

Lecture Notes
in
Population Genetics

Kent E. Holsinger
Department of Ecology & Evolutionary Biology, U-3043
University of Connecticut
Storrs, CT 06269-3043

© 2001-2012 Kent E. Holsinger

Creative Commons License

These notes are licensed under the Creative Commons Attribution-ShareAlike License. To view a copy of this license, visit <http://creativecommons.org/licenses/by-sa/3.0/> or send a letter to Creative Commons, 559 Nathan Abbott Way, Stanford, California 94305, USA.

Contents

Preface	vii
I The genetic structure of populations	1
1 Genetic transmission in populations	3
2 The Hardy-Weinberg Principle and estimating allele frequencies	7
3 Inbreeding and self-fertilization	19
4 Testing Hardy-Weinberg	27
5 Wahlund effect, Wright's F-statistics	33
6 Analyzing the genetic structure of populations	39
7 Analyzing the genetic structure of populations: a Bayesian approach	49
8 Analyzing the genetic structure of populations: individual assignment	57
9 Two-locus population genetics	61
II The genetics of natural selection	69
10 The Genetics of Natural Selection	71
11 Estimating viability	83

12 Selection at one locus with many alleles, fertility selection, and sexual selection	87
13 Selection Components Analysis	95
III Genetic drift	101
14 Genetic Drift	103
15 Mutation, Migration, and Genetic Drift	117
16 Selection and genetic drift	123
17 The Coalescent	127
IV Quantitative genetics	133
18 Introduction to quantitative genetics	135
19 Resemblance among relatives	147
20 Partitioning variance with WinBUGS	157
21 Evolution of quantitative traits	161
22 Selection on multiple characters	169
23 Mapping quantitative trait loci	177
24 Mapping Quantitative Trait Loci with R/qtl	185
25 Association mapping: the background from two-locus genetics	193
26 Association mapping: BAMD	199
V Molecular evolution	203
27 Introduction to molecular population genetics	205

28	The neutral theory of molecular evolution	215
29	Patterns of nucleotide and amino acid substitution	221
30	Detecting selection on nucleotide polymorphisms	227
31	Patterns of selection on nucleotide polymorphisms	235
32	Tajima's D , Fu's F_S , Fay and Wu's H , and Zeng et al.'s E	239
33	Evolution in multigene families	245
34	Analysis of mismatch distributions	253
VI	Phylogeography	261
35	Analysis of molecular variance (AMOVA)	263
36	Nested clade analysis	271
37	Statistical phylogeography	281
38	Fully coalescent-based approaches to phylogeography	285
39	Approximate Bayesian Computation	293
40	Population genomics	301

Preface

Acknowledgments

I've used various versions of these notes in my graduate course on population genetics <http://darwin.eeb.uconn.edu/eeb348> since 2001. Some of them date back even earlier than that. Several generations of students and teaching assistants have found errors and helped me to find better ways of explaining arcane concepts. In addition, the following people have found various errors and helped me to correct them.

Brian Cady
Rachel Prunier
Uzay Sezen
Robynn Shannon

Jennifer Steinbachs
Kathryn Theiss
Yufeng Wu

I am indebted to everyone who has found errors or suggested better ways of explaining concepts, but don't blame them for any errors that are left. Those are all mine.

Part I

The genetic structure of populations

Chapter 1

Genetic transmission in populations

Mendel's rules describe how genetic transmission happens between parents and offspring. Consider a monohybrid cross:

$$\begin{array}{c} A_1A_2 \times A_1A_2 \\ \downarrow \\ \frac{1}{4}A_1A_1 \quad \frac{1}{2}A_1A_2 \quad \frac{1}{4}A_2A_2 \end{array}$$

Population genetics describes how genetic transmission happens between a *population* of parents and a population of offspring. Consider the following data from the *Est-3* locus of *Zoarces viviparus*:¹

Maternal genotype	Genotype of offspring		
	A_1A_1	A_1A_2	A_2A_2
A_1A_1	305	516	
A_1A_2	459	1360	877
A_2A_2		877	1541

This table describes, empirically, the relationship between the genotypes of mothers and the genotypes of their offspring. We can also make some inferences about the genotypes of the fathers in this population, even though we didn't see them.

1. 305 out of 821 male gametes that fertilized eggs from A_1A_1 mothers carried the A_1 allele (37%).
2. 877 out of 2418 male gametes that fertilized eggs from A_2A_2 mothers carried the A_1 allele (36%).

¹from [12]

Question How many of the 2,696 male gametes that fertilized eggs from A_1A_2 mothers carried the A_1 allele?

Recall We don't know the paternal genotypes or we wouldn't be asking this question.

- There is no way to tell which of the 1360 A_1A_2 offspring received A_1 from their mother and which from their father.
- Regardless of what the genotype of the father is, half of the offspring of a heterozygous mother will be heterozygous.²
- Heterozygous offspring of heterozygous mothers contain no information about the frequency of A_1 among fathers, so we don't bother to include them in our calculations.

Rephrase How many of the 1336 homozygous progeny of heterozygous mothers received an A_1 allele from their father?

Answer 459 out of 1336 (34%)

New question How many of the offspring where the paternal contribution can be identified received an A_1 allele from their father?

Answer (305 + 459 + 877) out of (305 + 459 + 877 + 516 + 877 + 1541) or 1641 out of 4575 (36%)

An algebraic formulation of the problem

The above calculations tell us what's happening for this particular data set, but those of you who know me know that there has to be a little math coming to describe the situation more generally. Here it is:

Genotype	Number	Sex
A_1A_1	F_{11}	female
A_1A_2	F_{12}	female
A_2A_2	F_{22}	female
A_1A_1	M_{11}	male
A_1A_2	M_{12}	male
A_2A_2	M_{22}	male

²Assuming we're looking at data from a locus that has only two alleles. If there were four alleles at a locus, for example, *all* of the offspring might be heterozygous.

then

$$p_f = \frac{2F_{11}+F_{12}}{2F_{11}+2F_{12}+2F_{22}} \quad q_f = \frac{2F_{22}+F_{12}}{2F_{11}+2F_{12}+2F_{22}}$$
$$p_m = \frac{2M_{11}+M_{12}}{2M_{11}+2M_{12}+2M_{22}} \quad q_m = \frac{2M_{22}+M_{12}}{2M_{11}+2M_{12}+2M_{22}} \quad ,$$

where p_f is the frequency of A_1 in mothers and p_m is the frequency of A_1 in fathers.³

Since every individual in the population must have one father and one mother, the frequency of A_1 among offspring is the same in both sexes, namely

$$p = \frac{1}{2}(p_f + p_m) \quad ,$$

assuming that all matings have the same average fecundity and that the locus we're studying is autosomal.⁴

Question: Why do those assumptions matter?

Answer: If $p_f = p_m$, then the allele frequency among offspring is equal to the allele frequency in their parents, i.e., the allele frequency doesn't change from one generation to the next. This might be considered the First Law of Population Genetics: If no forces act to change allele frequencies between zygote formation and breeding, allele frequencies will not change.

Zero force laws

This is an example of what philosophers call a **zero force law**. Zero force laws play a very important role in scientific theories, because we can't begin to understand what a force does until we understand what would happen in the absence of any forces. Consider Newton's famous dictum:

An object in motion tends to remain in motion in a straight line. An object at rest tends to remain at rest.

or (as you may remember from introductory physics)⁵

$$F = ma \quad .$$

³ $q_f = 1 - p_f$ and $q_m = 1 - p_m$ as usual.

⁴And that there are enough offspring produced that we can ignore genetic drift. Have you noticed that I have a fondness for footnotes? You'll see a lot more before the semester is through, and you'll soon discover that most of my weak attempts at humor are buried in them.

⁵Don't worry if you're not good at physics. I'm probably worse. What I'm about to tell you is almost the only thing about physics I can remember.

If we observe an object accelerating, we can immediately infer that a force is acting on it, and we can infer something about the magnitude of that force. **However**, if an object is not accelerating we cannot conclude that no forces are acting. It might be that opposing forces act on the object in such a way that the resultant is no *net* force. Acceleration is a *sufficient* condition to infer that force is operating on an object, but it is not *necessary*.

What we might call the “First Law of Population Genetics” is analogous to Newton’s First Law of Motion:

If all genotypes at a particular locus have the same average fecundity and the same average chance of being included in the breeding population, allele frequencies in the population will remain constant.

For the rest of the semester we’ll be learning about the forces that cause allele frequencies to change and learning how to infer the properties of those forces from the changes that they induce. But you must always remember that while we can infer that some evolutionary force is present if allele frequencies change from one generation to the next, we *cannot* infer the absence of a force from a lack of allele frequency change.

Chapter 2

The Hardy-Weinberg Principle and estimating allele frequencies

To keep things relatively simple, we'll spend much of our time in this course talking about variation at a single genetic locus, even though alleles at many different loci are involved in expression of most morphological or physiological traits. We'll spend about three weeks in mid-October studying the genetics of quantitative variation, but until then you can assume that I'm talking about variation at a single locus unless I specifically say otherwise.

The genetic composition of populations

When I talk about the genetic composition of a population, I'm referring to three aspects of variation within that population:¹

1. The number of alleles at a locus.
2. The frequency of alleles at the locus.
3. The frequency of genotypes at the locus.

It may not be immediately obvious why we need both (2) and (3) to describe the genetic composition of a population, so let me illustrate with two hypothetical populations:

	A_1A_1	A_1A_2	A_2A_2
Population 1	50	0	50
Population 2	25	50	25

¹At each locus I'm talking about. Remember, I'm only talking about one locus at a time, unless I specifically say otherwise. We'll see why this matters when we get to two-locus genetics in a few weeks.

It's easy to see that the frequency of A_1 is 0.5 in both populations,² but the genotype frequencies are very different. In point of fact, we don't need both genotype and allele frequencies. We can always calculate allele frequencies from genotype frequencies, but we can't do the reverse unless ...

Derivation of the Hardy-Weinberg principle

We saw last time using the data from *Zoarces viviparus* that we can describe empirically and algebraically how genotype frequencies in one generation are related to genotype frequencies in the next. Let's explore that a bit further. To do so we're going to use a technique that is broadly useful in population genetics, i.e., we're going to construct a mating table. A mating table consists of three components:

1. A list of all possible genotype pairings.
2. The frequency with which each genotype pairing occurs.
3. The genotypes produced by each pairing.

Mating	Frequency	Offspring genotype		
		A_1A_1	A_1A_2	A_2A_2
$A_1A_1 \times A_1A_1$	x_{11}^2	1	0	0
A_1A_2	$x_{11}x_{12}$	$\frac{1}{2}$	$\frac{1}{2}$	0
A_2A_2	$x_{11}x_{22}$	0	1	0
$A_1A_2 \times A_1A_1$	$x_{12}x_{11}$	$\frac{1}{2}$	$\frac{1}{2}$	0
A_1A_2	x_{12}^2	$\frac{1}{4}$	$\frac{1}{2}$	$\frac{1}{4}$
A_2A_2	$x_{12}x_{22}$	0	$\frac{1}{2}$	$\frac{1}{2}$
$A_2A_2 \times A_1A_1$	$x_{22}x_{11}$	0	1	0
A_1A_2	$x_{22}x_{12}$	0	$\frac{1}{2}$	$\frac{1}{2}$
A_2A_2	x_{22}^2	0	0	1

Believe it or not, in constructing this table we've already made three assumptions about the transmission of genetic variation from one generation to the next:

Assumption #1 Genotype frequencies are the same in males and females, e.g., x_{11} is the frequency of the A_1A_1 genotype in both males and females.³

² $p_1 = 2(50)/200 = 0.5$, $p_2 = (2(25) + 50)/200 = 0.5$.

³It would be easy enough to relax this assumption, but it makes the algebra more complicated without providing any new insight, so we won't bother with relaxing it unless someone asks.

Assumption #2 Genotypes mate at random *with respect to their genotype at this particular locus*.

Assumption #3 Meiosis is fair. More specifically, we assume that there is no segregation distortion; no gamete competition; no differences in the developmental ability of eggs, or the fertilization ability of sperm.⁴ It may come as a surprise to you, but there are alleles at some loci in some organisms that subvert the Mendelian rules, e.g., the *t* allele in house mice, segregation distorter in *Drosophila melanogaster*, and spore killer in *Neurospora crassa*. A pair of paper describing the most recent work in *Neurospora* just appeared in July [33, 79].

Now that we have this table we can use it to calculate the frequency of each genotype in newly formed zygotes in the population,⁵ provided that we're willing to make three additional assumptions:

Assumption #4 There is no input of new genetic material, i.e., gametes are produced without mutation, and all offspring are produced from the union of gametes within this population, i.e., no migration from outside the population.

Assumption #5 The population is of infinite size so that the actual frequency of matings is equal to their expected frequency and the actual frequency of offspring from each mating is equal to the Mendelian expectations.

Assumption #6 All matings produce the same number of offspring, on average.

Taking these three assumptions together allows us to conclude that the frequency of a particular genotype in the pool of newly formed zygotes is

$$\sum (\text{frequency of mating})(\text{frequency of genotype produce from mating}) \quad .$$

So

$$\begin{aligned} \text{freq.}(A_1A_1 \text{ in zygotes}) &= x_{11}^2 + \frac{1}{2}x_{11}x_{12} + \frac{1}{2}x_{12}x_{11} + \frac{1}{4}x_{12}^2 \\ &= x_{11}^2 + x_{11}x_{12} + \frac{1}{4}x_{12}^2 \end{aligned}$$

⁴We are also assuming that we're looking at offspring genotypes at the zygote stage, so that there hasn't been any opportunity for differential survival.

⁵Not just the offspring from these matings

$$\begin{aligned}
&= (x_{11} + x_{12}/2)^2 \\
&= p^2 \\
\text{freq.}(A_1A_2 \text{ in zygotes}) &= 2pq \\
\text{freq.}(A_2A_2 \text{ in zygotes}) &= q^2
\end{aligned}$$

Those frequencies probably look pretty familiar to you. They are, of course, the familiar Hardy-Weinberg proportions. But we're not done yet. In order to say that these proportions will also be the genotype proportions of adults in the progeny generation, we have to make two more assumptions:

Assumption #7 Generations do not overlap.

Assumption #8 There are no differences among genotypes in the probability of survival.

The Hardy-Weinberg principle

After a single generation in which *all* eight of the above assumptions are satisfied

$$\text{freq.}(A_1A_1 \text{ in zygotes}) = p^2 \tag{2.1}$$

$$\text{freq.}(A_1A_2 \text{ in zygotes}) = 2pq \tag{2.2}$$

$$\text{freq.}(A_2A_2 \text{ in zygotes}) = q^2 \tag{2.3}$$

It's vital to understand the logic here.

1. If Assumptions #1–#8 are true, then equations 5.4–5.6 **must** be true.
2. If genotypes are in Hardy-Weinberg proportions, one or more of Assumptions #1–#8 may still be violated.
3. If genotypes are *not* in Hardy-Weinberg proportions, one or more of Assumptions #1–#8 **must** be false.
4. Assumptions #1–#8 are *sufficient* for Hardy-Weinberg to hold, but they are not *necessary* for Hardy-Weinberg to hold.

Point (3) is why the Hardy-Weinberg principle is so important. There isn't a population of any organism anywhere in the world that satisfies all 8 assumptions, even for a single generation.⁶ But *all* possible evolutionary forces within populations cause a violation of at least one of these assumptions. Departures from Hardy-Weinberg are one way in which we can detect those forces and estimate their magnitude.⁷

Estimating allele frequencies

Before we can determine whether genotypes in a population are in Hardy-Weinberg proportions, we need to be able to estimate the frequency of both genotypes and alleles. This is easy when you can identify all of the alleles within genotypes, but suppose that we're trying to estimate allele frequencies in the ABO blood group system in humans. Then we have a situation that looks like this:

Phenotype	A	AB	B	O
Genotype(s)	aa ao	ab	bb bo	oo
No. in sample	N_A	N_{AB}	N_B	N_O

Now we can't directly count the number of a , b , and o alleles. What do we do? Well, more than 50 years ago, some geneticists figured out how with a method they called "gene counting" [11] and that statisticians later generalized for a wide variety of purposes and called the EM algorithm [17]. It uses a trick you'll see repeatedly through this course. When we don't know something we want to know, we pretend that we know it and do some calculations with it. If we're lucky, we can fiddle with our calculations a bit to relate the thing that we pretended to know to something we actually do know so we can figure out what we wanted to know. Make sense? Probably not. But let's try an example.

If we knew p_a , p_b , and p_o , we could figure out how many individuals with the A phenotype have the aa genotype and how many have the ao genotype, namely

$$\begin{aligned} N_{aa} &= n_A \left(\frac{p_a^2}{p_a^2 + 2p_a p_o} \right) \\ N_{ao} &= n_A \left(\frac{2p_a p_o}{p_a^2 + 2p_a p_o} \right) . \end{aligned}$$

⁶There may be some that come reasonably close, but none that fulfill them *exactly*. There aren't any populations of infinite size, for example.

⁷Actually, there's a ninth assumption that I didn't mention. Everything I said here depends on the assumption that the locus we're dealing with is autosomal. We can talk about what happens with sex-linked loci, if you want. But again, mostly what we get is algebraic complications without a lot of new insight.

Obviously we could do the same thing for the B phenotype:

$$\begin{aligned} N_{bb} &= n_B \left(\frac{p_b^2}{p_b^2 + 2p_b p_o} \right) \\ N_{bo} &= n_B \left(\frac{2p_b p_o}{p_b^2 + 2p_b p_o} \right) . \end{aligned}$$

Notice that $N_{ab} = N_{AB}$ and $N_{oo} = N_O$ (lowercase subscripts refer to genotypes, uppercase to phenotypes). If we knew all this, then we could calculate p_a , p_b , and p_o from

$$\begin{aligned} p_a &= \frac{2N_{aa} + N_{ao} + N_{ab}}{2N} \\ p_b &= \frac{2N_{bb} + N_{bo} + N_{ab}}{2N} \\ p_o &= \frac{2N_{oo} + N_{ao} + N_{bo}}{2N} , \end{aligned}$$

where N is the total sample size.

Surprisingly enough we can actually estimate the allele frequencies by using this trick. Just take a guess at the allele frequencies. Any guess will do. Then calculate N_{aa} , N_{ao} , N_{bb} , N_{bo} , N_{ab} , and N_{oo} as described in the preceding paragraph.⁸ That's the **Expectation** part of the EM algorithm. Now take the values for N_{aa} , N_{ao} , N_{bb} , N_{bo} , N_{ab} , and N_{oo} that you've calculated and use them to calculate new values for the allele frequencies. That's the **Maximization** part of the EM algorithm. It's called "maximization" because what you're doing is calculating maximum-likelihood estimates of the allele frequencies, given the observed (and made up) genotype counts.⁹ Chances are your new values for p_a , p_b , and p_o won't match your initial guesses, but¹⁰ if you take these new values and start the process over and repeat the whole sequence several times, eventually the allele frequencies you get out at the end match those you started with. These are maximum-likelihood estimates of the allele frequencies.¹¹

Consider the following example:¹²

Phenotype	A	AB	AB	O
No. in sample	25	50	25	15

⁸Chances are N_{aa} , N_{ao} , N_{bb} , and N_{bo} won't be integers. That's OK. Pretend that there really are fractional animals or plants in your sample and proceed.

⁹If you don't know what maximum-likelihood estimates are, don't worry. We'll get to that in a moment.

¹⁰Yes, truth *is* sometimes stranger than fiction.

¹¹I should point out that this method *assumes* that genotypes are found in Hardy-Weinberg proportions.

¹²This is the default example available in the Java applet at <http://darwin.eeb.uconn.edu/simulations/em-abo.html>.

We'll start with the guess that $p_a = 0.33$, $p_b = 0.33$, and $p_o = 0.34$. With that assumption we would calculate that $25(0.33^2/(0.33^2 + 2(0.33)(0.34))) = 8.168$ of the A phenotypes in the sample have genotype aa , and the remaining 16.832 have genotype ao . Similarly, we can calculate that 8.168 of the B phenotypes in the population sample have genotype bb , and the remaining 16.823 have genotype bo . Now that we have a guess about how many individuals of each genotype we have,¹³ we can calculate a new guess for the allele frequencies, namely $p_a = 0.362$, $p_b = 0.362$, and $p_o = 0.277$. By the time we've repeated this process four more times, the allele frequencies aren't changing anymore. So the maximum likelihood estimate of the allele frequencies is $p_a = 0.372$, $p_b = 0.372$, and $p_o = 0.256$.

What is a maximum-likelihood estimate?

I just told you that the method I described produces “maximum-likelihood estimates” for the allele frequencies, but I haven't told you what a maximum-likelihood estimate is. The good news is that you've been using maximum-likelihood estimates for as long as you've been estimating anything, without even knowing it. Although it will take me awhile to explain it, the idea is actually pretty simple.

Suppose we had a sock drawer with two colors of socks, red and green. And suppose we were interested in estimating the proportion of red socks in the drawer. One way of approaching the problem would be to mix the socks well, close our eyes, take one sock from the drawer, record its color and replace it. Suppose we do this N times. We know that the number of red socks we'll get might be different the next time, so the number of red socks we get is a random variable. Let's call it K . Now suppose in our actual experiment we find k red socks, i.e., $K = k$. If we knew p , the proportion of red socks in the drawer, we could calculate the probability of getting the data we observed, namely

$$P(K = k|p) = \binom{N}{k} p^k (1 - p)^{(N-k)} \quad . \quad (2.4)$$

This is the *binomial probability distribution*. The part on the left side of the equation is read as “The probability that we get k red socks in our sample *given* the value of p .” The word “given” means that we're calculating the probability of our data conditional on the (unknown) value p .

Of course we don't know p , so what good does writing (2.4) do? Well, suppose we reverse the question to which equation (2.4) is an answer and call the expression in (2.4) the “likelihood of the data.” Suppose further that we find the value of p that makes the

¹³Since we're making these genotype counts up, we can also pretend that it makes sense to have fractional numbers of genotypes.

likelihood bigger than any other value we could pick.¹⁴ Then \hat{p} is the maximum-likelihood estimate of p .¹⁵

In the case of the ABO blood group that we just talked about, the likelihood is a bit more complicated

$$\binom{N}{N_A N_{AB} N_B N_O} (p_a^2 + 2p_a p_o)^{N_A} 2p_a p_b^{N_{AB}} (p_b^2 + 2p_b p_o)^{N_B} (p_o^2)^{N_O} \quad (2.5)$$

This is a *multinomial probability distribution*. It turns out that one way to find the values of p_a , p_b , and p_o is to use the EM algorithm I just described.¹⁶

An introduction to Bayesian inference

Maximum-likelihood estimates have a lot of nice features, but likelihood is a slightly backwards way of looking at the world. The likelihood of the data is the probability of the data, x , given parameters that we don't know, ϕ , i.e., $P(x|\phi)$. It seems a lot more natural to think about the probability that the unknown parameter takes on some value, given the data, i.e., $P(\phi|x)$. Surprisingly, these two quantities are closely related. Bayes' Theorem tells us that

$$P(\phi|x) = \frac{P(x|\phi)P(\phi)}{P(x)} \quad (2.6)$$

We refer to $P(\phi|x)$ as the *posterior distribution* of ϕ , i.e., the probability that ϕ takes on a particular value given the data we've observed, and to $P(\phi)$ as the *prior distribution* of ϕ , i.e., the probability that ϕ takes on a particular value *before* we've looked at any data. Notice how the relationship in (2.6) mimics the logic we use to learn about the world in everyday life. We start with some prior beliefs, $P(\phi)$, and modify them on the basis of data or experience, $P(x|\phi)$, to reach a conclusion, $P(\phi|x)$. That's the underlying logic of Bayesian inference.¹⁷

¹⁴Technically, we treat $P(K = k|p)$ as a function of p , find the value of p that maximizes it, and call that value \hat{p} .

¹⁵You'll be relieved to know that in this case, $\hat{p} = k/N$.

¹⁶There's another way I'd be happy to describe if you're interested, but it's a lot more complicated.

¹⁷If you'd like a little more information on why a Bayesian approach makes sense, you might want to take a look at my lecture notes from the Summer Institute in Statistical Genetics.

Estimating allele frequencies with two alleles

Let's suppose we've collected data from a population of *Desmodium cuspidatum*¹⁸ and have found 7 alleles coding for the *fast* allele at a enzyme locus encoding glucose-phosphate isomerase in a sample of 20 alleles. We want to estimate the frequency of the *fast* allele. The maximum-likelihood estimate is $7/20 = 0.35$, which we got by finding the value of p that maximizes

$$P(p|N, k) = \binom{N}{k} p^k (1-p)^{N-k} ,$$

where $N = 20$ and $k = 7$. A Bayesian uses the same likelihood, but has to specify a prior distribution for p . If we didn't know anything about the allele frequency at this locus in *D. cuspidatum* before starting the study, it makes sense to express that ignorance by choosing $P(p)$ to be a uniform random variable on the interval $[0, 1]$. That means we regarded all values of p as equally likely prior to collecting the data.¹⁹

Until a little over fifteen years ago it was necessary to do a bunch of complicated calculus to combine the prior with the likelihood to get a posterior. Since the early 1990s statisticians have used a simulation approach, Monte Carlo Markov Chain sampling, to construct numerical samples from the posterior. For the problems encountered in this course, we'll mostly be using the freely available software package WinBUGS to implement Bayesian analyses. For the problem we just encountered, here's the code that's needed to get our results:²⁰

```
model {  
  
  # likelihood  
  k ~ dbin(p, N)  
  
  # prior  
  p ~ dunif(0,1)  
  
}  
  
list(k = 7, n = 20)
```

¹⁸A few of you may recognize that I didn't choose that species entirely at random, even though the "data" are entirely fanciful.

¹⁹If we had prior information about the likely values of p , we'd pick a different prior distribution to reflect our prior information. See the Summer Institute notes for more information, if you're interested.

²⁰This code and other WinBUGS code used in the course can be found on the course web site by following the links associated with the corresponding lecture.

node	mean	sd	MC error	2.5%	median	97.5%	start	sample
p	0.3639	0.1004	0.001293	0.1841	0.3599	0.5697	1001	5000

Figure 2.1: Results of a WinBUGS analysis with the made up allele count data from *Desmodium cuspidatum*.

Running this in WinBUGS produces the results in Figure 2.1.

The column headings in Figure 2.1 should be fairly self-explanatory, except for the one labeled **MC error**.²¹ **mean** is the posterior mean. It's our best guess of the value for the frequency of the *fast* allele. **s.d.** is the posterior standard deviation. It's our best guess of the uncertainty associated with our estimate of the frequency of the *fast* allele. The 2.5%, 50%, and 97.5% columns are the percentiles of the posterior distribution. The [2.5%, 97.5%] interval is the 95% credible interval, which is analogous to the 95% confidence interval in classical statistics, except that we can say that there's a 95% chance that the frequency of the *fast* allele lies within this interval.²² Since the results are from a simulation, different runs will produce slightly different results. In this case, we have a posterior mean of about 0.36 (as opposed to the maximum-likelihood estimate of 0.35), and there is a 95% chance that p lies in the interval [0.18, 0.57].²³

Returning to the ABO example

Here's data from the ABO blood group:²⁴

Phenotype	A	AB	B	O	Total
Observed	862	131	365	702	2060

²¹If you're interested in what **MC error** means, ask. Otherwise, I don't plan to say anything about it.

²²If you don't understand why that's different from a standard confidence interval, ask me about it.

²³See the Summer Institute notes for more details on why the Bayesian estimate of p is different from the maximum-likelihood estimate. Suffice it to say that when you have a reasonable amount of data, the estimates are barely distinguishable.

²⁴This is almost the last time! I promise.

To estimate the underlying allele frequencies, p_A , p_B , and p_O , we have to remember how the allele frequencies map to phenotype frequencies:²⁵

$$\begin{aligned}\text{Freq}(A) &= p_A^2 + 2p_A p_O \\ \text{Freq}(AB) &= 2p_A p_B \\ \text{Freq}(B) &= p_B^2 + 2p_B p_O \\ \text{Freq}(O) &= p_O^2 \quad .\end{aligned}$$

Hers's the WinBUGS code we use to estimate the allele frequencies:

```
model {
  # likelihood
  pi[1] <- p.a*p.a + 2*p.a*p.o
  pi[2] <- 2*p.a*p.b
  pi[3] <- p.b*p.b + 2*p.b*p.o
  pi[4] <- p.o*p.o
  x[1:4] ~ dmulti(pi[],n)

  # priors
  a1 ~ dexp(1)
  b1 ~ dexp(1)
  o1 ~ dexp(1)
  p.a <- a1/(a1 + b1 + o1)
  p.b <- b1/(a1 + b1 + o1)
  p.o <- o1/(a1 + b1 + o1)

  n <- sum(x[])
}
```

```
list(x=c(862, 131, 365, 702))
```

The `dmulti()` is a multinomial probability, a simple generalization of the binomial probability to samples when there are more than two categories. The priors are some mumbo jumbo necessary to produce the rough equivalent of uniform $[0,1]$ priors with more than two alleles.²⁶ `sum()` is a built-in function that saves me the trouble of calculating the sample size and ensures that the `n` in `dmulti()` is consistent with the individual sample components.

²⁵Assuming genotypes are in Hardy-Weinberg proportions. We'll relax that assumption later.

²⁶It produces a Dirichlet(1,1,1), if you really want to know.

The image shows a screenshot of the WinBUGS 'Node statistics' window. The window title is 'Node statistics' and it contains a table with the following data:

node	mean	sd	MC error	2.5%	median	97.5%	start	sample
p.a	0.2813	0.007668	1.237E-4	0.2664	0.2812	0.2962	1001	5000
p.b	0.1293	0.005405	5.948E-5	0.119	0.1292	0.1404	1001	5000
p.o	0.5894	0.008327	1.242E-4	0.5734	0.5894	0.6059	1001	5000

Figure 2.2: Results of a WinBUGS analysis of the ABO data.

The `x=c()` produces a vector of counts arranged in the same order as the frequencies in `pi []`. The results are in Figure 2.2. Notice that the posterior means are very close to the maximum-likelihood estimates, but that we also have 95% credible intervals so that we have an assessment of how reliable the Bayesian estimates are. Getting them from a likelihood analysis is possible, but it takes a fair amount of additional work.

Chapter 3

Inbreeding and self-fertilization

Remember that long list of assumptions associated with derivation of the Hardy-Weinberg principle that I went over a couple of lectures ago? Well, we're about to begin violating assumptions to explore the consequences, but we're not going to violate them in order. We're first going to violate Assumption #2:

Genotypes mate at random with respect to their genotype at this particular locus.

There are many ways in which this assumption might be violated:

- Some genotypes may be more successful in mating than others — sexual selection.
- Genotypes that are different from one another may mate more often than expected — disassortative mating, e.g., self-incompatibility alleles in flowering plants, MHC loci in humans (the smelly t-shirt experiment) [95].
- Genotypes that are similar to one another may mate more often than expected — assortative mating.
- Some fraction of the offspring produced may be produced asexually.
- Individuals may mate with relatives — inbreeding.
 - self-fertilization
 - sib-mating
 - first-cousin mating
 - parent-offspring mating

– etc.

When there is sexual selection or disassortative mating genotypes differ in their chances of being included in the breeding population. As a result, allele and genotype frequencies will tend to change from one generation to the next. We'll talk a little about these types of departures from random mating when we discuss the genetics of natural selection in a few weeks, but we'll ignore them for now. In fact, we'll also ignore assortative mating, since its properties are fairly similar to those of inbreeding, and inbreeding is easier to understand.

Self-fertilization

Self-fertilization is the most extreme form of inbreeding possible, and it is characteristic of many flowering plants and some hermaphroditic animals, including freshwater snails and that darling of developmental genetics, *Caenorhabditis elegans*.¹ It's not too hard to figure out what the consequences of self-fertilization will be without doing any algebra.

- All progeny of homozygotes are themselves homozygous.
- Half of the progeny of heterozygotes are heterozygous and half are homozygous.

So you might expect that the frequency of heterozygotes would be halved every generation, and you'd be right. To see why, consider the following mating table:

Mating	frequency	Offspring genotype		
		A_1A_1	A_1A_2	A_2A_2
$A_1A_1 \times A_1A_1$	x_{11}	1	0	0
$A_1A_2 \times A_1A_2$	x_{12}	$\frac{1}{4}$	$\frac{1}{2}$	$\frac{1}{4}$
$A_2A_2 \times A_2A_2$	x_{22}	0	0	1

Using the same technique we used to derive the Hardy-Weinberg principle, we can calculate the frequency of the different offspring genotypes from the above table.

¹It may well be characteristic of many hermaphroditic animal parasites. You should also know that I just lied. The form of self-fertilization I'm going to describe actually isn't the most extreme form of selfing possible. That honor belongs to gametophytic self-fertilization in homosporous plants. The offspring of gametophytic self-fertilization are uniformly homozygous at every locus in the genome. For more information see [41]

$$x'_{11} = x_{11} + x_{12}/4 \quad (3.1)$$

$$x'_{12} = x_{12}/2 \quad (3.2)$$

$$x'_{22} = x_{22} + x_{12}/4 \quad (3.3)$$

I use the ' to indicate the next generation. Notice that in making this calculation I assume that all other conditions associated with Hardy-Weinberg apply (meiosis is fair, no differences among genotypes in probability of survival, no input of new genetic material, etc.). We can also calculate the frequency of the A_1 allele among offspring, namely

$$p' = x'_{11} + x'_{12}/2 \quad (3.4)$$

$$= x_{11} + x_{12}/4 + x_{12}/4 \quad (3.5)$$

$$= x_{11} + x_{12}/2 \quad (3.6)$$

$$= p \quad (3.7)$$

These equations illustrate two very important principles that are true with any system of strict inbreeding:

1. Inbreeding does not cause allele frequencies to change, but it will generally cause genotype frequencies to change.
2. Inbreeding reduces the frequency of heterozygotes relative to Hardy-Weinberg expectations. It need not eliminate heterozygotes entirely, but it is guaranteed to reduce their frequency.
 - Suppose we have a population of hermaphrodites in which $x_{12} = 0.5$ and we subject it to strict self-fertilization. Assuming that inbred progeny are as likely to survive and reproduce as outbred progeny, $x_{12} < 0.01$ in six generations and $x_{12} < 0.0005$ in ten generations.

Partial self-fertilization

Many plants reproduce by a mixture of outcrossing and self-fertilization. To a population geneticist that means that they reproduce by a mixture of selfing and random mating. Now I'm going to pull a fast one and derive the equations that determine how allele frequencies change from one generation to the next without using a mating table. To do so, I'm going

to imagine that our population consists of a mixture of two populations. In one part of the population all of the reproduction occurs through self-fertilization and in the other part all of the reproduction occurs through random mating. If you think about it for a while, you'll realize that this is equivalent to imagining that each plant reproduces some fraction of the time through self-fertilization and some fraction of the time through random mating. Let σ be the fraction of progeny produced through self-fertilization, then

$$x'_{11} = p^2(1 - \sigma) + (x_{11} + x_{12}/4)\sigma \quad (3.8)$$

$$x'_{12} = 2pq(1 - \sigma) + (x_{12}/2)\sigma \quad (3.9)$$

$$x'_{22} = q^2(1 - \sigma) + (x_{22} + x_{12}/4)\sigma \quad (3.10)$$

Notice that I use p^2 , $2pq$, and q^2 for the genotype frequencies in the part of the population that's mating at random. **Question:** Why can I get away with that?²

It takes a little more algebra than it did before, but it's not difficult to verify that the allele frequencies don't change between parents and offspring.

$$p' = \{p^2(1 - \sigma) + (x_{11} + x_{12}/4)\sigma\} + \{pq(1 - \sigma) + (x_{12}/4)\sigma\} \quad (3.11)$$

$$= p(p + q)(1 - \sigma) + (x_{11} + x_{12}/2)\sigma \quad (3.12)$$

$$= p(1 - \sigma) + p\sigma \quad (3.13)$$

$$= p \quad (3.14)$$

Because homozygous parents can always have heterozygous offspring (when they out-cross), heterozygotes are never completely eliminated from the population as they are with complete self-fertilization. In fact, we can solve for the *equilibrium* frequency of heterozygotes, i.e., the frequency of heterozygotes reached when genotype frequencies stop changing.³ By definition, an equilibrium for x_{12} is a value such that if we put it in on the right side of equation (3.9) we get it back on the left side, or in equations

$$\hat{x}_{12} = 2pq(1 - \sigma) + (\hat{x}_{12}/2)\sigma \quad (3.15)$$

$$\hat{x}_{12}(1 - \sigma/2) = 2pq(1 - \sigma) \quad (3.16)$$

$$\hat{x}_{12} = \frac{2pq(1 - \sigma)}{(1 - \sigma/2)} \quad (3.17)$$

²If you're being good little boys and girls and looking over these notes *before* you get to class, when you see **Question** in the notes, you'll know to think about that a bit, because I'm not going to give you the answer in the notes, I'm going to help you discover it during lecture.

³This is analogous to stopping the calculation and re-calculation of allele frequencies in the EM algorithm when the allele frequency estimates stop changing.

It's worth noting several things about this set of equations:

1. I'm using \hat{x}_{12} to refer to the equilibrium frequency of heterozygotes. I'll be using hats over variables to denote equilibrium properties throughout the course.⁴
2. I can solve for \hat{x}_{12} in terms of p because I know that p doesn't change. If p changed, the calculations wouldn't be nearly this simple.
3. The equilibrium is approached gradually (or asymptotically as mathematicians would say). A single generation of random mating will put genotypes in Hardy-Weinberg proportions (assuming all the other conditions are satisfied), but many generations may be required for genotypes to approach their equilibrium frequency with partial self-fertilization.

Inbreeding coefficients

Now that we've found an expression for \hat{x}_{12} we can also find expressions for \hat{x}_{11} and \hat{x}_{22} . The complete set of equations for the genotype frequencies with partial selfing are:

$$\hat{x}_{11} = p^2 + \frac{\sigma pq}{2(1 - \sigma/2)} \quad (3.18)$$

$$\hat{x}_{12} = 2pq - 2 \left(\frac{\sigma pq}{2(1 - \sigma/2)} \right) \quad (3.19)$$

$$\hat{x}_{22} = q^2 + \frac{\sigma pq}{2(1 - \sigma/2)} \quad (3.20)$$

Notice that all of those equations have a term $\sigma/(2(1 - \sigma/2))$. Let's call that f . Then we can save ourselves a little hassle by rewriting the above equations as:

$$\hat{x}_{11} = p^2 + fpq \quad (3.21)$$

$$\hat{x}_{12} = 2pq(1 - f) \quad (3.22)$$

$$\hat{x}_{22} = q^2 + fpq \quad (3.23)$$

Now you're going to have to stare at this a little longer, but notice that \hat{x}_{12} is the frequency of heterozygotes that we observe and $2pq$ is the frequency of heterozygotes we'd expect

⁴Unfortunately, I'll also be using hats to denote estimates of unknown parameters, as I did when discussing maximum-likelihood estimates of allele frequencies. I apologize for using the same notation to mean different things, but I'm afraid you'll have to get used to figuring out the meaning from the context. Believe me. Things are about to get a lot worse. Wait until I tell you how many different ways population geneticists use a parameter f that is commonly called the inbreeding coefficient.

under Hardy-Weinberg in this population if we were able to observe the genotype and allele frequencies without error. So

$$1 - f = \frac{\hat{x}_{12}}{2pq} \quad (3.24)$$

$$f = 1 - \frac{\hat{x}_{12}}{2pq} \quad (3.25)$$

$$= 1 - \frac{\text{observed heterozygosity}}{\text{expected heterozygosity}} \quad (3.26)$$

f is the inbreeding coefficient. When defined as $1 - (\text{observed heterozygosity})/(\text{expected heterozygosity})$ it can be used to measure the extent to which a particular population departs from Hardy-Weinberg expectations.⁵ When f is defined in this way, I refer to it as the *population inbreeding coefficient*.

But f can also be regarded as a function of a particular system of mating. With partial self-fertilization the population inbreeding coefficient when the population has reached equilibrium is $\sigma/(2(1 - \sigma/2))$. When regarded as the inbreeding coefficient predicted by a particular system of mating, I refer to it as the *equilibrium inbreeding coefficient*.

We'll encounter at least two more definitions for f once I've introduced idea of identity by descent.

Identity by descent

Self-fertilization is, of course, only one example of the general phenomenon of inbreeding—non-random mating in which individuals mate with close relatives more often than expected at random. We've already seen that the consequences of inbreeding can be described in terms of the inbreeding coefficient, f and I've introduced you to two ways in which f can be defined.⁶ I'm about to introduce you to one more.

Two alleles at a single locus are *identical by descent* if they are identical copies of the same allele in some earlier generation, i.e., both are copies that arose by DNA replication from the same ancestral sequence without any intervening mutation.

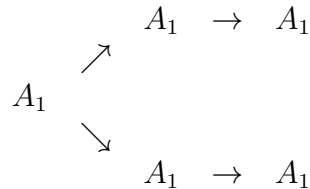
We're more used to classifying alleles by type than by descent. All though we don't usually say it explicitly, we regard two alleles as the "same," i.e., identical by type, if they

⁵ f can be negative if there are more heterozygotes than expected, as might be the case if cross-homozygote matings are more frequent than expected at random.

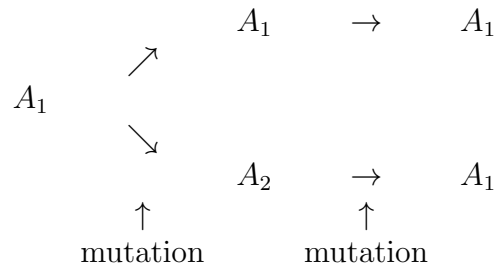
⁶See paragraphs above describing the population and equilibrium inbreeding coefficient.

have the same phenotypic effects. Whether or not two alleles are identical by descent, however, is a property of their genealogical history. Consider the following two scenarios:

Identity by descent



Identity by type



In both scenarios, the alleles at the end of the process are identical in type, i.e., they're both A_1 alleles. In the second scenario, however, they are identical in type only because one of the alleles has two mutations in its history.⁷ So alleles that are identical by descent will also be identical by type, but alleles that are identical by type need not be identical by descent.⁸

A third definition for f is the probability that two alleles *chosen at random* are identical by descent.⁹ Of course, there are several aspects to this definition that need to be spelled out more explicitly.

- In what sense are the alleles chosen at random, within an individual, within a particular population, within a particular set of populations?
- How far back do we trace the ancestry of alleles to determine whether they're identical by descent? Two alleles that are identical by type may not share a common ancestor if we trace their ancestry only 20 generations, but they may share a common ancestor if we trace their ancestry back 1000 generations and neither may have undergone any mutations since they diverged from one another.

⁷Notice that we could have had each allele mutate independently to A_2 .

⁸Systematists in the audience will recognize this as the problem of homoplasy.

⁹Notice that if we adopt this definition for f it can only take on values between 0 and 1. When used in the sense of a population or equilibrium inbreeding coefficient, however, f can be negative.

Let's imagine for a moment, however, that we've traced back the ancestry of all alleles in a particular population far enough to be able to say that if they're identical by type they're also identical by descent. Then we can write down the genotype frequencies in this population once we know f , where we define f as the probability that two alleles chosen at random in this population are identical by descent:

$$x_{11} = p^2(1 - f) + fp \tag{3.27}$$

$$x_{12} = 2pq(1 - f) \tag{3.28}$$

$$x_{22} = q^2(1 - f) + fq \quad . \tag{3.29}$$

It may not be immediately apparent, but you've actually seen these equations before in a different form. Since $p - p^2 = p(1 - p) = pq$ and $q - q^2 = q(1 - q) = pq$ these equations can be rewritten as

$$x_{11} = p^2 + fpq \tag{3.30}$$

$$x_{12} = 2pq(1 - f) \tag{3.31}$$

$$x_{22} = q^2 + fpq \quad . \tag{3.32}$$

You can probably see why population geneticists tend to play fast and loose with the definitions. *If* we ignore the distinction between identity by type and identity by descent, then the equations we used earlier to show the relationship between genotype frequencies, allele frequencies, and f (defined as a measure of departure from Hardy-Weinberg expectations) are identical to those used to show the relationship between genotype frequencies, allele frequencies, and f (defined as the probability that two randomly chosen alleles in the population are identical by descent).

Chapter 4

Testing Hardy-Weinberg

Because the Hardy-Weinberg principle tells us what to expect concerning the genetic composition of a sample when no evolutionary forces are operating, one of the first questions population geneticists often ask is “Are the genotypes in this sample present in the expected, i.e., Hardy-Weinberg, proportions?” We ask that question because we know that if the genotypes are *not* in Hardy-Weinberg proportions, at least one of the assumptions underlying derivation of the principle has been violated, i.e., that there is some evolutionary force operating on the population, and we know that we can use the magnitude and direction of the departure to say something about what those forces might be.

Of course we also know that the numbers in our sample may differ from expectation just because of random sampling error. For example, Table 4.1 presents data from a sample of 1000 English blood donors scored for MN phenotype. M and N are co-dominant, so that heterozygotes can be distinguished from the two homozygotes. Clearly the observed and expected numbers don’t look very different. The differences seem likely to be attributable purely to chance, but we need some way of assessing that “likeliness.”

Phenotype	Genotype	Observed Number	Expected Number
M	mm	298	294.3
MN	mn	489	496.3
N	nn	213	209.3

Table 4.1: Adapted from Table 2.4 in [37] (from [14])

Phenotype	A	AB	B	O	Total
Observed	862	131	365	702	2060

Table 4.2: Data on variation in ABO blood type.

Testing Hardy-Weinberg

One approach to testing the hypothesis that genotypes are in Hardy-Weinberg proportions is quite simple. We can simply do a χ^2 or G -test for goodness of fit between observed and predicted genotype (or phenotype) frequencies, where the predicted genotype frequencies are derived from our estimates of the allele frequencies in the population.¹ There's only one problem. To do either of these tests we have to know how many degrees of freedom are associated with the test. How do we figure that out? In general, the formula is

$$\begin{aligned} \text{d.f.} = & \quad (\# \text{ of categories in the data} - 1) \\ & - (\# \text{ number of parameters estimated from the data}) \end{aligned}$$

For this problem we have

$$\begin{aligned} \text{d.f.} = & \quad (\# \text{ of phenotype categories in the data} - 1) \\ & - (\# \text{ of allele frequencies estimated from the data}) \end{aligned}$$

In the ABO blood group we have 4 phenotype categories, and 3 allele frequencies. That means that a test of whether a particular data set has genotypes in Hardy-Weinberg proportions will have $(4 - 1) - (3 - 1) = 1$ degrees of freedom for the test. Notice that this also means that if you have completely dominant markers, like RAPDs or AFLPs, you can't determine whether genotypes are in Hardy-Weinberg proportions because you have 0 degrees of freedom available for the test.

An example

Table 4.2 exhibits data drawn from a study of phenotypic variation among individuals at the ABO blood locus:

¹If you're not familiar with the χ^2 or G -test for goodness of fit, consult any introductory statistics or biostatistics book, and you'll find a description. In fact, you probably don't have to go that far. You can probably find one in your old genetics textbook. Or you can just boot up your browser and head to Wikipedia: http://en.wikipedia.org/wiki/Goodness_of_fit.

The maximum-likelihood estimate of allele frequencies, assuming Hardy-Weinberg, is:²

$$\begin{aligned}p_a &= 0.281 \\p_b &= 0.129 \\p_o &= 0.590 \quad ,\end{aligned}$$

giving expected numbers of 846, 150, 348, and 716 for the four phenotypes. $\chi_1^2 = 3.8$, $0.05 < p < 0.1$.

A Bayesian approach

We saw last time how to use WinBUGS to provide allele frequency estimates from phenotypic data at the ABO locus.

```
model {
  # likelihood
  pi[1] <- p.a*p.a + 2*p.a*p.o
  pi[2] <- 2*p.a*p.b
  pi[3] <- p.b*p.b + 2*p.b*p.o
  pi[4] <- p.o*p.o
  x[1:4] ~ dmulti(pi[],n)

  # priors
  a1 ~ dexp(1)
  b1 ~ dexp(1)
  o1 ~ dexp(1)
  p.a <- a1/(a1 + b1 + o1)
  p.b <- b1/(a1 + b1 + o1)
  p.o <- o1/(a1 + b1 + o1)

  n <- sum(x[])
}
```

```
list(x=c(862, 131, 365, 702))
```

As you may recall, this produced the results in Figure 4.1.

²Take my word for it, or run the EM algorithm on these data yourself.

node	mean	sd	MC error	2.5%	median	97.5%	start	sample
p.a	0.2813	0.007668	1.237E-4	0.2664	0.2812	0.2962	1001	5000
p.b	0.1293	0.005405	5.948E-5	0.119	0.1292	0.1404	1001	5000
p.o	0.5894	0.008327	1.242E-4	0.5734	0.5894	0.6059	1001	5000

Figure 4.1: Results from WinBUGS analysis of the ABO data assuming genotypes are in Hardy-Weinberg proportions.

Now that we know about inbreeding coefficients and that they allow us to measure the departure of genotype frequencies from Hardy-Weinberg proportions, we can modify this a bit and estimate allele frequencies without assuming that genotypes are in Hardy-Weinberg proportions.

```

model {
  # likelihood
  pi[1] <- p.a*p.a + f*p.a*(1-p.a) + 2*p.a*p.o*(1-f)
  pi[2] <- 2*p.a*p.b*(1-f)
  pi[3] <- p.b*p.b + f*p.b*(1-p.b) + 2*p.b*p.o*(1-f)
  pi[4] <- p.o*p.o + f*p.o*(1-p.o)
  x[1:4] ~ dmulti(pi[],n)

  # priors
  a1 ~ dexp(1)
  b1 ~ dexp(1)
  o1 ~ dexp(1)
  p.a <- a1/(a1 + b1 + o1)
  p.b <- b1/(a1 + b1 + o1)
  p.o <- o1/(a1 + b1 + o1)

  f ~ dunif(0,1)

  n <- sum(x[])
}

list(x=c(862, 131, 365, 702))

```

This produces the results in Figure 4.2

node	mean	sd	MC error	2.5%	median	97.5%	start	sample
f	0.3991	0.1362	0.009539	0.07008	0.429	0.5907	1001	5000
p.a	0.3474	0.0258	0.001707	0.2905	0.3514	0.3885	1001	5000
p.b	0.1612	0.0139	8.712E-4	0.1321	0.1624	0.1861	1001	5000
p.o	0.4914	0.03765	0.002573	0.4327	0.4854	0.5758	1001	5000

Figure 4.2: Results from WinBUGS analysis of the ABO data relaxing the assumption that genotypes are in Hardy-Weinberg proportions.

Model	Dbar	Dhat	pD	DIC
$f > 0$	24.900	22.319	2.581	24.480
$f = 0$	27.827	25.786	2.041	29.869

Table 4.3: DIC calculations for the ABO example.

Notice that the allele frequency estimates have changed quite a bit and that the posterior mean of f is about 0.41. On first appearance, that would seem to indicate that we have lots of inbreeding in this sample. **BUT** it's a human population. It doesn't seem very likely that a human population is really that highly inbred.

Indeed, take a closer look at *all* of the information we have about that estimate of f . The 95% credible interval for f is between 0.06 and 0.55. That suggests that we don't have much information at all about f from these data.³ How can we tell if the model with inbreeding is better than the model that assumes genotypes are in Hardy-Weinberg proportions?

The Deviance Information Criterion

A widely used statistic for comparing models in a Bayesian framework is the Deviance Information Criterion. It can be calculated automatically in WinBUGS, just by clicking the right button. The results of the DIC calculations for our two models are summarized in Table 4.3.

Dbar and Dhat are measures of how well the model fits the data. Dbar is the posterior mean log likelihood, i.e., the average of the log likelihood values calculated from the parameters in each sample from the posterior. Dhat is the log likelihood at the posterior mean, i.e.,

³That shouldn't be too surprising, since any information we have about f comes indirectly through our allele frequency estimates.

the log likelihood calculated when all of the parameters are set to their posterior mean. pD is a measure of model complexity, roughly speaking the number of parameters in the model. DIC is a composite measure of how well the model does. It's a compromise between fit and complexity, and smaller DICs are preferred. A difference of more than 7-10 units is regarded as strong evidence in favor of the model with the smaller DIC.

In this case the difference in DIC values is about 5.5, so we have some evidence for $f > 0$ model for these data, even though they are from a human population. But the evidence is not very strong. This is consistent with the weak evidence for a departure from Hardy-Weinberg that was revealed in the χ^2 analysis.

Chapter 5

Wahlund effect, Wright's F-statistics

So far we've focused on inbreeding as one important way that populations may fail to mate at random, but there's another way in which virtually all populations and species fail to mate at random. Individuals tend to mate with those that are nearby. Even within a fairly small area, phenomena like nearest neighbor pollination in flowering plants or home-site fidelity in animals can cause mates to be selected in a geographically non-random way. What are the population genetic consequences of this form of non-random mating?

Well, if you think about it a little, you can probably figure it out. Since individuals that occur close to one another tend to be more genetically similar than those that occur far apart, the impacts of local mating will mimic those of inbreeding within a single, well-mixed population.

A numerical example

For example, suppose we have two subpopulations of green lacewings, one of which occurs in forests the other of which occurs in adjacent meadows. Suppose further that within each subpopulation mating occurs completely at random, but that there is no mating between forest and meadow individuals. Suppose we've determined allele frequencies in each population at a locus coding for phosphoglucose isomerase (*PGI*), which conveniently has only two alleles. The frequency of A_1 in the forest is 0.4 and in the meadow is 0.7. We can easily calculate the expected genotype frequencies within each population, namely

	A_1A_1	A_1A_2	A_2A_2
Forest	0.16	0.48	0.36
Meadow	0.49	0.42	0.09

Suppose, however, we were to consider a combined population consisting of 100 individuals from the forest subpopulation and 100 individuals from the meadow subpopulation. Then we'd get the following:¹

	A_1A_1	A_1A_2	A_2A_2
From forest	16	48	36
From meadow	49	42	9
Total	65	90	45

So the frequency of A_1 is $(2(65) + 90)/(2(65 + 90 + 45)) = 0.55$. Notice that this is just the average allele frequency in the two subpopulations, i.e., $(0.4 + 0.7)/2$. Since each subpopulation has genotypes in Hardy-Weinberg proportions, you might expect the combined population to have genotypes in Hardy-Weinberg proportions, but if you did you'd be wrong. Just look.

	A_1A_1	A_1A_2	A_2A_2
Expected (from $p = 0.55$)	$(0.3025)200$ 60.5	$(0.4950)200$ 99.0	$(0.2025)200$ 40.5
Observed (from table above)	65	90	45

The expected and observed don't match, even though there is random mating within both subpopulations. They don't match because there isn't random mating involving the combined population. Forest lacewings choose mates at random from other forest lacewings, but they never mate with a meadow lacewing (and *vice versa*). Our sample includes two populations that don't mix. This is an example of what's known as the *Wahlund effect* [94].

The algebraic development

You should know by now that I'm not going to be satisfied with a numerical example. I now feel the need to do some algebra to describe this situation a little more generally.

Suppose we know allele frequencies in k subpopulations.² Let p_i be the frequency of A_1 in the i th subpopulation. Then if we assume that all subpopulations contribute equally to combined population,³ we can calculate expected and observed genotype frequencies the way we did above:

¹If we ignore sampling error.

²For the time being, I'm going to assume that we know the allele frequencies without error, i.e., that we didn't have to estimate them from data. Next time we'll deal with real life, i.e., how we can detect the Wahlund effect when we have to *estimate* allele frequencies from data.

³We'd get the same result by relaxing this assumption, but the algebra gets messier, so why bother?

	A_1A_1	A_1A_2	A_2A_2
Expected	\bar{p}^2	$2\bar{p}\bar{q}$	\bar{q}^2
Observed	$\frac{1}{k} \sum p_i^2$	$\frac{1}{k} \sum 2p_iq_i$	$\frac{1}{k} \sum q_i^2$

where $\bar{p} = \sum p_i/k$ and $\bar{q} = 1 - \bar{p}$. Now

$$\frac{1}{k} \sum p_i^2 = \frac{1}{k} \sum (p_i - \bar{p} + \bar{p})^2 \quad (5.1)$$

$$= \frac{1}{k} \sum \left((p_i - \bar{p})^2 + 2\bar{p}(p_i - \bar{p}) + \bar{p}^2 \right) \quad (5.2)$$

$$= \frac{1}{k} \sum (p_i - \bar{p})^2 + \bar{p}^2 \quad (5.3)$$

$$= \text{Var}(p) + \bar{p}^2 \quad (5.4)$$

Similarly,

$$\frac{1}{k} \sum 2p_iq_i = 2\bar{p}\bar{q} - 2\text{Var}(p) \quad (5.5)$$

$$\frac{1}{k} \sum q_i^2 = \bar{q}^2 + \text{Var}(p) \quad (5.6)$$

Since $\text{Var}(p) \geq 0$ by definition, with equality holding only when all subpopulations have the same allele frequency, we can conclude that

- Homozygotes will be more frequent and heterozygotes will be less frequent than expected based on the allele frequency in the combined population.
- The magnitude of the departure from expectations is directly related to the magnitude of the variance in allele frequencies across populations, $\text{Var}(p)$.
- The effect will apply to *any* mixing of samples in which the subpopulations combined have different allele frequencies.⁴
- The same general phenomenon will occur if there are multiple alleles at a locus, although it is possible for one or a few heterozygotes to be *more* frequent than expected if there is positive covariance in the constituent allele frequencies across populations.⁵

⁴For example, if we combine samples from different years or across age classes of long-lived organisms, we may see a deficiency of heterozygotes in the sample purely as a result of allele frequency differences across years.

⁵If you're curious about this, feel free to ask, but I'll have to dig out my copy of Li [61] to answer. I don't carry those details around in my head.

- The effect is analogous to inbreeding. Homozygotes are more frequent and heterozygotes are less frequent than expected.⁶

To return to our earlier numerical example:

$$\text{Var}(p) = ((0.4 - 0.55)^2 + (0.7 - 0.55)^2) \quad (5.7)$$

$$= 0.0225 \quad (5.8)$$

	Expected		Observed	
A_1A_1	0.3025	+	0.0225	= 0.3250
A_1A_2	0.4950	-	$2(0.0225)$	= 0.4500
A_2A_2	0.2025	+	0.0225	= 0.2250

Wright's F -statistics

One limitation of the way I've described things so far is that $\text{Var}(p)$ doesn't provide a convenient way to compare population structure from different samples. $\text{Var}(p)$ can be much larger if both alleles are about equally common in the whole sample than if one occurs at a mean frequency of 0.99 and the other at a frequency of 0.01. Moreover, if you stare at equations (5.4)–(5.6) for a while, you begin to realize that they look a lot like some equations we've already encountered. Namely, if we were to define F_{st} ⁷ as $\text{Var}(p)/\bar{p}\bar{q}$, then we could rewrite equations (5.4)–(5.6) as

$$\frac{1}{k} \sum p_i^2 = \bar{p}^2 + F_{st}\bar{p}\bar{q} \quad (5.9)$$

$$\frac{1}{k} \sum 2p_iq_i = 2\bar{p}\bar{q}(1 - F_{st}) \quad (5.10)$$

$$\frac{1}{k} \sum q_i^2 = \bar{q}^2 + F_{st}\bar{p}\bar{q} \quad (5.11)$$

And it's not even completely artificial to define F_{st} the way I did. After all, the effect of geographic structure is to cause matings to occur among genetically similar individuals. It's rather like inbreeding. Moreover, the extent to which this local mating matters depends on the extent to which populations differ from one another. $\bar{p}\bar{q}$ is the maximum allele frequency

⁶And this is what we predicted when we started.

⁷The reason for the subscript will become apparent later. It's also *very* important to notice that I'm defining F_{ST} here in terms of the population parameters p and $\text{Var}(p)$. Again, we'll return to the problem of how to *estimate* F_{ST} from data next time.

variance possible, given the observed mean frequency. So one way of thinking about F_{st} is that it measures the amount of allele frequency variance in a sample relative to the maximum possible.⁸

There may, of course, be inbreeding within populations, too. But it's easy to incorporate this into the framework, too.⁹ Let H_i be the actual heterozygosity in individuals within subpopulations, H_s be the expected heterozygosity within subpopulations assuming Hardy-Weinberg within populations, and H_t be the expected heterozygosity in the combined population assuming Hardy-Weinberg over the whole sample.¹⁰ Then thinking of f as a measure of departure from Hardy-Weinberg and assuming that all populations depart from Hardy-Weinberg to the same degree, i.e., that they all have the same f , we can define

$$F_{it} = 1 - \frac{H_i}{H_t}$$

Let's fiddle with that a bit.

$$\begin{aligned} 1 - F_{it} &= \frac{H_i}{H_t} \\ &= \left(\frac{H_i}{H_s}\right) \left(\frac{H_s}{H_t}\right) \\ &= (1 - F_{is})(1 - F_{st}) \quad , \end{aligned}$$

where F_{is} is the inbreeding coefficient within populations, i.e., f , and F_{st} has the same definition as before.¹¹ H_t is often referred to as the genetic diversity in a population. So another way of thinking about $F_{st} = (H_t - H_s)/H_t$ is that it's the proportion of the diversity in the sample that's due to allele frequency differences among populations.

⁸I say "one way", because there are several other ways to talk about F_{st} , too. But we won't talk about them until later.

⁹At least it's easy once you've been shown how.

¹⁰Please remember that we're assuming we know those frequencies exactly. In real applications, of course, we'll *estimate* those frequencies from data, so we'll have to account for sampling error when we actually try to estimate these things. If you're getting the impression that I think the distinction between allele frequencies as *parameters*, i.e., the real allele frequency in the population, and allele frequencies as *estimates*, i.e., the sample frequencies from which we hope to estimate the parameters, is really important, you're getting the right impression.

¹¹It takes a fair amount of algebra to show that this definition of F_{st} is equivalent to the one I showed you before, so you'll just have to take my word for it.

Chapter 6

Analyzing the genetic structure of populations

We've now seen the principles underlying Wright's F -statistics. I should point out that Gustave Malécot developed very similar ideas at about the same time as Wright, but since Wright's notation stuck,¹ population geneticists generally refer to statistics like those we've discussed as Wright's F -statistics.²

Neither Wright nor Malécot worried too much about the problem of estimating F -statistics from data. Both realized that any inferences about population structure are based on a sample and that the characteristics of the sample may differ from those of the population from which it was drawn, but neither developed any explicit way of dealing with those differences. Wright develops some very *ad hoc* approaches in his book [102], but they have been forgotten, which is good because they aren't very satisfactory and they shouldn't be used. There are now three reasonable approaches available:³

1. Nei's G -statistics,
2. Weir and Cockerham's θ -statistics, and
3. A Bayesian analog of θ .⁴

¹Probably because he published in English and Malécot published in French.

²The Hardy-Weinberg proportions should probably be referred to as the Hardy-Weinberg-Castle proportions too, since Castle pointed out the same principle. For some reason, though, his demonstration didn't have the impact that Hardy's and Weinberg's did. So we generally talk about the Hardy-Weinberg principle.

³And as we'll soon see, I'm not too crazy about one of these three. To my mind, there are really only two approaches that anyone should consider.

⁴These is, as you have probably already guessed, my personal favorite. We'll talk about it next time.

Population	Genotype			\hat{p}
	A_1A_1	A_1A_2	A_2A_2	
Yackeyackine Soak	29	0	0	1.0000
Gnarlbine Rock	14	3	3	0.7750
Boorabbin	15	2	3	0.8000
Bullabulling	9	0	0	1.0000
Mt. Caudan	9	0	0	1.0000
Victoria Rock	23	5	2	0.8500
Yellowdine	23	3	4	0.8167
Wargangering	29	3	1	0.9242
Wagga Rock	5	0	0	1.0000
“Iron Knob Major”	1	0	0	1.0000
Rainy Rocks	0	1	0	0.5000
“Rainy Rocks Major”	1	0	0	1.0000

Table 6.1: Genotype counts at the $GOT - 1$ locus in *Isotoma petraea* (from [48]).

An example from *Isotoma petraea*

To make the differences in implementation and calculation clear, I’m going to use data from 12 populations of *Isotoma petraea* in southwestern Australia surveyed for genotype at $GOT-1$ [48] as an example throughout these discussions (Table 6.1).

Let’s ignore the sampling problem for a moment and calculate the F -statistics as if we had observed the population allele frequencies without error. They’ll serve as our baseline for comparison.

$$\begin{aligned}
 \bar{p} &= 0.8888 \\
 \text{Var}(p) &= 0.02118 \\
 F_{st} &= 0.2143 \\
 \text{Individual heterozygosity} &= (0.0000 + 0.1500 + 0.1000 + 0.0000 + 0.0000 + 0.1667 + 0.1000 \\
 &\quad + 0.0909 + 0.0000 + 0.0000 + 1.0000 + 0.0000)/12 \\
 &= 0.1340 \\
 \text{Expected heterozygosity} &= 2(0.8888)(1 - 0.8888) \\
 &= 0.1976 \\
 F_{it} &= 1 - \frac{\text{Individual heterozygosity}}{\text{Expected heterozygosity}}
 \end{aligned}$$

$$\begin{aligned}
&= 1 - \frac{0.1340}{0.1976} \\
&= 0.3221 \\
1 - F_{it} &= (1 - F_{is})(1 - F_{st}) \\
F_{is} &= \frac{F_{it} - F_{st}}{1 - F_{st}} \\
&= \frac{0.3221 - 0.2143}{1 - 0.2143} \\
&= 0.1372
\end{aligned}$$

Summary

Correlation of gametes due to inbreeding within subpopulations (F_{is}):	0.1372
Correlation of gametes within subpopulations (F_{st}):	0.2143
Correlation of gametes in sample (F_{it}):	0.3221

Why do I refer to them as the “correlation of gametes ...”? There are two reasons:

1. That’s the way Wright always referred to and interpreted them.
2. We can define indicator variables $x_{ijk} = 1$ if the i th allele in the j th individual of population k is A_1 and $x_{ijk} = 0$ if that allele is not A_1 . This may seem like a strange thing to do, but the Weir and Cockerham approach to F -statistics described below uses just such an approach. If we do this, then the definitions for F_{is} , F_{st} , and F_{it} follow directly.⁵

Notice that F_{is} could be negative, i.e., there could be an *excess* of heterozygotes within populations ($F_{is} < 0$). Notice also that we’re implicitly assuming that the extent of departure from Hardy-Weinberg proportions is the same in all populations. Equivalently, we can regard F_{is} as the *average* departure from Hardy-Weinberg proportions across all populations.

Statistical expectation and biased estimates

The concept of statistical expectation is actually quite an easy one. It is an arithmetic average, just one calculated from probabilities instead of being calculated from samples. So,

⁵See [96] for details.

for example, if $P(k)$ is the probability that we find k A_1 alleles in our sample, the *expected number* of A_1 alleles in our sample is just

$$\begin{aligned} E(k) &= \sum kP(k) \\ &= np \quad , \end{aligned}$$

where n is the total number of alleles in our sample and p is the frequency of A_1 in our sample.⁶

Now consider the expected value of our sample estimate of the population allele frequency, $\hat{p} = k/n$, where k now refers to the number of A_1 alleles we actually found.

$$\begin{aligned} E(\hat{p}) &= E\left(\sum(k/n)\right) \\ &= \sum(k/n)P(k) \\ &= (1/n)\left(\sum kP(k)\right) \\ &= (1/n)(np) \\ &= p \quad . \end{aligned}$$

Because $E(\hat{p}) = p$, \hat{p} is said to be an *unbiased estimate* of p . When an estimate is unbiased it means that if we were to repeat the sampling experiment an infinite number of times and to take the average of the estimates, the average of those values would be equal to the (unknown) parameter value.

What about estimating the frequency of heterozygotes within a population? The obvious estimator is $\tilde{H} = 2\hat{p}(1 - \hat{p})$. Well,

$$\begin{aligned} E(\tilde{H}) &= E(2\hat{p}(1 - \hat{p})) \\ &= 2\left(E(\hat{p}) - E(\hat{p}^2)\right) \\ &= ((n - 1)/n)2p(1 - p) \quad . \end{aligned}$$

Because $E(\tilde{H}) \neq 2p(1 - p)$, \tilde{H} is a *biased estimate* of $2p(1 - p)$. If we set $\hat{H} = (n/(n - 1))\tilde{H}$, however, \hat{H} is an unbiased estimator of $2p(1 - p)$.⁷

⁶ $P(k) = \binom{N}{k}p^k(1 - p)^{N-k}$. The algebra in getting from the first line to the second is a little complicated, but feel free to ask me about it if you're interested.

⁷If you're wondering how I got from the second equation for \tilde{H} to the last one, ask me about it or read the gory details section that follows.

If you've ever wondered why you typically divide the sum of squared deviations about the mean by $n - 1$ instead of n when estimating the variance of a sample, this is why. Dividing by n gives you a (slightly) biased estimator.

The gory details⁸

Starting where we left off above:

$$\begin{aligned} E(\tilde{H}) &= 2 \left((E\hat{p}) - E(\hat{p}^2) \right) \\ &= 2 \left(p - E \left((k/n)^2 \right) \right) \quad , \end{aligned}$$

where k is the number of A_1 alleles in our sample and n is the sample size.

$$\begin{aligned} E \left((k/n)^2 \right) &= \sum (k/n)^2 P(k) \\ &= (1/n)^2 \sum k^2 P(k) \\ &= (1/n)^2 \left(\text{Var}(k) + \bar{k}^2 \right) \\ &= (1/n)^2 \left(np(1-p) + n^2 p^2 \right) \\ &= p(1-p)/n + p^2 \quad . \end{aligned}$$

Substituting this back into the equation above yields the following:

$$\begin{aligned} E(\tilde{H}) &= 2 \left(p - \left(p(1-p)/n + p^2 \right) \right) \\ &= 2 \left(p(1-p) - p(1-p)/n \right) \\ &= (1 - 1/n) 2p(1-p) \\ &= ((n-1)/n) 2p(1-p) \quad . \end{aligned}$$

Corrections for sampling error

There are two sources of allele frequency difference among subpopulations in our sample: (1) real differences in the allele frequencies among our sampled subpopulations and (2) differences that arise because allele frequencies in our samples differ from those in the subpopulations from which they were taken.⁹

⁸Skip this part unless you are *really, really* interested in how I got from the second equation to the third equation in the last paragraph. This is more likely to confuse you than help unless you know that the variance of a binomial sample is $np(1-p)$ and that $E(k^2) = \text{Var}(p) + p^2$.

⁹There's actually a third source of error that we'll get to in a moment. The populations we're sampling from are the product of an evolutionary process, and since the populations aren't of infinite size, drift has

Nei's G_{st}

Nei and Chesser [67] described one approach to accounting for sampling error. So far as I've been able to determine, there aren't any currently supported programs¹⁰ that calculate the bias-corrected versions of G_{st} .¹¹ I calculated the results in Table 6.2 by hand.

The calculations are tedious, which is why you'll want to find some way of automating the calculations if you want to do them.¹²

$$\begin{aligned}H_i &= 1 - \frac{1}{N} \sum_{k=1}^N \sum_{i=1}^m X_{kii} \\H_s &= \frac{\tilde{n}}{\tilde{n} - 1} \left[1 - \sum_{i=1}^m \bar{x}_i^2 - \frac{H_I}{2\tilde{n}} \right] \\H_t &= 1 - \sum_{i=1}^m \bar{x}_i^2 + \frac{H_s}{\tilde{n}} - \frac{H_I}{2\tilde{n}N}\end{aligned}$$

where we have N subpopulations, $\bar{x}_i^2 = \sum_{k=1}^N x_{ki}^2/N$, $\bar{x}_i = \sum_{k=1}^N x_{ki}/N$, \tilde{n} is the harmonic mean of the population sample sizes, i.e., $\tilde{n} = \frac{1}{\frac{1}{N} \sum_{k=1}^N \frac{1}{n_k}}$, X_{kii} is the frequency of genotype A_iA_i in population k , x_{ki} is the frequency of allele A_i in population k , and n_k is the sample size from population k . Recall that

$$\begin{aligned}F_{is} &= 1 - \frac{H_i}{H_s} \\F_{st} &= 1 - \frac{H_s}{H_t} \\F_{it} &= 1 - \frac{H_i}{H_t} .\end{aligned}$$

Weir and Cockerham's θ

Weir and Cockerham [97] describe the fundamental ideas behind this approach. Weir and Hill [98] bring things up to date. Holsinger and Weir [44] provide a less technical overview.¹³

played a role in determining allele frequencies in them. As a result, if we were to go back in time and re-run the evolutionary process, we'd end up with a different set of real allele frequency differences. We'll talk about this more in just a moment when we get to Weir and Cockerham's statistics.

¹⁰**Popgene** estimates G_{st} , but I don't think it's been updated since 2000. **FSTAT** also estimates gene diversities, but the most recent version is from 2002.

¹¹There's a reason for this that we'll get to in a moment. It's alluded to in the last footnote.

¹²It is also one big reason why most people use Weir and Cockerham's θ . There's readily available software that calculates it for you.

¹³We also talk a bit more about how F -statistics can be used.

We'll be using the implementations from **GDA** and **WinARL** in this course. The most important difference between θ and G_{st} and the reason why G_{st} has fallen into disuse is that G_{st} ignores an important source of sampling error that θ incorporates.

In many applications, especially in evolutionary biology, the subpopulations included in our sample are not an exhaustive sample of all populations. Moreover, even if we have sampled from every population there is now, we may not have sampled from every population there ever was. And even if we've sampled from every population there ever was, we know that there are random elements in any evolutionary process. Thus, if we could run the clock back and start it over again, the genetic composition of the populations we have might be rather different from that of the populations we sampled. In other words, our populations are, in many cases, best regarded as a random sample from a much larger set of populations that could have been sampled.

Even more gory details¹⁴

Let $x_{mn,i}$ be an indicator variable such that $x_{mn,i} = 1$ if allele m from individual n is of type i and is 0 otherwise. Clearly, the sample frequency $\hat{p}_i = \frac{1}{2N} \sum_{m=1}^2 \sum_{n=1}^N x_{mn,i}$, and $E(\hat{p}_i) = p_i$, $i = 1 \dots A$. Assuming that alleles are sampled independently from the population

$$\begin{aligned} E(x_{mn,i}^2) &= p_i \\ E(x_{mn,i}x_{m'n',i}) &= E(x_{mn,i}x_{m'n',i}) = p_i^2 + \sigma_{x_{mn,i}x_{m'n',i}} \\ &= p_i^2 + p_i(1 - p_i)\theta \end{aligned}$$

where $\sigma_{x_{mn,i}x_{m'n',i}}$ is the intraclass covariance for the indicator variables and

$$\theta = \frac{\sigma_{p_i}^2}{p_i(1 - p_i)} \tag{6.1}$$

is the scaled among population variance in allele frequency in the populations from which this population was sampled. Using (6.1) we find after some algebra

$$\sigma_{\hat{p}_i}^2 = p_i(1 - p_i)\theta + \frac{p_i(1 - p_i)(1 - \theta)}{2N} .$$

A natural estimate for θ emerges using the method of moments when an analysis of variance is applied to indicator variables derived from samples representing more than one population.

¹⁴This is even worse than the last time. I include it for completeness only. I really don't expect anyone (unless they happen to be a statistician) to be able to understand these details.

Method	F_{is}	F_{st}	F_{it}
Direct	0.1372	0.2143	0.3221
Nei	0.3092	0.2395	0.4746
Weir & Cockerham	0.5398	0.0387	0.5577

Table 6.2: Comparison of Wright’s F -statistics when ignoring sampling effects with Nei’s G_{ST} and Weir and Cockerham’s θ .

Notation	
F_{it}	F
F_{is}	f
F_{st}	θ

Table 6.3: Equivalent notations often encountered in descriptions of population genetic structure.

Applying G_{st} and θ

If we return to the data that motivated this discussion, these are the results we get from analyses of the $GOT - 1$ data from *Isotoma petraea* (Table 6.1). But first a note on how you’ll see statistics like this reported in the literature. It can get a little confusing, because of the different symbols that are used. Sometimes you’ll see F_{is} , F_{st} , and F_{it} . Sometimes you’ll see f , θ , and F . And it will seem as if they’re referring to similar things. That’s because they are. They’re really just different symbols for the same thing (see Table 6.3). Strictly speaking the symbols in Table 6.3 are the *parameters*, i.e., values in the population that we try to estimate. We should put hats over any values estimated from data to indicate that they are estimates of the parameters, not the parameters themselves. But we’re usually a bit sloppy, and everyone knows that we’re presenting estimates, so we usually leave off the hats.

An example from Wright

Hierarchical analysis of variation in the frequency of the Standard chromosome arrangement of *Drosophila pseudoobscura* in the western United States (data from [18], analysis from [103]). Wright uses his rather peculiar method of accounting for sampling error. I

haven't gone back to the original data and used a more modern method of analysis.¹⁵

66 populations (demes) studied. Demes are grouped into eight regions. The regions are grouped into four primary subdivisions.

Results

Correlation of gametes within individuals relative to regions (F_{IR}):	0.0444
Correlation of gametes within regions relative to subdivisions (F_{RS}):	0.0373
Correlation of gametes within subdivisions relative to total (F_{ST}):	0.1478
Correlation of gametes in sample (F_{IT}):	0.2160

$$1 - F_{IT} = (1 - F_{IR})(1 - F_{RS})(1 - F_{ST})$$

Interpretation

There is relatively little inbreeding within regions ($F_{IR} = 0.04$) and relatively little genetic differentiation among regions within subdivisions ($F_{RS} = 0.04$). There is, however, substantial genetic differentiation among the subdivisions ($F_{ST} = 0.15$).

Thus, an explanation for the chromosomal diversity that predicted great local differentiation and little or no differentiation at a large scale would be inconsistent with these observations.

¹⁵Sounds like it might be a good project, doesn't it? We'll see.

Chapter 7

Analyzing the genetic structure of populations: a Bayesian approach

Our review of Nei's G_{st} and Weir and Cockerham's θ illustrated two important principles:

1. It's essential to distinguish *parameters* from *estimates*. *Parameters* are the things we're really interested in, but since we always have to make inferences about the things we're really interested in from limited data, we have to rely on *estimates* of those parameters.
2. This means that we have to identify the possible sources of sampling error in our estimates and to find ways of accounting for them. In the particular case of Wright's F -statistics we saw that, there are two sources of sampling error: the error associated with sampling only some individuals from a larger universe of individuals within populations (*statistical sampling*) and the error associated with sampling only some populations from a larger universe of populations (*genetic sampling*).¹

It shouldn't come as any surprise that there is a Bayesian way to do what I've just described. As I hope to convince you, there are some real advantages associated with doing so.

The Bayesian model

I'm not going to provide all of the gory details on the Bayesian model. If you're interested you can find most of them in my lecture notes from the Summer Institute in Statistical

¹The terms "statistical sampling" and "genetic sampling" are due to Weir [96].

Genetics last summer.² In fact, I'm only going to describe two pieces of the model.³ First, a little notation:

$$\begin{aligned}
 n_{11,i} &= \# \text{ of } A_1A_1 \text{ genotypes} \\
 n_{12,i} &= \# \text{ of } A_1A_2 \text{ genotypes} \\
 n_{22,i} &= \# \text{ of } A_2A_2 \text{ genotypes} \\
 i &= \text{population index} \\
 I &= \text{number of populations}
 \end{aligned}$$

These are the data we have to work with. The corresponding genotype frequencies are

$$\begin{aligned}
 x_{11,i} &= p_i^2 + fp_i(1 - p_i) \\
 x_{12,i} &= 2p_i(1 - p_i)(1 - f) \\
 x_{22,i} &= (1 - p_i)^2 + fp_i(1 - p_i)
 \end{aligned}$$

So we can express the likelihood of our sample as a product of multinomial probabilities

$$P(\mathbf{n}|\mathbf{p}, f) \propto \prod_{i=1}^I x_{11,i}^{n_{11,i}} x_{12,i}^{n_{12,i}} x_{22,i}^{n_{22,i}} .$$

To complete the Bayesian model, all we need are some appropriate priors. Specifically, we so far haven't done anything to describe the variation in allele frequency among populations. Suppose that the distribution of allele frequencies among populations is well-approximated by a Beta distribution. A Beta distribution is convenient for many reasons, and it is quite flexible. Don't worry about what the formula for a Beta distribution looks like. All you need to know is that it has two parameters and that if these parameters are π and θ , we can set things up so that

$$\begin{aligned}
 E(p_{ik}) &= \pi \\
 \text{Var}(p_{ik}) &= \pi(1 - \pi)\theta
 \end{aligned}$$

Thus π corresponds to \bar{p} and θ corresponds to F_{st} .⁴ Figure 7.1 illustrates the shape of the Beta distribution for different choices of π and θ . To complete the Bayesian model we need

²Or you can read Holsinger and Wallace [43], which I've linked to from the course web site.

³The good news is that to do the Bayesian analyses in this case, you don't have to write any WinBUGS code. You just have to get your data into a format that `Hickory` recognizes.

⁴For any of you who happen to be familiar with the usual parameterization of a Beta distribution, this parameterization corresponds to setting $\nu = ((1 - \theta)/\theta)\pi$ and $\omega = ((1 - \theta)/\theta)(1 - \pi)$.

Parameter	Mean (s.d.)	Credible Interval	
		2.5%	97.5%
f	0.52 (0.10)	0.32	0.70
$\theta^{(II)}$	0.19 (0.12)	0.03	0.50
G_{st}^B	0.10 (0.07)	0.02	0.30

Table 7.1: Results from analyzing the *Isotoma petraea* data with **Hickory**.

only to specify priors on π , f , and θ . In the absence of any prior knowledge about the parameters, a uniform prior on $[0,1]^5$ is a natural choice.

The *Isotoma petraea* example

If we put the data we've analyzed before into **Hickory**, we get the results shown in Table 7.1 f is the within-population inbreeding coefficient, F_{is} , $\theta^{(II)}$ is the Bayesian analog of Weir and Cockerham's θ ,⁶ and G_{st}^B is the Bayesian analog of Nei's G_{st} .

Hickory allows you to select from three models when running the data with codominant markers:

1. The `full model` estimates both f and $\theta^{(II)}$.
2. The `f = 0 model` constrains $f = 0$ and estimates $\theta^{(II)}$.
3. The `theta = 0 model` constrains $\theta^{(II)} = 0$ and estimates f .

As a result, we can use DIC comparisons among the models to determine whether we have evidence for inbreeding within populations ($f = 0$ versus $f \neq 0$) or for genetic differentiation among populations ($\theta^{(II)} = 0$ versus $\theta^{(II)} \neq 0$). If we do that with these data we get the results shown in Table 7.2.

The $f = 0$ has a much larger DIC than the full model, a difference of more than 20 units. Thus, we have strong evidence for inbreeding in these populations of *Isotoma petraea*.⁷ The $\theta^{(II)} = 0$ model also has a DIC substantially larger than the DIC for the full model, a difference of more than 10 units. Thus, we also have good evidence for genetic differentiation among these populations.⁸

⁵`dunif(0,1)` in WinBUGS notation

⁶This isn't quite true. If you're interested, ask me about it. If you're really interested, take a look at [85].

⁷Much stronger than the evidence we had for inbreeding in the ABO blood group data, by the way.

⁸It's important to remember that this conclusion applies *only* to the locus that we analyzed. Strong differentiation at this locus need not imply that there is strong differentiation at other loci.

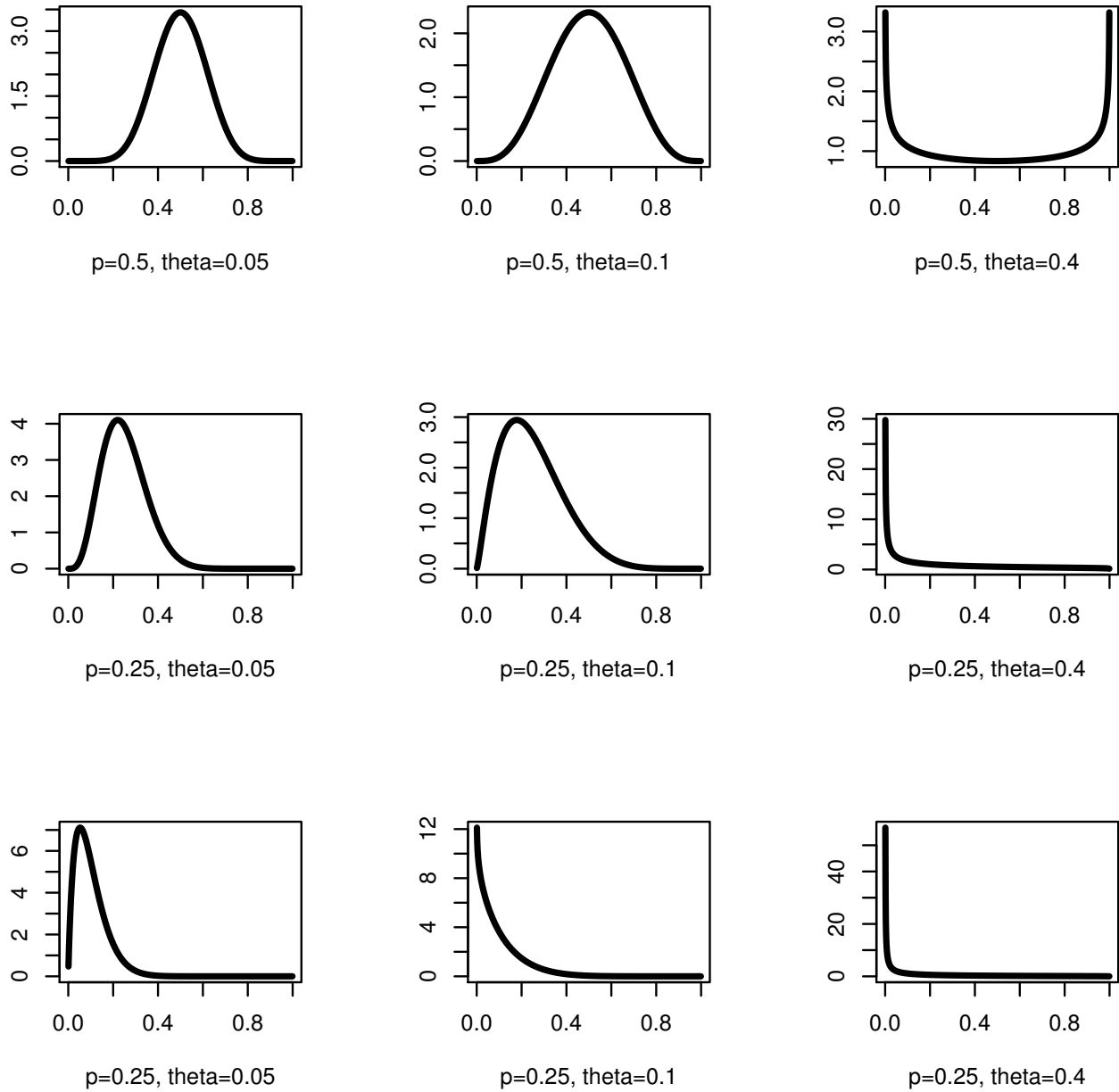


Figure 7.1: Shapes of the Beta distribution for different choices of π and θ . In the figure captions “p” corresponds to π , and “theta” corresponds to θ .

Model	DIC
Full	331.2
$f = 0$	355.4
$\theta^{(II)} = 0$	343.7

Table 7.2: DIC analysis of the *Isotoma petraea* data.

Method	F_{is}	F_{st}	
Direct	0.14	0.21	0.32
Nei	0.31	0.24	0.47
Weir & Cockerham	0.54	0.04	0.56
Bayesian	0.52 (0.32, 0.70)	0.19 (0.03, 0.50)	

Table 7.3: Comparison of F_{is} and F_{st} estimates calculated in different ways.

If we ask `Hickory` to keep log files of our analyses, we can also compare the estimates of f for the full and $\theta^{(II)} = 0$ models. The posterior means are 0.52 and 0.55, respectively, so they seem pretty similar. But we can do better than that. Since we have the full posterior distribution for f in both models, we can pick points at random from each, take their difference, and construct a 95% credible interval. Doing that we find that the 95% credible interval for $f_{full} - f_{\theta^{(II)}=0}$ is (-0.29, 0.24), meaning we have no evidence that the estimates are different. That may be a little surprising, since we strong evidence that $\theta^{(II)} \neq 0$ in these data, but it's also good news. It means that our estimate of within-population inbreeding is not much affected by the amount of differentiation among populations.⁹ What may be more surprising is that the estimates for $\theta^{(II)}$ with and without inbreeding are also very similar: 0.19 *versus* 0.22, respectively. Moreover, the 95% credible interval for $\theta_{full}^{(II)} - \theta_{f=0}^{(II)}$ is (-0.38, 0.33). Thus, even though we have strong evidence that $f > 0$, our estimate of $\theta^{(II)}$ is not strongly affected by what we think about f in these data.

One more thing we can do is to look at the posterior distribution of our parameter estimates and at the sample trace.¹⁰ (Figure 7.2). The result of an MCMC analysis, like the ones here or the ones in `WinBUGS`, is a large number of individual points. We can either fit a distribution to those points and display the results (the black lines in the figures), or we can use a non-parametric, kernel density estimate (the blue lines in the figures). The sample traces below show the values the chain took on at each point in the sampling process, and

⁹And when you say it that way, it sort of makes sense, doesn't it?

¹⁰You may have discovered that you can do this in `WinBUGS` too.

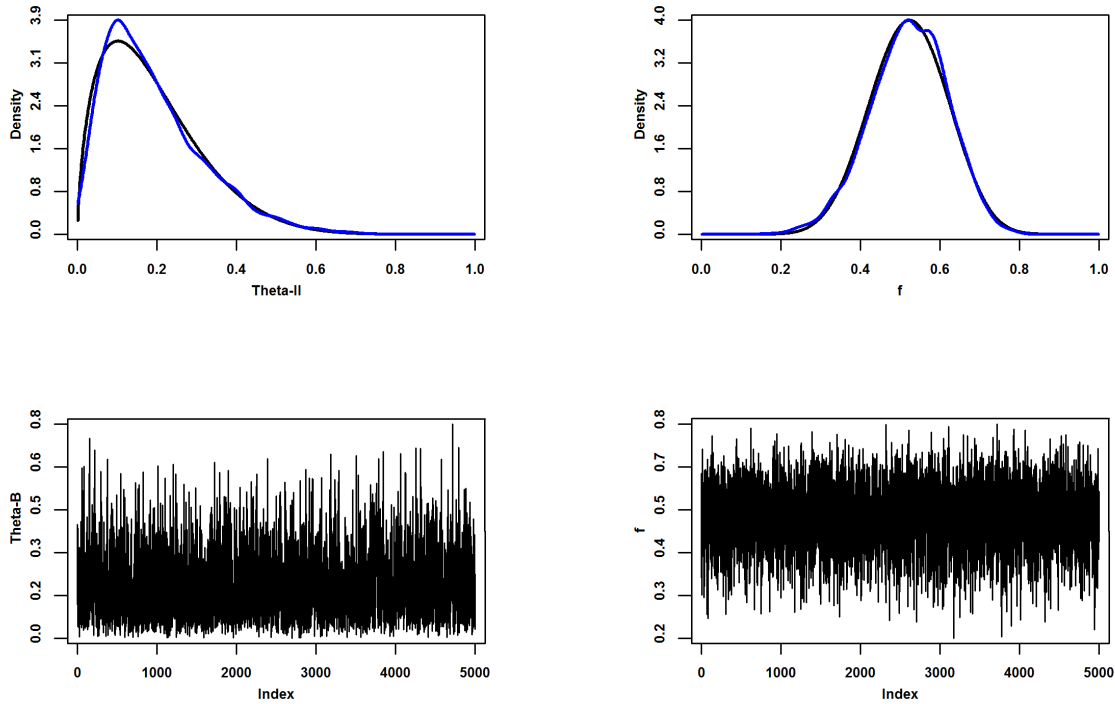


Figure 7.2: Posterior densities and traces for parameters of the full model applied to the *Isotoma petraea* data.

you can see that the values bounced around, which is good.

It's also useful to look back and think about the different ways we've used the data from *Isotoma petraea* (Table 7.3). Several things become apparent from looking at this table:

- The direct calculation is very misleading. A population that has only one individual sampled carries as much weight in determining F_{st} and F_{is} as populations with samples of 20-30 individuals.
- By failing to account for genetic sampling, Nei's statistics significantly underestimate F_{is} , while Weir & Cockerham's estimate is quite close to the Bayesian estimates.
- It's not illustrated here, but when a reasonable number of loci are sampled, say more than 8-10, the Weir & Cockerham estimates and the Bayesian estimates are quite

similar. But the Bayesian estimates allow for more convenient comparisons of different estimates, and the credible intervals don't depend either on asymptotic approximations or on bootstrapping across a limited collection of loci. The Bayesian approach can also be extended more easily to complex situations.

Chapter 8

Analyzing the genetic structure of populations: individual assignment

In the last 10-12 years a different approach to the analysis of genetic structure has emerged: analysis of individual assignment. Although the implementation details get a little hairy, the basic idea is fairly simple. Suppose we have genetic data on a series of individuals. Label the data we have for each individual x_i . Suppose that all individuals belong to one of K populations and let the genotype frequencies in population k be represented by γ_k . Then the likelihood that individual i comes from population k is just

$$P(i|k) = \frac{P(x_i|\gamma_k)}{\sum_k P(x_i|\gamma_k)} .$$

So if we can specify prior probabilities for γ_k , we can use Bayesian methods to estimate the posterior probability that individual i belongs to population k , and we can associate that assignment with some measure of its reliability.¹

Applying assignment to understand invasions

We'll use **Structure** to assess whether cultivated genotypes of *Berberis thunbergii* contribute to ongoing invasions in Connecticut and Massachusetts [63]. The first problem is to determine what K to use, because K doesn't necessarily have to equal the number of populations we sample from. Some populations may not be distinct from one another. There are a couple of ways to estimate K . The most straightforward is to run the analysis for a range of plausible values, repeat it 10-20 times for each value, calculate the mean "log probability of

¹You can find details in [77].

K	Mean L(K)
2	-2553.2
3	-2331.9
4	-2402.9
5	-2476.3

Table 8.1: Mean log probability of the data for $K = 2, 3, 4, 5$ in the *Berberis thunbergii* data (adapted from [63]).

the data” for each value of K , and pick the value of K that is the biggest, i.e., the least negative (Table 8.1). For the barberry data, $K = 3$ is the obvious choice.

Having determined that the data support $K = 3$, the results of the analysis are displayed in Figure 8.1. Each vertical bar corresponds to an individual in the sample, and the proportion of each bar that is of a particular color tells us the posterior probability that the individual belongs to the cluster with that color.

Figure 8.1 may not look terribly informative, but actually it is. Look at the labels beneath the figure. You’ll see that with the exception of individual 17 from Beaver Brook Park, all the of the individuals that are solid blue are members of the cultivated *Berberis thunbergii* var. *atropurpurea*. The solid red bar corresponds to *Berberis thunbergii* ‘Atropurpurea’, another modern cultivar. You’ll notice that individuals 1, 2, 18, and 19 from Beaver Brook Park and individual 1 from Bluff Point State Park fall into the same genotypic cluster as this cultivar. *Berberis* \times *ottawensis* is a hybrid cultivar whose parents are *Berberis thunbergii* and *Berberis vulgaris*, so it makes sense that individuals of this cultivar would be half blue and half red. The solid green bars are feral individuals from long-established populations. Notice that the cultivars are distinct from all but a few of the individuals in the long-established feral populations, suggesting that contemporary cultivars are doing relatively little to maintain the invasion in areas where it is already established.

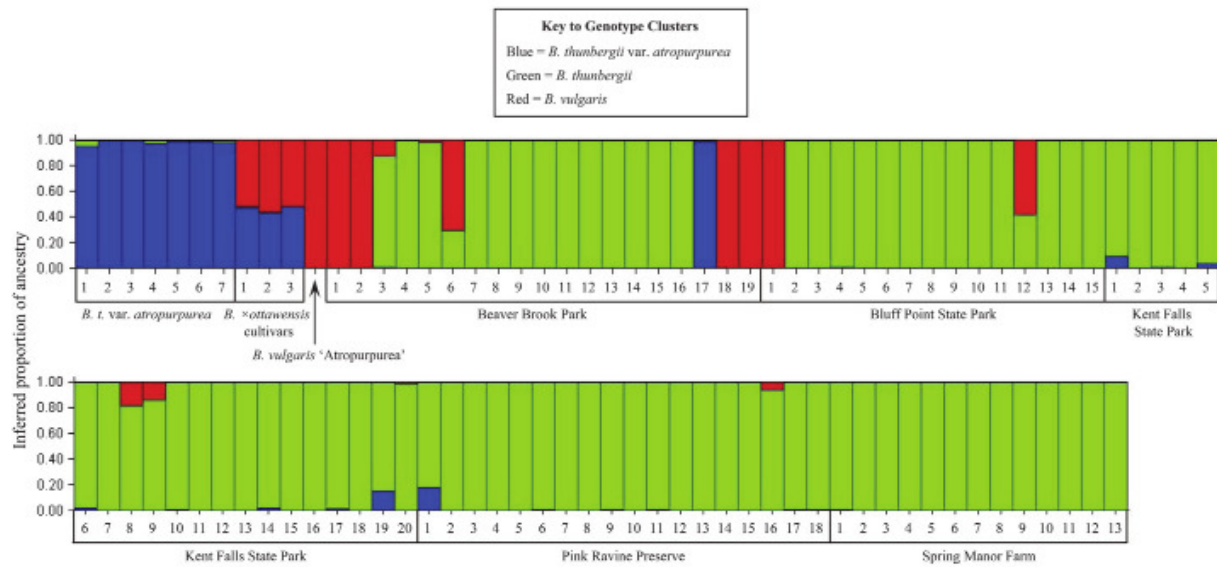


Figure 8.1: Analysis of AFLP data from *Berberis thunbergii* [63].

Chapter 9

Two-locus population genetics

So far in this course we've dealt only with variation at a single locus. There are obviously many traits that are governed by more than a single locus in whose evolution we might be interested. And for those who are concerned with the use of genetic data for forensic purposes, you'll know that forensic use of genetic data involves genotype information from multiple loci. I won't be discussing quantitative genetic variation for a few weeks, and I'm not going to say anything about how population genetics gets applied to forensic analyses, but I do want to introduce some basic principles of multilocus population genetics that are relevant to our discussions of the genetic structure of populations before moving on to the next topic. To keep things relatively simple *multilocus* population genetics will, for purposes of this lecture, mean *two-locus* population genetics.

Gametic disequilibrium

One of the most important properties of a two-locus system is that it is no longer sufficient to talk about allele frequencies alone, even in a population that satisfies all of the assumptions necessary for genotypes to be in Hardy-Weinberg proportions at each locus. To see why consider this. With two loci and two alleles there are four possible gametes:¹

Gamete	A_1B_1	A_1B_2	A_2B_1	A_2B_2
Frequency	x_{11}	x_{12}	x_{21}	x_{22}

If alleles are arranged randomly into gametes then,

$$x_{11} = p_1p_2$$

¹Think of drawing the Punnett square for a dihybrid cross, if you want.

$$\begin{aligned}
x_{12} &= p_1q_2 \\
x_{21} &= q_1p_2 \\
x_{22} &= q_1q_2 \quad ,
\end{aligned}$$

where $p_1 = \text{freq}(A_1)$ and $p_2 = \text{freq}(A_2)$. But alleles need not be arranged randomly into gametes. They may covary so that when a gamete contains A_1 it is more likely to contain B_1 than a randomly chosen gamete, or they may covary so that a gamete containing A_1 is less likely to contain B_1 than a randomly chosen gamete. This covariance could be the result of the two loci being in close physical association, but it doesn't have to be. Whenever the alleles covary within gametes

$$\begin{aligned}
x_{11} &= p_1p_2 + D \\
x_{12} &= p_1q_2 - D \\
x_{21} &= q_1p_2 - D \\
x_{22} &= q_1q_2 + D \quad ,
\end{aligned}$$

where $D = x_{11}x_{22} - x_{12}x_{21}$ is known as the *gametic disequilibrium*.² When $D \neq 0$ the alleles within gametes covary, and D measures *statistical* association between them. It does not (directly) measure the *physical* association. Similarly, $D = 0$ does not imply that the loci are unlinked, only that the alleles at the two loci are arranged into gametes independently of one another.

A little diversion

It probably isn't obvious why we can get away with only one D for all of the gamete frequencies. The short answer is:

There are four gametes. That means we need three parameters to describe the four frequencies. p_1 and p_2 are two. D is the third.

Another way is to do a little algebra to verify that the definition is self-consistent.

$$\begin{aligned}
D &= x_{11}x_{22} - x_{12}x_{21} \\
&= (p_1p_2 + D)(q_1q_2 + D) - (p_1q_2 - D)(q_1p_2 - D) \\
&= (p_1q_1p_2q_2 + D(p_1p_2 + q_1q_2) + D^2)
\end{aligned}$$

²You will sometimes see D referred to as the linkage disequilibrium, but that's misleading. Alleles at different loci may be non-randomly associated even when they are not linked.

$$\begin{aligned}
& - (p_1q_1p_2q_2 - D(p_1q_2 + q_1p_2) + D^2) \\
= & D(p_1p_2 + q_1q_2 + p_1q_2 + q_1p_2) \\
= & D(p_1(p_2 + q_2) + q_1(q_2 + p_2)) \\
= & D(p_1 + q_1) \\
= & D \quad .
\end{aligned}$$

Transmission genetics with two loci

I'm going to construct a reduced version of a mating table to see how gamete frequencies change from one generation to the next. There are ten different two-locus genotypes (if we distinguish coupling, A_1B_1/A_2B_2 , from repulsion, A_1B_2/A_2B_1 , heterozygotes as we must for these purposes). So a full mating table would have 100 rows. If we assume all the conditions necessary for genotypes to be in Hardy-Weinberg proportions apply, however, we can get away with just calculating the frequency with which any one genotype will produce a particular gamete.³

Genotype	Frequency	Gametes			
		A_1B_1	A_1B_2	A_2B_1	A_2B_2
A_1B_1/A_1B_1	x_{11}^2	1	0	0	0
A_1B_1/A_1B_2	$2x_{11}x_{12}$	$\frac{1}{2}$	$\frac{1}{2}$	0	0
A_1B_1/A_2B_1	$2x_{11}x_{21}$	$\frac{1}{2}$	0	$\frac{1}{2}$	0
A_1B_1/A_2B_2	$2x_{11}x_{22}$	$\frac{1-r}{2}$	$\frac{r}{2}$	$\frac{r}{2}$	$\frac{1-r}{2}$
A_1B_2/A_1B_2	x_{12}^2	0	1	0	0
A_1B_2/A_2B_1	$2x_{12}x_{21}$	$\frac{r}{2}$	$\frac{1-r}{2}$	$\frac{1-r}{2}$	$\frac{r}{2}$
A_1B_2/A_2B_2	$2x_{12}x_{22}$	0	$\frac{1}{2}$	0	$\frac{1}{2}$
A_2B_1/A_2B_1	x_{21}^2	0	0	1	0
A_2B_1/A_2B_2	$2x_{21}x_{22}$	0	0	$\frac{1}{2}$	$\frac{1}{2}$
A_2B_2/A_2B_2	x_{22}^2	0	0	0	1

Where do $\frac{1-r}{2}$ and $\frac{r}{2}$ come from?

Consider the coupling double heterozygote, A_1B_1/A_2B_2 . When recombination doesn't happen, A_1B_1 and A_2B_2 occur in equal frequency ($1/2$), and A_1B_2 and A_2B_1 don't occur at all. When recombination happens, the four possible gametes occur in equal frequency ($1/4$). So

³We're assuming random union of *gametes* rather than random mating of *genotypes*.

the recombination frequency,⁴ r , is half the crossover frequency,⁵ c , i.e., $r = c/2$. Now the results of crossing over can be expressed in this table:

Frequency	A_1B_1	A_1B_2	A_2B_1	A_2B_2
$1 - c$	$\frac{1}{2}$	0	0	$\frac{1}{2}$
c	$\frac{1}{4}$	$\frac{1}{4}$	$\frac{1}{4}$	$\frac{1}{4}$
Total	$\frac{2-c}{4}$	$\frac{c}{4}$	$\frac{c}{4}$	$\frac{2-c}{4}$
	$\frac{1-r}{2}$	$\frac{r}{2}$	$\frac{r}{2}$	$\frac{1-r}{2}$

Changes in gamete frequency

We can use this table as we did earlier to calculate the frequency of each gamete in the next generation. Specifically,

$$\begin{aligned}
 x'_{11} &= x_{11}^2 + x_{11}x_{12} + x_{11}x_{21} + (1 - r)x_{11}x_{22} + rx_{12}x_{21} \\
 &= x_{11}(x_{11} + x_{12} + x_{21} + x_{22}) - r(x_{11}x_{22} - x_{12}x_{21}) \\
 &= x_{11} - rD \\
 x'_{12} &= x_{12} + rD \\
 x'_{21} &= x_{21} + rD \\
 x'_{22} &= x_{22} - rD \quad .
 \end{aligned}$$

No changes in allele frequency

We can also calculate the frequencies of A_1 and B_1 after this whole process:

$$\begin{aligned}
 p'_1 &= x'_{11} + x'_{12} \\
 &= x_{11} - rD + x_{12} + rD \\
 &= x_{11} + x_{12} \\
 &= p_1 \\
 p'_2 &= p_2 \quad .
 \end{aligned}$$

Since each locus is subject to all of the conditions necessary for Hardy-Weinberg to apply at a single locus, allele frequencies don't change at either locus. Furthermore, genotype frequencies at each locus will be in Hardy-Weinberg proportions. But the two-locus gamete frequencies change from one generation to the next.

⁴The frequency of recombinant gametes in double heterozygotes.

⁵The frequency of cytological crossover during meiosis.

Population	Gamete frequencies				Allele frequencies		D
	A_1B_1	A_1B_2	A_2B_1	A_2B_2	p_{i1}	p_{i2}	
1	0.24	0.36	0.16	0.24	0.60	0.40	0.00
2	0.14	0.56	0.06	0.24	0.70	0.20	0.00
Combined	0.19	0.46	0.11	0.24	0.65	0.30	-0.005

Table 9.1: Gametic disequilibrium in a combined population sample.

Changes in D

You can probably figure out that D will eventually become zero, and you can probably even guess that how quickly it becomes zero depends on how frequent recombination is. But I'd be astonished if you could guess exactly how rapidly D decays as a function of r . It takes a little more algebra, but we can say precisely how rapid the decay will be.

$$\begin{aligned}
D' &= x'_{11}x'_{22} - x'_{12}x'_{21} \\
&= (x_{11} - rD)(x_{22} - rD) - (x_{12} + rD)(x_{21} + rD) \\
&= x_{11}x_{22} - rD(x_{11} + x_{12}) + r^2D^2 - (x_{12}x_{21} + rD(x_{12} + x_{21}) + r^2D^2) \\
&= x_{11}x_{22} - x_{12}x_{21} - rD(x_{11} + x_{12} + x_{21} + x_{22}) \\
&= D - rD \\
&= D(1 - r)
\end{aligned}$$

Notice that even if loci are unlinked, meaning that $r = 1/2$, D does not reach 0 immediately. That state is reached only asymptotically. The two-locus analogue of Hardy-Weinberg is that gamete frequencies will *eventually* be equal to the product of their constituent allele frequencies.

Population structure with two loci

You can probably guess where this is going. With one locus I showed you that there's a deficiency of heterozygotes in a combined sample even if there's random mating within all populations of which the sample is composed. The two-locus analog is that you can have gametic disequilibrium in your combined sample even if the gametic disequilibrium is zero in all of your constituent populations. Table 9.1 provides a simple numerical example involving just two populations in which the combined sample has equal proportions from each population.

The gory details

You knew that I wouldn't be satisfied with a numerical example, didn't you? You knew there had to be some algebra coming, right? Well, here it is. Let

$$\begin{aligned} D_i &= x_{11,i} - p_{1i}p_{2i} \\ D_t &= \bar{x}_{11} - \bar{p}_1\bar{p}_2 \quad , \end{aligned}$$

where $\bar{x}_{11} = \frac{1}{K} \sum_{k=1}^K x_{11,k}$, $\bar{p}_1 = \frac{1}{K} \sum_{k=1}^K p_{1k}$, and $\bar{p}_2 = \frac{1}{K} \sum_{k=1}^K p_{2k}$. Given these definitions, we can now calculate D_t .

$$\begin{aligned} D_t &= \bar{x}_{11} - \bar{p}_1\bar{p}_2 \\ &= \frac{1}{K} \sum_{k=1}^K x_{11,k} - \bar{p}_1\bar{p}_2 \\ &= \frac{1}{K} \sum_{k=1}^K (p_{1k}p_{2k} + D_k) - \bar{p}_1\bar{p}_2 \\ &= \frac{1}{K} \sum_{k=1}^K (p_{1k}p_{2k} - \bar{p}_1\bar{p}_2) + \bar{D} \\ &= \text{Cov}(p_1, p_2) + \bar{D} \quad , \end{aligned}$$

where $\text{Cov}(p_1, p_2)$ is the covariance in allele frequencies across populations and \bar{D} is the mean within-population gametic disequilibrium. Suppose $D_i = 0$ for all subpopulations. Then $\bar{D} = 0$, too (obviously). But that means that

$$D_t = \text{Cov}(p_1, p_2) \quad .$$

So if allele frequencies covary across populations, i.e., $\text{Cov}(p_1, p_2) \neq 0$, then there will be non-random association of alleles into gametes in the sample, i.e., $D_t \neq 0$, even if there is random association alleles into gametes within each population.⁶

Returning to the example in Table 9.1

$$\begin{aligned} \text{Cov}(p_1, p_2) &= 0.5(0.6 - 0.65)(0.4 - 0.3) + 0.5(0.7 - 0.65)(0.2 - 0.3) \\ &= -0.005 \\ \bar{x}_{11} &= (0.65)(0.30) - 0.005 \\ &= 0.19 \end{aligned}$$

⁶Well, duh! Covariation of allele frequencies across populations means that alleles are non-randomly associated across populations. What other result could you possibly expect?

$$\begin{aligned}\bar{x}_{12} &= (0.65)(0.7) + 0.005 \\ &= 0.46 \\ \bar{x}_{21} &= (0.35)(0.30) + 0.005 \\ &= 0.11 \\ \bar{x}_{22} &= (0.35)(0.70) - 0.005 \\ &= 0.24 \quad .\end{aligned}$$

Part II

The genetics of natural selection

Chapter 10

The Genetics of Natural Selection

So far in this course, we've focused on describing the pattern of variation within and among populations. We've talked about inbreeding, which causes *genotype* frequencies to change, although it leaves allele frequencies the same, and we've talked about how to describe variation among populations. But we haven't yet discussed any evolutionary processes that could lead to a change in allele frequencies within populations.¹

Let's return for a moment to the list of assumptions we developed when we derived the Hardy-Weinberg principle and see what we've done so far.

Assumption #1 Genotype frequencies are the same in males and females, e.g., x_{11} is the frequency of the A_1A_1 genotype in both males and females.

Assumption #2 Genotypes mate at random *with respect to their genotype at this particular locus*.

Assumption #3 Meiosis is fair. More specifically, we assume that there is no segregation distortion, no gamete competition, no differences in the developmental ability of eggs, or the fertilization ability of sperm.

Assumption #4 There is no input of new genetic material, i.e., gametes are produced without mutation, and all offspring are produced from the union of gametes within this population.

Assumption #5 The population is of infinite size so that the actual frequency of matings is equal to their expected frequency and the actual frequency of offspring from each mating is equal to the Mendelian expectations.

¹We mentioned migration and drift in passing, and I'm sure you all understand the rudiments of them, but we haven't yet discussed them in detail.

Assumption #6 All matings produce the same number of offspring, on average.

Assumption #7 Generations do not overlap.

Assumption #8 There are no differences among genotypes in the probability of survival.

The only assumption we've violated so far is Assumption #2, the random-mating assumption. We're going to spend the next several lectures talking about what happens when you violate Assumptions #3, #6, and #8. When any one of those assumptions is violated we have some form of natural selection going on.²

Components of selection

Depending on which of those three assumptions is violated and how it's violated we recognize that selection may happen in different ways and at different life-cycle stages.³

Assumption #3: *Meiosis is fair.* There are at least two ways in which this assumption may be violated.

- *Segregation distortion:* The two alleles are not equally frequent in gametes produced by heterozygotes. The *t*-allele in house mice, for example, is found in 95% of fertile sperm produced by heterozygous males.
- *Gamete competition:* Gametes may be produced in equal frequency in heterozygotes, but there may be competition among them to produce fertilized zygotes, e.g., sperm competition in animals, pollen competition in seed plants.

Assumption #6: *All matings produce the same number of progeny.*

- *Fertility selection:* The number of offspring produced may depend on maternal genotype (*fecundity selection*), paternal genotype (*virility selection*), or on both.

Assumption #8: *Survival does not depend on genotype.*

²As I alluded to when we first started talking about inbreeding, we can also have natural selection as a result of certain types of violations of assumption #2, e.g., sexual selection or disassortative mating. See below.

³To keep things *relatively* simple we're not even going to discuss differences in fitness that may be associated with different ages. We'll assume a really simple life-cycle in which there are non-overlapping generations. So we don't need to distinguish juveniles from adults.

- *Viability selection*: The probability of survival from zygote to adult may depend on genotype, and it may differ between sexes.

At this point you're probably thinking that I've covered all the possibilities. But by now you should also know me well enough to guess from the way I wrote that last sentence that if that's what you were thinking, you'd be wrong. There's one more way in which selection can happen that corresponds to violating

Assumption #2: *Individuals mate at random.*

- *Sexual selection*: Some individuals may be more successful at finding mates than others. Since females are typically the limiting sex (Bateman's principle), the differences typically arise either as a result of *male-male competition* or *female choice*.
- *Disassortative mating*: When individuals preferentially choose mates different from themselves, rare genotypes are favored relative to common genotypes. This leads to a form a frequency-dependent selection.

The genetics of viability selection

That's a pretty exhaustive (and exhausting) list of the ways in which selection can happen. Although we're going to focus our study of natural selection just on viability selection, it's important to remember that any or all of the other forms of selection may be operating simultaneously on the genes or the traits that we're studying, and the direction of selection due to these other components may be the same or different. We're going to focus on viability selection for two reasons:

1. The most basic properties of natural selection acting on other components of the life history are similar to those of viability selection. A good understanding of viability selection provides a solid foundation for understanding other types of selection.⁴
2. The algebra associated with understanding viability selection is a *lot* simpler than the algebra associated with understanding the other types of selection, and the dynamics are simpler and easier to understand.⁵

⁴There are some important differences, however, and I hope we have time to discuss a couple of them.

⁵Once you've seen what you're in for you may think I've lied about this. But if you really think I have, just ask me to illustrate some of the algebra necessary for understanding viability selection when males and females differ in fitness. That's about as simple an extension as you can imagine, and things start to get pretty complicated even then.

The basic framework

To understand the basics, we'll start with a numerical example using some data on *Drosophila pseudoobscura* that Theodosius Dobzhansky collected more than 50 years ago. You may remember that this species has chromosome inversion polymorphisms. Although these inversions involve many genes, they are inherited as if they were single Mendelian loci, so we can treat the karyotypes as single-locus genotypes and study their evolutionary dynamics. We'll be considering two inversion types the Standard inversion type, *ST*, and the Chiricahua inversion type, *CH*. We'll use the following notation throughout our discussion:

Symbol	Definition
N	number of individuals in the population
x_{11}	frequency of <i>ST/ST</i> genotype
x_{12}	frequency of <i>ST/CH</i> genotype
x_{22}	frequency of <i>CH/CH</i> genotype
w_{11}	fitness of <i>ST/ST</i> genotype, probability of surviving from egg to adult
w_{12}	fitness of <i>ST/CH</i> genotype
w_{22}	fitness of <i>CH/CH</i> genotype

The data look like this:⁶

Genotype	<i>ST/ST</i>	<i>ST/CH</i>	<i>CH/CH</i>
Number in eggs	41	82	27
	Nx_{11}	Nx_{12}	Nx_{22}
viability	0.6	0.9	0.45
	w_{11}	w_{12}	w_{22}
Number in adults	25	74	12
	$w_{11}x_{11}N$	$w_{12}x_{12}N$	$w_{22}x_{22}N$

Genotype and allele frequencies

It should be trivial for you by this time to calculate the genotype frequencies in eggs and adults. We'll be using the convention that genotype frequencies in eggs (or newly-formed zygotes) are the genotype frequencies *before selection* and that genotype frequencies in adults are the genotype frequencies *after selection*.

⁶Don't worry for the moment about how the viabilities were estimated.

$$\begin{aligned}
\text{freq}(ST/ST) \text{ before selection} &= \frac{41}{41 + 82 + 27} \\
&= 0.27 \\
\text{freq}(ST/ST) \text{ before selection} &= \frac{Nx_{11}}{Nx_{11} + Nx_{12} + Nx_{22}} \\
&= x_{11} \\
\text{freq}(ST/ST) \text{ after selection} &= \frac{25}{25 + 74 + 12} \\
&= 0.23 \\
\text{freq}(ST/ST) \text{ after selection} &= \frac{w_{11}x_{11}N}{w_{11}x_{11}N + w_{12}x_{12}N + w_{22}x_{22}N} \\
&= \frac{w_{11}x_{11}}{w_{11}x_{11} + w_{12}x_{12} + w_{22}x_{22}} \\
&= \frac{w_{11}x_{11}}{\bar{w}} \\
\bar{w} &= w_{11}x_{11} + w_{12}x_{12} + w_{22}x_{22}
\end{aligned}$$

where \bar{w} is the mean fitness, i.e., the average probability of survival in the population.

It is also trivial to calculate the allele frequencies before and after selection:

$$\begin{aligned}
\text{freq}(ST) \text{ before selection} &= \frac{2(41) + 82}{2(41 + 82 + 27)} \\
&= 0.55 \\
\text{freq}(ST) \text{ before selection} &= \frac{2(Nx_{11}) + Nx_{12}}{2(Nx_{11} + Nx_{12} + Nx_{22})} \\
&= x_{11} + x_{12}/2 \\
\text{freq}(ST) \text{ after selection} &= \frac{2(25) + 74}{2(25 + 74 + 12)} \\
&= 0.56 \\
\text{freq}(ST) \text{ after selection} &= \frac{2w_{11}x_{11}N + w_{12}x_{12}N}{2(w_{11}x_{11}N + w_{12}x_{12}N + w_{22}x_{22}N)} \\
&= \frac{2w_{11}x_{11} + w_{12}x_{12}}{2(w_{11}x_{11} + w_{12}x_{12} + w_{22}x_{22})} \\
p' &= \frac{w_{11}x_{11} + w_{12}x_{12}/2}{w_{11}x_{11} + w_{12}x_{12} + w_{22}x_{22}}
\end{aligned}$$

$$\begin{aligned}
x_{11} &= p^2, & x_{12} &= 2pq, & x_{22} &= q^2 \\
p' &= \frac{w_{11}p^2 + w_{12}pq}{w_{11}p^2 + w_{12}2pq + w_{22}q^2} \\
\bar{w} &= w_{11}x_{11} + w_{12}x_{12} + w_{22}x_{22} \\
&= p^2w_{11} + 2pqw_{12} + q^2w_{22}
\end{aligned}$$

If you're still awake, you're probably wondering⁷ why I was able to substitute p^2 , $2pq$, and q^2 for x_{11} , x_{12} , and x_{22} . Remember what I said earlier about what we're doing here. The *only* Hardy-Weinberg assumption we're violating is the one saying that all genotypes are equally likely to survive. Remember also that a single generation in which all of the conditions for Hardy-Weinberg is enough to establish the Hardy-Weinberg proportions. Putting those two observations together, it's not too hard to see that genotypes will be in Hardy-Weinberg proportions in newly formed zygotes. Viability selection will change that later in the life-cycle, but we restart every generation with genotypes in the familiar Hardy-Weinberg proportions, p^2 , $2pq$, and q^2 .

Selection acts on relative viability

Let's stare at the selection equation for awhile and see what it means.

$$p' = \frac{w_{11}p^2 + w_{12}pq}{\bar{w}} \quad . \quad (10.1)$$

Suppose, for example, that we were to divide the numerator and denominator of (10.1) by w_{11} .⁸ We'd then have

$$p' = \frac{p^2 + (w_{12}/w_{11})pq}{(\bar{w}/w_{11})} \quad . \quad (10.2)$$

Why did I bother to do that? Well, notice that we start with the same allele frequency, p , in the parental generation in both equations and that we end up with the same allele frequency in the offspring generation, p' , in both equations, but the fitnesses are different:

Equation	Fitnesses		
	A_1A_1	A_1A_2	A_2A_2
10.1	w_{11}	w_{12}	w_{22}
10.2	1	w_{12}/w_{11}	w_{22}/w_{11}

⁷Okay, "probably" is an overstatement. "May be" would have been a better guess.

⁸I'm dividing by 1, in case you hadn't noticed.

I could have, of course, divided the numerator and denominator by w_{12} or w_{22} instead and ended up with yet other sets of fitnesses that produce exactly the same change in allele frequency. This illustrates the following general principle:

The consequences of natural selection (in an infinite population) depend only on the *relative* magnitude of fitnesses, not on their *absolute* magnitude.

That means, for example, that in order to predict the outcome of viability selection, we don't have to know the probability that each genotype will survive, their *absolute viabilities*. We only need to know the probability that each genotype will survive relative to the probability that other genotypes will survive, their *relative viabilities*. As we'll see later, it's sometimes easier to estimate the relative viabilities than to estimate absolute viabilities.⁹

Marginal fitnesses

In case you haven't already noticed, there's almost always more than one way to write an equation.¹⁰ They're all mathematically equivalent, but they emphasize different things. In this case, it can be instructive to look at the difference in allele frequencies from one generation to the next, Δp :

$$\begin{aligned} \Delta p &= p' - p \\ &= \frac{w_{11}p^2 + w_{12}pq}{\bar{w}} - p \\ &= \frac{w_{11}p^2 + w_{12}pq - \bar{w}p}{\bar{w}} \\ &= \frac{p(w_{11}p + w_{12}q - \bar{w})}{\bar{w}} \\ &= \frac{p(w_1 - \bar{w})}{\bar{w}}, \end{aligned}$$

where w_1 is the *marginal fitness* of allele A_1 . To explain why it's called a *marginal* fitness, I'd have to teach you some probability theory that you probably don't want to learn.¹¹

⁹We'll also see when we get to studying the interaction between natural selection and drift that this statement is no longer true. To understand how drift and selection interact we have to know something about *absolute* viabilities.

¹⁰And you won't have noticed this and may not believe me when I tell you, but I'm *not* showing you every possible way to write these equations.

¹¹But remember this definition of marginal viability anyway. You'll see it return in a few weeks when we talk about the additive effect of an allele and about Fisher's Fundamental Theorem of Natural Selection.

Pattern	Description	Figure
Directional	$w_{11} > w_{12} > w_{22}$ or $w_{11} < w_{12} < w_{22}$	Figure 10.1
Disruptive	$w_{11} > w_{12}, w_{22} > w_{12}$	Figure 10.2
Stabilizing	$w_{11} < w_{12}, w_{22} < w_{12}$	Figure 10.3

Table 10.1: Patterns of viability selection at one locus with two alleles.

Fortunately, all you really need to know is that it corresponds to the probability that a randomly chosen A_1 allele in a newly formed zygote will survive into a reproductive adult.

Why do we care? Because it provides some (obvious) intuition on how allele frequencies will change from one generation to the next. If $w_1 > \bar{w}$, i.e., if the chances of a zygote carrying an A_1 allele of surviving to make an adult are greater than the chances of a randomly chosen zygote, then A_1 will increase in frequency. If $w_1 < \bar{w}$, A_1 will decrease in frequency. Only if $p = 0$, $p = 1$, or $w_1 = \bar{w}$ will the allele frequency not change from one generation to the next.

Patterns of natural selection

Well, all that algebra was lots of fun,¹² but what good did it do us? Not an enormous amount, except that it shows us (not surprisingly), that allele frequencies are likely to change as a result of viability selection, and it gives us a nice little formula we could plug into a computer to figure out exactly how. One of the reasons that it's useful¹³ to go through all of that algebra is that it's possible to make predictions about the consequences of natural selection simply by knowing the pattern of viability differences. What do I mean by pattern? Funny you should ask (Table 10.1).

Before exploring the consequences of these different patterns of natural selection, I need to introduce you to a very important result: Fisher's Fundamental Theorem of Natural Selection. We'll go through the details later when we get to quantitative genetics. For now all you need to know is that viability selection causes the mean fitness of the progeny generation to be greater than or equal to the mean fitness of the parental generation, with equality only at equilibrium, i.e.,

$$\bar{w}' \geq \bar{w} \quad .$$

¹²I'm kidding, in case you couldn't tell.

¹³If not exactly fun.

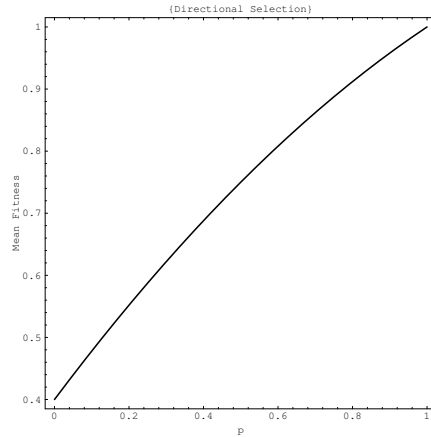


Figure 10.1: With *directional selection* (in this case $w_{11} > w_{12} > w_{22}$) viability selection leads to an ever increasing frequency of the favored allele. Ultimately, the population will be monomorphic for the homozygous genotype with the highest fitness.

How does this help us? Well, the best way to understand that is to illustrate how we can use Fisher's theorem to predict the outcome of natural selection when we know only the pattern of viability differences. Let's take each pattern in turn.

Directional selection

To use the Fundamental Theorem we plot \bar{w} as a function of p (Figure 10.1). The Fundamental Theorem now tells us that allele frequencies have to change from one generation to the next in such a way that $\bar{w}' > \bar{w}$, which can only happen if $p' > p$. So viability selection will cause the frequency of the A_1 allele to increase. Ultimately, the population will be monomorphic for the homozygous genotype with the highest fitness.¹⁴

Disruptive selection

If we plot \bar{w} as a function of p when $w_{11} > w_{12}$ and $w_{22} > w_{12}$, we see a very different pattern (Figure 10.2). Since the Fundamental Theorem tells us that $\bar{w}' \geq \bar{w}$, we know that

¹⁴A population is *monomorphic* at a particular locus when only one allele is present. If a population is monomorphic for allele A_1 , I might also say that allele A_1 is fixed in the population or that the population is fixed for allele A_1 .

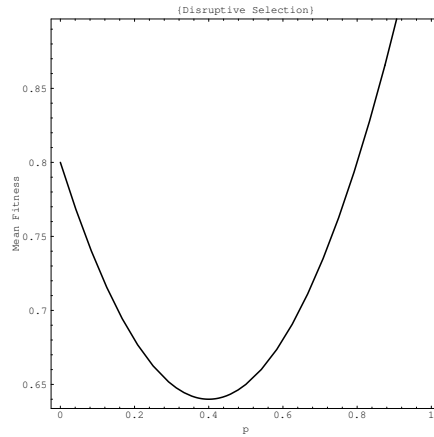


Figure 10.2: With *disruptive selection* ($w_{11} > w_{12} < w_{22}$) viability selection may lead either to an increasing frequency of the A allele or to a decreasing frequency. Ultimately, the population will be monomorphic for one of the homozygous genotypes. Which homozygous genotype comes to predominate, however, depends on the initial allele frequencies in the population.

if the population starts with an allele on one side of the bowl A_1 , will be lost. If it starts on the other side of the bowl, A_2 will be lost.¹⁵

Let's explore this example a little further. To do so, I'm going to set $w_{11} = 1 + s_1$, $w_{12} = 1$, and $w_{22} = 1 + s_2$.¹⁶ When fitnesses are written this way s_1 and s_2 are referred to as *selection coefficients*. Notice also with these definitions that the fitnesses of the homozygotes is greater than 1.¹⁷ Using these definitions and plugging them into (10.1),

$$\begin{aligned}
 p' &= \frac{p^2(1 + s_1) + pq}{p^2(1 + s_1) + 2pq + q^2(1 + s_2)} \\
 &= \frac{p(1 + s_1p)}{1 + p^2s_1 + q^2s_2} \quad . \quad (10.3)
 \end{aligned}$$

We can use equation (10.3) to find the equilibria of this system, i.e., the values of p such

¹⁵Strictly speaking, we need to know more than $\bar{w}' \geq \bar{w}$, but we do know the other things we need to know in this case. Trust me. Have I ever lied to you? (Don't answer that.)

¹⁶Why can I get away with this? Hint: Think about relative fitnesses.

¹⁷Which is why I gave you the relative fitness hint in the last footnote.

that $p' = p$.

$$\begin{aligned}
 p &= \frac{p(1 + s_1p)}{1 + p^2s_1 + q^2s_2} \\
 p(1 + p^2s_1 + q^2s_2) &= p(1 + s_1p) \\
 p((1 + p^2s_1 + q^2s_2) - (1 + s_1p)) &= 0 \\
 p(ps_1(p - 1) + q^2s_2) &= 0 \\
 p(-pq s_1 + q^2s_2) &= 0 \\
 pq(-ps_1 + qs_2) &= 0 \quad .
 \end{aligned}$$

So $p' = p$ if $\hat{p} = 0$, $\hat{q} = 0$, or $\hat{p}s_1 = \hat{q}s_2$.¹⁸ We can simplify that last one a little further, too.

$$\begin{aligned}
 \hat{p}s_1 &= \hat{q}s_2 \\
 \hat{p}s_1 &= (1 - \hat{p})s_2 \\
 \hat{p}(s_1 + s_2) &= s_2 \\
 \hat{p} &= \frac{s_2}{s_1 + s_2} \quad .
 \end{aligned}$$

Fisher's Fundamental Theorem tells us which of these equilibria matter. I've already mentioned that depending on which side of the bowl you start, you'll either lose the A_1 allele or the A_2 allele. But suppose you happen to start *exactly* at the bottom of the bowl. That corresponds to the equilibrium with $\hat{p} = s_2/(s_1 + s_2)$. What happens then?

Well, if you start *exactly* there, you'll stay there forever (in an infinite population). But if you start ever so slightly off the equilibrium, you'll move farther and farther away. It's what mathematicians call an *unstable equilibrium*. Any departure from that equilibrium gets larger and larger. For evolutionary purposes, we don't have to worry about a population getting to an unstable equilibrium. It never will. Unstable equilibria are ones that populations evolve away from.

When a population has only one allele present it is said to be *fixed* for that allele. Since having only one allele is also an equilibrium (in the absence of mutation), we can also call it a *monomorphic equilibrium*. When a population has more than one allele present, it is said to be *polymorphic*. If two or more alleles are present at an equilibrium, we can call it a *polymorphic equilibrium*. Thus, another way to describe the results of disruptive selection is to say that the monomorphic equilibria are stable, but the polymorphic equilibrium is not.¹⁹

¹⁸Remember that the "hats" can mean either the estimate of an unknown parameter or an equilibrium. The context will normally make it clear which meaning applies. In this case it should be pretty obvious that I'm talking about equilibria.

¹⁹Notice that a polymorphic equilibrium doesn't even exist when selection is directional.

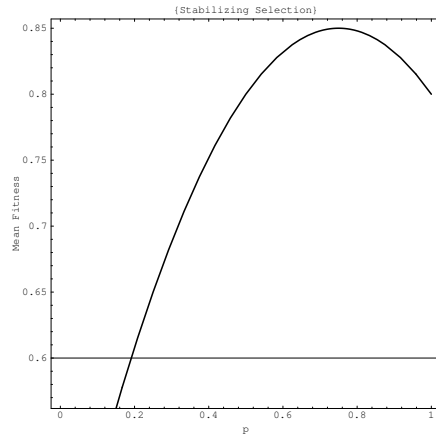


Figure 10.3: With *stabilizing selection* ($w_{11} < w_{12} > w_{22}$; also called balancing selection or heterozygote advantage) viability selection will lead to a stable polymorphism. All three genotypes will be present at equilibrium.

Stabilizing selection

If we plot \bar{w} as a function of p when $w_{11} < w_{12}$ and $w_{22} < w_{12}$, we see a third pattern. The plot is shaped like an upside down bowl (Figure 10.3).

In this case we can see that no matter what allele frequency the population starts with, the only way that $\bar{w}' \geq \bar{w}$ can hold is if the allele frequency changes in such a way that it gets close to the value where \bar{w} is maximized every generation. Unlike directional selection or disruptive selection, in which natural selection tends to eliminate one allele or the other, stabilizing selection tends to keep both alleles in the population. You'll also see this pattern of selection referred to as balancing selection, because the selection on each allele is “balanced” at the polymorphic equilibria.²⁰ We can summarize the results by saying that the monomorphic equilibria are unstable and that the polymorphic equilibrium is stable. By the way, if we write the fitness as $w_{11} = 1 - s_1$, $w_{12} = 1$, and $w_{22} = 1 - s_2$, then the allele frequency at the polymorphic equilibrium is $\hat{p} = s_2 / (s_1 + s_2)$.²¹

²⁰In fact, the marginal fitnesses are equal, i.e., $w_1 = w_2$.

²¹I'm not showing the algebra that justifies this conclusion on the off chance that you may want to test your understanding by verifying it yourself.

Chapter 11

Estimating viability

Being able to make predictions with known (or estimated) viabilities, doesn't do us a heck of a lot of good unless we can figure out what those viabilities are. Fortunately, figuring them out isn't too hard to do. If we know the number of individuals of each genotype before selection, it's really easy as a matter of fact. Consider that our data looks like this:

Genotype	A_1A_1	A_1A_2	A_2A_2
Number in zygotes	$n_{11}^{(z)}$	$n_{12}^{(z)}$	$n_{22}^{(z)}$
Viability	w_{11}	w_{12}	w_{22}
Number in adults	$n_{11}^{(a)} = w_{11}n_{11}^{(z)}$	$n_{12}^{(a)} = w_{12}n_{12}^{(z)}$	$n_{22}^{(a)} = w_{22}n_{22}^{(z)}$

In other words, estimating the absolute viability simply consists of estimating the probability that an individuals of each genotype that survive from zygote to adult. The maximum-likelihood estimate is, of course, just what you would probably guess:

$$w_{ij} = \frac{n_{ij}^{(a)}}{n_{ij}^{(z)}} \quad ,$$

Since w_{ij} is a probability and the outcome is binary (survive or die), you should be able to guess what kind of likelihood relates the observed data to the unseen parameter, namely, a binomial likelihood. In WinBUGS notation:¹

```
n.11.adult ~ dbin(w.11, n.11)
n.12.adult ~ dbin(w.12, n.11)
n.22.adult ~ dbin(w.22, n.11)
```

¹You knew you were going to see this again, didn't you?

Estimating relative viability

To estimate absolute viabilities, we have to be able to identify genotypes non-destructively. That's fine if we happen to be dealing with an experimental situation where we can do controlled crosses to establish known genotypes or if we happen to be studying an organism and a trait where we can identify the genotype from the phenotype of a zygote (or at least a very young individual) and from surviving adults.² What do we do when we can't follow the survival of individuals with known genotype? Give up?³

Remember that to make inferences about how selection will act, we only need to know *relative* viabilities, not absolute viabilities.⁴ We still need to know something about the genotypic composition of the population before selection, but it turns out that if we're only interested in relative viabilities, we don't need to follow individuals. All we need to be able to do is to score genotypes and estimate genotype frequencies before and after selection. Our data looks like this:

Genotype	A_1A_1	A_1A_2	A_2A_2
Frequency in zygotes	$x_{11}^{(z)}$	$x_{12}^{(z)}$	$x_{22}^{(z)}$
Frequency in adults	$x_{11}^{(a)}$	$x_{12}^{(a)}$	$x_{22}^{(a)}$

We also know that

$$\begin{aligned} x_{11}^{(a)} &= w_{11}x_{11}^{(z)}/\bar{w} \\ x_{12}^{(a)} &= w_{12}x_{12}^{(z)}/\bar{w} \\ x_{22}^{(a)} &= w_{22}x_{22}^{(z)}/\bar{w} \quad . \end{aligned}$$

Suppose we now divide all three equations by the middle one:

$$\begin{aligned} x_{11}^{(a)}/x_{12}^{(a)} &= w_{11}x_{11}^{(z)}/w_{12}x_{12}^{(z)} \\ 1 &= 1 \\ x_{22}^{(a)}/x_{12}^{(a)} &= w_{22}x_{22}^{(z)}/w_{12}x_{12}^{(z)} \quad , \end{aligned}$$

or, rearranging a bit

$$\begin{aligned} \frac{w_{11}}{w_{12}} &= \left(\frac{x_{11}^{(a)}}{x_{12}^{(a)}} \right) \left(\frac{x_{12}^{(z)}}{x_{11}^{(z)}} \right) \\ \frac{w_{22}}{w_{12}} &= \left(\frac{x_{22}^{(a)}}{x_{12}^{(a)}} \right) \left(\frac{x_{12}^{(z)}}{x_{22}^{(z)}} \right) \quad . \end{aligned}$$

²How many organisms and traits can you think of that satisfy this criterion? Any?

³Would I be asking the question if the answer were "Yes"?

⁴At least that's true until we start worrying about how selection and drift interact.

This gives us a complete set of relative viabilities.

Genotype	A_1A_1	A_1A_2	A_2A_2
Relative viability	$\frac{w_{11}}{w_{12}}$	1	$\frac{w_{22}}{w_{12}}$

If we use the maximum-likelihood estimates for genotype frequencies before and after selection, we obtain maximum likelihood estimates for the relative viabilities.⁵ If we use Bayesian methods to estimate genotype frequencies (including the uncertainty around those estimates), we can use these formulas to get Bayesian estimates of the relative viabilities (and the uncertainty around them).

An example

Let's see how this works with some real data from Dobzhansky's work on chromosome inversion polymorphisms in *Drosophila pseudoobscura*.⁶

Genotype	ST/ST	ST/CH	CH/CH	Total
Number in larvae	41	82	27	150
Number in adults	57	169	29	255

You may be wondering how the sample of adults can be larger than the sample of larvae. That's because to score an individual's inversion type, Dobzhansky had to kill it. The numbers in larvae are based on a sample of the population, and the adults that survived were not genotyped as larvae. As a result, all we can do is to estimate the relative viabilities.

$$\frac{w_{11}}{w_{12}} = \left(\frac{x_{11}^{(a)}}{x_{12}^{(a)}} \right) \left(\frac{x_{12}^{(z)}}{x_{11}^{(z)}} \right) = \left(\frac{57/255}{169/255} \right) \left(\frac{82/150}{41/150} \right) = 0.67$$

$$\frac{w_{22}}{w_{12}} = \left(\frac{x_{22}^{(a)}}{x_{12}^{(a)}} \right) \left(\frac{x_{12}^{(z)}}{x_{22}^{(z)}} \right) = \left(\frac{29/255}{169/255} \right) \left(\frac{82/150}{27/150} \right) = 0.52 \quad .$$

So it looks as if we have balancing selection, i.e., the fitness of the heterozygote exceeds that of either homozygote.

⁵If anyone cares, it's because of the invariance property of maximum-likelihood estimates. If you don't understand what that is, don't worry about it, just trust me.

⁶Taken from [19].

We can check to see whether this conclusion is statistically justified by comparing the observed number of individuals in each genotype category in adults with what we'd expect if all genotypes were equally likely to survive.

Genotype	ST/ST	ST/CH	CH/CH
Expected	$\left(\frac{41}{150}\right) 255$ 69.7	$\left(\frac{82}{150}\right) 255$ 139.4	$\left(\frac{27}{150}\right) 255$ 45.9
Observed	57	169	29

$\chi_2^2 = 14.82, P < 0.001$

So we have strong evidence that genotypes differ in their probability of survival.

We can also use our knowledge of how selection works to predict the genotype frequencies at equilibrium:

$$\frac{w_{11}}{w_{12}} = 1 - s_1$$

$$\frac{w_{22}}{w_{12}} = 1 - s_2 \quad .$$

So $s_1 = 0.33$, $s_2 = 0.48$, and the predicted equilibrium frequency of the ST chromosome is $s_2/(s_1 + s_2) = 0.59$.

Now all of those estimates are maximum-likelihood estimates. Doing these estimates in a Bayesian context is relatively straightforward and the details will be left as an exercise.⁷ In outline we simply

1. Estimate the genotype frequencies before and after selection as samples from a multinomial.
2. Apply the formulas above to calculate relative viabilities and selection coefficients.
3. Determine whether the 95% credible intervals for s_1 or s_2 overlap 0.⁸
4. Calculate the equilibrium frequency from $s_2/(s_1 + s_2)$, if $s_1 > 0$ and $s_2 > 0$. Otherwise, determine which fixation state will be approached.

In the end you then have not only viability estimates and their associated uncertainties, but a prediction about the ultimate composition of the population, associated with an accompanying level of uncertainty.

⁷In past years Problem #3 has consisted of making Bayesian estimates of viabilities from data like these and predicting the outcome of viability selection. That may or may not be part of Problem #3 this year. Stay tuned.

⁸Meaning that we don't have good evidence for selection either for or against the associated homozygotes, relative to the heterozygote.

Chapter 12

Selection at one locus with many alleles, fertility selection, and sexual selection

It's easy to extend the Hardy-Weinberg principle to multiple alleles at a single locus. In fact, we already did this when we were discussing the ABO blood group polymorphism. Just to get some notation out of the way, though, let's define x_{ij} as the frequency of genotype A_iA_j and p_i as the frequency of allele A_i . Then

$$x_{ij} = \begin{cases} p_i^2 & \text{if } i = j \\ 2p_i p_j & \text{if } i \neq j \end{cases}$$

Unfortunately, the simple principles we've learned for understanding selection at one locus with two alleles don't generalize completely to selection at one locus with many alleles (or even three).

- For one locus with two alleles, heterozygote advantage guarantees maintenance of a polymorphism.
- For one locus with multiple alleles, there are many different heterozygote genotypes. As a result, there is not a unique pattern identifiable as "heterozygote advantage," and selection may eliminate one or more alleles at equilibrium even if all heterozygotes have a higher fitness than all homozygotes.

Selection at one locus with multiple alleles

When we discussed selection at one locus with two alleles, I used the following set of viabilities:

$$\begin{array}{ccc} A_1A_1 & A_1A_2 & A_2A_2 \\ w_{11} & w_{12} & w_{22} \end{array}$$

You can probably guess where this is going. Namely, I'm going to use w_{ij} to denote the viability of genotype A_iA_j . What you probably wouldn't thought of doing is writing it as a matrix

$$\begin{array}{cc} & \begin{array}{cc} A_1 & A_2 \end{array} \\ \begin{array}{c} A_1 \\ A_2 \end{array} & \begin{array}{cc} w_{11} & w_{12} \\ w_{12} & w_{22} \end{array} \end{array}$$

Clearly we can extend an array like this to as many rows and columns as we have alleles so that we can summarize any pattern of viability selection with such a matrix. Notice that I didn't write both w_{12} and w_{21} , because (normally) an individual's fitness doesn't depend on whether it inherited a particular allele from its mom or its dad.¹

Marginal fitnesses and equilibria

After a little algebra it's possible to write down how allele frequencies change in response to viability selection:²

$$p'_i = \frac{p_i w_i}{\bar{w}} \quad ,$$

where $p_i = \sum p_i w_{ij}$ is the marginal fitness of allele i and $\bar{w} = \sum p_i^2 w_{ii} + \sum_i \sum_{j>i} 2p_i p_j w_{ij}$ is the mean fitness in the population.

It's easy to see³ that if the marginal fitness of an allele is less than the mean fitness of the population it will decrease in frequency. If its marginal fitness is greater than the mean fitness, it will increase in frequency. If its marginal fitness is equal to the mean fitness it won't change in frequency. So if there's a stable polymorphism, all alleles present at that equilibrium will have marginal fitnesses equal to the population mean fitness. And, since they're all equal to the same thing, they're also all equal to one another.

That's the only thing easy to say about selection with multiple alleles. To say anything more complete would require a lot of linear algebra. The only general conclusion I can

¹If it's a locus that's subject to genomic imprinting, it may be necessary to distinguish A_1A_2 from A_2A_1 . Isn't genetics fun?

²If you're ambitious (or a little weird), you might want to try to see if you can derive this yourself.

³At least it's easy to see if you've stared a lot at these things in the past.

mention, and I'll have to leave it pretty vague, is that for a complete polymorphism⁴ to be stable, none of the fitnesses can be too different from one another. Let's play with an example to illustrate what I mean.

An example

The way we always teach about sickle-cell anemia isn't entirely accurate. We talk as if there is a wild-type allele and the sickle-cell allele. In fact, there are at least three alleles at this locus in many populations where there is a high frequency of sickle-cell allele. In the wild-type, *A*, allele there is a glutamic acid at position six of the β chain of hemoglobin. In the most common sickle-cell allele, *S*, there is a valine in this position. In a rarer sickle-cell allele, *C*, there is a lysine in this position. The fitness matrix looks like this:

	<i>A</i>	<i>S</i>	<i>C</i>
<i>A</i>	0.976	1.138	1.103
<i>S</i>		0.192	0.407
<i>C</i>			0.550

There is a stable, complete polymorphism with these allele frequencies:

$$\begin{aligned} p_A &= 0.83 \\ p_S &= 0.07 \\ p_C &= 0.10 \end{aligned} .$$

If allele *C* were absent, *A* and *S* would remain in a stable polymorphism:

$$\begin{aligned} p_A &= 0.85 \\ p_S &= 0.15 \end{aligned}$$

If allele *A* were absent, however, the population would fix on allele *C*.⁵

The existence of a stable, complete polymorphism does not imply that all subsets of alleles could exist in stable polymorphisms. Loss of one allele as a result of random chance could result in a cascading loss of diversity.⁶

⁴A complete polymorphism is one in which all alleles are present.

⁵Can you explain why? Take a close look at the fitnesses, and it should be fairly obvious.

⁶The same thing can happen in ecological communities. Loss of a single species from a stable community may lead to a cascading loss of several more.

If the fitness of AS were 1.6 rather than 1.103, C would be lost from the population, although the $A - S$ polymorphism would remain.

Increasing the selection in favor of a heterozygous genotype may cause selection to eliminate one or more of the alleles not in that heterozygous genotype. This also means that if a genotype with a very high fitness in heterozygous form is introduced into a population, the resulting selection may eliminate one or more of the alleles already present.

Fertility selection

So far we've been talking about natural selection that occurs as a result of differences in the probability of survival, viability selection. There are, of course, other ways in which natural selection can occur:

- Heterozygotes may produce gametes in unequal frequencies, *segregation distortion*, or gametes may differ in their ability to participate in fertilization, *gametic selection*.
- Some genotypes may be more successful in finding mates than others, *sexual selection*.
- The number of offspring produced by a mating may depend on maternal and paternal genotypes, *fertility selection*.

In fact, most studies that have measured components of selection have identified far larger differences due to fertility than to viability. Thus, fertility selection is a very important component of natural selection in most populations of plants and animals. As we'll see a little later, it turns out that sexual selection is mathematically equivalent to a particular type of fertility selection. But before we get to that, let's look carefully at the mechanics of fertility selection.

Formulation of fertility selection

I introduced the idea of a fitness matrix earlier when we were discussing selection at one locus with more than two alleles. Even if we have only two alleles, it becomes useful to describe patterns of fertility selection in terms of a fitness matrix. Describing the matrix is easy. Writing it down gets messy. Each element in the table is simply the average number of offspring produced by a given mated pair. We write down the table with paternal genotypes in columns and maternal genotypes in rows:

Maternal genotype	Paternal genotype		
	A_1A_1	A_1A_2	A_2A_2
A_1A_1	$F_{11,11}$	$F_{11,12}$	$F_{11,22}$
A_1A_2	$F_{12,11}$	$F_{12,12}$	$F_{12,22}$
A_2A_2	$F_{22,11}$	$F_{22,12}$	$F_{22,22}$

Then the frequency of genotype A_1A_1 after one generation of fertility selection is:⁷

$$x'_{11} = \frac{x_{11}^2 F_{11,11} + x_{11}x_{12}(F_{11,12} + F_{12,11})/2 + (x_{12}^2/4)F_{12,12}}{\bar{F}}, \quad (12.1)$$

where \bar{F} is the mean fecundity of all matings in the population.⁸

It probably won't surprise you to learn that it's very difficult to say anything very general about how genotype frequencies will change when there's fertility selection. Not only are there nine different fitness parameters to worry about, but since genotypes are never guaranteed to be in Hardy-Weinberg proportion, all of the algebra has to be done on a system of three simultaneous equations.⁹ There are three weird properties that I'll mention:

1. \bar{F}' may be smaller than \bar{F} . Unlike selection on viabilities in which fitness evolved to the maximum possible value, there are situations in which fitness will evolve to the *minimum* possible value when there's selection on fertilities.¹⁰
2. A high fertility of heterozygote \times heterozygote matings is not sufficient to guarantee that the population will remain polymorphic.
3. Selection may prevent loss of either allele, but there may be no stable equilibria.

Conditions for protected polymorphism

There is one case in which it's fairly easy to understand the consequences of selection, and that's when one of the two alleles is very rare. Suppose, for example, that A_1 is very rare, then a little algebraic trickery¹¹ shows that

$$\begin{aligned} x'_{11} &\approx 0 \\ x'_{12} &\approx \frac{x_{12}(F_{12,22} + F_{22,12})/2}{F_{22,22}} \end{aligned}$$

⁷I didn't say it, but you can probably guess that I'm assuming that all of the conditions for Hardy-Weinberg apply, except for the assumption that all matings leave the same number of offspring, on average.

⁸As an exercise you might want to see if you can derive the corresponding equations for x'_{12} and x'_{22} .

⁹And you thought that dealing with one was bad enough!

¹⁰Fortunately, it takes rather weird fertility schemes to produce such a result.

¹¹The trickery isn't hard, just tedious. Justifying the trickery is a little more involved, but not too bad. If you're interested, drop by my office and I'll show you.

So A_1 will become more frequent if

$$(F_{12,22} + F_{22,12})/2 > F_{22,22} \quad (12.2)$$

Similarly, A_2 will become more frequent when it's very rare when

$$(F_{11,12} + F_{12,11})/2 > F_{11,11} \quad . \quad (12.3)$$

If both equation (12.2) and (12.3) are satisfied, natural selection will tend to prevent either allele from being eliminated. We have what's known as a *protected polymorphism*.

Conditions (12.2) and (12.3) are fairly easy to interpret intuitively: There is a protected polymorphism if the average fecundity of matings involving a heterozygote and the "resident" homozygote exceeds that of matings of the resident homozygote with itself.¹²

NOTE: It's entirely possible for neither inequality to be satisfied *and* for their to be a stable polymorphism. In other words, depending on where a population starts selection may eliminate one allele or the other or keep both segregating in the population in a stable polymorphism.¹³

Sexual selection

A classic example of sexual selection is the peacock's tail. The long, elaborate tail feathers do nothing to promote survival of male peacocks, but they are very important in determining which males attract mates and which don't. If you'll recall, when we originally derived the Hardy-Weinberg principle we said that the matings occurred randomly. Sexual selection is clearly an instance of non-random mating. Let's go back to our original mating table and see how we need to modify it to accommodate sexual selection.

¹²A "resident" homozygote is the one of which the populations is almost entirely composed when all but one allele is rare.

¹³Can you guess what pattern of fertilities is consistent with both a stable polymorphism and the *lack of* a protected polymorphism?

Mating	Frequency	Offspring genotype		
		A_1A_1	A_1A_2	A_2A_2
$A_1A_1 \times A_1A_1$	$x_{11}^f x_{11}^m$	1	0	0
A_1A_2	$x_{11}^f x_{12}^m$	$\frac{1}{2}$	$\frac{1}{2}$	0
A_2A_2	$x_{11}^f x_{22}^m$	0	1	0
$A_1A_2 \times A_1A_1$	$x_{12}^f x_{11}^m$	$\frac{1}{2}$	$\frac{1}{2}$	0
A_1A_2	$x_{12}^f x_{12}^m$	$\frac{1}{4}$	$\frac{1}{2}$	$\frac{1}{4}$
A_1A_2	$x_{12}^f x_{22}^m$	0	$\frac{1}{2}$	$\frac{1}{2}$
$A_2A_2 \times A_1A_1$	$x_{22}^f x_{11}^m$	0	1	0
A_1A_2	$x_{22}^f x_{12}^m$	0	$\frac{1}{2}$	$\frac{1}{2}$
A_2A_2	$x_{22}^f x_{22}^m$	0	0	1

What I've done is to assume that there is random mating in the populations *among those individuals that are included in the mating pool*. We'll assume that all females are mated so that $x_{ij}^f = x_{ij}$.¹⁴ We'll let the relative attractiveness of the male genotypes be a_{11} , a_{12} , and a_{22} . Then it's not too hard to convince yourself that

$$\begin{aligned} x_{11}^m &= \frac{x_{11}a_{11}}{\bar{a}} \\ x_{12}^m &= \frac{x_{12}a_{12}}{\bar{a}} \\ x_{22}^m &= \frac{x_{22}a_{22}}{\bar{a}} \end{aligned} ,$$

where $\bar{a} = x_{11}a_{11} + x_{12}a_{12} + x_{22}a_{22}$. A little more algebra and you can see that

$$x'_{11} = \frac{x_{11}^2 a_{11} + x_{11} x_{12} (a_{12} + a_{11}) / 2 + x_{12}^2 a_{12} / 4}{\bar{a}} \quad (12.4)$$

And we could derive similar equations for x'_{12} and x'_{22} . Now you're not likely to remember this, but equation (12.4) bears a striking resemblance to one you saw earlier, equation (12.1). In fact, sexual selection is equivalent to a particular type of fertility selection, in terms of how genotype frequencies will change from one generation to the next. Specifically, the fertility matrix corresponding to sexual selection on a male trait is:

	A_1A_1	A_1A_2	A_2A_2
A_1A_1	a_{11}	a_{12}	a_{22}
A_1A_2	a_{11}	a_{12}	a_{22}
A_2A_2	a_{11}	a_{12}	a_{22}

¹⁴There's a reason for doing this called Bateman's principle that we can discuss, if you'd like.

There are, of course, a couple of other things that make sexual selection interesting. First, traits that are sexually selected in males often come at a cost in viability, so there's a tradeoff between survival and reproduction that can make the dynamics complicated and interesting. Second, the evolution of a sexually selected trait involves two traits: the male characteristic that is being selected and a female preference for that trait. In fact the two tend to become associated so that the female preference evokes a sexually selected response in males, which evokes a stronger preference in females, and so on and so on. This is a process Fisher referred to as "runaway sexual selection."

Chapter 13

Selection Components Analysis

Consider the steps in a transition from one generation to the next, starting with a newly formed zygote:

- Zygote
- Adult — Survival from zygote to adult may differ between the sexes.
- Breeding population — Adult genotypes may differ in their probability of mating, and the differences may be different in males and females
- Newly formed zygotes

When the transition from one stage to the next depends on genotype, then selection has occurred at that stage. Thus, to determine whether selection is occurring we construct expectations of genotype or allele frequencies at one stage based on the frequencies at the immediately preceding stage assuming that no selection has occurred. Then we compare observed frequencies to those expected without selection. If they match, we have no evidence for selection. If they don't match, we do have evidence for selection.

As we've already seen, it's conceptually easy (if often experimentally difficult) to detect and measure selection if we can assay genotypes non-destructively at appropriate stages in the life-cycle. What if we can't? Well, there's a very nice approach known as *selection components analysis* that generalizes the approach to estimating relative viabilities that we've already seen [13].

The Data

Pregnant mothers are collected. One offspring from each mother is randomly selected and its genotype determined. In addition, the genotypes of a random sample of non-reproductive (“sterile”) females and adult males are determined. The data can be summarized as follows:

Mother	Offspring			Σ	“Sterile”	
	Females	Males	Females		Males	
A_1A_1	C_{11}	C_{12}	—	F_1	S_1	M_1
A_1A_2	C_{21}	C_{22}	C_{23}	F_2	S_2	M_2
A_2A_2	—	C_{32}	C_{33}	F_3	S_3	M_3
Total				F_0	S_0	M_0

Given the total sample size for mother-offspring pairs, “sterile” females, and males, how many free parameters are there? How many frequencies would we need to know to reproduce the data?

6	for mother-offspring pairs
2	for “sterile” females
2	for males
10	total

The Analysis

H_1 : **Half of the offspring from heterozygous mothers are also heterozygous.**

Under H_1

$$\gamma_{21} = (1/2)(F_2/F_0)(C_{21}/(C_{21} + C_{23}))$$

$$\gamma_{22} = (1/2)(F_2/F_0)$$

$$\gamma_{23} = (1/2)(F_2/F_0)(C_{23}/(C_{21} + C_{23}))$$

Under H_1 , γ_{22} can be predicted just from the frequency of heterozygous mothers in the sample. Thus, only 9 parameters are needed to describe the data under H_1 . Since 10 are required if we reject H_1 we can use a likelihood ratio test with one degree of freedom to see whether the above estimates provide an adequate description of the data.

If H_1 is rejected, we can conclude that there is either gametic selection or segregation distortion in A_1A_2 females.

H_2 : The frequency of transmitted male gametes is independent of the mother's genotype.

Under H_2

$$p_m = (C_{11} + C_{21} + C_{32}) / (F_0 - C_{22})$$

$$q_m = (C_{12} + C_{23} + C_{33}) / (F_0 - C_{22})$$

The expected frequency of the various mother-offspring combinations is

	A_1A_1	A_1A_2	A_2A_2
A_1A_1	$\phi_1 p_m$	$\phi_1 q_m$	—
A_1A_2	$(1/2)\phi_2 p_m$	$(1/2)\phi_2$	$(1/2)\phi_2 q_m$
A_2A_2	—	$\phi_3 p_m$	$\phi_3 q_m$

where $\phi_i = F_i/F_0$. Under H_2 only the female genotype frequencies and the male gamete frequencies are needed to describe the mother-offspring data. That's a total of $2 + 1 + 2 + 2 = 7$ frequencies needed to describe *all* of the data. Since H_1 needed 9, that gives us 2 degrees of freedom for our likelihood ratio test of H_2 given H_1 .

If H_2 is rejected, we can conclude that there is some form of non-random mating in the breeding population or female-specific selection of male gametes.

H_3 : The frequency of the transmitted male gametes is equal to the allele frequency in adult males.

Under H_3 the maximum likelihood estimates for p_m and q_m cannot be found explicitly, they are a complicated function of p_m and q_m as defined under H_2 and of M_1 , M_2 , and M_3 . Under H_3 , however, we no longer need to account separately for the gamete frequency in males, so a total of $2 + 2 + 2 = 6$ frequencies is needed to describe the data. Since H_2 needed 7, that gives us 1 degree of freedom for our likelihood ratio test of H_3 given H_2 .

If H_3 is rejected, we can conclude either that males differ in their ability to attract mates (i.e., there is sexual selection) or that male gametes differ in their ability to accomplish fertilization (e.g., sperm competition), or that there is segregation distortion in A_1A_2 males.

H_4 : The genotype frequencies of reproductive females are the same as those of “sterile” females.

Under H_4 the maximum likelihood estimates for the genotype frequencies in females are

$$\phi_i = (F_i + S_i)/(F_0 + S_0)$$

Under H_4 we no longer need to account separately for the genotype frequencies in “sterile” females, so a total of $2 + 2 = 4$ frequencies is needed to describe the data. Since H_3 needed 6, that gives us 2 degrees of freedom for our likelihood ratio test of H_4 given H_3 .

If H_4 is rejected, we can conclude that females differ in their ability to reproduce successfully.

H_5 : The genotype frequencies of adult females and adult males are equal.

Under H_5 the maximum likelihood estimates for the adult genotype frequencies can not be found explicitly. Instead, they are a complicated function of almost every piece of information that we have. Under H_5 , however, we no longer need to account separately for the genotype frequencies in females and males, so a total of 2 frequencies is needed to describe the data. Since H_4 needed 4, that gives us 2 degrees of freedom for our likelihood ratio test of H_5 given H_4 .

If H_5 is rejected we can conclude that the relative viabilities of the genotypes are different in the two sexes. (We have assumed implicitly throughout that the locus under study is an autosomal locus. Notice that rejection of H_5 is consistent with *no* selection in one sex.)

H_6 : **The genotype frequencies in the adult population are equal to those of the zygote population.**

Under H_6 the maximum-likelihood estimator for the allele frequency in the population is

$$p = \frac{((C_{11} + C_{21} + C_{32}) + 2(F_1 + S_1 + M_1) + (F_2 + S_2 + M_2))}{((F_0 - C_{21}) + F_0 + S_0 + M_0)}$$

Under H_6 the genotype frequencies in our original table can be summarized as follows:

Mother	A_1A_1	A_1A_2	A_2A_2	Σ	“Sterile” Females	Males
A_1A_1	p^3	p^2q	0	p^2	p^2	p^2
A_1A_2	p^2q	pq	pq^2	$2pq$	$2pq$	$2pq$
A_2A_2	0	pq^2	q^3	q^2	q^2	q^2

In short, under H_6 only one parameter, the allele frequency, is required to describe the entire data set. Since under H_5 needed two parameters, our likelihood ratio test of H_6 given H_5 will have one degree of freedom.

If H_6 is rejected, we can conclude that genotypes differ in their probability of survival from zygote to adult, i.e., that there is viability selection. If H_1 – H_6 are accepted, we have no evidence that selection is happening at any stage of the life cycle at this locus and no evidence of non-random mating with respect to genotype at this locus.

An example

This data is from a 2-allelic esterase polymorphism in our old friend *Zoarces viviparus*, the eelpout. The observations are in roman type in the table below. The numbers in italics are those expected if hypotheses H_1 – H_6 are accepted.

Mother	A_1A_1	A_1A_2	A_2A_2	Σ	“Sterile” Females	Males
	41	70	—	111	8	54
A_1A_1	<i>39.0</i>	<i>67.0</i>	—	<i>106.0</i>	<i>9.3</i>	<i>58.4</i>
	65	173	119	357	32	200
A_1A_2	<i>67.0</i>	<i>181.9</i>	<i>114.9</i>	<i>363.8</i>	<i>32.1</i>	<i>200.5</i>
	—	127	187	314	29	177
A_2A_2	—	<i>114.9</i>	<i>197.3</i>	<i>312.2</i>	<i>27.6</i>	<i>172.1</i>
Sum	106	370	306	782	69	431
	<i>106.0</i>	<i>363.8</i>	<i>312.2</i>	—	—	—

The results of the series of hypothesis tests is as follows:

Hypothesis	Degrees of freedom	χ^2	P	50% power point
H_1	1	0.34	>0.50	0.05
H_2	2	1.37	>0.50	≤ 0.09
H_3	1	0.98	>0.30	≤ 0.05
H_4	2	0.37	>0.50	≤ 0.10
H_5	2	0.22	>0.80	≤ 0.05
H_6	1	0.09	>0.70	0.03

We conclude from this analysis that there is no evidence of selection on the genetic variation at the esterase locus in *Zoarces viviparus* and that there is no evidence of non-random mating with respect to genotype at this locus. The power calculations increase our confidence that if there is selection happening, the differences among genotypes are on the order of just a few percent.

Part III
Genetic drift

Chapter 14

Genetic Drift

So far in this course we've talked about changes in genotype and allele frequencies as if they were completely deterministic. Given the current allele frequencies and viabilities, for example, we wrote down an equation describing how they will change from one generation to the next:

$$p' = \frac{p^2 w_{11} + pq w_{12}}{\bar{w}} \quad .$$

Notice that in writing this equation, we're claiming that we can predict the allele frequency in the next generation *without error*. But suppose the population is small, say 10 diploid individuals, and our prediction is that $p' = 0.5$. Then just as we wouldn't be surprised if we flipped a coin 20 times and got 12 heads, we shouldn't be surprised if we found that $p' = 0.6$. The difference between what we expect ($p' = 0.5$) and what we observe ($p' = 0.6$) can be chalked up to statistical sampling error. That sampling error is the cause of (or just another name for) *genetic drift*—the tendency for allele frequencies to change from one generation to the next in a finite population even if there is no selection.

A simple example

To understand in more detail what happens when there is genetic drift, let's consider the simplest possible example: a haploid population consisting of 2 individuals.¹ Suppose that there are initially two alleles in this population A_1 and A_2 . This implies that $p = q = 0.5$, but we'll ignore that numerical fact and simply imagine that the frequency of the A_1 allele is p .

We imagine the following scenario:

¹Notice that once we start talking about genetic drift, we have to specify the size of the population.

- Each individual in the population produces a very large number of offspring.
- Each offspring is an identical copy of its parent, i.e., A_1 begets A_1 and A_2 begets A_2 . In other words, there's no mutation.
- The next generation is constructed by picking two offspring at random from the very large number of offspring produced by these two individuals.

Then it's not too hard to see that

$$\begin{aligned} \text{Probability that both offspring are } A_1 &= p^2 \\ \text{Probability that one offspring is } A_1 \text{ and one is } A_2 &= 2pq \\ \text{Probability that both offspring are } A_2 &= q^2 \end{aligned}$$

Of course $p = 1$ if both offspring sampled are A_1 , $p = 1/2$ if one is A_1 and one is A_2 , and $p = 0$ if both are A_2 , so that set of equations is equivalent to this one:

$$P(p = 1) = p^2 \tag{14.1}$$

$$P(p = 1/2) = 2pq \tag{14.2}$$

$$P(p = 0) = q^2 \tag{14.3}$$

In other words, we can no longer predict with certainty what allele frequencies in the next generation will be. We can only assign probabilities to each of the three possible outcomes. Of course, in a larger population the amount of uncertainty about the allele frequencies will be smaller,² but there will be *some* uncertainty associated with the predicted allele frequencies unless the population is infinite.

The probability of ending up in any of the three possible states obviously depends on the current allele frequency. In probability theory we express this dependence by writing equations (14.1)–(14.3) as conditional probabilities:

$$P(p_1 = 1|p_0) = p_0^2 \tag{14.4}$$

$$P(p_1 = 1/2|p_0) = 2p_0q_0 \tag{14.5}$$

$$P(p_1 = 0|p_0) = q_0^2 \tag{14.6}$$

I've introduced the subscripts so that we can distinguish among various generations in the process. Why? Because if we can write equations (14.4)–(14.6), we can also write the

²More about that later.

following equations:³

$$\begin{aligned} P(p_2 = 1|p_1) &= p_1^2 \\ P(p_2 = 1/2|p_1) &= 2p_1q_1 \\ P(p_2 = 0|p_1) &= q_1^2 \end{aligned}$$

Now if we stare at those a little while, we⁴ begin to see some interesting possibilities. Namely,

$$\begin{aligned} P(p_2 = 1|p_0) &= P(p_2 = 1|p_1 = 1)P(p_1 = 1|p_0) + P(p_2 = 1|p_1 = 1/2)P(p_1 = 1/2|p_0) \\ &= (1)(p_0^2) + (1/4)(2p_0q_0) \\ &= p_0^2 + (1/2)p_0q_0 \\ P(p_2 = 1/2|p_0) &= P(p_2 = 1/2|p_1 = 1/2)P(p_1 = 1/2|p_0) \\ &= (1/2)(2p_0q_0) \\ &= p_0q_0 \\ P(p_2 = 0|p_0) &= P(p_2 = 0|p_1 = 0)P(p_1 = 0|p_0) + P(p_2 = 0|p_1 = 1/2)P(p_1 = 1/2|p_0) \\ &= (1)(q_0^2) + (1/4)(2p_0q_0) \\ &= q_0^2 + (1/2)p_0q_0 \end{aligned}$$

It takes more algebra than I care to show,⁵ but these equations can be extended to an arbitrary number of generations.

$$\begin{aligned} P(p_t = 1|p_0) &= p_0^2 + \left(1 - (1/2)^{t-1}\right) p_0q_0 \\ P(p_t = 1/2|p_0) &= p_0q_0(1/2)^{t-2} \\ P(p_t = 0|p_0) &= q_0^2 + \left(1 - (1/2)^{t-1}\right) p_0q_0 \end{aligned}$$

Why do I bother to show you these equations?⁶ Because you can see pretty quickly that as t gets big, i.e., the longer our population evolves, the smaller the probability that $p_t = 1/2$ becomes. In fact, it's not hard to verify two facts about genetic drift in this simple situation:

1. One of the two alleles originally present in the population is certain to be lost eventually.

³I know. I'm weird. I actually get a kick out of writing equations!

⁴Or at least the weird ones among us

⁵Ask me, if you're really interested.

⁶It's not just that I'm crazy.

2. The probability that A_1 is fixed is equal to its initial frequency, p_0 , and the probability that A_2 is fixed is equal to its initial frequency, q_0 .

Both of these properties are true in general for *any* finite population and *any* number of alleles.

1. Genetic drift will eventually lead to loss of all alleles in the population except one.⁷
2. The probability that any allele will eventually become fixed in the population is equal to its current frequency.

General properties of genetic drift

What I've shown you so far applies only to a haploid population with two individuals. Even I will admit that it isn't a very interesting situation. Suppose, however, we now consider a population with N diploid individuals. We can treat it as if it were a population of $2N$ haploid individuals using a direct analogy to the process I described earlier, and then things start to get a little more interesting.

- Each individual in the population produces a large number of gametes.
- Each gamete is an identical copy of its parent, i.e., A_1 begets A_1 and A_2 begets A_2 .
- The next generation is constructed by picking $2N$ gametes at random from the large number originally produced.

We can then write a general expression for how allele frequencies will change between generations. Specifically, the distribution describing the probability that there will be j copies of A_1 in the next generation given that there are i copies in this generation is

$$P(j \text{ } A_1 \text{ in offspring} \mid i \text{ } A_1 \text{ in parents}) = \binom{2N}{j} \left(\frac{i}{2N}\right)^j \left(1 - \frac{i}{2N}\right)^{2N-j},$$

i.e., a binomial distribution. I'll be astonished if any of what I'm about to say is apparent to any of you, but this equation implies three really important things. We've encountered two already:

⁷You obviously can't lose all of them unless the population becomes extinct.

- Allele frequencies will tend to change from one generation to the next purely as a result of sampling error. As a consequence, genetic drift will eventually lead to loss of all alleles in the population except one.
- The probability that any allele will eventually become fixed in the population is equal to its current frequency.
- The population has no memory.⁸ The probability that the offspring generation will have a particular allele frequency depends *only* on the allele frequency in the parental generation. It does not depend on how the parental generation came to have that allele frequency. This is exactly analogous to coin-tossing. The probability that you get a heads on the next toss of a fair coin is 1/2. It doesn't matter whether you've never tossed it before or if you've just tossed 25 heads in a row.⁹

Variance of allele frequencies between generations

For a binomial distribution

$$\begin{aligned}
 P(K = k) &= \binom{N}{k} p^k (1-p)^{N-k} \\
 \text{Var}(K) &= Np(1-p) \\
 \text{Var}(p) &= \text{Var}(K/N) \\
 &= \frac{1}{N^2} \text{Var}(K) \\
 &= \frac{p(1-p)}{N}
 \end{aligned}$$

Applying this to our situation,

$$\text{Var}(p_{t+1}) = \frac{p_t(1-p_t)}{2N}$$

$\text{Var}(p_{t+1})$ measures the amount of uncertainty about allele frequencies in the next generation, given the current allele frequency. As you probably guessed long ago, the amount of uncertainty is inversely proportional to population size. The larger the population, the smaller the uncertainty.

⁸Technically, we've described a Markov chain with a finite state space, but I doubt that you really care about that.

⁹Of course, if you've just tossed 25 heads in a row, you could be forgiven for having your doubts about whether the coin is actually fair.

If you think about this a bit, you might expect that a smaller variance would “slow down” the process of genetic drift—and you’d be right. It takes some pretty advanced mathematics to say how much the process slows down as a function of population size,¹⁰ but we can summarize the result in the following equation:

$$\bar{t} \approx -4N (p \log p + (1 - p) \log(1 - p)) \quad ,$$

where \bar{t} is the average time to fixation of one allele or the other and p is the current allele frequency.¹¹ So the average time to fixation of one allele or the other increases approximately linearly with increases in the population size.

Analogy to inbreeding

You may have noticed some similarities between drift and inbreeding. Specifically, both processes lead to a loss of heterozygosity and an increase in homozygosity. This analogy leads to a useful heuristic for helping us to understand the dynamics of genetic drift.

Remember our old friend f , the inbreeding coefficient? I’m going to re-introduce you to it in the form of the population inbreeding coefficient, the probability that two alleles chosen at random from a population are identical by descent. We’re going to study how the population inbreeding coefficient changes from one generation to the next as a result of reproduction in a finite population.¹²

$$\begin{aligned} f_{t+1} &= \text{Prob. ibd from preceding generation} \\ &\quad + (\text{Prob. not ibd from prec. gen.}) \times (\text{Prob. ibd from earlier gen.}) \\ &= \frac{1}{2N} + \left(1 - \frac{1}{2N}\right) f_t \end{aligned}$$

or, in general,

$$f_{t+1} = 1 - \left(1 - \frac{1}{2N}\right)^t (1 - f_0) \quad .$$

Summary

There are four characteristics of genetic drift that I think are particularly important for you to remember:

¹⁰Actually, we’ll encounter a way that isn’t quite so hard in a few lectures when we get to the coalescent.

¹¹Notice that this equation only applies to the case of one-locus with two alleles, although the principle applies to any number of alleles.

¹²Remember that I use the abbreviation ibd to mean identical by descent.

1. Allele frequencies tend to change from one generation to the next simply as a result of sampling error. We can specify a probability distribution for the allele frequency in the next generation, but we cannot predict the actual frequency with certainty.
2. There is no systematic bias to changes in allele frequency. The allele frequency is as likely to increase from one generation to the next as it is to decrease.
3. If the process is allowed to continue long enough without input of new genetic material through migration or mutation, the population will eventually become fixed for only one of the alleles originally present.¹³
4. The time to fixation on a single allele is directly proportional to population size, and the amount of uncertainty associated with allele frequencies from one generation to the next is inversely related to population size.

Effective population size

I didn't make a big point of it, but in our discussion of genetic drift so far we've assumed everything about populations that we assumed to derive the Hardy-Weinberg principle, *and* we've assumed that:

- We can model drift in a finite population as a result of sampling among haploid gametes rather than as a result of sampling among diploid genotypes. Since we're dealing with a finite population, this effectively means that the two gametes incorporated into an individual could have come from the same parent, i.e., self-fertilization occurs when there's random union of gametes in a finite, diploid population.
- Since we're sampling gametes rather than individuals, we're also implicitly assuming that there aren't separate sexes.¹⁴
- The number of gametes any individual has represented in the next generation is a binomial random variable.¹⁵
- The population size is constant.

¹³This will hold true even if there is strong selection for keeping alleles in the population. Selection can't prevent loss of diversity, only slow it down.

¹⁴How could there be separate sexes if there can be self-fertilization?

¹⁵More about this later.

How do we deal with the fact that one or more of these conditions will be violated in just about any case we're interested in?¹⁶ One way would be to develop all the probability models that incorporate that complexity and try to solve them. That's nearly impossible, except through computer simulations. Another, and by far the most common approach, is to come up with a conversion formula that makes our actual population seem like the "ideal" population that we've been studying. That's exactly what *effective population size* is.

The effective size of a population is the size of an ideal population that has the same properties with respect to genetic drift as our actual population does.

What does that phrase "same properties with respect to genetic drift" mean? Well there are two ways it can be defined.¹⁷

Variance effective size

You may remember¹⁸ that the variance in allele frequency in an ideal population is

$$\text{Var}(p_{t+1}) = \frac{p_t(1-p_t)}{2N} .$$

So one way we can make our actual population equivalent to an ideal population to make their allele frequency variances the same. We do this by calculating the variance in allele frequency for our actual population, figuring out what size of ideal population would produce the same variance, and pretending that our actual population is the same as an ideal population of the same size. To put that into an equation,¹⁹ let $\widehat{\text{Var}}(p)$ be the variance we calculate for our actual population. Then

$$N_e^{(v)} = \frac{p(1-p)}{2\widehat{\text{Var}}(p)}$$

is the *variance effective population size*, i.e., the size of an ideal population that has the same properties with respect to allele frequency variance as our actual population.

Inbreeding effective size

You may also remember that we can think of genetic drift as analogous to inbreeding. The probability of identity by descent within populations changes in a predictable way in relation

¹⁶OK, OK. They will probably be violated in *every* case we're interested in.

¹⁷There are actually more than two ways, but we're only going to talk about two.

¹⁸You probably won't, so I'll remind you

¹⁹As if that will make it any clearer. Does anyone actually read these footnotes?

to population size, namely

$$f_{t+1} = \frac{1}{2N} + \left(1 - \frac{1}{2N}\right) f_t \quad .$$

So another way we can make our actual population equivalent to an ideal population is to make them equivalent with respect to how f changes from generation to generation. We do this by calculating how the inbreeding coefficient changes from one generation to the next in our actual population, figuring out what size an ideal population would have to be to show the same change between generations, and pretending that our actual population is the same size at the ideal one. So suppose \hat{f}_t and \hat{f}_{t+1} are the actual inbreeding coefficients we'd have in our population at generation t and $t + 1$, respectively. Then

$$\begin{aligned} \hat{f}_{t+1} &= \frac{1}{2N_e^{(f)}} + \left(1 - \frac{1}{2N_e^{(f)}}\right) \hat{f}_t \\ &= \left(\frac{1}{2N_e^{(f)}}\right) (1 - \hat{f}_t) + \hat{f}_t \\ \hat{f}_{t+1} - \hat{f}_t &= \left(\frac{1}{2N_e^{(f)}}\right) (1 - \hat{f}_t) \\ N_e^{(f)} &= \frac{1 - \hat{f}_t}{2(\hat{f}_{t+1} - \hat{f}_t)} \quad . \end{aligned}$$

In many applications it's convenient to assume that $\hat{f}_t = 0$. In that case the calculation gets a lot simpler:

$$N_e^{(f)} = \frac{1}{2\hat{f}_{t+1}} \quad .$$

We also don't lose anything by doing so, because $N_e^{(f)}$ depends only on how much f *changes* from one generation to the next, not on its actual magnitude.

Comments on effective population sizes

Those are nice tricks, but there are some limitations. The biggest is that $N_e^{(v)} \neq N_e^{(f)}$ if the population size is changing from one generation to the next.²⁰ So you have to decide which of these two measures is more appropriate for the question you're studying.

²⁰It's even worse than that. When the population size is changing, it's not clear that any of the available adjustments to produce an effective population size are entirely satisfactory. Well, that's not entirely true either. Fu et al. [27] show that there is a reasonable definition in one simple case when the population size varies, and it happens to correspond to the solution presented below.

- $N_e^{(f)}$ is naturally related to the number of individuals in the parental populations. It tells you something about how the probability of identity by descent within a single population will change over time.
- $N_e^{(v)}$ is naturally related to the number of individuals in the offspring generation. It tells you something about how much allele frequencies in isolated populations will diverge from one another.

Examples

This is all pretty abstract. Let's work through some examples to see how this all plays out.²¹ In the case of separate sexes and variable population size, I'll provide a derivation of $N_e^{(f)}$. In the case of differences in the number of offspring left by individuals, I'll just give you the formula and we'll discuss some of the implications.

Separate sexes

We'll start by assuming that $\hat{f}_t = 0$ to make the calculations simple. So we know that

$$N_e^{(f)} = \frac{1}{2\hat{f}_{t+1}} \quad .$$

The first thing to do is to calculate \hat{f}_{t+1} . To do this we have to break the problem down into pieces.²²

- We assumed that $\hat{f}_t = 0$, so the only way for two alleles to be identical by descent is if they are identical copies of the *same* allele in the immediately preceding generation.
- Even if the numbers of reproductive males and reproductive females are different, every new offspring has exactly one father and one mother. Thus, the probability that the first gamete selected at random is female is just 1/2, and the probability that the second gamete selected is male is just 1/2.
- The probability that the second gamete selected is female given that the first one we selected was female is $(N - 1)/(2N - 1)$, because N out of the $2N$ alleles represented

²¹If you're interested in a comprehensive list of formulas relating various demographic parameters to effective population size, take a look at [16, p. 362]. They provide a pretty comprehensive summary and a number of derivations.

²²Remembering, of course, that \hat{f}_{t+1} is the probability that two alleles drawn at random are identical by descent.

among offspring came from females, and there are only $N - 1$ out of $2N - 1$ left after we've already picked one. The same logic applies for male gametes.

- Finally, the probability that one particular female gamete was chosen is $1/2N_f$, where N_f is the number of females in the population. Similarly the probability that one particular male gamete was chosen is $1/2N_m$, where N_m is the number of males in the population.

With those facts in hand, we're ready to calculate \hat{f}_{t+1} .

$$\begin{aligned} f_{t+1} &= \left(\frac{1}{2}\right) \left(\frac{N-1}{2N-1}\right) \left(\frac{1}{2N_f}\right) + \left(\frac{1}{2}\right) \left(\frac{N-1}{2N-1}\right) \left(\frac{1}{2N_m}\right) \\ &= \left(\frac{1}{2}\right) \left(\frac{N-1}{2N-1}\right) \left(\frac{1}{2N_f} + \frac{1}{2N_m}\right) \\ &\approx \left(\frac{1}{4}\right) \left(\frac{1}{2N_f} + \frac{1}{2N_m}\right) \end{aligned}$$

So,

$$N_e^{(f)} \approx \frac{4N_f N_m}{N_f + N_m} .$$

What does this all mean? Well, consider a couple of important examples. Suppose the numbers of females and males in a population are equal, $N_f = N_m = N/2$. Then

$$\begin{aligned} N_e^{(f)} &= \frac{4(N/2)(N/2)}{N/2 + N/2} \\ &= \frac{4N^2/4}{N} \\ &= N . \end{aligned}$$

The effective population size is equal to the actual population size if the sex ratio is 50:50. If it departs from 50:50, the effective population size will be smaller than the actual population size. Consider the extreme case where there's only one reproductive male in the population. Then

$$N_e^{(f)} = \frac{4N_f}{N_f + 1} . \tag{14.7}$$

Notice what this equation implies: The effective size of a population with only one reproductive male (or female) can *never* be bigger than 4, no matter how many mates that individual has and no matter how many offspring are produced.

Variable population size

The notation for this one gets a little more complicated, but the ideas are simpler than those you just survived. Since the population size is changing we need to specify the population size at each time step. Let N_t be the population size in generation t . Then

$$\begin{aligned} f_{t+1} &= \left(1 - \frac{1}{2N_t}\right) f_t + \frac{1}{2N_t} \\ 1 - f_{t+1} &= \left(1 - \frac{1}{2N_t}\right) (1 - f_t) \\ 1 - f_{t+K} &= \left(\prod_{i=1}^K \left(1 - \frac{1}{2N_{t+i}}\right)\right) (1 - f_t) \quad . \end{aligned}$$

Now if the population size were constant

$$\left(\prod_{i=1}^K \left(1 - \frac{1}{2N_{t+i}}\right)\right) = \left(1 - \frac{1}{2N_e^{(f)}}\right)^K \quad .$$

Dealing with products and powers is inconvenient, but if we take the logarithm of both sides of the equation we get something simpler:

$$\sum_{i=1}^K \log \left(1 - \frac{1}{2N_{t+i}}\right) = K \log \left(1 - \frac{1}{2N_e^{(f)}}\right) \quad .$$

It's a well-known fact²³ that $\log(1 - x) \approx -x$ when x is small. So if we assume that N_e and all of the N_t are large,²⁴ then

$$\begin{aligned} K \left(-\frac{1}{2N_e^{(f)}}\right) &= \sum_{i=1}^K -\frac{1}{2N_{t+i}} \\ \frac{K}{N_e^{(f)}} &= \sum_{i=1}^K \frac{1}{N_{t+i}} \\ N_e^{(f)} &= \left(\left(\frac{1}{K}\right) \sum_{i=1}^K \frac{1}{N_{t+i}}\right)^{-1} \end{aligned}$$

The quantity on the right side of that last equation is a well-known quantity. It's the *harmonic mean* of the N_t . It's another well-known fact²⁵ that the harmonic mean of a series

²³Well known to some of us at least.

²⁴So that there reciprocals are small

²⁵Are we ever going to run out of well-known facts? Probably not.

of numbers is always less than its arithmetic mean. This means that genetic drift may play a much more important role than we might have imagined, since the effective size of a population will be more influenced by times when it is small than by times when it is large.

Consider, for example, a population in which N_1 through N_9 are 1000, and N_{10} is 10.

$$N_e = \left(\left(\frac{1}{10} \right) \left(9 \left(\frac{1}{1000} \right) + \left(\frac{1}{10} \right) \right) \right)^{-1}$$

$$\approx 92$$

versus an arithmetic average of 901. So the population will behave with respect to the inbreeding associated with drift like a population a tenth of its arithmetic average size.

Variation in offspring number

I'm just going to give you this formula. I'm not going to derive it for you.²⁶

$$N_e^{(f)} = \frac{2N - 1}{1 + \frac{V_k}{2}} ,$$

where V_k is the variance in number of offspring among individuals in the population. Remember I told you that the number of gametes any individual has represented in the next generation is a binomial random variable in an ideal population? Well, if the population size isn't changing, that means that $V_k = 2(1 - 1/N)$ in an ideal population.²⁷ A little algebra should convince you that in this case $N_e^{(f)} = N$. It can also be shown (with more algebra) that

- $N_e^{(f)} < N$ if $V_k > 2(1 - 1/N)$ and
- $N_e^{(f)} > N$ if $V_k < 2(1 - 1/N)$.

That last fact is pretty remarkable. Conservation biologists try to take advantage of it to decrease the loss of genetic variation in small populations, especially those that are captive bred. If you can reduce the variance in reproductive success, you can substantially increase the effective size of the population. In fact, if you could reduce V_k to zero, then

$$N_e^{(f)} = 2N - 1 .$$

The effective size of the population would then be almost twice its actual size.

²⁶The details are in [16], if you're interested.

²⁷The calculation is really easy, and I'd be happy to show it to you if you're interested.

Chapter 15

Mutation, Migration, and Genetic Drift

So far in this course we've focused on single, isolated populations, and we've imagined that there isn't any mutation. We've also completely ignored the ultimate source of all genetic variation — mutation.¹ We're now going to study what happens when we consider multiple populations simultaneously and when we allow mutation to happen. Let's consider mutation first, because it's the easiest to understand.

Drift and mutation

Remember that in the absence of mutation

$$f_{t+1} = \left(\frac{1}{2N}\right) + \left(1 - \frac{1}{2N}\right) f_t \quad , \quad (15.1)$$

One way of modeling mutation is to assume that every time a mutation occurs it introduces a new allele into the population. This model is referred to as the *infinite alleles model*, because it implicitly assumes that there is potentially an infinite number of alleles. Under this model we need to make only one simple modification to equation (15.1). We have to multiply the expression on the right by the probability that neither allele mutated:

$$f_{t+1} = \left(\left(\frac{1}{2N}\right) + \left(1 - \frac{1}{2N}\right) f_t\right) (1 - \mu)^2 \quad , \quad (15.2)$$

¹Well, that's not quite true. We talked about multiple populations when we talked about the Wahlund effect and Wright's F_{ST} , but we didn't talk explicitly about any of the evolutionary processes associated with multiple populations.

where μ is the mutation rate, i.e., the probability that an allele in an offspring is different from the allele it was derived from in a parent. In writing down this expression, the reason this is referred to as an infinite alleles model becomes apparent: we are assuming that every time a mutation occurs it produces a new allele. The only way in which two alleles can be identical is if neither mutated.²

So where do we go from here? Well, if you think about it, mutation is always introducing new alleles that, by definition, are different from any of the alleles currently in the population. It stands to reason, therefore, that we'll never be in a situation where all of the alleles in a population are identical by descent as they would be in the absence of mutation. In other words we expect there to be an equilibrium between loss of diversity through genetic drift and the introduction of diversity through mutation.³ From the definition of an equilibrium,

$$\begin{aligned}
 \hat{f} &= \left(\left(\frac{1}{2N} \right) + \left(1 - \frac{1}{2N} \right) \hat{f} \right) (1 - \mu)^2 \\
 \hat{f} \left(1 - \left(1 - \frac{1}{2N} \right) (1 - \mu)^2 \right) &= \left(\frac{1}{2N} \right) (1 - \mu)^2 \\
 \hat{f} &= \frac{\left(\frac{1}{2N} \right) (1 - \mu)^2}{1 - \left(1 - \frac{1}{2N} \right) (1 - \mu)^2} \\
 &\approx \frac{1 - 2\mu}{2N \left(1 - \left(1 - \frac{1}{2N} \right) (1 - 2\mu) \right)} \\
 &= \frac{1 - 2\mu}{2N \left(1 - 1 + \frac{1}{2N} + 2\mu - \frac{2\mu}{2N} \right)} \\
 &= \frac{1 - 2\mu}{1 + 4N\mu - 2\mu} \\
 &\approx \frac{1}{4N\mu + 1}
 \end{aligned}$$

Since f is the probability that two alleles chosen at random are identical by descent within our population, $1 - f$ is the probability that two alleles chosen at random are *not*

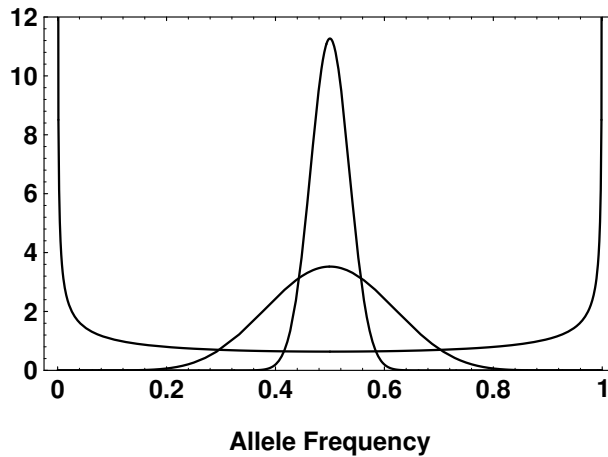
²Notice that we're also playing a little fast and loose with definitions here, since I've just described this in terms of identity by type when what the equation is written in terms of identity by descent. Can you see why it is that I can get away with this?

³Technically what the population reaches is not an equilibrium. It reaches a stationary distribution. At any point in time there is some probability that the population has a particular allele frequency. After long enough the probability distribution stops changing. That's when the population is at its stationary distribution. We often say that it's "reached stationarity." This is an example of a place where the inbreeding analogy breaks down a little.

identical by descent in our population. So $1 - f = 4N\mu / (4N\mu + 1)$ is a reasonable measure of the genetic diversity within the population. Notice that as N increases, the genetic diversity maintained in the population also increases. This shouldn't be too surprising. The rate at which diversity is lost declines as population size increases so larger populations should retain more diversity than small ones.⁴

A two-allele model with recurrent mutation

There's another way of looking at the interaction between drift and mutation. Suppose we have a set of populations with two alleles, A_1 and A_2 . Suppose further that the rate of mutation from A_1 to A_2 is equal to the rate of mutation from A_2 to A_1 .⁵ Call that rate μ . In the absence of mutation a fraction p_0 of the populations would fix on A_1 and the rest would fix on A_2 , where p_0 is the original frequency of A_1 . With recurrent mutation, no population will ever be permanently fixed for one allele or the other. Instead we see the following:



When $4N\mu < 1$ the stationary distribution of allele frequencies is bowl-shaped, i.e, most populations have allele frequencies near 0 or 1. When $4N\mu > 1$, the stationary distribution of

⁴Remember that if we're dealing with a non-ideal population, as we always are, you'll need to substitute N_e for N in this equation and others like it.

⁵We don't have to make this assumption, but relaxing it makes an already fairly complicated scenario even more complicated. If you're really interested, ask me about it.

allele frequencies is hump-shaped, i.e., most populations have allele frequencies near 0.5. In other words if the population is “small,” drift dominates the distribution of allele frequencies and causes populations to become differentiated. If the population is “large,” mutation dominates and keeps the allele frequencies in the different populations similar to one another. That’s what we mean when we say that a population is “large” or “small”. A population is “large” if evolutionary processes other than drift have a predominant influence on the outcome. It’s “small” if drift has a predominant role on the outcome.

A population is large with respect to the drift-mutation process if $4N\mu > 1$, and it is small if $4N\mu < 1$. Notice that calling a population large or small is really just a convenient shorthand. There isn’t much of a difference between the allele frequency distributions when $4N\mu = 0.9$ and when $4N\mu = 1.1$. Notice also that because mutation is typically rare, on the order of 10^{-5} or less per locus per generation for a protein-coding gene and on the order of 10^{-3} or less per locus for a microsatellite, a population must be pretty large ($> 25,000$ or > 250) to be considered large with respect to the drift-migration process. Notice also that whether the population is “large” or “small” will depend on the loci that you’re studying.

Drift and migration

I just pointed out that if populations are isolated from one another they will tend to diverge from one another as a result of genetic drift. Recurrent mutation, which “pushes” all populations towards the same allele frequency, is one way in which that tendency can be opposed. If populations are not isolated, but exchange migrants with one another, then migration will also oppose the tendency for populations to become different from one another. It should be obvious that there will be a tradeoff similar to the one with mutation: the larger the populations, the less the tendency for them to diverge from one another and, therefore, the more migration will tend to make them similar. To explore how drift and migration interact we can use an approach exactly analogous to what we used for mutation.

The model of migration we’ll consider is an extremely oversimplified one. It imagines that every allele brought into a population is different from any of the resident alleles.⁶ It also imagines that all populations receive the fraction of migrants. Because any immigrant allele is different, by assumption, from any resident allele we don’t even have to keep track of how far apart populations are from one another, since populations close by will be no more similar to one another than populations far apart. This is Wright’s island model of

⁶Sounds a lot like the infinite alleles model of mutation, doesn’t it? Just you wait. The parallel gets even more striking.

migration. Given these assumptions, we can write the following:

$$f_{t+1} = \left(\left(\frac{1}{2N} \right) + \left(1 - \frac{1}{2N} \right) f_t \right) (1 - m)^2 \quad . \quad (15.3)$$

That might look fairly familiar. In fact, it's identical to equation (15.2) except that there's an m in (15.3) instead of a μ . m is the migration rate, the fraction of individuals in a population that is composed of immigrants. More precisely, m is the *backward* migration rate. It's the probability that a randomly chosen individual in this generation *came from* a population different from the one in which it is currently found in the preceding generation. Normally we'd think about the *forward* migration rate, i.e., the probability that a randomly chosen individual with *go to* a different population in the next generation, but backwards migration rates turn out to be more convenient to work with in most population genetic models.⁷

It shouldn't surprise you that if equations (15.2) and (15.3) are so similar the equilibrium f under drift and migration is

$$\hat{f} \approx \frac{1}{4Nm + 1}$$

In fact, the two allele analog to the mutation model I presented earlier turns out to be pretty similar, too.

- If $2Nm > 1$, the stationary distribution of allele frequencies is hump-shaped, i.e., the populations tend not to diverge from one another.⁸
- If $2Nm < 1$, the stationary distribution of allele frequencies is bowl-shaped, i.e., the populations tend to diverge from one another.

Now there's a consequence of these relationships that's both surprising and odd. N is the population size. m is the fraction of individuals in the population that are immigrants. So Nm is the *number* of individuals in the population that are new immigrants in any generation. That means that if populations receive more than one new immigrant every other generation, on average, they'll tend not to diverge in allele frequency from one another.⁹ It doesn't make any difference if the populations have a million individuals a piece or ten. One new immigrant every other generation is enough to keep them from diverging.

With a little more reflection, this result is less surprising than it initially seems. After all in populations of a million individuals, drift will be operating very slowly, so it doesn't take

⁷I warned you weeks ago that population geneticists tend to think backwards.

⁸You read that right it's $2Nm$ not $4Nm$ as you might have expected from the mutation model. If you're *really* interested why there's a difference, I can show you. But the explanation isn't simple.

⁹In the sense that the stationary distribution of allele frequencies is hump-shaped.

a large proportion of immigrants to keep populations from diverging.¹⁰ In populations with only ten individuals, drift will be operating much more quickly, so it takes a large proportion of immigrants to keep populations from diverging.¹¹

¹⁰And one immigrant every other generation corresponds to a backwards migration rate of only 5×10^{-7} .

¹¹And one immigrant every other generation corresponds to a backwards migration rate of 5×10^{-2} .

Chapter 16

Selection and genetic drift

There are three basic facts about genetic drift that I really want you to remember, even if you forget everything else I've told you about it:

1. Allele frequencies tend to change from one generation to the next purely as a result of random sampling error. We can specify a probability distribution for the allele frequency in the next generation, but we cannot specify the numerical value exactly.
2. There is no systematic bias to the change in allele frequency, i.e., allele frequencies are as likely to increase from one generation to the next as to decrease.
3. Populations will eventually fix for one of the alleles that is initially present unless mutation or migration introduces new alleles.

Natural selection introduces a systematic bias in allele frequency changes. Alleles favored by natural selection *tend* to increase in frequency. Notice that word “tend.” It’s critical. Because there is a random component to allele frequency change when genetic drift is involved, we can’t say for sure that a selectively favored allele will increase in frequency. In fact, we can say that there’s a chance that a selectively favored allele *won’t* increase in frequency. There’s also a chance that a selectively *disfavored* allele will increase in frequency in spite of natural selection.

Loss of beneficial alleles

We’re going to confine our studies to our usual simple case: one locus, two alleles. We’re also going to consider a very simple form of directional viability selection in which the heterozygous genotype is exactly intermediate in fitness.

$$\begin{array}{ccc} A_1A_1 & A_1A_2 & A_2A_2 \\ 1 + s & 1 + \frac{1}{2}s & 1 \end{array}$$

After solving a reasonably complex partial differential equation, it can be shown that¹ the probability that allele A_1 ² is fixed, given that its current frequency is p is

$$P_1(p) = \frac{1 - e^{-2N_e s p}}{1 - e^{-2N_e s}} \quad (16.1)$$

Now it won't be immediately evident to you, but this equation actually confirms our intuition that even selectively favored alleles may sometimes be lost as a result of genetic drift. How does it do that? Well, it's not too hard to verify that $P_1(p) < 1$.³ The probability that the beneficial allele is fixed is less than one meaning that the probability it is lost is greater than zero, i.e., there's some chance it will be lost.

How big is the chance that a favorable allele will be lost? Well, consider the case of a newly arisen allele with a beneficial effect. If it's newly arisen, there is only one copy by definition. In a diploid population of N individuals that means that the frequency of this allele is $1/2N$. Plugging this into equation (16.1) above we find

$$\begin{aligned} P_1(p) &= \frac{1 - e^{-2N_e s(1/2N)}}{1 - e^{-2N_e s}} \\ &\approx 1 - e^{-N_e s(1/N)} \text{ if } 2N_e s \text{ is "large"} \\ &\approx s \left(\frac{N_e}{N} \right) \text{ if } s \text{ is "small."} \end{aligned}$$

In other words, most beneficial mutations are lost from populations unless they are *very* beneficial. If $s = 0.2$ in an ideal population, for example, a beneficial mutation will be lost about 80% of the time.⁴ Remember that in a strict harem breeding system with a single male $N_e \approx 4$ if the number of females with which the male breeds is large enough. Suppose that there are 99 females in the population. Then $N_e/N = 0.04$ and the probability that this beneficial mutation will be fixed is only 0.8%.

Notice that unlike what we saw with natural selection when we were ignoring genetic drift, the strength of selection⁵ affects the outcome of the interaction. The stronger selection is the more likely it is that the favored allele will be fixed. But it's also the case that the larger the population is, the more likely the favored allele will be fixed.⁶ Size *does* matter.

¹Remember, I told you that "it can be shown that" hides a *lot* of work.

²The beneficial allele.

³Unless $p = 1$.

⁴The exact calculation from equation (16.1) gives 82% for this probability.

⁵i.e., the magnitude of differences in relative viabilities

⁶Because the larger the population, the smaller the effect of drift.

s	N_e	
	4	100
0.001	1×10^{-2}	9×10^{-3}
0.01	1×10^{-2}	3×10^{-3}
0.1	7×10^{-3}	5×10^{-10}

Table 16.1: Fixation probabilities for a deleterious mutation as a function of effective population size and selection coefficient for a newly arisen mutant ($p = 0.01$).

Fixation of detrimental alleles

If drift can lead to the loss of beneficial alleles, it should come as no surprise that it can also lead to fixation of deleterious ones. In fact, we can use the same formula we've been using (equation (16.1)) if we simply remember that for an allele to be deleterious s will be negative. So we end up with

$$P_1(p) = \frac{1 - e^{2N_e s p}}{1 - e^{2N_e s}} \quad . \quad (16.2)$$

One implication of equation (16.2) that should not be surprising by now is that even a deleterious allele can become fixed. Consider our two example populations again, an ideal population of size 100 ($N_e = 100$) and a population with 1 male and 99 females ($N_e = 4$). Remember, the probability of fixation for a newly arisen allele with no effect on fitness is $1/2N = 5 \times 10^{-3}$ (Table 16.1).⁷

Conclusions

I'm not going to try to show you the formulas, but it shouldn't surprise you to learn that heterozygote advantage won't maintain a polymorphism indefinitely in a finite population. At best what it will do is to retard its loss.⁸ There are four properties of the interaction of drift and selection that I think you should take away from this brief discussion:

1. Most mutations, whether beneficial, deleterious, or neutral, are lost from the population in which they occurred.

⁷Because its probability of fixation is equal to its current frequency, i.e., $1/2N$. We'll return to this observation in a few weeks when we talk about the neutral theory of molecular evolution.

⁸In some cases it can actually accelerate its loss, but we won't discuss that unless you are really interested.

2. If selection against a deleterious mutation is weak or N_e is small,⁹ a deleterious mutation is almost as likely to be fixed as neutral mutants. They are “effectively neutral.”
3. If N_e is large, deleterious mutations are much less likely to be fixed than neutral mutations.
4. Even if N_e is large, most favorable mutations are lost.

⁹As with mutation and migration, what counts as large or small is determined by the product of N_e and s . If it's bigger than one the population is regarded as large, because selective forces predominate. If it's smaller than one, it's regarded as small, because drift predominates.

Chapter 17

The Coalescent

I've mentioned several times that population geneticists often look at the world backwards. Sometimes when they do, the result is very useful. Consider genetic drift, for example. So far we've been trying to predict what will happen in a population given a particular effective population size. But when we collect data we are often more interested in understanding the processes that produced any pattern we find than in predicting what will happen in the future. So let's take a backward look at drift and see what we find.

Reconstructing the genealogy of a sample of alleles

Specifically, let's keep track of the genealogy of alleles. In a finite population, two randomly chosen alleles will be identical by descent with respect to the immediately preceding generation with probability $1/2N_e$. That means that there's a chance that two alleles in generation t are copies of the same allele in generation $t - 1$. If the population size is constant, meaning that the number of alleles in the population is remaining constant, then there's also a chance that some alleles present in generation $t - 1$ will not have descendants in generation t . Looking backward, then, the number of alleles in generation $t - 1$ that have descendants in generation t is always less than or equal to the number of alleles in generation t . That means if we trace the ancestry of alleles in a sample back far enough, all of them will be descended from a single common ancestor. Figure 17.1 provides a simple schematic illustrating how this might happen.

Now take a look at Figure 17.1. Time runs from the top of the figure to the bottom, i.e., the current generation is represented by the circles in the bottom row of the figure. Each circle represents an allele. The eighteen alleles in our current sample are descended from only four alleles that were present in the populations ten generations ago. The other

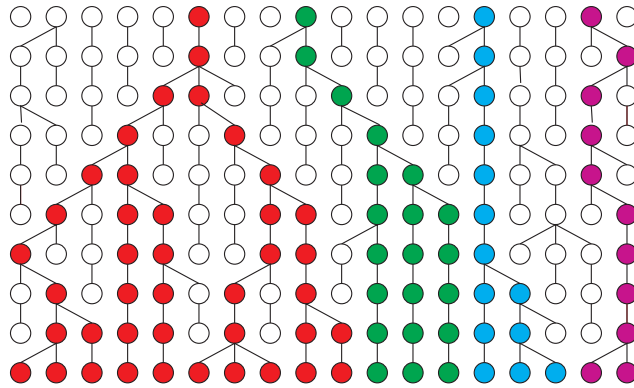


Figure 17.1: A schematic depiction of one possible realization of the coalescent process in a population with 18 haploid gametes. There are four coalescent events in the generation immediately preceding the last one illustrated, one involving three alleles.

fourteen alleles present in the population ten generations ago left no descendants. How far back in time we'd have to go before all alleles are descended from a single common ancestor depends on the effective size of the population, and how frequently two (or more) alleles are descended from the same allele in the preceding generation depends on the effective size of the population, too. But in any finite population the pattern will look something like the one I've illustrated here.

Mathematics of the coalescent: two alleles

J. F. C. Kingman developed a convenient and powerful way to describe how the time to common ancestry is related to effective population size [52, 53]. The process he describes is referred to as the *coalescent*, because it is based on describing the probability of *coalescent events*, i.e., those points in the genealogy of a sample of alleles where two alleles are descended from the same allele in the immediately preceding generation.¹ Let's consider a simple case, one that we've already seen, first, i.e., two alleles drawn at random from a single population.

The probability that two alleles drawn at random from a population are copies of the same allele in the preceding generation is also the probability that two alleles drawn at random

¹An important assumption of the coalescent is that populations are large enough that we can ignore the possibility that there is more than one coalescent event in a single generation. We also only allow coalescence between a pair of alleles, not three or more. In both ways the mathematical model of the process differs from the diagram in Figure 17.1.

from that population are identical by descent with respect to the immediately preceding generation. We know what that probability is,² namely

$$\frac{1}{2N_e^{(f)}} .$$

I'll just use N_e from here on out, but keep in mind that the appropriate population size for use with the coalescent is the inbreeding effective size. Of course, this means that the probability that two alleles drawn at random from a population are *not* copies of the same allele in the preceding generation is

$$1 - \frac{1}{2N_e} .$$

We'd like to calculate the probability that a coalescent event happened at a particular time t , in order to figure out how far back in the ancestry of these two alleles we have to go before they have a common ancestor. How do we do that?

Well, in order for a coalescent event to occur at time t , the two alleles must have *not* have coalesced in the generations preceding that.³ The probability that they did not coalesce in the first $t - 1$ generations is simply

$$\left(1 - \frac{1}{2N_e}\right)^{t-1} .$$

Then after having remained distinct for $t - 1$ generations, they have to coalesce in generation t , which they do with probability $1/2N_e$. So the probability that two alleles chosen at random coalesced t generations ago is

$$P(T = t) = \left(1 - \frac{1}{2N_e}\right)^{t-1} \left(\frac{1}{2N_e}\right) . \tag{17.1}$$

It's not too hard to show, once we know the probability distribution in equation (17.1), that the average time to coalescence for two randomly chosen alleles is $2N_e$.⁴

²Though you may not remember it.

³Remember that we're counting generations backward in time, so when I say that a coalescent event occurred at time t I mean that it occurred t generations ago.

⁴If you've had a little bit of probability theory, you'll notice that equation 17.1 shows that the coalescence time is a geometric random variable.

Mathematics of the coalescent: multiple alleles

It's quite easy to extend this approach to multiple alleles.⁵ We're interested in seeing how far back in time we have to go before all alleles are descended from a single common ancestor. We'll assume that we have m alleles in our sample. The first thing we have to calculate is the probability that any two of the alleles in our sample are identical by descent from the immediately preceding generation. To make the calculation easier, we assume that the effective size of the population is large enough that the probability of two coalescent events in a single generation is vanishingly small. We already know that the probability of a coalescence in the immediately preceding generation for two randomly chosen alleles is $1/2N_e$. But there are $m(m-1)/2$ different pairs of alleles in our sample. So the probability that one pair of these alleles is involved in a coalescent event in the immediately preceding generation is

$$\left(\frac{1}{2N_e}\right) \left(\frac{m(m-1)}{2}\right) .$$

From this it follows⁶ that the probability that the first coalescent event involving this sample of alleles occurred t generations ago is

$$P(T = t) = \left(1 - \left(\frac{1}{2N_e}\right) \left(\frac{m(m-1)}{2}\right)\right)^{t-1} \left(\frac{1}{2N_e}\right) \left(\frac{m(m-1)}{2}\right) . \quad (17.2)$$

So the mean time back to the first coalescent event is

$$\frac{2N_e}{m(m-1)/2} = \frac{4N_e}{m(m-1)} \text{ generations} .$$

But this is, of course, only the first coalescent event. We were interested in how long we have to wait until *all* alleles are descended from a single common ancestor. Now is where Kingman's sneaky trick comes in. After the first coalescent event, we now have $m-1$ alleles in our sample, instead of m . So the whole process starts over again with $m-1$ alleles instead of m . Since the time to the first coalescence depends only on the number of alleles in the sample and not on how long the first coalescence event took, we can calculate the average time until all coalescences have happened as

$$\bar{t} = \sum_{k=2}^m \bar{t}_k$$

⁵Okay, okay. What I should really have said is "It's not *too* hard to extend this approach to multiple alleles."

⁶Using logic just like what we used in the two allele case.

$$\begin{aligned}
&= \sum_{k=2}^m \frac{4N_e}{k(k-1)} \\
&= 4N_e \left(1 - \frac{1}{m}\right) \\
&\approx 4N_e
\end{aligned}$$

An example: Mitochondrial Eve

Cann et al. [10] sampled mitochondrial DNA from 147 humans of diverse racial and geographic origins. Based on the amount of sequence divergence they found among genomes in their sample and independent estimates of the rate of sequence evolution, they inferred that the mitochondria in their sample had their most recent common ancestor about 200,000 years ago. Because all of the most ancient lineages in their sample were from individuals of African ancestry, they also suggested that mitochondrial Eve lived in Africa. They used these arguments as evidence for the “Out of Africa” hypothesis for modern human origins, i.e., the hypothesis that anatomically modern humans arose in Africa about 200,000 years ago and displaced other members of the genus *Homo* in Europe and Asia as they spread. What does the coalescent tell us about their conclusion?

Well, we expect all mitochondrial genomes in the sample to have had a common ancestor about $2N_e$ generations ago. Why $2N_e$ rather than $4N_e$? Because mitochondrial genomes are haploid. Furthermore, since we all got our mitochondria from our mothers, N_e in this case refers to the effective number of *females*.

Given that a human generation is about 20 years, a coalescence time of 200,000 years implies that the mitochondrial genomes in the Cann et al. sample have their most recent common ancestor about 10,000 generations ago. If the effective number of females in the human populations is 5000, that’s exactly what we’d expect. While 5000 may sound awfully small, given that there are about 3 billion women on the planet now, remember that until the recent historical past (no more than 500 generations ago) the human population was small and humans lived in small hunter-gatherer groups, so an effective number of females of 5000 and a total effective size of 10,000 may not be unreasonable. If that’s true, then the geographical location of mitochondrial Eve need not tell us anything about the origin of modern human populations, because there had to be a coalescence somewhere. There’s no guarantee, from this evidence alone, that the Y-chromosome Adam would have lived in Africa, too. Having said that, my limited reading of the literature suggests that other data are consistent with the “Out of Africa” scenario. Y-chromosome polymorphisms, for example, are also consistent with the “Out of Africa” hypothesis [93]. Interestingly, dating of those polymorphisms suggests that Y-chromosome Adam left Africa 35,000 – 89,000 years ago.

The coalescent and F -statistics

Suppose we have a sample of alleles from a structured population. For alleles chosen randomly within populations let the average time to coalescence be \bar{t}_0 . For alleles chosen randomly from different populations let the average time to coalescence be \bar{t}_1 . If there are k populations in our sample, the average time to coalescence for two alleles drawn at random without respect to population is⁷

$$\bar{t} = \frac{k(k-1)\bar{t}_1 + k\bar{t}_0}{k} .$$

Slatkin [82] pointed out that F_{st} bears a simple relationship to average coalescence times within and among populations. Given these definitions of \bar{t} and \bar{t}_0 ,

$$F_{st} = \frac{\bar{t} - \bar{t}_0}{\bar{t}} .$$

So F_{st} measures the proportional increase in time to coalescence that is due to populations being separate from one another. One way to think about that relationship is this: the longer it has been, on average, since alleles in different populations diverged from a common ancestor, the greater the chances that they have become different. An implication of this relationship is that F -statistics, by themselves, can't tell us much of anything about patterns of migration among populations.

A given pattern of among-population relationships might reflect a migration-drift equilibrium, a sequence of population splits followed by genetic isolation, or any combination of the two. If we are willing to assume that populations in our sample have been exchanging genes long enough to reach stationarity in the drift-migration process, then F_{st} may tell us something about migration. If we are willing to assume that there's been no gene exchange among our populations, we can infer something about how recently they've diverged from one another. But unless we're willing to make one of those assumptions, we can't make any further progress.

⁷If you don't see why, don't worry about it. You can ask if you really care. We only care about \bar{t} for what follows anyway.

Part IV

Quantitative genetics

Chapter 18

Introduction to quantitative genetics

So far in this course we have dealt almost entirely with the evolution of characters that are controlled by simple Mendelian inheritance at a single locus. There are notes on the course website about gametic disequilibrium and how allele frequencies change at two loci simultaneously, but we didn't discuss them. In every other example we've considered we've imagined that we could understand something about evolution by examining the evolution of a single gene. That's the domain of classical population genetics.

For the next few weeks we're going to be exploring a field that's actually older than classical population genetics, although the approach we'll be taking to it involves the use of population genetic machinery. If you know a little about the history of evolutionary biology, you may know that after the rediscovery of Mendel's work in 1900 there was a heated debate between the "biometricians" (e.g., Galton and Pearson) and the "Mendelians" (e.g., de Vries, Correns, Bateson, and Morgan).

Biometricians asserted that the really important variation in evolution didn't follow Mendelian rules. Height, weight, skin color, and similar traits seemed to

- vary continuously,
- show blending inheritance, and
- show variable responses to the environment.

Since variation in such *quantitative traits* seemed to be more obviously related to organismal adaptation than the "trivial" traits that Mendelians studied, it seemed obvious to the biometricians that Mendelian geneticists were studying a phenomenon that wasn't particularly interesting.

Mendelians dismissed the biometricians, at least in part, because they seemed not to recognize the distinction between genotype and phenotype. It seemed to at least some of them that traits whose expression was influenced by the environment were, by definition, not inherited. Moreover, the evidence that Mendelian principles accounted for the inheritance of many discrete traits was incontrovertible.

Woltereck's [101] experiments on *Daphnia* helped to show that traits whose expression is environmentally influenced may also be inherited. He introduced the idea of a *norm of reaction* to describe the observation that the same genotype may produce different phenotypes in different environments. When you fertilize a plant, for example, it will grow larger and more robust than when you don't. The phenotype an organism expresses is, therefore, a product of *both* its genotype and its environment.

Nilsson-Ehle's [70] experiments on inheritance of kernel color in wheat showed how continuous variation and Mendelian inheritance could be reconciled. He demonstrated that what appeared to be continuous variation in color from red to white with blending inheritance could be understood as the result of three separate genes influencing kernel color that were inherited separately from one another. It was the first example of what's come to be known as *polygenic inheritance*. Fisher [25], in a paper that grew out of his undergraduate Honors thesis at Cambridge University, set forth the mathematical theory that describes how it all works. That's the theory of *quantitative genetics*, and it's what we're going to spend the next three weeks discussing.

An overview of where we're headed

Woltereck's ideas force us to realize that when we see a phenotypic difference between two individuals in a population there are three possible explanations for that difference:

1. The individuals have different genotypes.
2. The individuals developed in different environments.
3. The individuals have different genotypes *and* they developed in different environments.

This leads us naturally to think that phenotypic variation consists of two separable components, namely genotypic and environmental components.¹ Putting that into an equation

$$\text{Var}(P) = \text{Var}(G) + \text{Var}(E) \quad ,$$

¹We'll soon see that separating genotypic and environmental components is far from trivial.

where $\text{Var}(P)$ is the *phenotypic variance*, $\text{Var}(G)$ is the *genetic variance*, and $\text{Var}(E)$ is the environmental variance.² As we'll see in just a moment, we can also partition the genetic variance into components, the *additive genetic variance*, $\text{Var}(A)$, and the *dominance variance*, $\text{Var}(D)$.

There's a surprisingly subtle and important insight buried in that very simple equation: Because the expression of a quantitative trait is a result both of genes involved in that trait's expression and the environment in which it is expressed, it doesn't make sense to say of a particular individual's phenotype that genes are more important than environment in determining it. You wouldn't have a phenotype without both. What we might be able to say is that when we look at a particular population of organisms some fraction of the phenotypic differences among them is due to differences in the genes they carry and that some fraction is due to differences in the environment they have experienced.³

One important implication of this insight is that much of the “nature vs. nurture” debate concerning human intelligence or human personality characteristics. The intelligence and personality that you have is a product of both the genes you happen to carry and the environment that you experienced. Any differences between you and the person next to you probably reflect both differences in genes *and* differences in environment. Moreover, just because you have a genetic pre-disposition for a particular condition doesn't mean you're stuck with it.

Take phenylketonuria, for example. It's a condition in which individuals are homozygous for a deficiency that prevents them from metabolizing phenylalanine (<http://www.ncbi.nlm.nih.gov/pubmedhealth/PMH0002150/>). If individuals with phenylketonuria eat a normal diet, severe mental can result by the time an infant is one year old. But if they eat a diet that is very low in phenylalanine, their development is completely normal.

It's often useful to talk about how much of the phenotypic variance is a result of additive genetic variance or of genetic variance.

$$h_n^2 = \frac{\text{Var}(A)}{\text{Var}(P)}$$

is what's known as the *narrow-sense heritability*. It's the proportion of phenotypic variance that's attributable to differences among individuals in their additive genotype,⁴ much as F_{st} can be thought of as the proportion of genotypic diversity that attributable to differences

²Strictly speaking we should also include a term for the interaction between genotype and environment, but we'll ignore that for the time being.

³When I put it this way, I hope it's obvious that I'm neglecting genotype-environment interactions, and that I'm oversimplifying quite a bit.

⁴Don't worry about what I mean by *additive genotype*—yet. We'll get to it soon enough.

among populations. Similarly,

$$h_b^2 = \frac{\text{Var}(G)}{\text{Var}(P)}$$

is the *broad-sense heritability*. It's the proportion of phenotypic variance that's attributable to differences among individuals in their genotype. It is *not*, repeat *NOT*, a measure of how important genes are in determining phenotype. Every individual's phenotype is determined both by its genes and by its phenotype. It measures how much of the *difference* among individuals is attributable to differences in their genes.⁵ Why bother to make the distinction between narrow- and broad-sense heritability? Because, as we'll see, it's only the additive genetic variance that responds to natural selection.⁶ In fact,

$$R = h_n^2 S \quad ,$$

where R is the *response to selection* and S is the *selective differential*.

As you'll see in the coming weeks, there's a lot of stuff hidden behind these simple equations, including a lot of assumptions. But quantitative genetics is very useful. Its principles have been widely applied in plant and animal breeding for almost a century, and they have been increasingly applied in evolutionary investigations in the last thirty years. Nonetheless, it's useful to remember that quantitative genetics is a lot like a bikini. What it reveals is interesting, but what it conceals is crucial.

Partitioning the phenotypic variance

Before we worry about how to estimate any of those variance components I just mentioned, we first have to understand what they are. So let's start with some definitions (Table 18.1).⁷

You should notice something rather strange about Table 18.1 when you look at it. I motivated the entire discussion of quantitative genetics by talking about the need to deal with variation at many loci, and what I've presented involves only two alleles at a single locus. I do this for two reasons:

1. It's not too difficult to do the algebra with multiple alleles at one locus instead of only two, but it gets messy, doesn't add any insight, and I'd rather avoid the mess.

⁵As we'll see later it can do this only for the range of environments in which it was measured.

⁶Or at least only the additive genetic variance responds to natural selection when zygotes are found in Hardy-Weinberg proportions.

⁷Warning! There's a lot of algebra between here and the end. It's unavoidable. You can't possibly understand what additive genetic variance is without it. I'll try to focus on principles, but a lot of the algebra that follows *is* necessary. Sorry about that.

Genotype	A_1A_1	A_1A_2	A_2A_2
Frequency	p^2	$2pq$	q^2
Genotypic value	x_{11}	x_{12}	x_{22}
Additive genotypic value	$2\alpha_1$	$\alpha_1 + \alpha_2$	$2\alpha_2$

Table 18.1: Fundamental parameter definitions for quantitative genetics with one locus and two alleles.

2. Doing the algebra with multiple loci involves a *lot* of assumptions, which I'll mention when we get to applications, and the algebra is even worse than with multiple alleles.

Fortunately, the basic principles extend with little modification to multiple loci, so we can see all of the underlying logic by focusing on one locus with two alleles where we have a chance of understanding what the different variance components mean.

Two terms in Table 18.1 will almost certainly be unfamiliar to you: *genotypic value* and *additive genotypic value*. Of the two, *genotypic value* is the easiest to understand (Figure 18.1). It simply refers to the average phenotype associated with a given genotype.⁸ The *additive genotypic value* refers to the average phenotype associated with a given genotype, as would be inferred from the *additive effect* of the alleles of which it is composed. That didn't help much, did it? That's because I now need to tell you what we mean by the *additive effect* of an allele.⁹

The additive effect of an allele

In constructing Table 18.1 I used the quantities α_1 and α_2 , but I didn't tell you where they came from. Obviously, the idea should be to pick values of α_1 and α_2 that give additive genotypic values that are reasonably close to the genotypic values. A good way to do that is to minimize the squared deviation between the two, weighted by the frequency of the genotypes. So our first big assumption is that genotypes are in Hardy-Weinberg proportions.¹⁰

The objective is to find values for α_1 and α_2 that minimize:

$$a = p^2[x_{11} - 2\alpha_1]^2 + 2pq[x_{12} - (\alpha_1 + \alpha_2)]^2 + q^2[x_{22} - 2\alpha_2]^2 \quad .$$

⁸Remember. We're now considering traits in which the environment influences the phenotypic expression, so the same genotype can produce different phenotypes, depending on the environment in which it develops.

⁹Hold on. Things get even more interesting from here.

¹⁰As you should have noticed in Table 18.1.

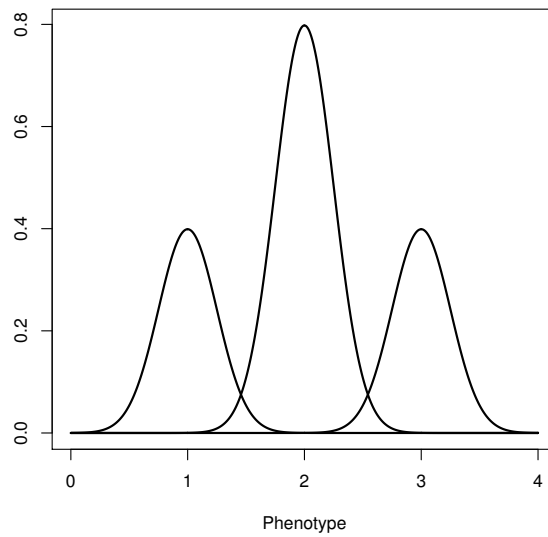


Figure 18.1: The phenotype distribution in a population in which the three genotypes at a single locus with two alleles occur in Hardy-Weinberg proportions and the alleles occur in equal frequency. The A_1A_1 genotype has a mean trait value of 1, the A_1A_2 genotype has a mean trait value of 2, and the A_2A_2 genotype has a mean trait value of 3, but each genotype can produce a range of phenotypes with the standard deviation of the distribution being 0.25 in each case.

To do this we take the partial derivative of a with respect to both α_1 and α_2 , set the resulting pair of equations equal to zero, and solve for α_1 and α_2 .¹¹

$$\begin{aligned}\frac{\partial a}{\partial \alpha_1} &= p^2\{2[x_{11} - 2\alpha_1][-2]\} + 2pq\{2[x_{12} - (\alpha_1 + \alpha_2)][-1]\} \\ &= -4p^2[x_{11} - 2\alpha_1] - 4pq[x_{12} - (\alpha_1 + \alpha_2)] \\ \frac{\partial a}{\partial \alpha_2} &= q^2\{2[x_{22} - 2\alpha_2][-2]\} + 2pq\{2[x_{12} - (\alpha_1 + \alpha_2)][-1]\} \\ &= -4q^2[x_{22} - 2\alpha_2] - 4pq[x_{12} - (\alpha_1 + \alpha_2)]\end{aligned}$$

Thus, $\frac{\partial a}{\partial \alpha_1} = \frac{\partial a}{\partial \alpha_2} = 0$ if and only if

$$\begin{aligned}p^2(x_{11} - 2\alpha_1) + pq(x_{12} - \alpha_1 - \alpha_2) &= 0 \\ q^2(x_{22} - 2\alpha_2) + pq(x_{12} - \alpha_1 - \alpha_2) &= 0\end{aligned}\tag{18.1}$$

Adding the equations in (18.1) we obtain (after a little bit of rearrangement)

$$[p^2x_{11} + 2pqx_{12} + q^2x_{22}] - [p^2(2\alpha_1) + 2pq(\alpha_1 + \alpha_2) + q^2(2\alpha_2)] = 0 \quad .\tag{18.2}$$

Now the first term in square brackets is just the mean phenotype in the population, \bar{x} . Thus, we can rewrite equation (18.2) as:

$$\begin{aligned}\bar{x} &= 2p^2\alpha_1 + 2pq(\alpha_1 + \alpha_2) + 2q^2\alpha_2 \\ &= 2p\alpha_1(p + q) + 2q\alpha_2(p + q) \\ &= 2(p\alpha_1 + q\alpha_2) \quad .\end{aligned}\tag{18.3}$$

Now divide the first equation in (18.1) by p and the second by q .

$$p(x_{11} - 2\alpha_1) + q(x_{12} - \alpha_1 - \alpha_2) = 0\tag{18.4}$$

$$q(x_{22} - 2\alpha_2) + p(x_{12} - \alpha_1 - \alpha_2) = 0 \quad .\tag{18.5}$$

Thus,

$$\begin{aligned}px_{11} + qx_{12} &= 2p\alpha_1 + q\alpha_1 + q\alpha_2 \\ &= \alpha_1(p + q) + p\alpha_1 + q\alpha_2 \\ &= \alpha_1 + p\alpha_1 + q\alpha_2 \\ &= \alpha_1 + \bar{x}/2 \\ \alpha_1 &= px_{11} + qx_{12} - \bar{x}/2 \quad .\end{aligned}$$

¹¹We won't bother with proving that the resulting estimates produce the minimum possible value of a . Just take my word for it. Or if you don't believe me and know a little calculus, take the second partials of a and evaluate it with the values of α_1 and α_2 substituted in. You'll find that the resulting matrix of partial derivatives, the Hessian matrix, is positive definite, meaning that we've found values that minimize the value of a .

Similarly,

$$\begin{aligned}
 px_{12} + qx_{22} &= 2q\alpha_2 + p\alpha_1 + p\alpha_2 \\
 &= \alpha_2(p + q) + p\alpha_1 + q\alpha_2 \\
 &= \alpha_2 + p\alpha_1 + q\alpha_2 \\
 &= \alpha_2 + \bar{x}/2 \\
 \alpha_2 &= px_{12} + qx_{22} - \bar{x}/2 \quad .
 \end{aligned}$$

α_1 is the additive effect of allele A_1 , and α_2 is the additive effect of allele A_2 . If we use these expressions, the additive genotypic values are as close to the genotypic values as possible, given the particular allele frequencies in the population.¹²

Components of the genetic variance

Let's assume for the moment that we can actually measure the genotypic values. Later, we'll relax that assumption and see how to use the resemblance among relatives to estimate the genetic components of variance. But it's easiest to see where they come from if we assume that the genotypic value of each genotype is known. If it is then, writing V_g for $\text{Var}(G)$

$$\begin{aligned}
 V_g &= p^2[x_{11} - \bar{x}]^2 + 2pq[x_{12} - \bar{x}]^2 + q^2[x_{22} - \bar{x}]^2 & (18.6) \\
 &= p^2[x_{11} - 2\alpha_1 + 2\alpha_1 - \bar{x}]^2 + 2pq[x_{12} - (\alpha_1 + \alpha_2) + (\alpha_1 + \alpha_2) - \bar{x}]^2 \\
 &\quad + q^2[x_{22} - 2\alpha_2 + 2\alpha_2 - \bar{x}]^2 \\
 &= p^2[x_{11} - 2\alpha_1]^2 + 2pq[x_{12} - (\alpha_1 + \alpha_2)]^2 + q^2[x_{22} - 2\alpha_2]^2 \\
 &\quad + p^2[2\alpha_1 - \bar{x}]^2 + 2pq[(\alpha_1 + \alpha_2) - \bar{x}]^2 + q^2[2\alpha_2 - \bar{x}]^2 \\
 &\quad + p^2[2(x_{11} - 2\alpha_1)(2\alpha_1 - \bar{x})] + 2pq[2(x_{12} - \{\alpha_1 + \alpha_2\})(\{\alpha_1 + \alpha_2\} - \bar{x})] \\
 &\quad + q^2[2(x_{22} - 2\alpha_2)(2\alpha_2 - \bar{x})] \quad . & (18.7)
 \end{aligned}$$

There are two terms in (18.7) that have a biological (or at least a quantitative genetic) interpretation. The term on the first line is the average squared deviation between the genotypic value and the additive genotypic value. It will be zero only if the effects of the alleles can be decomposed into strictly additive components, i.e., only if the phenotype of the heterozygote is exactly intermediate between the phenotype of the two homozygotes. Thus, it is a measure of how much variation is due to non-additivity (dominance) of allelic effects.

¹²If you've been paying close attention and you have a good memory, the expressions for α_1 and α_2 may look vaguely familiar. They look a lot like the expressions for marginal fitnesses we encountered when studying viability selection.

In short, the *dominance genetic variance*, V_d , is

$$V_d = p^2[x_{11} - 2\alpha_1]^2 + 2pq[x_{12} - (\alpha_1 + \alpha_2)]^2 + q^2[x_{22} - 2\alpha_2]^2 \quad . \quad (18.8)$$

Similarly, the term on the second line of (18.7) is the average squared deviation between the additive genotypic value and the mean genotypic value in the population. Thus, it is a measure of how much variation is due to differences between genotypes in their additive genotype. In short, the *additive genetic variance*, V_a , is

$$V_a = p^2[2\alpha_1 - \bar{x}]^2 + 2pq[(\alpha_1 + \alpha_2) - \bar{x}]^2 + q^2[2\alpha_2 - \bar{x}]^2 \quad . \quad (18.9)$$

What about the terms on the third and fourth lines of the last equation in 18.7? Well, they can be rearranged as follows:

$$\begin{aligned} & p^2[2(x_{11} - 2\alpha_1)(2\alpha_1 - \bar{x})] + 2pq[2(x_{12} - \{\alpha_1 + \alpha_2\})(\{\alpha_1 + \alpha_2\} - \bar{x})] \\ & \quad + q^2[2(x_{22} - 2\alpha_2)(2\alpha_2 - \bar{x})] \\ & = 2p^2(x_{11} - 2\alpha_1)(2\alpha_1 - \bar{x}) + 4pq[x_{12} - (\alpha_1 + \alpha_2)][(\alpha_1 + \alpha_2) - \bar{x}] \\ & \quad + 2q^2(x_{22} - 2\alpha_2)(2\alpha_2 - \bar{x}) \\ & = 4p^2(x_{11} - 2\alpha_1)[\alpha_1 - (p\alpha_1 + q\alpha_2)] \\ & \quad + 4pq[x_{12} - (\alpha_1 + \alpha_2)][(\alpha_1 + \alpha_2) - 2(p\alpha_1 + q\alpha_2)] \\ & \quad + 4q^2(x_{22} - 2\alpha_2)[\alpha_2 - (p\alpha_1 + q\alpha_2)] \\ & = 4p[\alpha_1 - (p\alpha_1 + q\alpha_2)][p(x_{11} - 2\alpha_1) + q(x_{12} - \{\alpha_1 + \alpha_2\})] \\ & \quad + 4q[\alpha_2 - (p\alpha_1 + q\alpha_2)][p(x_{11} - 2\alpha_1)p + q(x_{12} - \{\alpha_1 + \alpha_2\})] \\ & = 0 \end{aligned}$$

Where we have used the identities $\bar{x} = 2(p\alpha_1 + q\alpha_2)$ [see equation (18.3)] and

$$\begin{aligned} p(x_{11} - 2\alpha_1) + q(x_{12} - \alpha_1 - \alpha_2) & = 0 \\ q(x_{22} - 2\alpha_2) + p(x_{12} - \alpha_1 - \alpha_2) & = 0 \end{aligned}$$

[see equations (18.4) and (18.5)]. In short, we have now shown that the total genotypic variance in the population, V_g , can be subdivided into two components — the additive genetic variance, V_a , and the dominance genetic variance, V_d . Specifically,

$$V_g = V_a + V_d \quad ,$$

where V_g is given by the first line of (18.6), V_a by (18.9), and V_d by (18.8).

Genotype	A_1A_1	A_1A_2	A_2A_2
Genotypic value	0	1	2

Table 18.2: A set of perfectly additive genotypic values. Note that the genotypic value of the heterozygote is exactly halfway between the genotypic values of the two homozygotes.

An alternative expression for V_a

There's another way to write the expression for V_a when there are only two alleles at a locus. I show it here because it comes in handy some times.

$$\begin{aligned}
V_a &= p^2(2\alpha_1)^2 + 2pq(\alpha_1 + \alpha_2)^2 + q^2(2\alpha_2)^2 - 4(p\alpha_1 + q\alpha_2)^2 \\
&= 4p^2\alpha_1^2 + 2pq(\alpha_1 + \alpha_2)^2 + 4q^2\alpha_2^2 - 4(p^2\alpha_1^2 + 2pq\alpha_1\alpha_2 + q^2\alpha_2^2) \\
&= 2pq[(\alpha_1 + \alpha_2)^2 - 4\alpha_1\alpha_2] \\
&= 2pq[(\alpha_1^2 + 2\alpha_1\alpha_2 + \alpha_2^2) - 4\alpha_1\alpha_2] \\
&= 2pq[\alpha_1^2 - 2\alpha_1\alpha_2 + \alpha_2^2] \\
&= 2pq[\alpha_1 - \alpha_2]^2 \\
&= 2pq\alpha^2
\end{aligned}$$

An example: the genetic variance with known genotypes

We've been through a lot of algebra by now. Let's run through a couple of numerical examples to see how it all works. For the first one, we'll use the set of genotypic values in Table 18.2

For $p = 0.4$

$$\begin{aligned}
\bar{x} &= (0.4)^2(0) + 2(0.4)(0.6)(1) + (0.6)^2(2) \\
&= 1.20
\end{aligned}$$

$$\begin{aligned}
\alpha_1 &= (0.4)(0) + (0.6)(1) - (1.20)/2 \\
&= 0.0
\end{aligned}$$

$$\begin{aligned}
\alpha_2 &= (0.4)(1) + (0.6)(2) - (1.20)/2 \\
&= 1.0
\end{aligned}$$

Genotype	A_1A_1	A_1A_2	A_2A_2
Genotypic value	0	0.8	2

Table 18.3: A set of non-additive genotypic values. Note that the genotypic value of the heterozygote is closer to the genotypic value of A_1A_1 than it is to the genotypic value of A_2A_2 .

$$\begin{aligned}
V_g &= (0.4)^2(0 - 1.20)^2 + 2(0.4)(0.6)(1 - 1.20)^2 + (0.6)^2(2 - 1.20)^2 \\
&= 0.48 \\
V_a &= (0.4)^2[2(0.0) - 1.20]^2 + 2(0.4)(0.6)[(0.0 + 1.0) - 1.20]^2 + (0.6)^2[2(1.0) - 1.20]^2 \\
&= 0.48 \\
V_d &= (0.4)^2[0 - 2(0.0)]^2 + 2(0.4)(0.6)[1 - (0.0 + 1.0)]^2 + (0.6)^2[2 - 2(1.0)]^2 \\
&= 0.00 \quad .
\end{aligned}$$

For $p = 0.2$, $\bar{x} = 1.60$, $V_g = V_a = 0.32$, $V_d = 0.00$. You should verify for yourself that $\alpha_1 = 0$ and $\alpha_2 = 1$ for $p = 0.2$. If you are ambitious, you could try to prove that $\alpha_1 = 0$ and $\alpha_2 = 1$ for *any* allele frequency.

For the second example we'll use the set of genotypic values in Table 18.3.

For $p = 0.4$

$$\begin{aligned}
\bar{x} &= (0.4)^2(0) + 2(0.4)(0.6)(0.8) + (0.6)^2(2) \\
&= 1.104 \\
\alpha_1 &= (0.4)(0) + (0.6)(0.8) - (1.104)/2 \\
&= -0.072 \\
\alpha_2 &= (0.4)(0.8) + (0.6)(2) - (1.104)/2 \\
&= 0.968 \\
V_g &= (0.4)^2(0 - 1.104)^2 + 2(0.4)(0.6)(0.8 - 1.104)^2 + (0.6)^2(2 - 1.104)^2 \\
&= 0.5284 \\
V_a &= (0.4)^2[2(-0.072) - 1.104]^2 + 2(0.4)(0.6)[(-0.072 + 0.968) - 1.104]^2 \\
&\quad + (0.6)^2[2(0.968) - 1.104]^2 \\
&= 0.5192 \\
V_d &= (0.4)^2[0 - 2(-0.072)]^2 + 2(0.4)(0.6)[0.8 - (-0.072 + 0.968)]^2
\end{aligned}$$

$$\begin{aligned} & +(0.6)^2[2 - 2(0.968)]^2 \\ = & 0.0092 \quad . \end{aligned}$$

To test your understanding, it would probably be useful to calculate \bar{x} , α_1 , α_2 , V_g , V_a , and V_d for one or two other allele frequencies, say $p = 0.2$ and $p = 0.8$. Is it still true that α_1 and α_2 are independent of allele frequencies? If you are *really* ambitious you could try to prove that α_1 and α_2 are independent of allele frequencies if and only if $x_{12} = (x_{11} + x_{12})/2$, i.e., when heterozygotes are exactly intermediate.

Chapter 19

Resemblance among relatives

Just as individuals may differ from one another in phenotype because they have different genotypes, because they developed in different environments, or both, relatives may resemble one another more than they resemble other members of the population because they have similar genotypes, because they developed in similar environments, or both. In an experimental situation, we may be able to randomize individuals across environments. Under those circumstances any tendency for relatives to resemble one another more than non-relatives must be due to similarities in their genotypes.

Using this insight, we can develop a statistical technique that allows us to determine how much of the variance among individuals in phenotype is a result of genetic variance and how much is due to environmental variance. *Remember*, we can only ask about how much of the variability is due to genetic differences, and we can only do so *in a particular environment* and *with a particular set of genotypes*, and we can only do it when we *randomize genotypes across environments*.

An outline of the approach

The basic approach to the analysis is either to use a linear regression of offspring phenotype on parental phenotype, which as we'll see estimates h_n^2 , or to use a nested analysis of variance. One of the most complete designs is a full-sib, half-sib design in which each male sires offspring from several dams but each dam mates with only one sire.

The offspring of a single dam are full-sibs (they are nested within dams). Differences among the offspring of dams indicates that there are differences in maternal “genotype” in the trait being measured.¹

¹Assuming that we've randomized siblings across environments. If we haven't, siblings may resemble one

Maternal genotype	Frequency	Offspring genotype		
		A_1A_1	A_1A_2	A_2A_2
A_1A_1	p^2	p	q	0
A_1A_2	$2pq$	$\frac{p}{2}$	$\frac{1}{2}$	$\frac{q}{2}$
A_2A_2	q^2	0	p	q

Table 19.1: Half-sib family structure in a population with genotypes in Hardy-Weinberg proportions.

The offspring of different dams mated to a single sire are half-sibs. Differences among the offspring of sires indicates that there are differences in paternal “genotype” in the trait being measured.²

As we’ll see, this design has the advantage that it allows both additive and dominance components of the genetic variance to be estimated. It has the additional advantage that we don’t have to assume that the distribution of environments in the offspring generation is the same as it was in the parental generation.

The gory details

OK, so I’ve given you the basic idea. Where does it come from, and how does it work? Funny you should ask. The whole approach is based on calculations of the degree to which different relatives resemble one another. For these purposes we’re going to continue our focus on phenotypes influenced by one locus with two alleles, and we’ll do the calculations in detail only for half sib families. We start with something that may look vaguely familiar.³ Take a look at Table 19.1.

Note also that the probabilities in Table 19.1 are appropriate *only* if the progeny are from half-sib families. If the progeny are from full-sib families, we must specify the frequency of each of the nine possible matings (keeping track of the genotype of both mother and father) and the offspring that each will produce.⁴

another because of similarities in the environment they experienced, too.

²You’ll see the reason for the quotes around genotype in this paragraph and the last a little later. It’s a little more complex than what I’ve suggested.

³Remember our mother-offspring combinations with *Zoarces viviparus*?

⁴To check your understanding of all of this, you might want to try to produce the appropriate table.

Covariance of two random variables

Let p_{xy} be the probability that random variable X takes the value x and random variable Y takes the value y . Then the covariance between X and Y is:

$$\text{Cov}(X, Y) = \sum p_{xy}(x - \mu_x)(y - \mu_y) \quad ,$$

where μ_x is the mean of X and μ_y is the mean of Y .

Covariance between half-siblings

Here's how we can calculate the covariance between half-siblings: First, imagine selecting huge number of half-sibs pairs at random. The phenotype of the first half-sib in the pair is a random variable (call it S_1), as is the phenotype of the second (call it S_2). The mean of S_1 is just the mean phenotype in *all* the progeny taken together, \bar{x} . Similarly, the mean of S_2 is just \bar{x} . Now with one locus, two alleles we have three possible phenotypes: x_{11} (corresponding to the genotype A_1A_1), x_{12} (corresponding to the genotype A_1A_2), and x_{22} (corresponding to the genotype A_2A_2). So all we need to do to calculate the covariance between half-sibs is to write down all possible pairs of phenotypes and the frequency with which they will occur in our sample of randomly chosen half-sibs based on the frequencies in Table 19.1 above and the frequency of maternal genotypes. It's actually a bit easier to keep track of it all if we write down the frequency of each maternal genotype and the frequency with which each possible phenotypic combination will occur in her progeny.

$$\begin{aligned} \text{Cov}(S_1, S_2) &= p^2[p^2(x_{11} - \bar{x})^2 + 2pq(x_{11} - \bar{x})(x_{12} - \bar{x}) + q^2(x_{12} - \bar{x})^2] \\ &\quad + 2pq\left[\frac{1}{4}p^2(x_{11} - \bar{x})^2 + \frac{1}{2}p(x_{11} - \bar{x})(x_{12} - \bar{x}) + \frac{1}{2}pq(x_{11} - \bar{x})(x_{22} - \bar{x})\right. \\ &\quad \left. + \frac{1}{4}(x_{12} - \bar{x})^2 + \frac{1}{2}q(x_{12} - \bar{x})(x_{22} - \bar{x}) + \frac{1}{4}q^2(x_{22} - \bar{x})^2\right] \\ &\quad + q^2[p^2(x_{12} - \bar{x})^2 + 2pq(x_{12} - \bar{x})(x_{22} - \bar{x}) + q^2(x_{22} - \bar{x})^2] \\ &= p^2[p(x_{11} - \bar{x}) + q(x_{12} - \bar{x})]^2 \\ &\quad + 2pq\left[\frac{1}{2}p(x_{11} - \bar{x}) + \frac{1}{2}q(x_{12} - \bar{x}) + \frac{1}{2}p(x_{12} - \bar{x}) + \frac{1}{2}q(x_{22} - \bar{x})\right]^2 \\ &\quad + q^2[p(x_{12} - \bar{x}) + q(x_{22} - \bar{x})]^2 \\ &= p^2[px_{11} + qx_{12} - \bar{x}]^2 \\ &\quad + 2pq\left[\frac{1}{2}(px_{11} + qx_{12} - \bar{x}) + \frac{1}{2}(px_{12} + qx_{22} - \bar{x})\right]^2 \\ &\quad + q^2[px_{12} + qx_{22} - \bar{x}]^2 \end{aligned}$$

Genotype	A_1A_1	A_1A_2	A_2A_2
Phenotype	0	0.8	2

Table 19.2: An example of a non-additive relationship between genotypes and phenotypes.

Maternal genotype	Frequency	Offspring genotype		
		A_1A_1	A_1A_2	A_2A_2
A_1A_1	0.16	0.4	0.6	0.0
A_1A_2	0.48	0.2	0.5	0.3
A_2A_2	0.36	0.0	0.4	0.6

Table 19.3: Mother-offspring combinations (half-sib) for the numerical example in Table 19.2.

$$\begin{aligned}
&= p^2 \left[\alpha_1 - \frac{\bar{x}}{2} \right]^2 + 2pq \left[\frac{1}{2}(\alpha_1 - \frac{\bar{x}}{2}) + \frac{1}{2}(\alpha_2 - \frac{\bar{x}}{2}) \right]^2 + q^2 \left[\alpha_2 - \frac{\bar{x}}{2} \right]^2 \\
&= p^2 \left[\frac{1}{2}(2\alpha_1 - \bar{x}) \right]^2 + 2pq \left[\frac{1}{2}(\alpha_1 + \alpha_2 - \bar{x}) \right]^2 + q^2 \left[\frac{1}{2}(2\alpha_2 - \bar{x}) \right]^2 \\
&= \left(\frac{1}{4} \right) \left[p^2(2\alpha_1 - \bar{x})^2 + 2pq[(\alpha_1 + \alpha_2 - \bar{x})]^2 + q^2(2\alpha_2 - \bar{x})^2 \right] \\
&= \left(\frac{1}{4} \right) V_a
\end{aligned}$$

A numerical example

Now we'll return to an example we saw earlier (Table 19.2). This set of genotypes and phenotypes may look familiar. It is the same one we encountered earlier when we calculated additive and dominance components of variance. Let's assume that $p = 0.4$. Then we know that

$$\begin{aligned}
\bar{x} &= 1.104 \\
V_a &= 0.5192 \\
V_d &= 0.0092 \quad .
\end{aligned}$$

We can also calculate the numerical version of Table 19.1, which you'll find in Table 19.3.

So now we can follow the same approach we did before and calculate the numerical value of the covariance between half-sibs in this example:

$$\text{Cov}(S_1, S_2) = [(0.4)^2(0.16) + (0.2)^2(0.48)](0 - 1.104)^2$$

MZ twins (Cov_{MZ})	$V_a + V_d$
Parent-offspring (Cov_{PO}) ¹	$\left(\frac{1}{2}\right) V_a$
Full sibs (Cov_{FS})	$\left(\frac{1}{2}\right) V_a + \left(\frac{1}{4}\right) V_d$
Half sibs (Cov_{HS})	$\left(\frac{1}{4}\right) V_a$

¹One parent or mid-parent.

Table 19.4: Genetic covariances among relatives.

$$\begin{aligned}
& +[(0.6)^2(0.16) + (0.5)^2(0.48) + (0.4)^2(0.36)](0.8 - 1.104)^2 \\
& +[(0.3)^2(0.48) + (0.6)^2(0.36)](2 - 1.104)^2 \\
& +2[(0.4)(0.6)(0.16) + (0.2)(0.5)(0.48)](0 - 1.104)(0.8 - 1.104) \\
& +2(0.2)(0.3)(0.48)(0 - 1.104)(2 - 1.104) \\
& +2[(0.5)(0.3)(0.48) + (0.4)(0.6)(0.36)](0.8 - 1.104)(2.0 - 1.104) \\
= & 0.1298 \\
= & \left(\frac{1}{4}\right) 0.5192 \quad .
\end{aligned}$$

Covariances among relatives

Well, if we can do this sort of calculation for half-sibs, you can probably guess that it's also possible to do it for other relatives. I won't go through all of the calculations, but the results are summarized in Table 19.4

Estimating heritability

Galton introduced the term *regression* to describe the inheritance of height in humans. He noted that there is a tendency for adult offspring of tall parents to be tall and of short parents to be short, but he also noted that offspring tended to be less extreme than the parents. He described this as a “regression to mediocrity,” and statisticians adopted the term to describe a standard technique for describing the functional relationship between two variables.

Regression analysis

Measure the parents. Regress the offspring phenotype on: (1) the phenotype of one parent or (2) the mean of the parental phenotypes. In either case, the covariance between

the parental phenotype and the offspring genotype is $\left(\frac{1}{2}\right) V_a$. Now the regression coefficient between one parent and offspring, $b_{P \rightarrow O}$, is

$$\begin{aligned} b_{P \rightarrow O} &= \frac{\text{Cov}_{PO}}{\text{Var}(P)} \\ &= \frac{\left(\frac{1}{2}\right) V_a}{V_p} \\ &= \left(\frac{1}{2}\right) h_N^2 \quad . \end{aligned}$$

In short, the slope of the regression line is equal to one-half the narrow sense heritability. In the regression of offspring on mid-parent value,

$$\begin{aligned} \text{Var}(MP) &= \text{Var}\left(\frac{M + F}{2}\right) \\ &= \frac{1}{4} \text{Var}(M + F) \\ &= \frac{1}{4} (\text{Var}(M) + \text{Var}(F)) \\ &= \frac{1}{4} (2V_p) \\ &= \frac{1}{2} V_p \quad . \end{aligned}$$

Thus, $b_{MP \rightarrow O} = \frac{1}{2} V_a / \frac{1}{2} V_p = h_N^2$. In short, the slope of the regression line is equal to the narrow sense heritability.

Sib analysis

Mate a number of males (sires) with a number of females (dams). Each sire is mated to more than one dam, but each dam mates only with one sire. Do an analysis of variance on the phenotype in the progeny, treating sire and dam as main effects. The result is shown in Table 19.5.

Now we need some way to relate the variance components (σ_W^2 , σ_D^2 , and σ_S^2) to V_a , V_d , and V_e .⁵ How do we do that? Well,

$$V_p = \sigma_T^2 = \sigma_S^2 + \sigma_D^2 + \sigma_W^2 \quad .$$

⁵ σ_W^2 , σ_D^2 , and σ_S^2 are often referred to as the *observational* components of variance, because they are estimated from observations we make on phenotypic variation. V_a , V_d , and V_e are often referred to as the *causal* components of variance, because they represent the genetic and environmental influences on trait expression.

Source	d.f.	Mean square	Composition of mean square
Among sires	$s - 1$	MS_S	$\sigma_W^2 + k\sigma_D^2 + dk\sigma_s^2$
Among dams (within sires)	$s(d - 1)$	MS_D	$\sigma_W^2 + k\sigma_D^2$
Within progenies	$sd(k - 1)$	MS_W	σ_W^2

s = number of sires

d = number of dams per sire

k = number of offspring per dam

Table 19.5: Analysis of variance table for a full-sib analysis of quantitative genetic variation.

σ_S^2 estimates the variance among the means of the half-sib families fathered by each of the different sires or, equivalently, the covariance among half-sibs.⁶

$$\begin{aligned}\sigma_S^2 &= \text{Cov}_{HS} \\ &= \left(\frac{1}{4}\right) V_a \quad .\end{aligned}$$

Now consider the within progeny component of the variance, σ_W^2 . In general, it can be shown that *any* among group variance component is equal to the covariance among the members within the groups.⁷ Thus, a within group component of the variance is equal to the total variance minus the covariance within groups. In this case,

$$\begin{aligned}\sigma_W^2 &= V_p - \text{Cov}_{FS} \\ &= V_a + V_d + V_e - \left[\left(\frac{1}{2}\right) V_a + \left(\frac{1}{4}\right) V_d\right] \\ &= \left(\frac{1}{2}\right) V_a + \left(\frac{3}{4}\right) V_d + V_e \quad .\end{aligned}$$

There remains only σ_D^2 . Now $\sigma_W^2 = V_p - \text{Cov}_{FS}$, $\sigma_S^2 = \text{Cov}_{HS}$, and $\sigma_T^2 = V_p$. Thus,

$$\sigma_D^2 = \sigma_T^2 - \sigma_S^2 - \sigma_W^2$$

⁶To see why consider this is so, consider the following: The mean genotypic value of half-sib families with an A_1A_1 mother is $px_{11} + qx_{12}$; with an A_1A_2 mother, $px_{11}/2 + qx_{12}/2 + px_{12}/2 + qx_{22}/2$; with an A_2A_2 mother, $px_{12} + qx_{22}$. The equation for the variance of these means is identical to the equation for the covariance among half-sibs.

⁷With $x_{ij} = a_i + \epsilon_{ij}$, where a_i is the mean group effect and ϵ_{ij} is random effect on individual j in group i (with mean 0), $\text{Cov}(x_{ij}, x_{ik}) = E(a_i + \epsilon_{ij} - \mu)(a_i + \epsilon_{ik} - \mu) = E((a_i - \mu)^2 + a_i(\epsilon_{ij} + \epsilon_{ik}) + \epsilon_{ij}\epsilon_{ik}) = \text{Var}(A)$.

$$\begin{aligned}
&= V_p - \text{Cov}_{HS} - (V_p - \text{Cov}_{FS}) \\
&= \text{Cov}_{FS} - \text{Cov}_{HS} \\
&= \left[\left(\frac{1}{2}\right) V_a + \left(\frac{1}{4}\right) V_d \right] - \left(\frac{1}{4}\right) V_a \\
&= \left(\frac{1}{4}\right) V_a + \left(\frac{1}{4}\right) V_d \quad .
\end{aligned}$$

So if we rearrange these equations, we can express the genetic components of the phenotypic variance, the *causal* components of variance, as simple functions of the *observational* components of variance:

$$\begin{aligned}
V_a &= 4\sigma_S^2 \\
V_d &= 4(\sigma_D^2 - \sigma_S^2) \\
V_e &= \sigma_W^2 - 3\sigma_D^2 + \sigma_S^2 \quad .
\end{aligned}$$

Furthermore, the narrow-sense heritability is given by

$$h_N^2 = \frac{4\sigma_S^2}{\sigma_S^2 + \sigma_D^2 + \sigma_W^2} \quad .$$

An example: body weight in female mice

The analysis involves 719 offspring from 74 sires and 192 dams, each with one litter. The offspring were spread over 4 generations, and the analysis is performed as a nested ANOVA with the genetic analysis nested *within* generations. An additional complication is that the design was unbalanced, i.e., unequal numbers of progeny were measured in each sibship. As a result the degrees of freedom don't work out to be quite as simple as what I showed you.⁸ The results are summarized in Table 19.6.

Using the expressions for the composition of the mean square we obtain

$$\begin{aligned}
\sigma_W^2 &= MS_W \\
&= 2.19 \\
\sigma_D^2 &= \left(\frac{1}{k}\right) (MS_D - \sigma_W^2) \\
&= 2.47 \\
\sigma_S^2 &= \left(\frac{1}{dk'}\right) (MS_S - \sigma_W^2 - k'\sigma_D^2) \\
&= 0.48 \quad .
\end{aligned}$$

⁸What did you expect from real data? This example is extracted from Falconer and Mackay, pp. 169–170. See the book for details.

Source	d.f.	Mean square	Composition of mean square
Among sires	70	17.10	$\sigma_W^2 + k'\sigma_D^2 + dk'\sigma_s^2$
Among dams (within sires)	118	10.79	$\sigma_W^2 + k\sigma_D^2$
Within progenies	527	2.19	σ_W^2

$d = 2.33$
 $k = 3.48$
 $k' = 4.16$

Table 19.6: Quantitative genetic analysis of the inheritance of body weight in female mice (from Falconer and Mackay, pp. 169–170.)

Thus,

$$\begin{aligned}
 V_p &= 5.14 \\
 V_a &= 1.92 \\
 V_d + V_e &= 3.22 \\
 V_d &= (0.00\text{—}1.64) \\
 V_e &= (1.58\text{—}3.22)
 \end{aligned}$$

Why didn't I give a definite number for V_d after my big spiel above about how we can estimate it from a full-sib crossing design? Two reasons. First, if you plug the estimates for σ_D^2 and σ_S^2 into the formula above for V_d you get $V_d = 7.96$, $V_e = -4.74$, which is clearly impossible since V_d has to be less than V_p and V_e has to be greater than zero. It's a variance. Second, the experimental design confounds two sources of resemblance among full siblings: (1) genetic covariance and (2) environmental covariance. The full-sib families were all raised by the same mother in the same pen. Hence, we don't know to what extent their resemblance is due to a common natal environment.⁹ If we assume $V_d = 0$, we can estimate the amount of variance accounted for by exposure to a common natal environment, $V_{Ec} = 1.99$, and by environmental variation within sibships, $V_{Ew} = 1.23$.¹⁰ Similarly, if we assume $V_{Ew} = 0$, then $V_d = 1.64$ and $V_{Ec} = 1.58$. In any case, we can estimate the narrow sense heritability

⁹Notice that this doesn't affect our analysis of half-sib families, i.e., the progeny of different sires, since each father was bred with several females

¹⁰See Falconer for details.

as

$$\begin{aligned} h_N^2 &= \left(\frac{1.92}{5.14} \right) \\ &= 0.37 \quad . \end{aligned}$$

Chapter 20

Partitioning variance with WinBUGS

As you have already guessed, there is a way to partition variance with WinBUGS. One of the advantages of using WinBUGS for variance partitioning is that it is easy to get direct estimates of both the observational and the causal components of variance in a quantitative genetics design.¹

The data has a fairly straightforward structure, although it will look different from other WinBUGS data you've seen:

```
sire[] dam[] weight[]
1 1 99.2789343018915
1 1 92.8290811609914
1 1 92.1977611557311
1 1 96.2203721102781
1 1 97.435349690979
1 1 99.5683756647001
1 1 92.2818423839868
1 1 99.5110665900414
1 2 103.230303160319
. . .
. . .
7 52 102.100956942292
7 53 117.842216117441
```

¹I'll only describe the observational component here. Getting the causal components will be left as an exercise to be completed as part of Problem #4. The source code and data used for this example is available at <http://darwin.eeb.uconn.edu/eeb348/supplements-2006/partitioning-example.txt>.

```

7 53 124.174111459722
7 53 115.02591859104
7 53 111.885901553834
7 53 115.320308219732

```

The first column is the number of the sire involved in the mating, the second column is the number of the dam involved in the mating, and the last column is the weight of an individual offspring.

Recall that our objective is to describe how much of the overall variation is due to weight differences among individuals that share the same mother, how much is due to differences among mothers in the average weight of their offspring, and how much is due to differences among fathers in the average weight of their offspring. One way of approaching this is to imagine that the weight of each individual is determined as follows:

$$\text{weight}_i = \text{mean}(\text{weight}) + \text{sire contribution} + \text{dam contribution} + \text{error} \quad ,$$

where the sire contribution, the dam contribution, and error are normally distributed random variables with mean 0 and variances of V_s , V_d , and V_w , respectively. We can translate that into WinBUGS code like this:

```

# offspring
for (i in 1:300) {
  weight[i] ~ dnorm(mu[i], tau.within)
  mu[i] <- nu + alpha[sire[i]] + beta[dam[i]]
}

# sires
for (i in 1:6) {
  alpha[i] ~ dnorm(0.0, tau.sire)
}
alpha[7] <- -sum(alpha[1:6])

# dams
for (i in 1:52) {
  beta[i] ~ dnorm(0.0, tau.dam)
}
beta[53] <- -sum(beta[1:52])

# precisions

```

Variance component	Estimate (2.5%, 97.5%)
v.sire	8.18 (1.306, 31.29)
v.dam	4.571 (2.281, 7.974)
v.within	11.48 (9.599, 13.68)

Table 20.1: Estimated variance components from the full-sib data.

```
tau.sire <- 1.0/v.sire
tau.dam <- 1.0/v.dam
tau.within <- 1.0/v.within
```

Why the funny precisions, `tau.sire`, `tau.dam`, and `tau.within`? It turns out that it was easier for the people who wrote WinBUGS to describe normal distributions in terms of a mean and precision, where precision = 1/variance, than in terms of a mean and variance. Once we know that we need to work with precisions, then all we need to do is to put appropriate priors on `nu`, `v.within`, `v.dam`, and `v.sire`. The most convenient way to put priors on the variance components is to put a uniform distribution on the corresponding standard deviations.

```
# priors
nu ~ dnorm(0, 0.001)
sd.sire ~ dunif(0, 10)
sd.dam ~ dunif(0, 10)
sd.within ~ dunif(0, 10)

# observational components
v.sire <- sd.sire*sd.sire
v.dam <- sd.dam*sd.dam
v.within <- sd.within*sd.within
```

If you put this all together in a WinBUGS model, then WinBUGS will produce results like those shown in Table 20.1.²

²I would recommend using 10,000 iterations at each step of the analysis. This model takes a bit longer to settle down than the other ones we've seen.

Chapter 21

Evolution of quantitative traits

Let's stop and review quickly where we've come and where we're going. We started our survey of quantitative genetics by pointing out that our objective was to develop a way to describe the patterns of phenotypic resemblance among relatives. The challenge was that we wanted to do this for phenotypic traits that whose expression is influenced both by many genes and by the environment in which those genes are expressed. Beyond the technical, algebraic challenges associated with many genes, we have the problem that we can't directly associate particular genotypes with particular phenotypes. We have to rely on patterns of *phenotypic* resemblance to tell us something about how *genetic* variation is transmitted. Surprisingly, we've managed to do that. We now know that it's possible to:

- Estimate the additive effect of an allele.¹
- Partition the phenotypic variance into genotypic and environmental components and to partition the genotypic variance into additive and dominance components.²
- Estimate all of the variance components from a combination of appropriate crossing designs and appropriate statistical analyses.

¹Actually, we don't know this. You'll have to take my word for it that in certain breeding designs its possible to estimate not only the additive genetic variance and the dominance genetic variance, but also the actual additive effect of "alleles" that we haven't even identified. We'll see a more direct approach soon, when we get to quantitative trait locus analysis.

²I should point out that this is an oversimplification. I've mentioned that we typically assume that we can simply add the effects of alleles across loci, but if you think about how genes actually work in organisms, you realize that such additivity across loci isn't likely to be very common. Strictly speaking there are epistatic components to the genetic variance too, i.e., components of the genetic variance that have to do not with the interaction among alleles at a single locus (the dominance variance that we've already encountered), but with the interaction of alleles at different loci.

Now we're ready for the next step: applying all of these ideas to the evolution of a quantitative trait.

Evolution of the mean phenotype

We're going to focus on how the mean phenotype in a population changes in response to natural selection, specifically in response to viability selection. Before we can do this, however, we need to think a bit more carefully about the relationship between genotype, phenotype, and fitness. Let $F_{ij}(x)$ be the probability that genotype A_iA_j has a phenotype smaller than x .³ Then x_{ij} , the genotypic value of A_iA_j is

$$x_{ij} = \int_{-\infty}^{\infty} x dF_{ij}(x)$$

and the population mean phenotype is $p^2x_{11} + 2pqx_{12} + q^2x_{22}$. If an individual with phenotype x has fitness $w(x)$, then the fitness of an individual with genotype A_iA_j is

$$w_{ij} = \int_{-\infty}^{\infty} w(x) dF_{ij}(x)$$

and the mean fitness in the population is $\bar{w} = p^2w_{11} + 2pqw_{12} + q^2w_{22}$.

Now, there's a well known theorem from calculus known as Taylor's theorem. It says that for any function⁴ $f(x)$

$$f(x) = f(a) + \sum_{k=1}^{\infty} \left(\frac{(x-a)^k}{k!} \right) f^{(k)}(a) \quad .$$

Using this theorem we can produce an approximate expression describing how the mean phenotype in a population will change in response to selection. Remember that the mean phenotype, \bar{x} , depends both on the underlying genotypic values and on the allele frequency. So I'm going to write the mean phenotype as $\bar{x}(p)$ to remind us of that dependency.

$$\bar{x}(p') = \bar{x}(p) + (p' - p) \left(\frac{d\bar{x}}{dp} \right) + o(p^2)$$

$$\bar{x}(p) = p^2x_{11} + 2pqx_{12} + q^2x_{22}$$

³For those of you who have had probability theory, $F_{ij}(x)$ is the cumulative distribution for the probability density for phenotype associated with A_iA_j .

⁴Actually there are restrictions on the functions to which it applies, but we can ignore those restrictions for our purposes.

$$\begin{aligned}
\frac{d\bar{x}(p)}{dp} &= 2px_{11} + 2qx_{12} - 2px_{12} - 2qx_{22} \\
&= 2\{(px_{11} + qx_{12} - \bar{x}/2) + (px_{12} + qx_{22} - \bar{x}/2)\} \\
&= 2(\alpha_1 - \alpha_2)
\end{aligned}$$

$$\bar{x}(p') \approx \bar{x}(p) + (p' - p)(2(\alpha_1 - \alpha_2))$$

$$\Delta\bar{x} = (\Delta p)(2(\alpha_1 - \alpha_2))$$

Now you need to remember (from lo those many weeks ago) that

$$p' = \frac{p^2w_{11} + pqw_{12}}{\bar{w}} \ .$$

Thus,

$$\begin{aligned}
\Delta p &= p' - p \\
&= \frac{p^2w_{11} + pqw_{12}}{\bar{w}} - p \\
&= \frac{p^2w_{11} + pqw_{12} - p\bar{w}}{\bar{w}} \\
&= p \left(\frac{pw_{11} + qw_{12} - \bar{w}}{\bar{w}} \right) \ .
\end{aligned}$$

Now,⁵ let's do a linear regression of fitness on phenotype. After all, to make any further progress, we need to relate phenotype to fitness, so that we can use the relationship between phenotype and genotype to infer the change in allele frequencies, from which we will infer the change in mean phenotype.⁶ From our vast statistical knowledge, we know that the slope of this regression line is

$$\beta_1 = \frac{\text{Cov}(w, x)}{\text{Var}(x)}$$

and its intercept is

$$\beta_0 = \bar{w} - \beta_1\bar{x} \ .$$

⁵Since we're having so much fun with mathematics why should we stop here?

⁶Whew! That was a mouthful.

Let's use this regression equation to determine the fitness of each genotype. This is only an approximation to the true fitness,⁷ but it is adequate for many purposes.

$$\begin{aligned}
 w_{ij} &= \int_{-\infty}^{\infty} w(x) dF_{ij}(x) \\
 &\approx \int_{-\infty}^{\infty} (\beta_0 + \beta_1 x) dF_{ij}(x) \\
 &= \beta_0 + \beta_1 x_{ij} \\
 \bar{w} &= \beta_0 + \beta_1 \bar{x} \quad .
 \end{aligned}$$

If we substitute this into our expression for Δp above, we get

$$\begin{aligned}
 \Delta p &= p \left(\frac{pw_{11} + qw_{12} - \bar{w}}{\bar{w}} \right) \\
 &= p \left(\frac{p(\beta_0 + \beta_1 x_{11}) + q(\beta_0 + \beta_1 x_{12}) - (\beta_0 + \beta_1 \bar{x})}{\bar{w}} \right) \\
 &= p\beta_1 \left(\frac{px_{11} + qx_{12} - \bar{x}}{\bar{w}} \right) \\
 &= p\beta_1 \left(\frac{\alpha_1 - \bar{x}/2}{\bar{w}} \right) \\
 &= p\beta_1 \left(\frac{\alpha_1 - (p\alpha_1 + q\alpha_2)}{\bar{w}} \right) \\
 &= \frac{pq\beta_1(\alpha_1 - \alpha_2)}{\bar{w}} \quad .
 \end{aligned}$$

So where are we now?⁸ Let's substitute this result back into the equation for $\Delta \bar{x}$. When we do we get

$$\begin{aligned}
 \Delta \bar{x} &= (\Delta p) (2(\alpha_1 - \alpha_2)) \\
 &= \left(\frac{pq\beta_1(\alpha_1 - \alpha_2)}{\bar{w}} \right) (2(\alpha_1 - \alpha_2)) \\
 &= 2pq\alpha^2 \left(\frac{\beta_1}{\bar{w}} \right) \\
 &= V_a \left(\frac{\beta_1}{\bar{w}} \right) \quad .
 \end{aligned}$$

⁷Specifically, we are implicitly assuming that the fitnesses are adequately approximated by a linear function of our phenotypic measure.

⁸You don't have to tell me where you *wish* you were. I can reliably guess that it's not here.

This is great if we've done the regression between fitness and phenotype, but what if we haven't?⁹ Let's look at $\text{Cov}(w, x)$ in a little more detail.

$$\begin{aligned}
\text{Cov}(w, x) &= p^2 \int_{-\infty}^{\infty} xw(x)dF_{11}(x) + 2pq \int_{-\infty}^{\infty} xw(x)dF_{12}(x) \\
&\quad + q^2 \int_{-\infty}^{\infty} xw(x)dF_{22}(x) - \bar{x}\bar{w} \\
&= p^2 \left(\int_{-\infty}^{\infty} xw(x)dF_{11}(x) - x_{11}\bar{w} + x_{11}\bar{w} \right) \\
&\quad + 2pq \left(\int_{-\infty}^{\infty} xw(x)dF_{11}(x) - x_{12}\bar{w} + x_{12}\bar{w} \right) \\
&\quad + q^2 \left(\int_{-\infty}^{\infty} xw(x)dF_{22}(x) - x_{22}\bar{w} + x_{22}\bar{w} \right) \\
&\quad - \bar{x}\bar{w} \\
&= p^2 \left(\int_{-\infty}^{\infty} xw(x)dF_{11}(x) - x_{11}\bar{w} \right) \\
&\quad + 2pq \left(\int_{-\infty}^{\infty} xw(x)dF_{11}(x) - x_{12}\bar{w} \right) \\
&\quad + q^2 \left(\int_{-\infty}^{\infty} xw(x)dF_{22}(x) - x_{22}\bar{w} \right) \quad .
\end{aligned}$$

Now

$$\begin{aligned}
\int_{-\infty}^{\infty} xw(x)dF_{ij}(x) - x_{ij}\bar{w} &= \bar{w} \left(\int_{-\infty}^{\infty} \frac{xw(x)}{\bar{w}} dF_{ij}(x) - x_{ij} \right) \\
&= \bar{w}(x_{ij}^* - x_{ij}) \quad ,
\end{aligned}$$

where x_{ij}^* refers to the mean phenotype of A_iA_j after selection. So

$$\begin{aligned}
\text{Cov}(w, x) &= p^2\bar{w}(x_{11}^* - x_{11}) + 2pq\bar{w}(x_{12}^* - x_{12}) + q^2\bar{w}(x_{22}^* - x_{22}) \\
&= \bar{w}(\bar{x}^* - \bar{x}) \quad ,
\end{aligned}$$

where \bar{x}^* is the population mean phenotype after selection. In short,¹⁰ combining our equations for the change in mean phenotype and for the covariance of fitness and phenotype and remembering that $\beta_1 = \text{Cov}(w, x)/\text{Var}(x)$ ¹¹

$$\Delta\bar{x} = V_a \left(\frac{\bar{w}(\bar{x}^* - \bar{x})}{\bar{w}} \right)$$

⁹Hang on just a little while longer. We're almost there.

¹⁰We finally made it.

¹¹You also need to remember that $\text{Var}(x) = V_p$, since they're the same thing, the phenotypic variance.

Genotype	A_1A_1	A_1A_2	A_2A_2
Phenotype	1.303	1.249	0.948

Table 21.1: A simple example to illustrate response to selection in a quantitative trait.

$$= h_N^2(\bar{x}^* - \bar{x})$$

$\Delta\bar{x} = \bar{x}' - \bar{x}$ is referred to as the response to selection and is often given the symbol R . It is the change in population mean between the parental generation (before selection) and the offspring generation (before selection). $\bar{x}^* - \bar{x}$ is referred to as the selection differential and is often given the symbol S . It is the difference between the mean phenotype in the parental generation before selection and the mean phenotype in the parental generation after selection. Thus, we can rewrite our final equation as

$$R = h_N^2 S \quad .$$

This equation is often referred to as the *breeders equation*.

A Numerical Example

To illustrate how this works, let's examine the simple example in Table 21.1.

Given these phenotypes, $p = 0.25$, and $V_p = 0.16$, it follows that $\bar{x} = 1.08$ and $h_N^2 = 0.1342$. Suppose the mean phenotype after selection is 1.544. What will the phenotype be among the newly born progeny?

$$\begin{aligned}
 S &= \bar{x}^* - \bar{x} \\
 &= 1.544 - 1.08 \\
 &= 0.464 \\
 \Delta\bar{x} &= h_N^2 S \\
 &= (0.1342)(0.464) \\
 &= 0.06 \\
 \bar{x}' &= \bar{x} + \Delta\bar{x} \\
 &= 1.08 + 0.06 \\
 &= 1.14
 \end{aligned}$$

Genotype	A_1A_1	A_1A_2	A_2A_2
Frequency	p^2	$2pq$	q^2
Fitness	w_{11}	w_{12}	w_{22}
Additive fitness value	$2\alpha_1$	$\alpha_1 + \alpha_2$	$2\alpha_2$

Table 21.2: Fitnesses and additive fitness values used in deriving Fisher’s Fundamental Theorem of Natural Selection.

Fisher’s Fundamental Theorem of Natural Selection

Suppose the phenotype whose evolution we’re interested in following is fitness itself.¹² Then we can summarize the fitnesses as illustrated in Table 21.2.

Although I didn’t tell you this, a well-known fact about viability selection at one locus is that the change in allele frequency from one generation to the next can be written as

$$\Delta p = \left(\frac{pq}{2\bar{w}} \right) \left(\frac{d\bar{w}}{dp} \right) \quad .$$

Using our new friend, Taylor’s theorem, it follows immediately that

$$\bar{w}' = \bar{w} + (\Delta p) \left(\frac{d\bar{w}}{dp} \right) + \left(\frac{(\Delta p)^2}{2} \right) \left(\frac{d^2\bar{w}}{dp^2} \right) \quad .$$

Or, equivalently

$$\Delta\bar{w} = (\Delta p) \left(\frac{d\bar{w}}{dp} \right) + \left(\frac{(\Delta p)^2}{2} \right) \left(\frac{d^2\bar{w}}{dp^2} \right) \quad .$$

Recalling that $\bar{w} = p^2w_{11} + 2p(1-p)w_{12} + (1-p)^2w_{22}$ we find that

$$\begin{aligned} \frac{d\bar{w}}{dp} &= 2pw_{11} + 2(1-p)w_{12} - 2pw_{12} - 2(1-p)w_{22} \\ &= 2[(pw_{11} + qw_{12}) - (pw_{12} + qw_{22})] \\ &= 2[(pw_{11} + qw_{12} - \bar{w}/2) - (pw_{12} + qw_{22} - \bar{w}/2)] \\ &= 2[\alpha_1 - \alpha_2] \\ &= 2\alpha \quad , \end{aligned}$$

¹²The proof of the fundamental theorem that follows is due to C. C. Li [61]

where the last two steps use the definitions for α_1 and α_2 , and we set $\alpha = \alpha_1 - \alpha_2$. Similarly,

$$\begin{aligned}\frac{d^2\bar{w}}{dp^2} &= 2w_{11} - 2w_{12} - 2w_{12} + 2w_{22} \\ &= 2(w_{11} - 2w_{12} + w_{22})\end{aligned}$$

Now we can plug these back into the equation for $\Delta\bar{w}$:

$$\begin{aligned}\Delta\bar{w} &= \left\{ \left(\frac{pq}{2\bar{w}} \right) \left(\frac{d\bar{w}}{dp} \right) \right\} \left(\frac{d\bar{w}}{dp} \right) + \frac{\left\{ \left(\frac{pq}{2\bar{w}} \right) \left(\frac{d\bar{w}}{dp} \right) \right\}^2}{2} [2(w_{11} - 2w_{12} + w_{22})] \\ &= \left\{ \left(\frac{pq}{2\bar{w}} \right) (2\alpha) \right\} (2\alpha) + \left\{ \left(\frac{pq}{2\bar{w}} \right) (2\alpha) \right\}^2 (w_{11} - 2w_{12} + w_{22}) \\ &= \frac{2pq\alpha^2}{\bar{w}} + \frac{p^2q^2\alpha^2}{\bar{w}^2} (w_{11} - 2w_{12} + w_{22}) \\ &= \frac{V_a}{\bar{w}} \left\{ 1 + \frac{pq}{2\bar{w}} (w_{11} - 2w_{12} + w_{22}) \right\} \quad ,\end{aligned}$$

where the last step follows from the observation that $V_a = 2pq\alpha^2$. The quantity $\frac{pq}{2\bar{w}}(w_{11} - 2w_{12} + w_{22})$ is usually quite small, especially if selection is not too intense. So we are left with

$$\Delta\bar{w} \approx \frac{V_a}{\bar{w}} \quad .$$

Chapter 22

Selection on multiple characters

So far we've studied only the evolution of a single trait, e.g., height or weight. But organisms have many traits, and they evolve at the same time. How can we understand their simultaneous evolution? The basic framework of the quantitative genetic approach was first outlined by Russ Lande and Steve Arnold [59].

Let z_1, z_2, \dots, z_n be the phenotype of each character that we are studying. We'll use $\bar{\mathbf{z}}$ to denote the vector of these characters before selection and $\bar{\mathbf{z}}^*$ to denote the vector after selection. The selection differential, \mathbf{s} , is also a vector given by

$$\mathbf{s} = \bar{\mathbf{z}}^* - \bar{\mathbf{z}} \quad .$$

Suppose $p(\mathbf{z})$ is the probability that any individual has phenotype \mathbf{z} , and let $W(\mathbf{z})$ be the fitness (absolute viability) of an individual with phenotype \mathbf{z} . Then the mean absolute fitness is

$$\bar{W} = \int W(\mathbf{z})p(\mathbf{z})d\mathbf{z} \quad .$$

The relative fitness of phenotype \mathbf{z} can be written as

$$w(\mathbf{z}) = \frac{W(\mathbf{z})}{\bar{W}} \quad .$$

Using relative fitnesses the mean relative fitness, \bar{w} , is 1. Now

$$\bar{\mathbf{z}}^* = \int \mathbf{z}w(\mathbf{z})p(\mathbf{z})d\mathbf{z} \quad .$$

Recall that $Cov(X, Y) = E(X - \mu_x)(Y - \mu_y) = E(XY) - \mu_x\mu_y$. Consider

$$\begin{aligned} \mathbf{s} &= \bar{\mathbf{z}}^* - \bar{\mathbf{z}} \\ &= \int \mathbf{z}w(\mathbf{z})p(\mathbf{z})d\mathbf{z} - \bar{\mathbf{z}} \\ &= E(w, \mathbf{z}) - \bar{w}\bar{\mathbf{z}} \quad , \end{aligned}$$

where the last step follows since $\bar{w} = 1$ meaning that $\bar{w}\bar{z} = \bar{z}$. In short,

$$\mathbf{s} = \text{Cov}(w, z) \quad .$$

That should look familiar from our analysis of the evolution of a single phenotype.

If we assume that all genetic effects are additive, then the phenotype of an individual can be written as

$$\mathbf{z} = \mathbf{x} + \mathbf{e} \quad ,$$

where \mathbf{x} is the additive genotype and \mathbf{e} is the environmental effect. We'll denote by \mathbf{G} the matrix of genetic variances and covariances and by \mathbf{E} the matrix of environmental variances and covariances. The matrix of phenotype variances and covariances, \mathbf{P} , is then given by¹

$$\mathbf{P} = \mathbf{G} + \mathbf{E} \quad .$$

Now, if we're willing to assume that the regression of additive genetic effects on phenotype is linear² and that the environmental variance is the same for every genotype, then we can predict how phenotypes will change from one generation to the next

$$\begin{aligned} \bar{\mathbf{x}}^* - \bar{\mathbf{x}} &= \mathbf{GP}^{-1}(\bar{\mathbf{z}}^* - \bar{\mathbf{z}}) \\ \bar{\mathbf{z}}' - \bar{\mathbf{z}} &= \mathbf{GP}^{-1}(\bar{\mathbf{z}}^* - \bar{\mathbf{z}}) \\ \Delta\bar{\mathbf{z}} &= \mathbf{GP}^{-1}\mathbf{s} \end{aligned}$$

\mathbf{GP}^{-1} is the multivariate version of h_N^2 . This equation is also the multivariate version of the breeders equation.

But we have already seen that $\mathbf{s} = \text{Cov}(w, z)$. Thus,

$$\boldsymbol{\beta} = \mathbf{P}^{-1}\mathbf{s}$$

is a set of partial regression coefficients of relative fitness on the characters, i.e., the dependence of relative fitness on that character alone holding all others constant.

Note:

$$\begin{aligned} s_i &= \sum_{j=1}^n \beta_j P_{ij} \\ &= \beta_1 P_{i1} + \cdots + \beta_i P_{ii} + \cdots + \beta_n P_{in} \end{aligned}$$

is the total selective differential in character i , including the indirect effects of selection on other characters.

¹Assuming that there are no genotype \times environment interactions.

²And we were willing to do this when we were studying the evolution of only one trait, so why not do it now?

Character	Mean before selection	standard deviation		
head	0.880	0.034		
thorax	2.038	0.049		
scutellum	1.526	0.057		
wing	2.337	0.043		

	head	thorax	scutellum	wing
head	1.00	0.72	0.50	0.60
thorax		1.00	0.59	0.71
scutellum			1.00	0.62
wing				1.00

Character	s	s'	β	β'
head	-0.004	-0.11	-0.7 ± 4.9	-0.03 ± 0.17
thorax	-0.003	-0.06	$11.6 \pm 3.9^{**}$	$0.58 \pm 0.19^{**}$
scutellum	-0.16*	-0.28*	-2.8 ± 2.7	-0.17 ± 0.15
wing	-0.019**	-0.43**	$-16.6 \pm 4.0^{**}$	$-0.74 \pm 0.18^{**}$

Table 22.1: Selection analysis of pentastomid bugs on the shores of Lake Michigan.

An example: selection in a pentastomid bug

94 individuals were collected along shoreline of Lake Michigan in Parker County, Indiana after a storm. 39 were alive, 55 dead. The means of several characters before selection, the trait correlations, and the selection analysis are presented in Table 22.1.

The column labeled s is the selective differential for each character. The column labeled s' is the *standardized* selective differential, i.e., the change measured in units of standard deviation rather than on the original scale.³ A multiple regression analysis of fitness versus phenotype on the original scale gives estimates of β , the direct effect of selection on that trait. A multiple regression analysis of fitness versus phenotype on the transformed scale gives the standardized direct effect of selection, β' , on that trait.

Notice that the selective differential⁴ for the thorax measurement is negative, i.e., individuals that survived had larger thoraces than those that died. But the *direct* effect of selection on thorax is strongly positive, i.e., all other things being equal, an individual with a large

³To measure on this scale the data is simply transformed by setting $y_i = (x_i - \bar{x})/s$, where x_i is the raw score for the i th individual, \bar{x} is the sample mean for the trait, and s is its standard deviation.

⁴The cumulative effect of selection on the change in mean phenotype.

	body	tail
body	35.4606	11.3530
tail	11.3530	37.2973

Table 22.2: Genetic variance-covariance matrix for vertebral number in central Californian garter snakes.

was more likely to survive than one with a small thorax. Why the apparent contradiction? Because the thorax measurement is positively correlated with the wing measurement, and there's strong selection for decreased values of the wing measurement.

Cumulative selection gradients

Arnold [1] suggested an extension of this approach to longer evolutionary time scales. Specifically, he studied variation in the number of body vertebrae and the number of tail vertebrae in populations of *Thamnophis elegans* from two regions of central California. He found relatively little vertebral variation within populations, but there were considerable differences in vertebral number between populations on the coast side of the Coast Ranges and populations on the Central Valley side of the Coast Ranges. The consistent difference suggested that selection might have produced these differences, and Arnold attempted to determine the amount of selection necessary to produce these differences.

The data

Arnold collected pregnant females from two local populations in each of two sites in northern California 282 km apart from one another. Females were collected over a ten-year period and returned to the University of Chicago. Dam-offspring regressions were used to estimate additive genetic variances and covariances of vertebral number.⁵ Mark-release-recapture experiments in the California populations showed that females with intermediate numbers of vertebrae grow at the fastest rate, at least at the inland site, although no such relationship was found in males. The genetic variance-covariance matrix he obtained is shown in Table 22.2.

⁵1000 progeny from 100 dams.

The method

We know from Lande and Arnold's results that the change in multivariate phenotype from one generation to the next, $\Delta\bar{\mathbf{z}}$, can be written as

$$\Delta\bar{\mathbf{z}} = \mathbf{G}\beta \quad ,$$

where \mathbf{G} is the genotypic variance-covariance matrix, $\beta = \mathbf{P}^{-1}\mathbf{s}$ is the set of partial regression coefficients describing the direct effect of each character on relative fitness.⁶ If we are willing to assume that \mathbf{G} remains constant, then the total change in a character subject to selection for n generations is

$$\sum_{k=1}^n \Delta\bar{\mathbf{z}} = \mathbf{G} \sum_{k=1}^n \beta \quad .$$

Thus, $\sum_{k=1}^n \beta$ can be regarded as the cumulative selection differential associated with a particular observed change, and it can be estimated as

$$\sum_{k=1}^n \beta = \mathbf{G}^{-1} \sum_{k=1}^n \Delta\bar{\mathbf{z}} \quad .$$

The results

The overall difference in vertebral number between inland and coastal populations can be summarized as:

$$\begin{aligned} \text{body}_{\text{inland}} - \text{body}_{\text{coastal}} &= 16.21 \\ \text{tail}_{\text{inland}} - \text{tail}_{\text{coastal}} &= 9.69 \end{aligned}$$

Given the estimate of \mathbf{G} already obtained, this corresponds to a cumulative selection gradient between inland and coastal populations of

$$\begin{aligned} \beta_{\text{body}} &= 0.414 \\ \beta_{\text{tail}} &= 0.134 \end{aligned}$$

Applying the same technique to looking at the differences between populations within the inland site and within the coastal site we find cumulative selection gradients of

$$\begin{aligned} \beta_{\text{body}} &= 0.035 \\ \beta_{\text{tail}} &= 0.038 \end{aligned}$$

⁶ \mathbf{P} is the phenotypic variance-covariance matrix and \mathbf{s} is the vector of selection differentials.

for the coastal site and

$$\begin{aligned}\beta_{\text{body}} &= 0.035 \\ \beta_{\text{tail}} &= -0.004\end{aligned}$$

for the inland site.

The conclusions

“To account for divergence between inland and coastal California, we must invoke cumulative forces of selection that are 7 to 11 times stronger than the forces needed to account for differentiation of local populations.”

Furthermore, recall that the selection gradients can be used to partition the overall response to selection in a character into the portion due to the direct effects of that character alone and the portion due to the indirect effects of selection on a correlated character. In this case the overall response to selection in number of body vertebrae is given by

$$\mathbf{G}_{11}\beta_1 + \mathbf{G}_{12}\beta_2 \quad ,$$

where $\mathbf{G}_{11}\beta_1$ is the direct effect of body vertebral number and $\mathbf{G}_{12}\beta_2$ is the indirect effect of tail vertebral number. Similarly, the overall response to selection in number of tail vertebrae is given by

$$\mathbf{G}_{12}\beta_1 + \mathbf{G}_{22}\beta_2 \quad ,$$

where $\mathbf{G}_{22}\beta_2$ is the direct effect of tail vertebral number and $\mathbf{G}_{12}\beta_1$ is the indirect effect of body vertebral number. Using these equations it is straightforward to calculate that 91% of the total divergence in number of body vertebrae is a result of direct selection on this character. In contrast, only 51% of the total divergence in number of tail vertebrae is a result of direct selection on this character, i.e., 49% of the difference in number of tail vertebrae is attributable to indirect selection as a result of its correlation with number of body vertebrae.

The caveats

While the approach Arnold suggests is intriguing, there are a number of caveats that must be kept in mind in trying to apply it.

- This approach assumes that the \mathbf{G} matrix remains constant.
- This approach cannot distinguish strong selection that happened over a short period of time from weak selection that happened over a long period of time.

- This approach *assumes* that the observed differences in populations are the result of selection, but populations isolated from one another will diverge from one another even in the absence of selection simply as a result of genetic drift.
 - Small amount of differentiation between populations within sites could reflect relatively recent divergence of those populations from a common ancestral population.
 - Large amount of differentiation between populations from inland versus coastal sites could reflect a more ancient divergence from a common ancestral population.

Chapter 23

Mapping quantitative trait loci

So far in our examination of the inheritance and evolution of quantitative genetics, we've been satisfied with a purely statistical description of how the phenotypes of parents are related to the phenotypes of their offspring. We've made pretty good progress with that. We know how to partition the phenotypic variance into genetic and phenotypic components and how to partition the genetic variance into additive and dominance components. We know how to predict the degree of resemblance among relatives for any particular trait in terms of the genetic components of variance. We know how to predict how a trait will respond to natural selection.

That's not bad, but in the last 20-25 years the emergence of molecular technologies that allow us to identify large numbers of Mendelian markers has led to a new possibility. It is sometimes possible to identify the chromosomal location, at least roughly, of a few genes that have a large effect on the expression of a trait by associating variation in the trait with genotypic differences at loci that happen to be closely linked to those genes. A locus identified in this way is referred to as a *quantitative trait locus*, and the name given to the approach is QTL mapping.¹

The basic ideas behind QTL mapping are actually very simple, although the implementation of those ideas can be quite complex. In broad outline, this is the approach:

- Produce a set of progeny of known parentage. One common design involves first crossing a single pair of "inbred" parents that differ in expression of the quantitative trait of interest and then crossing the F_1 s, either among themselves to produce F_2 s (or recombinant inbred lines) or backcrossing them to one or both parents.
- Construct a linkage map for the molecular markers you're using. Ideally, you'll have a

¹These notes draw heavily on [64]

large enough number of markers to cover virtually every part of the genome.²

- Measure the phenotype and score the genotype at every marker locus of every individual in your progeny sample.
- Collate the data and analyze it in a computer package like **QTL Cartographer** to identify the position and effects of QTL associated with variation in the phenotypic trait you're interested in.

If that sounds like a lot of work, you're right. It is. But the results can be quite informative, because they allow you to say more about the genetic influences on the expression of the trait you're studying than a simple parent-offspring regression.

Thoday's Method³

Suppose there is a locus, Q , influencing the expression of a quantitative trait situated between two known marker loci, A and B .⁴ If we have inbred lines with different phenotypes, we can assume that one line has the genotype AQB/AQB and the other has the genotype aqb/aqb . The procedure for detecting the presence of Q is as follows:

1. Cross the inbred lines to form an F_1 . The genotype of all F_1 progeny will be AQB/aqb .
2. Intercross the F_1 's to form an F_2 and look at the progeny with recombinant genotypes, e.g., aB/ab .
3. If Q lies between A and B
 - (a) The phenotypes of progeny will fall into two distinct classes corresponding with the genotypes: aqB/aqb and aQB/aqb .⁵
 - (b) The recombination fraction between A and Q is related to the proportion of qq and Qq genotypes among the progeny.

²We'll talk a little later about how many markers are required.

³Primarily of historical interest, but it sets the stage for what is to follow.

⁴Of course, we don't know it's there when we start, but as we've done so many other times in this course, we'll assume that we know it's there and come back to how we find out where "there" is later.

⁵Actually there could be a third phenotypic class if there are two recombination events between a and b , i.e., aQB/aQb . Thoday's method assumes that the recombination fraction between A and B is small enough that double recombination events can be ignored, because if we don't ignore that possibility we must also admit that there will be some aqB/aQb genotypes that we can't distinguish from aQB/aqb genotypes.

Notice that in this last step we actually have a criterion for determining whether Q lies between A and B . Namely, if A and B are close enough in the linkage map that there is essentially no chance of double recombination between them, then we'll get the two phenotype classes referred to in recombinants between A and B . If Q lies outside this region,⁶ we'll get the two phenotype classes in 1:1 proportions and associated independently with genotype differences at the B locus.⁷

Genetic recombination and mapping functions

Genetic mapping is based on the idea that recombination is more likely between genes that are far apart on chromosomes than between genes that are close. If we have three genes A , B , and C arranged in that order on a chromosome, then

$$r_{AC} = r_{AB}(1 - r_{BC}) + (1 - r_{AB})r_{BC} \quad ,$$

where r_{AB} , r_{AC} , and r_{BC} are the recombination rates between A and B , A and C , and B and C , respectively.⁸

Haldane pointed out that this relationship implies another, namely that the probability that there are k recombination events between two loci m map units apart is given by the Poisson distribution:

$$p(m, k) = \frac{e^{-m} m^k}{k!} \quad .$$

Now to observe a recombination event between A and C requires that there be an odd number of recombination events between them (1, 3, 5, ...), i.e.,

$$\begin{aligned} r_{AC} &= \sum_{k=0}^{\infty} \frac{e^{-m} m^{(2k+1)}}{(2k+1)!} \\ &= \frac{1 - e^{-2m}}{2} \quad . \end{aligned}$$

⁶More specifically, if Q is not linked to B , or if Q has only a small effect on expression of the trait we're studying.

⁷This logic not only depends on ignoring the possibility of double crossovers, but also on assuming that individual loci influencing the quantitative trait have effects large enough that there will be two categories of offspring, corresponding to the difference between qq homozygotes and qQ heterozygotes.

⁸In practice this isn't quite true. Interference may cause the recombination rate between A and C to differ from this simple prediction. That's not much of a problem since we can just use a mapping function that corrects for this problem, but we'll ignore interference to keep things simple.

This leads to a natural definition of map units as

$$m = -\ln(1 - 2r)/2 \quad .$$

m calculated in this way gives the map distance in Morgans ($1M$). Map distances are more commonly expressed as centiMorgans, where $100cM = 1M$. Notice that when r is small, $r \approx m$, so the map distance in centiMorgans is approximately equal to the recombination frequency expressed as a percentage. There are several other mapping functions that can be chosen for an analysis. In particular, for analysis of real data investigators typically choose a mapping function that allows for interference in recombination. We don't have time to worry about those complications, so we'll use only the Haldane mapping function in our further discussions.

How many markers will you need?

If markers are randomly placed through the genome, then the average distance between a QTL and the closest marker is

$$E(m) = \frac{L}{2(n+1)} \quad ,$$

where L is the total map length and n is the number of markers employed. The upper 95% confidence limit for the distance is

$$\frac{L}{2} \left(1 - 0.05^{(1/n)}\right) \quad .$$

Since the human genome is 33M (3300cM), 110 random markers give an average distance of 14.9cM and an upper 95% confidence limit of 44.3cM, corresponding to recombination frequencies of 0.13 and 0.29, respectively. Since there are about 30,000 genes in the human genome, there are roughly 10 genes per centimorgan. So if your QTL is 44cm from the nearest marker, there are probably over 400 genes in the chromosomal segment you've identified.

If r_{MQ} is the recombination fraction between the nearest marker locus and the QTL of interest, the frequency of recombinant genotypes among F_2 progeny is $2r_{MQ}(1 - r_{MQ}) + r_{MQ}^2$. As you can see from the graph in Figure 23.1, there's a nearly linear relationship between recombination frequency and the frequency of recombinant phenotypes (p in the graph). Think about what that means. Having a really dense map with a lot of markers is great, because it will allow you to map your QTL very precisely, *if* you look at enough segregating

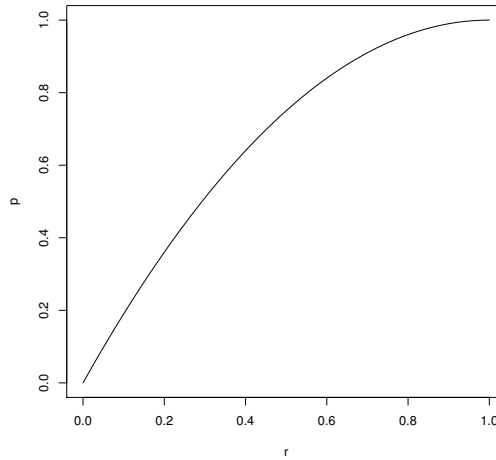


Figure 23.1: The relationship between recombination frequency, r , and the frequency of recombinant phenotypes, p , assuming a Haldane mapping function.

offspring to have a reasonable chance of picking up recombinants between them. With 3300cM in the human genome and roughly 3 GB of sequence to get within 1 MB of the actual QTL, you'd need one marker per centimorgan. To have 10 recombinants between markers bracketing the QTL, you'd need to analyze 1000 chromosomes.⁹

Analysis of an F_2 derived from inbred lines

An analysis of inbred lines uses the same basic design as Thoday, but takes advantage of more information.¹⁰ We start with two inbred lines M_1QM_2/M_1QM_2 and m_1qm_2/m_1qm_2 , make an F_1 , intercross them, and score the phenotype and marker genotype of each individual. Analysis of the data is based on calculating the frequency of each genotype at the Q locus as a function of the genotype at the marker loci and the recombination fractions between

⁹These calculations are moot in this context, of course. You can't ethically do a QTL study in humans. But the same principles apply to association mapping, which we'll get to in a couple of lectures.

¹⁰As I alluded to earlier, other breeding designs are possible, including backcrosses and recombinant inbred lines and analyses involving outbred parents. The principles are the same in every case, but the implementation is different.

the marker loci and Q .¹¹ For example,

$$\begin{aligned} P(M_1QM_2/M_1QM_2) &= ((1 - r_{1Q})(1 - r_{Q2})/2)^2 \\ P(M_1QM_2/M_1qM_2) &= 2((1 - r_{1Q})(1 - r_{Q2})/2)(r_{1Q}r_{Q2}/2) \\ P(M_1qM_2/M_1qM_2) &= (r_{1Q}r_{Q2}/2)^2 \quad . \end{aligned}$$

Because the frequency of $M_1M_2/M_1M_2 = ((1 - r_{12})/2)^2$, we can use Bayes' Theorem to write the conditional probabilities of getting each genotype as

$$\begin{aligned} P(QQ|M_1M_2/M_1M_2) &= \frac{(1-r_{1Q})^2(1-r_{Q2})^2}{(1-r_{12})^2} \\ P(Qq|M_1M_2/M_1M_2) &= \frac{2r_{1Q}r_{Q2}(1-r_{1Q})(1-r_{Q2})}{(1-r_{12})^2} \\ P(qq|M_1M_2/M_1M_2) &= \frac{r_{1Q}^2r_{Q2}^2}{(1-r_{12})^2} \quad . \end{aligned}$$

Clearly, if we wanted to we could right down similar expressions for the nine remaining marker genotype classes, but we'll stop here. You get the point.¹²

Now that we've got this we can write down the likelihood of getting our data, namely

$$L(x|M_j) = \sum_{k=1}^N \phi(x|\mu_{Q_k}, \sigma^2)P(Q_k|M_k) \quad ,$$

where N is the number of QTL genotypes considered, $\phi(x|\mu_{Q_k}, \sigma^2)$ is the probability of getting phenotype x given the mean phenotype, μ_{Q_k} , and variance, σ^2 , associated with Q_k , and $P(Q_k|M_k)$ is the probability of getting Q_k given the observed marker genotype. Fortunately, we don't have to do any of these calculations, all we do is to ask our good friend (QTL Cartographer) to do the calculations for us. It will scan the genome, and tell us how many QTL loci we are likely to have, where they are located relative to our known markers, and what the additive and dominance effects of the alleles are.

The Caveats

That's wonderful, isn't it? We have to do a little more work than for a traditional quantitative genetic analysis, i.e., we have to do a bunch of molecular genotyping in addition to all of the measurements we'd do for a quantitative genetic experiment anyway, but we now know how

¹¹You should be getting used to the idea now that we always assume we know something we don't and then backcalculate from what we do know to what we'd like to know.

¹²I should say, I *hope* you get the point.

how many genes are involved in the expression of our trait, where they are in the genetic map, and what their additive and dominance effects are. We can even tell something about how alleles at the different loci interact with one another. What more could you ask for? Well, there are a few things about QTL analyses to keep in mind.

- As currently implemented, QTL mapping procedures assume that the distribution of trait values around the genotype mean is normal, *with the same variance for all QTL genotypes*.¹³
- QTL mapping programs often estimate the effects of each locus individually. It's not at all easy to search simultaneously for the joint effects of two QTL loci, although it's not too hard to look at the combined effects of QTL loci first identified individually. Composite interval mapping, in which additional markers are included as cofactors in the analysis, partially addresses this limitation. Multiple interval mapping looks at several QTLs simultaneously and shows some promise, but as you may be able to imagine it's pretty hard to search for more than a few QTLs simultaneously.
- If some loci in the "high" line have "low" effects and vice versa, the effects of both loci (and possibly other loci) may be masked.
- Using this approach we can identify the QTL's that are important *in a particular cross*, but different crosses can identify different QTL's. Even the same cross may reveal different QTL's if the measurements are done in different environments. Methods to analyze several progeny sets simultaneously are only now being developed.

¹³I *know* you picked up on that when I said that the phenotypic variance associated with each QTL genotype was σ^2 . You were just too polite to point it out and interrupt me.

Chapter 24

Mapping Quantitative Trait Loci with R/qtl

There are two stages to making a QTL map for a particular trait (once you've scored tens or hundreds of marker loci in tens or hundreds of F_2 or backcross progeny):

1. Construct a genetic map of your markers.
2. Feed the genetic map, marker data, and phenotype data into QTL Cartographer and run the analysis.

Although constructing the genetic map of your markers is really important step, we're not going to talk about it further. It is, after all, simply an elaboration of classical Mendelian genetics.¹

The data²

We're going to focus on the second step. We'll be using data from Sugiyama et al., *Physiological Genomics* 10:512, 2002 as distributed from <http://www.rqtl.org/sug.csv> As you might have guessed from the extension on that file, the data is stored in a standard CSV file. Other formats are possible. They're described in the documentation, but we'll just deal with

¹Although it gets a *lot* more complicated when you're dealing with tens or hundreds of markers, and you don't even know which ones belong on which chromosomes!

²From here on out these notes depend heavily on the "shorter tour of R/qtl" at <http://www.rqtl.org/tutorials/>.

CSV files. To see what the data format looks like, open it up in your favorite spreadsheet program to take a look.

The data are from an intercross between two lines of inbred mice, BALB/cj and CBA/CaJ. Its format is fairly straightforward. After the header rows (lines 1-3), each line provides the data for one mouse. The first four columns contain phenotypic data: blood pressure, heart rate, body weight, and heart weight. The next two columns contain an indicator for the sex of the mouse³ and an individual ID. The remaining columns each correspond to markers (name on row 1, the chromosome on which they occur on row 2, and the map position on that chromosome on row 3). The two letters correspond to the alleles inherited from BALB/cj or CBA/CaJ, B and C respectively.⁴

Now it's time to get the data into R. To do so, type

```
sug <- read.cross(format="csv", file="sug.csv", genfile="sug.csv",
                  genotypes=c("CC", "CB", "BB"), alleles=c("C", "B"))
```

This reads data from `sug.csv` and puts it into the R object `sug`. Please note that R is case sensitive. If all goes well, you'll see this on your console after you hit return:

```
--Read the following data:
163 individuals
93 markers
6 phenotypes
--Cross type: f2
```

To get a sense of what's in the data simply type

```
summary(sug)
```

You should see

```
F2 intercross

No. individuals:      163

No. phenotypes:       6
Percent phenotyped: 95.1 95.7 99.4 99.4 100 100

No. chromosomes:     19
```

³All 1s. Only male mice are included in this data set.

⁴Since the parents were inbred lines, their genotypes were BB and CC.

```

Autosomes:      1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19

Total markers:  93
No. markers:    5 7 5 5 5 4 8 4 4 5 6 3 3 5 5 4 4 6 5
Percent genotyped: 98.3
Genotypes (%):  CC:23.9  CB:50.2  BB:26  not BB:0  not CC:0

```

You'll see that more than 95% of the 163 individuals have been phenotyped for each of the four traits and that more than 98% have been genotyped. You'll also see that all of the loci are autosomal. You can also `plot(sug)` to get a visual summary of the data (Figure 24.1).

The figure in the upper left shows which individuals (rows) are missing genotype information for particular markers (columns). There's so little missing genotype data in this data set, it doesn't show up. The next figure shows the location of markers on the genetic map, and the remainder summarize the distribution of phenotypes in the data set.⁵

The QTL analysis

Now we begin scanning the genome to locate QTLs. First, we have to calculate the probabilities of the three QTL genotypes in a grid across the genome. To do that we

```
sug <- calc.genoprob(sug, step=1)
```

Don't worry when you get nothing back except a command prompt. That's expected. What you've done is to calculate all of those genotype probabilities at a grid size of 1cM (`step=1`) across the genome and stored the results back into `sug`.

Now that we have the QTL genotype probabilities, we can run the QTL analysis and store the result

```
sug.em <- scanone(sug, pheno.col="bp")
```

Again, you'll just get the command prompt back.⁶ What you've done is to store the results in an object called `sug.em`.⁷ This analysis will locate only QTLs that influence blood pressure (`pheno.col='bp'`). If I'd wanted to analyze one of the other traits, I would have specified `pheno.col='hr'`, `pheno.col='bw'`, or `pheno.col='heart_wt'`. If I summarize the results to this point (`summary(sug.em)`), I'll get a report of the maximum LOD score on each chromosome.

⁵The histograms for sex and mouse_id obviously aren't very interesting

⁶Don't worry about the warning message. It's expected.

⁷I called it `sug.em` because `scanone()` is using the EM algorithm to obtain maximum-likelihood estimates of the parameters. Other algorithms are available, but we won't discuss them.

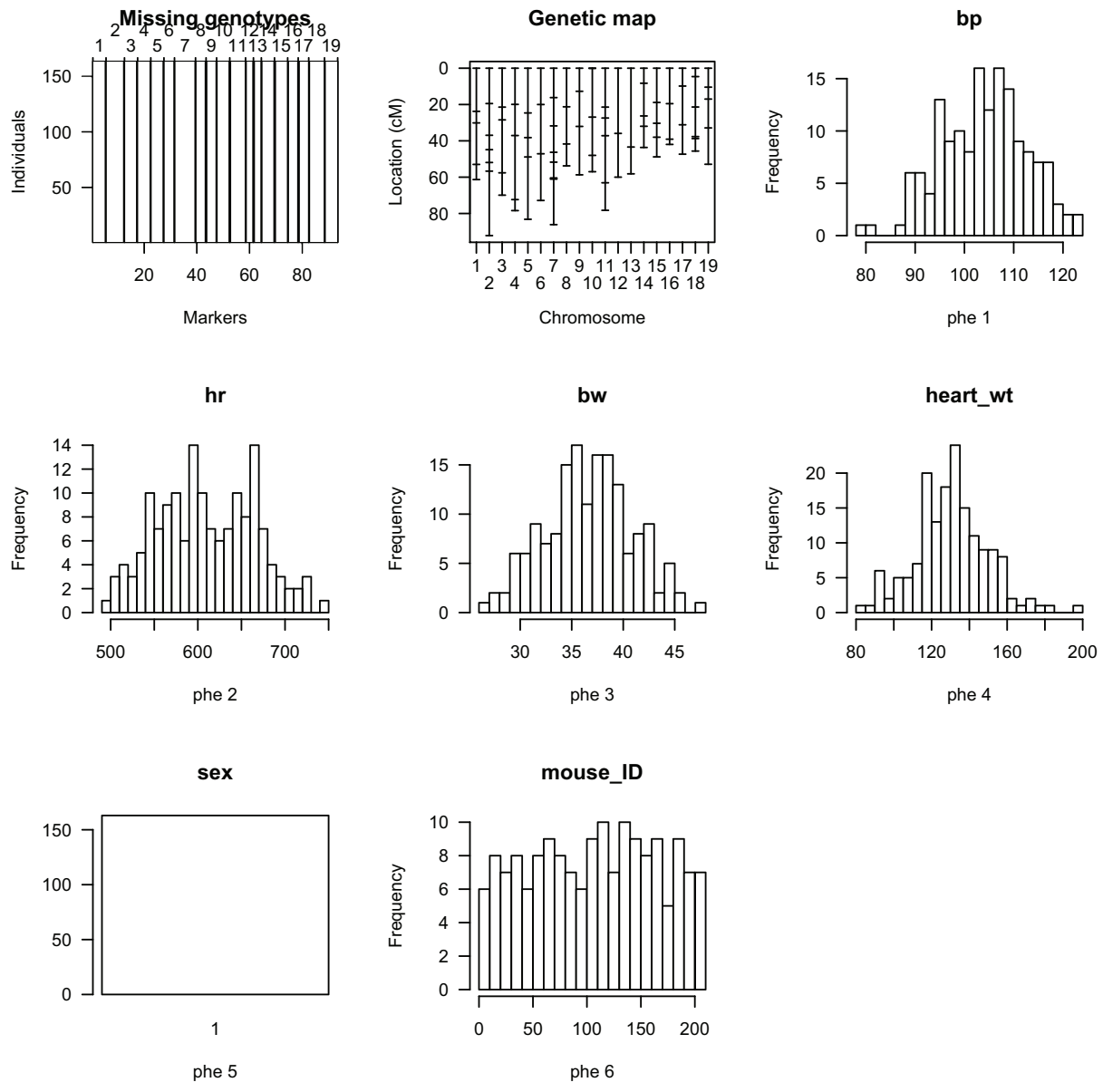


Figure 24.1: Results of `plot(sug)`

	chr	pos	lod
D1MIT36	1	76.73	1.449
c2.loc77	2	82.80	1.901
c3.loc42	3	52.82	1.393
c4.loc43	4	47.23	0.795
D5MIT223	5	86.57	1.312
c6.loc26	6	27.81	0.638
c7.loc45	7	47.71	6.109
c8.loc34	8	54.90	1.598
D9MIT71	9	27.07	0.769
c10.loc51	10	60.75	0.959
c11.loc34	11	38.70	2.157
D12MIT145	12	2.23	1.472
c13.loc20	13	27.26	1.119
D14MIT138	14	12.52	1.119
c15.loc8	15	11.96	5.257
c16.loc31	16	45.69	0.647
D17MIT16	17	17.98	1.241
D18MIT22	18	13.41	1.739
D19MIT71	19	56.28	0.402

To determine whether any of the markers are associated with the blood pressure phenotype more strongly than we would expect at random, we perform a permutation test.

```
sug.perm <- scanone(sug, pheno.col="bp", n.perm=5000)
```

You'll get a progress report as the permutations proceed, but be prepared to wait quite awhile. Each individual permutation reruns the entire `scanone()` analysis with phenotypes and genotypes randomized relative to one another. This gives us a distribution of LOD scores expected at random, and we'll use this to set a threshold that takes account of the multiple comparisons we make when we do separate likelihood-ratio tests at every potential QTL position in the genome.

With the permutations in hand, we can now summarize the results of the analysis and identify the position of QTLs for blood pressure (to a 1cM resolution).

```
summary(sug.em, perms=sug.perm, alpha=0.05, p.values=TRUE)
```

By specifying `alpha=0.05`, all peaks with a genome-adjusted p-value of less than 0.05 will be included in the summary. By specifying `p.values=TRUE` we ensure that only columns with genome-adjusted p-values are considered.⁸ The summary is very short and simple:

⁸Since we included all markers in our permutation test, this will simply include all columns.

```

      chr  pos  lod
c7.loc45   7 47.7 6.11
c15.loc8  15 12.0 5.26

```

It tells us that only two QTLs have a significant association with blood pressure, one on chromosome 7 at 47.7cM, the other on chromosome 15 at 12.0cM.

Finally, we visualize the effects of the QTLs on chromosome 7 and chromosome 15.

```

sug <- sim.geno(sug)
effectplot(sug, pheno.col="bp", mname1="7@47.7")
effectplot(sug, pheno.col="bp", mname2="15@12")
effectplot(sug, pheno.col="bp", mname1="7@47.7", mname2="15@12")
effectplot(sug, pheno.col="bp", mname1="15@12", mname2="7@47.7")

```

You'll see the results in Figure 24.2. The top two figures show the phenotypic means associated with markers on chromosome 7 and 15 respectively. The bottom two figures show how the phenotype depends on the genotype at both QTLs. The QTL on chromosome 15 (figure on the right) seems to have almost purely additive effects. The heterozygote is very close to intermediate between the two homozygotes. The QTL on chromosome 7, however, has substantial non-additive effects. Blood pressure of heterozygotes appears to be lower than that of either homozygote. The interaction plots suggests epistatic interactions between the loci. The lines aren't parallel.

We can get numerical estimates of the means and standard errors by changing those statements just a little.

```

print(effectplot(sug, pheno.col="bp", mname1="7@47.7", draw=FALSE))
$Means
D7MIT31.CC D7MIT31.CB D7MIT31.BB
 103.4679   101.3002   109.0165

$SEs
D7MIT31.CC D7MIT31.CB D7MIT31.BB
 1.4486284  0.9499457  1.0415715

print(effectplot(sug, pheno.col="bp", mname2="15@12", draw=FALSE))
$Means
D15MIT175.CC D15MIT175.CB D15MIT175.BB
 108.43902   104.70130   99.91892

```

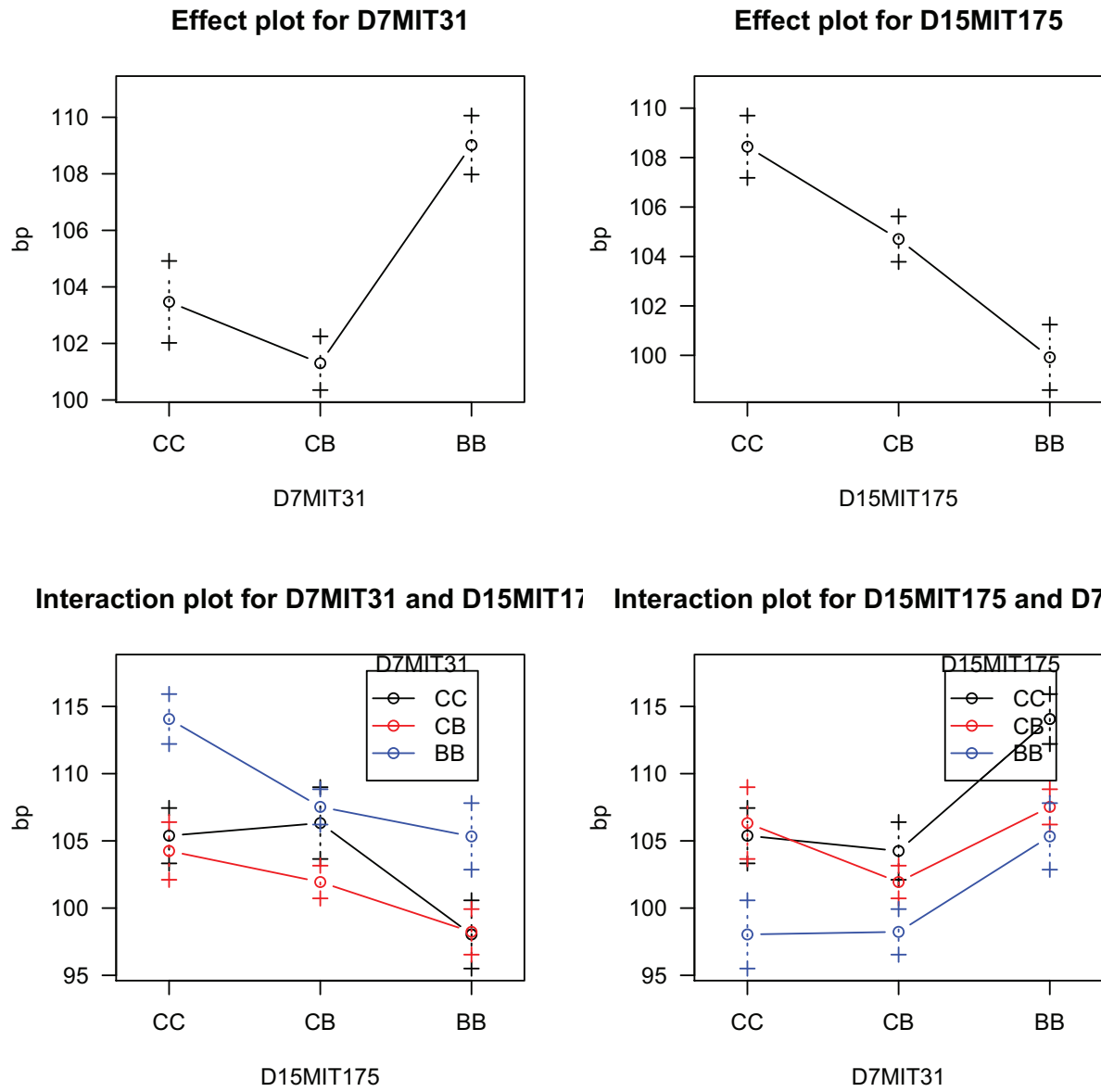


Figure 24.2: Effect plots for the QTLs on chromosome 7 and chromosome 15

\$SEs

```
D15MIT175.CC D15MIT175.CB D15MIT175.BB
      1.258112      0.918049      1.324373
```

We estimate additive and dominance effects associated with each marker from a linear regression

$$y_i = \beta_0 + \alpha_i a + \delta_i d$$

where $\alpha_i = (-1, 0, 1)$ and $\delta_i = (0, 1, 0)$ for genotypes CC, CB, and BB, respectively. With a and d estimated in this way we can specify the genotypic values as

CC	CB	BB
$\bar{x} - a$	$\bar{x} + d$	$\bar{x} + a$

Doing that regression may sound hard, but it's actually quite easy.

```
print(effectscan(sug, pheno.col="bp", draw=FALSE))
```

You'll get a very long table as a result. Here I just pull out the lines corresponding to the two QTLs we identified.

```
      chr  pos      a      d
D7MIT31   7 49.01  2.77491094 -4.950729964
D15MIT175 15  3.96 -4.26005274  0.522327047
```


Chapter 25

Association mapping: the background from two-locus genetics

QTL mapping is wonderful — provided that you're working with an organism where it's possible to design a breeding program and where the information derived from that breeding program is relevant to variation in natural populations. Think about it. If we do a QTL analysis based on segregation in an F_2 population derived from two inbred lines, all we really know is which loci are associated with phenotypic differences *between those two lines*. Typically what we really want to know, if we're evolutionary biologists, is which loci are associated with phenotypic differences *between individuals in the population we're studying*. That's where association mapping comes in. We look for statistical associations between phenotypes and genotypes across a whole population. We expect there to be such associations, if we have a dense enough map, because some of our marker loci will be closely linked to loci responsible for phenotypic variation.

A digression into two-locus population genetics¹

It's pretty obvious that if two loci are closely linked, alleles at those loci are likely to be closely linked, but let's take a closer look at exactly what that means.

One of the most important properties of a two-locus system is that it is no longer sufficient to talk about allele frequencies alone, even in a population that satisfies all of the assumptions necessary for genotypes to be in Hardy-Weinberg proportions at each locus. To see why

¹**Note:** We'll go over only a small part of this section in lecture. I'm providing all the details here so you can find them in the future if you ever need them.

consider this. With two loci and two alleles there are four possible gametes:²

Gamete	A_1B_1	A_1B_2	A_2B_1	A_2B_2
Frequency	x_{11}	x_{12}	x_{21}	x_{22}

If alleles are arranged randomly into gametes then,

$$\begin{aligned} x_{11} &= p_1p_2 \\ x_{12} &= p_1q_2 \\ x_{21} &= q_1p_2 \\ x_{22} &= q_1q_2 \quad , \end{aligned}$$

where $p_1 = \text{freq}(A_1)$ and $p_2 = \text{freq}(A_2)$. But alleles need not be arranged randomly into gametes. They may covary so that when a gamete contains A_1 it is more likely to contain B_1 than a randomly chosen gamete, or they may covary so that a gamete containing A_1 is less likely to contain B_1 than a randomly chosen gamete. This covariance could be the result of the two loci being in close physical association, but it doesn't have to be. Whenever the alleles covary within gametes

$$\begin{aligned} x_{11} &= p_1p_2 + D \\ x_{12} &= p_1q_2 - D \\ x_{21} &= q_1p_2 - D \\ x_{22} &= q_1q_2 + D \quad , \end{aligned}$$

where $D = x_{11}x_{22} - x_{12}x_{21}$ is known as the *gametic disequilibrium*.³ When $D \neq 0$ the alleles within gametes covary, and D measures *statistical* association between them. It does not (directly) measure the *physical* association. Similarly, $D = 0$ does not imply that the loci are unlinked, only that the alleles at the two loci are arranged into gametes independently of one another.

A little diversion

It probably isn't obvious why we can get away with only one D for all of the gamete frequencies. The short answer is:

²Think of drawing the Punnett square for a dihybrid cross, if you want.

³You will sometimes see D referred to as the linkage disequilibrium, but that's misleading. Alleles at different loci may be non-randomly associated even when they are not linked.

There are four gametes. That means we need three parameters to describe the four frequencies. p_1 and p_2 are two. D is the third.

Another way is to do a little algebra to verify that the definition is self-consistent.

$$\begin{aligned}
 D &= x_{11}x_{22} - x_{12}x_{21} \\
 &= (p_1p_2 + D)(q_1q_2 + D) - (p_1q_2 - D)(q_1p_2 - D) \\
 &= \left(p_1q_1p_2q_2 + D(p_1p_2 + q_1q_2) + D^2 \right) \\
 &\quad - \left(p_1q_1p_2q_2 - D(p_1q_2 + q_1p_2) + D^2 \right) \\
 &= D(p_1p_2 + q_1q_2 + p_1q_2 + q_1p_2) \\
 &= D(p_1(p_2 + q_2) + q_1(q_2 + p_2)) \\
 &= D(p_1 + q_1) \\
 &= D \quad .
 \end{aligned}$$

Transmission genetics with two loci

I'm going to construct a reduced version of a mating table to see how gamete frequencies change from one generation to the next. There are ten different two-locus genotypes (if we distinguish coupling, A_1B_1/A_2B_2 , from repulsion, A_1B_2/A_2B_1 , heterozygotes as we must for these purposes). So a full mating table would have 100 rows. If we assume all the conditions necessary for genotypes to be in Hardy-Weinberg proportions apply, however, we can get away with just calculating the frequency with which any one genotype will produce a particular gamete.⁴

Genotype	Frequency	Gametes			
		A_1B_1	A_1B_2	A_2B_1	A_2B_2
A_1B_1/A_1B_1	x_{11}^2	1	0	0	0
A_1B_1/A_1B_2	$2x_{11}x_{12}$	$\frac{1}{2}$	$\frac{1}{2}$	0	0
A_1B_1/A_2B_1	$2x_{11}x_{21}$	$\frac{1}{2}$	$0\frac{1}{2}$	0	
A_1B_1/A_2B_2	$2x_{11}x_{22}$	$\frac{1-r}{2}$	$\frac{r}{2}$	$\frac{r}{2}$	$\frac{1-r}{2}$
A_1B_2/A_1B_2	x_{12}^2	0	1	0	0
A_1B_2/A_2B_1	$2x_{12}x_{21}$	$\frac{r}{2}$	$\frac{1-r}{2}$	$\frac{1-r}{2}$	$\frac{r}{2}$
A_1B_2/A_2B_2	$2x_{12}x_{22}$	0	$\frac{1}{2}$	0	$\frac{1}{2}$
A_2B_1/A_2B_1	x_{21}^2	0	0	1	0
A_2B_1/A_2B_2	$2x_{21}x_{22}$	0	0	$\frac{1}{2}$	$\frac{1}{2}$
A_2B_2/A_2B_2	x_{22}^2	0	0	0	1

⁴We're assuming random union of *gametes* rather than random mating of *genotypes*.

Where do $\frac{1-r}{2}$ and $\frac{r}{2}$ come from?

Consider the coupling double heterozygote, A_1B_1/A_2B_2 . When recombination doesn't happen, A_1B_1 and A_2B_2 occur in equal frequency ($1/2$), and A_1B_2 and A_2B_1 don't occur at all. When recombination happens, the four possible gametes occur in equal frequency ($1/4$). So the recombination frequency,⁵ r , is half the crossover frequency,⁶ c , i.e., $r = c/2$. Now the results of crossing over can be expressed in this table:

Frequency	A_1B_1	A_1B_2	A_2B_1	A_2B_2
$1 - c$	$\frac{1}{2}$	0	0	$\frac{1}{2}$
c	$\frac{1}{4}$	$\frac{1}{4}$	$\frac{1}{4}$	$\frac{1}{4}$
Total	$\frac{2-c}{4}$	$\frac{c}{4}$	$\frac{c}{4}$	$\frac{2-c}{4}$
	$\frac{1-r}{2}$	$\frac{r}{2}$	$\frac{r}{2}$	$\frac{1-r}{2}$

Changes in gamete frequency

We can use the mating table table as we did earlier to calculate the frequency of each gamete in the next generation. Specifically,

$$\begin{aligned}
 x'_{11} &= x_{11}^2 + x_{11}x_{12} + x_{11}x_{21} + (1-r)x_{11}x_{22} + rx_{12}x_{21} \\
 &= x_{11}(x_{11} + x_{12} + x_{21} + x_{22}) - r(x_{11}x_{22} - x_{12}x_{21}) \\
 &= x_{11} - rD \\
 x'_{12} &= x_{12} + rD \\
 x'_{21} &= x_{21} + rD \\
 x'_{22} &= x_{22} - rD \quad .
 \end{aligned}$$

No changes in allele frequency

We can also calculate the frequencies of A_1 and B_1 after this whole process:

$$\begin{aligned}
 p'_1 &= x'_{11} + x'_{12} \\
 &= x_{11} - rD + x_{12} + rD \\
 &= x_{11} + x_{12} \\
 &= p_1 \\
 p'_2 &= p_2 \quad .
 \end{aligned}$$

⁵The frequency of recombinant gametes in double heterozygotes.

⁶The frequency of cytological crossover during meiosis.

Since each locus is subject to all of the conditions necessary for Hardy-Weinberg to apply at a single locus, allele frequencies don't change at either locus. Furthermore, genotype frequencies at each locus will be in Hardy-Weinberg proportions. But the two-locus gamete frequencies change from one generation to the next.

Changes in D

You can probably figure out that D will eventually become zero, and you can probably even guess that how quickly it becomes zero depends on how frequent recombination is. But I'd be astonished if you could guess exactly how rapidly D decays as a function of r . It takes a little more algebra, but we can say precisely how rapid the decay will be.

$$\begin{aligned}
 D' &= x'_{11}x'_{22} - x'_{12}x'_{21} \\
 &= (x_{11} - rD)(x_{22} - rD) - (x_{12} + rD)(x_{21} + rD) \\
 &= x_{11}x_{22} - rD(x_{11} + x_{12}) + r^2D^2 - (x_{12}x_{21} + rD(x_{12} + x_{21}) + r^2D^2) \\
 &= x_{11}x_{22} - x_{12}x_{21} - rD(x_{11} + x_{12} + x_{21} + x_{22}) \\
 &= D - rD \\
 &= D(1 - r)
 \end{aligned}$$

Notice that even if loci are unlinked, meaning that $r = 1/2$, D does not reach 0 immediately. That state is reached only asymptotically. The two-locus analogue of Hardy-Weinberg is that gamete frequencies will *eventually* be equal to the product of their constituent allele frequencies.

D in a finite population

In the absence of mutation, D will eventually decay to 0, although the course of that decay isn't as regular as what I've just shown [40]. If we allow recurrent mutation at both loci, however, where



then it can be shown [74] that the expected value of $D^2/p_1(1-p_1)p_2(1-p_2)$ is

$$\begin{aligned}
 \frac{E(D^2)}{E(p_1(1-p_1)p_2(1-p_2))} &= \frac{1}{3 + 4N_e(r + \mu_1 + \nu_1 + \mu_2 + \nu_2) - \frac{2}{(2.5 + N_e(r + \mu_1 + \nu_1 + \mu_2 + \nu_2) + N_e(\mu_1 + \nu_1 + \mu_2 + \nu_2))}} \\
 &\approx \frac{1}{3 + 4N_e r} \quad .
 \end{aligned}$$

Bottom line: In a finite population, we don't expect D to go to 0, and the magnitude of D^2 is inversely related to amount of recombination between the two loci. The less recombination there is between two loci, i.e., the smaller r is, the larger the value of D^2 we expect.

This has all been a long way⁷ of showing that our initial intuition is correct. If we can detect a statistical association between a marker locus and a phenotypic trait, it suggests that the marker locus and a locus influence expression of the trait are physically linked. So how do we detect such an association and why do I say that it *suggests* the loci are physically linked?⁸ Funny you should ask. That's the subject of the next set of notes.

⁷OK. You can say it. A *very* long way.

⁸Rather than concluding definitely that they *are* physically linked.

Chapter 26

Association mapping: BAMD

We’ve now seen that a naïve, locus-by-locus approach to identifying associations between marker loci and SNPs could be misleading, both because we have to correct for multiple comparisons¹ and, more importantly, because we need to account for the possibility that loci are statistically associated simply because there is genetic substructure within the sample. Stephens and Balding [86] outline one set of Bayesian approaches to dealing with both of these problems. We’ll focus on the problem of accounting for population structure, using the approach implemented in BAMD, an R package similar to R/qt1.

The statistical model

BAMD² uses a multiple regression approach to investigate the relationship between genotypes at a marker locus and phenotypes. Specifically, they use a “mixed-model” that allows the residual variances and covariances to be specified in ways that reflect the underlying population structure. Suppose y_i is the phenotype of the i th individual in our sample and $\mathbf{y} = (y_1, \dots, y_I)$. Then the statistical model is:³

$$\mathbf{y} = \mathbf{X}\beta + \mathbf{Z}\gamma + \epsilon \quad ,$$

where \mathbf{X} is a matrix describing how each individual is assigned to a particular genetic grouping,⁴ β is a vector of coefficients describing the mean phenotype associated with individuals

¹Strictly speaking, we didn’t see this in the context of association mapping, but we encountered it in our discussion of QTL mapping.

²And a couple of other packages we won’t discuss, TASSEL and EMMAX.

³Hang on. This looks pretty complicated, but it’s really not as bad as it looks.

⁴For example, you could use STRUCTURE to identify genetic groupings in your data. Then row i of \mathbf{X} would correspond to the posterior probability that individual i is assigned to each of the groupings you

belonging to that grouping, \mathbf{Z} is a matrix in which element ij is the genotype of individual i at locus j ,⁵ γ is a vector of coefficients describing the effect of different genotypes at each locus,⁶ and ϵ is a vector of residuals.

In a typical regression problem, we'd assume $\epsilon \sim N(0, \sigma^2 \mathbf{I})$, where \mathbf{I} is the identity matrix. Translating that to English,⁷ we'd typically assume that the errors are independently distributed with a mean of 0 and a variance of σ^2 . In some applications, that's not a good assumption. Some of the individuals included in the analysis are related to one another. Fortunately, if you know (or can estimate) that degree of relationship, **BAMD** can help you out. If \mathbf{R} is a matrix in which element ij indicates the degree of relationship between individual i and j ,⁸ then we simply⁹ let $\epsilon \sim N(0, \sigma^2 \mathbf{R})$. Now we allow the residual errors to be correlated when individuals are related and to be uncorrelated when they are not.

There's only one more piece of the model that you need to understand in order to interpret the output. If I tell you that **BAMD** is an acronym for **B**ayesian **A**ssociation with **M**issing **D**ata, you can probably guess that the last piece has something to do with prior distributions. Here's what you need to know. We will, obviously, have to place prior distributions on β , γ , and σ^2 . We don't need to talk much about the priors on β or σ^2 . We simply assume $\beta_j \sim \text{uniform}$, and we use a standard prior for variance parameters.¹⁰ The prior for γ is, however, a bit more complicated.

The covariates in \mathbf{X} reflect aspects of the experimental design, even if the elements of \mathbf{X} are inferred from a **STRUCTURE** analysis.¹¹ They are, to some degree at least, imposed by how we collected our samples of individuals. In contrast, the covariates reflected in \mathbf{Z} represent genotypes selected at random from within those groups. Moreover, the set of marker loci we chose isn't the only possible set we could have chosen. As a result we have to think of both the genotypes we chose and the coefficients associated with them as being samples from some underlying distribution.¹² Specifically, we assume $\gamma_k \sim N(0, \sigma^2 \phi^2)$, where ϕ^2 is simply

identify.

⁵**BAMD** is intended for the analysis of SNP data. Thus, the genotypes can be scored as 1, 2, or 3. Which homozygote is associated with genotype 1 doesn't affect the results, only the sign of the associated coefficient.

⁶These are the coefficients we're really interested in. They tell us the magnitude of the affect associated with a particular locus. In the implementation we're using, the relationship between genotype and phenotype is assumed to be strictly additive, since heterozygotes are perfectly intermediate.

⁷Or at least translating it to something *closer* to English.

⁸Individuals are perfectly related to themselves, so $r_{ii} = 1$. Unrelated individuals have $r_{ij} = 0$.

⁹It's simple because the authors of **BAMD** included this possibility in their code. All you have to do is to specify \mathbf{R} . **BAMD** will take care of the rest.

¹⁰If you must know, we use $1/\sigma^2 \sim G(a, b)$, where G stands for the Gamma distribution and a and b are its parameters.

¹¹Some people like to call these "fixed" effects.

¹²People who like to refer to \mathbf{X} as fixed effects like to refer to these as "random" effects.

a positive constant that “adjusts” the variance of γ_k relative to the residual variance. Then we just put a standard prior on ϕ^2 .¹³

The good news is that once you’ve got your data into the right format, BAMD will take care of all of the calculations for you. It will give you samples from the posterior distribution of β , γ , σ^2 , and ϕ^2 , from which you can derive the posterior mean, the posterior standard deviation, and the credible intervals.

What about the “Missing Data” part of the name?

There’s one more thing that BAMD does for us behind the scenes. In any real association analysis data set, every individual is likely to be missing data at one or more loci. That’s a problem. If we’re doing a multiple regression, we can’t include sample points where there are missing data, but if we dropped every individual for which we couldn’t score one or more SNPs, we wouldn’t have any data left. So what do we do? We “impute” the missing data, i.e., we use the data we do have to guess what the data would have been if we’d been able to observe it. BAMD does this in a very sophisticated and reliable way. As a result, we’re able to include every individual in our analysis and make use of all the data we’ve collected.¹⁴

¹³You may be able to guess, if you’ve been reading footnotes, that we use $1/\phi^2 \sim G(c, d)$.

¹⁴If you’re interested in why we can get away with what seems like making up data, stop by and talk to me. It involves a lot more statistics than I want to get into here.

Part V

Molecular evolution

Chapter 27

Introduction to molecular population genetics

The study of evolutionary biology is commonly divided into two components: study of the *processes* by which evolutionary change occurs and study of the *patterns* produced by those processes. By “pattern” we mean primarily the pattern of phylogenetic relationships among species or genes.¹ Studies of evolutionary processes often don’t often devote too much attention to evolutionary patterns, except insofar as it is often necessary to take account of evolutionary history in determining whether or not a particular feature is an adaptation. Similarly, studies of evolutionary pattern sometimes try not to use any knowledge of evolutionary processes to improve their guesses about phylogenetic relationships, because the relationship between process and pattern can be tenuous.² Those who take this approach argue that invoking a particular evolutionary process seems often to be a way of making sure that you get the pattern you want to get from the data.

Or at least that’s the way it was in evolutionary biology when evolutionary biologists were concerned primarily with the evolution of morphological, behavioral, and physiological traits and when systematists used primarily anatomical, morphological, and chemical features (but not proteins or DNA) to describe evolutionary patterns. With the advent of molecular biology after the Second World War and its application to an increasing diversity of organisms in the late 1950s and early 1960s, that began to change. Goodman [31] used

¹In certain cases it may make sense to talk about a phylogeny of populations within species, but in many cases it doesn’t. We’ll discuss this further when we get to phylogeography in a couple of weeks.

²One way of justifying a strict parsimony approach to cladistics is by arguing (a) that by minimizing character state changes on a tree you’re merely trying to find a pattern of character changes as consistent as possible with the data you’ve gathered and (b) that evolutionary processes should be invoked only to explain that pattern, not to construct it.

the degree of immunological cross-reactivity between serum proteins as an indication of the evolutionary distance among primates. Zuckerkandl and Pauling [105] proposed that after species diverged, their proteins diverged according to a “molecular clock,” suggesting one that molecular similarities could be used to reconstruct evolutionary history. In 1966, Harris [35] and Lewontin and Hubby [45, 60] showed that human populations and populations of *Drosophila pseudoobscura* respectively, contained surprising amounts of genetic diversity.

In this course, we’ll focus on advances made in understanding the processes of molecular evolution and pay relatively little attention to the ways in which inferences about evolutionary patterns can be made from molecular data. Up to this point in the course we’ve completely ignored evolutionary pattern. As you’ll see in what follows, however, any discussion of molecular evolution, even if it focuses on understanding the processes, cannot avoid some careful attention to the pattern.

Types of data

Before we delve any further into our study of molecular evolution, it’s probably useful to back up a bit and talk a bit about the types of data that are available to molecular evolutionists. We’ve already encountered a couple of these (microsatellites and SNPs), but there are a variety of important categories into which we can group data used for molecular evolutionary analyses. Even though studies of molecular evolution in the last 10-15 years have focused on data derived from DNA sequence or copy number variation, modern applications of molecular markers evolved from earlier applications. Those markers had their limitations, but analyses of them also laid the groundwork for most or all of what’s going on in analyses of molecular evolution today. Thus, it’s useful to remind everyone what those groups are and to agree on some terminology for the ones we’ll say something about. Let’s talk first about the physical basis of the underlying data. Then we’ll talk about the laboratory methods used to reveal variation.

The physical basis of molecular variation

With the exception of RNA viruses, the hereditary information in all organisms is carried in DNA. Ultimately, differences in any of the molecular markers we study (and of genetically-based morphological, behavioral, or physiological traits) is associated with some difference in the physical structure of DNA, and molecular evolutionists study a variety of its aspects.

Nucleotide sequence A difference in nucleotide sequence is the most obvious way in which two homologous stretches of DNA may differ. The differences may be in translated

portions of protein genes (exons), portions of protein genes that are transcribed but not translated (e.g., introns, 5' or 3' untranslated regions), non-transcribed functional regions (e.g., promoters), or regions without apparent function.

Protein sequence Because of redundancy in the genetic code, a difference in nucleotide sequence at a protein-coding locus may or may not result in proteins with a different amino acid sequence. **Important note:** Don't forget that some loci code for RNA that has an immediate function without being translated to a protein, e.g., ribosomal RNA and various small nuclear RNAs.

Secondary, tertiary, and quaternary structure Differences in amino acid sequence may or may not lead to a different distribution of α -helices and β -sheets, to a different three-dimensional structure, or to different multisubunit combinations.

Imprinting At certain loci in some organisms the expression pattern of a particular allele depends on whether that allele was inherited from the individual's father or its mother.

Expression Functional differences among individuals may arise because of differences in the patterns of gene expression, even if there are no differences in the primary sequences of the genes that are expressed.³

Sequence organization Particular genes may differ between organisms because of differences in the position and number of introns. At the whole genome level, there may be differences in the amount and kind of repetitive sequences, in the amount and type of sequences derived from transposable elements, in the relative proportion of G-C relative to A-T, or even in the identity and arrangement of genes that are present. In microbial species, only a subset of genes are present in all strains. For example, in *Streptococcus pneumoniae* the "core genome" contains only 73% of the loci present in one fully sequenced reference strain [71]. Similarly, a survey of 20 strains of *Escherichia coli* and one of *E. fergusonii*, *E. coli*'s closest relative, identified only 2000 homologous loci that were present in all strains out of 18,000 orthologous loci identified [92]

Copy number variation Even within diploid genomes, there may be substantial differences in the number of copies of particular genes. In humans, for example, 76 copy-number polymorphisms (CNPs) were identified in a sample of only 20 individuals, and individuals differed from one another by an average of 11 CNPs. [81].

³Of course, differences in expression must ultimately be the result of a DNA sequence (or at least a methylation difference) difference somewhere, e.g., in a promoter sequence or the locus encoding a promoter or repressor protein, if it is a genetic difference or the result of an epigenetic modification of the sequence, e.g., by methylation.

It is worth remembering that in nearly all eukaryotes there are two different genomes whose characteristics may be analyzed: the nuclear genome and the mitochondrial genome. In plants there is a third: the chloroplast genome. In some protists, there may be even more, because of secondary or tertiary endosymbiosis. The mitochondrial and chloroplast genomes are typically inherited only through the maternal line, although some instances of biparental inheritance are known.

Revealing molecular variation

The diversity of laboratory techniques used to reveal molecular variation is even greater than the diversity of underlying physical structures. Various techniques involving direct measurement of aspects of DNA sequence variation are by far the most common today, so I'll mention only the techniques that have been most widely used.

Immunological distance Some molecules, notably protein molecules, induce an immune response in common laboratory mammals. The extent of cross-reactivity between an antigen raised to humans and chimps, for example, can be used as a measure of evolutionary distance. The immunological distance between humans and chimps is smaller than it is between humans and orangutans, suggesting that humans and chimps share a more recent common ancestor.

DNA-DNA hybridization Once repetitive sequences of DNA have been “subtracted out”,⁴ the rate and temperature at which DNA species from two different species anneal reflects the average percent sequence divergence between them. The percent sequence divergence can be used as a measure of evolutionary distance. Immunological distances and DNA-DNA hybridization were once widely used to identify phylogenetic relationships among species. Neither is now widely used in molecular evolution studies.

Isozymes Biochemists recognized in the late 1950s that many soluble enzymes occurred in multiple forms within a single individual. Population geneticists, notably Hubby and Lewontin, later recognized that in many cases, these different forms corresponded to different alleles at a single locus, *allozymes*. Allozymes are relatively easy to score in most macroscopic organisms, they are typically co-dominant (the allelic composition of heterozygotes can be inferred), and they allow investigators to identify both variable and non-variable loci.⁵ Patterns of variation at allozyme loci may not be representative

⁴See below for a description of some of these repetitive sequences.

⁵Classical Mendelian genetics, and quantitative genetics too for that matter, depend on genetic variation in traits to identify the presence of a gene.

of genetic variation that does not result from differences in protein structure or that are related to variation in proteins that are insoluble.

RFLPs In the 1970s molecular geneticists discovered restriction enzymes, enzymes that cleave DNA at specific 4, 5, or 6 base pair sequences, the *recognition site*. A single nucleotide change in a recognition site is usually enough to eliminate it. Thus, presence or absence of a restriction site at a particular position in a genome provides compelling evidence of an underlying difference in nucleotide sequence at that position.

RAPDs, AFLPs, ISSRs With the advent of the polymerase chain reaction in the late 1980s, several related techniques for the rapid assessment of genetic variation in organisms for which little or no prior genetic information was available. These methods differ in details of how the laboratory procedures are performed, but they are similar in that they (a) use PCR to amplify anonymous stretches of DNA, (b) generally produce larger amounts of variation than allozyme analyses of the same taxa, and (c) are bi-allelic, dominant markers. They have the advantage, relative to allozymes, that they sample more or less randomly through the genome. They have the disadvantage that heterozygotes cannot be distinguished from dominant homozygotes, meaning that it is difficult to use them to obtain information about levels of within population inbreeding.⁶

Microsatellites Satellite DNA, highly repetitive DNA associated with heterochromatin, had been known since biochemists first began to characterize the large-scale structure of genomes in DNA-DNA hybridization studies. In the mid-late 1980s several investigators identified smaller repetitive units dispersed throughout many genomes. Microsatellites, which consist of short (2-6) nucleotide sequences repeated many times, have proven particularly useful for analyses of variation within populations since the mid-1990s. Because of high mutation rates at each locus, they commonly have many alleles. Moreover, they are typically co-dominant, making them more generally useful than dominant markers. Identifying variable microsatellite loci is more laborious than identifying AFLPs, RAPDs, or ISSRs.

Nucleotide sequence The advent of automated sequencing has greatly increased the amount of population-level data available on nucleotide sequences. Nucleotide sequence data has an important advantage over most of the types of data discussed so

⁶To be fair, it is possible to distinguish heterozygotes from homozygotes with AFLPs, if you are **very** careful with your PCR technique [49]. That being said, few people are careful enough with their PCR to be able to score AFLPs reliably as codominant markers, and I am unaware of anyone who has done so outside of a controlled breeding program.

far: allozymes, RFLPs, AFLPs, RAPDs, and ISSRs may all hide variation. Nucleotide sequence differences need not be reflected in any of those markers. On the other hand, each of those markers provides information on variation at several or many, independently inherited loci. Nucleotide sequence information reveals differences at a location that rarely extends more than 2-3kb. Of course, as next generation sequencing techniques become less expensive and more widely available, we will see more and more examples of nucleotide sequence variation from many loci within individuals.

Single nucleotide polymorphisms In organisms that are genetically well-characterized it may be possible to identify a large number of single nucleotide positions that harbor polymorphisms. These SNPs potentially provide high-resolution insight into patterns of variation within the genome. For example, the HapMap project has identified approximately 3.2M SNPs in the human genome, or about one every kb [15].

As you can see from these brief descriptions, each of the markers reveals different aspects of underlying hereditary differences among individuals, populations, or species. There is no single “best” marker for evolutionary analyses. Which is best depends on the question you are asking. In many cases in molecular evolution, the interest is intrinsically in the evolution of the molecule itself, so the choice is based not on what those molecules reveal about the organism that contains them but on what questions about which molecules are the most interesting.

Divergence of nucleotide sequences

Underlying much of what we’re going to discuss in this last part of the course is the idea that we should be able to describe the degree of difference between nucleotide sequences, proteins, or anything else as a result of some underlying evolutionary processes. To illustrate the principle, let’s start with nucleotide sequences and develop a fairly simple model that describes how they become different over time.⁷

Let q_t be the probability that two homologous nucleotides are identical after having been evolving for t generations independently since the gene in which they were found was replicated in their common ancestor. Let λ be the probability of a substitution occurring at this nucleotide position in either of the two genes during a small time interval, Δt . Then

$$q_{t+\Delta t} = (1 - \lambda\Delta t)^2 q_t + 2(1 - \lambda\Delta t) \left(\frac{1}{3}\lambda\Delta t\right) (1 - q_t) + o(\Delta t^2)$$

⁷By now you should realize that when I write that somethin is “fairly simple”, I mean that it’s fairly simple to someone who’s comfortable with mathematics.

$$\begin{aligned}
&= (1 - 2\lambda\Delta t)q_t + \left(\frac{2}{3}\lambda\Delta t\right)(1 - q_t) + o(\Delta t^2) \\
q_{t+\Delta t} - q_t &= \frac{2}{3}\lambda\Delta t - \frac{8}{3}\lambda\Delta tq_t + o(\Delta t^2) \\
\frac{q_{t+\Delta t} - q_t}{\Delta t} &= \frac{2}{3}\lambda - \frac{8}{3}\lambda q_t + o(\Delta t) \\
\lim_{\Delta t \rightarrow 0} \frac{q_{t+\Delta t} - q_t}{\Delta t} = \frac{dq_t}{dt} &= \frac{2}{3}\lambda - \frac{8}{3}\lambda q_t \\
q_t &= 1 - \frac{3}{4}\left(1 - e^{-8\lambda t/3}\right)
\end{aligned}$$

The expected number of nucleotide substitutions separating the two sequences at any one position since they diverged is $d = 2\lambda t$.⁸ Thus,

$$\begin{aligned}
q_t &= 1 - \frac{3}{4}\left(1 - e^{-4d/3}\right) \\
d &= -\frac{3}{4}\ln\left[1 - \frac{4}{3}(1 - q_t)\right]
\end{aligned}$$

This is the simplest model of nucleotide substitution possible—the Jukes-Cantor model. It assumes

- that mutations are equally likely at all positions and
- that mutation among all nucleotides is equally likely.

Let's examine the second of those assumptions first. Observed differences between nucleotide sequences shows that some types of substitutions, i.e., transitions ($A \iff G$, $T \iff C$), occur much more frequently than others, i.e., transversions ($A \iff G$, $A \iff C$, $T \iff A$, $T \iff G$). There are a variety of different substitution models corresponding to different assumed patterns of mutation: Kimura 2 parameter (K2P), Felsenstein 1984 (F84), Hasegawa-Kishino-Yano 1985 (HKY85), Tamura and Nei (TrN), and generalized time-reversible (GTR). The GTR is, as its name suggests, the most general

⁸The factor 2 is there because λt substitutions are expected on each branch. In fact you will usually see the equation for q_t written as $q_t = 1 - (3/4)(1 - e^{-4\alpha t/3})$, where $\alpha = 2\lambda$. α is also referred to as the substitution rate, but it refers to the rate of substitution between the two sequences, not to the rate of substitution between each sequence and their common ancestor. If mutations are neutral λ equals the mutation rate, while α equals twice the mutation rate.

time-reversible model. It allows substitution rates to differ between each pair of nucleotides. That's why it's general. It still requires, however, that the substitution rate be the same in both directions. That's what it means to say that it's time reversible. While it is possible to construct a model in which the substitution rate differs depending on the direction of substitution, it leads to something of a paradox: with non-reversible substitution models the distance between two sequences A and B depends on whether we measure the distance from A to B or from B to A .

There are two ways in which the rate of nucleotide substitution can be allowed to vary from position to position—the phenomenon of among-site rate variation. First, we expect the rate of substitution to depend on codon position in protein-coding genes. The sequence can be divided into first, second, and third codon positions and rates calculated separately for each of those positions. Second, we can assume *a priori* that there is a distribution of different rates possible and that this distribution is described by one of the standard distributions from probability theory. We then imagine that the substitution rate at any given site is determined by a random draw from the given probability distribution. The gamma distribution is widely used to describe the pattern of among-site rate variation, because it can approximate a wide variety of different distributions (Figure 27.1).⁹

The mean substitution rate in each curve above is 0.1. The curves differ only in the value of a parameter, α , called the “shape parameter.” The shape parameter gives a nice numerical description of how much rate variation there is, except that it's backwards. The larger the parameter, the less among-site rate variation there is.

⁹And, to be honest, because it is mathematically convenient to work with.

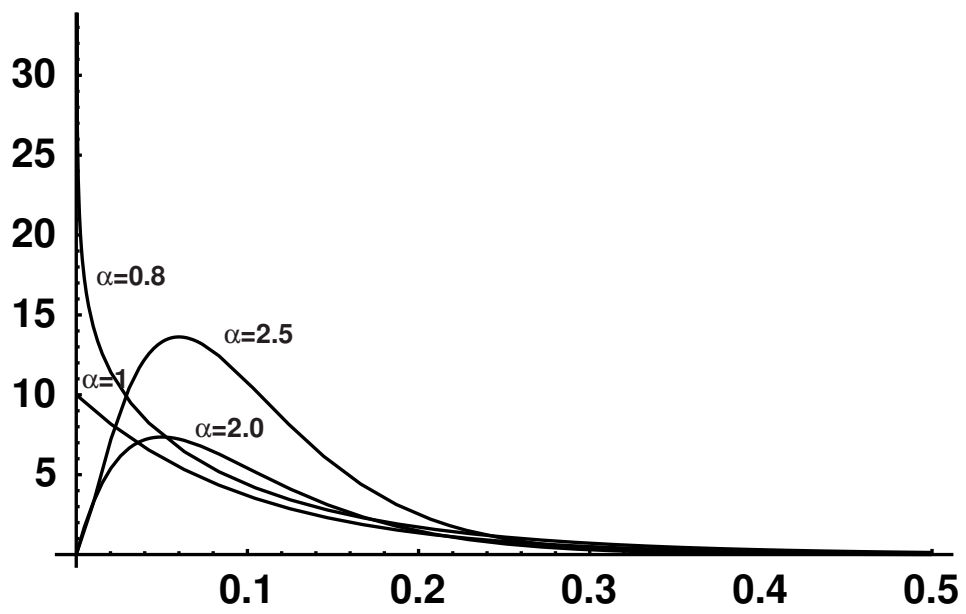


Figure 27.1: Examples of a gamma distribution.

Chapter 28

The neutral theory of molecular evolution

I didn't make a big deal of it in what we just went over, but in deriving the Jukes-Cantor equation I used the phrase "substitution rate" instead of the phrase "mutation rate." As a preface to what is about to follow, let me explain the difference.

- *Mutation rate* refers to the rate at which changes are incorporated into a nucleotide sequence during the process of replication, i.e., the probability that an allele differs from the copy of that in its parent from which it was derived. *Mutation rate* refers to the rate at which mutations arise.
- An allele substitution occurs when a newly arisen allele is incorporated into a population, e.g., when a newly arisen allele becomes fixed in a population. *Substitution rate* refers to the rate at which allele substitutions occur.

Mutation rates and substitution rates are obviously related related—substitutions can't happen unless mutations occur, after all—, but it's important to remember that they refer to different processes.

Early empirical observations

By the early 1960s amino acid sequences of hemoglobins and cytochrome *c* for many mammals had been determined. When the sequences were compared, investigators began to notice that the number of amino acid differences between different pairs of mammals seemed to be roughly proportional to the time since they had diverged from one another, as inferred

from the fossil record. Zuckerkandl and Pauling [105] proposed the *molecular clock hypothesis* to explain these results. Specifically, they proposed that there was a constant rate of amino acid substitution over time. Sarich and Wilson [78, 100] used the molecular clock hypothesis to propose that humans and apes diverged approximately 5 million years ago. While that proposal may not seem particularly controversial now, it generated enormous controversy at the time, because at the time many paleoanthropologists interpreted the evidence to indicate humans diverged from apes as much as 30 million years ago.

One year after Zuckerkandl and Pauling's paper, Harris [35] and Hubby and Lewontin [45, 60] showed that protein electrophoresis could be used to reveal surprising amounts of genetic variability within populations. Harris studied 10 loci in human populations, found three of them to be polymorphic, and identified one locus with three alleles. Hubby and Lewontin studied 18 loci in *Drosophila pseudoobscura*, found seven to be polymorphic, and five that had three or more alleles.

Both sets of observations posed real challenges for evolutionary geneticists. It was difficult to imagine an evolutionary mechanism that could produce a constant rate of substitution. It was similarly difficult to imagine that natural selection could maintain so much polymorphism within populations. The “cost of selection,” as Haldane called it would simply be too high.

Neutral substitutions and neutral variation

Kimura [50] and King and Jukes [51] proposed a way to solve both empirical problems. If the vast majority of amino acid substitutions are selectively neutral, then substitutions will occur at approximately a constant rate (assuming that mutation rates don't vary over time) and it will be easy to maintain lots of polymorphism within populations because there will be no cost of selection. I'll develop both of those points in a bit more detail in just a moment, but let me first be precise about what the neutral theory of molecular evolution actually proposes. More specifically, let me first be precise about what it does *not* propose. I'll do so specifically in the context of protein evolution for now, although we'll broaden the scope later.

- *The neutral theory asserts that alternative alleles at variable protein loci are selectively neutral.* This does *not* mean that the locus is unimportant, only that the alternative alleles found at this locus are selectively neutral.
 - Glucose-phosphate isomerase is an essential enzyme. It catalyzes the first step of glycolysis, the conversion of glucose-6-phosphate into fructose-6-phosphate.

- Natural populations of many, perhaps most, populations of plants and animals are polymorphic at this locus, i.e., they have two or more alleles with different amino acid sequences.
 - The neutral theory asserts that the alternative alleles are selectively neutral.
- By *selectively neutral* we do *not* mean that the alternative alleles have no effect on physiology or fitness. We mean that the selection among different genotypes at this locus is sufficiently weak that the pattern of variation is determined by the interaction of mutation, drift, mating system, and migration. This is roughly equivalent to saying that $N_e s < 1$, where N_e is the effective population size and s is the selection coefficient on alleles at this locus.
 - Experiments in *Colias* butterflies, and other organisms have shown that different electrophoretic variants of GPI have different enzymatic capabilities and different thermal stabilities. In some cases, these differences have been related to differences in individual performance.
 - If populations of *Colias* are large and the differences in fitness associated with differences in genotype are large, i.e., if $N_e s > 1$, then selection plays a predominant role in determining patterns of diversity at this locus, i.e., the neutral theory of molecular evolution would not apply.
 - If populations of *Colias* are small or the differences in fitness associated with differences in genotype are small, or both, then drift plays a predominant role in determining patterns of diversity at this locus, i.e., the neutral theory of molecular evolution applies.

In short, the neutral theory of molecular really asserts only that observed amino acid substitutions and polymorphisms are *effectively* neutral, not that the loci involved are unimportant or that allelic differences at those loci have no effect on fitness.

The rate of molecular evolution

We're now going to calculate the rate of molecular evolution, i.e., the rate of allelic substitution, under the hypothesis that mutations are selectively neutral. To get that rate we need two things: the rate at which new mutations occur and the probability with which new mutations are fixed. In a word equation

$$\begin{aligned} \# \text{ of substitutions/generation} &= (\# \text{ of mutations/generation}) \times (\text{probability of fixation}) \\ \lambda &= \mu_0 p_0 \quad . \end{aligned}$$

Surprisingly,¹ it's pretty easy to calculate both μ_0 and p_0 from first principles.

In a diploid population of size N , there are $2N$ gametes. The probability that any one of them mutates is just the mutation rate, μ , so

$$\mu_0 = 2N\mu \quad . \quad (28.1)$$

To calculate the probability of fixation, we have to say something about the dynamics of alleles in populations. Let's suppose that we're dealing with a single population, to keep things simple. Now, you have to remember a little of what you learned about the properties of genetic drift. If the current frequency of an allele is p_0 , what's the probability that it is eventually fixed? p_0 . When a new mutation occurs there's only one copy of it,² so the frequency of a newly arisen mutation is $1/2N$ and

$$p_0 = \frac{1}{2N} \quad . \quad (28.2)$$

Putting (28.1) and (28.2) together we find

$$\begin{aligned} \lambda &= \mu_0 p_0 \\ &= (2N\mu) \left(\frac{1}{2N} \right) \\ &= \mu \quad . \end{aligned}$$

In other words, if mutations are selectively neutral, the substitution rate is equal to the mutation rate. Since mutation rates are (mostly) governed by physical factors that remain relatively constant, mutation rates should remain constant, implying that substitution rates should remain constant if substitutions are selectively neutral. In short, if mutations are selectively neutral, we expect a molecular clock.

Diversity in populations

Protein-coding genes consist of hundreds or thousands of nucleotides, each of which could mutate to one of three other nucleotides.³ That's not an infinite number of possibilities, but it's pretty large.⁴ It suggests that we could treat every mutation that occurs as if it

¹Or perhaps not.

²By definition. It's new.

³Why three when there are four nucleotides? Because if the nucleotide at a certain position is an A, for example, it can only *change* to a C, G, or T.

⁴If a protein consists of 400 amino acids, that's 1200 nucleotides. There are $4^{1200} \approx 10^{720}$ different sequences that are 1200 nucleotides long.

were completely new, a mutation that has never been seen before and will never be seen again. Does that description ring any bells? Does the infinite alleles model sound familiar? It should, because it exactly fits the situation I've just described.

Having remembered that this situation is well described by the infinite alleles model, I'm sure you'll also remember that we can calculate the equilibrium inbreeding coefficient for the infinite alleles model, i.e.,

$$f = \frac{1}{4N_e\mu + 1} \quad .$$

What's important about this for our purposes, is that to the extent that the infinite alleles model is appropriate for molecular data, then f is the frequency of homozygotes we should see in populations and $1 - f$ is the frequency of heterozygotes. So in large populations we should find more diversity than in small ones, which is roughly what we do find. Notice, however, that here we're talking about heterozygosity at individual nucleotide positions,⁵ not heterozygosity of haplotypes.

Conclusions

In broad outline then, the neutral theory does a pretty good job of dealing with at least some types of molecular data. I'm sure that some of you are already thinking, "But what about third codon positions *versus* first and second?" or "What about the observation that histone loci evolve much more slowly than interferons or MHC loci?" Those are good questions, and those are where we're going next. As we'll see, molecular evolutionists have elaborated the framework extensively⁶ in the last thirty years, but these basic principles underlie every investigation that's conducted. That's why I wanted to spend a fair amount of time going over the logic and consequences. Besides, it's a rare case in population genetics where the fundamental mathematics that lies behind some important predictions are easy to understand.⁷

⁵Since the mutation rate we're talking about applies to individual nucleotide positions.

⁶That mean's they've made it more complicated.

⁷It's the concepts that get tricky, not the algebra, or at least that's what I think.

Chapter 29

Patterns of nucleotide and amino acid substitution

So I've just suggested that the neutral theory of molecular evolution explains quite a bit, but it also ignores quite a bit.¹ The derivations we did assumed that all substitutions are equally likely to occur, because they are selectively neutral. That isn't plausible. We need look no further than sickle cell anemia to see an example of a protein polymorphism in which a single amino acid difference has a very large effect on fitness. Even reasoning from first principles we can see that it doesn't make much sense to think that all nucleotide substitutions are created equal. Just as it's unlikely that you'll improve the performance of your car if you pick up a sledgehammer, open its hood, close your eyes, and hit something inside, so it's unlikely that picking a random amino acid in a protein and substituting it with a different one will improve the function of the protein.²

The genetic code

Of course, not all nucleotide sequence substitutions lead to amino acid substitutions in protein-coding genes. There is redundancy in the genetic code. Table 29.1 is a list of the codons in the universal genetic code.³ Notice that there are only two amino acids, methionine

¹I won't make my bikini joke, because it doesn't conceal as much as quantitative genetics. But still the "pure" version of the neutral theory of molecular evolution makes a *lot* of simplifying assumptions.

²Obviously it happens sometimes. If it didn't, there wouldn't be any adaptive evolution. It's just that, on average, mutations are more likely to decrease fitness than to increase it.

³By the way, the "universal" genetic code is not universal. There are at least eight, but all of them have similar redundancy properties.

Codon	Amino Acid	Codon	Amino Acid	Codon	Amino Acid	Codon	Amino Acid
UUU	Phe	UCU	Ser	UAU	Tyr	UGU	Cys
UUC	Phe	UCC	Ser	UAC	Tyr	UGC	Cys
UUA	Leu	UCA	Ser	UAA	Stop	UGA	Stop
UUG	Leu	UCG	Ser	UAG	Stop	UGG	Trp
CUU	Leu	CCU	Pro	CAU	His	CGU	Arg
CUC	Leu	CCC	Pro	CAC	His	CGC	Arg
CUA	Leu	CCA	Pro	CAA	Gln	CGA	Arg
CUG	Leu	CCG	Pro	CAG	Gln	CGG	Arg
AUU	Ile	ACU	Thr	AAU	Asn	AGU	Ser
AUC	Ile	ACC	Thr	AAC	Asn	AGC	Ser
AUA	Ile	ACA	Thr	AAA	Lys	AGA	Arg
AUG	Met	ACG	Thr	AAG	Lys	AGG	Arg
GUU	Val	GCU	Ala	GAU	Asp	GGU	Gly
GUC	Val	GCC	Ala	GAC	Asp	GGC	Gly
GUA	Val	GCA	Ala	GAA	Glu	GGA	Gly
GUG	Val	GCG	Ala	GAG	Glu	GGG	Gly

Table 29.1: The universal genetic code.

and tryptophan, that have a single codon. All the rest have at least two. Serine, arginine, and leucine have six.

Moreover, most of the redundancy is in the third position, where we can distinguish 2-fold from 4-fold redundant sites (Table 29.2). 2-fold redundant sites are those at which either one of two nucleotides can be present in a codon for a single amino acid. 4-fold redundant sites are those at which any of the four nucleotides can be present in a codon for a single amino acid. In some cases there is redundancy in the first codon position, e.g, both AGA and CGA are codons for arginine. Thus, many nucleotide substitutions at third positions do not lead to amino acid substitutions, and some nucleotide substitutions at first positions do not lead to amino acid substitutions. But every nucleotide substitution at a second codon position leads to an amino acid substitution. Nucleotide substitutions that do not lead to amino acid substitutions are referred to as *synonymous substitutions*, because the codons involved are synonymous, i.e., code for the same amino acid. Nucleotide substitutions that

Codon	Amino Acid	Redundancy
CCU	Pro	4-fold
CCC		
CCA		
CCG		
AAU	Asn	2-fold
AAC	Lys	2-fold
AAA		
AAG		

Table 29.2: Examples of 4-fold and 2-fold redundancy in the 3rd position of the universal genetic code.

do lead to amino acid substitutions are *non-synonymous substitutions*.

Rates of synonymous and non-synonymous substitution

By using a modification of the simple Jukes-Cantor model we encountered before, it is possible to make separate estimates of the number of synonymous substitutions and of the number of non-synonymous substitutions that have occurred since two sequences diverged from a common ancestor. If we combine an estimate of the *number* of differences with an estimate of the *time of divergence* we can estimate the rates of synonymous and non-synonymous substitution (number/time). Table 29.3 shows some representative estimates for the rates of synonymous and non-synonymous substitution in different genes studied in mammals.

Two very important observations emerge after you've looked at this table for awhile. The first won't come as any shock. The rate of non-synonymous substitution is generally lower than the rate of synonymous substitution. This is a result of my "sledgehammer principle." Mutations that change the amino acid sequence of a protein are more likely to reduce that protein's functionality than to increase it. As a result, they are likely to lower the fitness of individuals carrying them, and they will have a lower probability of being fixed than those mutations that do not change the amino acid sequence.

The second observation is more subtle. Rates of non-synonymous substitution vary by more than two orders of magnitude: 0.02 substitutions per nucleotide per billion years in ribosomal protein S14 to 3.06 substitutions per nucleotide per billion years in γ -interferon,

Locus	Non-synonymous rate	Synonymous rate
Histone		
H4	0.00	3.94
H2	0.00	4.52
Ribosomal proteins		
S17	0.06	2.69
S14	0.02	2.16
Hemoglobins & myoglobin		
α -globin	0.56	4.38
β -globin	0.78	2.58
Myoglobin	0.57	4.10
Interferons		
γ	3.06	5.50
$\alpha 1$	1.47	3.24
$\beta 1$	2.38	5.33

Table 29.3: Representative rates of synonymous and non-synonymous substitution in mammalian genes (from [62]). Rates are expressed as the number of substitutions per 10^9 years.

while rates of synonymous substitution vary only by a factor of two (2.16 in ribosomal protein S14 to 4.52 in histone H2). If synonymous substitutions are neutral, as they probably are to a first approximation,⁴ then the rate of synonymous substitution should equal the mutation rate. Thus, the rate of synonymous substitution should be approximately the same at every locus, which is roughly what we observe. But proteins differ in the degree to which their physiological function affects the performance and fitness of the organisms that carry them. Some, like histones and ribosomal proteins, are intimately involved with chromatin or translation of messenger RNA into protein. It's easy to imagine that just about any change in the amino acid sequence of such proteins will have a detrimental effect on its function. Others, like interferons, are involved in responses to viral or bacterial pathogens. It's easy to imagine not only that the selection on these proteins might be less intense, but that some amino acid substitutions might actually be favored by natural selection because they enhance resistance to certain strains of pathogens. Thus, the probability that a non-synonymous substitution will be fixed is likely to vary substantially among genes, just as we

⁴We'll see that they may not be completely neutral a little later, but at least it's reasonable to believe that the intensity of selection to which they are subject is less than that to which non-synonymous substitutions are subject.

observe.

Revising the neutral theory

So we've now produced empirical evidence that many mutations are *not* neutral. Does this mean that we throw the neutral theory of molecular evolution away? Hardly. We need only modify it a little to accommodate these new observations.

- *Most non-synonymous substitutions are deleterious.* We can actually generalize this assertion a bit and say that most mutations that affect function are deleterious. After all, organisms have been evolving for about 3.5 billion years. Wouldn't you expect their cellular machinery to work pretty well by now?
- *Most molecular variability found in natural populations is selectively neutral.* If most function-altering mutations are deleterious, it follows that we are unlikely to find much variation in populations for such mutations. Selection will quickly eliminate them.
- *Natural selection is primarily purifying.* Although natural selection for variants that improve function is ultimately the source of adaptation, even at the molecular level, most of the time selection is simply eliminating variants that are less fit than the norm, not promoting the fixation of new variants that increase fitness.
- *Alleles enhancing fitness are rapidly incorporated.* They do not remain polymorphic for long, so we aren't likely to find them when they're polymorphic.

As we'll see, even these revisions aren't entirely sufficient, but what we do from here on out is more to provide refinements and clarifications than to undertake wholesale revisions.

Chapter 30

Detecting selection on nucleotide polymorphisms

At this point, we've refined the neutral theory quite a bit. Our understanding of how molecules evolve now recognizes that some substitutions are more likely than others, but we're still proceeding under the assumption that most nucleotide substitutions are neutral or detrimental. So far we've argued that variation like what Hubby and Lewontin [45, 60] found is not likely to be maintained by natural selection. But we have strong evidence that heterozygotes for the sickle-cell allele are more fit than either homozygote in human populations where malaria is prevalent. That's an example where selection is acting to maintain a polymorphism, not to eliminate it. Are there other examples? How could we detect them?

In the 1970s a variety of studies suggested that a polymorphism in the locus coding for alcohol dehydrogenase in *Drosophila melanogaster* might not only be subject to selection but that selection may be acting to maintain the polymorphism. As DNA sequencing became more practical at about the same time,¹ population geneticists began to realize that comparative analyses of DNA sequences at protein-coding loci could provide a powerful tool for unraveling the action of natural selection. Synonymous sites within a protein-coding sequence provide a powerful standard of comparison. Regardless of

- the demographic history of the population from which the sequences were collected,
- the length of time that populations have been evolving under the sample conditions and whether it has been long enough for the population to have reached a drift-migration-mutation-selection equilibrium, or

¹It was still *vastly* more laborious than it is now.

- the actual magnitude of the mutation rate, the migration rate, or the selection coefficients

the synonymous positions within the sequence provide an internal control on the amount and pattern of differentiation that should be expected when substitutions.² Thus, if we see different patterns of nucleotide substitution at synonymous and non-synonymous sites, we can infer that selection is having an effect on amino acid substitutions.

Nucleotide sequence variation at the *Adh* locus in *Drosophila melanogaster*

Kreitman [56] took advantage of these ideas to provide additional insight into whether natural selection was likely to be involved in maintaining the polymorphism at *Adh* in *Drosophila melanogaster*. He cloned and sequenced 11 alleles at this locus, each a little less than 2.4kb in length.³ If we restrict our attention to the coding region, a total of 765bp, there were 6 distinct sequences that differed from one another at between 1 and 13 sites. Given the observed level of polymorphism within the gene, there should be 9 or 10 amino acid differences observed as well, but only one of the nucleotide differences results in an amino acid difference, the amino acid difference associated with the already recognized electrophoretic polymorphism. Thus, there is significantly less amino acid diversity than expected if nucleotide substitutions were neutral, consistent with my assertion that most mutations are deleterious and that natural selection will tend to eliminate them. In other words, another example of the “sledgehammer principle.”

Does this settle the question? Is the *Adh* polymorphism another example of allelic variants being neutral or selected against? Would I be asking these questions if the answer were “Yes”?

Kreitman and Aguadé

A few years after Kreitman [56] appeared, Kreitman and Aguadé [57] published an analysis in which they looked at levels of nucleotide diversity in the *Adh* region, as revealed through analysis of RFLPs, in *D. melanogaster* and the closely related *D. simulans*. Why the comparative approach? Well, Kreitman and Aguadé recognized that the neutral theory

²Ignoring, for the moment, the possibility that there may be selection on codon usage.

³Think about how the technology has changed since then. This work represented a major part of his Ph.D. dissertation, and the results were published as an article in *Nature*.

	5' flanking	<i>Adh</i> locus	3' flanking
Diversity ¹			
Observed	9	14	2
Expected	10.8	10.8	3.4
Divergence ²			
Observed	86	48	31
Expected	55	76.9	33.1

¹Number of polymorphic sites within *D. melanogaster*

²Number of nucleotide differences between *D. melanogaster* and *D. simulans*

Table 30.1: Diversity and divergence in the *Adh* region of *Drosophila* (from [57]).

of molecular evolution makes two predictions that are related to the underlying mutation rate:

- If mutations are neutral, the substitution rate is equal to the mutation rate.
- If mutations are neutral, the diversity within populations should be about $4N_e\mu/(4N_e\mu + 1)$.

Thus, if variation at the *Adh* locus in *D. melanogaster* is selectively neutral, the amount of divergence between *D. melanogaster* and *D. simulans* should be related to the amount of diversity within each. What they found instead is summarized in Table 30.1.

Notice that there is substantially less divergence at the *Adh* locus than would be expected, based on the average level of divergence across the entire region. That's consistent with the earlier observation that most amino acid substitutions are selected against. On the other hand, there is *more* nucleotide diversity within *D. melanogaster* than would be expected based on the levels of diversity seen in across the entire region. What gives?

Time for a trip down memory lane. Remember something called “coalescent theory?” It told us that for a sample of neutral genes from a population, the expected time back to a common ancestor for all of them is about $4N_e$ for a nuclear gene in a diploid population. That means there's been about $4N_e$ generations for mutations to occur. Suppose, however, that the electrophoretic polymorphism were being maintained by natural selection. Then we might well expect that it would be maintained for a lot longer than $4N_e$ generations. If so, there would be a lot more time for diversity to accumulate. Thus, the excess diversity could be accounted for if there is balancing selection at ADH.

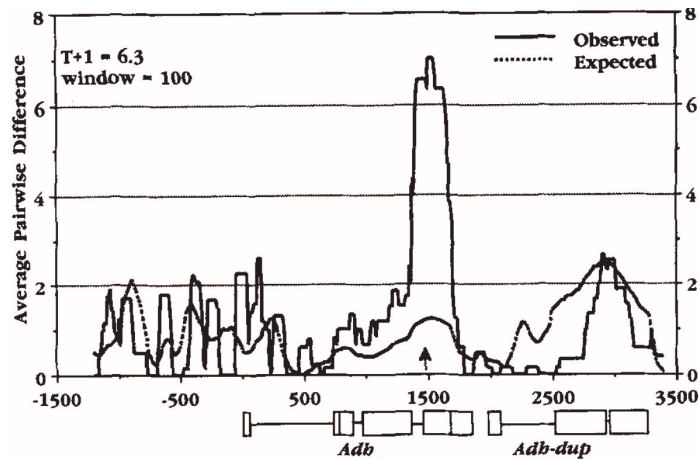


Figure 30.1: Sliding window analysis of nucleotide diversity in the *Adh-Adh-dup* region of *Drosophila melanogaster*. The arrow marks the position of the single nucleotide substitution that distinguishes *Adh-F* from *Adh-S* (from [58])

Kreitman and Hudson

Kreitman and Hudson [58] extended this approach by looking more carefully within the region to see where they could find differences between observed and expected levels of nucleotide sequence diversity. They used a “sliding window” of 100 silent base pairs in their calculations. By “sliding window” what they mean is that first they calculate statistics for bases 1-100, then for bases 2-101, then for bases 3-102, and so on until they hit the end of the sequence. It’s rather like walking a chromosome for QTL mapping, and the results are rather pretty (Figure 30.1).

To me there are two particularly striking things about this figure. First, the position of the single nucleotide substitution responsible for the electrophoretic polymorphism is clearly evident. Second, the excess of polymorphism extends for only a 200-300 nucleotides in each direction. That means that the rate of recombination *within* the gene is high enough to randomize the nucleotide sequence variation farther away.⁴

⁴Think about what that means for association mapping. In organisms with a large effective population size, associations due to physical linkage may fall off *very* rapidly, meaning that you would have to have a *very* dense map to have a hope of finding associations.

Detecting selection in the human genome

I've already mentioned the HapMap project [15], a collection of genotype data at roughly 3.2M SNPs in the human genome. The data in phase II of the project were collected from four populations:

- Yoruba (Ibadan, Nigeria)
- Japanese (Tokyo, Japan)
- Han Chinese (Beijing, China)
- ancestry from northern and western Europe (Utah, USA)

We expect genetic drift to result in allele frequency differences among populations, and we can summarize the extent of that differentiation at each locus with F_{ST} . If all HapMap SNPs are selectively neutral,⁵ then all loci should have the same F_{ST} within the bounds of statistical sampling error and the evolutionary sampling due to genetic drift. A scan of human chromosome 7 reveals both a lot of variation in individual-locus estimates of F_{ST} and a number of loci where there is substantially more differentiation among populations than is expected by chance (Figure 30.2). At very fine genomic scales we can detect even more outliers (Figure 30.3), suggesting that human populations have been subject to divergent selection pressures at many different loci [32].

⁵And unlinked to sites that are under selection.

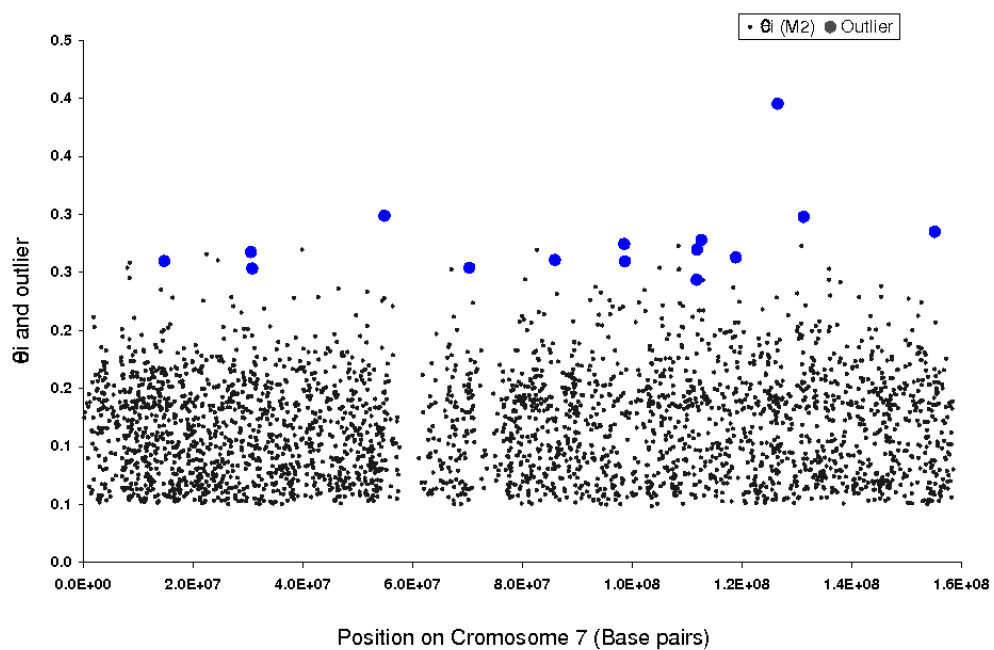


Figure 30.2: Single-locus estimates of F_{ST} along chromosome 7 in the HapMap data set. Blue dots denote outliers. Adjacent SNPs in this sample are separated, on average, by about 52kb. (from [32])

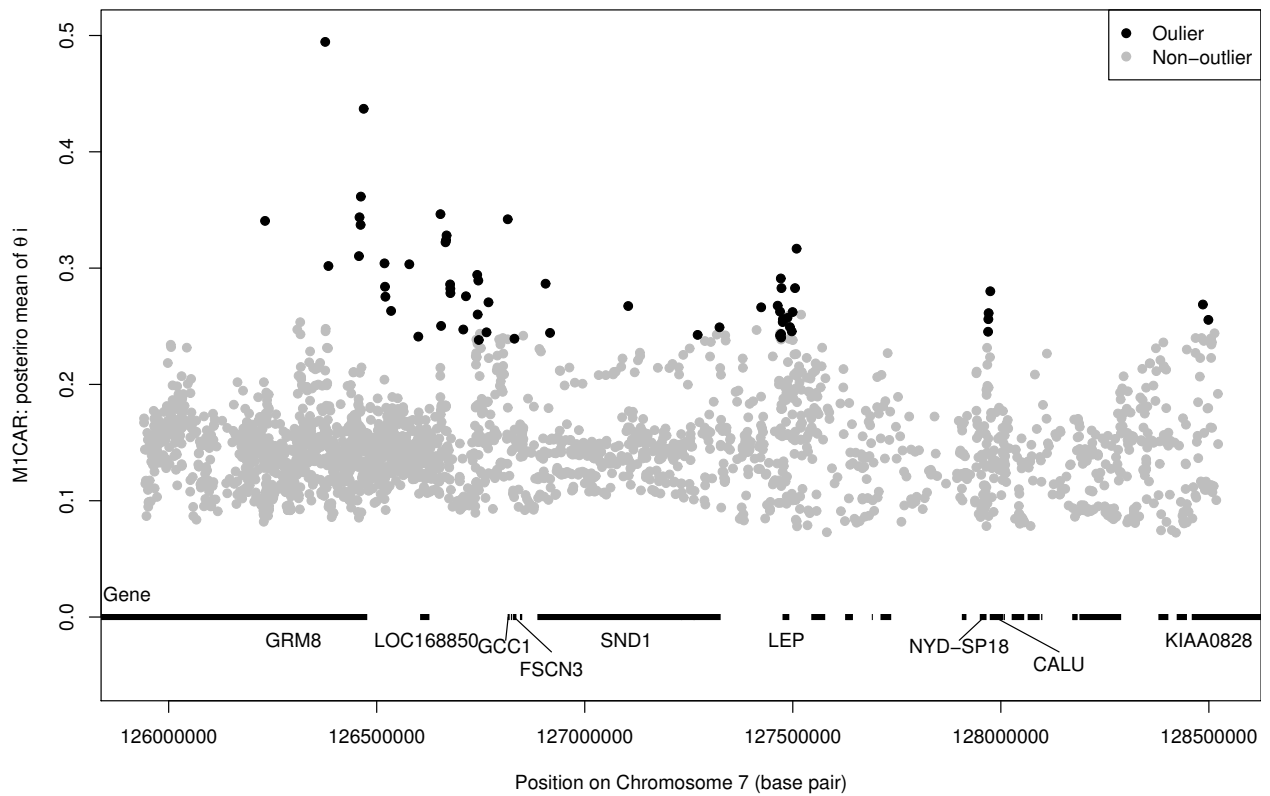


Figure 30.3: Single-locus estimates of F_{ST} along a portion of chromosome 7 in the HapMap data set. Black dots denote outliers. Solid bars refer to previously identified genes. Adjacent SNPs in this sample are separated, on average, by about 1kb. (from [32])

Chapter 31

Patterns of selection on nucleotide polymorphisms

We've now seen one good example of natural selection acting to maintain diversity at the molecular level, but that example involves only a pair of alleles. Let's examine how selection operates on a more complex polymorphism involving many alleles and several loci, specifically the polymorphisms at the major histocompatibility complex (MHC) loci of vertebrates.

MHC molecules are responsible for cellular immune responses in vertebrates. They are expressed on all nucleated cells in vertebrates, and they present intracellularly processed "foreign" antigens to T cell receptor lymphocytes. When the MHC + antigen complex is recognized, a cytotoxic reaction is triggered killing cells presenting the antigen. It's been known for many years that the genes are highly polymorphic.¹ Although plausible adaptive scenarios for that variation existed, a competing hypothesis had been that MHC loci were "hypervariable" not because of selection for diversity, but because of an unusually high mutation rate.

Patterns of amino acid substitution at MHC loci

Hughes and Nei [46] recognized that these hypotheses could be distinguished by comparing rates of synonymous and non-synonymous substitution at MHC loci. The results are summarized in Table 31.1. Notice that they distinguished among three functional regions within the protein and calculated statistics separately for each one:

- codons in the *antigen recognition site*,

¹They were discovered as a result of investigations into rejection of transplanted organs and tissues. They are the loci governing acceptance/rejection of transplants in vertebrates.

Locus	ARS		α_1 and α_2		α_3	
	K_s	K_a	K_s	K_a	K_s	K_a
Human						
HLA-A	3.5	13.3***	2.5	1.6	9.5	1.6**
HLA-B	7.1	18.1**	6.9	2.4	1.5	0.5
HLA-C	3.8	8.8	10.4	4.8	2.1	1.0
Mean	4.7	14.1***	5.1	2.4	5.8	1.1**
Mouse						
H2-K	15.0	22.9	8.7	5.8	2.3	4.0
H2-L	11.4	19.5	8.8	6.8	0.0	2.5**
Mean	13.2	21.2*	8.8	6.3	1.2	3.6**

Table 31.1: Rates of synonymous and non-synonymous substitution for loci in the MHC complex of humans and mice (modified from [62] and based on [46]). ARS refers to the antigen recognition site. Significant differences between K_s and K_a are denoted as: * ($P < 0.05$), ** ($P < 0.01$), and *** ($P < 0.001$).

- the remaining codons in the extracellular domain involved in presenting the antigen on the cell surface (the α_1 and α_2 domains), and
- codons in the extracellular domain that are not directly involved in presenting the antigen on the cell surface (the α_3 domain).

Hughes and Nei argue that the unusually low value of K_s in the α_3 domain of H2-L in mice is due to interlocus genetic exchange. If we discount that set of data as unreliable, a clear pattern emerges.

- In the part of the MHC molecule that is not directly involved in presenting antigen, α_3 in humans, the rate of non-synonymous substitution is significantly lower than the rate of synonymous substitution, i.e., there is selection *against* amino acid substitutions.²
- In the parts of the MHC molecule that presents antigens, α_1 and α_2 , the rate of synonymous and non-synonymous substitution is indistinguishable, except within the antigen recognition site where there are *more* non-synonymous than synonymous substitutions, i.e., there is selection *for* amino acid substitutions.

²No surprise there. That's the "sledgehammer principle in operation.

It's worth spending a little time thinking about what I mean when I say that there is selection *for* or *against* amino acid substitutions.

- Everything we know about DNA replication and mutation tells us that mutations arise independently of any fitness effect they have.
- Since the substitution rate is the product of the mutation rate and the probability of fixation, if some substitutions occur at a slower rate than neutral substitutions, they must have a lower probability of fixation, and the only way that can happen is if there is natural selection *against* those substitutions.
- Similarly, if some substitutions occur at a higher rate than neutral substitutions, they must have a higher probability of fixation, i.e., there is natural selection *for* those substitutions.

In a later paper Hughes et al. [47] took these observations even further. They subdivided the antigen recognition site into the binding cleft, the T-cell-receptor-directed residues, and the outward-directed residues. They found that the rate of non-synonymous substitution is much higher in the binding cleft than in other parts of the antigen recognition site and that nucleotide substitutions that change the charge of the associated amino acid residue are even more likely to be incorporated than those that are charge-conservative. In short, we have very strong evidence that natural selection is promoting diversity in the antigen binding capacity of MHC molecules.

Notice, however, that this selection for diversity is combined with overall conservatism in amino acid substitutions. Across the protein as a whole, most non-synonymous substitutions are selected *against*. Of course, it is that small subset of amino acids where non-synonymous substitutions are selected *for* that are responsible for adaptive responses to new pathogens.

Chapter 32

Tajima's D , Fu's F_S , Fay and Wu's H , and Zeng et al.'s E

So far we've been comparing rates of synonymous and non-synonymous substitution to detect the effects of natural selection on molecular polymorphisms. Tajima [87] proposed a method that builds on the foundation of the neutral theory of molecular evolution in a different way. I've already mentioned the infinite alleles model of mutation several times. When thinking about DNA sequences a closely related approximation is to imagine that every time a mutation occurs, it occurs at a different site.¹ If we do that, we have an *infinite sites* model of mutation.

Tajima's D

When dealing with nucleotide sequences in a population context there are two statistics of potential interest:

- The *number* of nucleotide positions at which a polymorphism is found or, equivalently, the number of segregating sites, k .
- The average per nucleotide diversity, π , where π is estimated as

$$\pi = \sum x_i x_j \delta_{ij} / N \quad .$$

¹Of course, we know this isn't true. Multiple substitutions *can* occur at any site. That's why the percent difference between two sequences isn't equal to the number of substitutions that have happened at any particular site. We're simply assuming that the sequences we're comparing are closely enough related that nearly all mutations have occurred at different positions.

In this expression, x_i is the frequency of the i th haplotype, δ_{ij} is the number of nucleotide sequence differences between haplotypes i and j , and N is the total length of the sequence.²

The quantity $4N_e\mu$ comes up a lot in mathematical analyses of molecular evolution. Population geneticists, being a lazy bunch, get tired of writing that down all the time, so they invented the parameter $\theta = 4N_e\mu$ to save themselves a little time.³ Under the infinite-sites model of DNA sequence evolution, it can be shown that

$$\begin{aligned} E(\pi) &= \theta \\ E(k) &= \theta \sum_i^{n-1} \frac{1}{i} \quad , \end{aligned}$$

where n is the number of haplotypes in your sample.⁴ This suggests that there are two ways to estimate θ , namely

$$\begin{aligned} \hat{\theta}_\pi &= \hat{\pi} \\ \hat{\theta}_k &= \frac{k}{\sum_i^{n-1} \frac{1}{i}} \quad , \end{aligned}$$

where $\hat{\pi}$ is the average heterozygosity at nucleotide sites in our sample and k is the observed number of segregating sites in our sample.⁵ If the nucleotide sequence variation among our haplotypes is neutral and the population from which we sampled is in equilibrium with respect to drift and mutation, then $\hat{\theta}_\pi$ and $\hat{\theta}_k$ should be statistically indistinguishable from one another. In other words,

$$\hat{D} = \hat{\theta}_\pi - \hat{\theta}_k$$

²I lied, but you must be getting used to that by now. This isn't quite the way you estimate it. To get an unbiased estimate of pi, you have to multiply this equation by $n/(n-1)$, where n is the number of haplotypes in your sample. And, of course, if you're Bayesian you'll be even a little more careful. You'll estimate x_i using an appropriate prior on haplotype frequencies and you'll estimate the probability that haplotypes i and j are different at a randomly chosen position given the observed number of differences and the sequence length. That probability will be close to δ_{ij}/N , but it won't be identical.

³This is *not* the same θ we encountered when discussing F -statistics. Weir and Cockerham's θ is a different beast. I know it's confusing, but that's the way it is. When reading a paper, the context should make it clear which conception of θ is being used. Another thing to be careful of is that sometimes authors think of θ in terms of a haploid population. When they do, it's $2N_e\mu$. Usually the context makes it clear which definition is being used, but you have to remember to pay attention to be sure.

⁴The "E" refers to expectation. It is the average value of a random variable. $E(\pi)$ is read as "the expectation of π ".

⁵If your memory is really good, you may recognize that those estimates are method of moments estimates, i.e., parameter estimates obtained by equating sample statistics with their expected values.

should be indistinguishable from zero. If it is either negative or positive, we can infer that there's some departure from the assumptions of neutrality and/or equilibrium. Thus, \hat{D} can be used as a test statistic to assess whether the data are consistent with the population being at a neutral mutation-drift equilibrium. Consider the value of D under following scenarios:

Neutral variation If the variation is neutral and the population is at a drift-mutation equilibrium, then \hat{D} will be statistically indistinguishable from zero.

Overdominant selection Overdominance will allow alleles belonging to the different classes to become quite divergent from one another. δ_{ij} within each class will be small, but δ_{ij} between classes will be large and both classes will be in intermediate frequency, leading to large values of θ_π . There won't be a similar tendency for the *number* of segregating sites to increase, so θ_k will be relatively unaffected. As a result, \hat{D} will be positive.

Population bottleneck If the population has recently undergone a bottleneck, then π will be little affected unless the bottleneck was prolonged and severe.⁶ k , however, may be substantially reduced. Thus, \hat{D} should be positive.

Purifying selection If there is purifying selection, mutations will occur and accumulate at silent sites, but they aren't likely ever to become very common. Thus, there are likely to be lots of segregating sites, but not much heterozygosity, meaning that $\hat{\theta}_k$ will be large, $\hat{\theta}_\pi$ will be small, and \hat{D} will be negative.

Population expansion Similarly, if the population has recently begun to expand, mutations that occur are unlikely to be lost, increasing $\hat{\theta}_k$, but it will take a long time before they contribute to heterozygosity, $\hat{\theta}_\pi$. Thus, \hat{D} will be negative.

In short, \hat{D} provides a different avenue for insight into the evolutionary history of a particular nucleotide sequence. But interpreting it can be a little tricky.

$\hat{D} = 0$: We have no evidence for changes in population size or for any particular pattern of selection at the locus.⁷

$\hat{D} < 0$: The population size may be increasing or we may have evidence for purifying selection at this locus.

⁶Why? Because most of the heterozygosity is due to alleles of moderate to high frequency, and those are not the ones likely to be lost in a bottleneck. See the Appendix32 for more details.

⁷Please remember that the failure to detect a difference from 0 could mean that your sample size is too small to detect an important effect. If you can't detect a difference, you should try to assess what values of D are consistent with your data and be appropriately circumspect in your conclusions.

$\hat{D} > 0$: The population may have suffered a recent bottleneck (or be decreasing) or we may have evidence for overdominant selection at this locus.

If we have data available for more than one locus, we may be able to distinguish changes in population size from selection at any particular locus. After all, all loci will experience the same demographic effects, but we might expect selection to act differently at different loci, especially if we choose to analyze loci with different physiological function.

A quick search in Google Scholar reveals that the paper in which Tajima described this approach [87] has been cited over 5300 times. Clearly it has been widely used for interpreting patterns of nucleotide sequence variation. Although it is a very useful statistic, Zeng et al. [104] point out that there are important aspects of the data that Tajima's D does not consider. As a result, it may be less powerful, i.e., less able to detect departures from neutrality, than some alternatives.

Fu's F_S

Fu [29] proposes a different statistic based on the infinite sites model of mutation. He suggests estimating the probability of observing a random sample with a number of alleles equal to or smaller than the observed value under given the observed level of diversity and the assumption that all of the alleles are selectively neutral. If we call this probability \hat{S} , then

$$F_S = \ln \left(\frac{\hat{S}}{1 - \hat{S}} \right) .$$

A negative value of F_S is evidence for an excess number of alleles, as would be expected from a recent population expansion or from genetic hitchhiking. A positive value of F_S is evidence for a deficiency of alleles, as would be expected from a recent population bottleneck or from overdominant selection. Fu's simulations suggest that F_S is a more sensitive indicator of population expansion and genetic hitchhiking than Tajima's D . Those simulations also suggest that the conventional P-value of 0.05 corresponds to a P-value from the coalescent simulation of 0.02. In other words, F_S should be regarded as significant if $P < 0.02$.

Fay and Wu's H

Let ξ_i be the number of sites at which a sequence occurring i times in the sample differs from the sequence of the most recent common ancestor for all the sequences. Fu [28] showed that

$$E(\xi_i) = \frac{\theta}{i} .$$

Remember that i is the number of times this haplotype occurs in the sample. Using this result, we can rewrite $\hat{\theta}_\pi$ and $\hat{\theta}_k$ as

$$\begin{aligned}\hat{\theta}_\pi &= \binom{n}{2}^{-1} \sum_{i=1}^{n-1} i(n-i)\hat{\xi}_i \\ \hat{\theta}_k &= \frac{1}{a_n} \sum_{i=1}^{n-1} \hat{\xi}_i\end{aligned}$$

There are also at least two other statistics that could be used to estimate θ from these data:

$$\begin{aligned}\theta_H &= \binom{n}{2}^{-1} \sum_{i=1}^{n-1} i^2 \hat{\xi}_i \\ \theta_L &= \frac{1}{n-1} \sum_{i=1}^{n-1} i \hat{\xi}_i .\end{aligned}$$

Notice that to estimate θ_H or θ_L , you'll need information on the sequence of an ancestral haplotype. To get this you'll need an outgroup. As we've already seen, we can get estimates of θ_π and θ_k without an outgroup.

Fay and Wu [24] suggest using the statistic

$$H = \hat{\theta}_\pi - \theta_H$$

to detect departures from neutrality. So what's the difference between Fay and Wu's H and Tajima's D ? Well, notice that there's an i^2 term in θ_H . The largest contributions to this estimate of θ are coming from alleles in relatively high frequency, i.e., those with lots of copies in our sample. In contrast, intermediate-frequency alleles contribute most to estimates of θ_π . Thus, H measures departures from neutrality that are reflected in the difference between high-frequency and intermediate-frequency alleles. In contrast, D measures departures from neutrality that are reflected in the difference between low-frequency and intermediate frequency alleles. Thus, while D is sensitive to population expansion (because the number of segregating sites responds more rapidly to changes in population size than the nucleotide heterozygosity), H will not be. As a result, combining both tests may allow you to distinguish population expansion from purifying selection.

Zeng et al.'s E

So if we can use D to compare estimates of θ from intermediate- and low-frequency variants and H to compare estimates from intermediate- and high-frequency variants, what about

comparing estimates from high-frequency and low-frequency variants? Funny you should ask, Zeng et al. [104] suggest looking at

$$E = \theta_L - \theta_k \quad .$$

E doesn't put quite as much weight on high frequency variants as H ,⁸ but it still provides a useful contrast between estimates of θ derived from high-frequency variants and low-frequency variants. For example, suppose a new favorable mutation occurs and sweeps to fixation. All alleles other than those carrying the new allele will be eliminated from the population. Once the new variant is established, neutral variation will begin to accumulate. The return to neutral expectations after such an event, however, happens much more rapidly in low frequency variants than in high-frequency ones. Thus, a negative E may provide evidence of a recent selective sweep at the locus being studied. For similar reasons, it will be a sensitive indicator of recent population expansion.

Appendix

I noted earlier that π will be little affected by a population bottleneck unless it is prolonged and severe. Here's one way of thinking about it that might make that counterintuitive assertion a little clearer.

Remember that π is defined as $\pi = \sum x_i x_j \delta_{ij} / N$. Unless one haplotype in the population happens to be very divergent from all other haplotypes in the population, the magnitude of π will be approximately equal to the average difference between any two nucleotide sequences times the probability that two randomly chosen sequences represent different haplotypes. Thus, we can treat haplotypes as alleles and ask what happens to heterozygosity as a result of a bottleneck. Here we recall the relationship between identity by descent and drift, and we pretend that homozygosity is the same thing as identity by descent. If we do, then the heterozygosity after a bottleneck is

$$H_t = \left(1 - \frac{1}{2N_e}\right)^t H_0 \quad .$$

So consider a *really* extreme case: a population reduced to one male and one female for 5 generations. $N_e = 2$, so $H_5 \approx 0.24H_0$, so the population would retain roughly 24% of its original diversity even after such a bottleneck. Suppose it were less severe, say, five males and five females for 10 generations, then $N_e = 10$ and $H_{10} \approx 0.6$.

⁸Because it has an i rather than an i^2 in its formula

Chapter 33

Evolution in multigene families

We now know a lot about the dynamics of nucleotide substitutions within existing genes, but we've neglected one key component of molecular evolution. We haven't talked about where new genes come from. It's important to understand this phenomenon because, after all, new metabolic functions are likely to arise only when there are new genes that can perform them. It's not likely that an existing gene can adopt a new function while continuing to serve its old one.

Fundamentally the source of new genes is the *duplication* of existing genes and their *divergence* in function. As we'll see in a moment, for example, genes coding for myoglobin and hemoglobin in mammals are descendants of a single common ancestor. That's the duplication. Myoglobin is involved in oxygen metabolism in muscle, while hemoglobin is involved in oxygen transport in blood. That's the divergence. Although there are many interesting things to say about the processes by which duplication and divergence occur, we're going to focus on the pattern of nucleotide sequence evolution that arises as a result.

Globin evolution

I've just pointed out the distinction between myoglobin and hemoglobin. You may also remember that hemoglobin is a multimeric protein consisting of four subunits, 2 α subunits and 2 β subunits. What you may not know is that in humans there are actually two types of α hemoglobin and four types of β hemoglobin, each coded by a different genetic locus (see Table 33.1). The five α -globin loci (α_1 , α_2 , ζ , and two non-functional pseudogenes) are found in a cluster on chromosome 16. The six β -globin loci (ϵ , γ_G , γ_A , δ , β , and a pseudogene) are found in a cluster on chromosome 11. The myoglobin locus is on chromosome 22.

Not only do we have all of these different types of globin genes in our bodies, they're all

Developmental stage	α globin	β globin
Embryo	ζ	ϵ
	α	ϵ
Fetus	α	β
	α	γ
Adult	α	β
	α	δ

Table 33.1: Human hemoglobins arranged in developmental sequence. Adult hemoglobins composed of 2α and 2δ subunits typically account for less than 3% of hemoglobins in adults (<http://sickle.bwh.harvard.edu/hbsynthesis.html>).

related to one another. Comparative sequence analysis has shown that vertebrate myoglobin and hemoglobins diverged from one another about 450 million years ago. Figure 33.1 shows a phylogenetic analysis of globin genes from humans, mice, and a variety of Archaea. Focus your attention on the part of the tree that has human and mouse sequences. You'll notice two interesting things:

- Human and mouse neuroglobins (Ngb) are more closely related to one another than they are to other globins, even those from the same species. The same holds true for cytoglobins (Cyg) and myoglobins (Mb).
- Within the hemoglobins, only mouse β -globin (Mouse HbB) is misplaced. All other α - and β -globins group with the corresponding mouse and human loci.

This pattern is exactly what we expect as a result of duplication and divergence. Up to the time that a gene becomes duplicated, its evolutionary history matches the evolutionary history of the organisms containing it. Once there are duplicate copies, each follows an independent evolutionary history. Each traces the history of speciation and divergence. And over long periods duplicate copies of the same gene share more recent common ancestry with copies of the same gene in a different species than they do with duplicate genes in the same genome.

A history of duplication and divergence in multigene families makes it important to distinguish between two classes of related loci: those that represent the same locus in different species and between which divergence is a result of species divergence are *orthologs*. Those that represent different loci and between which divergence occurred after duplication of an ancestral gene are *paralogs*. The β -globin loci of humans and chickens are orthologous. The α - and β -globin loci of any pair of taxa are paralogous.

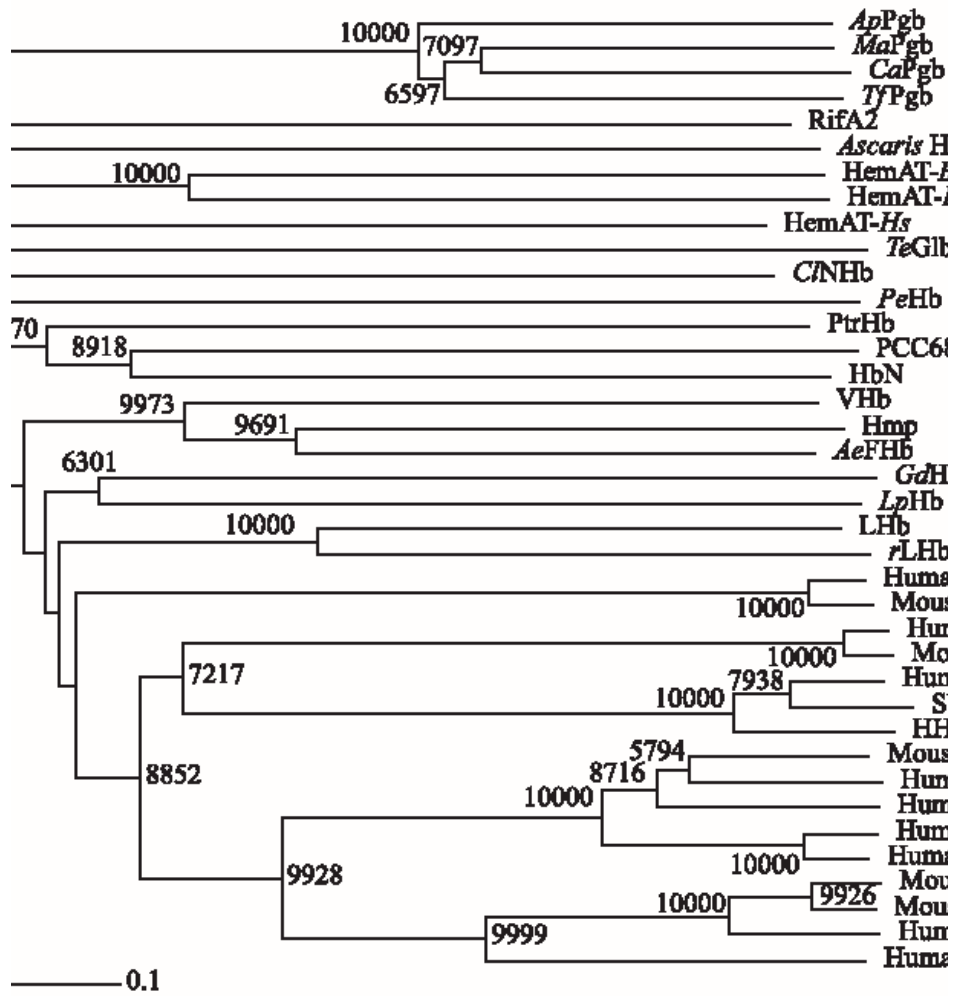


Figure 33.1: Evolution of globin genes in Archaea and mammals (from [26]).

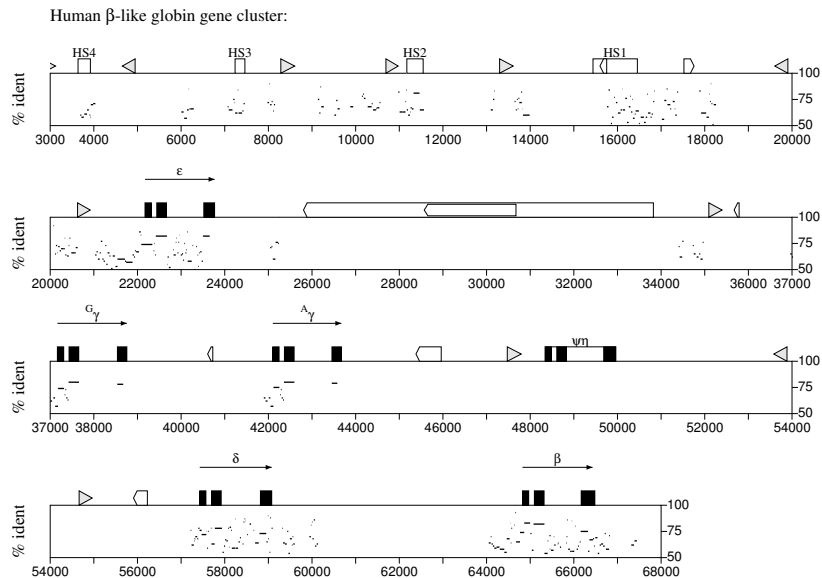


Figure 33.2: Structure of the human β -globin gene cluster. % identity refers to similarity to the mouse β -globin sequence. From <http://globin.cse.psu.edu/html/pip/betaglobin/iplot.ps> (retrieved 28 Nov 2006).

As multigene families go, the globin family is relatively simple and easy to understand. There are only about a dozen loci involved, one isolated locus (myoglobin) and two clusters of loci (α - and β -globins). You'll find a diagram of the β -globin cluster in Figure 33.2. As you can see the β -globins are not only evolutionarily related to one another they occur relatively close to one another on chromosome 11 in humans.

Other families are far more complex. Class I and class II MHC loci, for example are part of the same multigene family. Moreover, immunoglobulins, T-cell receptors, and, and MHC loci are part of a larger superfamily of genes, i.e., all are ultimately derived from a common ancestral gene by duplication and divergence. Table 33.2 lists a few examples of multigene families and superfamilies in the human genome and the number of proteins produced.

Protein family domain	Number of proteins
Actin	61
Immunoglobulin	381
Fibronectin type I	5
Fibronectin type II	11
Fibronectin type III	106
Histone	
H2A/H2B/H3/H4	75
Homeobox	160
Immunoglobulin	381
MHC Class I	18
MHC Class II α	5
MHC Class II β	7
T-cell receptor α	16
T-cell receptor β	15
T-cell receptor γ	1
T-cell receptor δ	1
Zinc finger, C2H2	564
Zinc finger, C3HC4	135

Table 33.2: A few gene families from the human genome (adapted from [73, 22]).

Concerted evolution

Although the patterns of gene relationships produced through duplication and divergence can be quite complex, the processes are relatively easy to understand. In some multigene families, however, something quite different seems to be going on. In many plants and animals, genes encoding ribosomal RNAs are present in hundreds of copies and arranged end to end in long tandem arrays in one or a few places in the genome (Figure 33.3). Brown et al. [9] compared the ribosomal RNA of *Xenopus laevis* and *X. mulleri* and found a surprising pattern. There was little or no detectable variation among copies of the repeat units within either species, in spite of substantial divergence between them. This pattern can't be explained by purifying selection. Members of the gene family presumably diverged before *X. laevis* and *X. mulleri* diverged. Thus, we would expect more divergence among copies *within* species than *between* species, i.e., the pattern we see in the globin family. Explaining this pattern requires some mechanism that causes different copies of the repeat to be homogenized within each species while allowing the repeats to diverge between species. The phenomenon is referred to as concerted evolution.

Two mechanisms that can result in concerted evolution have been widely studied: unequal crossing over and gene conversion. Both depend on misalignments during prophase. These misalignments allow a mutation that occurs in one copy of a tandemly repeated gene array to “spread” to other copies of the gene array. Tomoko Ohta and Thomas Nagylaki have provided exhaustive mathematical treatments of the process [66, 72]. We'll follow Ohta's treatment, but keep it fairly simple and straightforward. First some notation:

- f = P(two alleles at same locus are ibd)
- c_1 = P(two alleles at different loci in same chromosome are ibd)
- c_2 = P(two alleles at different loci in different chromosomes are ibd)
- μ = mutation rate
- n = no. of loci in family
- λ = rate of gene conversion

Now remember that for the infinite alleles model

$$f = \frac{1}{4N_e\mu + 1} \quad ,$$

and f is the probability that neither allele has undergone mutation. By analogy

$$g = \frac{1}{4N_e\lambda + 1} \quad ,$$

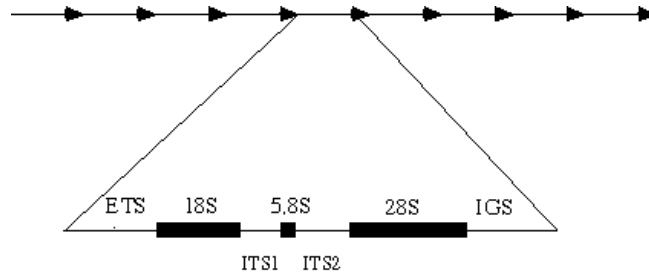


Figure 33.3: Diagrammatic representation of ribosomal DNA in vascular plant genomes (from Muir & Schlötterer, 1999 <http://webdoc.sub.gwdg.de/ebook/y/1999/whichmarker/m11/Chap11.htm>).

where g is the probability that two alleles at a homologous position are ibd in the sense that neither has ever moved from that position in the array. Thus, for our model

$$\begin{aligned}
 f &= P(\text{neither has moved})P(\text{ibd}) \\
 &\quad + P(\text{one has moved})P(\text{ibd anyway}) \\
 &= \left(\frac{1}{4N_e\lambda + 1}\right) \left(\frac{1}{4N_e\mu + 1}\right) + \left(\frac{4N_e\lambda}{4N_e\lambda + 1}\right) c_2 \\
 &\approx \frac{4N_e\lambda c_2 + 1}{4N_e\lambda + 4N_e\mu + 1} \\
 c_1 = c_2 &= \frac{\lambda}{\lambda + (n - 1)\mu} \quad .
 \end{aligned}$$

Notice that $(n - 1)\mu$ is approximately the number of mutations that occur in a single array every generation. Consider two possibilities:

- *Gene conversion occurs much more often than mutation:* $\lambda \gg (n - 1)\mu$.

Under these conditions $c_2 \approx 1$ and $f \approx 1$. In short, all copies of alleles at every locus in the array are virtually identical—concerted evolution.

- *Gene conversion occurs much less often than mutation:* $\lambda \ll (n - 1)\mu$.

Under these conditions $c_2 \approx 0$ and $f \approx \frac{1}{4N_e\mu + 1}$. In short, copies of alleles at different loci are almost certain to be different from one another, and the diversity at any single locus matches neutral expectations—non-concerted evolution.

Chapter 34

Analysis of mismatch distributions

Introduction

Remember when we were talking about Tajima's D ?¹ I pointed out that $\hat{\theta}_\pi$, the estimate of $4N_e\mu$ derived from nucleotide sequence diversity is less sensitive to demographic changes than θ_k , the estimate of $4N_e\mu$ derived from the number of segregating sites in the sample. I went on to argue that in a rapidly expanding population, mutation will not have “built up” the level of nucleotide diversity we'd expect based on the number of segregating sites, so that $\hat{D} = \theta_\pi - \theta_k$ will be negative. In a population that's suffered a recent bottleneck, on the other hand, there will be more nucleotide diversity than we'd expect on the basis of the number of segregating sites, so that \hat{D} will be positive.

Figures 1–3 may help you to visualize what's going on. We get to revisit our old friend the coalescent. Figure 34.1 shows the genealogical relationships among a set of alleles sampled from two different populations that exchange genes every other generation, on average that haven't changed size. The four different coalescent trees correspond to four different loci. The red and green dots correspond to the different populations from which the alleles were collected.

Looking at Figure 34.1 isn't particularly revealing by itself, except that it shows how much variability there is in coalescent history among loci, even when the demographic parameters. What's more interesting is to compare those trees with similar trees generated when the populations have undergone either a recent expansion (Figure 34.2) or a recent contraction (Figure 34.3). As you can see, when populations have undergone a recent expansion, all of the branches are relatively long. When they've undergone a recent bottleneck, on the other hand, all of the branches are quite short.

¹Don't answer that. I don't think I want to know the answer.

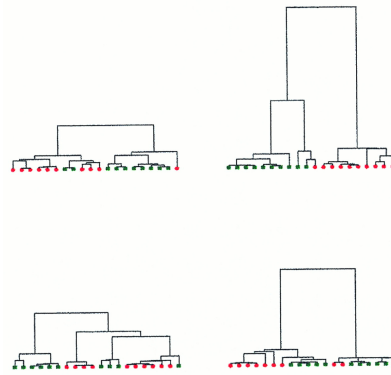


Figure 34.1: Four simulated coalescent trees for a sample of alleles from two populations of constant size that exchange genes every other generation on average (from [34]).

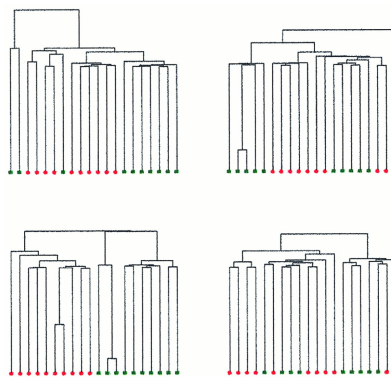


Figure 34.2: Four simulated coalescent trees for a sample of alleles from two populations that have undergone a recent expansion and exchange genes every other generation on average (from [34]).

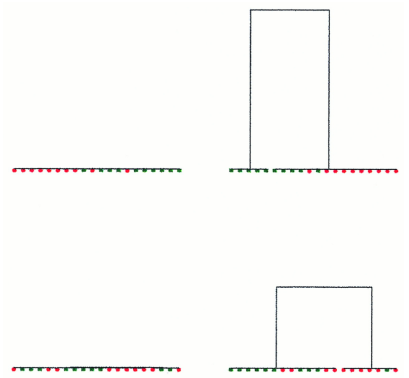


Figure 34.3: Four simulated coalescent trees for a sample of alleles from two populations that have undergone a recent contraction and exchange genes every other generation on average (from [34]).

Mismatch distributions

Since the amount of sequence difference between alleles depends on the length of time since they diverged, these observations suggest that we might be able to learn more about the recent demographic history of populations by looking not just at a summary statistic like θ_π or θ_k , but at the whole distribution of sequence differences. In fact, as Figure 34.4 and Figure 34.5 show, the differences are quite dramatic.

Harpending et al. [34] used this approach to analyze nucleotide sequence diversity in a sample of 636 mtDNA sequences. Their analysis focused on 411 positions in the first hypervariable segment of human mitochondrial DNA (Figure 34.6). The large excess of low-frequency variants suggests that the human population has undergone a recent population expansion. There is, of course, the possibility that purifying selection on the mitochondrion could explain the pattern, so they also analyzed sequence variation on the Y chromosome and found the same pattern. Patterns of variation at a variety of other loci are also compatible with the hypothesis of a recent expansion of human populations.

Estimating population parameters from mismatch distributions

Well, if we can detect recent population expansion (or contraction) in the characteristics of the mismatch distribution, maybe we can estimate some characteristics of the expansion.

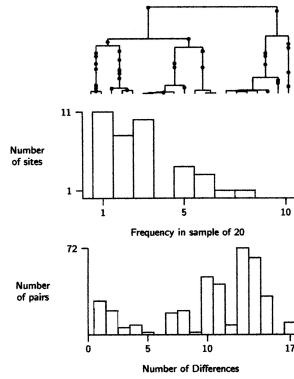


Figure 34.4: A gene tree (top), the frequency with which different haplotypes are found (middle), and the mismatch distribution (bottom) for a sample from a population of constant size (from [34]).

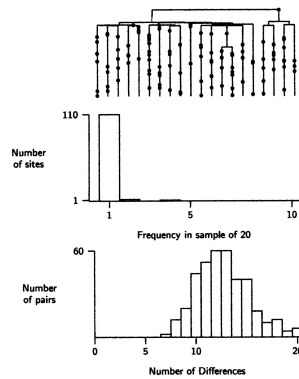


Figure 34.5: A gene tree (top), the frequency with which different haplotypes are found (middle), and the mismatch distribution (bottom) for a sample from a population that has undergone a recent expansion (from [34]).

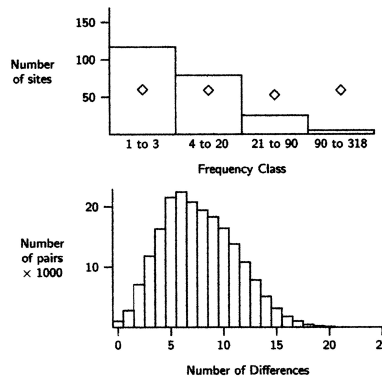


Figure 34.6: Mismatch distribution in a sample of 636 mtDNA sequences. The diamonds indicate expected values in the case of constant population size (from [34])

Suppose, for example, we consider a really simple example (Figure 34.7) where the population had some constant effective size, N_0 , and underwent an instantaneous expansion to a new effective size, N_1 , some unknown time t ago. We'll also assume that the mutation rate is μ ,² and we'll assume that we're dealing with a haploid population, e.g., human mitochondrial DNA.

You can probably already guess that the properties of the mismatch distribution under this simple model depends only on the products $N_0\mu$ and $N_1\mu$, so we'll define new parameters $\theta_0 = 2N_0\mu$ and $\theta_1 = 2N_1\mu$ to save ourselves a little bit of time. Similarly, we'll let $\tau = 2\mu t$. Given these assumptions, it's possible to calculate the mismatch distribution.³ Fortunately, you don't have to calculate it yourself. **Arlequin** will take care of that for you.⁴ Unfortunately, Schneider and Excoffier [80] show that of the three parameters we could estimate using this model, only τ is estimated with a reasonable degree of reliability.

DNA sequence data from hypervariable region 1 (mtDNA) in a sample from Senegalese Mandenka is distributed with **Arlequin**. If we estimate the parameters of demographic expansion, we get the results in Table 34.1. The sequence in question is 406 nucleotides long. If we assume that mutations occur at a rate of 2×10^{-6} per nucleotide per generation, then $\mu \approx 8 \times 10^{-4}$, so $t \approx 3875$ generations. In other words, according to these data the

²Since we're dealing with a neutral locus, the substitution rate is equal to the mutation rate. This is also the mutation rate for the entire stretch of DNA we're looking at. In other words, if the per nucleotide mutation rate is 10^{-9} and our DNA sequence is a thousand nucleotides long $\mu = 10^{-9} \times 10^3 = 10^{-6}$.

³You may be astonished to learn that I'm not going to give you a formula for the distribution. If you're interested, you can find it in [80].

⁴And give you bootstrapped confidence intervals to boot!

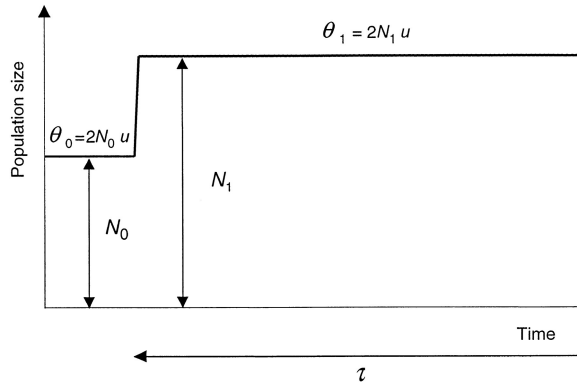


Figure 34.7: The simple demographic scenario underlying estimation of population expansion parameters from the mismatch distribution (from [80]).

Parameter	Mean	(99% CI)
θ_0	1.4	(1.4, 15.4)
θ_1	23.0	(0.0, 9.2)
τ	6.2	(10.8, 99999)

Table 34.1: Parameters of demographic expansion based on mtDNA sequence data distributed with *Arlequin*.

expansion took place roughly 75,000 years ago.⁵

Arlequin also reports two statistics that we can use to assess whether our model is working well: the sum of squared deviations (Ssd) and Harpending's raggedness index⁶

$$r = \sum_{i=1}^{d+1} (x_i - x_{i-1})^2 \quad ,$$

where x_i is the frequency of haplotypes that differ at i positions and d is the maximum number of observed differences. In these data

$$\begin{aligned} P(\text{Expected Ssd} \geq \text{Observed Ssd}) &= 0.83 \\ P(\text{Expected } r \geq \text{Observed } r) &= 0.96 \quad . \end{aligned}$$

⁵If we assume that I got the mutation rate roughly right.

⁶A test of Tajima's D will tell us whether we have evidence that there might have been an expansion. In this case, $\hat{D} = -1.14$ and $P = 0.10$, so we have only weak evidence of any departure from a drift-mutation equilibrium.

The large value of Harpending's r , in particular, suggests that the model doesn't provide a particularly good fit to these data.⁷

⁷Which is consistent with our analysis of Tajima's D that there isn't much evidence for departure from a drift-mutation equilibrium.

Part VI
Phylogeography

Chapter 35

Analysis of molecular variance (AMOVA)

We've already encountered π , the nucleotide diversity in a population, namely

$$\pi = \sum_{ij} x_i x_j \delta_{ij} \quad ,$$

where x_i is the frequency of the i th haplotype and δ_{ij} is the fraction of nucleotides at which haplotypes i and j differ.¹ It shouldn't come to any surprise to you that just as there is interest in partitioning diversity within and among populations when we're dealing with simple allelic variation, i.e., Wright's F -statistics, there is interest in partitioning diversity within and among populations when we're dealing with nucleotide sequence or other molecular data. We'll see later that AMOVA can be used very generally to partition variation when there is a distance we can use to describe how different alleles are from one another, but for now, let's stick with nucleotide sequence data for the moment and think of δ_{ij} simply as the fraction of nucleotide sites at which two sequences differ.

¹When I introduced nucleotide diversity before, I defined δ_{ij} as the *number* of nucleotides that differ between haplotypes i and j . It's a little easier for what follows if we think of it as the *fraction* of nucleotides at which they differ instead.

Analysis of molecular variation (AMOVA)

The notation now becomes just a little bit more complicated. We will now use x_{ik} to refer to the frequency of the i th haplotype in the k th population. Then

$$x_{i\cdot} = \frac{1}{K} \sum_{k=1}^K x_{ik}$$

is the mean frequency of haplotype i across all populations, where K is the number of populations. We can now define

$$\begin{aligned}\pi_t &= \sum_{ij} x_i \cdot x_j \cdot \delta_{ij} \\ \pi_s &= \frac{1}{K} \sum_{k=1}^K \sum_{ij} x_{ik} x_{jk} \delta_{ij} \quad ,\end{aligned}$$

where π_t is the nucleotide sequence diversity across the entire set of populations and π_s is the average nucleotide sequence diversity within populations. Then we can define

$$\Phi_{st} = \frac{\pi_t - \pi_s}{\pi_t} \quad , \quad (35.1)$$

which is the direct analog of Wright's F_{st} for nucleotide sequence diversity. Why? Well, that requires you to remember stuff we covered eight or ten weeks ago.

To be a bit more specific, refer back to <http://darwin.eeb.uconn.edu/eeb348/lecture-notes/wahlund/node4.html>. If you do, you'll see that we defined

$$F_{IT} = 1 - \frac{H_i}{H_t} \quad ,$$

where H_i is the average heterozygosity in individuals and H_t is the expected panmictic heterozygosity. Defining H_s as the average panmictic heterozygosity within populations, we then observed that

$$\begin{aligned}1 - F_{IT} &= \frac{H_i}{H_t} \\ &= \frac{H_i}{H_s} \frac{H_s}{H_t} \\ &= (1 - F_{IS})(1 - F_{ST}) \\ 1 - F_{ST} &= \frac{1 - F_{IT}}{1 - F_{IS}}\end{aligned}$$

$$\begin{aligned}
F_{ST} &= \frac{(1 - F_{IS}) - (1 - F_{IT})}{1 - F_{IS}} \\
&= \frac{(H_i/H_s) - (H_i/H_t)}{H_i/H_s} \\
&= 1 - \frac{H_s}{H_t} \quad .
\end{aligned}$$

In short, another way to think about F_{ST} is

$$F_{ST} = \frac{H_t - H_s}{H_t} \quad . \quad (35.2)$$

Now if you compare equation (35.1) and equation (35.2), you'll see the analogy.

So far I've motivated this approach by thinking about δ_{ij} as the fraction of sites at which two haplotypes differ and π_s and π_t as estimates of nucleotide diversity. But nothing in the algebra leading to equation (35.1) requires that assumption. Excoffier et al. [23] pointed out that other types of molecular data can easily be fit into this framework. We simply need an appropriate measure of the “distance” between different haplotypes or alleles. Even with nucleotide sequences the appropriate δ_{ij} may reflect something about the mutational pathway likely to connect sequences rather than the raw number of differences between them. For example, the distance might be a Jukes-Cantor distance or a more general distance measure that accounts for more of the properties we know are associated with nucleotide substitution. The idea is illustrated in Figure 35.1. Once we have δ_{ij} for all pairs of haplotypes or alleles in our sample, we can use the ideas lying behind equation (35.1) to partition diversity — the average distance between randomly chosen haplotypes or alleles — into within and among population components.² This procedure for partitioning diversity in molecular markers is referred to as an analysis of molecular variance or AMOVA (by analogy with the ubiquitous statistical procedure analysis of variance, ANOVA). Like Wright's F -statistics, the analysis can include several levels in the hierarchy.

An AMOVA example

Excoffier et al. [23] illustrate the approach by presenting an analysis of restriction haplotypes in human mtDNA. They analyze a sample of 672 mitochondrial genomes representing two

²As with F -statistics, the actual estimation procedure is more complicated than I describe here. Standard approaches to AMOVA use method of moments calculations analogous to those introduced by Weir and Cockerham for F -statistics [97]. Bayesian approaches are possible, but they are not yet widely available (meaning, in part, that I know how to do it, but I haven't written the necessary software yet).

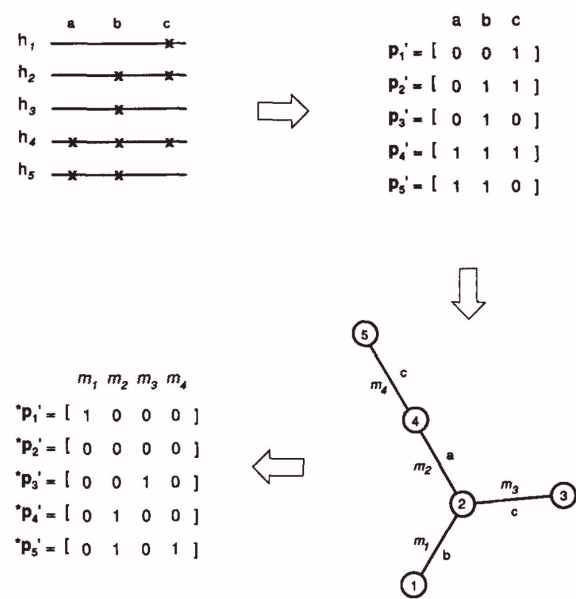


Figure 35.1: Converting raw differences in sequence (or presence and absence of restriction sites) into a minimum spanning tree and a mutational measure of distance for an analysis of molecular variance (from [23]).

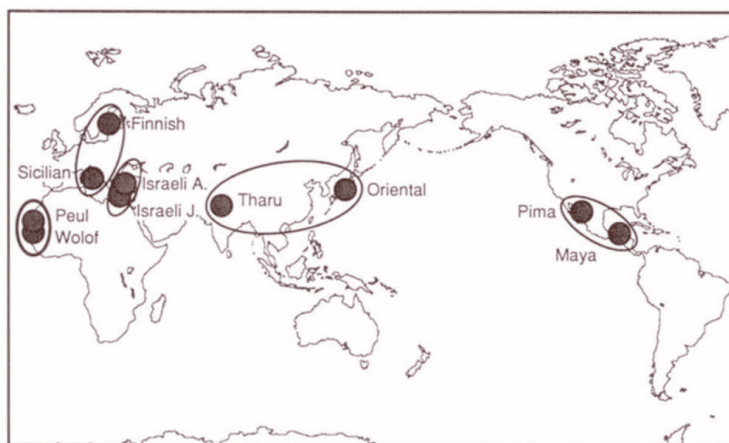


Figure 35.2: Locations of human mtDNA samples used in the example analysis (from [23]).

populations in each of five regional groups (Figure 35.2). They identified 56 haplotypes in that sample. A minimum spanning tree illustrating the relationships and the relative frequency of each haplotype is presented in Figure 35.3.

It's apparent from the figure that haplotype 1 is very common. In fact, it is present in substantial frequency in every sampled population. An AMOVA using the minimum spanning network in Figure 35.3 to measure distance produces the results shown in Table 35.1. Notice that there is relatively little differentiation among populations within the same geographical region ($\Phi_{SC} = 0.044$). There is, however, substantial differentiation among regions ($\Phi_{CT} = 0.220$). In fact, differences among populations in different regions is responsible for nearly all of the differences among populations ($\Phi_{ST} = 0.246$). Notice also that Φ -statistics follow the same rules as Wright's F -statistics, namely

$$\begin{aligned}
 1 - \Phi_{ST} &= (1 - \Phi_{SC})(1 - \Phi_{CT}) \\
 0.754 &= (0.956)(0.78) \quad ,
 \end{aligned}$$

within the bounds of rounding error.³

³There wouldn't be any rounding error if we had access to the raw data.

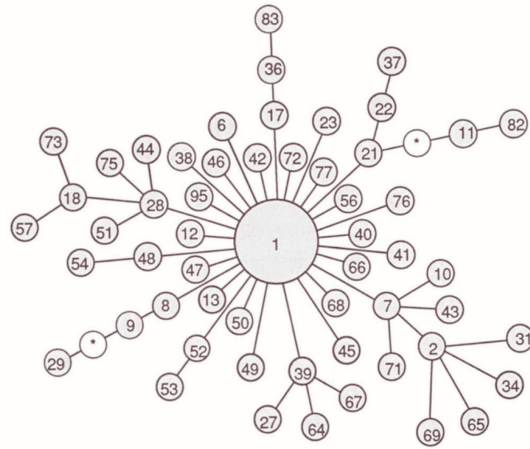


Figure 35.3: Minimum spanning network of human mtDNA samples in the example. The size of each circle is proportional to its frequency (from [23]).

Component of differentiation	Φ -statistics
Among regions	$\Phi_{CT} = 0.220$
Among populations within regions	$\Phi_{SC} = 0.044$
Among all populations	$\Phi_{ST} = 0.246$

Table 35.1: AMOVA results for the human mtDNA sample (from [23]).

An extension

As you may recall,⁴ Slatkin [83] pointed out that there is a relationship between coalescence time and F_{st} . Namely, if mutation is rare then

$$F_{ST} \approx \frac{\bar{t} - \bar{t}_0}{\bar{t}} \quad ,$$

where \bar{t} is the average time to coalescence for two genes drawn at random without respect to population and \bar{t}_0 is the average time to coalescence for two genes drawn at random from the same populations. Results in [42] show that when δ_{ij} is linearly proportional to the time since two sequences have diverged, Φ_{ST} is a good estimator of F_{ST} when F_{ST} is thought of as a measure of the relative excess of coalescence time resulting from dividing a species into several population. This observation suggests that the combination of haplotype frequency differences and evolutionary distances among haplotypes may provide insight into the evolutionary relationships among populations of the same species.

⁴Look back at <http://darwin.eeb.uconn.edu/eeb348/lecture-notes/coalescent/node6.html> for the details.

Chapter 36

Nested clade analysis

In a very influential paper Avise et al. [2] introduced the term “phylogeography” to refer to evolutionary studies lying at the interface of population genetics and phylogenetic systematics. An important property of molecular sequences is that the degree of difference among them contains information about their relatedness. Avise et al. proposed combining information derived from the phylogenetic relationship of molecular sequences with information about where the sequences were collected from to infer something about the biogeography of relationships among populations within species. Figure 36.1 provides an early and straightforward example.

The data are from bowfins, *Amia calva*, and consist of mtDNA haplotypes detected by restriction site mapping. There are two highly divergent groups of haplotypes separated from one another by a minimum of four restriction site differences. Moreover, the two sets of haplotypes are found in areas that are geographically disjunct. Haplotypes 1-9 are found exclusively in the eastern portion of the range, while haplotypes 10-13 are found exclusively in the western part of the range. This pattern suggests that the populations of bowfin in the two geographical regions have had independent evolutionary histories for a relatively long period of time. Interestingly, this disjunction between populations west and east of the Appalachian River is shared by a number of other species, as are disjunctions between the Atlantic and Gulf coasts, the west and east sides of the Tombigbee River, the west and east sides of the Appalachian mountains, and the west and east sides of the Mississippi River [84].

Early analyses often provided very clear patterns, like the one in bowfins. As data accumulated, however, it became clear that in some species it was necessary to account for differences in frequency, not just presence *versus* absence of particular haplotypes. We saw this in the application of AMOVA to mtDNA haplotype variation in humans. These approaches have two critical things in common:

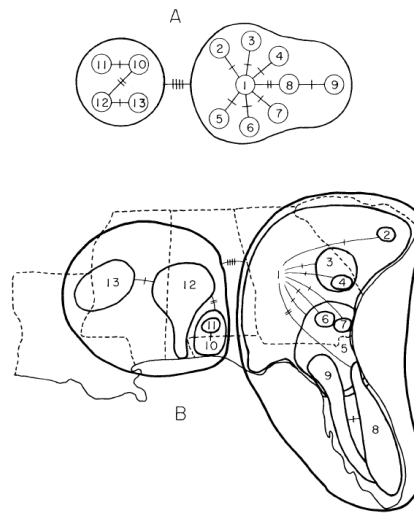


Figure 36.1: A phylogeographic analysis of 75 bowfins, *Amia calva*, sampled from the southeastern United States. **A.** A parsimony network connecting the 13 mtDNA haplotypes identified from the sample. **B.** The geographical distribution of the haplotypes.

- Haplotype networks are constructed as minimum-spanning (parsimony) networks without consideration as to whether assuming a parsimonious reconstruction of among haplotype differences is reasonable.¹
- The relationship between geographical distributions and haplotypes contains information about the history of those distributions, but there is no formal way to assess different interpretations of that history.

Nested-clade analysis (NCA) has become a widely used technique for phylogeographic analysis because it provides methods intended to assess each of those concerns [89].² In broad outline the ideas are pretty simple:

- Use statistical parsimony to construct a statistically supportable haplotype network.

¹For those of you who are familiar with molecular phylogenetics as it is usually applied, there's another important difference. Not only are these parsimony networks. They are networks in which some haplotypes are regarded as ancestral to others, i.e., haplotypes appear not only at the tips of a tree, but also at the nodes.

²It continues to produce networks in which some haplotypes are ancestral to others, but in this context, such an approach is reasonable. Ask me about it if you're interested in why it's reasonable.

- Identify nested clades, test for an association between geography and haplotype distribution, and work through an inference key to identify the processes that could have produced the association.

As we'll see, implementing these simple ideas poses some challenges.³

Statistical parsimony

Templeton et al. [90] lay out the theory and procedures involved in statistical parsimony in great detail. As with NCA in general, the details get a little complicated. We'll get to those complications soon enough, but again as with NCA in general the basic ideas are pretty simple:

- Evaluate the limits of parsimony, i.e., the number of mutational steps that can be reliably inferred without having to worry about multiple substitutions.
- Construct “the set of parsimonious and non-parsimonious cladograms that is consistent with these limits” (p. 619).⁴

So why use parsimony? Within species the time for substitutions to occur is relatively short. As a result, it may be reasonable to assume that we don't have to worry about multiple substitutions having occurred, at least between those haplotypes that are the most closely related. To “identify the limits of parsimony” we first estimate $\theta = 4N_e\mu$ from our data. Then we plug it into a formula that allows us to assess the probability that the difference between two randomly drawn haplotypes in our sample is the result of more than one substitution.⁵ If that probability is small, say less than 5%, we can connect all of the haplotypes into a parsimonious network, i.e., one that involves only single substitutions between haplotypes (some of which may be hypothetical and unobserved).

More likely than not, we won't be able to connect all of the haplotypes parsimoniously, but there's still a decent chance that we'll be able to identify subsets of the haplotypes for which the assumption of parsimonious change is reasonable. Templeton et al. [90] suggest the following procedure to construct a haplotype network:

³When we talk about statistical phylogeography, you'll see an alternative approach to addressing the concerns NCA was intended to address.

⁴Makes you wonder a little about why it's called statistical parsimony if some of the reconstructed cladograms aren't parsimonious, doesn't it?

⁵If you're interested, you can find the formula for restriction site differences in equation (1), p. 620.

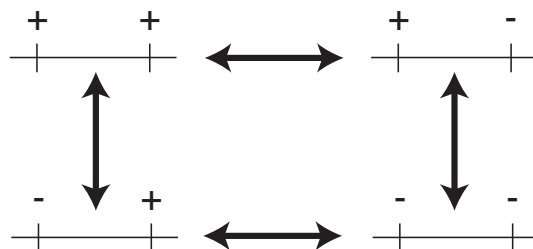


Figure 36.2: An example of four haplotypes connected in a single-step network showing that two paths are possible between haplotypes that differ in two positions.

Step 1: Estimate P_1 the probability that haplotype pairs differing by a single change are the result of a single substitution. If $P_1 > 0.95$, as is likely, connect all pairs of haplotypes that differ by a single change. There may be ambiguities in the reconstruction, including loops. Keep these in the network (Figure 36.2).

Step 2: Identify the products of recombination by inspecting the 1-step network to determine if postulating recombination between a pair of sequences can remove ambiguity identified in step 1.

Step 3: Augment j by one and estimate P_j . If $P_j > 0.95$, join $j - 1$ -step networks into a j -step network by connecting the two haplotypes that differ by j steps. Repeat until either all haplotypes are included in a single network or you are left with two or more non-overlapping networks.

Step 4: If you have two or more networks left to connect, estimate the smallest number of non-parsimonious changes that will occur with greater than 95% probability, and connect the networks.

Refer to Templeton et al. [90] for details of the calculations. Figure 36.3 provides an example of the kind of network that may result from this analysis.

Nested clade analysis

Once we have constructed the haplotype network, we're then faced with the problem of identifying nested clades. Templeton et al. [88] propose the following algorithm to construct a unique set of nested clades:

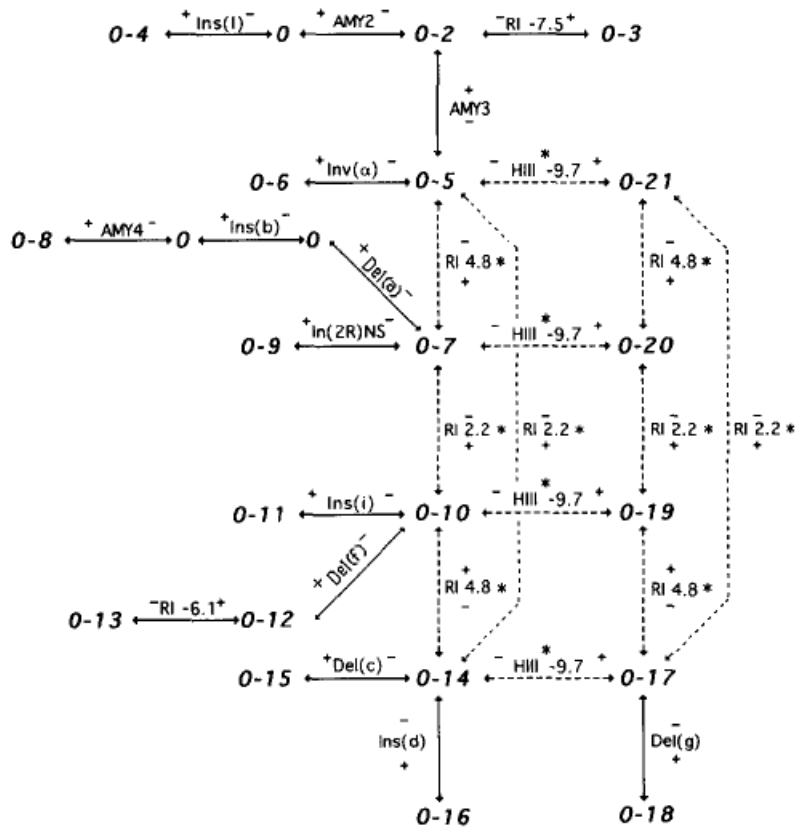


Figure 36.3: Statistical parsimony network for the *Amy* locus of *Drosophila melanogaster*.

- Step 1.** Each haplotype in the sample comprises a 0-step clade, i.e., each copy of a particular haplotype in the sample is separated by zero evolutionary steps from other copies of the same haplotype. “Tip” haplotypes are those that are connected to only one other haplotype. “Interior” haplotypes are those that are connected to two or more haplotypes. Set $j = 0$
- Step 2.** Pick a tip haplotype that is not part of any $j + 1$ -step network.
- Step 3.** Identify the interior haplotype with which it is connected by $j + 1$ mutational steps.
- Step 4.** Identify all tip haplotypes connected to that interior haplotype by $j + 1$ mutational steps.
- Step 5.** The set of all such tip and interior haplotypes constitutes a $j + 1$ -step clade.
- Step 6.** If there are tip haplotypes remaining that are not part of a $j + 1$ -step clade, return to step 2.
- Step 7.** Identify internal j -step clades that are not part of a $j + 1$ step clade and are separated by $j + 1$ steps.
- Step 8.** Designate these clades as “terminal” and return to step 2.
- Step 9.** Increment j by one and return to step 2.

That sounds fairly complicated, but if you look at the example in Figure 36.4, you’ll see that it isn’t all *that* horrible.

This algorithm produces a set of nested clades, i.e., a 1-step clade is contained within a 2-step clade, a 2-step clade is contained within a 3-step clade, and so on. Once such sets of nested clades have been identified, we can calculate statistics related to the geographical distribution of each clade in the sample. Templeton et al. [91] define two statistics that are used in an inferential key (the most recent version of the key is in [89]; see Figure 36.5):

Clade distance The average distance of each haplotype in the particular clade from the center of its geographical distribution. “Distance” may be the great circle distance or it might be the distance measured along a presumed dispersal corridor. The clade distance for clade X is symbolized $D_c(X)$, and it measures how far this clade has spread.

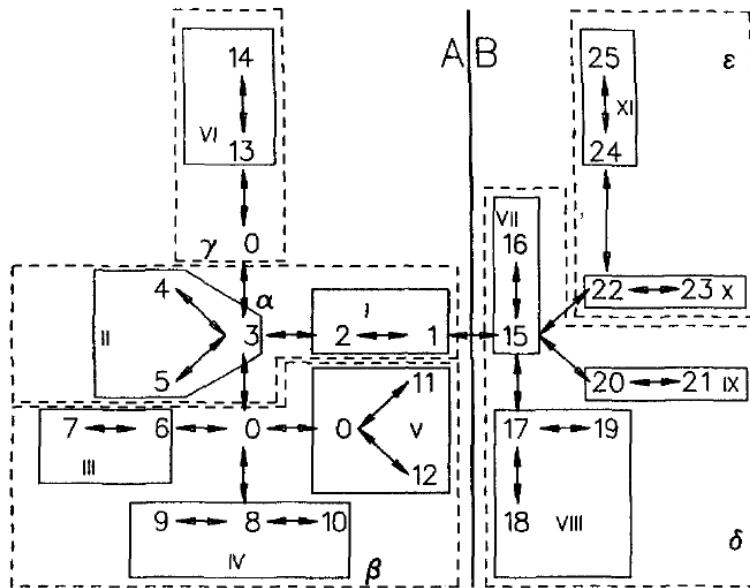


Figure 36.4: Nesting of haplotypes at the *Adh* locus in *Drosophila melanogaster*.

Nested clade distance The average distance of the center of distribution for this haplotype from the center of distribution for the haplotype within which it is nested. So if clade X is nested within clade Y , we calculate $D_n(X)$ by determining the geographic center of clades X and clade Y and measuring the distance between those centers. $D_n(X)$ measures how far the clade has changed position relative to the clade from which it originated.

Once you've calculated these distances, you randomly permute the clades across sample locations. This shuffles the data randomly, keeping the number of haplotypes and the sample size per location the same as in the original data set. For each of these permutations, you calculate $D_c(X)$ and $D_n(X)$. If the observed clade distance, the observed nested clade difference, or both are significantly different from expected by chance, then you have evidence of (a) geographical expansion of the clade (if $D_c(X)$ is greater than null expectation) or (b) a range-shift (if $D_n(X)$ is greater than null expectation). Using these kinds of statistics, you run your data set through Templeton's inference key to reach a conclusion. For example, applying this procedure to data from *Ambystoma tigrinum* (Figure 36.6), Templeton et al. [91] construct the scenario in Figure 36.7.

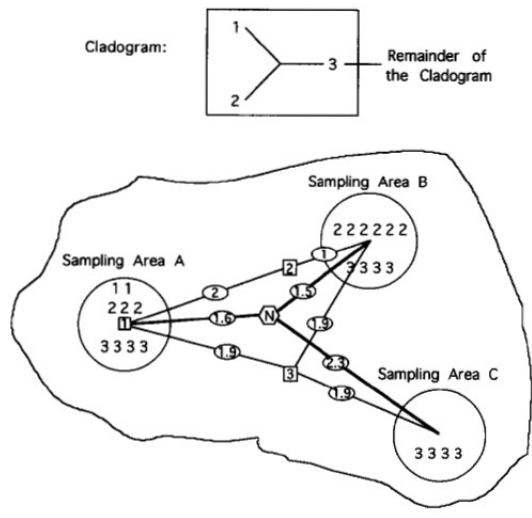


Figure 36.5: Each number corresponds to a haplotype in the sample. Haplotypes 1 and 2 are “tip” haplotypes. Haplotype 3 is an interior haplotype. The numbers in square boxes illustrate the center for each 0-step clade (a haplotype). The hexagonal “N” represents the center for the clade containing 1, 2, and 3. Numbers in ovals are the distances from the center of each collecting area to the clade center. $D_c(1) = 0$, $D_c(2) = (3/9)(2) + (6/9)(1) = 1.33$, $D_c(3) = (4/12)(1.9) + (4/12)(1.9) + (4/12)(1.9) = 1.9$. $D_n(1) = 1.6$, $D_n(2) = (3/9)(1.6) + (6/9)(1.5) = 1.53$, $D_n(3) = (4/12)(1.6) + (4/12)(1.5) + (4/12)(2.3) = 1.8$.

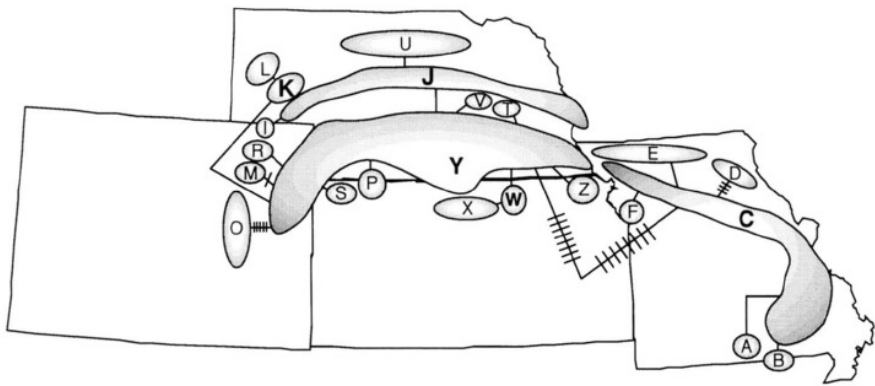


Figure 36.6: Geographic distribution of mtDNA haplotypes in *Ambystoma tigrinum*.

Clade	Chain of inference	Inference
Haplotypes nested in 1-1	1-2-3-5-6-13-14 NO	Range expansion, but cannot discriminate between contiguous range expansion and long-distance colonization
Haplotypes nested in 1-2	1-2-3-4 NO	Restricted gene flow via isolation by distance
One-step clades nested in 2-1	1-2-3-4 NO	Restricted gene flow via isolation by distance
One-step clades nested in 2-2	1-2-11-12 NO	Contiguous range expansion
Two-step clades nested in 3-2	1-2-3-4 NO	Restricted gene flow via isolation by distance
>Four-step clades nested in entire cladogram	1-2-3-5-9 NO and associated with longest branch length	Allopatric fragmentation

Figure 36.7: Inference key for *Ambystoma tigrinum*.

Chapter 37

Statistical phylogeography

Nested clade analysis represented the earliest attempt to develop a formal approach to using an estimate of phylogenetic relationships among haplotypes to infer something both about the biogeographic history of the populations in which they are contained and the evolutionary processes associated with the pattern of diversification implied by the phylogenetic relationships among haplotypes and their geographic distribution. The statistical parsimony part of NCA depends heavily on coalescent theory for calculating the “limits” of parsimony. As a result, NCA combines aspects of pure phylogenetic inference — parsimony — with aspects of pure population genetics — coalescent theory — to develop a set of inferences about the phylogeographic history of populations within species. But well before NCA was developed, Pamilo and Nei [76] pointed out that the phylogenetic relationships of a single gene might be different from those of the populations from which the samples were collected.

Gene trees *versus* population trees

There are several reasons why *gene trees* might not match *population trees*.

- It could simply be a problem of estimation. Given a particular set of gene sequences, we *estimate* a phylogenetic relationship among them. But our estimate could be wrong. In fact, given the astronomical number of different trees possible with 50 or 60 distinct sequences, it’s virtually certain to be wrong somewhere. We just don’t know where. It could be that if we had the right gene tree it would match the species tree.
- There might have been a hybridization event in the past so that the phylogenetic history of the gene we’re studying is different from that of the populations from which we sampled. Hybridization is especially likely to have a large impact if the locus for

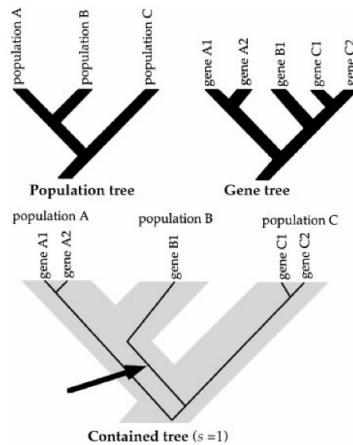


Figure 37.1: Discordance between gene and population trees as a result of ancestral polymorphism (from [54]).

which we have information is uniparentally inherited, e.g., mitochondrial or chloroplast DNA. A single hybridization event in the distant past in which the maternal parent was from a different population will give mtDNA or cpDNA a very different phylogeny than nuclear genes that underwent a lot of backcrossing after the hybridization event.

- If the ancestral population was polymorphic at the time the initial split occurred alleles that are more distantly related might, by chance, end up in the same descendant population (see Figure 37.1)

As Pamilo and Nei showed, it's possible to calculate the probability of discordance between the gene tree and the population tree using some basic ideas from coalescent theory. That leads to a further refinement, using coalescent theory directly to examine alternative biogeographic hypotheses.

Phylogeography of montane grasshoppers

Lacey Knowles studied grasshoppers in the genus *Melanopus*. She collected 1275bp of DNA sequence data from cytochrome oxidase I (COI) from 124 individuals of *M. oregonensis* and two outgroup species. The specimens were collected from 15 “sky-island” sites in the northern Rocky Mountains (see Figure 37.2; [54]). Two alternative hypotheses had been proposed to

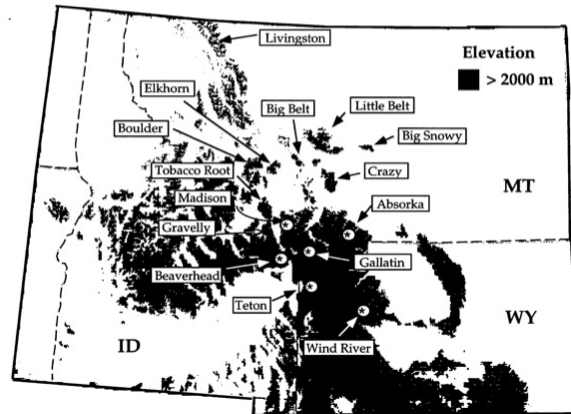


Figure 37.2: Collection sites for *Melanopus oregonensis* in the northern Rocky Mountains (from [54]).

describe the evolutionary relationships among these grasshoppers (refer to Figure 37.3 for a pictorial representation):

- **Widespread ancestor:** The existing populations might represent independently derived remnants of a single, widespread population. In this case all of the populations would be equally related to one another.
- **Multiple glacial refugia:** Populations that shared the same refugium will be closely related while those that were in different refugia will be distantly related.

As is evident from Figure 37.3, the two hypotheses have very different consequences for the coalescent history of alleles in the sample. Since the interrelationships between divergence times and time to common ancestry differ so markedly between the two scenarios, the pattern of sequence differences found in relation to the geographic distribution will differ greatly between the two scenarios.

Using techniques described in Knowles and Maddison [55], Knowles simulated gene trees under the widespread ancestor hypothesis. She then placed them within a population tree representing the multiple glacial refugia hypothesis and calculated a statistic, s , that measures the discordance between a gene tree and the population tree that contains it. This gave her a distribution of s under the widespread ancestor hypothesis. She compared the s estimated from her actual data with this distribution and found that the observed value of

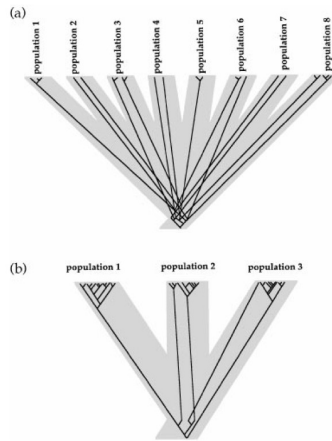


Figure 37.3: Pictorial representations of the “widespread ancestor” (top) and “glacial refugia” (bottom) hypotheses (from [54]).

s was only $1/2$ to $1/3$ the size of the value observed in her simulations.¹ In short, Knowles presented strong evidence that her data are not consistent with the widespread ancestor hypothesis.

¹The discrepancy was largest when divergence from the widespread ancestor was assumed to be very recent.

Chapter 38

Fully coalescent-based approaches to phylogeography

Last time we saw an early example of using coalescent theory to distinguish between two scenarios describing the history of populations. In the example we considered, Knowles [54] compared two scenarios, the “widespread ancestor” and the “multiple glacial refugia” scenarios. To make the comparison she simulated data under the “widespread ancestor” hypothesis, collected the samples into a multiple-refuge tree, and calculated a statistic that measures the discrepancy between the gene trees and the population trees. Her observed gene tree was far less discordant than the simulated trees, leading her to conclude that her grasshoppers had been dispersed among multiple refugia in the past rather than being the remnants of a single, widespread ancestral population. As I mentioned, one limitation of the approach Knowles [54] takes is that it requires the investigator to identify alternative scenarios before beginning the analysis, and it can only identify which of the scenarios is more likely than the others with which it is compared. It cannot determine whether there are other scenarios that are even more likely. Another approach is to back off a bit, specify a particular process that we are interested in and to use what we know about that process to try and estimate its properties.

Coalescent-based estimates of migration rate

A few years before Knowles [54] appeared Beerli and Felsenstein [7, 8] proposed a coalescent-based method to estimate migration rates among populations. As with other analytical methods we’ve encountered in this course, the details can get pretty hairy, but the basic idea is (relatively) simple.

Recall that in a single population we can describe the coalescent history of a sample without too much difficulty. Specifically, given a sample of n alleles in a diploid population with effective size N_e , the probability that the first coalescent event took place t generations ago is

$$P(t|n, N_e) = \left(\frac{n(n-1)}{4N_e}\right) \left(1 - \frac{n(n-1)}{4N_e}\right)^{t-1}. \quad (38.1)$$

Now suppose that we have a sample of alleles from K different populations. To keep things (relatively) simple, we'll imagine that we have a sample of n alleles from every one of these populations and that every population has an effective size of N_e . In addition, we'll imagine that there is migration among populations, but again we'll keep it really simple. Specifically, we'll assume that the probability that a given allele in our sample from one population had its ancestor in a different population in the immediately preceding generation is m .¹ Under this simple scenario, we can again construct the coalescent history of our sample. How? Funny you should ask.

We start by using the same logic we used to construct equation (38.1). Specifically, we ask "What's the probability of an 'event' in the immediately preceding generation?" The complication is that there are two kinds of events possible: (1) a coalescent event and (2) a migration event. As in our original development of the coalescent process, we'll assume that the population sizes are large enough that the probability of two coalescent events in a single time step is so small as to be negligible. In addition, we'll assume that the number of populations and the migration rates are small enough that the probability of more than one event of either type is so small as to be negligible. That means that all we have to do is to calculate the probability of either a coalescent event or a migration event and combine them to calculate the probability of an event. It turns out that it's easiest to calculate the probability that there *isn't* an event first and then to calculate the probability that there is an event as one minus that.

We already know that the probability of a coalescent event in population k , is

$$P_k(\text{coalescent}|n, N_e) = \frac{n(n-1)}{4N_e},$$

so the probability that there is *not* a coalescent event in any of our K populations is

$$P(\text{no coalescent}|n, N_e, K) = \left(1 - \frac{n(n-1)}{4N_e}\right)^K.$$

¹In other words, m is the backwards migration rate, the probability that a gene in one population came from another population in the preceding generation. This is the same migration rate we encountered weeks ago when we discussed the balance between drift and migration.

If m is the probability that there was a migration event in a particular population than the probability that there is *not* a migration event involving any of our nK alleles² is

$$P(\text{no migration}|m, K) = (1 - m)^{nK} \quad .$$

So the probability that there *is* an event of some kind is

$$P(\text{event}|n, m, N_e, K) = 1 - P(\text{no coalescent}|n, N_e, K)P(\text{no migration}|m, K) \quad .$$

Now we can calculate the time back to the first event

$$P(\text{event at } t|n, m, N_e, K) = P(\text{event}|n, m, N_e, K) (1 - P(\text{event}|n, m, N_e, K))^{t-1} \quad .$$

We can then use Bayes theorem to calculate the probability that the event was a coalescence or a migration and the populations involved. Once we've done that, we have a new population configuration and we can start over. We continue until all of the alleles have coalesced into a single common ancestor, and then we have the complete coalescent history of our sample.³ That's roughly the logic that Beerli and Felsenstein use to construct coalescent histories for a sample of alleles from a set of populations—except that they allow effective population sizes to differ among populations and they allow migration rates to differ among all pairs of populations. As if that weren't bad enough, now things start to get even more complicated.

There are lots of different coalescent histories possible for a sample consisting of n alleles from each of K different populations, even when we fix m and N_e . Worse yet, given any one coalescent history, there are a lot of different possible mutational histories possible. In short, there are a lot of different possible sample configurations consistent with a given set of migration rates and effective population size. Nonetheless, some combinations of m and N_e will make the data more likely than others. In other words, we can construct a likelihood for our data:

$$P(\text{data}|m, N_e) \propto f(n, m, N_e, K) \quad ,$$

where $f(n, m, N_e, K)$ is some very complicated function of the probabilities we derived above. In fact, the function is so complicated, we can't even write it down. Beerli and Felsenstein, being very clever people, figured out a way to simulate the likelihood, and **Migrate** provides a (relatively) simple way that you can use your data to estimate m and N_e for a set of populations. In fact, **Migrate** will allow you to estimate pairwise migration rates among all populations in your sample, and since it can simulate a likelihood, if you put priors on

² K populations each with n alleles

³This may not seem very simple, but just think about how complicated it would be if I allowed every population to have a different effective size and if I allowed each pair of populations to have different migration rates between them.

the parameters you're interested in, i.e., m and N_e , you can get Bayesian estimates of those parameters rather than maximum likelihood estimates, including credible intervals around those estimates so that you have a good sense of how reliable your estimates are.⁴

There's one further complication I need to mention, and it involves a lie I just told you. **Migrate** can't give you estimates of m and N_e . Remember how every time we've dealt with drift and another process we always end up with things like $4N_e m$, $4N_e \mu$, and the like. Well, the situation is no different here. What **Migrate** can actually estimate are the two parameters $4N_e m$ and $\theta = 4N_e \mu$.⁵ How did μ get in here when I only mentioned it in passing? Well, remember that I said that once the computer has constructed a coalescent history, it has to apply mutations to that history. Without mutation, all of the alleles in our sample would be identical to one another. Mutation is what produces the diversity. So what we get from **Migrate** isn't the fraction of a population that's composed of migrants. Rather, we get an estimate of how much migration contributes to local population diversity relative to mutation. That's a pretty interesting estimate to have, but it may not be everything that we want.

There's a further complication to be aware of. Think about the simulation process I described. All of the alleles in our sample are descended from a single common ancestor. That means we are implicitly assuming that the set of populations we're studying have been around long enough and have been exchanging migrants with one another long enough that we've reached a drift-mutation-migration equilibrium. If we're dealing with a relatively small number of populations in a geographically limited area, that may not be an unreasonable assumption, but what if we're dealing with populations of crickets spread across all of the northern Rocky Mountains? And what if we haven't sampled all of the populations that exist?⁶ In many circumstances, it may be more appropriate to imagine that populations diverged from one another at some time in the not too distant past, have exchanged genes since their divergence, but haven't had time to reach a drift-mutation-migration equilibrium. What do we do then?

⁴If you'd like to see a comparison of maximum likelihood and Bayesian approaches, Beerli [5] provides an excellent overview.

⁵Depending on the option you pick when you run **Migrate** you can either get θ and $4N_e m$ or θ and $M = m/\mu$.

⁶Beerli [6] discusses the impact of "ghost" populations. He concludes that you have to be careful about which populations you sample, but that you don't necessarily need to sample every population. Read the paper for the details.

Divergence and migration

Nielsen and Wakely [69] consider the simplest generalization of Beerli and Felsenstein [7, 8] you could imagine (Figure 38.1). They consider a situation in which you have samples from only two populations and you're interested in determining both how long ago the populations diverged from one another and how much gene exchange there has been between the populations since they diverged. As in **Migrate** mutation and migration rates are confounded with effective population size, and the relevant parameters become:

- θ_a , which is $4N_e\mu$ in the ancestral population.
- θ_1 , which is $4N_e\mu$ in the first population.
- θ_2 , which is $4N_e\mu$ in the second population.
- M_1 , which is $2N_em$ in the first population, where m is the fraction of the first population composed of migrants from the second population.
- M_2 , which is $2N_em$ in the second population.
- T , which is the time since the populations diverged. Specifically, if there have been t units since the two populations diverged, $T = t/2N_1$, where N_1 is the effective size of the first population.

Given that set of parameters, you can probably imagine that you can calculate the likelihood of the data for a given set of parameters.⁷ Once you can do that you can either obtain maximum-likelihood estimates of the parameters by maximizing the likelihood, or you can place prior distributions on the parameters and obtain Bayesian estimates from the posterior distribution. Either way, armed with estimates of θ_a , θ_1 , θ_2 , M_1 , M_2 , and T you can say something about: (1) the effective population sizes of the two populations relative to one another and relative to the ancestral population, (2) the relative frequency with which migrants enter each of the two populations from the other, and (3) the time at which the two populations diverged from one another. Keep in mind, though, that the estimates of M_1 and M_2 confound local effective population sizes with migration rates. So if $M_1 > M_2$, for example, it does not mean that the fraction of migrants incorporated into population 1 exceeds the fraction incorporated into population 2. It means that the impact of migration has been felt more strongly in population 1 than in population 2.

⁷As with **Migrate**, you can't calculate the likelihood explicitly, but you can approximate it numerically. See [69] for details.

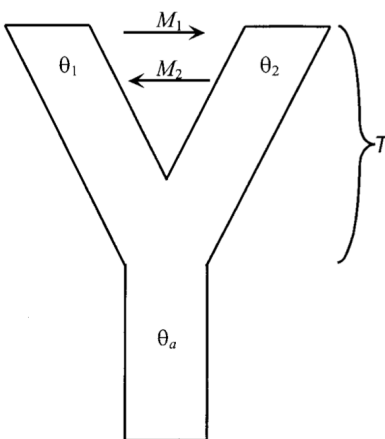


Figure 38.1: The simple model developed by Nielsen and Wakeley [69]. θ_a is $4N_e\mu$ in the ancestral population; θ_1 and θ_2 are $4N_e\mu$ in the descendant populations; M_1 and M_2 are $2N_e m$, where m is the backward migration rate; and T is the time since divergence of the two populations.

An example

Orti et al. [75] report the results of phylogenetic analyses of mtDNA sequences from 25 populations of threespine stickleback, *Gasterosteus aculeatus*, in Europe, North America, and Japan. The data consist of sequences from a 747bp fragment of cytochrome *b*. Nielsen and Wakeley [69] analyze these data using their approach. Their analyses show that “[a] model of moderate migration and very long divergence times is more compatible with the data than a model of short divergence times and low migration rates.” By “very long divergence times” they mean $T > 4.5$, i.e., $t > 4.5N_1$. Focusing on populations in the western (population 1) and eastern Pacific (population 2), they find that the maximum likelihood estimate of M_1 is 0, indicating that there is little if any gene flow from the eastern Pacific (population 2) into the western Pacific (population 1). In contrast, the maximum likelihood estimate of M_2 is about 0.5, indicating that one individual is incorporated into the eastern Pacific population from the western Pacific population every other generation. The maximum-likelihood estimates of θ_1 and θ_2 indicate that the effective size of the population eastern Pacific population is about 3.0 times greater than that of the western Pacific population.

Extending the approach to multiple populations

A couple of years ago, Jody Hey announced the release of **IMa2**. Building on work described in Hey and Nielsen [38, 39], **IMa2** allows you to estimate relative divergence times, relative effective population sizes, and relative pairwise migration rates for more than two populations at a time. That flexibility comes at a cost, of course. In particular, you have to specify the phylogenetic history of the populations before you begin the analysis.

Chapter 39

Approximate Bayesian Computation

Just when you thought it was safe to go back into the water, I'm going to complicate things even further.¹ The Nielsen-Wakely-Hey [69, 38, 39] approach is *very* flexible and *very* powerful, but even it doesn't cover all possible scenarios. It allows for non-equilibrium scenarios in which the populations from which we sampled diverged from one another at different times, but suppose that we think our populations have dramatically increased in size over time (as in humans) or dramatically changed their distribution (as with an invasive species). Is there a way to use genetic data to gain some insight into those processes? Would I be asking that question if the answer were "No"?

An example

Let's change things up a bit this time and start with an example of a problem we'd like to solve first. Once you see what the problem is, then we can talk about how we might go about solving it. Let's talk about the case of the cane toad, *Bufo marinus*, in Australia.

You may know that the cane toad is native to the American tropics. It was purposely introduced into Australia in 1935 as a biocontrol agent, where it has spread across an area of more than 1 million km². Its range is still expanding in northern Australia and to a lesser extent in eastern Australia (Figure 39.1).² Estoup et al. [21] Collected microsatellite data from 30 individuals in each of 19 populations along roughly linear transects in the northern and eastern expansion areas.

¹Look on the bright side. The semester is nearly over. Besides, you need to know a little about approximate Bayesian computation in order to write up your final problem.

²All of this information is from the introduction to [21].

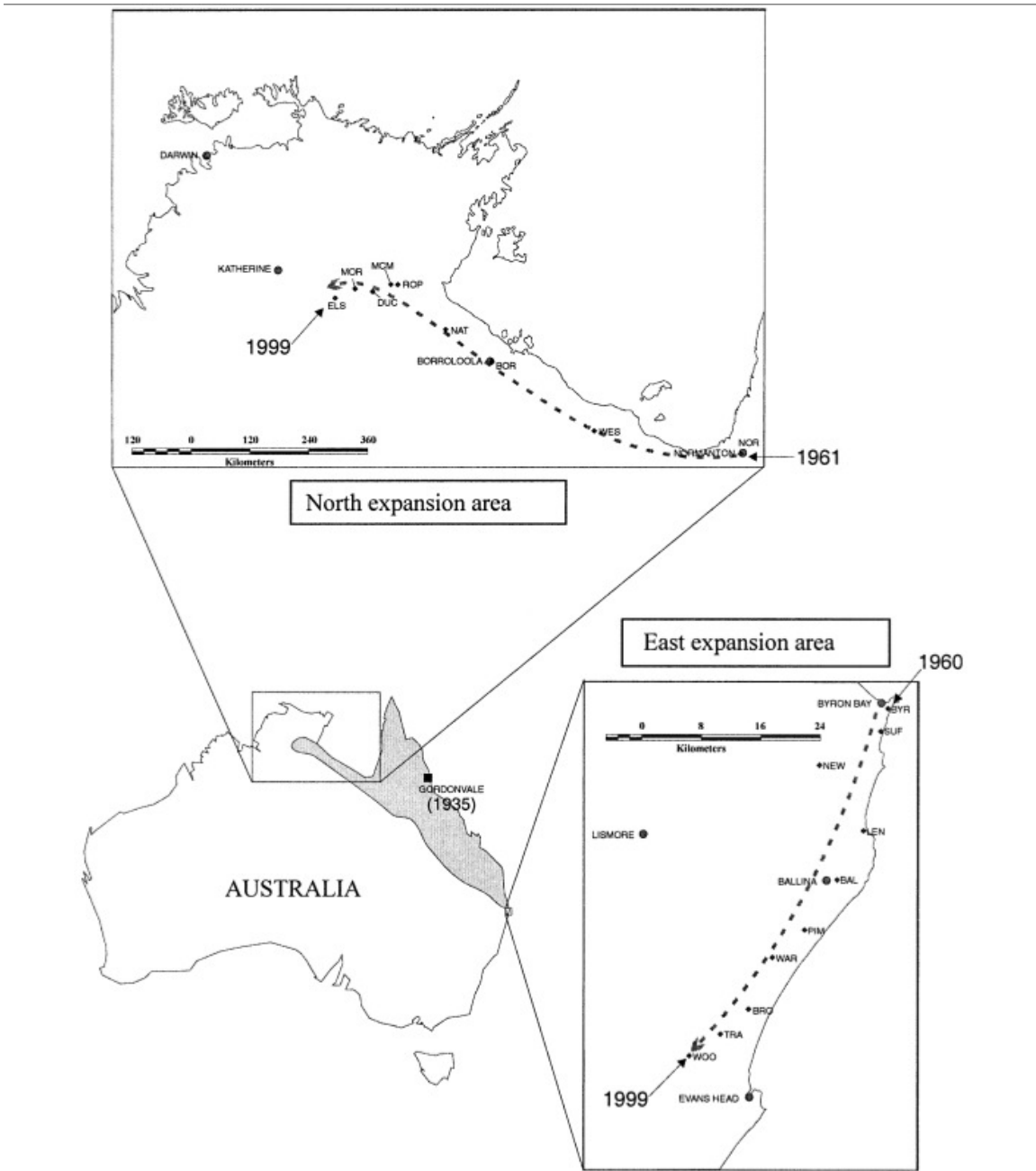


Figure 39.1: Maps showing the expansion of the cane toad population in Australia since its introduction in 1935 (from [21]).

With these data they wanted to distinguish among five possible scenarios describing the geographic spread:

- **Isolation by distance:** As the expansion proceeds, each new population is founded by or immigrated into by individuals with a probability proportional to the distance from existing populations.
- **Differential migration and founding:** Identical to the preceding model except that the probability of founding a population may be different from the probability of immigration into an existing population.
- **“Island” migration and founding:** New populations are established from existing populations without respect to the geographic distances involved, and migration occurs among populations without respect to the distances involved.
- **Stepwise migration and founding with founder events:** Both migration and founding of populations occurs only among immediately adjacent populations. Moreover, when a new population is established, the number of individuals involved may be very small.
- **Stepwise migration and founding without founder events:** Identical to the preceding model except that when a population is founded its size is assumed to be equal to the effective population size.

That’s a pretty complex set of scenarios. Clearly, you could use `Migrate` or `IMa2` to estimate parameters from the data Estoup et al. [21] report, but would those parameters allow you to distinguish those scenarios? Not in any straightforward way, at least. Neither `Migrate` nor `IMa2` distinguishes between founding and migration events for example. And with `IMa2` we’d have to specify the relationships among our sampled populations before we could make any of the calculations. But how would we arrive at the hypothesis of relationships to use? So what do we do?

Approximate Bayesian Computation

Well, in principle we could take an approach similar to what `Migrate` and `IMa2` use. Let’s start by reviewing what we did last time³ with `Migrate` and `IMa2`. In both cases, we knew how to simulate data given a set of mutation rates, migration rates, local effective population

³More accurately, what Peter Beerli, Joe Felsenstein, Rasmus Nielsen, John Wakeley, and Jody Hey did

sizes, and times since divergence. Let's call that whole, long string of parameters ϕ and our big, complicated data set X . If we run enough simulations, we can keep track of how many of those simulations produce data identical to the data we collected. With those results in hand, we can estimate $P(X|\phi)$, the likelihood of the data, as the fraction of simulations that produce data identical to the data we collected.⁴ In principle, we could take the same approach in this, much more complicated, situation. But the problem is that there are an astronomically large number of different possible coalescent histories and different allelic configurations possible with any one population history both because the population histories being considered are pretty complicated and because the coalescent history of every locus will be somewhat different from the coalescent history at other loci. As a result, the chances of getting *any* simulated samples that match our actual samples is virtually nil, and we can't estimate $P(X|\phi)$ in the way we have so far.

Approximate Bayesian computation is an approach that allows us to get around this problem. It was introduced by Beaumont et al. [4] precisely to allow investigators to get approximate estimates of parameters and data likelihoods in a Bayesian framework. Again, the details of the implementation get pretty hairy,⁵ but the basic idea is relatively straightforward.⁶

1. Calculate “appropriate” summary statistics for your data set, e.g., pairwise estimates of ϕ_{ST} (possibly one for every locus if you're using microsatellite or SNP data), estimates of within population diversity, counts of the number of segregating sites (for nucleotide sequence data, both within each population and across the entire sample). Call that set of summary statistics S .
2. Specify a prior distribution for the unknown parameters, ϕ .
3. Pick a random set of parameter values, ϕ' from the prior distribution and simulate a data set for that set of parameter values.
4. Calculate the same summary statistics for the simulated data set as you calculated for your actual data. Call that set of statistics S' .
5. Calculate the Euclidean distance between S and S' . Call it δ . If it's less than some value you've decided on, δ^* , keep track of S' and the associated ϕ' and δ . Otherwise, throw all of them away and forget you ever saw them.

⁴The actual implementation is a bit more involved than this, but that's the basic idea.

⁵You're welcome to read the Methods in [4], and feel free to ask questions if you're interested.

⁶OK. This maybe calling it “relatively straightforward” is misleading. Even this simplified outline is fairly complicated, but compared to some of what you've already survived in this course, it may not look too awful.

6. Return to step 2 and repeat until you have accepted a large number of pairs of S' and ϕ' .

Now you have a bunch of S' 's and a bunch of ϕ' 's that produced them. Let's label them S_i and ϕ_i , and let's remember what we're trying to do. We're trying to estimate ϕ for our real data. What we have from our real data is S . So far it seems as if we've worked our computer pretty hard, but we haven't made any progress.

Here's where the trick comes in. Suppose we fit a regression to the data we've simulated

$$\phi_i = \alpha + S_i\beta + \epsilon \quad ,$$

where α is an intercept, β is a vector of regression coefficients relating each of the summary statistics to ϕ , and ϵ is an error vector.⁷ Now we can use that regression relationship to predict what ϕ should be in our real data, namely

$$\phi = \alpha + S\beta \quad .$$

If we throw in some additional bells and whistles, we can approximate the posterior distribution of our parameters. With that we can get not only a point estimate for ϕ , but also credible intervals for all of its components.

Back to the real world⁸

OK. So now we know how to do ABC, how do we apply it to the cane toad data. Well, using the additional bells and whistles I mentioned, we end up with a whole distribution of δ for each of the scenarios we try. The scenario with the smallest δ provides the best fit of the model to the data. In this case, that corresponds to model 4, the stepwise migration with founder model, although it is only marginally better than model 1 (isolation by distance) and model 2 (isolation by distance with differential migration and founding) in the northern expansion area (Figure 39.2).

Of course, we also have estimates for various parameters associated with this model:

- N_{e_s} : the effective population size when the population is stable.

⁷I know what you're thinking to yourself now. This doesn't sound very simple. Trust me. It is as simple as I can make it. The actual procedure involves local linear regression. I'm also not telling you how to go about picking δ or how to pick "appropriate" summary statistics. There's a fair amount of "art" involved in that.

⁸Or at least something resembling the real world

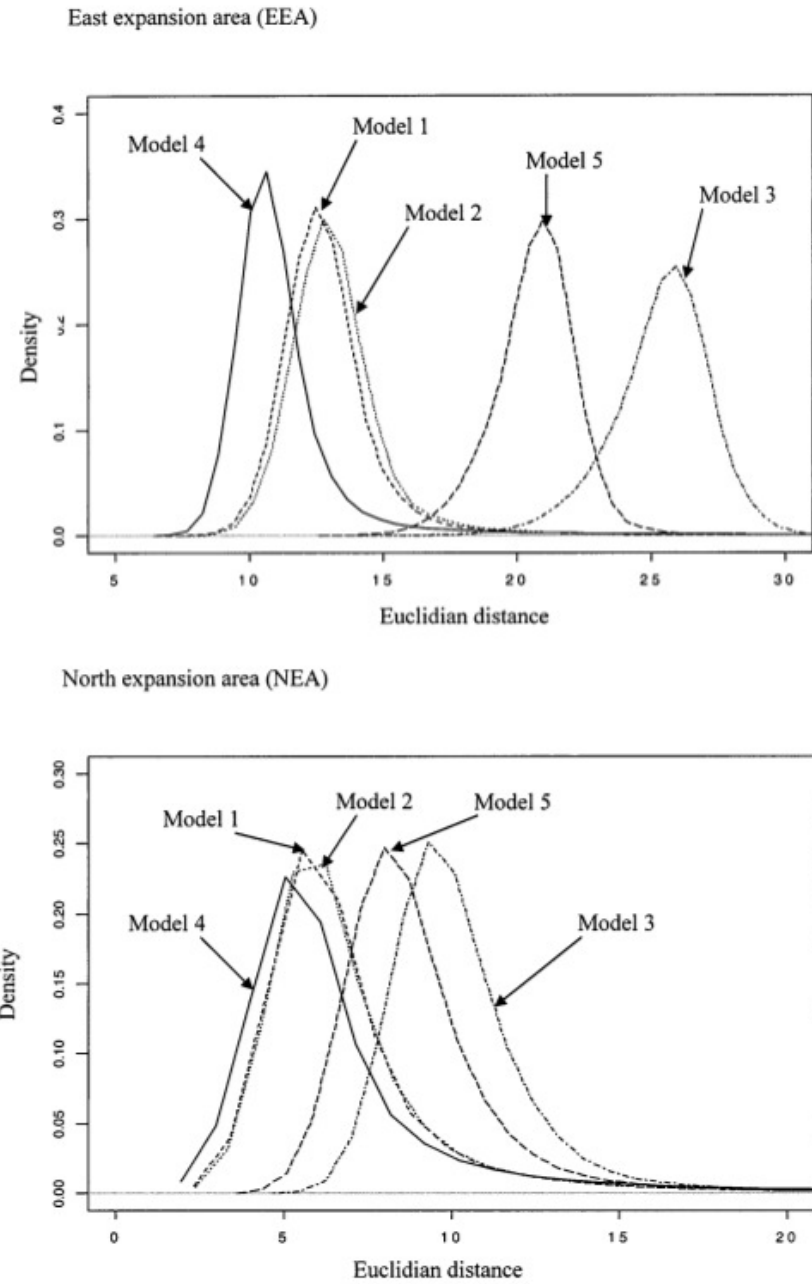


Figure 39.2: Posterior distribution of δ for the five models considered in Estoup et al. [21].

- N_{ef} : the effective population size when a new population is founded.
- F_R : the founding ratio, N_{es}/N_{ef} .
- m : the migration rate.
- $N_{es}m$: the effective number of migrants per generation.

The estimates are summarized in Table 39.1. Although the credible intervals are fairly broad,⁹ There are a few striking features that emerge from this analysis.

- Populations in the northern expansion area are larger, than those in the eastern expansion region. Estoup et al. [21] suggest that this is consistent with other evidence suggesting that ecological conditions are more homogeneous in space and more favorable to cane toads in the north than in the east.
- A smaller number of individuals is responsible for founding new populations in the east than in the north, and the ratio of “equilibrium” effective size to the size of the founding population is bigger in the east than in the north. (The second assertion is only weakly supported by the results.)
- Migration among populations is more limited in the east than in the north.

As Estoup et al. [21] suggest, results like these could be used to motivate and calibrate models designed to predict the future course of the invasion, incorporating a balance between gene flow (which can reduce local adaptation), natural selection, drift, and colonization of new areas.

Limitations of ABC

If you’ve learned anything by now, you should have learned that there is no perfect method. An obvious disadvantage of ABC relative to either `Migrate` or `IMa2` is that it is much more computationally intensive.

- Because the scenarios that can be considered are much more complex, it simply takes a long time to simulate all of the data.

⁹And notice that these are 90% credible intervals, rather than the conventional 95% credible intervals, which would be even broader.

Parameter	area	mean (5%, 90%)
N_{e_s}	east	744 (205, 1442)
	north	1685 (526, 2838)
N_{e_f}	east	78 (48, 118)
	north	311 (182, 448)
F_R	east	10.7 (2.4, 23.8)
	north	5.9 (1.6, 11.8)
m	east	0.014 (6.0×10^{-6} , 0.064)
	north	0.117 (1.4×10^{-4} , 0.664)
$N_{e_s}m$	east	4.7 (0.005, 19.9)
	north	188 (0.023, 883)

Table 39.1: Posterior means and 90% credible intervals for parameters of model 4 in the eastern and northern expansion areas of *Bufo marinus*.

- In the last few years, one of the other disadvantages—that you had to know how to do some moderately complicated scripting to piece together several different packages in order to run analysis—has become less of a problem. `popABC` (<http://code.google.com/p/popabc/>) and `DIYABC` (<http://www1.montpellier.inra.fr/CBGP/diyabc/>) make it *relatively* easy¹⁰ to perform the simulations.
- Selecting an appropriate set of summary statistics isn't easy, and it turns out that which set is most appropriate may depend on the value of the parameters that you're trying to estimate and the which of the scenarios that you're trying to compare is closest to the actual scenario applying to the populations from which you collected the data. Of course, if you knew what the parameter values were and which scenario was closest to the actual scenario, you wouldn't need to do ABC in the first place.
- In the end, ABC allows you to compare a small number of evolutionary scenarios. It can tell you which of the scenarios you've imagined provides the best combination of fit to the data and parsimonious use of parameters (if you choose model comparison statistics that include both components), but it takes additional work to determine whether the model is adequate, in the sense that it does a good job of explaining the data. Moreover, even if you determine that the model is adequate, you can't exclude the possibility that there are other scenarios that might be equally adequate—or even better.

¹⁰Emphasis on “relatively”.

Chapter 40

Population genomics

In the past few years, the development of high-throughput methods for genomic sequencing (next-generation sequencing: NGS) have drastically changed how many biologists work. It is now possible to produce so much data so rapidly that simply storing and processing the data poses great challenges, so great that *Nature Reviews Genetics* recently published an overview of approaches to analysis, interpretation, reproducibility, and accessibility of NGS data [68]. Their review didn't even discuss the new challenges that face population geneticists and evolutionary biologists as they start to take advantage of those tools, nor did it discuss the promise these data hold for providing new insight into long-standing, but the challenges and the promise are at least as great as those they do describe.

To some extent the most important opportunity provided by NGS sequencing is simply that we now have a lot more data to answer the same questions. For example, using a technique like RAD sequencing [3], it is now possible to identify thousands of polymorphic SNP markers in non-model organisms, even if you don't have a reference sequence available. As we've seen several times this semester, the variance associated with drift is enormous. Every SNP identified through RAD sequencing is likely to be independently inherited, thus the amount and pattern of variation at each locus will represent an independent sample from the underlying evolutionary process. As a result, we should be able to get much better estimates of fundamental parameters like $\theta = 4N_e\mu$, $M = 4N_em$, and $R = 4N_er$ and to have much greater power to discriminate among different evolutionary scenarios. Willing et al. [99], for example, present simulations suggesting that accurate estimates of F_{ST} are possible with sample sizes as small as 4–6 individuals per population, so long as the number of markers used for inference is greater than 1000.

Next-generation phylogeography

The American pitcher plant mosquito *Wyeomyia smithii* has been extensively studied for many years. It's a model organism for ecology, but its genome has not been sequenced. An analysis of *COI* from 20 populations and two outgroups produced the set of relationships you see in Figure 40.1 [20]. As you can see, this analysis allows us to distinguish a northern group of populations from a southern group of populations, but it doesn't provide us any reliable insight into finer scale relationships.

Using the same set of samples, the authors used RAD sequencing to identify 3741 SNPs. That's more than 20 times the number of variable sites found in *COI*, 167. Not surprisingly, the large number of additional sites allowed the authors to produce a much more highly resolved phylogeny 40.2. With this phylogeny it's easy to see that southern populations are divided into two distinct groups, those from North Carolina and those from the Gulf Coast. Similarly, the northern group of populations is subdivided into those from the Appalachians in North Carolina, those from the mid-Atlantic coast, and those from further north. The glacial history of North America means that both the mid-Atlantic the populations farther north must have been derived from one or more southern populations after the height of the last glaciation. Given the phylogenetic relationships recovered here, it seems clear that they were all derived from a population or populations the Appalachian of North Carolina.

That's the promise of NGS for population genetics. What are the challenges? Funny you should ask.

Estimates of nucleotide diversity¹

Beyond the simple challenge of dealing with all of the short DNA fragments that emerge from high-throughput sequencing, there are at least two challenges that don't arise with data obtained in more traditional ways.

1. Most studies involve "shotgun" sequencing of entire genomes. In large diploid genomes, this leads to variable coverage. At sites where coverage is low, there's a good chance that all of the reads will be derived from only one of the two chromosomes present, and a heterozygous individual will be scored as homozygous. "Well," you might say, "let's just throw away all of the sites that don't have at least $8\times$ coverage."² That would

¹This section draws heavily on [65]

²If both chromosomes have an equal probability of being sequenced, the probability that one of them is missed with $8\times$ coverage is $(1/2)^8 = 1/256$.

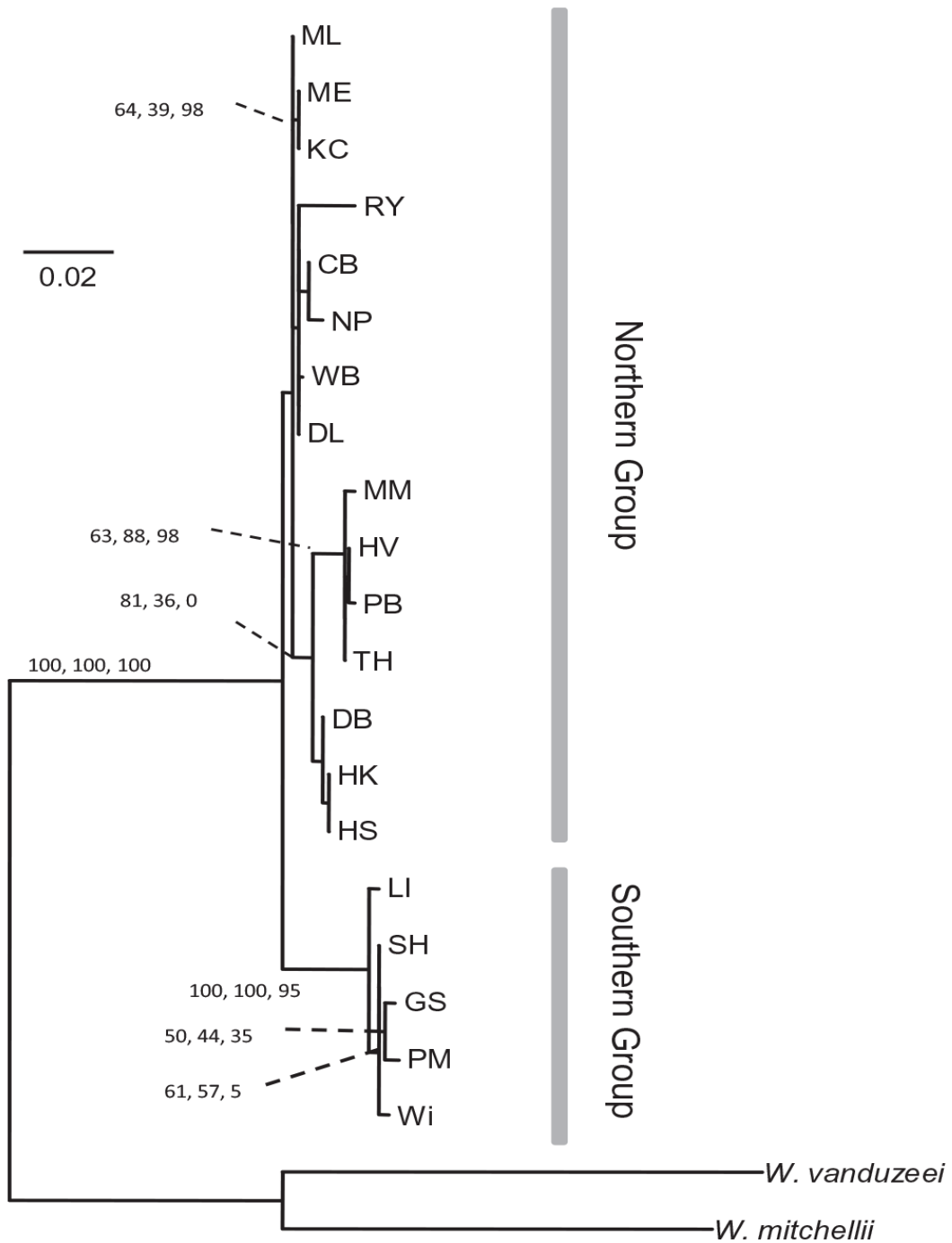


Figure 40.1: Maximum-likelihood phylogenetic tree depicting relationships among populations of *W. smithii* relative to the outgroups *W. vanduzeei* and *W. mitchellii* (from [20]).

work, but you would also be throwing out a lot of potentially valuable information.³ It seems better to develop an approach that lets us use *all* of the data we collect.

2. Sequencing errors are more common with high-throughput methods than with traditional methods, and since so much data is produced, it's not feasible to go back and resequence apparent polymorphisms to see if they reflect sequencing error rather than real differences. Quality scores can be used, but they only reflect the quality of the reads from the sequencing reaction, not errors that might be introduced during sample preparation. Again, we might focus on high-coverage sites and ignore "polymorphisms" associated with single reads, but we'd be throwing away a lot of information.

A better approach than setting arbitrary thresholds and throwing away data is to develop an explicit model of how errors can arise during sequencing and to use that model to interpret the data we've collected. That's precisely the approach that Lynch [65] adopts. Here's how it works assuming that we have a sample from a single, diploid individual:

- Any particular site will have a sequence profile, (n_1, n_2, n_3, n_4) , corresponding to the number of times an A, C, G, or T was observed. $n = n_1 + n_2 + n_3 + n_4$ is the depth of coverage for that site.
- Let ϵ be the probability of a sequencing error at any site, and assume that all errors are equiprobable, e.g., there's no tendency for an A to be miscalled as a C rather than a T when it's miscalled.
- If the site in question were homozygous, the probability of getting our observed sequence profile is:⁴

$$P(n_1, n_2, n_3, n_4 | \text{homozygous}, \epsilon) = \sum_{i=1}^4 \left(\frac{p_i^2}{\sum_{j=1}^4 p_j^2} \right) \binom{n}{n_i} (1 - \epsilon)^{n_i} \epsilon^{n-n_i} \quad .$$

- If the site in question were heterozygous, the probability of getting our observed sequence profile is a bit more complicated. Let k_1 be the number of reads from the first chromosome and k_2 be the number of reads from the second chromosome ($n = k_1 + k_2$). Then

$$\begin{aligned} P(k_1, k_2) &= \binom{n}{k_1} \left(\frac{1}{2}\right)^{k_1} \left(\frac{1}{2}\right)^{k_2} \\ &= \binom{n}{k_1} \left(\frac{1}{2}\right)^n \quad . \end{aligned}$$

³It's valuable information, providing you know how to deal with in properly.

⁴This expression looks a little different from the one in [65], but I'm pretty sure it's equivalent.

Now consider the ordered genotype $x_i x_j$, where x_i refers to the nucleotide on the first chromosome and x_j refers to the nucleotide on the second chromosome. The probability of getting our observed sequence profile from this genotype is:

$$P(n_1, n_2, n_3, n_4 | x_i, x_j, k_1, k_2) = \sum_{l=1}^4 \sum_{m=0}^{k_1} \binom{k_1}{m} (1 - \delta_{il})^m \delta_{il}^{k_1 - m} \binom{k_2}{n_i - m} (1 - \delta_{jl})^{n_i - m} \delta_{jl}^{k_2 - (n_i - m)} ,$$

where

$$\delta_{il} = \begin{cases} 1 - \epsilon & \text{if } i = l \\ \epsilon & \text{if } i \neq l \end{cases} .$$

We can use Bayes' Theorem⁵ to get

$$P(n_1, n_2, n_3, n_4 | x_i, x_j, \epsilon) = P(n_1, n_2, n_3, n_4 | x_i, x_j, k_1, k_2, \epsilon) P(k_1, k_2) ,$$

and with that in hand we can get

$$P(n_1, n_2, n_3, n_4 | \text{heterozygous}, \epsilon) = \sum_{i=1}^4 \sum_{j=1}^4 \left(\frac{x_i x_j}{1 - \sum_{l=1}^4 p_l^2} \right) P(n_1, n_2, n_3, n_4 | x_i, x_j)$$

- Let π be the probability that any site is heterozygous. Then the probability of getting our data is:

$$P(n_1, n_2, n_3, n_4 | \pi, \epsilon) = \pi P(n_1, n_2, n_3, n_4 | \text{heterozygous}, \epsilon) + (1 - \pi) P(n_1, n_2, n_3, n_4 | \text{homozygous}, \epsilon) .$$

- What we've just calculated is the probability of the configuration we observed at a particular site. The probability of our data is just the product of this probability across all of the sites in our sample:

$$P(\text{data} | \pi, \epsilon) = \prod_{s=1}^S P(n_1^{(s)}, n_2^{(s)}, n_3^{(s)}, n_4^{(s)} | \pi, \epsilon) ,$$

where the superscript (s) is used to index each site in the data.

- What we now have is the likelihood of the data in terms of ϵ , which isn't very interesting since it's just the average sequencing error rate in our sample, and π , which is interesting, because it's the genome-wide nucleotide diversity. Now we "simply" maximize that likelihood, and we have maximum-likelihood estimates of both parameters. Alternatively, we could supply priors for ϵ and π and use MCMC to get Bayesian estimates of ϵ and π .

⁵Ask me for details if you're interested.

Taxon	$4N_e\mu$	$4N_e\mu$ (low coverage)	ϵ
<i>Cionia intestinalis</i>	0.0111	0.012	0.00113
<i>Daphnia pulex</i>	0.0011	0.0012	0.00121

Table 40.1: Estimates of nucleotide diversity and sequencing error rate in *Cionia intestinalis* and *Daphnia pulex* (results from [36]).

Notice that this genome-wide estimate of nucleotide diversity is obtained from a sample derived from a single diploid individually. Lynch [65] develops similar methods for estimating gametic disequilibrium as a function of genetic distance for a sample from a single diploid individual. He also extends that method to samples from a pair of individuals, and he describes how to estimate mutation rates by comparing sequences derived from individuals in mutation accumulation lines with consensus sequences.⁶

Haubold et al. [36] describe a program implementing these methods. Recall that under the infinite sites model of mutation $\pi = 4N_e\mu$. They analyzed data sets from the sea squirt *Cionia intestinalis* and the water flea *Daphnia pulex* (Table 40.1). Notice that the sequencing error rate in *D. pulex* is indistinguishable from the nucleotide diversity.

Next-generation AMOVA⁷

What we’ve discussed so far gets us estimates of some population parameters ($4N_e\mu$, $4N_e r$), but they’re derived from the sequences in a single diploid individual. That’s not much of a population sample, and it certainly doesn’t tell us anything about how different populations are from one another. Gompert and Buerkle [30] describe an approach to estimate statistics very similar to Φ_{ST} from AMOVA. Since they take a Bayesian approach to developing their estimates, they refer to approach as BAMOVA, Bayesian models for analysis of molecular variance. They propose several related models.

- **NGS-individual model:** This model assumes that sequencing errors are negligible.⁸ Under this model, the only trick is that we may or may not pick up both sequences

⁶Mutation accumulation lines are lines propagated through several (sometimes up to hundreds) of generations in which population sizes are repeatedly reduced to one or a few individuals, allowing drift to dominate the dynamics and alleles to “accumulate” with little regard to their fitness effects.

⁷This section depends heavily on [30]

⁸Or that they’ve already been corrected. We don’t care *how* they might have been corrected. We care only that we can assume that the reads we get from a sequencing run faithfully reflect the sequences present on each of the chromosomes.

from a heterozygote. The probability of not seeing both sequences in a heterozygote is related to the depth of coverage.

- **NGS-population model:** In some NGS experiments, investigators pool all of the samples from a population into a single sample. Again, Gompert and Buerkle assume that sequencing errors are negligible. Here we assume that the number of reads for one of two alleles at a particular SNP site in a sample is related to the underlying allele frequency at that site. Roughly speaking, the likelihood of the data at that site is then

$$P(x_i|p_i, n_i, k_i) = \binom{n_i}{k_i} p_i^{k_i} (1 - p_i)^{n_i - k_i} \quad ,$$

where p_i is the allele frequency at this site, n_i is the sample size, and k_i is the count of one of the alleles in the sample. The likelihood of the data is just the product across the site-specific likelihoods.⁹

Then, as we did way back when we used a Bayesian approach to estimate F_{ST} [43], we put a prior on the p_i and the parameters of this prior are defined in terms of Φ_{ST} (among other things).¹⁰ They also propose a method for detecting SNP loci¹¹ that have unusually large or small values of Φ_{ST} .

BAMOVA example

Gompert and Buerkle [30] used data derived from two different human population data sets:

- 316 fully sequenced genes in an African population and a population with European ancestry. With these data, they didn't have to worry about the sequencing errors that their model neglects and they could simulate pooled samples allowing them to compare estimates derived from pooled versus individual-level data.
- 12,649 haplotype regions and 11,866 genes derived from 597 individuals across 33 widely distributed human populations.

In analysis of the first data set, they estimated $\Phi_{ST} = 0.08$. Three loci were identified as having unusually high values of Φ_{ST} .

⁹The actual model they use is a bit more complicated than this, but the principles are the same.

¹⁰Again, the actual model is a bit more complicated than what I'm describing here, but the principle is the same.

¹¹Or sets of SNP loci that are parts of a single contig.

- **HSD11B2**: $\Phi_{ST} = 0.32(0.16, 0.48)$. Variants at this locus are associated with an inherited form of high blood pressure and renal disease. A microsatellite in an intron of this locus is weakly associated with type 1 diabetes.
- **FOXA2**: $\Phi_{ST} = 0.32(0.12, 0.51)$. This gene is involved in regulation of insulin sensitivity.
- **POLG2**: $\Phi_{ST} = 0.33(0.18, 0.48)$. This locus was identified as a target of selection in another study.

In analysis of the 33-population data set, they found similar values of Φ_{ST} on each chromosome, ranging from 0.083 (0.075, 0.091) on chromosome 22 to 0.11 (0.10, 0.12) on chromosome 16. Φ_{ST} for the X chromosome was higher: 0.14 (0.13, 0.15). They detected 569 outlier loci, 518 were high outliers and 51 were low outliers. Several of the loci they detected as outlier had been previously identified as targets of selection. The loci they identified as candidates for balancing selection have not been suggested before as targets of such selection.

Literature cited

- [1] S J Arnold. Quantitative genetics and selection in natural populations: microevolution of vertebral numbers in the garter snake *Thamnophis elegans*. In B S Weir, E J Eisen, M M Goodman, and G Namkoong, editors, *Proceedings of the Second International Conference on Quantitative Genetics*, pages 619–636. Sinauer Associates, Sunderland, MA, 1988.
- [2] J C Avise, J Arnold, R M Ball, E Bermingham, T Lamb, J E Neigel, C A Reeb, and N C Saunders. Intraspecific phylogeography: the mitochondrial DNA bridge between population genetics and systematics. *Annual Review of Ecology & Systematics*, 18:489–522, 1987.
- [3] Nathan A Baird, Paul D Etter, Tressa S Atwood, Mark C Currey, Anthony L Shiver, Zachary A Lewis, Eric U Selker, William A Cresko, and Eric A Johnson. Rapid SNP discovery and genetic mapping using sequenced RAD markers. *PLoS ONE*, 3(10):e3376, 2008.
- [4] Mark A Beaumont, Wenyang Zhang, and David J Balding. Approximate Bayesian computation in population genetics, 2002.
- [5] P Beerli. Comparison of Bayesian and maximum-likelihood estimation of population genetic parameters. *Bioinformatics*, 22:341–345, 2006.
- [6] Peter Beerli. Effect of unsampled populations on the estimation of population sizes and migration rates between sampled populations, 2004.
- [7] Peter Beerli and Joseph Felsenstein. Maximum-likelihood estimation of migration rates and effective population numbers in two populations using a coalescent approach, 1999.
- [8] Peter Beerli and Joseph Felsenstein. Maximum likelihood estimation of a migration matrix and effective population sizes in n subpopulations by using a coalescent approach, 2001.

- [9] D D Brown, P C Wensink, and E Jordan. Comparison of the ribosomal DNA's of *Xenopus laevis* and *Xenopus mulleri*: the evolution of tandem genes. *J. Mol. Biol.*, 63:57–73, 1972.
- [10] R L Cann, M Stoneking, and A C Wilson. Mitochondrial DNA and human evolution. *Nature*, 325:31–36, 1987.
- [11] R Ceppellini, M Siniscalco, and C A B Smith. The estimation of gene frequencies in a random-mating population. *Annals of Human Genetics*, 20:97–115, 1955.
- [12] F B Christiansen. Studies on selection components in natural populations using population samples of mother-offspring combinations. *Hereditas*, 92:199–203, 1980.
- [13] F B Christiansen and O Frydenberg. Selection component analysis of natural polymorphisms using population samples including mother-offspring combinations. *Theoretical Population Biology*, 4:425–445, 1973.
- [14] T E Cleghorn. MNSs gene frequencies in English blood donors. *Nature*, 187:701, 1960.
- [15] The International HapMap Consortium. A second generation human haplotype map of over 3.1 million SNPs. *Nature*, 449(7164):851–861, 2007.
- [16] J F Crow and M Kimura. *An Introduction to Population Genetics Theory*. Burgess Publishing Company, Minneapolis, Minn., 1970.
- [17] A P Dempster, N M Laird, and D B Rubin. Maximum likelihood from incomplete data. *Journal of the Royal Statistical Society Series B*, 39:1–38, 1977.
- [18] T Dobzhansky and C Epling. *Contributions to the genetics, taxonomy, and ecology of Drosophila pseudoobscura and its relatives. Publication 554*. Carnegie Institution of Washington, Washington, DC, 1944.
- [19] Th. Dobzhansky. Genetics of natural populations. XIV. A response of certain gene arrangements in the third chromosome of *Drosophila pseudoobscura* to natural selection. *Genetics*, 32:142–160, 1947.
- [20] Kevin Emerson, Clayton Merz, Julian Catchen, Paul A Hohenlohe, William Cresko, William Bradshaw, and Christina Holzapfel. Resolving postglacial phylogeography using high-throughput sequencing. *Proceedings of the National Academy of Sciences of the United States of America*, 107(37):16196–16200, 2010.

- [21] Arnaud Estoup, Mark A Beaumont, Florent Sennedot, Craig Moritz, and Jean-Marie Cornuet. Genetic analysis of complex demographic scenarios: spatially expanding populations of the cane toad, *Bufo marinus*, 2004.
- [22] J C et al. Venter. The sequence of the human genome. *Science*, 291:1304–1351, 2001.
- [23] L Excoffier, P E Smouse, and J M Quattro. Analysis of molecular variance inferred from metric distances among DNA haplotypes: application to human mitochondrial DNA restriction data. *Genetics*, 131(2):479–491, 1992.
- [24] J C Fay and C.-I. Wu. Hitchhiking under positive Darwinian selection. *Genetics*, 155:1405–1413, 2000.
- [25] R A Fisher. The correlation between relatives on the supposition of Mendelian inheritance. *Transactions of the Royal Society of Edingurgh*, 52:399–433, 1918.
- [26] T A K Freitas, S Hou, E M Dioum, J A Saito, J Newhouse, G Gonzalez, M.-A. Gilles-Gonzalez, and M Alam. Ancestral hemoglobins in Archaea. *Proceedings of the National Academy of Sciences USA*, 101:6675–6680, 2004.
- [27] R Fu, A E Gelfand, and K E Holsinger. Exact moment calculations for genetic models with migration, mutation, and drift. *Theoretical Population Biology*, 63:231–243, 2003.
- [28] Y X Fu. Statistical properties of segregating sites. *Theoretical Population Biology*, 48:172–197, 1995.
- [29] Y.-X. Fu. Statistical tests of neutrality of mutations against population growth, hitchhiking, and background selection. *Genetics*, 147:915–925, 1997.
- [30] Zachariah Gompert and C Alex Buerkle. A hierarchical Bayesian model for next-generation population genomics. *Genetics*, 187(3):903–917, March 2011.
- [31] M Goodman. Immunocytochemistry of the primates and primate evolution. *Annals of the New York Academy of Sciences*, 102:219–234, 1962.
- [32] Feng Guo, Dipak K Dey, and Kent E Holsinger. A Bayesian hierarchical model for analysis of SNP diversity in multilocus, multipopulation samples. *Journal of the American Statistical Association*, 104(485):142–154, March 2009.
- [33] Thomas M Hammond, David G Rehard, Hua Xiao, and Patrick K T Shiu. Molecular dissection of *Neurospora* Spore killer meiotic drive elements. *Proceedings of the National Academy of Sciences of the United States of America*, 109(30):12093–12098, July 2012.

- [34] Henry C Harpending, Mark A Batzer, Michael Gurven, Lynn B Jorde, Alan R Rogers, and Stephen T Sherry. Genetic traces of ancient demography. *Proceedings of the National Academy of Sciences of the United States of America*, 95(4):1961–1967, 1998.
- [35] H Harris. Enzyme polymorphisms in man. *Proceedings of the Royal Society of London, Series B*, 164:298–310, 1966.
- [36] Bernhard Haubold, Peter Pfaffelhuber, and MICHAEL LYNCH. mlRho - a program for estimating the population mutation and recombination rates from shotgun-sequenced diploid genomes. *Molecular Ecology*, 19:277–284, March 2010.
- [37] P W Hedrick. *Genetics of Populations*. Jones and Bartlett Publishers, Sudbury, MA, 2nd ed. edition, 2000.
- [38] Jody Hey and Rasmus Nielsen. Multilocus methods for estimating population sizes, migration rates and divergence time, with applications to the divergence of *Drosophila pseudoobscura* and *D. persimilis*, 2004.
- [39] Jody Hey and Rasmus Nielsen. Integration within the Felsenstein equation for improved Markov chain Monte Carlo methods in population genetics. *Proceedings of the National Academy of Sciences*, 104(8):2785–2790, 2007.
- [40] W G Hill and A Robertson. Linkage disequilibrium in finite populations. *Theoretical and Applied Genetics*, 38:226–231, 1968.
- [41] K E Holsinger. The population genetics of mating system evolution in homosporous plants. *American Fern Journal*, pages 153–160, 1990.
- [42] K E Holsinger and R J Mason-Gamer. Hierarchical analysis of nucleotide diversity in geographically structured populations. *Genetics*, 142(2):629–639, 1996.
- [43] K E Holsinger and L E Wallace. Bayesian approaches for the analysis of population structure: an example from *Platanthera leucophaea* (Orchidaceae). *Molecular Ecology*, 13:887–894, 2004.
- [44] Kent E. Holsinger and Bruce S. Weir. Genetics in geographically structured populations: defining, estimating, and interpreting F_{ST} . *Nature Reviews Genetics*, 10:639–650, 2009.
- [45] J L Hubby and R C Lewontin. A molecular approach to the study of genic heterozygosity in natural populations. I. The number of alleles at different loci in *Drosophila pseudoobscura*. *Genetics*, 54:577–594, 1966.

- [46] A L Hughes and M Nei. Pattern of nucleotide substitution at major histocompatibility complex class I loci reveals overdominant selection. *Nature*, 335:167–170, 1988.
- [47] A L Hughes, T Ota, and M Nei. Positive Darwinian selection promotes charge profile diversity in the antigen-binding cleft of class I major histocompatibility complex molecules. *Molecular Biology & Evolution*, 7(6):515–524, 1990.
- [48] S H James, A P Wylie, M S Johnson, S A Carstairs, and G A Simpson. Complex hybridity in *Isotoma petraea* V. Allozyme variation and the pursuit of hybridity. *Heredity*, 51(3):653–663, 1983.
- [49] R C Jansen, H Geerlings, A J VanOeveren, and R C VanSchaik. A comment on codominant scoring of AFLP markers. *Genetics*, 158(2):925–926, 2001.
- [50] M Kimura. Evolutionary rate at the molecular level. *Nature*, 217:624–626, 1968.
- [51] J L King and T L Jukes. Non-Darwinian evolution. *Science*, 164:788–798, 1969.
- [52] J F C Kingman. On the genealogy of large populations. *Journal of Applied Probability*, 19A:27–43, 1982.
- [53] J F C Kingman. The coalescent. *Stochastic Processes and their Applications*, 13:235–248, 1982.
- [54] L Knowles. Did the Pleistocene glaciations promote divergence? Tests of explicit refugial models in montane grasshoppers, 2001.
- [55] L Knowles and Wayne P Maddison. Statistical phylogeography, 2002.
- [56] M Kreitman. Nucleotide polymorphism at the alcohol dehydrogenase locus of *Drosophila melanogaster*. *Nature*, 304:412–417, 1983.
- [57] M Kreitman and M Aguadé. Excess polymorphism at the alcohol dehydrogenase locus in *Drosophila melanogaster*. *Genetics*, 114:93–110, 1986.
- [58] M Kreitman and R R Hudson. Inferring the evolutionary history of the *Adh* and *Adh-dup* loci in *Drosophila melanogaster* from patterns of polymorphism and divergence. *Genetics*, 127:565–582, 1991.
- [59] R Lande and S J Arnold. The measurement of selection on correlated characters. *Evolution*, 37:1210–1226, 1983.

- [60] R C Lewontin and J L Hubby. A molecular approach to the study of genic heterozygosity in natural populations. II. Amount of variation and degree of heterozygosity in natural populations of *Drosophila pseudoobscura*. *Genetics*, 54:595–609, 1966.
- [61] C C Li. *First Course in Population Genetics*. Boxwood Press, Pacific Grove, CA, 1976.
- [62] W.-H. Li. *Molecular Evolution*. Sinauer Associates, Sunderland, MA, 1997.
- [63] J D Lubell, M H Brand, J M Lehrer, and K E Holsinger. Detecting the influence of ornamental *Berberis thunbergii* var. *atropurpurea* in invasive populations of *Berberis thunbergii* (Berberidaceae) using AFLP. *American Journal of Botany*, 95(6):700–705, 2008.
- [64] M Lynch and B Walsh. *Genetics and Analysis of Quantitative Traits*. Sinauer Associates, Sunderland, MA, 1998.
- [65] MICHAEL LYNCH. Estimation of nucleotide diversity, disequilibrium coefficients, and mutation rates from high-coverage genome-sequencing projects. *Molecular biology and evolution*, 25(11):2409–2419, November 2008.
- [66] T Nagylaki. Evolution of multigene families under interchromosomal gene conversion. *Proceedings of the National Academy of Sciences USA*, 81:3796–3800, 1984.
- [67] M Nei and R K Chesser. Estimation of fixation indices and gene diversities. *Annals of Human Genetics*, 47:253–259, 1983.
- [68] Anton Nekrutenko and James Taylor. Next-generation sequencing data interpretation: enhancing reproducibility and accessibility. *Nature Publishing Group*, 13(9):667–672, September 2012.
- [69] Rasmus Nielsen and J Wakeley. Distinguishing migration from isolation: a Markov chain Monte Carlo approach, 2001.
- [70] H Nilsson-Ehle. Kreuzungsuntersuchungen an Hafer und Weizen. *Lunds Universitets Arsskrift*, n.s. 2, vo(no. 2), 1909.
- [71] Caroline Obert, Jack Sublett, Deepak Kaushal, Ernesto Hinojosa, Theresa Barton, Elaine I Tuomanen, and Carlos J Orihuela. Identification of a Candidate *Streptococcus pneumoniae* Core Genome and Regions of Diversity Correlated with Invasive Pneumococcal Disease. *Infection and Immunity*, 74(8):4766–4777, 2006.

- [72] T Ohta. Some models of gene conversion for treating the evolution of multigene families. *Genetics*, 106:517–528, 1984.
- [73] T Ohta. Gene families: multigene families and superfamilies. In *Encyclopedia of the Human Genome*. Macmillan Publishers Ltd., London, 2003.
- [74] Tomoko Ohta and Motoo Kimura. Linkage disequilibrium at steady state determined by random genetic drift and recurrent mutation. *Genetics*, 63:229–238, 1969.
- [75] Guillermo Orti, Michael A Bell, Thomas E Reimchen, and Axel Meyer. Global survey of mitochondrial DNA sequences in the threespine stickleback: evidence for recent migrations. *Evolution*, 48(3):608–622, 1994.
- [76] P Pamilo and M Nei. Relationships between gene trees and species trees. *Molecular Biology and Evolution*, 5(5):568–583, 1988.
- [77] Jonathan Pritchard, Matthew Stephens, and Peter Donnelly. Inference of Population Structure Using Multilocus Genotype Data, 2000.
- [78] V M Sarich and A C Wilson. Immunological time scale for hominid evolution. *Science*, 158:1200–1203, 1967.
- [79] Sven J Saupe. A fungal gene reinforces Mendel’s laws by counteracting genetic cheating. *Proceedings of the National Academy of Sciences of the United States of America*, 109(30):11900–11901, July 2012.
- [80] Stefan Schneider and Laurent Excoffier. Estimation of Past Demographic Parameters From the Distribution of Pairwise Differences When the Mutation Rates Vary Among Sites: Application to Human Mitochondrial DNA. *Genetics*, 152(3):1079–1089, 1999.
- [81] Jonathan Sebat, B Lakshmi, Jennifer Troge, Joan Alexander, Janet Young, Par Lundin, Susanne Maner, Hillary Massa, Megan Walker, Maoyen Chi, Nicholas Navin, Robert Lucito, John Healy, James Hicks, Kenny Ye, Andrew Reiner, T Conrad Gilliam, Barbara Trask, Nick Patterson, Anders Zetterberg, and Michael Wigler. Large-Scale Copy Number Polymorphism in the Human Genome. *Science*, 305(5683):525–528, 2004.
- [82] Montgomery Slatkin. Inbreeding coefficients and coalescence times. *Genetical Research*, 58:167–175, 1991.
- [83] Montgomery Slatkin. Inbreeding coefficients and coalescence times. *Genetical Research*, 58:167–175, 1991.

- [84] Douglas E Soltis, Ashley B Morris, Jason S McLachlan, Paul S Manos, and Pamela S Soltis. Comparative phylogeography of unglaciated eastern North America. *Molecular Ecology*, 15(14):4261–4293, 2006.
- [85] S Song, D K Dey, and K E Holsinger. Differentiation among populations with migration, mutation, and drift: implications for genetic inference. *Evolution*, 60:1–12, 2006.
- [86] Matthew Stephens and D Balding. Bayesian statistical methods for genetic association studies. *Nature Reviews Genetics*, 10:681–690, 2009.
- [87] F Tajima. Statistical method for testing the neutral mutation hypothesis by DNA polymorphism. *Genetics*, 123:585–595, 1989.
- [88] A R Templeton, E Boerwinkle, and C F Sing. A cladistic analysis of phenotypic associations with haplotypes inferred from restriction endonuclease mapping. I. Basic theory and an analysis of alcohol dehydrogenase activity in *Drosophila*. *Genetics*, 117:343–351, 1987.
- [89] Alan R Templeton. Statistical phylogeography: methods of evaluating and minimizing inference errors, 2004.
- [90] Alan R Templeton, Keith A Crandall, and Charles F Sing. A cladistic analysis of phenotypic associations with haplotypes inferred from restriction endonuclease mapping and DNA sequence data. III. Cladogram estimation. *Genetics*, 132(2):619–633, 1992.
- [91] Alan R Templeton, Eric Routman, and Christopher A Phillips. Separating population structure from population history: a cladistic analysis of the geographical distribution of mitochondrial DNA haplotypes in the tiger salamander, *Ambystoma tigrinum*. *Genetics*, 140(2):767–782, 1995.
- [92] Marie Touchon, Claire Hoede, Olivier Tenailon, Valérie Barbe, Simon Baeriswyl, Philippe Bidet, Edouard Bingen, Stéphane Bonacorsi, Christiane Bouchier, Odile Bouvet, Alexandra Calteau, Hélène Chiapello, Olivier Clermont, Stéphane Cruveiller, Antoine Danchin, Médéric Diard, Carole Dossat, Meriem El Karoui, Eric Frapy, Louis Garry, Jean Marc Ghigo, Anne Marie Gilles, James Johnson, Chantal Le Bougénéec, Mathilde Lescat, Sophie Mangenot, Vanessa Martinez-Jéhanne, Ivan Matic, Xavier Nassif, Sophie Oztas, Marie Agnès Petit, Christophe Pichon, Zoé Rouy, Claude Saint Ruf, Dominique Schneider, Jérôme Turret, Benoit Vacherie, David Vallenet, Claudine Médigue, Eduardo P C Rocha, and Erick Denamur. Organised Genome Dynamics in

- the *Escherichia coli* Species Results in Highly Diverse Adaptive Paths. *PLoS Genet*, 5(1):e1000344, 2009.
- [93] Peter A Underhill, Peidong Shen, Alice A Lin, Li Jin, Giuseppe Passarino, Wei H Yang, Erin Kauffman, Batsheva Bonne-Tamir, Jaume Bertranpetit, Paolo Francalacci, Muntaser Ibrahim, Trefor Jenkins, Judith R Kidd, S Qasim Mehdi, Mark T Seielstad, R Spencer Wells, Alberto Piazza, Ronald W Davis, Marcus W Feldman, L Luca Cavalli-Sforza, and Peter J Oefner. Y chromosome sequence variation and the history of human populations. *Nature Genetics*, 26(3):358–361, 2000.
- [94] S Wahlund. Zusammensetzung von Population und Korrelationserscheinung vom Standpunkt der Vererbungslehre aus betrachtet. *Hereditas*, 11:65–106, 1928.
- [95] C Wedekind, T Seebeck, F Bettens, and A J Paepke. MHC-dependent mate preferences in humans. *Proceedings of the Royal Society of London, Series B*, 260:245–249, 1995.
- [96] B S Weir. *Genetic Data Analysis II*. Sinauer Associates, Sunderland, MA, 1996.
- [97] B S Weir and C C Cockerham. Estimating F -statistics for the analysis of population structure. *Evolution*, 38:1358–1370, 1984.
- [98] B S Weir and W G Hill. Estimating F -statistics. *Annual Review of Genetics*, 36:721–750, 2002.
- [99] Eva-Maria Willing, Christine Dreyer, and Cock van Oosterhout. Estimates of Genetic Differentiation Measured by F_{ST} Do Not Necessarily Require Large Sample Sizes When Using Many SNP Markers. *PLoS ONE*, 7(8):e42649, August 2012.
- [100] A C Wilson and V M Sarich. A molecular time scale for human evolution. *Proceedings of the National Academy of Sciences U.S.A.*, 63:1088–1093, 1969.
- [101] R Woltereck. Weiterer experimentelle Untersuchungen über Artveränderung, Speziell über das Wesen Quantitativer Arunterschiede bei Daphiden. *Versuchung Deutsch Zoologische Gesellschaft*, 19:110–172, 1909.
- [102] Sewall Wright. *Evolution and the Genetics of Populations*, volume 2. University of Chicago Press, Chicago, IL, 1969.
- [103] Sewall Wright. *Evolution and the Genetics of Populations.*, volume 4. University of Chicago Press, Chicago, IL, 1978.

- [104] K Zeng, Y.-X. Fu, S Shi, and C.-I. Wu. Statistical tests for detecting positive selection by utilizing high-frequency variants. *Genetics*, 174:1431–1439, 2006.
- [105] E Zuckerkandl and L Pauling. Evolutionary divergence and convergence in proteins. In V Bryson and H J Vogel, editors, *Evolving Genes and Proteins*, pages 97–166. Academic Press, New York, NY, 1965.

Index

- additive effect, 139, 142
 - Hardy-Weinberg assumption, 139
- additive genetic variance, 143
- additive genotypic value, 138
- Adh, 228
 - balancing selection, 229
 - purifying selection, 228
- alcohol dehydrogenase, 228
- allele fixation, 79, 81
- allele frequency distribution, 50
- allele genealogy, 127
- Ambystoma*
 - tigrinum*, 278, 279
- Amia*
 - calva*, 272
- among-site rate variation, 212
 - shape parameter, 212
- AMOVA, 264
 - example, 265
- ancestral polymorphism, 282
- Approximate Bayesian Computation, 295
 - limitations, 299
 - regression, 297
- association mapping
 - BAMD priors, 200
 - linear mixed model, 199
 - relatedness, 200
- assortative mating, 19
- BAMOVA, 307
 - example, 308
- Bayesian inference, 14
- Beta distribution, 50
- biometricians, 135
- blending inheritance, 135
- breeders equation, 166, 169, 170
- Bufo*
 - marinus*, 293, 297
- Cionia*
 - intestinalis*, 307
- clade distance, 276
- coalescent, 127
 - balancing selection, 229
 - diverging populations, 289
 - estimating migration, 289
 - estimating migration rates, 288
 - F*-statistics, 132
 - migration, 285
 - mitochondrial Eve, 131
 - multiple alleles, 130
 - time to coalescence, 130
 - two alleles, 128
- coalescent events, 128
- components of selection, 72
- components of variance
 - causal, 152, 154
 - observational, 152, 154
- covariance, 149
 - half-siblings, 149
- cumulative selection gradient, 172
 - caveats, 174

Daphnia
 pulex, 307
 Deviance Information Criterion, 31, 51
 directional selection, 79
 disassortative mating, 73
 disruptive selection, 80
 diversity-divergence, 229
 dominance genetic variance, 143
Drosophila
 melanogaster, 228, 274
 pseudoobscura, 74

E-matrix, 170
 effective neutrality, 217
 effectively neutral, 126
 EM algorithm, 11
 environmental variance, 137
 equilibrium, 22, 81
 monomorphic, 81
 polymorphic, 81
 unstable, 81
 estimate, 49
 evolutionary pattern, 205
 evolutionary process, 205

F-statistics, 36, 39
 G_{st} , 44
 coalescent, 132
 notation, 46
 outliers, 231
 Weir and Cockerham, 44
 Fay and Wu's H , 242
 fecundity selection, 72
 fertility selection, 72
 First law of population genetics, 6
 Fisher's Fundamental Theorem of Natural Selection, 79
 fitness
 regression on phenotype, 164, 165
 Fu's F_S , 242
 full-sib analysis, 147, 152
 fundamental theorem of natural selection, 167

G-matrix, 170
 gamete competition, 72
 gametic disequilibrium, 194
 drift, 197
 gene tree, 281
 genetic code, 222
 redundancy, 223
 genetic composition of populations, 7
 genetic drift, 103
 allele frequency variance, 107
 binomial distribution, 106
 effective population size, 110
 effective population size, limitations, 111
 effective population size, separate sexes, 112
 effective population size, variable population size, 114
 effective population size, variation in offspring number, 115
 effectively neutral, 126
 fixation of deleterious alleles, 125
 fixation probability, 124
 fixation time, 108
 ideal population, 109
 inbreeding analogy, 108
 inbreeding effective size, 110
 loss of beneficial alleles, 123
 migration, 120
 mutation, 117
 mutation, recurrent, 119
 mutation, stationary distribution, 119
 mutation, stationary distribution, 118

- population size, 120
- properties, 105, 106, 108
- properties with selection, 125
- uncertainty in allele frequencies, 104
- variance effective size, 110
- genetic variance, 137
 - additive, 144
 - components, 142
 - dominance, 143
- genotypic value, 138, 162
 - additive, 138
 - fitness, 162
- geographic structure, 33
- Hardy-Weinberg assumptions, 8
- Hardy-Weinberg principle, 10
- Hardy-Weinberg proportions
 - multiple alleles, 87
- heritability, 152, 154
 - broad sense, 138
 - narrow sense, 138
- identity by type, 24
- identity by descent, 24
- immunological distance, 206
- inbreeding, 19
 - consequences, 21
 - partial self-fertilization, 21
 - self-fertilization, 20
 - types, 19
- inbreeding coefficient, 23
 - equilibrium, 24
 - population, 24
- individual assignment, 57
 - application, 57
- Jukes-Cantor distance, 210
 - assumptions, 211
- linkage disequilibrium, 194
- marginal fitness, 78
- mating table, 8
 - self-fertilization, 20
- maximum-likelihood estimates, 13
- MCMC sampling, 15
- mean fitness, 75
- Melanopus*, 283
- Mendelians, 135
- MHC
 - conservative and non-conservative substitutions, 237
 - synonymous and non-synonymous substitutions, 236
- MHC polymorphism, 235
- migration
 - estimating, 288
- migration rate
 - backward, 121
 - forward, 121
- molecular clock, 206, 216
 - derivation, 217
- molecular variation
 - markers, 208
 - physical basis, 206
- monomorphic, 79
- mother-offspring pairs, 148
- mutation
 - infinite alleles model, 117, 219
 - infinite sites model, 239
- mutation rate, 215
- natural selection, 72
 - components of selection, 72
 - disassortative mating, 73
 - fertility selection, 72, 90
 - fertility selection, fertility matrix, 90

- fertility selection, properties, 91
- fertility selection, protected polymorphism, 91
- gamete competition, 72
- multiple alleles, marginal viability, 88
- patterns, 78
- segregation distortion, 72
- sexual selection, 73, 92
- viability selection, 73
- nature vs. nurture, 137
- nested clade analysis, 272
 - clade distance, 276
 - constructing nested clades, 274
 - nested clade distance, 277
 - statistical parsimony, 273
- neutral alleles, 216
- neutral theory
 - effective neutrality, 217
 - modifications, 225
- next-generation sequencing, 301
 - estimating F_{ST} , 301
 - estimating nucleotide diversity, 302
 - partitioning diversity, 307
 - phylogeography, 302
- non-synonymous substitutions, 223
- norm of reaction, 136
- nucleotide diversity, 239, 263
 - partitioning, 264
- nucleotide substitutions
 - selection against, 237
 - selection for, 237
- P -matrix, 170
- parameter, 49
- parent-offspring regression, 147, 151
- phenotypic plasticity, 136
- phenotypic variance
 - partitioning, 137
- phenylketonuria, 137
- Φ_{st} , 264
- phylogeography, 271
- polygenic inheritance, 136
- population tree, 281
- QTL mapping
 - caveats, 182
 - inbred lines, 181
 - outline, 177
- quantitative traits, 135
- quantitative trait locus, 177
- R/qtl, 185
 - data format, 186
 - estimating QTL effects, 190
 - identifying QTLs, 189
 - permutation test, 189
 - QTL analysis, 187
 - visualizing QTL effects, 190
- RAD sequencing, 301
- recombination frequency, 196
- relative fitness, 77
- response to selection, 138, 162, 166, 169
- sampling
 - genetic, 49
 - statistical, 49
- sampling error, 43
- segregating sites, 239
- segregation distortion, 72
- selection
 - directional selection, 79
 - disruptive, 80
 - multivariate example, 171
 - stabilizing, 82
- selection coefficient, 80
- selection differential, 166, 169
- selection equation, 76

self-fertilization, 20
 partial, 21
 sexual selection, 19, 73
 sledgehammer principle, 221, 223, 228, 236
 stabilizing selection, 82
 statistical expectation, 41
 statistical parsimony, 273
 example, 274
 haplotype network, 273
 statistical phylogeography
 example, 284
 substitution rate, 215
 substitution rates, 223
 synonymous substitutions, 223

 Tajima's D , 239
 interpretation, 241
 TCS parsimony, 273
 testing Hardy-Weinberg, 27
 Bayesian approach, 29
 goodness of fit, 28
 two-locus genetics
 decay of disequilibrium, 197
 drift, 197
 gamet frequencies, 194
 gametic disequilibrium, 194
 Hardy-Weinberg, 197
 recombination, 196
 transmission, 195

 unbiased estimate, 42

 viability
 absolute, 77
 estimating absolute, 83
 estimating relative, 84
 relative, 77
 viability selection, 73
 genetics, 73

 virility selection, 72

 Wahlund effect, 33
 properties, 35
 theory, 34
 WinBUGS, 15
 Wyeomyia
 smithii, 302

 Zeng et al.'s E , 243
 zero force laws, 5
 Zoarces viviparus, 3