# Chapter 5

## **Population Structure**

### Introduction

(1) Statistical and genetic sampling :

Population genetic theory depends on the concept of replicated populations that are maintained under the same conditions, but will differ because of genetic sampling. The derived variance for statistics of interest should consider both types of variation.

- (2) Between-population variation cannot be estimated with a sample from a single population. Different loci sometimes may be regarded as playing the role of separate populations.
- (3) The distinction between statistical and genetic sampling can also be phrased in terms of fixed and random effects.

## **Fixed populations**

Every member of the population has an equal chance of being sampled, and individuals are sampled independently, the genotypic counts are multinomially distributed.

Different populations for the same species are compared simply by comparing frequencies.

**Contingency Tables:** With *v* alleles at a locus, the genotypic counts in each of *r* samples are arranged in a  $v(v+1)/2 \times r$  contingency table and a chi-square statistic with  $[v(v+1)/2-1] \times r$  degrees of freedom is calculated.

#### Numerical Resampling:

Bootstrapping: Two populations can be judged to have different allele frequencies if the estimated frequencies have nonoverlapping confidence intervals.

Chebyshev's inequality:

$$\Pr(\left|\widetilde{p}_{i_{l}} - p_{i_{l}}\right| > k\sigma) \leq \frac{1}{k^{2}}$$
$$\Pr(\left|\widetilde{p}_{i_{l}} - p_{i_{ll}}\right| > K\sigma\sqrt{2}) \leq \frac{1}{K^{2}}$$

#### **Permutation Tests:**

$$\Pr(\{n_{i_{I}}\},\{n_{i_{II}}\}) = \left(\frac{(2n_{I})!}{\prod n_{i_{I}}!}\prod p_{i_{I}}^{n_{i_{I}}}\right) \left(\frac{(2n_{II})!}{\prod n_{i_{II}}!}\prod p_{i_{II}}^{n_{i_{II}}}\right)$$
$$\Pr(\{n_{i_{I}}\},\{n_{i_{II}}\}|\{n_{i_{I}}+n_{i_{II}}\}) = \frac{(2n_{I})!(2n_{II})!}{(2n_{I}+2n_{II})!}\frac{\prod (n_{i_{I}}+n_{i_{II}})!}{\prod n_{i_{I}}!\prod n_{i_{II}}!}$$

	Sample I	Sample II
Allele A	$n_{A_I} = 8$	$n_{A_{II}}=2$
Allele a	$n_{a_{I}} = 2$	$n_{a_{II}} = 8$
	$2n_I = 10$	$2n_{II} = 10$

**Table 5.1**Possible allocations of 10 A and 10 a alleles into two samples of size 10alleles, along with conditional probabilities.

Number of $A$ alleles		
Sample I	Sample II	Probability
0	10	0.0000
1	9	0.0005
2	8	0.0110
3	7	0.0779
4	6	0.2387
5	• 5	0.3438
6	4	0.2387
7	3	0.0779
8	2	0.0110
9	1	0.0005
10	0	0.0000

#### F statistics:

$$F_{st} = \frac{\sum_{i} (\tilde{P}_{i} - \overline{P})^{2}}{\overline{P}(1 - \overline{P})}$$

In the fixed-population framework,  

$$X^2 = (r-1)\overline{n}F_{st}$$

#### Analysis of Variance:

$$X_{ij} = \begin{cases} 1 & \text{if allele is } A \\ 0 & \text{if allele is not } A \end{cases}$$

**Table 5.2** Analysis of variance layout for variable indicating allele A in fixed populations.

Source	d.f.	Sum of Squares	Expected Mean Square <sup>*</sup>
Between	r-1	$\sum_{i} \frac{x_{i.}^2}{n_i} - \frac{x_{}^2}{n}$	$\frac{1}{r-1}\sum_i k_1 p_{Ai}(1-p_{Ai})$
p op allotions		$=\sum_i n_i (\tilde{p}_{Ai} - \tilde{p}_{A.})^2$	$+\frac{1}{r-1}\sum_{i}n_{i}(p_{Ai}-\bar{p}_{A.})^{2}$
Within	$\sum_{i}(n_i-1)$	$\sum_{i=1}^{r} \sum_{j=1}^{n_i} x_{ij}^2 - \sum_i \frac{x_{i.}^2}{n_i}$	$\frac{1}{\sum_{i}(n_{i}-1)}\sum_{i}k_{2}p_{Ai}(1-p_{Ai})$
populations	= n r	$=\sum_{i}n_{i}\tilde{p}_{Ai}(1-\tilde{p}_{Ai})$	

$${}^{*}k_{1} = \left(1 - \frac{n_{i}}{\sum_{i} n_{i}}\right), \ k_{2} = (n_{i} - 1)$$

### **Random Populations**

- The action of evolutionary forces (genetic sampling) causes different alleles in a population to be dependent, or related, and will result in intraspecific differentiation. Even though individuals, or alleles, may be sampled randomly, the process of taking expectations must recognize that they are dependent through their shared ancestry.
- The differentiation is quantified with the F statistics of Wright (1951), or the analogous measures of Cockerham (1963,1973).



Haploid Data

There is only one F statistic:  $\theta$ 

Shared ancestry means that the expected value of a squared sample frequency from a sample of size  $n_i$  is

$$E(\tilde{P}_{Ai}^{2}) = P_{A}^{2} + P_{A}(1 - P_{A})\theta + \frac{1}{n_{i}}P_{A}(1 - P_{A})(1 - \theta)$$

Between population differentiation goes hand in hand with relatedness of alleles within populations. As individuals become more related within populations, the independent populations are expected to become more differentiated.

**Table 5.3** Analysis of variance layout for variable indicating allele A in random populations.

Source	d.f.	Sum of Squares	Expected Mean Square <sup>*</sup>
Between	r - 1	$\sum_{i=1}^{r} \frac{x_{i.}^2}{n_i} - \frac{x_{}^2}{n_i}$	$p_A(1-p_A)[(1-\theta)+n_c\theta]$
1-1		$=\sum_{i=1}^{r}n_{i}(\tilde{p}_{Ai}-\tilde{p}_{A.})^{2}$	$= \sigma_G^2 + n_c \sigma_P^2$
Within populations	$\sum_{i=1}^{r} (n_i - 1)$ $= n r$	$\sum_{i=1}^{r} \sum_{j=1}^{n_i} x_{ij}^2 - \sum_i \frac{x_{i.}^2}{n_i}$ $= \sum_i n_i \tilde{p}_{Ai} (1 - \tilde{p}_{Ai})$	$p_A(1-p_A)(1-\theta)$ $= \sigma_G^2$

$${}^{*}n_{c} = \frac{1}{r-1} \left( \sum_{i=1}^{r} n_{i} - \frac{\sum_{i} n_{i}^{2}}{\sum_{i} n_{i}} \right)$$

$$p_{A}(1-p_{A})(1-\theta) \stackrel{\text{a}}{=} \frac{1}{\sum_{i} n_{i} - 1} \sum_{i} n_{i} \tilde{p}_{Ai}(1-\tilde{p}_{Ai}) = \text{MSG}$$

$$p_{A}(1-p_{A})(1-\theta) + n_{c}p_{A}(1-p_{A})\theta \stackrel{\text{a}}{=} \frac{1}{r-1} \sum_{i} n_{i}(\tilde{p}_{Ai} - \tilde{p}_{A}) = \text{MSP}$$

$$\sigma_{G}^{2} = p_{A}(1-p_{A})(1-\theta) \qquad \sigma_{P}^{2} = p_{A}(1-p_{A})\theta$$

$$\hat{\sigma}_{G}^{2} = \text{MSG} \qquad \hat{\sigma}_{P}^{2} = \frac{1}{n_{c}} (\text{MSP-MSG})$$

#### An estimate of $\theta$ from ANOVA is

$$\hat{\theta} = \frac{\hat{\sigma}_P^2}{\hat{\sigma}_P^2 + \hat{\sigma}_G^2} = \frac{\text{MSP} - \text{MSG}}{\text{MSP} + (n_c - 1)\text{MSG}}$$

For a large sample,

$$\hat{\theta} = \frac{S_A^2}{\tilde{P}_{A\bullet}(1 - \tilde{P}_{A\bullet})} = F_{ST}$$

For multiple loci and alleles, see pg174.

### Some comments on $\theta$ (pg 174 ~ 176)

#### Diploid Data

Table 5.4 Analysis of variance layout for genotypic data in random populations.

Source	d.f.	Sum of Squares	Expected Mean Square <sup>*</sup>
Between populations	r - 1	$2\sum_{i=1}^{r} n_i (\tilde{p}_{Ai} - \tilde{p}_{A.})^2$ = 2(r - 1)\bar{n}s_A^2	$p_A(1-p_A) \left[ (1-F) + 2(F-\theta) + 2n_c \theta \right]$
			$=\sigma_G^2+2\sigma_I^2+2n_c\sigma_P^2$
Individuals in populations	$\sum_{i=1}^{r} (n_i - 1)$	$\sum_{i=1}^{r} n_i (\tilde{p}_{Ai} + \tilde{P}_{AAi} - 2\tilde{p}_{Ai}^2)$	$p_A(1-p_A)\left[(1-F)\right.$
	= n r	$=2r\bar{n}\tilde{p}_{A.}(1-\tilde{p}_{A.})-\frac{1}{2}r\bar{n}\tilde{H}_{A.}$	$+2(F-\theta)]$
		$-2(r-1)\bar{n}s_A^2$	$= \sigma_G^2 + 2\sigma_I^2$
Alleles in individuals	$\sum_{i=1}^{r} n_i$	$\sum_{i=1}^{r} n_i (\tilde{p}_{Ai} - \tilde{P}_{AAi})$	$p_A(1-p_A)(1-F)$
	= n.	$=\frac{1}{2}r\bar{n}\tilde{H}_{A}$	$=\sigma_G^2$

Estimation of F,  $\theta$ , f (pg177~179)

$$\hat{F} = \frac{\hat{\sigma}_P^2 + \hat{\sigma}_I^2}{\hat{\sigma}_P^2 + \hat{\sigma}_I^2 + \hat{\sigma}_G^2}$$

$$= 1 - \frac{2n_c \text{MSG}}{\text{MSP} + (n_c - 1)\text{MSI} + n_c \text{MSG}}$$

$$= 1 - \frac{S_3}{S_2}$$

$$\hat{F} = \frac{\hat{\sigma}_P^2}{\hat{\sigma}_P^2 + \hat{\sigma}_I^2 + \hat{\sigma}_G^2}$$

$$= \frac{\text{MSP} - \text{MSI}}{\text{MSP} + (n_c - 1)\text{MSI} + n_c \text{MSG}}$$

$$= \frac{S_1}{S_2}$$

$$\hat{F} = 1 - \frac{\tilde{H}_A}{\tilde{p}_A (1 - \tilde{p}_A)} \qquad \hat{\theta} = \frac{S_A^2}{\tilde{p}_A (1 - \tilde{p}_A)}$$

$$\hat{f} = \frac{\hat{F} - \hat{\theta}}{1 - \hat{\theta}}$$

Effects of Evolutionary Forces (on  $\theta$  )

When different mating system and mutation or migration process are specified, the expected

variance of  $\theta$  can be estimated by the methods of Cockerham and Weir (1983) Under random mating,  $F = \theta$  and f = 0.

- Any avoidance of mating between relatives will cause  $F < \theta$  and f < 0.
- Different patterns of differences for the two estimates of F and  $\theta$  <u>at different loci</u> indicate that there are forces other than nonrandom mating affecting these loci.
- The effects of selection on the F statistics were detailed by Cockerham (1973).

If forces such as mutation are involved, ....(pg 180 ~183)

$$\theta = \frac{1}{1+4N\mu}$$
  
Migration :  $\theta = \frac{1}{1+4Nm}$ 

### **Population Subdivision**

 
 Table 5.5
 Analysis of variance layout for a three-level sampling hierarchy in random populations.

Source	d.f.	M.S.	Expected M.S.*
Between populations	r-1	MSP	$(1 - F) + 2(F - \theta_{\rm S}) + 2n_{c1}(\theta_{\rm S} - \theta_{\rm P}) + 2n_{c2}\theta_{\rm P}$
Subpopulations in populations	$\sum_{i=1}^{r} (s_i - 1)$ = s r	MSS	$\begin{array}{l} (1-F)+2(F-\theta_{\rm S}) \\ +2n_{c3}(\theta_{\rm S}-\theta_{\rm P}) \end{array}$
Individuals in subpopulations	$\sum_{i=1}^{r} \sum_{j=1}^{s_i} (n_{ij} - 1) = n_{} - s.$	MSI	$(1-F)+2(F- heta_{ m S})$
Alleles in individuals	$\sum_{i=1}^{r} \sum_{j=1}^{s_i} n_{ij}$ $= n_{}$	MSG	(1 - F)

\*To be multiplied by  $p_A(1-p_A)$ .  $n_{c1}, n_{c2}, n_{c3}$  are defined in the text.

Three-Level Hierarchy Four-Level Hierarchy

### **Genetic Distance**

Geometric distance and Genetic distance

Geometric Distances :

$$d_{P\theta} = P_{\sqrt{\sum_{u=1}^{v} (P_u - q_u)^2}}$$

Coancestry as Distance

$$\theta_{t+1} = \frac{1}{2N} + (1 - \frac{1}{2N})\theta_t$$
  
If  $\theta_0 = 0$ ,  $\theta_t = 1 - (1 - \frac{1}{2N})^t$ 

will give information about t, the time since the populations diverged.

Specifically

$$d = -\ln(1-\theta) = -t\ln(1-\frac{1}{2N}) \approx \frac{t}{2N}$$

is an appropriate distance for divergence due to drift. For equal sample sizes from two populations in which

the frequencies of allele  $A_u$  at locus l, the estimator becomes

$$\hat{\theta} = \frac{\sum_{l} \left\{ \frac{1}{2} \sum_{u} (\tilde{p}_{lu_{1}} - \tilde{p}_{lu_{2}})^{2} - \frac{1}{2(2n-1)} \left[ 2 - \sum_{u} (\tilde{p}_{lu_{1}}^{2} + \tilde{p}_{lu_{2}}^{2}) \right] \right\}}{\sum_{l} (1 - \sum_{u} \tilde{p}_{lu_{1}} \tilde{p}_{lu_{2}})}$$

Nei's Genetic Distance

$$I = \frac{\sum_{l} \sum_{u} \widetilde{P}_{lu_{1}}^{2} - \widetilde{P}_{lu_{2}}^{2}}{\sqrt{\sum_{l} \sum_{u} \widetilde{P}_{lu_{1}}^{2} \sum_{l} \sum_{u} \widetilde{P}_{lu_{2}}^{2}}}$$

Nei's distance is appropriate for long-term evolution when populations diverge because of drift and mutation. The distance is proportional to the time since divergence in the special case of the infinite alleles mutation model and equilibrium in the ancestry The Coancestry distance population. İS appropriate divergence due to drift only, and no assumptions need to be made about the ancestral population.

Variance of Distance Estimates

Summary

Homework: Exercise 5.3 (pg 200)