

## § 2: HWE:

• Hardy-Weinberg equilibrium:

$p_A$  = allele frequency of A

Under HWE,

$$\begin{aligned} P_{AA} &= p_A^2 \\ P_{Aa} &= 2p_A(1-p_A) \\ P_{aa} &= (1-p_A)^2 \end{aligned}$$

I.e. Genotypes are  $\text{Binom}(2, p_A)$

Recall Binomial conditions

- ① Success and failure occur
- ② See  $n$  (here  $n=2$ ) outcomes
- ③ outcomes are independent ("random mating")
- ④ Infinite population size

• Departures from HWE can be because of

- ① No random mating (e.g. pop structure, inbreeding)
- ② Genotyping errors
- ③ Small population sizes

• Deviations from HWE

$$P_{AA} = p_A^2 + p_A(1-p_A)f$$

$$P_{Aa} = 2p_A(1-p_A)(1-f)$$

$$P_{aa} = (1-p_A)^2 + p_A(1-p_A)f$$

$$f=0 \Rightarrow \text{HWE}$$

•  $f$  = "inbreeding frequency" =  $F_{IS}$  (measure of relative mating)  
↑ "fixation index" =  $F_{ST}$  (measure of population differentiation)  
2 different models

$$x_j = \begin{cases} 1 & \text{if allele is A} \\ 0 & \text{if allele is a} \end{cases}$$

$$\text{Var}(x_j) = p_A(1-p_A)$$

$$\begin{aligned} \text{cov}(x_j, x_j') &= E[x_j x_j'] - E[x_j] E[x_j'] \\ &= p_{AA} - p_A^2 \\ &= f \underbrace{p_A(1-p_A)}_{= \text{var}(x_j)} \end{aligned}$$

$$\Rightarrow f = \text{cov}(x_j, x_j')$$

• Disequilibrium coefficient  $f$ :

$$D = p_{AA} - p_A^2 \quad (\text{difference b/t genotype freq. and what is expected at HWE})$$

$$p_{AA} = p_A^2 + D_A$$

$$p_{Aa} = 2p_A(1-p_A) - 2D_A$$

$$p_{aa} = (1-p_A)^2 + D_A$$

• Note:  $\max\{-p_A^2, -(1-p_A)^2\} \leq D_A \leq p_A(1-p_A)$

proof,

$$0 \leq p_{AA} \leq p_A \quad \text{since max occurs at}$$

$$\begin{aligned} 0 \leq p_A^2 + D_A \leq p_A \\ \Rightarrow -p_A^2 \leq D_A \leq p_A - p_A^2 \end{aligned}$$



Do same thing for  $p_{Aa}$

Recall  $n_{AA}$   $n_{Au}$   $n_{uA}$   $n_{uu}$   $n = n_{AA} + n_{Au} + n_{uA} + n_{uu}$

$$\hat{p}_A = \frac{n_{AA} + n_{Au}}{n}$$

$$\hat{p}_{AA} = \frac{n_{AA}}{n}$$

$$\hat{D}_A = \hat{p}_{AA} - \hat{p}_A^2$$

Can get mean and variance of the sampling dist of  $\hat{D}_A$

By properties of MLE

$$\hat{D}_A \sim N\left(\underbrace{E(\hat{D}_A)}_{\approx 0 \text{ under } H_0: D_A = 0}, \text{var}(\hat{D}_A)\right)$$

$$\text{so } \frac{\hat{D}_A}{\sqrt{\text{var}(\hat{D}_A)}} \approx N(0, 1)$$

↑ compare to  $N(0, 1)$

$$\text{under Null: } E(\hat{D}_A) \approx 0$$

$$\text{var}(\hat{D}_A) \approx \frac{1}{n} p_A^2 (1-p_A)^2$$

$$\chi^2_A = \frac{n \hat{D}_A^2}{\hat{p}_A (1-\hat{p}_A)} \approx \chi^2_1 \quad \text{b/c } \frac{\hat{D}_A}{\frac{1}{n} p_A (1-p_A)} \approx N(0, 1)$$

Approach 2:  $\chi^2$  - goodness-of-fit

Observed	$n \hat{p}_A$	$n p_A$	$n p_A$
Expected	$n \hat{p}_A$	$2n \hat{p}_A (1 - \hat{p}_A)$	$n (1 - \hat{p}_A)^2$
observed - expected	$n \hat{D}_A$	$-2n \hat{D}_A$	$n \hat{D}_A$

$$\chi^2_A = \sum_{\text{categories}} \frac{(\text{observed} - \text{expected})^2}{\text{expected}} \sim \chi^2_1$$

2 parameters under full  $\leftarrow p_A + p_B$   
 1 parameter under null  $\leftarrow p_A$   
 $2 - 1$

$$= \frac{(n \hat{D}_A)^2}{n \hat{p}_A^2} + \frac{(-2n \hat{D}_A)^2}{2n \hat{p}_A (1 - \hat{p}_A)} + \frac{(n \hat{D}_A)^2}{n (1 - \hat{p}_A)^2}$$

$$= \frac{n \hat{D}_A^2}{\hat{p}_A^2 (1 - \hat{p}_A)^2}$$

• Approach 3: Exact tests

① Calculate  $P_r(n_{AA}, n_{Au}, n_{aA} | n_A, n_a)$

for all possible values of  $n_{AA}, n_{Au}, n_{aA}$  given

$$\begin{array}{l} n_A = n_{Au} + 2n_{AA} \\ n_a = n_{aA} + 2n_{aa} \end{array}$$

↑ fixed
↑ vary these

Ex:  $n_A = 5, n_a = 5$

observed	$n_{AA}$	$n_{Au}$	$n_{aA}$	$P_r(n_{AA}, n_{Au}, n_{aA}   n_A, n_a)$
→ 2	2	1	2	$p_1$
1	1	3	1	$p_2$
0	0	5	0	$p_3$

② Sort probabilities and add up smaller values

$$p_3 < p_1 < p_2$$

$$p_1 + p_3 = p\text{-value} \quad \text{b/c we saw } p_1$$

$$P_r(n_{AA}, n_{Au}, n_{aA}) = \frac{n!}{n_{AA}! n_{Au}! n_{aA}!} (p_A^2)^{n_{AA}} (2p_A(1-p_A))^{n_{Au}} [(1-p_A)^2]^{n_{aA}}$$

$$P_r(n_A, n_a) = \frac{(2n)!}{n_A! n_a!} (p_A)^{n_A} (1-p_A)^{n_a}$$

only under HWF  
 $n_A \sim \text{Bin}(2n, p_A)$

$$(n_{AA}, n_{Au}, n_{aA}) \sim \text{Mult}(n, p_A^2, 2p_A(1-p_A), (1-p_A)^2)$$

$$P_r(n_{AA}, n_{Au}, n_{aA} | n_A, n_a) = \frac{P_r(n_{AA}, n_{Au}, n_{aA})}{P_r(n_A, n_a)}$$

$$= \frac{n! n_{AA}! n_{Au}! n_{aA}! 2^{n_{Au}}}{n_{AA}! n_{Au}! n_{aA}! (2n)!} \left[ \frac{(p_A^2)^{n_{AA}} (p_A(1-p_A))^{n_{Au}} [(1-p_A)^2]^{n_{aA}}}{p_A^{n_A} (1-p_A)^{n_a}} \right]$$

$$= p_A^{2n_{AA} + n_{Au} - n_A} (1-p_A)^{2n_{aA} + n_{Au} - n_a}$$

$$= p_A^0 (1-p_A)^0$$

$$= 1 //$$

• So probabilities do not depend on  $p_A$

If  $n$  is large, too computationally expensive to find all values of  $n_{AA}, n_{Aa}, n_{aa}$

So take  $n_A$  alleles at  $n_a$  alleles, randomly pair, find  $n_{AA}, n_{Aa}, n_{aa}$ , calculate prob of this  
 ↑ repeat many times  
 ↑ proportion tables less probable  $\approx$  p-value

Approach 4: Likelihood ratio test (ch genome, "G" tests, but no else)

$$L_0 = \max_{p_A} P(n_{AA}, n_{Aa}, n_{aa} \mid p_A, 2p_A(1-p_A), (1-p_A)^2)$$

$$L_1 = \max_{p_{AA}, p_{Aa}, p_{aa}} P(n_{AA}, n_{Aa}, n_{aa} \mid p_{AA}, p_{Aa}, p_{aa})$$

Under  $H_0$ :  $-2 \log \left( \frac{L_0}{L_1} \right) \sim \chi^2_1$   
 ↑ 2 parameters in alt - 1 parameter in  $H_0$   
 $= -2(\log(L_0) - \log(L_1))$