

Stat 2911 Lecture Notes

Class 13, 2017

Uri Keich

© Uri Keich, The University of  
Sydney

Estimating the minor allele frequency in a Hardy-Weinberg (HW) equilibrium, Comparing estimators: MSE and unbiased estimators (Rice 8.4, 8.5)

## Hardy-Weinberg equilibrium

(Rice p.273)

If gene frequencies are at an equilibrium, the genotypes AA, Aa, aa occur in the population with corresponding frequencies  $(1-\theta)^2$ ,  $2\theta(1-\theta)$ ,  $\theta^2$  where  $\theta$  is the frequency of the (minor) allele "a".

How come?

Ignoring the diploid (2-chromosome copy) structure of the parents genome (imagine artificially creating an embryo by randomly drawing individual chromosomes from the population),

$$\begin{array}{ccc}
 \text{"M"} & \text{"F"} & P(M \rightarrow a, F \rightarrow a) = \\
 \downarrow & \downarrow & \\
 a/A & a/A & P(M \rightarrow a, F \rightarrow A) = P(M \rightarrow A, F \rightarrow a) =
 \end{array}$$

Having established those frequencies, note that taking the diploid structure into account we have:

$$P(M \rightarrow a) = \theta^2 + \frac{1}{2} \cdot 2 \cdot \theta(1-\theta) = \theta,$$

which, reassuringly, is the same as when we ignored the diploidy.

Regardless, we will next assume the HW model applies.

We will analyze a sample from a Chinese population of HK in 1937 which surveys the genotype of the MN erythrocyte (red blood cell) antigenic group. For this gene, genotype=phenotype, or *wysiwyG*, (alleles are not dominant/recessive)

The sample:

blood type	M	MN	N	
genotype	AA	Aa	aa	
count	342	500	187	(1029 in total)
	$x_1$	$x_2$	$x_3$	

Goal: estimate the frequency  $\theta$  of the minor allele "a" in the population.

Model: the sample is taken from a multinomial distribution with  $P = ((1-\theta)^2, 2\theta(1-\theta), \theta^2)$  and  $n = 1029$   
 $\underline{X} = (\underline{x}_1, \underline{x}_2, \underline{x}_3)$

This is different from estimating  $\Theta = (\theta_1, \theta_2, \theta_3)$  given a multinomial sample: our parameters here depend on a single  $\theta$ .

We begin with estimating  $\theta$  using the method of moments. First, we express  $\theta$  in terms of moments of the dist.

What is the dist. of  $\bar{X}_3$ ?

$$\bar{X}_3 \sim \text{Binom}(n, \theta)$$

$$\Rightarrow E(\bar{X}_3) = \quad \Rightarrow \theta = \sqrt{E(\bar{X}_3)/n}$$

$$\Rightarrow \tilde{\theta} = \sqrt{\hat{E}(\bar{X}_3)/n} = \sqrt{\bar{X}_3/n} \stackrel{\text{in our case}}{=} \sqrt{\frac{187}{4029}} \approx 0.4263.$$

Does something look odd about  $\tilde{\theta}$ ?

Clearly,  $\tilde{\theta}$  does not take into account all available data!

Can we do better?

Imagine estimating a Bernoulli( $p$ ) using only half the data!

What is the proportion of the "a" allele in the sample?

$$\frac{X_1 \cdot 0 + X_2 \cdot 1 + X_3 \cdot 2}{2n}$$

Extending this to a random sample and taking expectations:

$$\begin{aligned} E\left[\frac{\bar{X}_1 \cdot 0 + \bar{X}_2 \cdot 1 + \bar{X}_3 \cdot 2}{2n}\right] &= \frac{1}{2n} [E(\bar{X}_2) + 2E(\bar{X}_3)] \\ &= \frac{1}{2n} [2n\theta(1-\theta) + 2n\theta^2] = \theta \end{aligned}$$

$$\bar{X}_2 \sim \text{Binom}(n, 2\theta(1-\theta))$$

as expected

$$\Rightarrow \theta = \frac{E(\bar{X}_2) + 2E(\bar{X}_3)}{2n} \Rightarrow \hat{\theta} = \frac{\bar{X}_2 + 2\bar{X}_3}{2n} \stackrel{\text{in our case}}{\approx} 0.4247$$

Both  $\tilde{\theta}$  and  $\hat{\theta}$  are moment estimators for this problem.

Which one is better?

Intuitively,  $\hat{\theta}$  because it uses all available data.  
(why?)

MLE

Recall the multinomial pmf:  $P_{nr}(n_1, \dots, n_r) = \binom{n}{n_1 \dots n_r} p_1^{n_1} \dots p_r^{n_r}$

In our case  $p = ((1-\theta)^2, 2\theta(1-\theta), \theta^2)$  and  $n=1029$  so

$$\ell(\theta) = \underbrace{x_1 \log[(1-\theta)^2]}_{2\log(1-\theta)} + \underbrace{x_2 \log[2\theta(1-\theta)]}_{\log 2 + \log \theta + \log(1-\theta)} + \underbrace{x_3 \log(\theta^2)}_{2\log \theta} + \log(x_1 x_2 x_3)$$

$$\begin{aligned}\ell'(\theta) &= -\frac{2x_1}{1-\theta} + \frac{x_2}{\theta} - \frac{x_2}{1-\theta} + \frac{2x_3}{\theta} \\ &= \frac{(x_2+2x_3)(1-\theta) - (2x_1+x_2)\theta}{\theta(1-\theta)} \\ &= \frac{(x_2+2x_3) - 2\theta(x_1+x_2+x_3)}{\theta(1-\theta)} \\ &= \frac{(x_2+2x_3)/2n - \theta}{\theta(1-\theta)/2n}\end{aligned}$$

Sgn  $\ell'$ :

$\nearrow$	$\circ$	$\searrow$
$+$	$ $	$-$

$\theta \in (0, 1)$

if  $x_2+2x_3 > 0$

$$\Rightarrow \hat{\theta}_{MLE} = \frac{x_2+2x_3}{2n} \quad (\text{same if } x_2+2x_3 = 0)$$

Same as  $\hat{\theta}$ , boosting our confidence in  $\hat{\theta}$  but how do we rigorously compare estimators?

Note that even a poor estimator might "get lucky" for any particular sample.

For example, given a Poisson( $\lambda$ ) sample  $x_1, \dots, x_n$  we can estimate  $\hat{\lambda} = x_1$  instead of  $\bar{x}$ . Still, if  $\lambda=2$  and  $x_1=2$  it would probably do better than  $\bar{x}$ . Hence, to compare estimators in a meaningful way we need to consider their performance on all possible samples.

So far our setup was: given a single sample  $x_1, \dots, x_n$  from a dist.  $F(\theta)$  we construct an estimate of  $\theta$  by  $\hat{\theta} = \varphi(x_1, \dots, x_n)$ , e.g.,  $\hat{\theta} = \bar{x}$ . To consider all possible samples we recall our assumption the sample  $x_1, \dots, x_n$  is a particular realization of the  $F(\theta)$ -distributed RVs  $\bar{X}_1, \dots, \bar{X}_n$ , so we define the RV

$$\textcircled{M} = \varphi(\bar{X}_1, \dots, \bar{X}_n).$$

Example. Our estimate for the Bernoulli( $\theta$ ) case was  $\bar{x}$ . If we now consider the sample consisting of the ind.  $\bar{X}_1, \dots, \bar{X}_n$  Bernoulli( $\theta$ ) RVs, then our estimator  $\bar{x}$  corresponds to the RV  $\textcircled{M} = \frac{1}{n} \sum_i \bar{X}_i$ . Strictly speaking,  $\textcircled{M} = \textcircled{M}(n)$

How can we gauge the quality of an estimator?

What does an estimator need to do?

The estimator  $\hat{\theta}$  needs to estimate the unknown  $\theta$ .

When is it doing a good job?

The error  $\hat{\theta} - \theta$  needs to be small.

Which samples should we consider?

All, each weighted by its probability, and in order to avoid cancellations between positive and negative errors we need to take  $|\cdot|$  or  $(\cdot)^2$  of the error.

The mean square error (MSE) of an estimator  $\hat{\theta}$  is defined as,

$$\begin{aligned} \text{MSE}(\hat{\theta}) &:= E_{\theta} \left[ (\hat{\theta} - \theta)^2 \right] \\ &= V(\hat{\theta} - \theta) + [E(\hat{\theta} - \theta)]^2 \\ &= V(\hat{\theta}) + \underbrace{[E(\hat{\theta}) - \theta]}_{\text{bias}}^2 \quad \text{both positive} \end{aligned}$$

An estimator is called unbiased if  $E_{\theta}(\hat{\theta}) = \theta$  (bias = 0).

For unbiased estimators,  $\text{MSE}(\hat{\theta}) = V(\hat{\theta})$ .

## Examples

1)  $X_1, \dots, X_n$  are ind. Bernoulli( $\theta$ ) RVs.

$$\hat{\theta} = \hat{\theta}_n = \frac{1}{n} \sum_{i=1}^n \bar{X}_i$$

$$\begin{aligned} E(\hat{\theta}) &= E\left(\frac{1}{n} \sum_{i=1}^n \bar{X}_i\right) \\ &= \frac{1}{n} \sum_{i=1}^n E(\bar{X}_i) \end{aligned}$$

$$= \theta$$

$\Rightarrow \hat{\theta}$  is estimator

$$\begin{aligned} \text{MSE}(\hat{\theta}_n) &= V(\hat{\theta}_n) = V\left(\frac{1}{n} \sum_{i=1}^n \bar{X}_i\right) \\ &\stackrel{?}{=} \frac{1}{n^2} \sum_{i=1}^n V(\bar{X}_i) \\ &= \frac{\theta(1-\theta)}{n} \end{aligned}$$

Do you see a problem?

$\theta$  is, of course, unknown, but we can:

(i) use  $\theta(1-\theta) \leq \frac{1}{4}$  for  $\theta \in [0, 1]$  to conclude that  $\text{MSE}(\hat{\theta}) \leq \frac{1}{4n}$

(ii) estimate  $\text{MSE}(\hat{\theta})$  by plugging in  $\hat{\theta} = \bar{x}$

Regardless, clearly  $\text{MSE}(\hat{\theta}_n) \xrightarrow{n} \dots$  (does it make sense?)