

## 1.5.1 Testing Hardy-Weinberg Proportions (HWE)

n	AA	AB	BB	l1	q	l2	2(l1 - l2)
100	36	48	16	-101.33	0.6	-101.33	0
100	30	60	10	-89.79	0.6	-93.01	6.5
100	45	30	25	-106.71	0.6	-113.81	14.2

Consider the above three samples, each of 100 individuals. Each has 120 A alleles, so the MLE of q is 0.6, but different genotypic counts in each genotype.

- With count  $n_c$  and probability  $p_c$  for class  $c$ ,  
 $\lambda = \text{const} + \sum_c n_c \log(p_c)$  with  $\sum_c p_c = 1$

## 1.5.2 Testing ABO blood group models

- |           | factor freq. |       | phenotype frequencies |         |       |                |
|-----------|--------------|-------|-----------------------|---------|-------|----------------|
|           | A            | B     | A                     | B       | AB    | O              |
| Data      |              |       | 0.422                 | 0.206   | 0.078 | 0.294          |
| H1 theory | p            | q     | p(1-q)                | (1-p)q  | pq    | (1-p)(1-q)     |
| H1 fitted | 0.500        | 0.284 | 0.358                 | 0.142   | 0.142 | 0.358          |
| H2 theory | p            | q     | p(p+2r)               | q(q+2r) | 2pq   | r <sup>2</sup> |
| H2 fitted | 0.295        | 0.155 | 0.411                 | 0.194   | 0.091 | 0.303          |
- Bernstein reported ABO blood types on a sample of 502 individuals: 42.2% type A, 20.6% type B, 7.8% type AB and 29.4% type O. (Did he drop 2 individuals?)
  - For the general model, log-likelihood is  
 $\lambda = 502 ( .422 \log.422 + .206 \log.206 + .078 \log.078 + .294 \log.294 ) = -626.71$

## Testing HWE ctd.

- With no constraints, MLE of  $p_c$  is  $n_c/n$ , and maximized value of the log-likelihood is  
 $\lambda_1 = \sum_c n_c \log(n_c/n) = \sum_c n_c \log(n_c) - n \log(n)$ .
- Assuming HWE, estimated genotype frequencies are  $(q^2, 2q(1-q), (1-q)^2) = (0.36, 0.48, 0.16)$  and  
 $\lambda_2 = n_1 \log 0.36 + n_2 \log 0.48 + n_3 \log 0.16$ .
- Now, if HWE is true,  $2 \log \Lambda = 2(\lambda_1 - \lambda_2)$  is approximately chi-squared ( $\chi^2$ ) with 1 df. If HWE is not true  $\Lambda$  will tend to be larger.
- In our three examples, the values are 0, 6.5 and 14.2 (see the table). What do we conclude?

## Testing H1: indep factor model

- H1: A and B are independently inherited factors.
- Frequency of individuals having the factor A is 0.500 and of B is 0.284.
- Bernstein observed that independence of the factors would give an AB frequency of  $0.500 \times 0.284 = 0.142$  much larger than the 0.078 observed.
- Under H1 the estimated frequencies are as shown in Table, and the log-likelihood is  
 $\lambda_1 = 502 ( .422 \log.358 + .206 \log.142 + .078 \log.142 + .294 \log.358 ) = -647.50$
- Twice the log-likelihood difference is 41.58, and would be the value of a  $\chi$ -squared random variable with 1 df if H1 were true.
- Clearly, H1 is rejected.

## Testing H2: Bernstein's method

- Under H2: A and B are the two non-null alleles of a single system. We assume HWE.
- If the three alleles A, B and O have frequencies p, q and r ( $p+q+r = 1$ ), then the frequencies of the four blood types are  $p^2+2pr$  etc, as shown in the Table.
- Bernstein pointed out that  $\alpha = \text{freq-type-A} + \text{freq-type-O} = (p+r)^2$ , or  $1-\sqrt{\alpha} = (1-p-r) = q$ .
- Similarly if  $\beta = \text{freq-type-B} + \text{freq-type-O} = (q+r)^2$  or  $1-\sqrt{\beta} = p$ , and  $\sqrt{(\text{freq-type-O})} = r$ .
- The sum of these three numbers should be one.
- For his data,  $(1 - \sqrt{0.422+0.294}) + (1 - \sqrt{0.206+0.294}) + \sqrt{0.294} = 0.99$ , which is close to one, suggesting a good fit.

## A note about degrees of freedom

- For the general hypothesis, with no constraints, the df is the number of categories less 1.
- We always lose 1, as the sum of the probabilities is always fixed to be 1.
- Here, Total number of categories = 4
- Lose 1 degree of freedom for fixed total:  $4-1 = 3$
- Lose 1 for each parameter estimated:  
Under H1 we estimate p and q.  
Under H2, we estimate p, q, and r, but  $r = 1-p-q$ .  
So in each case we estimate 2 parameters.
- So there is  $3-2 = 1$  df to test each of H1 and H2.

## Testing H2 using log-likelihoods

- Or, as with H1, we may perform a likelihood ratio test.
- Finding the MLEs of the parameters p, q and r is not simple. In fact, we shall see later that these MLEs are  $p = 0.2945$  and  $q = 0.1547$ , with the resulting fitted frequencies given in the table.
- The log-likelihood is  
 $\lambda^2 = 502 ( .422 \log .4114 + .206 \log .1942 + .078 \log .0911 + .294 \log .3033 ) = -627.52$
- Twice the log-likelihood difference between this and the general alternative is now only 1.62. Again, this is the value of a  $\chi^2$  random variable with 1 df if H2 is true.
- H2 is not rejected.

### 1.6.1 Gene counting: case of recessive allele A

- The three genotypes are AA, AB and BB, with counts say  $t_i$ , ( $i=1,2,3$ ). Now,  $n_1 = t_1$ , but the counts of AB and BB are unobservable since B is dominant to A.
- If counts,  $t_2$  and  $t_3$ , were known, then the number of A alleles is  $(2 t_1 + t_2)$ , and the MLE of q would be  $(2 t_1 + t_2)/2n$ . (M-step)
- Now  $P(\text{AB} | \text{AB or BB}) = \{2 q(1-q)\}/\{1 - q^2\} = 2q/\{1+q\}$
- So  $E(t_2 | t_2 + t_3 = 64) = 64 \{2q\}/\{1+q\}$  (E-step)
- The EM-algorithm implements the sequence of iterates shown in the Table. Starting from any value (e.g.  $q=0.5$ ), the proportion  $2q/(1+q)$  is computed, and the 64 individuals of dominant phenotype divided into the expected numbers of AB and BB, respectively.
- Then a new value of q is estimated as  $(2 t_1 + t_2)/2n$ .

## Table of gene-counting iterates

current	current	recess- ive	dominant phenotype		new estimate
q	$2q/(1+q)$	t1=36	t2+t3 = 64		q
0.5	0.667	36	42.67	21.33	0.573
0.573	0.729	36	46.64	17.36	0.593
0.593	0.745	36	47.66	16.34	0.598
0.598	0.749	36	47.91	16.09	0.600
0.600	0.750	36	48.00	16.00	0.600

Theoretical result:  $\lambda(\hat{\theta}) \geq \lambda(\theta^*)$ .

Thus the EM algorithm for finding MLEs alternates E-steps and M-steps. The likelihood is non-decreasing over the process. Where the likelihood surface is unimodal, convergence to the MLE is assured, although it may be slow. Where computable, evaluate the (log)-likelihood to assess convergence.

For multinomial data, let  $n_c$  be actual data-counts, and  $m_{c^*}$  complete-data counts for idealized data. So  $\lambda^* = \sum_{c^*} m_{c^*} \log(p_{c^*})$ , and finding the ECDLL just means finding

$$E(m_{c^*} | n_c) = n_c \Pr(c^* | c, \theta^*) = n_c \frac{p_{c^*}(\theta^*)}{\sum_{c^*} p_{c^*}(\theta^*)}$$

## 1.6.2 EM for multinomial data

EM algorithm for multinomial data

In latent variable problems, suppose the actual data are  $Y$ , and the ideal data that would make the problem easy are  $(Y, X)$ . The complete-data log-likelihood is

$$\lambda^* = \log \Pr((Y, X) = (y, x)).$$

The actual log-likelihood to be maximized is

$$\lambda = \log \Pr(Y = y) = \log \left( \sum_x \Pr((Y, X) = (y, x)) \right).$$

E-step (expectation):

At the current estimate  $\theta^*$  compute ECDLL

$$H_y(\theta; \theta^*) = E_{\theta^*}(\log P_{\theta}(X, Y) | Y = y)$$

M-step (maximization):

Maximize  $H_y(\theta; \theta^*)$  w.r.t.  $\theta$  to obtain a new estimate  $\bar{\theta}$ .

## 1.6.3 The ABO log-likelihood

- Observed counts are  $Y = (NA, NB, NAB, NO)$  with frequencies  $p(p+2r)$ ,  $q(q+2r)$ ,  $2pq$ , and  $r^2$ .
- $\lambda = \sum_{i} \{ \text{obs counts} \} Y_i \log P(Y_i)$
- Complete-data count are genotype count  $X = (nAA, nAO, \dots)$
- $\lambda^* = \sum_{i} \{ \text{all counts} \} X_i \log P(X_i)$
- Do not confuse  $\lambda$  and  $\lambda^*$ .  $\lambda^*$  is just a tool that lets us maximize  $\lambda$ .
- Compute  $E(\lambda^* | Y)$ . -- in the multinomial case this just involves imputing the "hidden" counts -- but only because  $\lambda^*$  is a linear function of these counts.

## 1.6.4 Estimation of ABO allele frequencies

- The EM algorithm is one of the easiest ways to find the MLEs of the ABO blood group allele frequencies.
- See table, next page.
- E-step: partition the A phenotypes into expected counts of AA and AO genotypes, and similarly B into BB and BO
- $P(\text{AO} \mid \text{blood type A}) = \frac{2pr}{p^2 + 2pr} = \frac{2r}{p + 2r}$   
 $P(\text{BO} \mid \text{blood type B}) = \frac{2qr}{q^2 + 2qr} = \frac{2r}{q + 2r}$ .
- M-step: new estimates of p and q are  
 $p = (2P(\text{AA}) + P(\text{AO}) + P(\text{AB}))/2$ , and  
 $q = (2P(\text{BB}) + P(\text{BO}) + P(\text{AB}))/2$ .
- Note p does not change monotonely, but  $\lambda$  does. Note we are interested in the current value of  $\lambda$ , not of  $\lambda^*$ .

Current values				Phenotype A		Phenotype B	
p	q	2r/ (p +2r)	2r/ (q +2r)	P(type A) =0.422		P(type B) = 0.206	
				AA	AO	BB	BO
0.3	0.3	0.73	0.73	0.115	0.307	0.056	0.150
0.308	0.170	0.77	0.86	0.096	0.326	0.029	0.177
0.298	0.156	0.79	0.87	0.091	0.331	0.026	0.180
0.295	0.155	0.79	0.88	0.089	0.333	0.025	0.181
	.....	ph AB	phen O	New values		$\lambda$	
	.....	0.078	0.294	p	q	-687.00	
	.....	0.078	0.294	0.308	0.170	-629.00	
	.....	0.078	0.294	0.298	0.156	-627.57	
	.....	0,078	0.294	0.295	0.155	-627.53	
	.....	0.078	0.294	0.295	0.155	-627.52	