

# LDSP - Linear Detection of Selection in Pooled data

*Hussein Al-Asadi, Matthew Stephens*

## 1 Introduction

### Phase I - Calculate effective coverage using haplotypic information:

The intuition is that data at each SNP are binomial counts, which help estimate the frequency of a SNP in a pool. But by combining information across multiple corrected SNPs, you can improve the estimate.

**Phase II - detect selection using effective coverage :** fit a linear model.

## 2 Phase I - calculating effective coverage

### 2.1 The Prior

Let  $y = (y_1, y_2, \dots, y_p)'$  denote the vector of allele frequencies in the study sample. Let  $E[y_i] = \mu_i$  and  $M$  denote the  $(2m) \times p$  panel (i.e.  $2m$  haplotypes and  $p$  SNPs). As in (Wen & Stephens, 2010), we assume

$$y^{true}|M \sim N_p(\mu, \Sigma) \quad (1)$$

(Wen & Stephens, 2010) derive estimates for  $\mu$  and  $\Sigma$ :

$$\hat{\mu} = (1 - \theta)f^{panel} + \frac{\theta}{2}1 \quad (2)$$

$$\hat{\Sigma} = (1 - \theta)^2 S + \frac{\theta}{2}(1 - \frac{\theta}{2})I \quad (3)$$

and  $S$  is obtained from  $\Sigma^{panel}$ , specifically,

$$S_{i,j} = \begin{cases} \Sigma_{i,j}^{panel} & i = j \\ e^{-\frac{\rho_{i,j}}{2m}} \Sigma_{i,j}^{panel} & i \neq j \end{cases} \quad (4)$$

$\rho_{i,j} = -4Nc_{i,j}d_{i,j}$  where  $d_{i,j}$  is the physical distance between markers  $i$  and  $j$ ,  $N$  is the effective diploid population size,  $c_{i,j}$  is the average rate of crossover per unit physical distance, per meiosis, between sites  $i$  and  $j$  (so that  $c_{i,j}d_{i,j}$  is the genetic distance between sites  $i$  and  $j$ ).

and,

$$\theta = \frac{(\sum_{i=1}^{2m-1} \frac{1}{i})^{-1}}{2m + (\sum_{i=1}^{2m-1} \frac{1}{i})^{-1}} \quad (5)$$

and,

$$\Sigma_{i,j}^{panel} = \begin{cases} f_i^{panel}(1 - f_i^{panel}) & i = j \\ f_{ij}^{panel} - f_i^{panel} f_j^{panel} & i \neq j \end{cases} \quad (6)$$

where  $f_{ij}^{panel}$  is the panel frequency of the haplotype “1-1” consisting of loci  $i$  and loci  $j$

## 2.2 The likelihood

Let  $(n_i^0, n_i^1)$  denote the counts of 0 and 1 alleles at SNP  $i$ ,  $n_i = n_i^0 + n_i^1$ , and  $y_i^{true}$  is the population frequency of the SNP  $i$  “1” allele.

$$\begin{aligned} n_i^1 | y_i^{true} &\sim \text{Bin}(n_i, y_i^{true}) \sim N(n_i y_i, n_i y_i^{true}(1 - y_i^{true})) \\ \implies \frac{n_i^1}{n_i} | y_i^{true} &\sim N(y_i^{true}, \frac{y_i^{true}(1 - y_i^{true})}{n_i}) \end{aligned} \quad (7)$$

let  $y_i^{obs} = \frac{n_i^1}{n_i}$  and replace  $y_i^{true}$  with  $y_i^{obs}$  in the variance (a common simplification). Therefore our equation becomes,

$$y_i^{obs} | y_i^{true} \sim N(y_i^{true}, \frac{y_i^{obs}(1 - y_i^{obs})}{n_i}) \quad (8)$$

$y_1^{obs} | y_1^{true}, y_2^{obs} | y_2^{true}, \dots, y_p^{obs} | y_p^{true}$  are independent therefore we can write,

$$y^{obs} | y^{true} \sim N_p(y^{true}, \text{diag}(\epsilon_1, \dots, \epsilon_p)) \quad (9)$$

where  $\epsilon_i = \frac{y_i^{obs}(1 - y_i^{obs})}{n_i}$

### 2.2.1 Avoiding $y_i^{obs} = 0$

If the coverage is low, then a frequency estimate can be zero (i.e.  $\frac{n_i^1}{n_i} = 0$ ) which will introduce complications when we must invert matrices. Therefore we make the following modification,

$$y_i^{obs} = \frac{n_i^1 + \frac{1}{2}}{n_i + 1} \quad (10)$$

which has nice properties.

### 2.2.2 Incorporating base quality scores

## 2.3 The Posterior

In the distribution of  $y^{true}$ , we assumed that the panel and study individuals are from the sample population, and the parameters  $\theta$  and  $\rho$  are estimated without

error. Deviations from these assumptions will cause over-dispersion: the true allele frequencies will lie further from their expected values than the model predicts. To allow this, we modify equation 1 by introducing an over-dispersion parameter  $\sigma^2$ .

$$y^{true}|M \sim N_p(\hat{\mu}, \sigma^2 \hat{\Sigma}) \quad (11)$$

We estimate  $\sigma^2$  by maximizing the multivariate normal likelihood:

$$y^{obs}|M \sim N_p(\hat{\mu}, \sigma^2 \hat{\Sigma} + \text{diag}(\epsilon_1, \dots, \epsilon_p)) \quad (12)$$

To obtain the distribution for the true frequencies conditional on the observed data, we use Bayes theorem

$$P(y^{true}|y^{obs}, M) \propto P(y^{obs}|y^{true})P(y^{true}|M)$$

Let,

$$\bar{\Sigma} = \left( \frac{\hat{\Sigma}^{-1}}{\sigma^2} + \text{diag}\left(\frac{1}{\epsilon_1}, \dots, \frac{1}{\epsilon_p}\right) \right)^{-1} \quad (13)$$

and,

$$\bar{\mu} = \bar{\Sigma} \left( \frac{\hat{\Sigma}^{-1}}{\sigma^2} \hat{\mu} + \text{diag}\left(\frac{1}{\epsilon_1}, \dots, \frac{1}{\epsilon_p}\right) y^{obs} \right) \quad (14)$$

since the normal is in the conjugate family,

$$y^{true}|y^{obs}, M \sim N_p(\bar{\mu}, \bar{\Sigma}) \quad (15)$$

Therefore a natural point estimate for  $y^{true}$  is  $\bar{\mu}$ .

Note: we can also calculate effective coverage here (just simply use the reverse mapping of bin -i, normal approximation)

### 2.3.1 Avoiding prior mean bias

*Note: for simplicity we now assume the overdispersion parameter  $\sigma^2 = 1$*

As mentioned above, we assume the the panel and sample individuals are drawn from the same population. This is the never the case in reality but in some applications, the frequencies of alleles of interest have changed significantly but the correlation structure has changed slightly (i.e. very little recombination between nearby SNPs). Therefore we would just like to use the information from SNP correlations.  $L(y_i^{true})$  will do the job.

$$L(y_i^{true}) = P(y^{obs}|y_i^{true}, M) \propto \frac{P(y_i^{true}|y^{obs}, M)}{P(y_i^{true}|M)} \quad (16)$$

We showed above,

$$y_i^{true}|y^{obs}, M \sim N(\bar{\mu}_i, \bar{\Sigma}_{ii}) \quad (17)$$

and,

$$y_i^{true}|M \sim N(\hat{\mu}_i, \hat{\Sigma}_{ii}) \quad (18)$$

Thus,

$$P(y^{obs}|y_i^{true}, M) \propto e^{-\frac{(y_i^{true} - \bar{\mu}_i)^2}{2\bar{\Sigma}_{ii}} + \frac{(y_i^{true} - \hat{\mu}_i)^2}{2\hat{\Sigma}_{ii}}} \quad (19)$$

Completing the square, we can see

$$\frac{-(y_i^{true} - \bar{\mu}_i)^2}{2\bar{\Sigma}_{ii}} + \frac{(y_i^{true} - \hat{\mu}_i)^2}{2\hat{\Sigma}_{ii}} = \frac{1}{2} \left( \frac{1}{\hat{\Sigma}_{ii}} - \frac{1}{\bar{\Sigma}_{ii}} \right) (y_i^{true} - \frac{\hat{\mu}_i \bar{\Sigma}_{ii} - \bar{\mu}_i \hat{\Sigma}_{ii}}{\bar{\Sigma}_{ii} - \hat{\Sigma}_{ii}})^2 + K \quad (20)$$

Let,

$$\tilde{\mu}_i = \frac{\hat{\mu}_i \bar{\Sigma}_{ii} - \bar{\mu}_i \hat{\Sigma}_{ii}}{\bar{\Sigma}_{ii} - \hat{\Sigma}_{ii}} \quad (21)$$

and,

$$\tilde{\sigma}_i^2 = \frac{1}{-\frac{1}{\bar{\Sigma}_{ii}} + \frac{1}{\hat{\Sigma}_{ii}}} \quad (22)$$

$$L(y_i^{true}) \propto e^{-\frac{(y_i^{true} - \tilde{\mu}_i)^2}{2\tilde{\sigma}_i^2}} \quad (23)$$

with  $\tilde{\mu}_i$  being the MLE of  $y_i^{true}$ .

## 2.4 Unequal contributions of individuals to the pool

## 2.5 Computing the Effective Coverage

To calculate effective coverage ( $n_e$ ) and the effective proportion ( $p_e$ ), we approximate the normal likelihood with a binomial likelihood.

### 2.5.1 Simply taking the reverse mapping of the well known binomial to normal transformation

$$p_e = \tilde{\mu}_i \quad (24)$$

and,

$$\frac{p_e(1-p_e)}{n_e} = \tilde{\sigma}_i^2 \quad (25)$$

Then,

$$n_e = \frac{\tilde{\mu}_i(1-\tilde{\mu}_i)}{\tilde{\sigma}_i^2} \quad (26)$$

### 2.5.2 Using the Taylor expansion

Let,

$$f(p) = \log l(p) = \log(p^{n_1}(1-p)^{n-n_1}) \quad (27)$$

Taking the Taylor expansion of  $f(p)$  around its maximum ( $\hat{p}$ ) we get,

$$f(p) \approx f(\hat{p}) + \frac{(p-\hat{p})^2}{2} f''(\hat{p}) \quad (28)$$

Therefore,

$$e^{f(p)} = p^{n_1}(1-p)^{n-n_1} \approx C e^{\frac{-(p-\hat{p})^2}{2 f''(\hat{p})}} \quad (29)$$

working back from equation 23 let

$$\tilde{\mu}_i = \hat{p} \quad (30)$$

and,

$$\tilde{\sigma}_i^2 = \frac{-1}{f''(\hat{p})} = \frac{(1-\hat{p})\hat{p}^2}{n_1} \quad (31)$$

where  $\hat{p} = \frac{n_1}{n}$ . Solving the equations for  $n$ :

$$n = \frac{\tilde{\mu}_i(1-\tilde{\mu}_i)}{\tilde{\sigma}_i^2} \quad (32)$$

which is the same as above!

## 3 Phase II - estimating $\beta$

Let  $f_{i,k,j}$  denote the frequency of the  $j$ th SNP in population  $i$  and replicate  $k$ . Then,

$$\log\left(\frac{1-f_{i,k,j}}{f_{i,k,j}}\right) = \mu_j + \beta_j g_i + \epsilon \quad (33)$$

where  $\epsilon \sim N(0, \sigma_d^2)$ ,  $\sigma_d^2$  is the variance due to drift,  $\mu_j$  is the frequency of the  $j$ th SNP in the founding population and

$$g_i = \begin{cases} -1 & i = 0 \\ 0 & i = 1 \\ 1 & i = 2 \end{cases}$$

The intuition here is that sites with large  $\beta$  coefficients are under selection.

### 3.1 Finding the casual SNP

## 4 Computational Issues

### 4.1 Calculating the inverse of the covariance matrix

$\hat{\Sigma}$  is singular when SNPs are perfectly correlated, to fix this we... (look at notes)

We modify 13 which is now probably non-singular?

$$\bar{\Sigma} = \left( \frac{\hat{\Sigma}^+}{\sigma^2} + \text{diag}\left(\frac{1}{\epsilon_1}, \dots, \frac{1}{\epsilon_p}\right) \right)^{-1} \quad (34)$$