

# Selection regimes through hierarchical MK approaches

Jesus Murga Moreno

October 22, 2019

## 1 Estimation of slightly deleterious mutations

Here we present  $b$  estimator, which calculate the fraction/excess of slightly deleterious mutations at selected sites. This estimation bring us a direct evidence of the strength of purifying selection at selected sites, altogether with the estimation of effectively neutral mutations and strongly deleterious mutations (defined as deleterious mutations not segregating at samples because of fitness). In addition  $b$  represents the variable biasing  $P_i/P_0$  ratio. Then  $b$  is define as the excess of slightly deleterious mutations normalized by the total number of sites.

$$b = \sum_{j=0}^1 \frac{P_{wd(j)}}{P_{0(j)}} \cdot \frac{m_0}{m_i} \quad (1)$$

Assuming beneficial alleles get fixed quickly and barely contribute to polymorphism, the site frequency spectrum of selected sites could be defined by one-side gamma distribution as described in REFS, where mutations will be effectively neutral or deleterious (weakly or strongly). EXPLANATION WHY GAMMA. The distribution is defined by two parameters:  $k$  and  $a$ , governing shape and mean of the distribution. The kurtosis will determinate the number of slightly deleterious mutations segregating on selected sites.

### 1.1 $b$ through $\alpha$ MK

$\alpha$  will be underestimated depending on the presence and proportion of slightly deleterious. In the absence of slightly deleterious mutations in that frequency,  $\alpha$  should be approximately the asymptotic value in a given frequency:

$$\alpha_{(j)} \approx \alpha_a \quad (2)$$

Under this assumption we could redefine the  $P_{i(j)}/P_{0(j)}$  ratio, where the slightly deleterious mutations ( $P_{wd(j)}$ ) tend to be 0 in that frequency.

$$\frac{P_{i(j)} - P_{wd(j)}}{P_{0(j)}} \quad (3)$$

$$P_{wd(j)} \rightarrow 0$$

At frequencies where slightly deleterious mutations are segregating, we could estimate  $P_{wd(j)}$  from equation (2), assuming again  $\alpha_{(j)}$  should be  $\alpha_a$ . Then we redefined the standard  $\alpha$  with equation (3) to include the expected count of slightly deleterious mutations biasing  $\alpha_{(j)}$ .

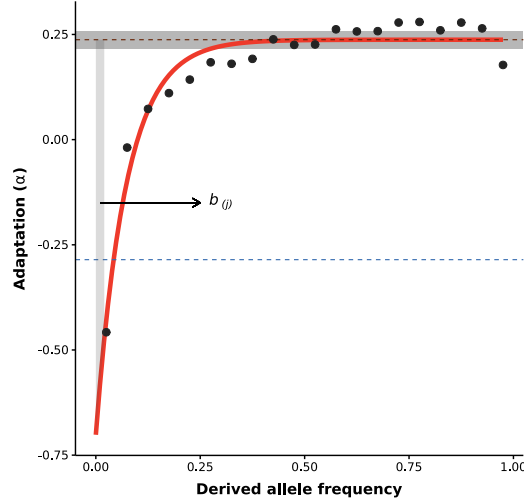
$$\alpha_a = 1 - \frac{(P_{i(j)} - P_{wd(j)}) \cdot D_0}{P_{0(j)} \cdot D_i} \quad (4)$$

We understand  $P_{wd(j)}$  as the expected number of slightly deleterious mutations biasing  $\alpha$  at a given frequency, distorting the frequency spectra of selected sites and therefore underestimating  $\alpha_{(j)}$ .

$$P_{wd(j)} = P_{i(j)} - \frac{(1 - \alpha_a) \cdot P_0 \cdot D_i}{D_0} \quad (5)$$

We just expect slightly deleterious at low frequency, although it depends on several important factors as the DFE, the effective population size ( $N_e$ ) or recent demography events. Nevertheless, we estimated  $P_{wd}$  at the range  $\alpha_0 \rightarrow \alpha_1$ , following the  $\alpha_a$  estimation at the  $\alpha$ MK methodology fitting  $\alpha_{(j)}$ . At the moment  $\alpha_{(j)}$  is within the confidence intervals estimated at  $\alpha$ MK,  $\alpha_{(j)}$  will range around  $\alpha_a$  values. At these frequencies we don't expect slightly deleterious mutations and depending on  $\alpha_{(j)}$  values we could estimate positive or negatives values of  $P_{wd(j)}$  if  $\alpha_{(j)}$  is above or below the asymptote. However, once  $\alpha_{(j)}$  is at the confidence intervals,  $P_{wd(j)}$  values are minimal and  $\alpha_{(j)}$  fluctuation is canceling the total count at  $P_{wd(cilow)} \rightarrow P_{wd(1)}$ . In these way we avoid several assumptions of the DFE at selected sites, not considering where slightly deleterious mutations should segregate *a priori*, maximizing the estimation.

$b_{(j)}$  definition could be easily visualize as the area under the asymptotic curve. At frequencies where  $P_{(i)}$  have an excess of slightly deleterious mutation, the difference between  $\alpha_{(j)}$  and  $\alpha_a$  should be greater than  $\alpha_{(j+1)}$ , assuming slightly deleterious mutation tend to segregate low along the site frequency spectrum.



## 1.2 $b$ through $eMK$

### 1.3 Testing $b$ estimator

To test the sensitivity and accuracy of  $b$  we perform multiple simulations through SLiM software. Forward in-time simulations allow us to compare the real excess of slightly deleterious mutations with our estimator. In order to compare both, we estimate the real excess of slightly deleterious mutations (*true*  $b$ ) dividing the spectrum of selected mutations ( $P_i$ ) (estimated from a gamma distribution) through their fitness coefficients ( $s$ ) on: (1) effectively neutral mutations ( $-1 < N_e s < 1$ ), (2) slightly deleterious mutations ( $-10 < N_e s < 1$ ) and (3) strongly deleterious mutations ( $N_e s < -10$ ). The *true*  $b$  is defined like in equation (1), where  $P_{wd}$  is the sum of slightly and strongly deleterious mutations. Although we don't expect mutations segregating at the range  $N_e s < -10$ , it will depend on the DFE and  $N_e$ , since DFE goes  $0 \rightarrow -\infty$ . We observe some situations where linkage could sweep some of these mutations despite of their low fitness coefficients, although the total number are minimal and tend to segregate at really low frequencies. In this way we compare our estimators with the real excess of slightly deleterious mutations.

We performed a total of 13 simulations. Simulations were performed over a 10Mb with two genomic element (synonymous and non-synonymous sites) simulating a typical coding sequence structure, in a proportion of 1/2 (eg. 004) with a mutation rate of  $1e-9$ , a recombination rate of  $1e-7$  and

a dominance coefficient of 0.5 over  $2e5$  generations with a  $10N_e$  burnin period. We defined synonymous sites as neutral sites introducing only neutral mutations with fixed probability. On the other hand non-synonymous sites were used as selected sites. The spectra of selected sites were defined using a one-side gamma distribution with a  $k$  parameter controlling the shape and  $a$  parameter controlling the mean ( $a$  parameter refers to selection coefficient mean). We check  $b$  values changing in several ways the DFE of selected sites (table 1). In this way we increase or decrease the number of deleterious sites segregating at  $P_i$  using leptokurtics and platykurtic distributions and changes in the strength of selection. Simulations were divided in two categories: neutral and adaptive, where we introduce the presence of adaptive mutations and test the presence of high  $\alpha$  values on  $b$  estimation. In addition we replicate the estimation from two different effective population size ( $N_e = 100$  and  $N_e = 1000$ ). We based all the scenarios at the same SLiM recipe from which we change determinate parameters to test effects on  $b$ .

#### 1.4 $b$ estimation on presence of recent positive selection

Linkage and recent beneficial alleles segregating on selected sites could lead to underestimate  $\alpha$  too. A new methodology combining *abc* frameworks and  $\alpha$ MK approaches developed by Uricchio et al. is able to re-estimate  $\alpha$  values taking into account these variables. Because of *abc* scripts are currently set up to run on the Stanford cluster, we couldn't apply  $b$  estimation at genetics scenarios where adaptation is lead by selective sweeps. At these sites frequency spectrum will be biased due to the proportion of weakly beneficial (determined by  $\alpha_w$  in the methodology) and deleterious alleles. Taking into account  $\alpha_w$  on the site frequency spectrum, we would expect similar results on accuracy and sensitivity since  $b$  only depends on the asymptotic value of  $\alpha$ . In any case, at these kind of complex linkage scenarios would require an extensive exploration to determine possibly errors in  $b$  estimations.