

$S'_B$ : a modified version of the  $S_B$  statistic that is more robust to low sample sizes

Fernando Racimo

October 2021

In Refoyo-Martinez et al. (2020), we defined the  $S_B$  statistic for a specific branch  $k$  of an admixture graph as:

$$S_B = \frac{((\mathbf{p} - \bar{p}\mathbf{1})^T \mathbf{b}_k)^2}{\bar{p}(1 - \bar{p}) \mathbf{b}_k^T \hat{\mathbf{F}} \mathbf{b}_k} \quad (1)$$

Here,  $\mathbf{p}$  is the vector of sample allele frequencies across populations,  $\hat{\mathbf{F}}$  is an estimate of the genome-wide allele frequency covariance matrix, and  $\bar{p}$  is the mean allele frequency among populations. The elements of the branch vector  $\mathbf{b}_k$  are the ancestry contributions of that branch to each of the populations in the leaves of the graph.

This statistic makes use of the sample allele frequencies  $\mathbf{p}$  as an approximation the true population allele frequencies  $\check{\mathbf{p}}$ , which may be particularly poor if a SNP under study has sequence data from a few individuals in a given population.

Conditional on the population allele frequency for a population  $j$ , the sample allele frequency is binomially distributed:

$$p_j | \check{p}_j \sim \text{Bin}(2n_j, \check{p}_j) \quad (2)$$

where  $n_j$  is the number of diploid individuals for which there is reliable genotype data at a particular SNP of interest in population  $j$ . We can approximate the above equation using a Normal distribution:

$$p_j | \check{p}_j \sim \text{Normal}(\check{p}_j, \check{q}_j) \quad (3)$$

where  $\check{q}_j$  is equal to  $\frac{\check{p}_j(1-\check{p}_j)}{2n_j}$ . Like  $\check{p}_j$ ,  $\check{q}_j$  will also not be known, and here we approximate it as  $q_j = \frac{p_j(1-p_j)}{2n_j}$ .

$$p_j | \check{p}_j \sim \text{Normal}(\check{p}_j, q_j) \quad (4)$$

Conditional on knowing the population allele frequencies for all populations, the sample allele frequencies for each population are independent of each other. In vector notation:

$$\mathbf{p} | \check{\mathbf{p}} \sim \text{MVN}(\check{\mathbf{p}}, \text{diag}(\mathbf{q})) \quad (5)$$

The population frequencies are, in turn, assumed to depend on some population-wide ancestral allele frequency  $\mathbf{e}$ , as in Refoyo-Martinez et al. (2020):

$$\check{\mathbf{p}} \sim \text{MVN}(\mathbf{e}\mathbf{1}, \mathbf{e}(\mathbf{1} - \mathbf{e})\mathbf{F}) \quad (6)$$

where  $\mathbf{1}$  is a vector of ones. If we make one further approximation and treat the variance of the conditional distribution as a constant that is not dependent on the mean, we can marginalize the population allele frequencies, and obtain:

$$\mathbf{p} \sim \text{MVN}(\mathbf{e}\mathbf{1}, \text{diag}(\mathbf{q}) + \mathbf{e}(\mathbf{1} - \mathbf{e})\mathbf{F}) \quad (7)$$

We then mean-center the vector  $\mathbf{p}$ :

$$\mathbf{y} = \mathbf{p} - \mathbf{e}\mathbf{1} \sim \text{MVN}(0, \text{diag}(\mathbf{q}) + \mathbf{e}(\mathbf{1} - \mathbf{e})\mathbf{F}) \quad (8)$$

We multiply the mean-centered vector by the branch vector  $b_k$  for a branch of interest and obtain:

$$\mathbf{y}^T \mathbf{b} \sim \text{Normal}(0, \mathbf{b}_k^T \text{diag}(\mathbf{q}) \mathbf{b}_k + e(1-e) \mathbf{b}_k^T \mathbf{F} \mathbf{b}_k) \quad (9)$$

Finally, we derive a statistic that follows a chi-squared distribution under neutrality:

$$\frac{((\mathbf{p} - \bar{p}\mathbf{1})^T \mathbf{b}_k)^2}{\mathbf{b}_k^T \text{diag}(\mathbf{q}) \mathbf{b}_k + e(1-e) \mathbf{b}_k^T \mathbf{F} \mathbf{b}_k} \sim \chi_1^2 \quad (10)$$

If we use the mean sample frequency across populations  $\bar{p}$  as an estimate of the ancestral frequency  $e$ , and also use the empirical covariance matrix  $\hat{\mathbf{F}}$  as an estimate of the true covariance matrix  $\mathbf{F}$ , we can obtain a statistic that penalizes sites in which the number of sampled individuals for a given branch's subtended populations is low:

$$S'_B = \frac{((\mathbf{p} - \bar{p}\mathbf{1})^T \mathbf{b}_k)^2}{\mathbf{b}_k^T \text{diag}(\mathbf{q}) \mathbf{b}_k + \bar{p}(1-\bar{p}) \mathbf{b}_k^T \hat{\mathbf{F}} \mathbf{b}_k} \quad (11)$$