

# On the heterozygosity of an admixed population

Simina M. Boca\*, Lucy Huang<sup>†</sup> and Noah A. Rosenberg<sup>‡</sup>

October 12, 2019

A population is termed *admixed* if its members possess recent ancestry from two or more separate sources. As a result of the fusion of source populations with different genetic variants, admixed populations can exhibit high levels of genetic variation, reflecting contributions of their multiple ancestral groups. For a model of an admixed population derived from  $K$  source groups, we obtain a relationship between its level of genetic variation, as measured by heterozygosity, and its proportions of admixture from the various source populations. We show that the heterozygosity of the admixed population is at least as great as that of the least heterozygous source population, and that it potentially exceeds the heterozygosities of *all* of the source populations. The admixture proportions that maximize the heterozygosity possible for an admixed population formed from a specified set of source populations are also obtained under specific conditions. We examine the special case of  $K = 2$  source populations in detail, characterizing the maximal admixture in terms of the heterozygosities of the two source populations and the value of  $F_{ST}$  between them. In this case, the heterozygosity of the admixed population exceeds the maximal heterozygosity of the source groups if the divergence between them, measured by  $F_{ST}$ , is large enough, namely above a certain bound that is a function of the heterozygosities of the source groups. We present applications to simulated data as well as to data from human admixture scenarios, providing results useful for interpreting the properties of genetic variability in admixed populations.

**Keywords.** Admixture, allele frequencies, heterozygosity, population genetics

**Mathematics subject classification.** 15A63, 92D10

---

\*Innovation Center for Biomedical Informatics, Department of Oncology, Department of Biostatistics, Bioinformatics and Biomathematics, Georgetown University Medical Center, Washington, DC 20007, United States. Corresponding author, email: [smb310@georgetown.edu](mailto:smb310@georgetown.edu).

<sup>†</sup>Bioinformatics Graduate Program, University of Michigan, Ann Arbor, MI 48109, United States.

<sup>‡</sup>Department of Biology, Stanford University, Stanford, CA 94305, United States.

## 1 Introduction

Admixed populations are populations that possess ancestry from multiple source groups. They result from the fusion of populations that have long been separated, in processes such as long-distance migration and hybrid-zone formation at population boundaries.

Several features of ancestry and allele frequencies are characteristic of admixed populations (Chakraborty, 1986; Long, 1991; Verdu & Rosenberg, 2011; Gravel, 2012). In an admixed population, the values of allele frequencies are typically intermediate between those of the various sources. Unlike in a mixture that pools individuals taken from separate populations, in an admixed population, alleles from different sources cooccur within individuals. The contributions from the source populations are each large enough that most members of an admixed population have ancestry in more than one source group.

In admixed populations, the history of mating among populations is recent enough that time has not yet eroded differences among admixed individuals in their relative proportions of ancestry. This feature of high levels of variability in admixture proportions has been central to studies of admixed populations. Investigations of such phenomena as the timing and contributions of the source populations (Verdu & Rosenberg, 2011; Gravel, 2012), the effect of admixture levels on assortative mating patterns (Risch *et al.*, 2009; Zou *et al.*, 2015), and the genetic basis of traits in admixed populations (Buerkle & Lexer, 2008; Zhu *et al.*, 2008) all make use of variation in levels of admixture levels across admixed individuals.

A second aspect of variability in admixed populations is potentially of interest: the variability of alleles as captured by genetic diversity measures. The effect of admixture in contributing to increased genetic diversity, however, is not simple. For example, in a study of the genetics of populations founded by relatively small groups, Mooney *et al.* (2018) examined genetic diversity in admixed and non-admixed populations, some of which were regarded as founder populations. Mooney *et al.* (2018) observed that genetic diversity was relatively high in multiple admixed populations of Latin America. This pattern was observed even for populations that, on the basis of small population size and past history of isolation, might have been expected to have relatively low levels of genetic diversity.

Here, to deepen understanding of the relationship between admixture and genetic variability, we focus in admixed populations on levels of genetic diversity computed from allele

frequencies, rather than on variability among individuals in admixture proportions. For a model of an admixed population with  $K$  source groups, we derive a relationship between genetic diversity, as measured by heterozygosity, and proportions of admixture drawn from the various source populations. The model is the same model we have previously used to examine the genetic differentiation between admixed populations and their source groups, as measured by  $F_{ST}$  (Boca & Rosenberg, 2011). We show that for all values of the admixture contributions from the source populations, the heterozygosity of the admixed population is greater than or equal to the smallest of the source population heterozygosities. We further examine the maximal values of the heterozygosity of the admixed population over the space of possible admixture proportions. We consider in more detail special cases of the admixture model with  $K = 2$  and  $K = 3$  source populations, providing explicit results for  $K = 2$  in terms of relatively few parameters. Finally, we use simulations and example analyses from admixed human populations to illustrate the mathematical results.

## 2 Notation and model

We consider a model with  $K$  source populations and an admixed population arising from these sources. A single locus is considered, with  $J$  distinct alleles. In Sections 2.1, 2.2, and 2.3, respectively, we define the expected heterozygosity, fixation index, and allele-frequency dot product statistics that we use in our analysis. In Section 2.4, we introduce the admixture model. Notation is summarized in Table 1.

### 2.1 Expected heterozygosity

We denote by  $p_{kj}$  the frequency of allelic type  $j$ ,  $1 \leq j \leq J$ , in source population  $k$ ,  $1 \leq k \leq K$ , with  $0 \leq p_{kj} \leq 1$ . The expected heterozygosity is a measure of genetic diversity, giving the probability that two alleles randomly drawn from the population differ in type.

**Definition 1.** The *expected heterozygosity* in a population for a given locus with  $J$  distinct alleles is defined as  $H = 1 - \sum_{j=1}^J p_j^2$ , where  $p_j$  is the frequency of allelic type  $j$ .

We denote by  $H_k$  the expected heterozygosity of source population  $k$  at a locus. We have  $0 \leq H_k < 1$ , with  $H_k = 0$  if and only if source population  $k$  has only a single allelic type of nonzero frequency. We refer to expected heterozygosity simply as heterozygosity.

## 2.2 Fixation index

The fixation index  $F_{ST}$  is a measure of genetic divergence among a set of subpopulations. In its general form, it is computed from  $H_S$ , the mean of the heterozygosities of the subpopulations, and  $H_T$ , the heterozygosity of a population formed by pooling the subpopulations into a single “total” population.

**Definition 2.** The *fixation index*,  $F_{ST}$  is defined as  $F_{ST} = (H_T - H_S)/H_T$ , where  $H_T$  is the heterozygosity of the total population and  $H_S$  is the mean heterozygosity across subpopulations.

The fixation index can be regarded as a measure of genetic divergence between two populations, with  $F_{k\ell}$  denoting the value of  $F_{ST}$  between source populations  $k$  and  $\ell$ . We assume that the two subpopulations have the same contribution to the overall population, so that they are weighted equally in producing the total population. We also assume that when pooled together, they produce a polymorphic population. In other words, we disallow the case in which there is some allelic type  $j$  for which  $p_{kj} = p_{\ell j} = 1$ .

For this pairwise scenario,  $H_S = (H_k + H_\ell)/2$ ,  $H_T = 1 - \sum_{j=1}^J (\frac{p_{kj} + p_{\ell j}}{2})^2$ , and

$$F_{k\ell} = \frac{[1 - \sum_{j=1}^J (\frac{p_{kj} + p_{\ell j}}{2})^2] - \frac{H_k + H_\ell}{2}}{1 - \sum_{j=1}^J (\frac{p_{kj} + p_{\ell j}}{2})^2}. \quad (1)$$

We can observe by the Cauchy-Schwarz inequality that  $0 \leq F_{k\ell} \leq 1$ , with  $F_{k\ell} = 0$  requiring  $p_{kj} = p_{\ell j}$  for all  $j$ .  $F_{k\ell} = 1$  requires  $H_S = H_k = H_\ell = 0$ .

## 2.3 Allele frequency dot product

We will have occasion to use a quantity,  $C_{k\ell}$ , the probability that, when randomly drawing one allele from population  $k$  and one allele from population  $\ell$ , the two alleles differ in type. For population  $k$ , let  $\underline{p}_k$  denote a  $J \times 1$  column vector of its allele frequencies.  $C_{k\ell}$  can then be written as 1 minus the dot product of the allele frequency vectors of populations  $k$  and  $\ell$ :

$$C_{k\ell} = 1 - \underline{p}_k' \cdot \underline{p}_\ell = 1 - \sum_{j=1}^J p_{kj} p_{\ell j}. \quad (2)$$

Note that this quantity can be viewed as a generalization of heterozygosity to two populations, as  $H_k = C_{kk}$ . Because we exclude the case in which populations  $k$  and  $\ell$  are fixed for

the same allelic type,  $C_{k\ell}$  strictly exceeds 0, so that  $0 < C_{k\ell} \leq 1$ . The upper bound of 1 is achieved if populations  $k$  and  $\ell$  share no allelic types in common.

We can rewrite eq. 1 as

$$F_{k\ell} = \frac{2C_{k\ell} - H_k - H_\ell}{2C_{k\ell} + H_k + H_\ell}. \quad (3)$$

If  $F_{k\ell} < 1$ , then we can solve for  $C_{k\ell}$ :

$$C_{k\ell} = \left( \frac{H_k + H_\ell}{2} \right) \left( \frac{1 + F_{k\ell}}{1 - F_{k\ell}} \right). \quad (4)$$

Recall that  $F_{k\ell} = 1$  implies  $H_k = H_\ell = 0$ , so that populations  $k$  and  $\ell$  each have only a single allelic type with nonzero frequency. We have excluded the case in which the two populations are fixed for the same allelic type; hence, they must be fixed for different allelic types, and  $C_{k\ell} = 1$ .

By the Cauchy-Schwarz inequality,  $1 - \sqrt{(1 - H_k)(1 - H_\ell)} \leq C_{k\ell} \leq 1$  (Mehta *et al.*, 2019, eq. 7). Equality in the lower bound requires  $p_{kj} = p_{\ell j}$  for all  $j$ , and hence  $H_k = H_\ell$ . Rewriting this inequality with eq. 4, we obtain the allowable space of  $F_{k\ell}$  given  $H_k, H_\ell \in [0, 1)$ :

$$F_{k\ell} \in \left[ \frac{2 - H_k - H_\ell - 2\sqrt{(1 - H_k)(1 - H_\ell)}}{2 + H_k + H_\ell - 2\sqrt{(1 - H_k)(1 - H_\ell)}}, \frac{2 - H_k - H_\ell}{2 + H_k + H_\ell} \right]. \quad (5)$$

The lower limit is achieved if and only if the two populations  $k$  and  $\ell$  are identical, with  $H_k = H_\ell$  and  $p_{kj} = p_{\ell j}$  for all  $j$ . The upper limit is achieved if and only if populations  $k$  and  $\ell$  share no allelic types in common.

Appendix A of Mehta *et al.* (2019) shows that given  $H_k$  and  $H_\ell$  in  $[0, 1)$ , if the number of distinct alleles  $J$  is not fixed, then we can choose allele frequency vectors  $\underline{p}_k$  and  $\underline{p}_\ell$  such that each  $C_{k\ell}$  in value in  $[1 - \sqrt{(1 - H_k)(1 - H_\ell)}, 1]$  is achievable. The lower bound is achievable only if  $H_k = H_\ell$ . Hence, each value in the interval in eq. 5 for  $F_{ST}$  is also achievable by some pair of allele frequency vectors  $\underline{p}_k$  and  $\underline{p}_\ell$ , the lower bound only if  $H_k = H_\ell$ .

## 2.4 Admixture model

In our  $K$ -source-population model,  $K \geq 2$ , we follow Section 2.2 in assuming that no two populations are fixed for the same allelic type. We now make a stronger assumption that no two populations are identical, so that for each  $(k, \ell)$ , some  $j$  exists for which  $p_{kj} \neq p_{\ell j}$ .

Following a commonly used approach for describing variation in an admixed population, we treat allele frequencies in the admixed population as linear combinations of those of the

source populations (e.g. Pritchard *et al.*, 2000; Boca & Rosenberg, 2011). It is convenient to assume that no source population can have its vector of allele frequencies written as the linear combination of vectors of allele frequencies of other source populations; otherwise, an admixed population would not have a unique representation as a linear combination of sources. We thus assume that not only are no two source populations identical, no source population can be described as an admixture of two or more of the other sources.

Note that the assumption that no population is a linear combination of the others also excludes linear combinations with one or more negative coefficients. In addition, because the maximal number of vectors of length  $J$  that can be linearly independent is  $J$ , the assumption implies that  $J \geq K$ . A succinct way of describing the linear independence assumption is that if we define the  $J \times K$  matrix of allele frequencies in the source populations,

$$P = \begin{pmatrix} p_{11} & p_{21} & \cdots & p_{K1} \\ p_{12} & p_{22} & \cdots & p_{K2} \\ \cdots & \cdots & \cdots & \cdots \\ p_{1J} & p_{2J} & \cdots & p_{KJ} \end{pmatrix} = (\underline{p}_1, \underline{p}_2, \dots, \underline{p}_K), \quad (6)$$

then we assume that  $P$  has rank  $K$ .

For the admixed population generated from the  $K$  source populations, we denote by  $\gamma_k$  the admixture fraction for source population  $k$ , so that for each  $k$  with  $1 \leq k \leq K$ , fraction  $\gamma_k$  of the ancestry of the admixed population,  $0 \leq \gamma_k \leq 1$ , derives from source population  $k$ . We denote by  $\underline{\gamma}$  the  $K \times 1$  column vector of admixture fractions. This vector represents a point in the simplex  $\Delta^{K-1}$ , the set of all vectors of  $K$  nonnegative entries with  $\sum_{k=1}^K \gamma_k = 1$ .

The frequency of allele  $j$  in the admixed population is denoted by  $\bar{p}_j$ . According to the linear combination assumption, we have

$$\bar{p}_j = \sum_{k=1}^K \gamma_k p_{kj}. \quad (7)$$

### 3 General case: $K$ source populations

Our goal is to study the heterozygosity of the admixed population. Using Definition 1 with eq. 7, we compute the heterozygosity for the admixed population, which we denote by  $H_{\text{adm}}$ :

$$H_{\text{adm}} = 1 - \sum_{j=1}^J \bar{p}_j^2 = 1 - \sum_{j=1}^J \left( \sum_{k=1}^K \gamma_k P_{kj} \right)^2. \quad (8)$$

The heterozygosity of the admixed population can be written in terms of the heterozygosities of the source populations and the dot products of the allele frequencies. Using eq. 4 and the formula for  $H_{\text{adm}}$  in eq. 8, we have:

$$H_{\text{adm}} = \sum_{k=1}^K \gamma_k^2 H_k + 2 \sum_{k=1}^{K-1} \sum_{\ell=k+1}^K \gamma_k \gamma_\ell C_{k\ell} \quad (9)$$

$$= \sum_{k=1}^K \gamma_k^2 H_k + \sum_{k=1}^{K-1} \sum_{\ell=k+1}^K \gamma_k \gamma_\ell (H_k + H_\ell) \left( \frac{1 + F_{k\ell}}{1 - F_{k\ell}} \right). \quad (10)$$

The last simplification can be made only for  $F_{k\ell} \neq 1$ ; if  $F_{k\ell} = 1$ , then eq. 9 is used, or, as noted after eq. 4,  $(H_k + H_\ell)(1 + F_{k\ell})/(1 - F_{k\ell})$  is understood to equal 2.

With the formula for  $H_{\text{adm}}$  established, we now explore how  $H_{\text{adm}}$  varies in relation to the admixture fractions  $\underline{\gamma}$ . Given the allele frequencies  $P$ , we determine how small and how large  $H_{\text{adm}}$  can be over the space of possible values of  $\underline{\gamma}$ . We write  $H_m$  for the smallest heterozygosity among the source populations,  $H_m = \min_{k \in \{1, 2, \dots, K\}} H_k$ , and  $H_M$  for the largest heterozygosity among the source populations,  $H_M = \max_{k \in \{1, 2, \dots, K\}} H_k$ .

#### 3.1 Minimum of $H_{\text{adm}}$ in terms of the ancestry proportions

For the minimum of  $H_{\text{adm}}$  over vectors  $(\gamma_1, \gamma_2, \dots, \gamma_K)$ , we can immediately observe from the form of eq. 10 that for a fixed set of source population allele frequencies  $P$ ,  $H_{\text{adm}}$  is minimized as a function of the admixture fractions when the admixed population consists of only one of the source populations.

**Proposition 3.** The minimum of  $H_{\text{adm}}$  as a function of the ancestry proportions  $\underline{\gamma}$  is  $H_m = \min_{k \in \{1, 2, \dots, K\}} H_k$ , the smallest heterozygosity among the source populations, and it is obtained when the admixed population consists solely of that source population.

*Proof.* To obtain this result, we use eq. 10 and the fact that  $H_k \geq H_m$  for all  $k$ :

$$\begin{aligned} H_{\text{adm}} &= \sum_{k=1}^K \gamma_k^2 H_k + \sum_{k=1}^{K-1} \sum_{\ell=k+1}^K \gamma_k \gamma_\ell (H_k + H_\ell) \left( \frac{1 + F_{k\ell}}{1 - F_{k\ell}} \right) \\ &\geq \sum_{k=1}^K \gamma_k^2 H_m + \sum_{k=1}^{K-1} \sum_{\ell=k+1}^K 2\gamma_k \gamma_\ell H_m \\ &= \left( \sum_{k=1}^K \gamma_k \right)^2 H_m = H_m. \end{aligned}$$

Because equality is achieved when  $\gamma_m = 1$  and  $\gamma_k = 0$  for all  $k \neq m$ , we have shown that the minimal value of  $H_{\text{adm}}$  as a function of the ancestry proportions is  $H_m$ .  $\square$

The result applies whether or not  $H_1, H_2, \dots, H_K$  are mutually distinct. If two or more of  $H_1, H_2, \dots, H_K$  are tied for the minimal heterozygosity  $H_m$ , then the minimum of  $H_{\text{adm}}$  is achieved at each vector associated with complete ancestry from one of the minimally heterozygous populations.

A consequence of Proposition 3 is that if all  $K$  populations have the same heterozygosity  $H_m$ , then  $H_{\text{adm}} > H_m$  for all ancestry vectors  $\underline{\gamma}$  with two or more nonzero entries. In particular, note that  $F_{k\ell} > 0$  for each  $(k, \ell)$ ,  $k \neq \ell$ , by the assumption that each pair of source populations has distinct allele frequencies. Hence,  $(H_k + H_\ell)(1 + F_{k\ell})/(1 - F_{k\ell}) > 2H_m$  for each  $(k, \ell)$ ,  $k \neq \ell$ . Because at least one product  $\gamma_k \gamma_\ell$  is positive, the inequality  $\gamma_k \gamma_\ell (H_k + H_\ell)(1 + F_{k\ell})/(1 - F_{k\ell}) \geq 2\gamma_k \gamma_\ell H_m$  is strict for at least one  $(k, \ell)$ , so that  $H_{\text{adm}} > (\sum_{k=1}^K \gamma_k)^2 H_m = H_m$ . This same reasoning shows that if two or more populations are tied with heterozygosity  $H_m$ , then  $H_{\text{adm}} > H_m$  for each  $\underline{\gamma}$  with two or more nonzero entries.

### 3.2 Maximum of $H_{\text{adm}}$ in terms of the ancestry proportions

To obtain the maximum of  $H_{\text{adm}}$  over the space of values of  $\underline{\gamma}$ , we write eq. 9 as a quadratic form in terms of the ancestry proportions,

$$H_{\text{adm}}(\underline{\gamma}) = \underline{\gamma}' A \underline{\gamma}.$$



Here,  $\underline{\gamma}'$  represents the transpose of the column vector  $\underline{\gamma}$  and  $A$  is the  $K \times K$  symmetric matrix with the  $H_k$  on the diagonal and the  $C_{k\ell}$  off the diagonal:

$$A = \begin{pmatrix} H_1 & C_{12} & \dots & C_{1K} \\ C_{12} & H_2 & \dots & C_{2K} \\ \dots & \dots & \dots & \dots \\ C_{1K} & C_{2K} & \dots & H_K \end{pmatrix} = \underline{\mathbf{1}}\underline{\mathbf{1}}' - \begin{pmatrix} \sum_{j=1}^J p_{1j}^2 & \sum_{j=1}^J p_{1j}p_{2j} & \dots & \sum_{j=1}^J p_{1j}p_{Kj} \\ \sum_{j=1}^J p_{1j}p_{2j} & \sum_{j=1}^J p_{2j}^2 & \dots & \sum_{j=1}^J p_{2j}p_{Kj} \\ \dots & \dots & \dots & \dots \\ \sum_{j=1}^J p_{1j}p_{Kj} & \sum_{j=1}^J p_{2j}p_{Kj} & \dots & \sum_{j=1}^J p_{Kj}^2 \end{pmatrix} = \underline{\mathbf{1}}\underline{\mathbf{1}}' - P'P, \quad (11)$$

where  $P$  is the  $J \times K$  allele frequency matrix (eq. 6) and  $\underline{\mathbf{1}}$  is a  $K \times 1$  vector of ones.

Maximizing  $H_{\text{adm}}$  in terms of  $\underline{\gamma}$  is equivalent to finding  $\max_{\underline{\gamma} \in \Delta^{K-1}} \underline{\gamma}'A\underline{\gamma}$  subject to  $\underline{\mathbf{1}}'\underline{\gamma} = 1$ . We denote by  $\underline{\gamma}_{\text{arg max}}$  the location of the maximal value of  $H_{\text{adm}}$ . We first observe that  $\underline{\gamma}_{\text{arg max}}$  is sometimes interior to the simplex, and that it sometimes lies at a vertex.

**Proposition 4.** Consider the case of  $K$  source populations,  $K \geq 2$ .

- (i) There exists some collection of source population allele frequencies  $P$  and some collection of admixture proportions  $\underline{\gamma}$  for which the heterozygosity of the admixed population exceeds the heterozygosity  $H_M$  of the most heterozygous source population.
- (ii) There exists some collection of source population allele frequencies  $P$  for which *no* collection of admixture proportions  $\underline{\gamma}$  produces an admixed population with heterozygosity greater than the heterozygosity  $H_M$  of the most heterozygous source population.

*Proof.* (i) Consider  $K$  populations, each with different allele frequencies, but identical heterozygosity:  $\underline{p}_k \neq \underline{p}_\ell$  for  $k \neq \ell$  but  $H_k = H$  for  $k = 1, 2, \dots, K$ . Suppose that a locus has  $K + 1$  distinct alleles, and that the allele frequencies are  $\underline{p}_1 = (\frac{1}{2}, \frac{1}{2}, 0, 0, \dots, 0)$ ,  $\underline{p}_2 = (\frac{1}{2}, 0, \frac{1}{2}, 0, \dots, 0)$ ,  $\dots$ ,  $\underline{p}_K = (\frac{1}{2}, 0, 0, \dots, 0, \frac{1}{2})$ . By eq. 9,  $H_{\text{adm}} = \frac{3}{4} - \frac{1}{4} \sum_{k=1}^K \gamma_k^2$ , which is minimized if and only if  $\sum_{k=1}^K \gamma_k^2 = 1$ , or  $\underline{\gamma} = \underline{e}_k$  for some  $k$ . The minimal value of  $H_{\text{adm}}$  is thus  $\frac{1}{2}$ , all other values of the admixture proportions resulting in  $H_{\text{adm}} > H = \frac{1}{2}$ .

(ii) Consider  $K$  populations and a locus with  $K$  distinct alleles. Suppose that the number of distinct alleles at the locus is  $k$  for population  $k$ , with  $\underline{p}_k = (\frac{1}{k}, \dots, \frac{1}{k})$ . Hence,  $H_k = 1 - \frac{1}{k}$  and, in particular,  $H_1 < \dots < H_K$ . We show that  $H_{\text{adm}} \leq H_K$  irrespective of  $\underline{\gamma}$ .

By eq. 9,

$$H_{\text{adm}} = 1 - \left( \gamma_1 + \frac{\gamma_2}{2} + \dots + \frac{\gamma_K}{K} \right)^2 - \dots - \left( \frac{\gamma_K}{K} \right)^2.$$

By the Cauchy-Schwarz inequality:

$$\left[ \left( \gamma_1 + \frac{\gamma_2}{2} + \dots + \frac{\gamma_K}{K} \right)^2 + \dots + \left( \frac{\gamma_K}{K} \right)^2 \right] K \geq \left( \gamma_1 + \frac{\gamma_2}{2} 2 + \dots + \frac{\gamma_K}{K} K \right)^2 = \left( \sum_{k=1}^K \gamma_k \right)^2 = 1.$$

Thus,  $H_{\text{adm}} \leq 1 - \frac{1}{K} = H_K$ .  $\square$

Note that it is trivial to see that in general,  $\max_{\underline{\gamma} \in \Delta^{K-1}} H_{\text{adm}}(\underline{\gamma}) \geq \max\{H_1, \dots, H_K\}$ : the  $K$  source populations simply correspond to the  $K$  vertices of the simplex. This result that the maximal  $H_{\text{adm}}$  is at least as great as the heterozygosity of the most heterozygous source population immediately implies  $\max_{\underline{\gamma} \in \Delta^{K-1}} H_{\text{adm}}(\underline{\gamma}) \geq \sum_{k=1}^K H_k / K$ .

Having established that the maximum can be at a vertex or an interior point of the simplex, we now provide a general theorem. The theorem gives the location of the maximum when it lies in the interior of  $\Delta^{K-1}$ , rather than on the boundary, assuming a condition applies on the allele frequencies. The proof is in Appendix 1.

**Theorem 5.** Suppose that  $\underline{1}'(P'P)^{-1}\underline{1} \neq 1$ . Suppose also that  $\frac{A^{-1}\underline{1}}{\underline{1}'A^{-1}\underline{1}} \in \Delta^{K-1}$ . Then the maximum of  $H_{\text{adm}}$  as a function of the ancestry proportions  $\underline{\gamma} \in \Delta^{K-1}$  is attained at  $\underline{\gamma}_{\text{arg max}} = \underline{\gamma}^*$ , where:

$$\underline{\gamma}^* = \frac{A^{-1}\underline{1}}{\underline{1}'A^{-1}\underline{1}} = \frac{(P'P)^{-1}\underline{1}}{\underline{1}'(P'P)^{-1}\underline{1}}.$$

The maximum is equal to:

$$H_{\text{adm}}(\underline{\gamma}^*) = \frac{1}{\underline{1}'A^{-1}\underline{1}} = 1 - \frac{1}{\underline{1}'(P'P)^{-1}\underline{1}}.$$

If  $\frac{A^{-1}\underline{1}}{\underline{1}'A^{-1}\underline{1}} \notin \Delta^{K-1}$ , then  $\underline{\gamma}_{\text{arg max}}$  lies on the boundary of the set  $\{\underline{\gamma} : \underline{1}'\underline{\gamma} = 1 \text{ and } \underline{\gamma} \in \Delta^{K-1}\}$ .

The following corollary, also proven in Appendix 1, further describes the possible locations of the maximum of  $H_{\text{adm}}$ . Note that if the maximum is not at  $\underline{\gamma}^*$ , then it lies at a point that has some elements equal to 0, with the nonzero subvector having a similar form to  $\underline{\gamma}^*$ , but in a lower number of dimensions. Thus, the maximum can occur in a scenario in which the admixture involves only a strict subset of the source populations.

Consider a nonempty subset  $\mathcal{S} \subset \{1, 2, \dots, K\}$ . Define by  $A_{\mathcal{S}}$  the  $|\mathcal{S}| \times |\mathcal{S}|$  matrix that has diagonal terms  $H_k$  for each  $k \in \mathcal{S}$  and off-diagonal terms  $C_{k\ell}$  for each distinct  $k, \ell \in \mathcal{S}$ . Additionally, denote by  $P_{\mathcal{S}}$  the matrix consisting of the columns of  $P$  corresponding to the subset  $\mathcal{S}$ .  $P_{\mathcal{S}}$  contains the allele frequencies for the source populations in  $\mathcal{S}$ .

**Corollary 6.** Suppose that  $\underline{1}'(P'_{\mathcal{S}}P_{\mathcal{S}})^{-1}\underline{1} \neq 1$  for all nonempty  $\mathcal{S} \subset \{1, 2, \dots, K\}$ . Then the maximum of  $H_{\text{adm}}$  as a function of the ancestry proportions  $\underline{\gamma} \in \Delta^{K-1}$  is attained at a point that has nonzero elements for some nonempty subset of the source populations  $\mathcal{S}^* \subset \{1, 2, \dots, K\}$ . The nonzero subvector of ancestry proportions at the location of the maximum is equal to  $\underline{\gamma}_{\mathcal{S}^*} = \frac{A_{\mathcal{S}^*}^{-1}\underline{1}}{\underline{1}'A_{\mathcal{S}^*}^{-1}\underline{1}}$ .

In particular, note that  $\underline{\gamma}_{\text{arg max}} = \underline{\gamma}^*$  corresponds to  $\mathcal{S}^* = \{1, 2, \dots, K\}$ : all source populations contribute nonzero admixture fractions. The  $K$  vertices of the simplex  $\Delta^{K-1}$  correspond to the cases of  $\mathcal{S}^* = \{k\}$ , at which only one source population contributes.  $\mathcal{S}$  has  $2^K - 1$  nonempty subsets.

## 4 $K = 2$ source populations

With some general results established for the case of arbitrary  $K$ , we now focus on the simplest case, with  $K = 2$  source populations contributing to the admixed population.

We continue to exclude the scenario in which the allele frequencies for the two source populations are identical, so that we assume  $\underline{p}_1 \neq \underline{p}_2$ . Noting that  $\gamma_2 = 1 - \gamma_1$ , we can consider  $H_{\text{adm}}$  in terms of a single admixture coefficient  $\gamma_1$ , the admixture fraction of the first population, with  $\gamma_1 \in [0, 1]$ . Using eqs. 9 and 10 with this substitution, we obtain:

$$H_{\text{adm}} = \gamma_1^2 H_1 + (1 - \gamma_1)^2 H_2 + 2\gamma_1(1 - \gamma_1)C_{12} \quad (12)$$

$$= \gamma_1^2 H_1 + (1 - \gamma_1)^2 H_2 + \gamma_1(1 - \gamma_1)(H_1 + H_2) \frac{1 + F_{12}}{1 - F_{12}} \quad (13)$$

$$= \gamma_1^2(H_1 + H_2 - 2C_{12}) - 2\gamma_1(H_2 - C_{12}) + H_2. \quad (14)$$

In particular, we note that from eq. 13 that  $H_{\text{adm}}$  is increasing as a function of  $F_{12}$ .

From eq. 14, we can see that  $H_{\text{adm}}$  is concave down as a function of  $\gamma_1$ . We have  $d^2 H_{\text{adm}}/d\gamma_1^2 = 2(H_1 + H_2 - 2C_{12})$ . By Definition 1 and eq. 2,  $2(H_1 + H_2 - 2C_{12}) = -2 \sum_{j=1}^J (p_{1j} - p_{2j})^2$ . Because  $\underline{p}_1 \neq \underline{p}_2$ ,  $p_{1j} \neq p_{2j}$  for at least one choice of  $j$ , and hence  $d^2 H_{\text{adm}}/d\gamma_1^2 < 0$ . By symmetry,  $H_{\text{adm}}$  is also concave down as a function of  $\gamma_2$ .

To illustrate eq. 13, for fixed values of  $H_1$  and  $H_2$ , Figure 1 plots  $H_{\text{adm}}$  as a function of  $\gamma_1$  for a variety of values of  $F_{12}$ . The figure illustrates the concave-down quadratic nature of  $H_{\text{adm}}$  as a function of  $\gamma_1$ . We observe that for each value of  $F_{12}$  considered, the minimum of  $H_{\text{adm}}$  occurs at  $(\gamma_1, \gamma_2) = (0, 1)$ , reflecting the result of Proposition 3 that the minimum occurs when the admixed population consists solely of the less heterozygous source population. In accord with the fact that in eq. 13,  $H_{\text{adm}}$  increases for fixed  $H_1$ ,  $H_2$ , and  $\gamma_1$  with increasing  $F_{12}$ , the value at the maximum increases with increasing  $F_{12}$ . The location of the maximum lies at a value of  $\gamma_1 \geq \frac{1}{2}$ , decreasing with increasing  $F_{12}$ . This location has a pattern where for larger values of  $F_{12}$ , it lies interior to the unit interval, and for smaller values of  $F_{12}$ , it occurs when the admixed population consists solely of the more heterozygous source population. We now consider this pattern in more detail.

#### 4.1 Minimum and maximum of $H_{\text{adm}}$ in terms of the ancestry proportions

Applying the general results from Section 3.1 describing the minimum and maximum of  $H_{\text{adm}}$  as a function of  $\underline{\gamma}$ , by Proposition 3, the minimum of  $H_{\text{adm}}$  is simply  $\min\{H_1, H_2\}$ . As shown in the following proposition, the maximum can occur in one of three locations.

**Proposition 7.** Consider two source populations with distinct allele frequencies,  $\underline{p}_1 \neq \underline{p}_2$ . As a function of  $\gamma_1$ ,  $H_{\text{adm}}$  is maximized at  $\gamma_1 = \gamma_1^*$ , where  $\gamma_1^*$  takes one of three forms.

(i) If  $H_1 < C_{12}$  and  $H_2 < C_{12}$ , then  $\gamma_1^* \in (0, 1)$  satisfies

$$\gamma_1^* = \frac{C_{12} - H_2}{2(C_{12} - H_S)} = \frac{1}{2} + \frac{H_1 - H_2}{8(H_T - H_S)}, \quad (15)$$

and  $H_{\text{adm}}$  has maximum equal to

$$H_{\text{adm}}(\gamma_1^*) = \frac{C_{12}^2 - H_1 H_2}{2(C_{12} - H_S)} = H_T + \frac{(H_1 - H_2)^2}{16(H_T - H_S)}. \quad (16)$$

(ii) If  $H_1 < C_{12}$  and  $H_2 \geq C_{12}$ , then  $\gamma_1^* = 0$  and  $H_{\text{adm}}$  has maximum  $H_2$ .

(iii) If  $H_1 \geq C_{12}$  and  $H_2 < C_{12}$ , then  $\gamma_1^* = 1$  and  $H_{\text{adm}}$  has maximum  $H_1$ .

An elementary proof appears in Appendix 2. We can see that these three cases capture all possible values of  $(H_1, H_2, C_{12})$ . By the Cauchy-Schwarz inequality,  $(1 - C_{12})^2 \leq (1 - H_1)(1 - H_2)$ , with equality requiring  $\underline{p}_1 = \underline{p}_2$ . Hence, with  $\underline{p}_1 \neq \underline{p}_2$  assumed, either  $1 - C_{12} < 1 - H_1$  and  $1 - C_{12} \geq 1 - H_2$  (case ii),  $1 - C_{12} < 1 - H_2$  and  $1 - C_{12} \geq 1 - H_1$  (case iii), or both  $1 - C_{12} < 1 - H_1$  and  $1 - C_{12} < 1 - H_2$  (case i).

Note that the locations specified in Proposition 7 accord with those in Theorem 5 and Corollary 6. For  $K = 2$ ,

$$A = \begin{pmatrix} H_1 & C_{12} \\ C_{12} & H_2 \end{pmatrix}.$$

The result of Theorem 5 gives  $\underline{\gamma}^* = \frac{A^{-1}\mathbf{1}}{\mathbf{1}'A^{-1}\mathbf{1}} = \left( \frac{C_{12}-H_2}{2(C_{12}-H_S)}, \frac{C_{12}-H_1}{2(C_{12}-H_S)} \right)$ . The locations in Corollary 6 are  $\gamma_1^* = \frac{A_1^{-1}}{A_1^{-1}} = 1$  and  $\gamma_2^* = 0$ , and  $\gamma_1^* = 0$  and  $\gamma_2^* = \frac{A_2^{-1}}{A_2^{-1}} = 1$ .

In accord with the observation in Figure 1 that the maximal  $H_{\text{adm}}$  lies at a value of  $\gamma_1 \geq \frac{1}{2}$  in an example with  $H_1 \geq H_2$ , the proposition demonstrates  $\gamma_1^* \geq \frac{1}{2}$  if and only if  $H_1 \geq H_2$ . To obtain this result, suppose  $H_1 \geq H_2$ . If case (i) applies, then because  $H_T > H_S$ ,  $\gamma_1^* \geq \frac{1}{2}$ . Case (ii) cannot apply because  $H_1 < C_{12}$ ,  $H_2 \geq C_{12}$ , and  $H_1 \geq H_2$  cannot hold simultaneously. In case (iii),  $\gamma_1^* = 1 \geq \frac{1}{2}$ . For the reverse direction, if  $H_1 < H_2$  and case (i) or case (ii) applies, then  $\gamma_1^* < \frac{1}{2}$ . Case (iii) cannot apply because  $H_1 \geq C_{12}$ ,  $H_2 < C_{12}$ , and  $H_1 < H_2$  cannot hold simultaneously.

We also have  $H_{\text{adm}}(\gamma_1^*) \geq H_T$  in all three cases, and  $H_{\text{adm}}(\gamma_1^*) = H_T$  if  $H_1 = H_2$ . To obtain this result, we see that  $H_{\text{adm}}(\gamma_1^*) \geq H_T$  in case (i). In case (ii),  $H_2 > H_T = \frac{H_1+H_2+2C_{12}}{4}$  because  $H_2 > H_1$  and  $H_2 \geq C_{12}$ . In case (iii),  $H_1 > H_T$  because  $H_1 > H_2$  and  $H_1 \geq C_{12}$ . Note that if  $H_1 = H_2$ , then case (i) applies, producing  $H_{\text{adm}}(\gamma_1^*) = H_T$ .

We can succinctly describe the region where  $\gamma_1^*$  lies interior to  $(0, 1)$ .

**Corollary 8.** Consider two source populations with distinct allele frequencies,  $\underline{p}_1 \neq \underline{p}_2$ .  $\gamma_1^*$  lies in  $(0, 1)$  if and only if the following inequality holds:

$$F_{12} > \frac{|H_1 - H_2|}{2(H_1 + H_2) + |H_1 - H_2|}. \quad (17)$$

This corollary is proven in Appendix 2. Note that if  $H_1 + H_2$  is fixed, then the right-hand side of eq. 17 increases with  $|H_1 - H_2|$ , from a minimum of 0 when  $H_1 = H_2$  to a supremum of  $\frac{1}{3}$  as  $|H_1 - H_2|$  approaches  $H_1 + H_2$ . Thus, in accord with the observation in Section 3.1 that  $H_{\text{adm}} > H$  for all nontrivial admixtures of equal-heterozygosity source populations, the maximal  $H_{\text{adm}}$  exceeds  $\max\{H_1, H_2\}$  over a broader range of  $F_{12}$  values if  $|H_1 - H_2|$  is small rather than large. Moreover, if  $F_{12} > \frac{1}{3}$ , then eq. 17 necessarily holds. Hence, irrespective of  $H_1$  and  $H_2$ , if the source populations are distant enough that  $F_{12} > \frac{1}{3}$ , then the maximal heterozygosity exceeds the heterozygosities of the source populations.

## 4.2 Special case of $J = 2$ alleles

In the case with  $K = 2$  source populations in which the locus has only  $J = 2$  allelic types, it is possible to make further simplifications, as the results can be stated in terms of frequencies of one specific allele. We substitute  $p_{12} = 1 - p_{11}$  and  $p_{22} = 1 - p_{21}$ .

**Proposition 9.** Consider two source populations with distinct allele frequencies,  $\underline{p}_1 \neq \underline{p}_2$ . For a biallelic locus,  $H_{\text{adm}}$  is maximized at  $\gamma_1 = \gamma_1^*$ , where  $\gamma_1^*$  takes one of three forms.

(i) If  $p_{11} > \frac{1}{2} > p_{21}$  or  $p_{21} > \frac{1}{2} > p_{11}$ , then  $\gamma_1^* \in (0, 1)$  satisfies

$$\gamma_1^* = \frac{1 - 2p_{21}}{2(p_{11} - p_{21})}, \quad (18)$$

and  $H_{\text{adm}}$  has maximum equal to

$$H_{\text{adm}}(\gamma_1^*) = \frac{1}{2}. \quad (19)$$

(ii) If  $\frac{1}{2} \geq p_{21} > p_{11}$  or  $p_{11} > p_{21} \geq \frac{1}{2}$ , then  $\gamma_1^* = 0$  and  $H_{\text{adm}}$  has maximum  $H_2$ .

(iii) If  $\frac{1}{2} \geq p_{11} > p_{21}$  or  $p_{21} > p_{11} \geq \frac{1}{2}$ , then  $\gamma_1^* = 1$  and  $H_{\text{adm}}$  has maximum  $H_1$ .

*Proof.* We apply Proposition 7 with  $J = 2$ . Substituting  $p_{12} = 1 - p_{11}$  and  $p_{22} = 1 - p_{21}$  in eqs. 15 and 16, we obtain  $C_{12} - H_2 = (p_{11} - p_{21})(1 - 2p_{21})$ ,  $C_{12} - H_1 = (p_{21} - p_{11})(1 - 2p_{11})$ ,  $C_{12} - H_S = (p_{11} - p_{21})^2$ , and  $C_{12}^2 - H_1H_2 = (p_{11} - p_{21})^2$ . Thus, because  $p_{11} = p_{21}$  is not permitted, the quantities in eqs. 15 and 16 reduce to those of eqs. 18 and 19, respectively.

To complete the application of Proposition 7 to  $K = 2$ , note that case (i) of Proposition 7 occurs when  $(p_{11} - p_{21})(1 - 2p_{21}) > 0$  and  $(p_{21} - p_{11})(1 - 2p_{11}) > 0$ . The first of this pair of inequalities requires both  $p_{11} - p_{21} > 0$  and  $1 - 2p_{21} > 0$ , so that  $p_{11} > p_{21}$  and  $\frac{1}{2} > p_{21}$ , or both  $p_{11} - p_{21} < 0$  and  $1 - 2p_{21} < 0$ , so that  $p_{11} < p_{21}$  and  $\frac{1}{2} < p_{21}$ . The second inequality requires both  $p_{21} - p_{11} > 0$  and  $1 - 2p_{11} > 0$ , so that  $p_{21} > p_{11}$  and  $\frac{1}{2} > p_{11}$ , or both  $p_{21} - p_{11} < 0$  and  $1 - 2p_{11} < 0$ , so that  $p_{21} < p_{11}$  and  $\frac{1}{2} < p_{11}$ . Thus, the conditions of case (i) of Proposition 7 obtain if and only if  $p_{11} > \frac{1}{2} > p_{21}$  or  $p_{21} > \frac{1}{2} > p_{11}$ .

Similarly, using the expressions for  $H_1$ ,  $H_2$ , and  $C_{12}$  when  $K = 2$ , the conditions of case (ii) of Proposition 7 are equivalent to  $\frac{1}{2} \geq p_{21} > p_{11}$  or  $p_{11} > p_{21} \geq \frac{1}{2}$ . The conditions of case (iii) are equivalent to  $\frac{1}{2} \geq p_{11} > p_{21}$  or  $p_{21} > p_{11} \geq \frac{1}{2}$ .  $\square$

The unit square representing the possible values of the location of the maximum appears in Figure 2. The square has six nonoverlapping regions: in Proposition 9, each of the three

cases corresponds to two disjoint subsets of  $[0, 1]^2$ . A smooth gradient exists for the regions in case (i). However, an abrupt transition occurs at the line  $p_{21} = p_{11}$  between the case-(ii) regions where  $\gamma_1^* = 0$  and the case-(iii) regions where  $\gamma_1^* = 1$ . Note that the  $p_{21} = p_{11}$  line is disallowed, as it corresponds to the two populations having the same allele frequencies.

## 5 Simulations

We illustrate a number of properties of  $H_{\text{adm}}$  by simulating population sets for different values of  $K$  and  $J$ . Given a value of  $K$ , we generated allele frequency vectors for the  $K$  source populations from independent and identically distributed symmetric multivariate  $J$ -dimensional Dirichlet distributions with a common concentration parameter  $\alpha = 1$ . This distribution corresponds to a uniform distribution on the simplex  $\Delta^{J-1}$ . Mathematical results concerning  $H_{\text{adm}}$  under the Dirichlet distribution on allele frequencies appear in Appendix 3.

First, for  $K = 2$  and  $K = 3$ , we assessed the probability that the maximal  $H_{\text{adm}}$  over possible admixture vectors  $\underline{\gamma}$  occurs interior to the simplex  $\Delta^{K-1}$ , rather than on its boundary. This computation gives the probability that the heterozygosity-maximizing admixture vector contains nonzero contributions from all  $K$  source populations. We considered  $2 \leq J \leq 30$  for  $K = 2$  and  $3 \leq J \leq 30$  for  $K = 3$ , recalling the condition  $J \geq K$  for the  $K$  allele frequency vectors to be linearly independent.

For each  $(K, J)$ , we ran 10,000 simulation replicates. In each replicate, to determine the location of the maximum, we applied Theorem 5 and Corollary 6 to identify the locations specified for each choice  $\mathcal{S}$  of the nonempty subset of the  $K$  populations with nonzero allele frequencies. Among these  $2^{K-1}$  locations, excluding those outside the simplex  $\Delta^{K-1}$ , we identified the point with the largest  $H_{\text{adm}}$ . Note that in each replicate, we observed that the  $\underline{1}'(P_{\mathcal{S}}'P_{\mathcal{S}})^{-1}\underline{1} \neq 1$  condition of Corollary 6 was satisfied for each  $\mathcal{S}$ .

Figure 3 finds that, for both  $K = 2$  and  $K = 3$ , the maximum of  $H_{\text{adm}}$  is increasingly likely to be in the interior of the simplex as the number of distinct alleles,  $J$ , increases. For  $K = 3$ , we also observe that the probability that  $H_{\text{adm}}$  is maximized on an edge, corresponding to nonzero contributions from two of the three source populations, exceeds the probability that it is maximized at a vertex, with only one contributing source population.

Next, we assessed the probability  $\mathbb{P}[H_{\text{adm}} > \max\{H_1, \dots, H_K\}]$  in a scenario in which both

the allele frequency vectors  $\underline{p}_k$  and the admixture fractions  $\underline{\gamma}$  were chosen from independent Dirichlet distributions. We simulated the  $\underline{p}_k$  as before, additionally simulating  $\underline{\gamma}$  from a  $K$ -dimensional symmetric Dirichlet- $(1, 1, \dots, 1)$  distribution. For each  $(K, J)$  with  $K = 2, 3, 4, 5$  and  $J = 2, 3, \dots, 30$ , we simulated 50,000 replicate populations. Note that here, unlike in Section 2.4, we impose no restrictions on linear combinations of allele frequency vectors from the source populations, so that it is not necessarily true that  $J \geq K$ .

The fraction of replicates with  $\mathbb{P}[H_{\text{adm}} > \max\{H_1, \dots, H_K\}]$  appears in Figure 4. We see that this fraction increases as  $K$  increases, indicating that for an admixture involving more populations, the probability is larger that the admixed population has greater heterozygosity than all source populations. This probability also increases with increasing  $J$ .

For the special case of  $K = 2$  and  $J = 2$ , Proposition 15 in Appendix 3 obtains the probability analytically,  $\mathbb{P}[H_{\text{adm}} > \max\{H_1, H_2\}] = 1 - \log 2 \approx 0.307$ . In accord with this result, the  $K = 2$  curve in Figure 4 begins near  $(2, 0.307)$ .

Figure 5 provides further detail on the behavior of  $H_{\text{adm}}$  in the  $K = 2$  case by graphing  $H_{\text{adm}}$  versus  $\gamma_1$  for 10 simulation replicates chosen at random for each of three values of  $J$ . The figure illustrates that  $H_{\text{adm}}$  is a concave-down quadratic polynomial in  $\gamma_1$ , as in eq. 14. Averaging across replicates, by examining the panels of the figure from left to right, we can also observe that  $\mathbb{E}[H_{\text{adm}}]$  increases as a function of  $J$ , as in Corollary 14 of Appendix 3. For  $J = 2$ , as in Proposition 9, the possible values of  $H_{\text{adm}}$  at the maximum are  $H_1, H_2$ , and  $\frac{1}{2}$ .

## 6 Application to data

We illustrate the mathematical results using data from human populations, following Boca & Rosenberg (2011) in considering data from Wang *et al.* (2008) on 678 microsatellite loci typed in 160 Europeans, 463 Native Americans, 123 Africans, and 249 individuals from admixed Mestizo populations. To represent admixed Mestizo populations under our model, we used sample allele frequencies for the Europeans and Native Americans as source populations in the  $K = 2$  case, also including sample allele frequencies for the Africans for  $K = 3$ . As in Boca & Rosenberg (2011), we treated allele frequencies, heterozygosities, and  $F_{ST}$  values computed from the data as parametric values rather than estimates.



## 6.1 $K = 2$ source populations

We selected 20 loci at random from Wang *et al.* (2008) for illustration, choosing the same loci as in our study on  $F_{ST}$  and admixture (Boca & Rosenberg, 2011). Treating  $\gamma_1$  as the fraction of European ancestry and  $1 - \gamma_1$  as the fraction of Native American ancestry in an admixed population, for each locus, the plot for  $H_{adm}$  versus  $\gamma_1$  appears in Figure 6. Following Proposition 3, the minimal value of  $H_{adm}$  lies either at  $\gamma_1 = 0$  or at  $\gamma_1 = 1$  for all the loci. For 12 of the 20 loci, the maximum of  $H_{adm}$  lies in the interior of the unit interval for  $\gamma_1$ . Seven loci have the maximum at  $\gamma_1 = 1$ , representing membership in the more heterozygous European population, and only one locus has the maximum at  $\gamma_1 = 0$ , representing membership in the less heterozygous Native American population.

Examining all 678 loci, 52% have the maximum in the interior, 39% at  $\gamma_1 = 1$ , and 8% at  $\gamma_1 = 0$ . That more loci have the maximum at  $\gamma_1 = 1$  than at  $\gamma_1 = 0$  is expected from the fact that European populations generally tend to have greater heterozygosity than Native American populations (e.g. Pemberton *et al.*, 2013).

The Dirichlet model in Corollary 14 in Appendix 3 and Figures 3 and 5 predicts a dependence of the location of the maximum on the number of distinct alleles of a locus, with the probability that the maximum lies in the interior increasing with the number of distinct alleles. The data produce a trend in the same direction as this prediction. The mean numbers of distinct alleles are 9.44, 10.41, and 10.74, for the loci with  $\gamma_1^*$  at 0, 1, and in  $(0, 1)$ , respectively (one-way ANOVA,  $P = 0.01$ ,  $F$  test, 2 df). The mean number of distinct alleles for the loci with the maximum on either boundary is 10.24, smaller than the mean of 10.74 for those with the mean in the interior ( $P = 0.04$ , two-tailed  $t$  test).

## 6.2 Comparison of predicted $H_{adm}$ to observed $H_{adm}$

Next, we compare predicted and observed  $H_{adm}$  values for the 678 loci for the admixed Mestizo population. In this approach, we used estimated locus-wise values of  $\gamma_1$  in the Mestizo population together with locus-wise heterozygosities in the European and Native American populations to “predict” locus-wise Mestizo heterozygosities. The prediction is compared to the observed heterozygosity value to examine if our formulas for the heterozygosity of an admixed population are reflected in actual heterozygosities in an admixed group.

This computation follows a similar computation of Boca & Rosenberg (2011). The estimated admixture fractions, computed for the same data, are taken from Schroeder *et al.* (2009), who obtained them by a maximum likelihood approach (Millar, 1987) that does not take into account the source population heterozygosities. Using these estimates, locus-wise heterozygosity estimates in the source populations, and locus-wise  $F_{ST}$  values calculated from the allele frequencies in the source populations, we predicted  $H_{\text{adm}}$  with eq. 13.

The predicted and observed  $H_{\text{adm}}$  values for individual loci are compared in Figure 7. In general, the observation closely matches the prediction (Figure 7A), with the correlation between the observed and predicted  $H_{\text{adm}}$  values equaling 0.978 (Figure 7B). For 56% of the 678 loci, the prediction provides an underestimate of the observed value.

### 6.3 $K = 3$ source populations

We now consider the European, Native American, and African populations as the source populations, using  $\gamma_1$  for the proportion of European ancestry,  $\gamma_2$  for Native American ancestry, and  $\gamma_3$  for African ancestry. We select 3 loci for illustration, choosing the same ones as in a similar analysis of Boca & Rosenberg (2011).

Plots for  $H_{\text{adm}}$  over the unit simplex for  $(\gamma_1, \gamma_2, \gamma_3)$  appear in Figure 8. Each plot depicts  $H_{\text{adm}}$  as a function of  $(\gamma_1, \gamma_2, \gamma_3)$  for a specific locus. The three panels show the possible locations of the maximal value of  $H_{\text{adm}}$ : in the first panel, the maximum lies in the interior of the simplex; in the second panel, at a vertex, and in the third panel, on an edge.

Considering all 678 loci, 14% have the maximum in the interior of the region, with  $\gamma_1 > 0$ ,  $\gamma_2 > 0$ , and  $\gamma_3 > 0$ . The fractions with the maximum on an edge are 19% for a maximum on the edge with  $\gamma_1 = 0$ , 26% on the  $\gamma_2 = 0$  edge, and 5% on the  $\gamma_3 = 0$  edge. The fractions with the maximum at a vertex are 27% for the vertex  $(0, 0, 1)$ , 2% for  $(0, 1, 0)$ , and 5% for  $(1, 0, 0)$ . The observations that  $(0, 0, 1)$  is the vertex with the largest number of maxima and  $(1, 0, 1)$  is the edge with the most maxima accord with the fact that African populations have generally higher heterozygosity than European populations, which in turn have higher heterozygosity than Native American populations (e.g. Pemberton *et al.*, 2013).

## 7 Discussion

We have considered the heterozygosity  $H_{\text{adm}}$  of an admixed population in terms of the admixture fractions of the source populations, and their heterozygosities and  $F_{ST}$  values at a locus. We have derived formulas describing  $H_{\text{adm}}$  in relation to these quantities (eqs. 8-10). In particular, we showed that  $H_{\text{adm}}$  is minimized over the set of possible admixture coefficient vectors when the admixed population consists of only one of the source populations (Proposition 3); hence, an admixed population is at least as heterozygous as the least heterozygous source population. The maximal  $H_{\text{adm}}$  has a more complex characterization, as it has the possibility of being either greater than or equal to the heterozygosity of the most heterozygous source population (Proposition 4).

In studying the possible locations of the maximal  $H_{\text{adm}}$  for a fixed set of source populations, we found that the maximum can lie either in the interior of the region describing the allowable values of the admixture fractions—in which case all source populations contribute to the admixed population—or on the boundary, where one or more source populations does not contribute to the admixed population (Propositions 4-6, Figures 1-3). Simulations under a Dirichlet model for allele frequencies suggest that the maximal value of  $H_{\text{adm}}$  lies with increasing frequency in the interior of the allowable region as  $K$  and  $J$  increase (Figure 4).

For  $K = 2$  source populations, we obtained further results, in particular showing that  $H_{\text{adm}}$  is a concave-down quadratic polynomial in the admixture coefficient  $\gamma_1$  (eqs. 12-14). We obtained an analytical expression for the maximal heterozygosity of an admixture of a specific pair of source populations in terms of  $H_1, H_2$ , and the  $F_{ST}$  value between the two populations (Proposition 7). For fixed values of  $H_1, H_2$ , and the admixture fraction  $\gamma_1$ ,  $H_{\text{adm}}$  is increasing as a function of  $F_{ST}$  (eq. 13, Figure 1). If  $H_1 > H_2$ , then the admixture fraction in source population 1 that maximizes  $H_{\text{adm}}$  is greater than  $\frac{1}{2}$  (Proposition 7), meaning that at the maximal heterozygosity of the admixed population, the contribution of the more heterozygous source population exceeds that of the less heterozygous one. Interestingly, for the  $K = 2$  case with  $J = 2$  allelic types, if the location of the maximal value lies in  $(0, 1)$ , then the heterozygosity at the maximum is always  $\frac{1}{2}$  (Proposition 9 and Figure 5): irrespective of the allele frequencies of the source populations, a linear combination  $(\gamma_1, \gamma_2)$  always exists so that the admixed population has allele frequencies of  $\frac{1}{2}$  for both alleles.

For  $K = 2$  source populations, a key result is that the maximal value of  $H_{\text{adm}}$  exceeds the larger of the two source population heterozygosities if and only if  $F_{ST}$  exceeds a bound defined by those heterozygosities (Corollary 8). Thus, with all other quantities equal, combining source populations that are more divergent rather than less divergent is more likely to lead to an admixed population with heterozygosity exceeding those of the source populations.

In human data, we observed that for heterozygosities and  $F_{ST}$  values for putative sources of Mestizo populations, the maximal  $H_{\text{adm}}$  was more likely to be in the interior of the unit simplex or on an edge rather than at a vertex (Figures 6 and 8). This result indicates that the heterozygosities and  $F_{ST}$  values of these populations lie in a parameter range for which admixed populations are frequently more heterozygous than all their source populations. Examining the heterozygosities of 267 worldwide populations in Table S20 of Pemberton *et al.* (2013), the 13 Mestizo populations all have heterozygosities exceeding all 29 Native American populations, and 4 have heterozygosities exceeding all 8 European populations. Interestingly, the top 10 most heterozygous populations among the 267 include all five admixed populations involving a source population from the high-heterozygosity region of Africa: a Cape Mixed Ancestry population from South Africa, and four African-American populations. Thus, our mathematical results predicting that admixed populations often exceed all their source populations in heterozygosity are reflected in admixed human groups.

For  $K = 2$ , our model successfully predicted the heterozygosities in an admixed population from the source population heterozygosities, the  $F_{ST}$  between the source populations, and the estimated admixture coefficient  $\hat{\gamma}_1$  for one of the source populations (Figure 7). Because  $H_{\text{adm}}$  is not necessarily monotonic in the admixture fraction  $\gamma_1$ , however, the reverse problem of using  $H_{\text{adm}}$  to estimate  $\gamma_1$  is problematic—unlike for the monotonically varying  $F_{ST}$  between an admixed population and one of the source populations (Boca & Rosenberg, 2011, Theorem 3). Given a value of  $H_{\text{adm}}$ , source population heterozygosities  $H_1$  and  $H_2$ , and  $F_{ST}$  between the source populations, two solutions to eq. 13 might exist for  $\gamma_1$ —so that although  $H_{\text{adm}}$  can be predicted from  $\gamma_1$ , it is inadvisable to proceed in the reverse direction to estimate the admixture coefficient  $\gamma_1$  from the heterozygosity of an admixed population.

Our approach has followed the study of  $F_{ST}$  and admixture from Boca & Rosenberg (2011), and it shares similar limitations. For example, the model assumes source population

allele frequencies are known rather than estimated, and it considers only population-level rather than individual-level admixture. It also relies on patterns of variation from a single time point and does not incorporate mechanistic evolutionary processes. Despite these limitations, the observed  $H_{\text{adm}}$  values and those predicted under our model are strongly correlated (Figure 7B), indicating that the model captures key population features relevant to the relationship between admixture and heterozygosity. Thus, the empirical results suggest that assessing this relationship in the mathematical formulations we have presented can be useful for understanding the genetics of admixed populations.

**Acknowledgments.** Support was provided by NIH grant HG005855 and NSF grant BCS-1515127.

## Appendix 1. Proofs for arbitrary $K$ : Theorem 5 and Corollary 6

For the proof of Theorem 5, we first show (i) that  $P'P$  and  $A$  are both invertible under the conditions stated in the theorem, and that:

$$\frac{1}{\underline{1}'A^{-1}\underline{1}} = 1 - \frac{1}{\underline{1}'(P'P)^{-1}\underline{1}}.$$

We then (ii) use constrained optimization via Lagrange multipliers to obtain the maximum of  $\underline{\gamma}'A\underline{\gamma}$  subject to  $\underline{1}'\underline{\gamma} = 1$ . This step consists of the first derivative test to find a stationary point, coupled with the second derivative test, in Lemma 10, to show that the stationary point defines a local maximum. Finally, we (iii) show that this means that the overall maximum is either at the local maximum  $\underline{\gamma}^*$  as described in the statement of the theorem or on the boundary of the set  $\{\underline{\gamma} : \underline{1}'\underline{\gamma} = 1 \text{ and } \underline{\gamma} \in \Delta^{K-1}\}$ .

*Proof of Theorem 5 (i)* Because  $P$  is a  $J \times K$  matrix with column rank  $K$ , the  $K \times K$  matrix  $P'P$  is positive definite. As a positive definite matrix,  $P'P$  is invertible and  $(P'P)^{-1}$  is also positive definite (Graybill, 1976, pp. 21-22).

To show that  $A = \underline{1}\underline{1}' - P'P$  is invertible, we use the Sherman-Morrison formula for the inverse of a rank-one update of an invertible matrix (Horn & Johnson, 2012, pp. 18-19). This formula states that for an invertible square  $n \times n$  matrix  $X$  and  $n \times 1$  column vectors  $\underline{y}$  and  $\underline{z}$ ,  $X + \underline{y}\underline{z}'$  is invertible if and only if  $1 + \underline{z}'X^{-1}\underline{y} \neq 0$ , with:

$$(X + \underline{y}\underline{z}')^{-1} = X^{-1} - \frac{X^{-1}\underline{y}\underline{z}'X^{-1}}{1 + \underline{z}'X^{-1}\underline{y}}.$$

Because we assumed  $\underline{1}'(P'P)^{-1}\underline{1} \neq 1$ , the Sherman-Morrison formula applies with  $-(P'P)$  in the role of  $X$ , and  $K \times 1$  column vectors  $\underline{1}$  in the role of  $\underline{y}$  and  $\underline{z}$ .  $A$  has inverse:

$$A^{-1} = \frac{(P'P)^{-1}\underline{1}\underline{1}'(P'P)^{-1}}{\underline{1}'(P'P)^{-1}\underline{1} - 1} - (P'P)^{-1}. \quad (20)$$

Left-multiplying by  $\underline{1}'$  and right-multiplying by  $\underline{1}$ , we obtain

$$\frac{1}{\underline{1}'A^{-1}\underline{1}} = 1 - \frac{1}{\underline{1}'(P'P)^{-1}\underline{1}}.$$

Because  $(P'P)^{-1}$  is positive definite,  $\underline{1}'(P'P)^{-1}\underline{1} > 0$  by definition, and because  $\underline{1}'(P'P)^{-1}\underline{1} \neq 1$  by assumption, we conclude that  $\frac{1}{\underline{1}'A^{-1}\underline{1}}$  is always defined.

(ii) To maximize  $\underline{\gamma}'A\underline{\gamma}$  subject to  $\underline{1}'\underline{\gamma} = 1$ , we use Lagrange multipliers. Let  $f(\underline{\gamma}) = \underline{\gamma}'A\underline{\gamma}$ , and let  $g(\underline{\gamma}) = \underline{1}'\underline{\gamma}$ . The Lagrange function is defined as:

$$\Lambda(\underline{\gamma}, \lambda) = f(\underline{\gamma}) + \lambda[g(\underline{\gamma}) - 1].$$

Denoting by  $\underline{0}$  is a column vector of length  $K$ , we solve a system of equations for  $\underline{\gamma}$  and  $\lambda$ ,

$$\left( \frac{\delta\Lambda(\underline{\gamma}, \lambda)}{\delta\underline{\gamma}}, \frac{\delta\Lambda(\underline{\gamma}, \lambda)}{\delta\lambda} \right) = (\underline{0}, 0). \quad (21)$$

Eq. 21 includes  $K$  equations  $\delta\Lambda(\underline{\gamma}, \lambda)/\delta\gamma_k = 0$  for  $1 \leq k \leq K$ .

$A$  is symmetric, so we have

$$\begin{aligned} \frac{\delta f(\underline{\gamma})}{\delta\underline{\gamma}} &= \frac{\delta(\underline{\gamma}'A\underline{\gamma})}{\delta\underline{\gamma}} = (A + A')\underline{\gamma} = 2A\underline{\gamma} \\ \frac{\delta g(\underline{\gamma})}{\delta\underline{\gamma}} &= \underline{1}. \end{aligned}$$

For the derivatives of the Lagrange function, we have:

$$\left( \frac{\delta\Lambda(\underline{\gamma}, \lambda)}{\delta\underline{\gamma}}, \frac{\delta\Lambda(\underline{\gamma}, \lambda)}{\delta\lambda} \right) = (2A\underline{\gamma} + \lambda\underline{1}, \underline{1}'\underline{\gamma} - 1).$$

Setting the derivatives with respect to  $\underline{\gamma}$  to  $\underline{0}$  leads to:

$$(\underline{\gamma}, \lambda) = \left( -\frac{\lambda}{2}A^{-1}\underline{1}, -\frac{2}{\underline{1}'A^{-1}\underline{1}} \right).$$

Hence, the solution for  $\underline{\gamma}$  is:

$$\underline{\gamma}^* = \frac{A^{-1}\underline{1}}{\underline{1}'A^{-1}\underline{1}}.$$

Because  $\underline{\gamma}'A\underline{\gamma}$  is a differentiable function of  $\underline{\gamma}$ , its maximum on  $\Delta^{K-1}$  can occur either on the boundary or at a critical point. The following lemma shows that the critical point  $\underline{\gamma}^* = \frac{A^{-1}\underline{1}}{\underline{1}'A^{-1}\underline{1}}$  is a local maximum.

**Lemma 10.** The critical point  $\underline{\gamma}^* = \frac{A^{-1}\underline{1}}{\underline{1}'A^{-1}\underline{1}}$  is a local maximum of  $H_{\text{adm}}$  seen as a function of  $\underline{\gamma}$  on  $\Delta^{K-1}$ , under the conditions stated in Theorem 5.

*Proof.* To show that  $\underline{\gamma}^*$  is a local maximum, we use the second derivative test for constrained optimization (e.g. Magnus & Neudecker, 2007, p. 155). This test considers the bordered Hessian matrix, representing the matrix of second derivatives of the Lagrange function  $\Lambda$  with respect to  $\lambda$  and the components of  $\underline{\gamma}$ :

$$F = \begin{pmatrix} \frac{\delta^2 \Lambda}{\delta \lambda^2} & \left( \frac{\delta^2 \Lambda}{\delta \gamma \delta \lambda} \right)' \\ \frac{\delta^2 \Lambda}{\delta \gamma \delta \lambda} & \frac{\delta^2 \Lambda}{\delta \gamma^2} \end{pmatrix} = \begin{pmatrix} 0 & \left( \frac{\delta g}{\delta \gamma} \right)' \\ \frac{\delta g}{\delta \gamma} & \frac{\delta^2 \Lambda}{\delta \gamma^2} \end{pmatrix} = \begin{pmatrix} 0 & \underline{1}' \\ \underline{1} & 2A \end{pmatrix}.$$

We must consider the principal minors—determinants of matrices in the upper-left corner—of  $F$ . We denote the upper-left corner matrix of size  $r \times r$  of  $F$  by  $F_r$ , for  $r = 2, 3, \dots, K$ . The principal minors are the  $\det(F_r)$ . Using the definition of  $A$  from eq. 11, we obtain

$$F_r = \begin{pmatrix} 0 & 1 & 1 & \dots & 1 \\ 1 & 2H_1 & 2C_{12} & \dots & 2C_{1r} \\ 1 & 2C_{12} & 2H_2 & \dots & 2C_{2r} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & 2C_{1r} & 2C_{2r} & \dots & 2H_r \end{pmatrix}.$$

A sufficient condition for the critical point to be a local maximum is for  $(-1)^r \det(F_r) > 0$  for each  $r$  (Magnus & Neudecker, 2007, p. 155). We now show that this condition is satisfied.

Using the fact that multiplying a row or column of a matrix by a scalar multiplies the determinant by that scalar, we multiply rows 2 through  $r + 1$  by  $-1$  and get

$$\det(F_r) = \det \begin{pmatrix} 0 & 1 & 1 & \dots & 1 \\ 1 & 2H_1 & 2C_{12} & \dots & 2C_{1r} \\ 1 & 2C_{12} & 2H_2 & \dots & 2C_{2r} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & 2C_{1r} & 2C_{2r} & \dots & 2H_r \end{pmatrix} = (-1)^r \det \begin{pmatrix} 0 & 1 & 1 & \dots & 1 \\ -1 & -2H_1 & -2C_{12} & \dots & -2C_{1r} \\ -1 & -2C_{12} & -2H_2 & \dots & -2C_{2r} \\ \vdots & \vdots & \vdots & \vdots & \dots \\ -1 & -2C_{1r} & -2C_{2r} & \dots & -2H_r \end{pmatrix}$$

Using the fact that adding a multiple of a row or column to another row does not change the determinant, we add  $-2$  times the first column to each of the remaining columns. We

also multiply the first column by  $-1$ . We then have

$$(-1)^r \det(F_r) = (-1)^{2r+1} \det \begin{pmatrix} 0 & \underline{1}_r' \\ \underline{1}_r & 2M_r \end{pmatrix} = - \det \begin{pmatrix} 0 & \underline{1}_r' \\ \underline{1}_r & 2M_r \end{pmatrix}, \quad (22)$$

where  $M_r$  is the  $r \times r$  matrix consisting of the upper-left corner of matrix  $P'P$ , and  $\underline{1}_r$  is the column vector of length  $r$  consisting of 1s.

We now apply a result for the determinant of partitioned matrices (Graybill, 1976, pp. 19-20). If  $W$  is invertible, then

$$\det \begin{pmatrix} X & Y \\ Z & W \end{pmatrix} = \det(W) \det(X - YW^{-1}Z).$$

Applying this result to eq. 22, we obtain

$$\begin{aligned} (-1)^r \det(F_r) &= - \det(2M_r) \det(-\underline{1}_r'(2M_r)^{-1}\underline{1}_r) \\ &= -[2^r \det(M_r)] \left[ \left( -\frac{1}{2} \right) \underline{1}_r' M_r^{-1} \underline{1}_r \right] \\ &= 2^{r-1} \det(M_r) (\underline{1}_r' M_r^{-1} \underline{1}_r). \end{aligned}$$

Because  $P'P$  is positive definite,  $M_r$  is also positive definite. To demonstrate this result, note that because  $\underline{x}'P'P\underline{x} > 0$  for each nonzero column vector  $\underline{x}$ ,  $\underline{x}'P'P\underline{x} > 0$  for each nonzero  $\underline{x}$  with  $x_k = 0$  for  $k > r$ . Because  $M_r$  is positive definite,  $\det(M_r) > 0$  and  $M_r^{-1}$  is also positive definite, leading to  $\underline{1}_r' M_r^{-1} \underline{1}_r > 0$ . We conclude

$$(-1)^r \det(F_r) > 0,$$

so that the critical point is the location of a local maximum.  $\square$

*Concluding the proof of Theorem 5.* Returning to part (iii) of the proof, following Lemma 10, if  $\underline{\gamma}^* = \frac{A^{-1}\underline{1}}{\underline{1}'A^{-1}\underline{1}}$  is interior to the simplex  $\Delta^{K-1}$ , then  $H_{\text{adm}}$  is maximal at  $\underline{\gamma} = \underline{\gamma}^*$ , with maximum  $H(\underline{\gamma}) = \frac{1}{\underline{1}'A^{-1}\underline{1}}$ . This value is the reciprocal of the sum of the elements of  $A^{-1}$ . If  $\underline{\gamma}^*$  is not interior to  $\Delta^{K-1}$ , then the maximum lies on the boundary of  $\Delta^{K-1}$ .

Finally, we note that  $\underline{\gamma}^* = \frac{(P'P)^{-1}\underline{1}}{\underline{1}'(P'P)^{-1}\underline{1}}$  by using eq. 20.  $\square$

*Proof of Corollary 6.* In Theorem 5, the maximum of  $H_{\text{adm}}$  occurs either in the interior of the simplex  $\Delta^{K-1}$  or on its boundary,  $\{\underline{\gamma} : \underline{1}'\underline{\gamma} = 1 \text{ and } \underline{\gamma} \in \Delta^{K-1}\}$ .



The boundary of the simplex is the union of  $K$  faces, which are themselves  $(K - 2)$ -simplices. If the maximum lies on the boundary of  $\Delta^{K-1}$ , then without loss of generality, we can permute the labels of the source populations so that  $\gamma_K = 0$ .

We drop column  $K$  from matrix  $P$  and apply Theorem 5 with this new  $J \times (K - 1)$  matrix,  $P_{\{1, \dots, K-1\}}$ , which has rank  $K - 1$ . By assumption,  $\underline{1}'(P'_{\{1, \dots, K-1\}}P_{\{1, \dots, K-1\}})^{-1}\underline{1} \neq 1$ .

We then apply Theorem 5 to  $P_{\{1, \dots, K-1\}}$ . The maximum of  $H_{\text{adm}}$  occurs either at the point  $\gamma_{\mathcal{S}}$ , where  $\mathcal{S} = \{1, 2, \dots, K - 1\}$ , or on the boundary of the set  $\{\underline{\gamma} : \underline{1}'\underline{\gamma} = 1 \text{ and } \underline{\gamma} \in \Delta^{K-2}\}$ .

We repeat this method of descent, decrementing the dimension (and permuting population labels without loss of generality) until we reach the case of only two source populations. A final application of Theorem 5 then finds that  $H_{\text{adm}}$  is maximized either interior to the 1-simplex—the line connecting vertices  $(1, 0)$  and  $(0, 1)$ —or at one of these vertices.  $\square$

## Appendix 2. Proofs for $K = 2$ : Proposition 7 and Corollary 8

*Proof of Proposition 7.* We maximize the quadratic polynomial in eqs. 12-14 over  $\gamma \in [0, 1]$ . The maximum occurs at the unique critical point or on the boundary of the interval.

Setting the derivative of eq. 14 with respect to  $\gamma_1$  to 0, we find that the critical point is

$$(\gamma_1^*, H_{\text{adm}}) = \left( \frac{C_{12} - H_2}{2(C_{12} - H_S)}, \frac{C_{12}^2 - H_1 H_2}{2(C_{12} - H_S)} \right). \quad (23)$$

Because the leading coefficient of eq. 14 is negative for  $\underline{p}_1 \neq \underline{p}_2$ , the critical point is a maximum. Hence, if  $(C_{12} - H_2)/[2(C_{12} - H_S)] \in (0, 1)$ , then the maximum of  $H_{\text{adm}}$  on the interval  $[0, 1]$  lies at  $\gamma_1 = (C_{12} - H_2)/[2(C_{12} - H_S)]$ . Otherwise, the maximum lies either at  $\gamma_1 = 0$ , in which case it equals  $H_2$ , or at  $\gamma_1 = 1$ , in which case it equals  $H_1$ .

The conditions describing the location of the maximum can be written in terms of  $H_1$ ,  $H_2$ , and  $C_{12}$ . Because the denominator of  $\gamma_1^*$  in eq. 23 is always positive for  $\underline{p}_1 \neq \underline{p}_2$  (Section 4),  $\gamma_1^* \in (0, 1)$  becomes equivalent to  $C_{12} > H_1$  and  $C_{12} > H_2$ , the former inequality arising from the condition  $\gamma_1^* < 1$  and the latter from the condition  $\gamma_1^* > 0$ .

If the requirement  $C_{12} > H_1$  and  $C_{12} > H_2$  for  $\gamma_1^* \in (0, 1)$  fails, then the maximum occurs on the boundary of the unit interval. We have  $H_{\text{adm}}(0) = H_2$  and  $H_{\text{adm}}(1) = H_1$ . Thus, the maximum lies at  $\gamma_1 = 0$  if  $H_2 > H_1$  and at  $\gamma_1 = 1$  if  $H_1 > H_2$ .

If  $C_{12} > H_1$  and  $C_{12} > H_2$  do not both hold, then one of them must hold, as we showed in Section 4 that  $2C_{12} > H_1 + H_2$ . Combining the fact that either  $C_{12} > H_1$  or  $C_{12} > H_2$

holds with the observation that  $H_2 > H_1$  leads to a maximum at  $\gamma_1 = 0$  and  $H_1 > H_2$  leads to a maximum at  $\gamma_1 = 1$ , we complete the characterization of the three cases.  $\square$

Note that alternative expressions in terms of  $H_1$ ,  $H_2$ , and  $F_{12}$  can be derived by noting that  $H_S = \frac{1}{2}(H_1 + H_2)$ ,  $H_1H_2 = H_S^2 - [(H_1 - H_2)/2]^2$  and  $C_{12} = H_S(1 + F_{12})/(1 - F_{12})$ , the latter relationship simply restating eq. 4 (recalling  $C_{12} = 1$  for  $F_{12} = 1$ ). Thus, we have

$$\gamma_1^* = \frac{C_{12} - H_2}{2(C_{12} - H_S)} = \frac{1}{2} + \frac{H_1 - H_2}{4 \frac{F_{12}}{1-F_{12}}(H_1 + H_2)}, \quad (24)$$

$$H_{\text{adm}}(\gamma^*) = \frac{C_{12}^2 - H_1H_2}{2(C_{12} - H_S)} = \frac{H_1 + H_2}{2(1 - F_{12})} + \frac{(H_1 - H_2)^2}{8 \frac{F_{12}}{1-F_{12}}(H_1 + H_2)}. \quad (25)$$

Another formulation uses the heterozygosity of a population formed by equal admixture of populations 1 and 2, or  $H_T$ . Because  $F_{12} = 1 - H_S/H_T$  by eq. 1,  $F_{12}/(1 - F_{12}) = (H_T - H_S)/H_S$ . Using this relationship in eqs. 24 and 25:

$$\gamma_1^* = \frac{1}{2} + \frac{H_1 - H_2}{8(H_T - H_S)},$$

$$H_{\text{adm}}(\gamma^*) = H_T + \frac{(H_1 - H_2)^2}{16(H_T - H_S)}.$$

*Proof of Corollary 8.* We restate the condition  $0 < (C_{12} - H_2)/[2(C_{12} - H_S)] < 1$  as

$$0 < \frac{1}{2} + \frac{\left(\frac{H_1 - H_2}{2}\right)}{2 \frac{F_{12}}{1-F_{12}}(H_1 + H_2)} < 1.$$

Subtracting  $\frac{1}{2}$  from both sides and multiplying by 2, an equivalent condition is

$$-1 < \frac{(H_1 - H_2)}{2 \frac{F_{12}}{1-F_{12}}(H_1 + H_2)} < 1,$$

or, equivalently,  $|H_1 - H_2|/[2 \frac{F_{12}}{1-F_{12}}(H_1 + H_2)] < 1$ . We rearrange this last expression to obtain the desired result.  $\square$

### Appendix 3: Dirichlet model for allele frequencies

We first provide results concerning  $H_{\text{adm}}$  in the case that the  $K$  source populations have independently and identically distributed (IID) allele frequency vectors. Next, we specify these IID vectors to be Dirichlet distributions.

## IID allele frequency vectors

We begin by examining the expected values of  $H_k$  and  $H_{\text{adm}}$ .

**Proposition 11.** Suppose the allele frequency vectors  $\underline{p}_k$  are independently and identically distributed for  $1 \leq k \leq K$ . Then  $\mathbb{E}[H_{\text{adm}}] = \mathbb{E}[H_1] + (1 - \sum_{k=1}^K \gamma_k^2)(\sum_{j=1}^J \text{Var}[p_{1j}])$ .

*Proof.* We use eq. 8:

$$\mathbb{E}[H_{\text{adm}}] = 1 - \sum_{k=1}^K \gamma_k^2 \left[ \sum_{j=1}^J \mathbb{E}[p_{kj}^2] \right] - 2 \sum_{k=1}^{K-1} \sum_{\ell=k+1}^K \gamma_k \gamma_\ell \left[ \sum_{j=1}^J \mathbb{E}[p_{kj} p_{\ell j}] \right].$$

Using the IID assumption and simplifying by noting that  $1 = (\sum_{k=1}^K \gamma_k)^2 = (\sum_{k=1}^K \gamma_k^2) + (2 \sum_{k=1}^{K-1} \sum_{\ell=k+1}^K \gamma_k \gamma_\ell)$ , we have

$$\begin{aligned} \mathbb{E}[H_{\text{adm}}] &= 1 - \left( \sum_{k=1}^K \gamma_k^2 \right) \left( \sum_{j=1}^J \mathbb{E}[p_{1j}^2] \right) - 2 \sum_{k=1}^{K-1} \sum_{\ell=k+1}^K \gamma_k \gamma_\ell \left[ \sum_{j=1}^J (\mathbb{E}[p_{1j}])^2 \right] \\ &= 1 - \sum_{j=1}^J \mathbb{E}[p_{1j}^2] + \sum_{j=1}^J \mathbb{E}[p_{1j}^2] \left( 1 - \sum_{k=1}^K \gamma_k^2 \right) - \sum_{j=1}^J (\mathbb{E}[p_{1j}])^2 \left( 1 - \sum_{k=1}^K \gamma_k^2 \right), \end{aligned}$$

from which the result follows.  $\square$

An immediate corollary of Proposition 11 is that under the IID assumption,  $H_{\text{adm}}$  has expectation greater than or equal to the expectation of the heterozygosity of each of the source populations.

**Corollary 12.** Suppose the allele frequency vectors  $\underline{p}_k$  are independently and identically distributed for  $1 \leq k \leq K$ . Then  $\mathbb{E}[H_{\text{adm}}] \geq \mathbb{E}[H_k]$ .

A second corollary results from the Cauchy-Schwarz inequality, by which  $\sum_{k=1}^K \gamma_k^2 \geq \frac{1}{K}$ , with equality if and only if  $(\gamma_1, \gamma_2, \dots, \gamma_K) = (\frac{1}{K}, \frac{1}{K}, \dots, \frac{1}{K})$ .

**Corollary 13.** Suppose the allele frequency vectors  $\underline{p}_k$  are independently and identically distributed for  $1 \leq k \leq K$ . Considering all admixture vectors  $\underline{\gamma} \in \Delta^{K-1}$ ,  $\mathbb{E}[H_{\text{adm}}]$  is maximized at  $\underline{\gamma} = (\frac{1}{K}, \frac{1}{K}, \dots, \frac{1}{K})$ , and has maximal value  $\mathbb{E}[H_1] + (1 - \frac{1}{K}) \sum_{j=1}^J \text{Var}[p_{1j}]$ .

## IID allele frequency vectors from a symmetric Dirichlet distribution

We now further assume that the independently and identically distributed allele frequency vectors follow a symmetric multivariate Dirichlet distribution. This distribution is frequently used for allele frequency distributions (Balding & Nichols, 1995; Pritchard *et al.*, 2000; Huelsenbeck & Andolfatto, 2007), and it is a natural probability distribution to assume for allelic types with the same marginal distributions.

The  $J$ -dimensional Dirichlet- $(\alpha_1, \alpha_2, \dots, \alpha_J)$  distribution is defined over the open unit  $(J - 1)$ -simplex  $\Delta^{J-1}$  and has concentration parameters  $\alpha_j > 0$ . The means and variances for the individual allele frequencies are (Lange, 1997; Kotz *et al.*, 2000, chapter 49):

$$\begin{aligned}\mathbb{E}[p_{kj}] &= \frac{\alpha_j}{J\bar{\alpha}}, \\ \text{Var}[p_{kj}] &= \frac{\alpha_j(J\bar{\alpha} - \alpha_j)}{J^2\bar{\alpha}^2(J\bar{\alpha} + 1)},\end{aligned}\tag{26}$$

where  $\bar{\alpha} = \frac{1}{J} \sum_{j=1}^J \alpha_j$ .

The symmetric Dirichlet distribution assumes  $\alpha_1 = \alpha_2 = \dots = \alpha_J = \bar{\alpha}$ , leading to:

$$\begin{aligned}\mathbb{E}[p_{kj}] &= \frac{1}{J}, \\ \text{Var}[p_{kj}] &= \frac{J - 1}{J^2(J\bar{\alpha} + 1)}.\end{aligned}\tag{27}$$

Making these substitutions in Proposition 11, we obtain the expectation of  $H_{\text{adm}}$  under the assumption that the allele frequency vectors follow independent Dirichlet distributions.

**Corollary 14.** Suppose the allele frequency vectors  $\underline{p}_k$  are independently and identically distributed for  $1 \leq k \leq K$ , all with symmetric multivariate Dirichlet distributions with concentration parameter  $\bar{\alpha}$ . Then

$$\begin{aligned}\mathbb{E}[H_k] &= \left(1 - \frac{1}{J}\right) \left(1 - \frac{1}{J\bar{\alpha} + 1}\right), \\ \mathbb{E}[H_{\text{adm}}] &= \left(1 - \frac{1}{J}\right) \left(1 - \frac{1}{J\bar{\alpha} + 1} \sum_{k=1}^K \gamma_k^2\right).\end{aligned}$$

This corollary implies that both  $\mathbb{E}[H_k]$  and  $\mathbb{E}[H_{\text{adm}}]$  are increasing functions of  $J$  and  $\bar{\alpha}$ .

The next proposition considers the special case of  $K = 2$  and  $J = 2$ , further specifying a uniform distribution for  $\gamma_1$ .

**Proposition 15.** Consider  $K = 2$  and  $J = 2$ . Suppose that the values of  $p_{11}$  and  $p_{21}$  are independently chosen from a uniform-[0,1] distribution. Suppose also that  $\gamma_1$  is also chosen from a uniform-[0,1] distribution. Then  $\mathbb{P}[H_{\text{adm}}(\gamma_1) > \max\{H_1, H_2\}] = 1 - \log 2 \approx 0.307$ .

*Proof.* Using Proposition 9, we identify the regions of the unit square for  $(p_{11}, p_{21})$  in which  $\max_{\gamma_1 \in (0,1)} H_{\text{adm}}(\gamma_1) > \max\{H_1, H_2\}$ . These regions are  $\{(p_{11}, p_{21}) \mid \frac{1}{2} < p_{11} < 1, 0 < p_{21} < \frac{1}{2}\}$  and  $\{(p_{11}, p_{21}) \mid 0 < p_{11} < \frac{1}{2}, \frac{1}{2} < p_{21} < 1\}$ .

Within those regions, we must determine the portion of the unit interval for  $\gamma_1$  in which  $H_{\text{adm}}(\gamma_1) > \max\{H_1, H_2\}$ .  $H_{\text{adm}}(\gamma_1)$  is a quadratic function of  $\gamma_1$ . We ignore the set of zero volume with  $H_1 = H_2$ . In the regions for  $(p_{11}, p_{21})$  in which  $\max_{\gamma_1 \in (0,1)} H_{\text{adm}}(\gamma_1) > \max\{H_1, H_2\}$  and  $H_2 > H_1$ , the interval for  $\gamma_1$  in which  $H_{\text{adm}}(\gamma_1) > H_1$  is  $(0, \frac{1-2p_{21}}{p_{11}-p_{21}})$ . In the regions for  $(p_{11}, p_{21})$  in which  $\max_{\gamma_1 \in (0,1)} H_{\text{adm}}(\gamma_1) > \max\{H_1, H_2\}$  and  $H_1 > H_2$ , the interval for  $\gamma_1$  in which  $H_{\text{adm}}(\gamma_1) > H_1$  is  $(\frac{p_{21}-1+p_{11}}{p_{21}-p_{11}}, 1)$ .

The desired probability is the volume within the unit cube for  $(p_{11}, p_{21}, \gamma_1)$  of the regions in which  $H_{\text{adm}}(\gamma_1) > \max\{H_1, H_2\}$ . The volume is

$$\begin{aligned} & \int_{1/2}^1 \int_{1-p_{11}}^{1/2} \int_0^{\frac{1-2p_{21}}{p_{11}-p_{21}}} 1 \, d\gamma_1 \, dp_{21} \, dp_{11} + \int_{1/2}^1 \int_0^{1-p_{11}} \int_{\frac{p_{21}-1+p_{11}}{p_{21}-p_{11}}}^1 1 \, d\gamma_1 \, dp_{21} \, dp_{11} \\ & \int_0^{1/2} \int_{1-p_{11}}^1 \int_{\frac{p_{21}-1+p_{11}}{p_{21}-p_{11}}}^1 1 \, d\gamma_1 \, dp_{21} \, dp_{11} + \int_0^{1/2} \int_{1/2}^{1-p_{11}} \int_0^{\frac{1-2p_{21}}{p_{11}-p_{21}}} 1 \, d\gamma_1 \, dp_{21} \, dp_{11} \\ & = 4 \frac{1 - \log 2}{4}. \end{aligned}$$

□

## References

- Balding, D. J. and Nichols, R. A. 1995. A method for quantifying differentiation between populations at multi-allelic loci and its implications for investigating identity and paternity, *Genetics* **96**, 3–12.
- Boca, S. M. and Rosenberg, N. A. 2011. Mathematical properties of  $F_{st}$  between admixed populations and their parental source populations, *Theor. Pop. Biol.* **80**, 208–216.
- Buerkle, C. A. and Lexer, C. 2008. Admixture as the basis for genetic mapping, *Trends Ecol. Evol.* **23**, 686–694.
- Chakraborty, R. 1986. Gene admixture in human populations: Models and predictions, *Yrbk. Phys. Anthropol.* **29**, 1–43.
- Gravel, S. 2012. Population genetics models of local ancestry, *Genetics* **191**, 607–619.
- Graybill, F. A. 1976. “Theory and application of the linear model”, Duxbury, Pacific Grove, CA.
- Horn, R. A. and Johnson, C. R. 2012. “Matrix analysis”, Cambridge University Press, New York, NY.

- Huelsenbeck, J. P. and Andolfatto, P. 2007. Inference of population structure under a Dirichlet process model, *Genetics* **175**, 1787–1802.
- Kotz, S., Balakrishnan, N., and Johnson, N. L. 2000. “Continuous Multivariate Distributions. Volume 1: Models and Applications”, Wiley, New York.
- Lange, K. 1997. “Mathematical and Statistical Methods for Genetic Analysis”, Springer, New York.
- Long, J. C. 1991. The genetic structure of admixed populations, *Genetics* **127**, 417–428.
- Magnus, J. R. and Neudecker, H. 2007. “Matrix differential calculus with applications in statistics and econometrics”, John Wiley & Sons, Chichester, UK, 3rd edition.
- Mehta, R. S., Feder, A. F., Boca, S. M., and Rosenberg, N. A. 2019. The relationship between haplotype-based FST and haplotype length., *Genetics* **213**, 281–295.
- Millar, R. B. 1987. Maximum likelihood estimation of mixed stock fishery composition, *Can. J. Fish. Aquat. Sci.* **44**, 583–590.
- Mooney, J. A., Huber, C. D., Service, S., Sul, J. H., Marsden, C. D., Zhang, Z., Sabatti, C., Ruiz-Linares, A., Bedoya, G., Costa Rica/Colombia Consortium for Genetic Investigation of Bipolar Endophenotypes, Freimer, N., and Lohmueller, K. E. 2018. Understanding the hidden complexity of Latin American population isolates, *Am. J. Hum. Genet.* **103**, 707–726.
- Pemberton, T. J., DeGiorgio, M., and Rosenberg, N. A. 2013. Population structure in a comprehensive genomic data set on human microsatellite variation, *G3: Genes, Genomes, Genetics* **3**, 891–907.
- Pritchard, J. K., Stephens, M., and Donnelly, P. 2000. Inference of population structure using multilocus genotype data, *Genetics* **155**, 945–959.
- Risch, N., Choudhry, S., Via, M., Basu, A., Sebro, R., Eng, C., Beckman, K., Thyne, S., Chapela, R., Rodriguez-Santana, J. R., Rodriguez-Cintron, W., Avila, P. C., Ziv, E., and Burchard, E. G. 2009. Ancestry-related assortative mating in Latino populations, *Genome Biol.* **10**, R132.
- Schroeder, K. B., Jakobsson, M., Crawford, M. H., Schurr, T. G., Boca, S. M., Conrad, D. F., Tito, R. Y., Osipova, L. P., Tarskaia, L. A., Zhadanov, S. I., Wall, J. D., Pritchard, J. K., Malhi, R. S., Smith, D. G., and Rosenberg, N. A. 2009. Haplotypic background of a private allele at high frequency in the Americas, *Mol. Biol. Evol.* **26**, 995–1016.
- Verdu, P. and Rosenberg, N. A. 2011. A general mechanistic model for admixture histories of hybrid populations, *Genetics* **189**, 1413–1426.
- Wang, S., Ray, N., Rojas, W., Parra, M. V., Bedoya, G., Gallo, C., Poletti, G., Mazzotti, G., Hill, K., Hurtado, A. M., Camrena, B., Nicolini, H., Klitz, W., Barrantes, R., Molina, J. A., Freimer, N. B., Bortolini, M. C., Salzano, F. M., Petzl-Erler, M. L., Tsuneto, L. T., Dipierri, J. E., Alfaro, E. L., Bailliet, G., Bianchi, N. O., Llop, E., Rothhammer, F., Excoffier, L., and Ruiz-Linares, A. 2008. Geographic patterns of genome admixture in Latin American Mestizos, *PLoS Genet.* **4**, e1000037.
- Zhu, X., Tang, H., and Risch, N. 2008. Admixture mapping and the role of population structure for localizing disease genes, *Adv. Genet.* **60**, 547–569.
- Zou, J. Y., Park, D. S., Burchard, E. G., Torgerson, D. G., Pino-Yanes, M., Song, Y. S., Sankararaman, S., Halperin, E., and Zaitlen, N. 2015. Genetic and socioeconomic study of mate choice in Latinos reveals novel assortment patterns, *Proc. Natl. Acad. Sci. USA* **112**, 13621–13626.

Table 1: Notation

| Type of quantity    | Symbol               | Description  |
|---------------------|----------------------|--|
| Indices             | $j = 1, \dots, J$    | Index over alleles   |
|                     | $k = 1, \dots, K$    | Index over source populations  |
| Allele frequencies  | $p_{kj}$             | Frequency of allelic type $j$ in population $k$  |
|                     | $\underline{p}_k$    | $J \times 1$ vector of allele frequencies for population $k$   |
|                     | $\overline{P}$       | $J \times K$ matrix of allele frequencies in the source populations  |
|                     | $\bar{p}_j$          | Frequency of allelic type $j$ in the admixed population  |
| Admixture fractions | $\gamma_k$           | Admixture fraction for population $k$  |
|                     | $\underline{\gamma}$ | $K \times 1$ vector of admixture fractions   |
| Heterozygosities    | $\overline{H}_k$     | Heterozygosity for population $k$ ; probability that two alleles drawn from population $k$ differ in type      |
|                     | $H_{\text{adm}}$     | Heterozygosity for the admixed population  |
|                     | $C_{k\ell}$          | Probability that an allele drawn from population $k$ and an allele drawn from population $\ell$ differ in type |
| Fixation index      | $F_{k\ell}$          | Fixation index $F_{ST}$ between populations $k$ and $\ell$   |

Figure 1:  $H_{\text{adm}}$  versus  $\gamma_1$  for fixed values of  $H_1$  and  $H_2$ . We choose  $H_1 = 0.727$  and  $H_2 = 0.628$ ; the horizontal lines represent  $H_{\text{adm}} = H_1$  and  $H_{\text{adm}} = H_2$ . Eq. 13 is plotted for multiple values of  $F_{12}$ , considering the allowable range of  $F_{12}$  values in  $[0.003, 0.192]$  as specified by eq. 5. The red curve, which plots  $(\gamma_1, H_{\text{adm}})$  in terms of  $H_1$ ,  $H_2$ , and  $F_{12}$  in the form of eqs. 24 and 25, indicates the maxima of  $H_{\text{adm}}$  as  $F_{12}$  varies, with black dots specifying the maxima for the specific plotted values of  $F_{12}$ . The shaded region corresponds to the region where  $\gamma_1^* \in (0, 1)$ , as specified by Corollary 8; the value  $F_{12} \approx 0.034$  gives the boundary of this region. The values chosen for  $H_1$  and  $H_2$  are, respectively, the mean heterozygosities across 8 European and 29 Native American populations, based on population-wise estimates in Table S20 of Pemberton *et al.* (2013). The value of  $\gamma_1$  can be viewed as the fraction of European ancestry in an admixed population and  $1 - \gamma_1$  can be considered the fraction of Native American ancestry.

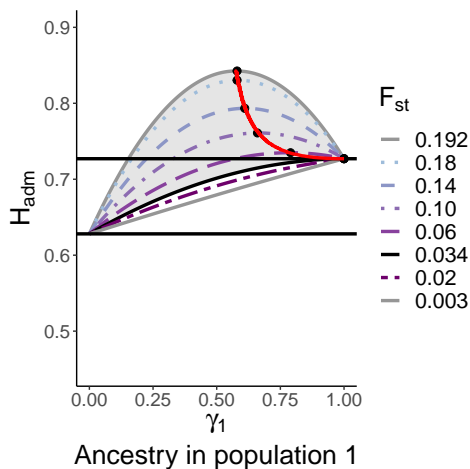




Figure 2: The admixture coefficient  $\gamma_1$  that maximizes  $H_{\text{adm}}$  in the case of  $K = 2$  source populations and  $J = 2$  allelic types. The plot shows the unit square for  $(p_{11}, p_{21})$ . In the red regions, the maximizing value of  $\gamma_1$  lies in  $(0, 1)$ , whereas in the white and gray regions, it lies on one or the other boundary. The figure depicts the result of Proposition 9.

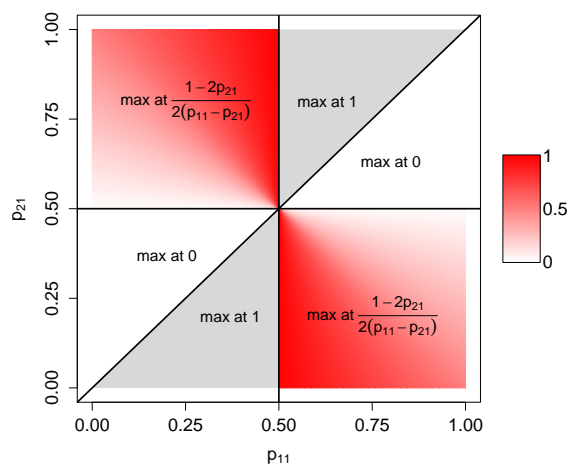


Figure 3: Location of the maximum of  $H_{\text{adm}}$  in simulation replicates. (A)  $K = 2$ . (B)  $K = 3$ . The location  $\underline{\gamma}_{\text{arg max}}$  can be in the interior of the simplex  $\Delta^{K-1}$ , corresponding to nontrivial admixture of all source groups, or on the boundary of the simplex. For  $K = 3$ , it can be on an edge, corresponding to admixture of two of three source populations, and for both  $K = 2$  and  $K = 3$ , it can be at a vertex, corresponding to membership in only one source population. For each  $(K, J)$ , points plotted are based on 10,000 simulations with independently and identically distributed Dirichlet- $(1, 1, \dots, 1)$  distributions for the allele frequency vectors  $\underline{p}_k$  in the  $K$  populations.

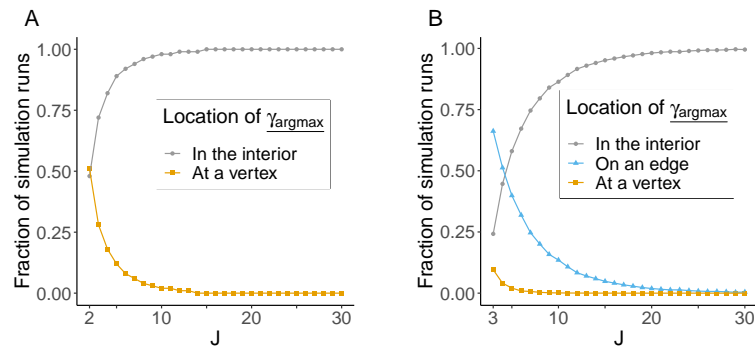


Figure 4: The fraction of simulation replicates for which  $H_{\text{adm}} > \max\{H_1, \dots, H_K\}$ , for various values of  $K$  and  $J$ . For each  $(K, J)$ , points plotted are based on 50,000 simulation replicates with independently and identically distributed Dirichlet-(1, 1, ..., 1) distributions for the allele frequency vectors  $\underline{p}_k$  in the  $K$  populations, and a Dirichlet-(1, 1, ..., 1) distribution for the admixture coefficient vector  $\underline{\gamma}$ .

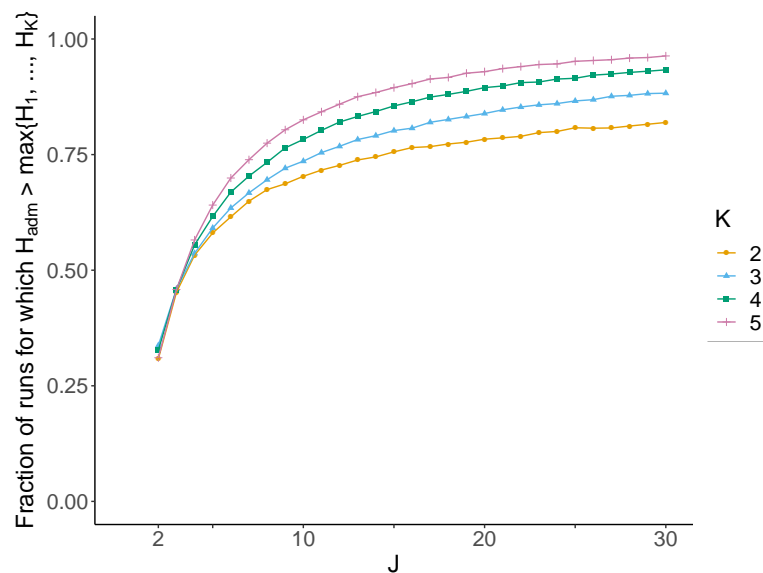


Figure 5:  $H_{\text{adm}}$  versus  $\gamma_1$  for 10 simulation replicates for  $K = 2$  source populations, for each of three values of the number of allelic types  $J$ . For each replicate, allele frequency vectors  $p_k$  in the two populations are simulated according to Dirichlet- $(1, 1, \dots, 1)$  distributions, and  $H_{\text{adm}}$  is plotted as a function of  $\gamma_1$  according to eq. 8. The maximum of  $H_{\text{adm}}$  is indicated by a black circle in each replicate. The red dashed lines represent the expected values of  $H_{\text{adm}}$  according to Corollary 14 in Appendix 3.

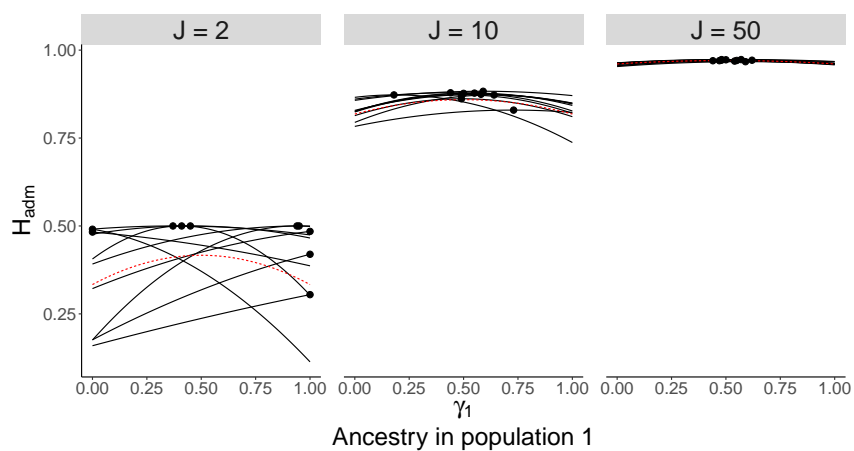


Figure 6:  $H_{\text{adm}}$  versus  $\gamma_1$  for 20 random loci from Wang *et al.* (2008). The two source populations providing the allele frequencies are the European and Native American populations, with  $\gamma_1$  corresponding to membership in the European population.  $H_{\text{adm}}$  is plotted according to eq. 8. Circles indicate the location of the maximum along each curve.

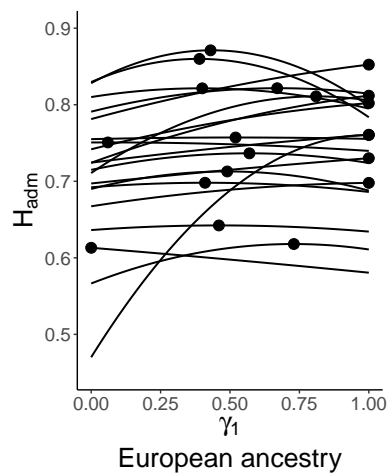


Figure 7: Predicted and observed  $H_{\text{adm}}$ . (A) The predicted and observed  $H_{\text{adm}}$  values for an admixed Mestizo population are plotted against the locus-wise estimated European admixture fraction  $\hat{\gamma}_1$  in the Mestizo population, estimated by maximum likelihood. The prediction is based on eq. 8, using European and Native American allele frequencies estimated from Wang *et al.* (2008) as  $p_1$  and  $p_2$ , respectively, together with the maximum likelihood estimate of  $\gamma_1$ . The observation is based on  $H_{\text{adm}}$  computed from Definition 1, inserting estimated allele frequencies from Wang *et al.* (2008) for the Mestizo population. (B) The observed  $H_{\text{adm}}$  value is plotted against the predicted  $H_{\text{adm}}$  value. The identity line is shown in gray. In both panels, each point represents one of the 678 loci used. The correlation coefficient between the predicted and observed  $H_{\text{adm}}$  values is 0.978.

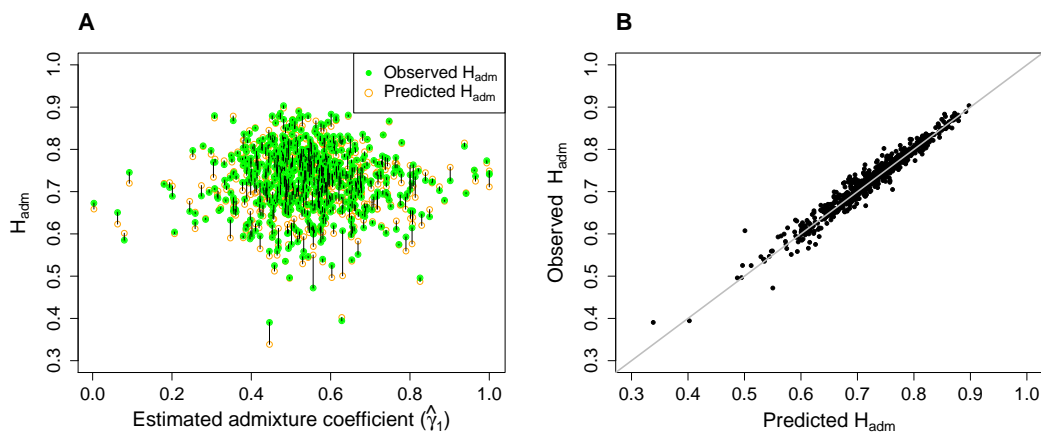
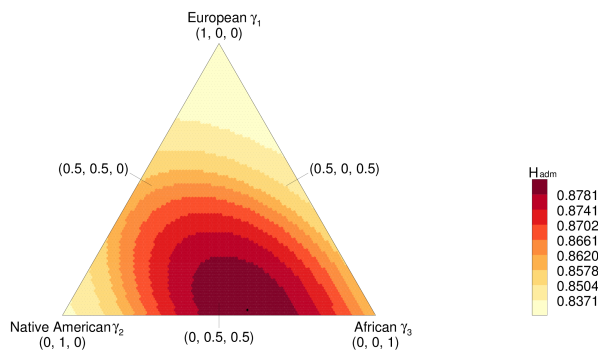
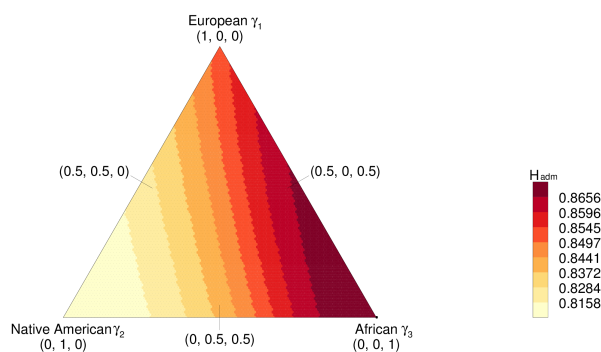


Figure 8:  $H_{adm}$  versus  $(\gamma_1, \gamma_2, \gamma_3)$  for three loci. The loci are from Wang *et al.* (2008) and have 14, 14, and 8 distinct alleles, respectively. The value of  $H_{adm}$  is computed from eq. 8. Black circles indicate the maximum  $H_{adm}$ . (A) Locus D2S1399: the maximum lies in the interior of the region. (B) Locus GATA101G01: the maximum lies at the  $(0, 0, 1)$  vertex. (C) Locus GATA146D07: the maximum lies on the  $\gamma_2 = 0$  edge.

A



B



C

