

Algorithms and equations utilized in poppr version 1.1.0.99

Zhian N. Kamvar¹ and Niklaus J. Grünwald^{1,2}

1) Department of Botany and Plant Pathology, Oregon State University, Corvallis, OR

2) Horticultural Crops Research Laboratory, USDA-ARS, Corvallis, OR

May 24, 2014

Abstract

This vignette is focused on simply explaining the different algorithms utilized in calculations such as the index of association and different distance measures. Many of these are previously described in other papers and it would be prudent to cite them properly if they are used.

Contents

1	Mathematical representation of data in <i>adegenet</i> and <i>poppr</i>	1
2	The Index of Association	2
3	Genetic distances	4
3.1	Distances that assume genetic drift	4
3.1.1	Nei's 1978 Distance	4
3.1.2	Edwards' angular distance	4
3.1.3	Reynolds' coancestry distance	4
3.2	Distances without assumptions	5
3.2.1	Rogers' distance	5
3.2.2	Provesti's absolute genetic distance	5
3.3	Bruvo's distance (stepwise mutation for microsatellites)	5
3.3.1	Special cases of Bruvo's distance	6
3.4	Example	6

1 Mathematical representation of data in *adegenet* and *poppr*

The sections dealing with the index of association and genetic distances will be based on the same data structure, a matrix with samples in rows and alleles in columns. The number of columns is equal to the total number of alleles observed in the data set. Much of this description is derived from *adegenet*'s `dist.genpop` manual page.

Let **A** be a table containing allelic frequencies with t samples¹ (rows) and m alleles (columns).

The above statement describes the table present in `genind` or `genpop` object where, instead of having the number of columns equal the number of loci, the number of columns equals the number of observed alleles in the entire data set.

¹populations or individuals

Let ν be the number of loci. The locus j gets $m(j)$ alleles.

$$m = \sum_{j=1}^{\nu} m(j) \quad (1)$$

So, if you had a data set with 5 loci that had 2 alleles each, your table would have ten columns. Of course, codominant loci like microsatellites have varying numbers of alleles.

For the row i and the modality k of the variable j , notice the value

$$a_{ijk} (1 \leq i \leq t, 1 \leq j \leq \nu, 1 \leq k \leq m(j)) \quad (2)$$

$$a_{ij\cdot} = \sum_{k=1}^{m(j)} a_{ijk} \quad (3)$$

$$p_{ijk} = \frac{a_{ijk}}{a_{ij\cdot}} \quad (4)$$

The above couple of equations are basically defining the allele counts (a_{ijk}) and frequency (p_{ijk}). Remember that i is individual, j is locus, and k is allele. The following continues to describe properties of the frequency table used for analysis:

$$p_{ij\cdot} = \sum_{k=1}^{m(j)} p_{ijk} = 1 \quad (5)$$

The sum of all allele frequencies for a single population (or individual) at a single locus is one.

$$p_{i\cdot\cdot} = \sum_{j=1}^{\nu} p_{ij\cdot} = \nu \quad (6)$$

The sum of all allele frequencies over all loci is equal to the number of loci.

$$p_{\dots} = \sum_{j=1}^{\nu} p_{i\cdot\cdot} = t\nu \quad (7)$$

The the sum of the entire table is the sum of all loci multiplied by the number of populations (or individuals).

2 The Index of Association

The index of association was originally developed by A.H.D. Brown analyzing population structure of wheat and has been widely used as a tool to detect clonal reproduction within populations [2, 4]. Populations whose members are undergoing sexual reproduction, whether it be selfing or out-crossing, will produce gametes via meiosis, and thus have a chance to shuffle alleles in the next generation. Populations whose members are undergoing clonal reproduction, however, generally do so via mitosis.

The most likely mechanism, therefor for a change in genotype for a clonal organism is via mutation. The rate of mutation varies from species to species, but it is rarely sufficiently high to approximate a random shuffling of alleles. The index of association is a calculation based on the ratio of the variance of the raw

number of differences between individuals and the sum of those variances over each locus [4]. It can also be thought of as the observed variance over the expected variance. If both variances are equal, then the index is zero after subtracting one (from Maynard-Smith, 1993 [4]):

$$I_A = \frac{V_O}{V_E} - 1 \quad (8)$$

Any sort of marker can be used for this analysis as it only counts differences between pairs of samples. This can be thought of as a distance whose maximum is equal to the number of loci multiplied by the ploidy of the sample. This is calculated using an absolute genetic distance.

Remember that in *poppr*, genetic data is stored in a table where the rows represent samples and the columns represent potential allelic states grouped by locus. Notice also that the sum of the rows all equal one. *Poppr* uses this to calculate distances by simply taking the sum of the absolute values of the differences between rows.

The calculation for the distance between two individuals at a single locus with k allelic states and a ploidy of l is as follows²:

$$d(a, b) = \frac{l}{2} \sum_{j=1}^k |p_{ajk} - p_{bjk}| \quad (9)$$

To find the total number of differences between two individuals over all loci, you just take d over ν loci, a value we'll call D :

$$D(a, b) = \sum_{i=1}^{\nu} d_i \quad (10)$$

An interesting observation: $D(a, b)/(l\nu)$ is Provesti's distance.

These values are calculated over all possible combinations of individuals in the data set, $\binom{n}{2}$ after which you end up with $\binom{n}{2} \cdot \nu$ values of d and $\binom{n}{2}$ values of D . Calculating the observed variances is fairly straightforward (modified from Agapow and Burt, 2001) [1]:

$$V_O = \frac{\sum_{i=1}^{\binom{n}{2}} D_i^2 - \frac{(\sum_{i=1}^{\binom{n}{2}} D_i)^2}{\binom{n}{2}}}{\binom{n}{2}} \quad (11)$$

Calculating the expected variance is the sum of each of the variances of the individual loci. The calculation at a single locus, j is the same as the previous equation, substituting values of D for d [1]:

$$var_j = \frac{\sum_{i=1}^{\binom{n}{2}} d_i^2 - \frac{(\sum_{i=1}^{\binom{n}{2}} d_i)^2}{\binom{n}{2}}}{\binom{n}{2}} \quad (12)$$

The expected variance is then the sum of all the variances over all ν loci [1]:

$$V_E = \sum_{j=1}^{\nu} var_j \quad (13)$$

Now you can plug the sums of equations (11) and (13) into equation (8) to get the index of association. Of course, Agapow and Burt showed that this index increases steadily with the number of loci, so they came

²Individuals with Presence / Absence data will have the $l/2$ term dropped.

up with an approximation that is widely used, \bar{r}_d [1]. For the derivation, see the manual for *multilocus*. The equation is as follows, utilizing equations (11), (12), and (13) [1]:

$$\bar{r}_d = \frac{V_O - V_E}{2 \sum_{j=1}^{m\nu} \sum_{k \neq j}^m \sqrt{\text{var}_j \cdot \text{var}_k}} \quad (14)$$

3 Genetic distances

Genetic distances are great tools for analyzing diversity in populations as they are the basis for creating dendrograms with bootstrap support and also for AMOVA. This section will simply present different genetic distances along with a few notes about them. Most of these distances are derived from the *ade4* and *adegenet* packages, where they were implemented as distances between populations. *Popp* extends the implementation to individuals as well (with the exception of Bruvo’s distance).

Table 1: Distance measures and their respective assumptions

Method	Function	Assumption	Euclidean
Provesti	<code>provesti.dist</code> <code>diss.dist</code>	-	No
Nei	<code>nei.dist</code>	Infinite Alleles Genetic Drift	No
Edwards	<code>edwards.dist</code>	Genetic Drift	Yes
Reynolds	<code>reynolds.dist</code>	Genetic Drift	Yes
Rogers	<code>rogers.dist</code>	-	Yes
Bruvo	<code>bruvo.dist</code>	Stepwise Mutation	No

3.1 Distances that assume genetic drift

3.1.1 Nei’s 1978 Distance

$$D_{Nei}(a, b) = -\ln\left(\frac{\sum_{k=1}^{\nu} \sum_{j=1}^{m(k)} p_{ajk} p_{bjk}}{\sqrt{\sum_{k=1}^{\nu} \sum_{j=1}^{m(k)} (p_{ajk})^2} \sqrt{\sum_{k=1}^{\nu} \sum_{j=1}^{m(k)} (p_{bjk})^2}}\right) \quad (15)$$

Note: if comparing individuals in *poppr*, those that do not share any alleles normally receive a distance of ∞ . As you cannot draw a dendrogram with infinite branch lengths, these are converted to an order of magnitude higher distance than the largest observed less than ∞ .

3.1.2 Edwards’ angular distance

$$D_2(a, b) = \sqrt{1 - \frac{1}{\nu} \sum_{k=1}^{\nu} \sum_{j=1}^{m(k)} \sqrt{p_{ajk} p_{bjk}}} \quad (16)$$

3.1.3 Reynolds’ coancestry distance

$$D_3(a, b) = \sqrt{\frac{\sum_{k=1}^{\nu} \sum_{j=1}^{m(k)} (p_{ajk} - p_{bjk})^2}{2 \sum_{k=1}^{\nu} (1 - \sum_{j=1}^{m(k)} p_{ajk} p_{bjk})}} \quad (17)$$

3.2 Distances without assumptions

3.2.1 Rogers' distance

$$D_4(a, b) = \frac{1}{\nu} \sum_{k=1}^{\nu} \sqrt{\frac{1}{2} \sum_{j=1}^{m(k)} (p_{ajk} - p_{bjk})^2} \quad (18)$$

3.2.2 Provesti's absolute genetic distance

$$D_P(a, b) = \frac{1}{2\nu} \sum_{k=1}^{\nu} \sum_{j=1}^{m(k)} |p_{ajk} - p_{bjk}| \quad (19)$$

Note: for AFLP data, the 2 is dropped.

3.3 Bruvo's distance (stepwise mutation for microsatellites)

Bruvo's distance between two individuals calculates the minimum distance across all combinations of possible pairs of alleles at a single locus and then averaging that distance across all loci [3]. The distance between each pair of alleles is calculated as³[3]:

$$m_x = 2^{-|x|} \quad (20)$$

$$d_a = 1 - m_x \quad (21)$$

Where x is the number of steps between each allele. So, let's say we were comparing two haploid ($k = 1$) individuals with alleles 228 and 244 at a locus that had a tetranucleotide repeat pattern (CATG) ^{n} . The number of steps for each of these alleles would be $228/4 = 57$ and $244/4 = 61$, respectively. The number of steps between them is then $|57 - 61| = 4$. Bruvo's distance at this locus between these two individuals is then $1 - 2^{-4} = 0.9375$. For samples with higher ploidy (k), there would be k such distances of which we would need to take the sum [3].

$$s_i = \sum_{a=1}^k d_a \quad (22)$$

Unfortunately, it's not as simple as that since we do not assume to know phase. Because of this, we need to take all possible combinations of alleles into account. This means that we will have k^2 values of d_a , when we only want k . How do we know which k distances we want? We will have to invoke parsimony for this and attempt to take the minimum sum of the alleles, of which there are $k!$ possibilities [3]:

$$d_l = \frac{\left(\min_{i \dots k!} s_i \right)}{k} \quad (23)$$

Finally, after all of this, we can get the average distance over all loci [3].

$$D = \frac{\sum_{i=1}^l d_i}{l} \quad (24)$$

This is calculated over all possible combinations of individuals and results in a lower triangle distance matrix over all individuals.

³Notation presented unmodified from Bruvo et al, 2004

3.3.1 Special cases of Bruvo's distance

As shown in the above section, ploidy is irrelevant with respect to calculation of Bruvo's distance. However, since it makes a comparison between all alleles at a locus, it only makes sense that the two loci need to have the same ploidy level. Unfortunately for polyploids, it's often difficult to fully separate distinct alleles at each locus, so you end up with genotypes that appear to have a lower ploidy level than the organism [3].

To help deal with these situations, Bruvo has suggested three methods for dealing with these differences in ploidy levels [3]:

- Infinite Model - The simplest way to deal with it is to count all missing alleles as infinitely large so that the distance between it and anything else is 1. Aside from this being computationally simple, it will tend to inflate distances between individuals.
- Genome Addition Model - If it is suspected that the organism has gone through a recent genome expansion, the missing alleles will be replaced with all possible combinations of the observed alleles in the shorter genotype. For example, if there is a genotype of [69, 70, 0, 0] where 0 is a missing allele, the possible combinations are: [69, 70, 69, 69], [69, 70, 69, 70], and [69, 70, 70, 70]. The resulting distances are then averaged over the number of comparisons.
- Genome Loss Model - This is similar to the genome addition model, except that it assumes that there was a recent genome reduction event and uses the observed values in the full genotype to fill the missing values in the short genotype. As with the Genome Addition Model, the resulting distances are averaged over the number of comparisons.
- Combination Model - Combine and average the genome addition and loss models.

As mentioned above, the infinite model is biased, but it is not nearly as computationally intensive as either of the other models. The reason for this is that both of the addition and loss models requires replacement of alleles and recalculation of Bruvo's distance. The number of replacements required is equal to the multiset coefficient: $\binom{n}{k} = \binom{n-k+1}{k}$ where n is the number of potential replacements and k is the number of alleles to be replaced. So, for the example given above, The genome addition model would require $\binom{2}{2} = 3$ calculations of Bruvo's distance, whereas the genome loss model would require $\binom{4}{2} = 10$ calculations.

To reduce the number of calculations and assumptions otherwise, Bruvo's distance will be calculated using the largest observed ploidy. This means that when comparing [69,70,71,0] and [59,60,0,0], they will be treated as triploids.

3.4 Example

As these distances can affect data in different ways, it might be important to see what kind of trees they produce. To demonstrate, we will use 5 diploid samples at a single locus demonstrating a range of possibilities:

Genotype
1/1
1/2
2/3
3/4
4/4

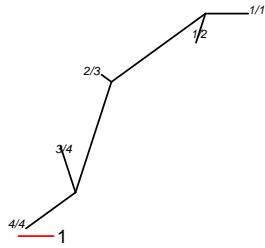
Table 2: Table of genotypes to be used for analysis

```
library(poppr)
dat.df <- data.frame(Genotype = c("1/1", "1/2", "2/3", "3/4", "4/4"))
dat <- as.genclone(df2genind(dat.df, sep = "/", ind.names = dat.df[[1]]))
```

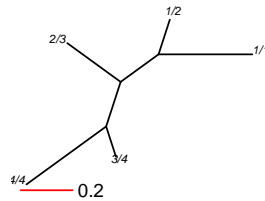
We will now compute the distances, construct a neighbor joining dendrogram with the package *ape*, and plot it.

```
distances <- c("Nei", "Rogers", "Edwards", "Reynolds", "Provesti")
dists <- lapply(distances, function(x) {
  DISTFUN <- match.fun(paste(tolower(x), "dist", sep = "."))
  DISTFUN(dat)
})
names(dists) <- distances
dists$Bruvo <- bruvo.dist(dat, replen = 1)
library(ape)
par(mfrow = c(2, 3))
x <- lapply(names(dists), function(x) {
  plot(nj(dists[[x]]), main = x, type = "unrooted")
  add.scale.bar(lcol = "red")
})
```

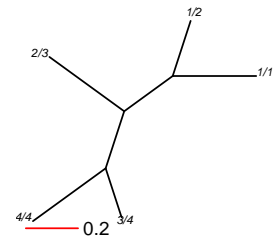
Nei



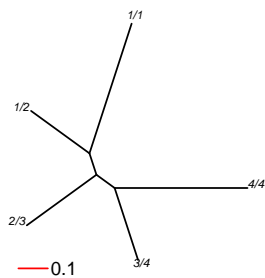
Rogers



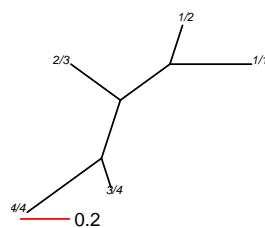
Edwards



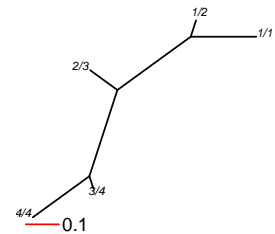
Reynolds



Provesti



Bruvo



References

- [1] Paul-Michael Agapow and Austin Burt. Indices of multilocus linkage disequilibrium. *Molecular Ecology Notes*, 1(1-2):101–102, 2001.
- [2] A.H.D. Brown, M.W. Feldman, and E. Nevo. Multilocus structure of natural populations of hordeum spontaneum. *Genetics*, 96(2):523–536, 1980.
- [3] Ruzica Bruvo, Nicolaas K. Michiels, Thomas G. D’Souza, and Hinrich Schulenburg. A simple method for the calculation of microsatellite genotype distances irrespective of ploidy level. *Molecular Ecology*, 13(7):2101–2106, 2004.
- [4] J M Smith, N H Smith, M O’Rourke, and B G Spratt. How clonal are bacteria? *Proceedings of the National Academy of Sciences*, 90(10):4384–4388, 1993.