

Documentation for EMALAM

CAROLA SOPHIA HEINZEL^(1,*), PETER PFAFFELHUBER⁽¹⁾

(1) Department of Mathematical Stochastics, Albert-Ludwigs-University Freiburg,
79104 Freiburg im Breisgau, Germany

(*) Corresponding author: Carola.Heinzel@stochastik.uni-freiburg.de

August 20, 2024

1 Introduction and Overview

To estimate the ancestry of individuals, we can use the Admixture Model and a maximum likelihood estimator. We proved that this estimator is usually not unique. EMALAM calculates some other maximum likelihood estimators (MLEs). The user determines which other estimators EMALAM calculates.

Here, we provide an explanation how to use the software EMALAM. Briefly, this software has as an input the estimated IAs and the estimated allele frequencies. To estimate them, we could e.g. use STRUCTURE [5]. Additionally, EMALAM requires some additional information from the reader that is specified in section 3. The output are MLEs, i.e. estimated individuals ancestries (IAs) and estimated allele frequencies. The user defines which MLEs EMALAM calculates, e.g. the MLEs with the maximal admixture.

First, we explain how EMALAM works. Then, we describe the decisions the user has to make to use EMALAM and the input of EMALAM. Finally, we describe the interpretation of the output and a possibility to depict the results.

2 Explanation how EMALAM works

We use `scipy.minimize` to minimize our objective function under some constraints. To explain the method, we need some notation. This is an often considered problem, see e.g. [4, 2, 3]. Here, we solved this problem with a new method, EMALAM.

We write $\sigma(\hat{q})$ for all permutations of the matrix \hat{q} and $\sigma(\hat{q})_k$ for the k th permutation of

23 this matrix. For example, let $\hat{q} = \begin{pmatrix} 0.2 & 0.3 & 0.5 \\ 0.9 & 0.1 & 0 \end{pmatrix}$, i.e. $N = 2, K = 3$. Then,

$$\sigma(\hat{q})_1 = \begin{pmatrix} 0.2 & 0.5 & 0.3 \\ 0.9 & 0 & 0.1 \end{pmatrix}$$

$$\sigma(\hat{q})_2 = \begin{pmatrix} 0.2 & 0.3 & 0.5 \\ 0.9 & 0.1 & 0 \end{pmatrix}$$

$$\sigma(\hat{q})_3 = \begin{pmatrix} 0.3 & 0.2 & 0.5 \\ 0.1 & 0.9 & 0 \end{pmatrix}$$

$$\sigma(\hat{q})_4 = \begin{pmatrix} 0.3 & 0.5 & 0.2 \\ 0.1 & 0 & 0.9 \end{pmatrix}$$

$$\sigma(\hat{q})_5 = \begin{pmatrix} 0.5 & 0.2 & 0.3 \\ 0 & 0.9 & 0.1 \end{pmatrix}$$

$$\sigma(\hat{q})_6 = \begin{pmatrix} 0.5 & 0.3 & 0.2 \\ 0 & 0.1 & 0.9 \end{pmatrix}.$$

24 Additionally, the matrix $S_K(a)$ is defined as in the corresponding paper. The a emphasizes
 25 that this matrix depends on the parameters $a = (a_\ell)_{\ell=1,\dots,2(K-1)}$. We define the matrix A and
 26 the vector b so that $Aa \leq b$ make sure that the estimated IAs are between 0 and 1 and that they
 27 sum up to one, i.e. $\sum_{k=1}^K \hat{q}_{i,k} = 1, \hat{q}_{i,k} \in [0, 1]$. Here, the vector a contains the same parameters
 28 as the matrix $S_K(a)$, but with an other format.

29 EMALAM minimizes or maximizes a function $f_{obj}(a)$, that the user chooses (see 3), with
 30 respect to a and the constraints

31 $C1 \quad Aa \leq b$

32 $C2 \quad S_K^{-1}(a)\hat{p} \in [0, 1]^{K \times J \times M}$

33 $C3 \quad \sum_{j=1}^J (S_K^{-1}(a)\hat{p})_{k,j,m} = 1 \text{ for } J \geq 3.$

34 $C4 \quad |\hat{q}S_K - \hat{q}| \leq |\hat{q}S_K(a) - \sigma(\hat{q})_k| \quad \forall k = 1, \dots, K.$

35 Here, condition C1 makes sure that we can interpret $\hat{q}S_K(a)$ as the IAs, i.e.

$$\hat{q}S_K(a) \geq 0, \sum_{k=1}^K (\hat{q}_i, S_K(a))_k = 1.$$

36 Additionally, the conditions C2 and C3 make sure that the estimated allele frequencies are
 37 between 0 and 1 ($S_K^{-1}(a)\hat{p} \in [0, 1]^{K \times J \times M}$) and that it holds $\sum_{j=1}^J p_{k,j,m} = 1$. Condition C4 is
 38 optional and consequences that the output of EMALAM does not belong to label switching, if
 39 we consider the most similar IAs. However, please note that the running of EMALAM with
 40 condition C4 takes much longer than without C4. Hence, we recommend to use C4 only for
 41 small number of individuals and small number of markers.

3 Application Decisions for the User

information) There are two different decisions that the user can make: The choice for the function that should be maximized and the choice for the measure of similarity to take label switching into account. We explain the choice concerning the function that should be maximized first.

3.1 Choice of the Target Function

Researchers can choose five different objective functions $f_{obj}(a)$. We discuss them in more detail now.

- (I) Maximize and minimize the estimated IA for individual ind in the input data of the estimated IAs in every population (`poss` = "P1" in the code). In this case the user also has to specify the index of the individual. Here, we define

$$f_{obj}(a) = \tilde{q}_{ind,k}.$$

We minimize and maximize every the estimated ancestry of every population for this individual ind .

- (II) Maximize the admixture of the estimated IAs (`poss` = "P2" in the code), i.e. we maximize the entropy

$$f_{obj}(a) = - \sum_{i=1}^N \sum_{k=1}^K \tilde{q}_{i,k} \ln(\tilde{q}_{i,k}).$$

- (III) Minimize the admixture of the estimated IAs (`poss` = "P3" in the code). We minimize the entropy. This consequences the minimal admixture.

- (IV) Maximize the ancestries for a specific population $k_{specific}$ (`poss` = "P4" in the code), i.e. we minimize

$$f_{obj}(a) = - \sum_{i=1}^N \tilde{q}_{i,k_{specific}}.$$

- (V) Minimize the ancestries for a specific population $k_{specific}$ (`poss` = "P5" in the code). Specifically, we minimize

$$f_{obj}(a) = \sum_{i=1}^N \tilde{q}_{i,k_{specific}}.$$

To summarize, we can apply EMALAM to single individuals or to whole populations as described above to find the most extreme optima.

3.2 Choice of the Definition for Label Switching

Additionally, EMALAM can take label switching into account. Therefore, there are different possibilities to define the similarity of the different MLEs. Here, we use the euclidean norm. This means that we choose the labels for the populations in order to minimize the euclidean norm between the different estimators for the IA. Let us consider a simple example for this.

Let $\hat{q}_{1,1} = 0.4, \hat{q}_{1,2} = 0.6, \hat{p}_{1,1,1} = 0.9, \hat{p}_{2,1,1} = 0.2$ be the output of STRUCTURE for $K = 2, M = 1, N = 1$ (Figure 1). Furthermore, we have the output $\tilde{q}_{1,2} = 0.7, \tilde{q}_{1,1} = 0.3, \tilde{p}_{1,1,1} = 0.6, \tilde{p}_{2,1,1} = 0.2$ for an other run of STRUCTURE. Now, we consider two possibilities to depict these opportunities, i.e to avoid label switching:

- (i) Minimize the differences between the allele frequencies (second column in Figure 1).
- (ii) Minimize the differences between the IAs (first column in Figure 1).

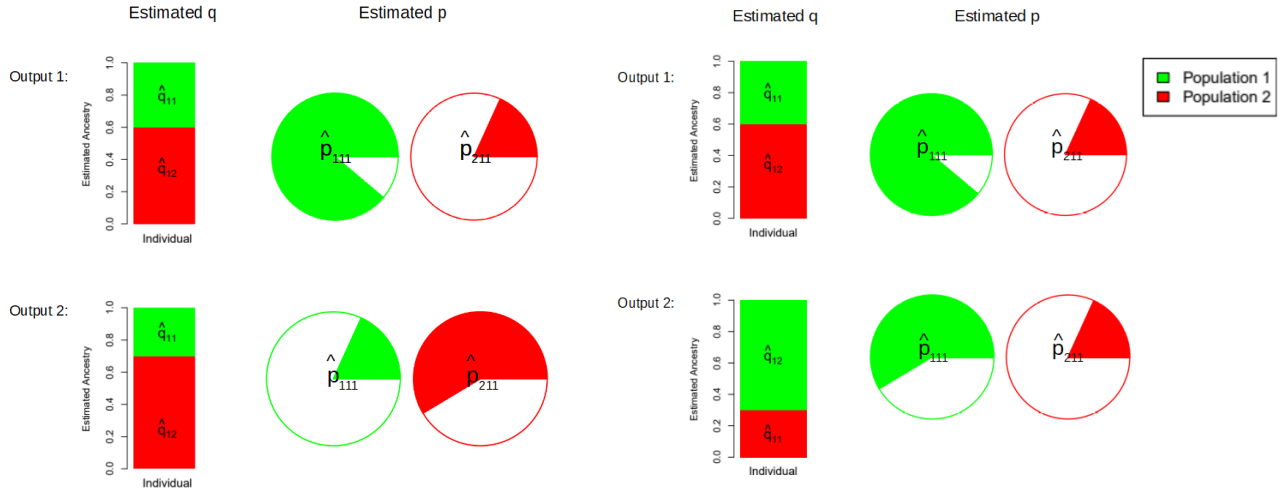


Figure 1: Example Label Switching

EMALAM uses possibility (ii). Specifically, EMALAM uses the assignment of population labels to the estimated IAs and allele frequencies with the smallest euclidean norm between the estimated IAs for the $K!$ different possibilities.

However, it is easy to change this in the function `constraint4` or otherwise, apply pong to take label switching into account.

4 Input of EMALAM

There is some information that the user has to provide to use EMALAM. These information are

- (i) `file_path_q`: Estimated IAs. We provide an example for the input file, called `q-CEU_IBS_TSI_K2`. Specifically, for the estimated IAs, the file has to contain N rows and K columns. The rows represent the individuals and the columns represent the populations. The required format can e.g. be created by applying the code in `Extract_q-p.R`

- 87 (ii) **file_path_p**: Estimated allele frequencies. The input file has K columns and the number
 88 of columns depend on M and J_m . For marker m , we only exclude $J_m - 1$ estimated allele
 89 frequencies from the data. The reason for this is that we know the allele frequencies
 90 for the other alleles as $\sum_{j=1}^J p_{k,j,m} = 1$ has to hold. The example file for this is called
 91 **p_CEU_IBS_TSI_K2**. The required format can e.g. be created by applying the code in
 92 **Extract_q_p.R** to the output of STRUCTURE for $J = 2$. For J arbitrary, they can be
 93 extracted with the code **Extract_P_J_arbitrary.R**. Please note that we can exclude the
 94 allele frequencies that are the same in every population from the input for EMALAM.
 95 Our code does this, if the allele 0 has either the frequency 0 or 1.
- 96 (iii) **simi**: Either 1 (takes label switching into account) or 0 (does not take label switching
 97 into account, i.e. we have to apply a software as e.g. pong [1] to the output of EMALAM
 98 afterwards). The details are described in condition C4. If **simi** = 1, C4 is an additional
 99 constraint of the minimization problem. However, this is only recommended for small
 100 number of individuals and markers since the run time is very long for **simi** = 1.
- 101 (iv) **file_path_pJ**: The sum of the estimated allele frequencies per marker. For marker
 102 with more than two alleles, we also calculate the sum of the allele frequencies. The
 103 file **p_CEU_IBS_TSI_K2_J** contains an example for this. We describe the details below.
- 104 (v) **poss**: Definition of the function $f_{opt}(a)$, i.e. of the most extreme values as described
 105 above. Here, the user chooses between (I), (II), (III), (IV) or (V). Alternatively, you can
 106 also define the function by yourself.
- 107 (vi) **names**: The names of the output file (including the directory where they should be saved).
 108 The first entry is the name for the ancestries and the second entry stands for the allele
 109 frequencies. If the user chooses **poss** = "P1", then the output is $2K$ different files, named
 110 **names[j]_i**, $i = 1, \dots, 2K$, $j = 0, 1$.
- 111 (vii) $k_{specific}$: Population that is considered in P4 or P5.
- 112 (viii) n_{trial} : The number of different initial values for **scipy.minimize** that are used. The
 113 default value is 10. It might occur in rare cases that the function **scipy.minimize** does
 114 not convergence to an optimal point. In this case, please try other initial values.

115 Let us consider (iv) in more detail. For just bi-allelic marker, we do not need this input.
 116 Otherwise, the Input format is a .txt file with K columns. It can be created by applying
 117 **Extract_P_J_arbitrary.R** to the output of STRUCTURE. Let us consider an example: The

118 estimated allele frequencies in STRUCTURE are

```

Locus1 :
2alleles
0.0% missing data
1(0.244)0.2080.103
0(0.756)0.7920.897
Locus2 :
3alleles
0.0% missing data
1(0.7)0.8320.650
0(0.267)0.1680.350
2(0.023)0.030 0.036.

```

Then, the output would be

0.168 + 0.030 0.350 +0.036

Table 1: Example Input for EMALAM

119

120 5 Interpretation of the Output

121 We also provide an example output. The files are called
122 q_K2_P1_0.txt (for the IAs),
123 p_K2_P1_0.txt (for the allele frequencies),
124 and have the same format as the input. Additionally, the likelihood is also the same as the
125 one for the input estimators. However, this output is the most extreme one in the sense that
126 the user chose. Please note that there might occur slightly negative values or values that are
127 slightly bigger than 1. To avoid this, the user can choose an other threshold for \tilde{q}, \tilde{p} , e.g.
128 $\tilde{q} \in [0.001, 0.999]$.

129 6 Depiction of the Results

130 You can use the code Create_Figures.py to depict the different estimated IAs. On the x-axis,
131 you have the individuals and on the y-axis, you have the estimated IAs. Figure 2 is an example
132 for the depiction of the results.

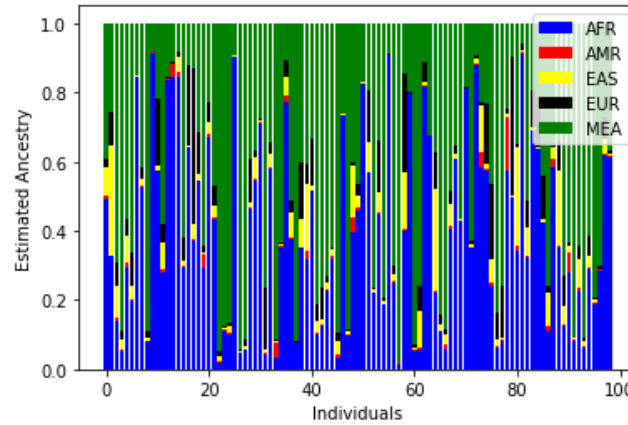


Figure 2: Example Depiction of the estimated IAs

(Fig:ex)

133 The list

$$data = [[0.326, 0.004, 0.318, 0.097, 0.254],$$

$$[0.14, 0.008, 0.094, 0.065, 0.693],$$

$$[0.053, 0.004, 0.027, 0.027, 0.889],$$

$$[0.296, 0.01, 0.078, 0.05, 0.567]]$$

134 is an example for the input of this code. However, we present an other example with more
 135 individuals in Figure 2.

136 References

- [behr2016](#) [1] Aaron A Behr, Katherine Z Liu, Gracie Liu-Fang, Priyanka Nakka, and Sohini Ramachandran. Pong: fast analysis and visualization of latent clusters in population genetic data. *Bioinformatics*, 32(18):2817–2823, 2016.
- [kraft1988](#) [2] Dieter Kraft. A software package for sequential quadratic programming. *Forschungsbericht-Deutsche Forschungs- und Versuchsanstalt für Luft- und Raumfahrt*, 1988.
- [lalee1998](#) [3] Marucha Lalee, Jorge Nocedal, and Todd Plantenga. On the implementation of an algorithm for large-scale equality constrained optimization. *SIAM Journal on Optimization*, 8(3):682–706, 1998.
- [powell1994](#) [4] Michael JD Powell. *A direct search optimization method that models the objective and constraint functions by linear interpolation*. Springer, 1994.
- [pritchard2000](#) [5] Jonathan K Pritchard, Matthew Stephens, and Peter Donnelly. Inference of population structure using multilocus genotype data. *Genetics*, 155(2):945–959, 2000.