

# ALStructure Manual

February 18, 2018

---

alstructure

*Main function for execution of the ALStructure algorithm*

---

## Description

Computes global ancestry estimates under the admixture model given a SNP data matrix  $\mathbf{X}$ . This function is based on the ALStructure algorithm from (Cabreros and Storey 2017).

## Usage

```
alstructure(X, d_hat = NULL, svd_method = "base", tol = 1e-05,  
            max_iters = 1000, order_method = "ave_admixture", P_init, Q_init)
```

## Arguments

<code>X</code>	The $m \times n$ SNP data matrix.
<code>d_hat</code>	Estimate of the latent space dimension $d$ . If left blank, this is estimated by the function <code>estimate_d()</code>
<code>svd_method</code>	One of "vanilla" or "trunc." If "vanilla" is chosen, the base <code>svd()</code> function is used. If "trunc" is used, the truncated svd algorithm from the <code>lfa</code> package is used.
<code>tol</code>	The convergence criterion. If $RMSE(\hat{\mathbf{Q}}_t - \hat{\mathbf{Q}}_{t+1}) < tol$ , then the algorithm halts
<code>max_iters</code>	The maximum number of iterations (repetitions of steps (6) and (7) in Algorithm 1) to be executed
<code>order_method</code>	One of "ave_admixture" or "var_explained." If "ave_admixture," the $d$ populations are ordered by decreasing average admixture accross samples (i.e. $1/n \sum_j q_{ij}$ ). If "var_explained", the $d$ populations are ordered by decreasing variation explained. Specifically, we compute a modified version of the eigen- $R^2$ statistic from (L. S. Chen and Storey 2008). The statistic is modified in the following ways: 1) we treat rows of $\mathbf{Q}$ as the response variables 2) we regress each row of $\mathbf{Q}$ on the eigenvectors of $\mathbf{G}$ rather than the eigenvectors of the data matrix itself 3) we take the weighted average only over the top $d$ eigenvectors. and columns of $\mathbf{P}$ are ordered by amount of variation explained by each row of $\mathbf{Q}$ by the function <code>order_Q</code>
<code>P_init</code>	Optional initialization of $\mathbf{P}$ . Only available for cALS method.
<code>Q_init</code>	Optional initialization of $\mathbf{Q}$ . Only available for cALS method.

**Value**

A list with the following elements:

**P\_hat** : The estimated  $\mathbf{P}$  matrix. Each column of  $\mathbf{P}$  may be interpreted as a vector of allele frequencies for a specific ancestral population.

**Q\_hat** : The estimated  $\mathbf{Q}$  matrix. Each column of  $\mathbf{Q}$  may be interpreted as the admixture proportions of a specific individual.

**rowspace** : a list with the following elements:

**vectors** : The top  $d$  eigenvectors of the matrix  $\mathbf{G}$  sorted by decreasing eigenvalue. These vectors approximate the subspace spanned by the rows of  $\mathbf{Q}$ .

**values** : The top  $d$  eigenvalues of the matrix  $\mathbf{G}$  sorted by decreasing eigenvalue.

**References**

Cabreros, I., and J. D. Storey. 2017. “A Nonparametric Estimator of Population Structure Unifying Admixture Models and Principal Components Analysis.” *BioRxiv*. Cold Spring Harbor Laboratory. doi:10.1101/240812.

Hao, W., M. Song and J. D. Storey. 2015. “lfa: Logistic Factor Analysis for Categorical Data.” R package version 1.8.0, <https://github.com/StoreyLab/lfa>.

---

estimate\_d

*Estimate the latent space dimension*

---

**Description**

Estimates the dimension of the rowspace of  $\mathbf{Q}$  (equivalently, the rank of the matrix  $\mathbf{F}$ ). This estimate  $\hat{d}$  is based on the estimator from (Leek 2011), page 6.

**Usage**

estimate\_d(X)

**Arguments**

X                      the  $m \times n$  SNP data matrix

**Value**

an estimate  $\hat{d}$  of the dimension of the latent space dimension.

**References**

Leek, J. T. 2011. “Asymptotic conditional singular value decomposition for high-dimensional genomic data.” *Biometrics* 67 (4): 344–52.

---

estimate_F	<i>Estimates the individual-specific allele frequency matrix <math>F</math></i>
------------	---

---

## Description

Estimates the  $m \times n$  individual-specific allele frequency matrix  $F$  using the method of Latent Subspace Estimation (X. Chen and Storey 2015) as described in (Cabreros and Storey 2017) in Section 2.3.

## Usage

```
estimate_F(X, d, svd_method = "base")
```

## Arguments

<code>X</code>	The $m \times n$ SNP data matrix
<code>d</code>	The rank of $F$ . This can be estimated using the function <code>d_estimate()</code> .
<code>svd_method</code>	One of "base" or "truncated_svd." If "base" is chosen, the base <code>svd()</code> function is used. If "truncated_svd" is used, the truncated svd algorithm from the <code>lfa</code> package is used.

## Value

A list with the following elements:

**F\_hat** : The  $m \times n$  matrix  $\hat{F}$ .

**rowspace** : a list with the following elements:

**vectors** : The top  $d$  eigenvectors of the matrix  $G$  sorted by decreasing eigenvalue. These vectors approximate the subspace spanned by the rows of  $Q$ .

**values** : The top  $d$  eigenvalues of the matrix  $G$  sorted by decreasing eigenvalue.

## References

Cabreros, I., and J. D. Storey. 2017. "A Nonparametric Estimator of Population Structure Unifying Admixture Models and Principal Components Analysis." *BioRxiv*. Cold Spring Harbor Laboratory. doi:10.1101/240812.

Chen, X., and J. D. Storey. 2015. "Consistent Estimation of Low-Dimensional Latent Structure in High-Dimensional Data." *ArXiv E-Prints*, October.

---

factor_F	<i>A fast algorithm for factoring <math>\hat{\mathbf{F}}</math>.</i>
----------	--

---

### Description

An algorithm for finding approximate factors  $\hat{\mathbf{P}}$  and  $\hat{\mathbf{Q}}$  of the matrix  $\hat{\mathbf{F}}$  that obey constraints of the admixture model:

1.  $p_{ij} \in [0, 1] \forall (i, j)$
2.  $q_{ij} \geq 0 \forall (i, j)$  and  $\sum_i q_{ij} = 1 \forall j$

This algorithm is described in Algorithm 2 in (Cabreros and Storey 2017). While it lacks the theoretical guarantees of the cALS algorithm, it is much faster.

### Usage

```
factor_F(F_hat, d, tol = 1e-05, max_iters = 1000)
```

### Arguments

F_hat	The estimate of the matrix $\mathbf{F}$ to be factored
d	The dimension of the latent space. This can be estimated by the function <code>d_estimate</code> .
tol	The convergence criterion. If $RMSE(\hat{\mathbf{Q}}_t - \hat{\mathbf{Q}}_{t+1}) < tol$ , then the algorithm halts
max_iters	The maximum number of iterations (repetitions of steps (6) and (7) in Algorithm 1) to be executed

### Value

A list with the following elements:

**P\_hat** : The estimated  $\mathbf{P}$  matrix. Each column of  $\mathbf{P}$  may be interpreted as a vector of allele frequencies for a specific ancestral population.

**Q\_hat** : The estimated  $\mathbf{Q}$  matrix. Each column of  $\mathbf{Q}$  may be interpreted as the admixture proportions of a specific individual.

@references Cabreros, I., and J. D. Storey. 2017. "A Nonparametric Estimator of Population Structure Unifying Admixture Models and Principal Components Analysis." *BioRxiv*. Cold Spring Harbor Laboratory. doi:10.1101/240812.

---

lse	<i>Estimates the latent subspace</i>
-----	--------------------------------------

---

### Description

Estimates the rowspace of  $\mathbf{Q}$  using the method of latent subspace estimation. The function returns the top  $d$  eigenvalues and vectors of the matrix

$$\mathbf{G} = \frac{1}{m} \mathbf{X}^T \mathbf{X} - \mathbf{D}$$

where the matrix  $\mathbf{D}$  is a diagonal matrix with each diagonal entry  $d_{ii}$  an estimate of the average of the variances of the random variables in the  $i$  column of  $\mathbf{X}$ . As is proven in (X. Chen and Storey 2015), the span of the top  $d$  eigenvectors of  $\mathbf{G}$  span the same space as the rows of  $\mathbf{Q}$ . The eigenvectors are returned in order of decreasing eigenvalue.

### Usage

```
lse(X, d, svd_method = "base")
```

### Arguments

$\mathbf{X}$	The $m \times n$ SNP data matrix
$d$	The rank of $\mathbf{F}$ . This can be estimated using the function <code>d_estimate()</code> . When $d = n$ , all eigenvectors of $\mathbf{G}$ are returned.
<code>svd_method</code>	One of "base" or "truncated_svd." If "base" is chosen, the base <code>svd()</code> function is used. If "truncated_svd" is used, the truncated svd algorithm from the <code>lfa</code> package is used.

### Value

A list with the following elements:

**vectors** : The top  $d$  eigenvectors of the matrix  $\mathbf{G}$  sorted by decreasing eigenvalue. These vectors approximate the subspace spanned by the rows of  $\mathbf{Q}$ .

**values** : The top  $d$  eigenvalues of the matrix  $\mathbf{G}$  sorted by decreasing eigenvalue.

### References

Chen, X., and J. D. Storey. 2015. "Consistent Estimation of Low-Dimensional Latent Structure in High-Dimensional Data." ArXiv E-Prints, October.

---

order_pops	<i>Orders the <math>d</math> populations</i>
------------	--

---

### Description

Orders the  $d$  populations according to one of two methods: "ave\_admixture" or "var\_explained." Function returns matrices  $P$  and  $Q$  with permuted columns and rows according to the determined ordering of populations.

### Usage

```
order_pops(P, Q, method = "ave_admixture", Q_space = NULL)
```

### Arguments

$P$	The $m \times d$ loadings matrix with columns to be ordered
$Q$	The $d \times n$ admixture matrix with rows to be ordered
method	One of "ave_admixture" or "var_explained." If "ave_admixture," the $d$ populations are ordered by decreasing average admixture accross samples (i.e. $1/n \sum_j q_{ij}$ ). If "var_explained", the $d$ populations are ordered by decreasing variation explained. Specifically, we compute a modified version of the eigen- $R^2$ statistic from (L. S. Chen and Storey 2008). The statistic is modified in the following ways: 1) we treat rows of $Q$ as the response variables 2) we regress each row of $Q$ on the eigenvectors of $G$ rather than the eigenvectors of the data matrix itself 3) we take the weighted average only over the top $d$ eigenvectors.
Q_space	Only required for "var_explained" methodThe list containing the top $d$ eigenvectors and their corresponding eigenvalues of the $G$ . These may be obtained through the function lse.

### Value

A list with the following elements:

**P\_ordered** : The permuted  $P$

**Q\_ordered** : The permuted  $Q$

**perm\_mat** : The permutation matrix  $A$  such that  $PA^T = P_{\text{ord}}$  and  $AQ = Q_{\text{ord}}$ .

### References

Chen, L. S., and J. D. Storey. 2008. "Eigen-R2 for dissecting variation in high-dimensional studies." *Bioinformatics* 24 (19): 2260–2.

---

simulate_admixture	<i>Simulates data from the PSD model</i>
--------------------	--

---

### Description

Creates a data frame that contains the parameters of the admixture model ( $\mathbf{F}$ ,  $\mathbf{P}$ ,  $\mathbf{Q}$ ) as well as a single draw  $\mathbf{X}$  such that  $x_{ij} \sim \text{Bernoulli}(f_{ij})$ . The  $\mathbf{Q}$  matrix is drawn from the Dirichlet distribution with parameter  $\alpha$ , and the  $\mathbf{P}$  matrix is simulated from the Balding-Nichols model. The parameter  $\alpha$  is supplied by the user. The  $m \times 2$  matrix of Balding-Nichols parameters is optional. If BN\_params is not supplied, the parameters are derived from a random sample of estimated Balding-Nichols parameters from the Human Genomes Diversity Project (HGDP) dataset. The Balding-Nichols parameter estimates are provided by (Gopalan et al. 2016), and included in this package in the object hgdpBN.

### Usage

```
simulate_admixture(m, n, d, alpha, BN_params = NA, seed = NA)
```

### Arguments

m	number of SNPs
n	number of individuals
d	number of groups
alpha	dirichlet parameter; $\text{length}(\alpha) = d$
BN_params	a $m \times 2$ matrix of parameters. The first column contains $F_{ST}$ for each SNP, while the second column contains the allele frequency.

### Value

a list with the following elements:

$\mathbf{P}$  : the  $m \times d$  matrix of loadings  
 $\mathbf{Q}$  : the  $d \times n$  matrix of latent admixture components  
 $\mathbf{F}$  : the  $m \times n$  matrix  $PQ$   
 $\mathbf{X}$  : a random draw such that  $x_{ij} \sim \text{Bernoulli}(f_{ij}, 2)$ .

### References

Gopalan, P., W. Hao, D. M. Blei, and J. D. Storey. 2016. "Scaling probabilistic models of genetic variation to millions of humans." *Nature Genetics* 48 (12): 1587–90.

# Index

`alstructure`, [1](#)

`estimate_d`, [2](#)

`estimate_F`, [3](#)

`factor_F`, [4](#)

`lse`, [5](#)

`order_pops`, [6](#)

`simulate_admixture`, [7](#)