Identifying groups 00000 000000 Exploring group diversity 000 0000 000000

Multivariate analysis of genetic data — exploring group diversity —

Thibaut Jombart, Marie-Pauline Beugin

MRC Centre for Outbreak Analysis and Modelling Imperial College London

Genetic data analysis with PR~Statistics, Millport Field Station 18 Aug 2016

Exploring group diversity 000 0000 000000

Outline

Introduction

Identifying groups Hierarchical clustering K-means

Exploring group diversity

Aggregating data Optimizing group differences Discriminant Analysis of Principal Components

Exploring group diversity 000 0000 000000

Outline

Introduction

Identifying groups Hierarchical clustering K-means

Exploring group diversity

Aggregating data Optimizing group differences Discriminant Analysis of Principal Components

Exploring group diversity 000 0000 000000

Genetic data: introducing group data



• How to identify groups?

• How to explore group diversity?

Exploring group diversity 000 0000 000000

Genetic data: introducing group data



- How to identify groups?
- How to explore group diversity?

Identifying groups

Exploring group diversity 000 0000 000000

Outline

Introduction

Identifying groups Hierarchical clustering K-means

Exploring group diversity

Aggregating data Optimizing group differences Discriminant Analysis of Principal Components

Exploring group diversity 000 0000 000000

Hierarchical clustering: a variety of algorithms

- single linkage
- complete linkage
- UPGMA
- Ward
- ...

Exploring group diversity 000 0000 000000

Rationale

1. compute pairwise genetic distances D (or similarities)

- 2. group the closest pair(s) together
- 3. (optional) update ${\bf D}$
- 4. return to 2) until no new group can be made

Exploring group diversity 000 0000 000000

- 1. compute pairwise genetic distances D (or similarities)
- 2. group the closest pair(s) together
- 3. (optional) update \mathbf{D}
- 4. return to 2) until no new group can be made

Exploring group diversity 000 0000 000000

- 1. compute pairwise genetic distances D (or similarities)
- 2. group the closest pair(s) together
- 3. (optional) update ${\bf D}$
- 4. return to 2) until no new group can be made

Exploring group diversity 000 0000 000000

- 1. compute pairwise genetic distances D (or similarities)
- 2. group the closest pair(s) together
- 3. (optional) update \mathbf{D}
- 4. return to 2) until no new group can be made

Identifying groups

Exploring group diversity 000 0000 00000



Identifying groups

Exploring group diversity 000 0000 000000



- single linkage: $D_{k,g} = \min(D_{k,i}, D_{k,j})$
- complete linkage: $D_{k,g} = \max(D_{k,i}, D_{k,j})$

• UPGMA:
$$D_{k,g} = \frac{D_{k,i} + D_{k,j}}{2}$$

Exploring group diversity 000 0000 000000



- single linkage: $D_{k,g} = \min(D_{k,i}, D_{k,j})$
- complete linkage: $D_{k,g} = \max(D_{k,i}, D_{k,j})$

• UPGMA:
$$D_{k,g} = \frac{D_{k,i} + D_{k,j}}{2}$$

Identifying groups

Exploring group diversity 000 0000 000000



- single linkage: $D_{k,g} = \min(D_{k,i}, D_{k,j})$
- complete linkage: $D_{k,g} = \max(D_{k,i}, D_{k,j})$

• UPGMA:
$$D_{k,g} = \frac{D_{k,i} + D_{k,j}}{2}$$

1

Exploring group diversity 000 0000 000000



- single linkage: $D_{k,g} = \min(D_{k,i}, D_{k,j})$
- complete linkage: $D_{k,g} = \max(D_{k,i}, D_{k,j})$

• UPGMA:
$$D_{k,g} = \frac{D_{k,i} + D_{k,j}}{2}$$

Exploring group diversity 000 0000 000000



Exploring group diversity 000 0000 000000

K-means underlying model

ANOVA model:

total var. = (var. between groups) + (var. within groups)



Exploring group diversity 000 0000 000000

K-means rationale

Find groups which minimize *within group var*. (equally: maximize *between group var*.).

In other words:

Identify a partition $\mathcal{G} = \{g_1, \ldots, g_k\}$ solving:

$$\arg\min_{\mathcal{G}=\{g_1,\ldots,g_k\}}\sum_k\sum_{i\in g_k}\|\mathbf{x}_i-\boldsymbol{\mu}_k\|^2$$

with:

- $\mathbf{x}_i \in \mathbb{R}^p$: vector of allele frequencies of individual i
- $\mu_k \in \mathbb{R}^p$: vector of means allele frequencies of group k

Exploring group diversity 000 0000 000000

K-means rationale

Find groups which minimize *within group var.* (equally: maximize *between group var.*).

In other words:

Identify a partition $\mathcal{G} = \{g_1, \ldots, g_k\}$ solving:

$$\arg\min_{\mathcal{G}=\{g_1,\ldots,g_k\}}\sum_k\sum_{i\in g_k}\|\mathbf{x}_i-\boldsymbol{\mu}_k\|^2$$

with:

- $\mathbf{x}_i \in \mathbb{R}^p$: vector of allele frequencies of individual i
- $\mu_k \in \mathbb{R}^p$: vector of means allele frequencies of group k

Exploring group diversity 000 0000 000000

K-means rationale

Find groups which minimize *within group var.* (equally: maximize *between group var.*).

In other words:

Identify a partition $\mathcal{G} = \{g_1, \ldots, g_k\}$ solving:

$$\arg\min_{\mathcal{G}=\{g_1,\ldots,g_k\}}\sum_k\sum_{i\in g_k}\|\mathbf{x}_i-\boldsymbol{\mu}_k\|^2$$

with:

- $\mathbf{x}_i \in \mathbb{R}^p$: vector of allele frequencies of individual i
- $\mu_k \in \mathbb{R}^p$: vector of means allele frequencies of group k

Exploring group diversity 000 0000 000000

K-means algorithm

- 1. select random group means (μ_k , $k=1,\ldots,K$)
- 2. assign each individual \mathbf{x}_i to the closest group $\longrightarrow g_k$
- 3. update group means μ_k
- 4. go back to 2) until convergence (groups no longer change)

Exploring group diversity 000 0000 000000

K-means algorithm

- 1. select random group means (μ_k , $k=1,\ldots,K$)
- 2. assign each individual \mathbf{x}_i to the closest group $\longrightarrow g_k$
- 3. update group means μ_k
- 4. go back to 2) until convergence (groups no longer change)

Exploring group diversity 000 0000 000000

K-means algorithm

- 1. select random group means (μ_k , $k=1,\ldots,K$)
- 2. assign each individual \mathbf{x}_i to the closest group $\longrightarrow g_k$
- 3. update group means μ_k
- 4. go back to 2) until convergence (groups no longer change)

Exploring group diversity 000 0000 000000

K-means algorithm

- 1. select random group means (μ_k , $k = 1, \ldots, K$)
- 2. assign each individual \mathbf{x}_i to the closest group $\longrightarrow g_k$
- 3. update group means μ_k
- 4. go back to 2) until convergence (groups no longer change)

Exploring group diversity 000 0000 00000

K-means algorithm



Exploring group diversity 000 0000 000000

K-means: limitations and extensions

Limitations

- slower for large numbers of alleles (e.g. 100,000)
- K-means does not identify the number of clusters (K)

Extension

- run K-means after dimension reduction using PCA
- try increasing values of K
- use Bayesian Information Criterion (BIC) for model selection

Exploring group diversity 000 0000 000000

K-means: limitations and extensions

Limitations

- slower for large numbers of alleles (e.g. 100,000)
- K-means does not identify the number of clusters (K)

Extension

- run K-means after dimension reduction using PCA
- try increasing values of K
- use Bayesian Information Criterion (BIC) for model selection

Exploring group diversity 000 0000 000000

Genetic clustering using K-means & BIC

(Jombart et al. 2010, BMC Genetics)

Simulated data: island model with 6 populations



Performances:

- K-means ≥ STRUCTURE on simulated data (various island and stepping stone models)
- orders of magnitude faster (seconds vs hours/days)

Exploring group diversity 000 0000 000000

Genetic clustering using K-means & BIC

(Jombart et al. 2010, BMC Genetics)

Simulated data: island model with 6 populations



Performances:

- K-means ≥ STRUCTURE on simulated data (various island and stepping stone models)
- orders of magnitude faster (seconds vs hours/days)

Identifying groups 00000 000000 Exploring group diversity

Outline

Introduction

Identifying groups Hierarchical clustering K-means

Exploring group diversity

Aggregating data Optimizing group differences Discriminant Analysis of Principal Components

Exploring group diversity ••• ••• •••• ••••

Why identifying clusters is not the whole story Example of cattle breeds diversity (30 microsatellites, 704 individuals).

Group membership probabilities:



Important to assess the relationships between clusters.

Exploring group diversity ••• ••• •••• ••••

Why identifying clusters is not the whole story Example of cattle breeds diversity (30 microsatellites, 704 individuals).

Group membership probabilities:



Multivariate analysis:



Important to assess the relationships between clusters.

Exploring group diversity ••• ••• •••• ••••

Why identifying clusters is not the whole story Example of cattle breeds diversity (30 microsatellites, 704 individuals).

Group membership probabilities:



Multivariate analysis:



Important to assess the relationships between clusters.

Identifying groups 00000 000000 Exploring group diversity OOO OOO OOOOOO

Aggregating data by groups



 \longrightarrow multivariate analysis of group allele frequencies.

Exploring group diversity

Analysing group data

Available methods:

- Principal Component Analysis (PCA) of allele frequency table
- Genetic distance between populations → Principal Coordinates Analysis (PCoA)
- Correspondance Analysis (CA) of allele counts

Criticism:

- Lose individual information
- Neglect within-group diversity
- CA: possible artefactual outliers

Exploring group diversity

Analysing group data

Available methods:

- Principal Component Analysis (PCA) of allele frequency table
- Genetic distance between populations → Principal Coordinates Analysis (PCoA)
- Correspondance Analysis (CA) of allele counts

Criticism:

- Lose individual information
- Neglect within-group diversity
- CA: possible artefactual outliers

Identifying groups 00000 000000 Exploring group diversity

Multivariate analysis: reminder



Identifying groups 00000 000000 Exploring group diversity

Multivariate analysis: reminder



Identifying groups 00000 000000 Exploring group diversity

Multivariate analysis: reminder



Identifying groups 00000 000000 Exploring group diversity

Multivariate analysis: reminder



Exploring group diversity

But total variance may not reflect group differences



Need to optimize different criteria.

Exploring group diversity

But total variance may not reflect group differences



Need to optimize different criteria.

Exploring group diversity

But total variance may not reflect group differences



Need to optimize different criteria.

Exploring group diversity

Optimizing different criteria

Similar approaches to PCA can be used to optimize different quantities:

• PCA: total variance

- Between-group PCA: variance between groups
- Within-group PCA: variance within groups
- **Discriminant Analysis**: variance *between* groups / variance *within* groups

Exploring group diversity

Optimizing different criteria

Similar approaches to PCA can be used to optimize different quantities:

- PCA: total variance
- Between-group PCA: variance between groups
- Within-group PCA: variance within groups
- **Discriminant Analysis**: variance *between* groups / variance *within* groups

Exploring group diversity

Optimizing different criteria

Similar approaches to PCA can be used to optimize different quantities:

- PCA: total variance
- Between-group PCA: variance between groups
- Within-group PCA: variance within groups
- **Discriminant Analysis**: variance *between* groups / variance *within* groups

Exploring group diversity

Optimizing different criteria

Similar approaches to PCA can be used to optimize different quantities:

- PCA: total variance
- Between-group PCA: variance *between* groups
- Within-group PCA: variance within groups
- **Discriminant Analysis**: variance *between* groups / variance *within* groups

Exploring group diversity

From PCA to DA: increasing group differentiation



Discriminant Analysis: limitations and extensions

Limitations:

- DA requires less variables (alleles) than observations (individuals)
- DA requires uncorrelated variables (no frequencies, no linkage disequilibrium)

Discriminant Analysis of Principal Components (DAPC)¹:

- data orthogonalisation/reduction using PCA before DA
- overcomes limitations of DA
- group membership probabilities, group prediction

¹ Jombart et al. 2010, *BMC Genetics*

Discriminant Analysis: limitations and extensions

Limitations:

- DA requires less variables (alleles) than observations (individuals)
- DA requires uncorrelated variables (no frequencies, no linkage disequilibrium)

Discriminant Analysis of Principal Components (DAPC)¹:

- data orthogonalisation/reduction using PCA before DA
- overcomes limitations of DA
- group membership probabilities, group prediction

¹ Jombart et al. 2010, BMC Genetics

Identifying groups 00000 000000 Exploring group diversity

Rationale of DAPC



Exploring group diversity

PCA of seasonal influenza (A/H3N2) data

Data: seasonal influenza (A/H3N2), 500 HA segments.



Little temporal evolution, burst of diversity in 2002??

Exploring group diversity

PCA of seasonal influenza (A/H3N2) data

Data: seasonal influenza (A/H3N2), 500 HA segments.



Little temporal evolution, burst of diversity in 2002??

Exploring group diversity

DAPC of seasonal influenza (A/H3N2) data



Strong temporal signal, originality of 2006 isolates (new alleles).

Exploring group diversity

DAPC of seasonal influenza (A/H3N2) data



Strong temporal signal, originality of 2006 isolates (new alleles).

Exploring group diversity

Other features

DAPC can be used to:

- provides group assignment probabilities
- can use supplementary individuals
- can predict group membership of new data
- can be used for variable selection



Identifying groups 00000 000000 Exploring group diversity

Time to get your hands dirty (again)!



The pdf of the practical is online:

```
http://adegenet.r-forge.r-project.org/
```

or

```
\mathsf{Google} \rightarrow \mathsf{adegenet} \rightarrow \mathsf{documents} \rightarrow \text{``GDAR August 2016''}
```