# Multivariate analysis of genetic data
## — exploring group diversity —

Thibaut Jombart

MRC Centre for Outbreak Analysis and Modelling
Imperial College London

*Population genomics in* ®
Lausanne
23 Aug 2016

# Outline

Introduction

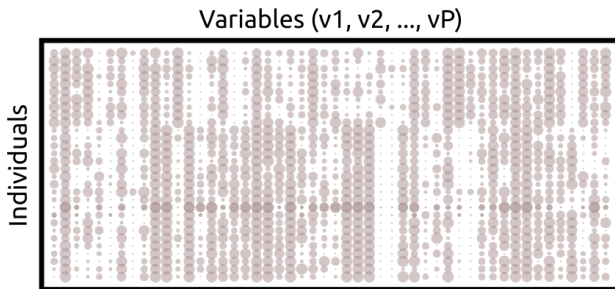Identifying groups using $K$-means clustering

Exploring group diversity
    Aggregating data
    Optimizing group differences
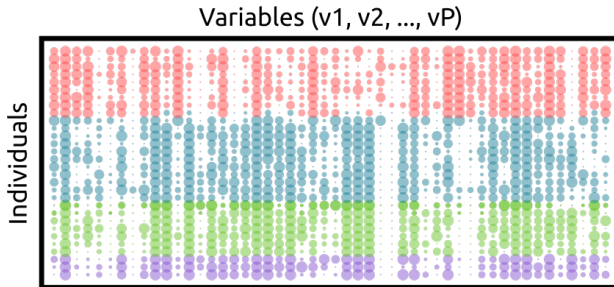    Discriminant Analysis of Principal Components

# Outline

## Genetic data: introducing group data

Variables (v1, v2, ..., vP)



Individuals

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

Introduction

Identifying groups using $K$-means clustering
oooooo

Exploring group diversity
ooo
oooo
oooooo

## Genetic data: introducing group data



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

# Outline

Introduction
○

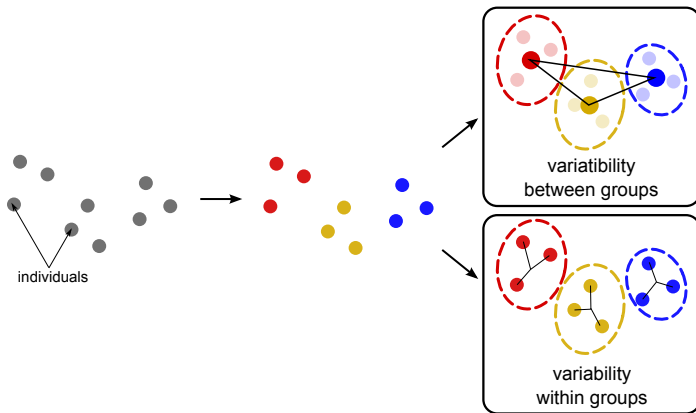Identifying groups using $K$-means clustering
●○○○○○

Exploring group diversity
○○○
○○○○
○○○○○○

# K-means underlying model

ANOVA model:

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

Introduction
○

Identifying groups using $K$-means clustering
○●○○○○

Exploring group diversity
○○○
○○○○
○○○○○○

# 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 $i$
- $\boldsymbol{\mu}_k \in \mathbb{R}^p$: vector of means allele frequencies of group $k$

Introduction
○

Identifying groups using $K$-means clustering
○●○○○○

Exploring group diversity
○○○
○○○○
○○○○○○

# 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 $i$
- $\boldsymbol{\mu}_k \in \mathbb{R}^p$: vector of means allele frequencies of group $k$

# 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 $i$
- $\boldsymbol{\mu}_k \in \mathbb{R}^p$: vector of means allele frequencies of group $k$

# K-means algorithm

The K-mean problem is solved by the following algorithm:

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

# K-means algorithm

The K-mean problem is solved by the following algorithm:

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

# K-means algorithm
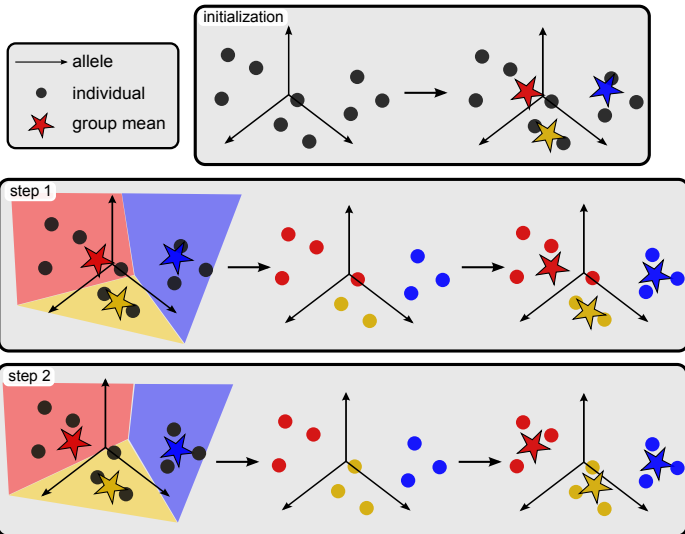
The K-mean problem is solved by the following algorithm:

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

# K-means algorithm

The K-mean problem is solved by the following algorithm:

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

Introduction
○

Identifying groups using $K$-means clustering
○○○●○○

Exploring group diversity
○○○
○○○○
○○○○○○

# K-means algorithm

# K-means: limitations and extensions

## Limitations

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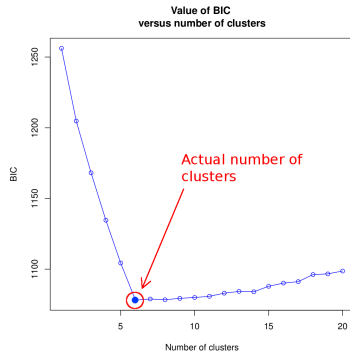- K-means does not identify the number of clusters ($K$)

## Extension

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

Introduction

Identifying groups using $K$-means clustering
○○○○●○

Exploring group diversity
○○○
○○○○
○○○○○○

# 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

# Genetic clustering using K-means & BIC

(Jombart *et al.* 2010, *BMC Genetics*)

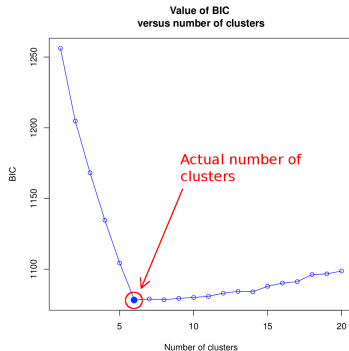Simulated data: island model with 6 populations



**Value of BIC versus number of clusters**

Actual number of clusters

Performances:

- K-means ≥ STRUCTURE on simulated data (various island and stepping stone models)

- orders of magnitude faster (seconds vs hours/days)

Introduction
○

Identifying groups using $K$-means clustering
○○○○○●

Exploring group diversity
○○○
○○○○
○○○○○○

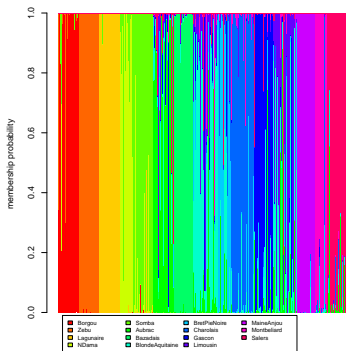# Genetic clustering using K-means & BIC

(Jombart *et al.* 2010, *BMC Genetics*)

Simulated data: island model with 6 populations



Performances:

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

# Outline

Introduction

Identifying groups using $K$-means clustering

Exploring group diversity
  Aggregating data
  Optimizing group differences
  Discriminant Analysis of Principal Components

# Why identifying clusters is not the whole story

Example of cattle breeds diversity (30 microsatellites, 704 individuals).
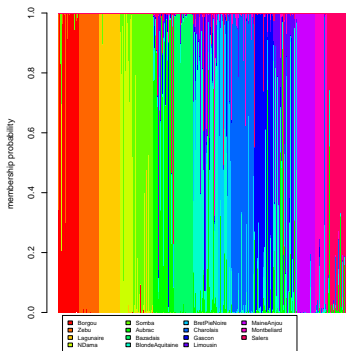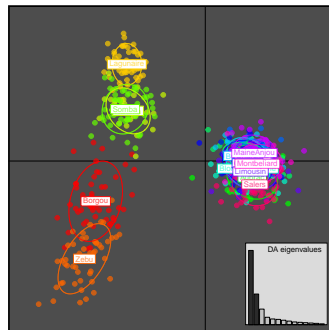
Group membership probabilities:



Important to assess the relationships between clusters.

Introduction
○

Identifying groups using $K$-means clustering
○○○○○○

Exploring group diversity
●○○
○○○○
○○○○○○

# Why identifying clusters is not the whole story

Example of cattle breeds diversity (30 microsatellites, 704 individuals).

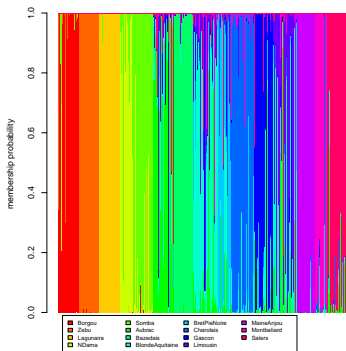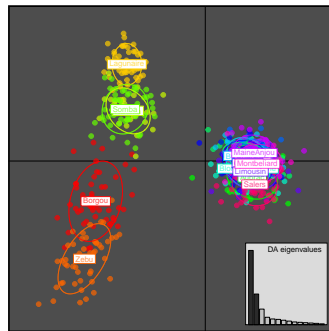Group membership probabilities:



Multivariate analysis:



Important to assess the relationships between clusters.

# Why identifying clusters is not the whole story

Example of cattle breeds diversity (30 microsatellites, 704 individuals).
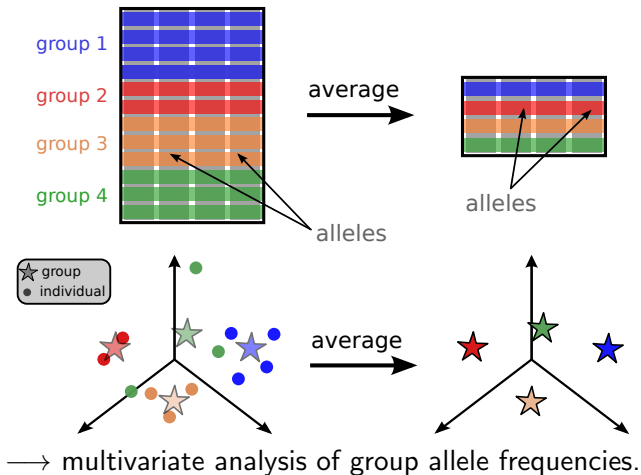
Group membership probabilities:

Multivariate analysis:



**Important to assess the relationships between clusters.**

# Aggregating data by groups



$\longrightarrow$ multivariate analysis of group allele frequencies.

Introduction    Identifying groups using $K$-means clustering    Exploring group diversity
o               oooooo                                           oo●
                                                                 oooo
                                                                 oooooo

# Analysing group data

Available methods:

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

Criticism:

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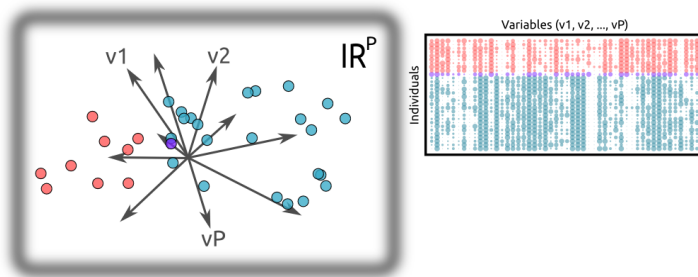- CA: possible artefactual outliers

Introduction
Identifying groups using $K$-means clustering
Exploring group diversity

# Analysing group data

Available methods:

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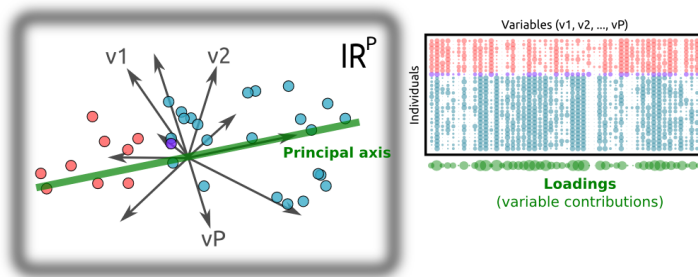- Correspondance Analysis (CA) of allele counts

Criticism:

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

Introduction
○

Identifying groups using $K$-means clustering
○○○○○○

Exploring group diversity
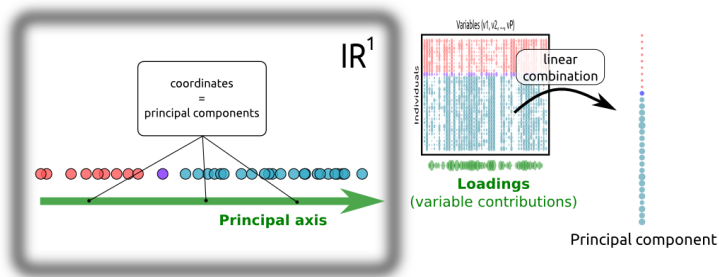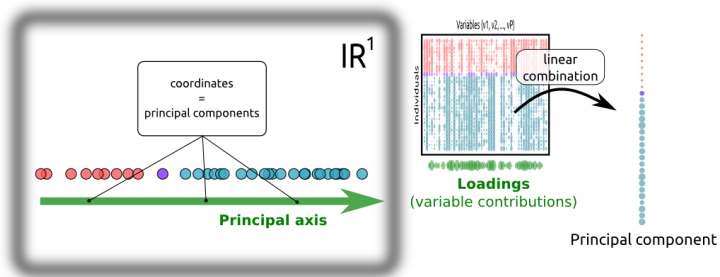○○○
●○○○
○○○○○○

# Multivariate analysis: reminder



Find principal components with *maximum total variance.*

# Multivariate analysis: reminder



Find principal components with *maximum total variance.*
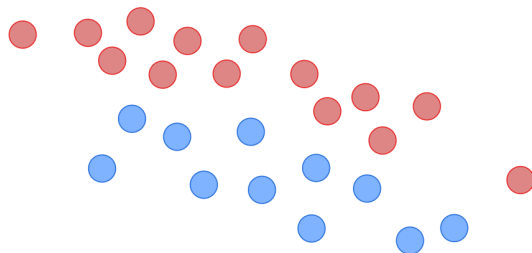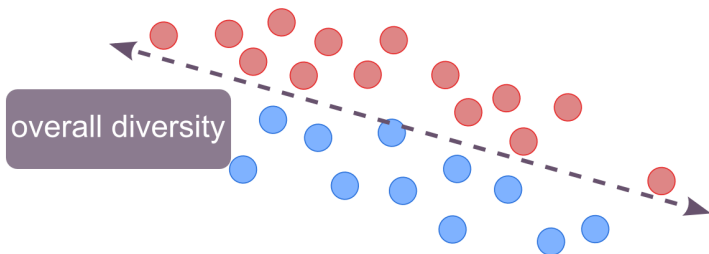
# Multivariate analysis: reminder



Find principal components with *maximum total variance.*

Introduction
○

Identifying groups using $K$-means clustering
○○○○○○

Exploring group diversity
○○○
●○○○
○○○○○○

# Multivariate analysis: reminder



Find principal components with *maximum total variance*.

Introduction
○

Identifying groups using $K$-means clustering
○○○○○○

Exploring group diversity
○○○
○●○○
○○○○○○

# But total variance may not reflect group differences



Need to optimize different criteria.

Introduction
○

Identifying groups using $K$-means clustering
○○○○○○

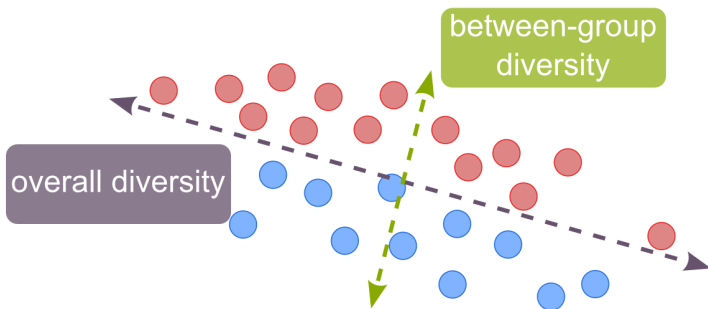Exploring group diversity
○○○
○●○○
○○○○○○

# But total variance may not reflect group differences



overall diversity

Need to optimize different criteria.

# But total variance may not reflect group differences



between-group diversity

overall diversity

**Need to optimize different criteria.**

# 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

# Optimizing different criteria

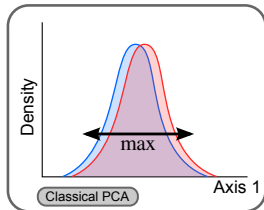Similar approaches to PCA can be used to optimize different quantities:

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

# Optimizing different criteria

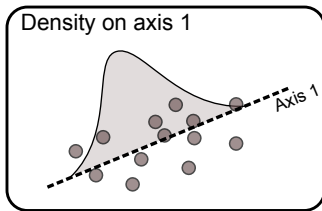Similar approaches to PCA can be used to optimize different quantities:

- **PCA**: *total* variance
- **Between-group PCA**: variance *between* groups
- **Within-group PCA**: variance *within* groups
- **Discriminant Analysis**: variance *between* groups / variance *within* groups
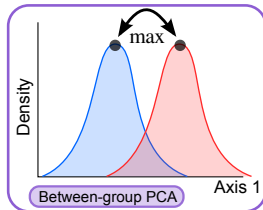
# Optimizing different criteria

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
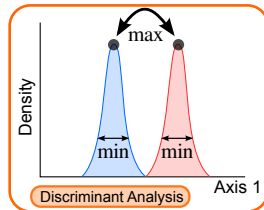
Introduction
o

Identifying groups using $K$-means clustering
oooooo

Exploring group diversity
ooo
ooo●
oooooo

# From PCA to DA: increasing group differentiation



Density on axis 1

Axis 1

Max. total diversity

Max. diversity between groups

Max. separation of groups

## 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]:
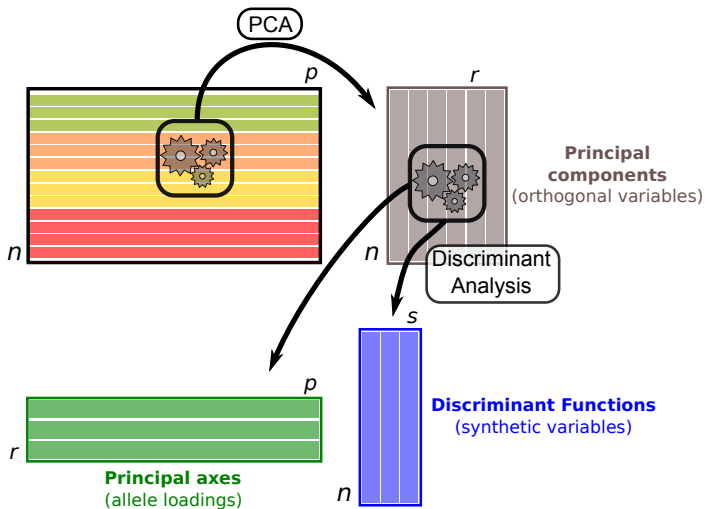
- data orthogonalisation/reduction using PCA before DA

- overcomes limitations of DA

- group membership probabilities, group prediction

---

[1] Jombart et al. 2010, *BMC Genetics*

Introduction
○

Identifying groups using $K$-means clustering
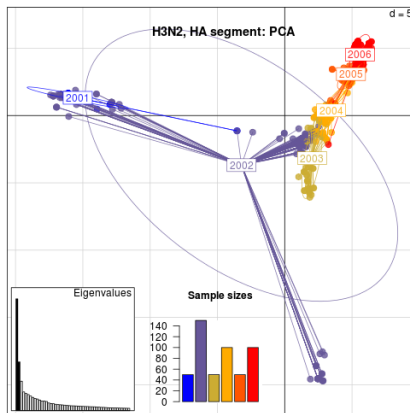○○○○○○

Exploring group diversity
○○○
○○○○
○●○○○○

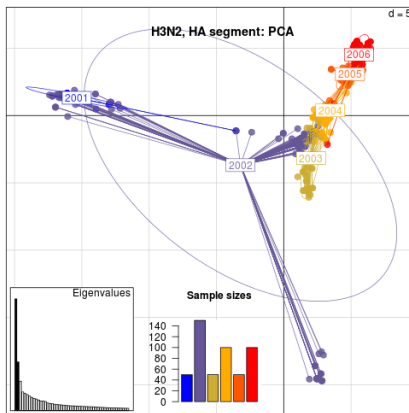# Rationale of DAPC

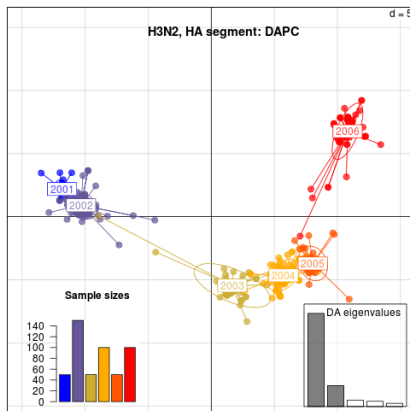# PCA of seasonal influenza (A/H3N2) data

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



Little temporal evolution, burst of diversity in 2002??

# PCA of seasonal influenza (A/H3N2) data

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



Little temporal evolution, burst of diversity in 2002??

Introduction
○

Identifying groups using $K$-means clustering
○○○○○○

Exploring group diversity
○○○
○○○○
○○○●○○

# DAPC of seasonal influenza (A/H3N2) data



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

Introduction
○

Identifying groups using $K$-means clustering
○○○○○○

Exploring group diversity
○○○
○○○○
○○○●○○

# DAPC of seasonal influenza (A/H3N2) data



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

# 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

Introduction

Identifying groups using $K$-means clustering
○○○○○○

Exploring group diversity
○○○
○○○○
○○○○○●

# Coming next: group diversity, funny plants, and alien abductions



The pdf of the practical is online:

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

or

Google → adegenet → documents → "Lausanne August 2016"