Imperial College London





## Exploring genetic diversity: multivariate approaches

Thibaut Jombart – Imperial College London

t.jombart@imperial.ac.uk

Multivariate analysis of genetic data



I) Introduction: data and methods

2) Multivariate analysis in a nutshell

3) Some examples

4) Spatial genetic structures

5) In practice



#### I) Introduction: data and methods

2) Multivariate analysis in a nutshell

3) Some examples

4) Spatial genetic structures

5) In practice

Multivariate analysis of genetic data

## Genetic data (genetic markers)



Different ways of exploiting this rich information

Multivariate analysis of genetic data

## Two complementary approaches



- One or a small number of models (possible comparisons)
- Parameter estimation based on the data

## Two complementary approaches



- No explicit model, little underlying assumptions
- Description of specific features (e.g. genetic diversity)

## Different types of methods

#### Model-based approaches

- Bayesian clustering (e.g. STRUCTURE, BAPS)
- > Phylogenetic trees (e.g. PhyML, BEAST)

#### Exploratory approaches

- > Distance-based trees (e.g. UPGMA, NJ)
- Multivariate methods (e.g. PCA, PCoA)

## **Different methods**

- Model-based approaches
  - Bayesian clustering (e.g. STRUCTURE, BAPS)
  - > Phylogenetic trees (e.g. PhyML, BEAST)

- Exploratory approaches
  - Distance-based trees (e.g. UPGMA, NJ)
  - > Multivariate methods (e.g. PCA, PCoA)



I) Introduction: data and methods

#### 2) Multivariate analysis in a nutshell

3) Some examples

4) Spatial genetic structures

5) In practice

Multivariate analysis of genetic data

## Genetic data (again)



This information can be considered from the geometric point of view.

## Multivariate analysis – rationale (1/3)



Which directions are the most informative?

Multivariate analysis of genetic data

## Multivariate analysis – rationale (2/3)



Multivariate analysis of genetic data

## Multivariate analysis – rationale (2/3)



Summarise the genetic diversity among individuals / populations

## Multivariate analysis – rationale (3/3)



Multivariate analysis of genetic data

## In practice, lots of methods are used

#### To name a few (used in genetics):

- Principal Component Analysis (PCA)
  - > centred / not centred / fancy centring
  - > scaled / not scaled / fancy scaling
  - > transformed for compositional data
- Principal Coordinates Analysis (PCoA), aka (Metric) Multidimensional Scaling (MDS)
   > many genetic distances
- Correspondence Analysis (CA)
- Discriminant Analysis (DA)
- ... (review in Jombart et al. 2009, Heredity 102, 330-341)

### Content

I) Introduction: data and methods

2) Multivariate analysis in a nutshell

3) Some examples

4) Spatial genetic structures

5) In practice

Multivariate analysis of genetic data

## Multivariate analyses: some examples (1/6)

Getting a picture of the genetic diversity



## Multivariate analyses: some examples (2/6)



Multivariate analysis of genetic data

## Multivariate analyses: some examples (3/6)

#### Mapping the genetic differentiation



Multivariate analysis of genetic data

## Multivariate analyses: some examples (4/6)

#### Mapping the genetic differentiation (again)



Data

matrix

Multivariate analysis of genetic data

## Multivariate analyses: some examples (5/6)

#### Studying hybridization





(Paul et al. 2010, Biological Invasions)



Multivariate analysis of genetic data

## Multivariate analyses: some examples (6/6)

#### Analysing the temporal evolution of pathogens



Multivariate analysis of genetic data

## A wide range of possible applications

- Many different data and questions
- Wide range of existing methods
- Adaptability to specific problems

   (e.g. non-even sampling, different ploidy levels,
   heterogeneous variation amongst loci)
- Datasets are growing bigger and more complex: more multivariate analyses to come...

### Content

I) Introduction: data and methods

2) Multivariate analysis in a nutshell

3) Some examples

4) Spatial genetic structures

5) In practice

Multivariate analysis of genetic data

## Why look for spatial genetic structures?

Most genetic models predict that genetic diversity should be spatially structured:







## Taking spatial information into account

• Usual multivariate methods do not use spatial information.



• They can reveal `obvious' spatial patterns, but will overlook more subtle structures.

• To seek spatial genetic structures, we must find the part of the **genetic variability related to spatial proximity** between individuals/populations.

## Usual multivariate analyses (recall)



Multivariate analysis of genetic data

#### Spatial Principal Component Analysis (sPCA) (Jombart *et al.* 2009)



Multivariate analysis of genetic data

### Content

I) Introduction: data and methods

2) Multivariate analysis in a nutshell

3) Some examples

4) Spatial genetic structures

#### 5) In practice

#### Population genetic software – the scary picture



(Excoffier & Heckel 2006, Nature Reviews Genetics)

`In a perfect world, research teams would be able to develop analysis tools to address their specific problem, but in practice they have to make their data fit the available tools, leading to obvious discrepancies between the initial goals and the results'

# Taking genetic markers into the field of (multivariate) statistics

#### **Population genetic software:**

- very few multivariate methods
- no plasticity
- poor data interoperability



- many (most) multivariate methods
- total plasticity
- tons of statistical methods (tests, modelling, Monte-Carlo)
- great graphics
- great interoperability (e.g., GIS)
- programming language
- free software



changed directly with the other programs.

## The adegenet package for (1/3) (Jombart 2008)

#### **Purpose:**

- take genetic markers into a suitable format
- adapt multivariate methods to genetic markers
- provide advanced data handling
- provide standard population genetics tools
- implement novel methods (*e.g.*, sPCA, DAPC, seqTrack)

## The adegenet package for (2/3)





#### Where to get information:

- reference: Jombart (2008) Bioinformatics 24: 1403-1405
- *adegenet* website: http://adegenet.r-forge.r-project.org/ (tutorials and manuals for download)
- *adegenet* forum: adegenet\_forum@lists.r\_forge.r\_project.org
- Here and now!