

CS 260: Foundations of Data Science

Prof. Sara Mathieson

Fall 2023



HVERFORD
COLLEGE

- Piazza counts as participation (both asking AND **answering!**)
- Tomorrow (Friday): **zoom office hours** 3:30-5pm
- **Lab 7 and project proposal** (both short)
 - Due Wednesday Nov 8

Outline for November 2

- Discuss final project
- Dimensionality reduction
- PCA for data visualization

Outline for November 2

- Discuss final project
- Dimensionality reduction
- PCA for data visualization

Timeline and Logistics

- November 8: project proposal due
- November 8 - December 7: working on projects
- December 7, 12, 14: oral project presentations during class
- December 22: github repos must be finalized

Outline for a typical project:

- Find a dataset (see project writeup)
- Run an algorithm we've discussed on the dataset
- Try to do a comparison
 - run the algorithm in multiple ways
 - different data pre-processing
 - try a different algorithm
- Evaluate, interpret, and visualize the results

Project Proposal

- **Title** and **names** of both partners
 - Pair work is required!
- A **dataset** (what is n ? what is p ?)
- An **algorithm** or set of algorithms you will develop and/or apply to this dataset
- A **scientific question** you are trying to answer
 - “Will Naive Bayes or logistic regression perform better on my dataset?”
 - “How will pre-processing a dataset or subsampling features affect the results?”
- A way to **evaluate, interpret, and visualize** the results
- **References**

Project Group Options

- If you would like a random partner, please email me ASAP!
- If you **really** prefer to work individually or in a group of 3, email me ASAP!

Final Project Deliverables

- Main deliverable: presentation
 - In class Dec 7, 12, 14 (last 1.5 weeks of classes)
 - 10 min per pair
 - Peer feedback
- On git (by Dec 22)
 - Lab Notebook (in README.md)
 - Project Code
 - Slides

Project Lab Notebook

- As you accept your git repo, start creating a “lab notebook” in your **README.md**
- This should say:
 - who was working (which partner)
 - date
 - how long
 - briefly what what accomplished

Sara: 03-07-18 (2hrs)

- now averaging the Markov chain, fixed all the results
- combined ancestral 1000 genomes still running (need to start similar for SGDP)
- started new runs with filtering to only have selected alleles in the “selected pop” and only have ancestral alleles in the “reference panel”

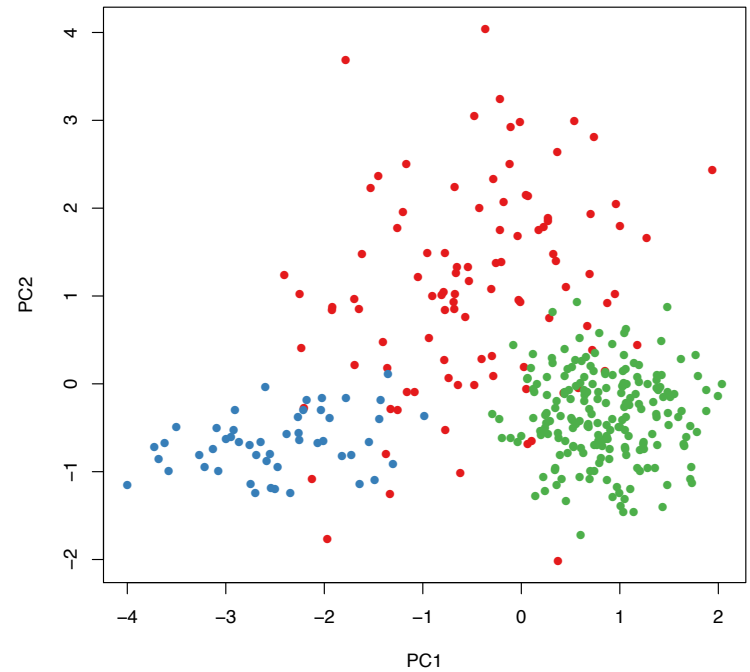
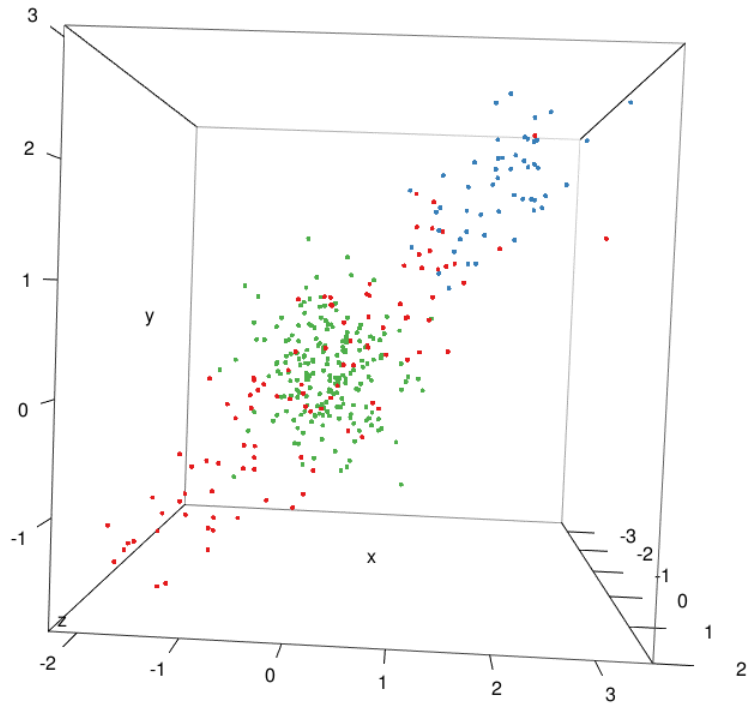
Outline for November 2

- Discuss final project
- Dimensionality reduction
- PCA for data visualization

Principal Components Analysis (PCA)

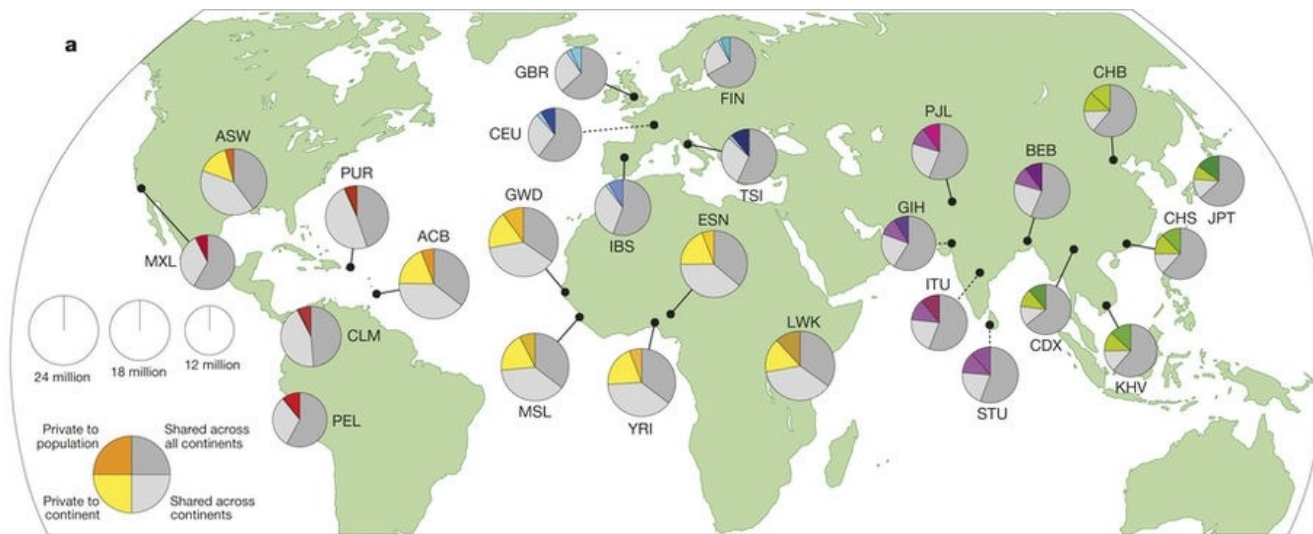
- Transforms p -dimensional data so that the new first dimension explains as much of the variation as possible, the new second explains as much of the remaining variation as possible, and so on
- Typically, we look at the first few dimensions of the transformed data and use as a means of dimensionality reduction and visualization
- PCA is a linear transformation
- PCA is often used for:
 - Data visualization
 - Infer qualitative relationships between groups

Principal component analysis

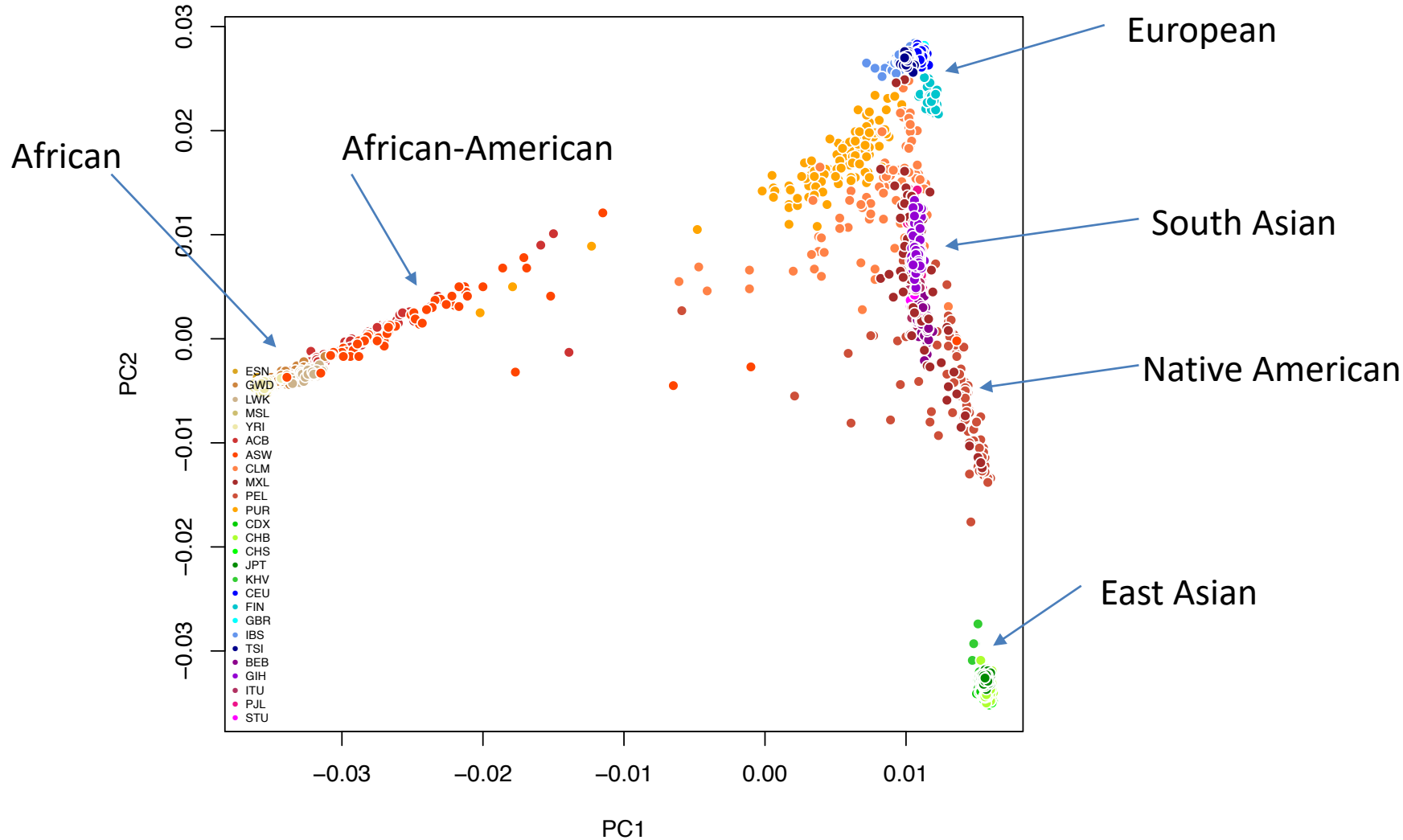


The 1000 Genomes project

- Whole-genome **sequence data** from 2504 individuals from 26 populations
- A catalog of human genetic variation, useful as a reference or **imputation** panel
- Completely public. Download from <ftp://ftp-trace.ncbi.nih.gov/1000genomes/>

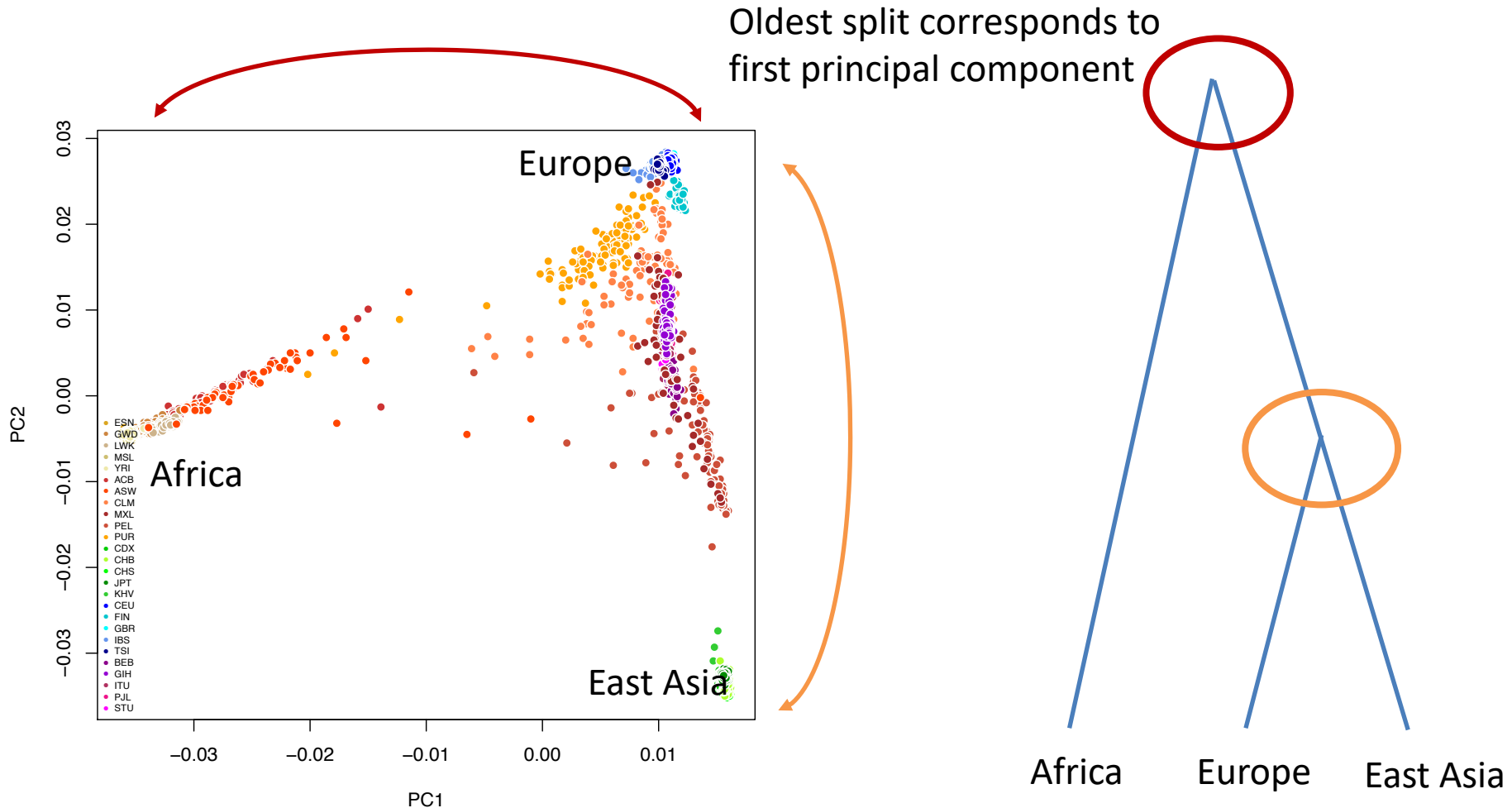


Global population structure



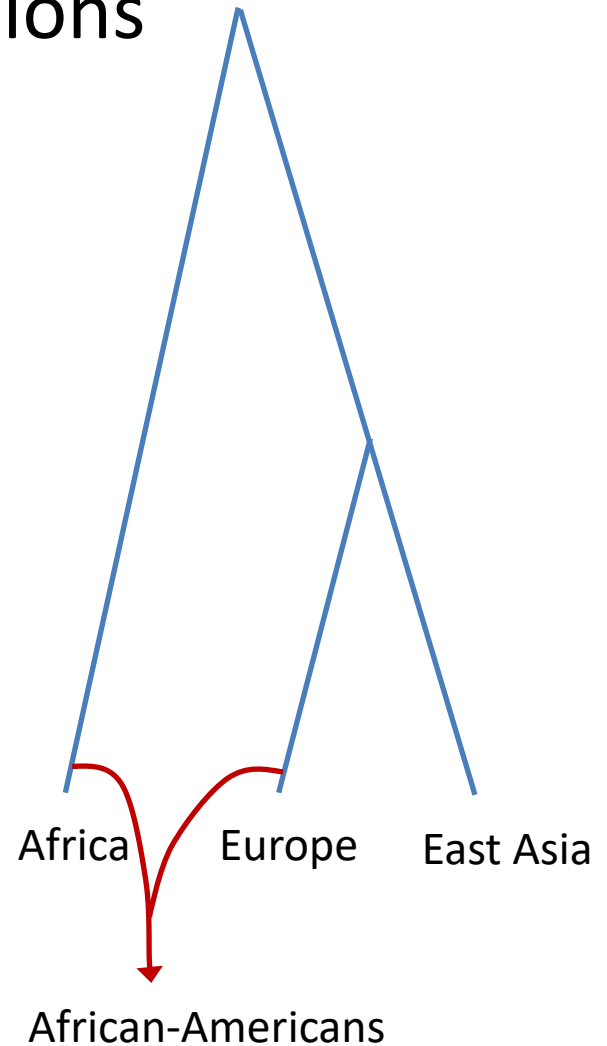
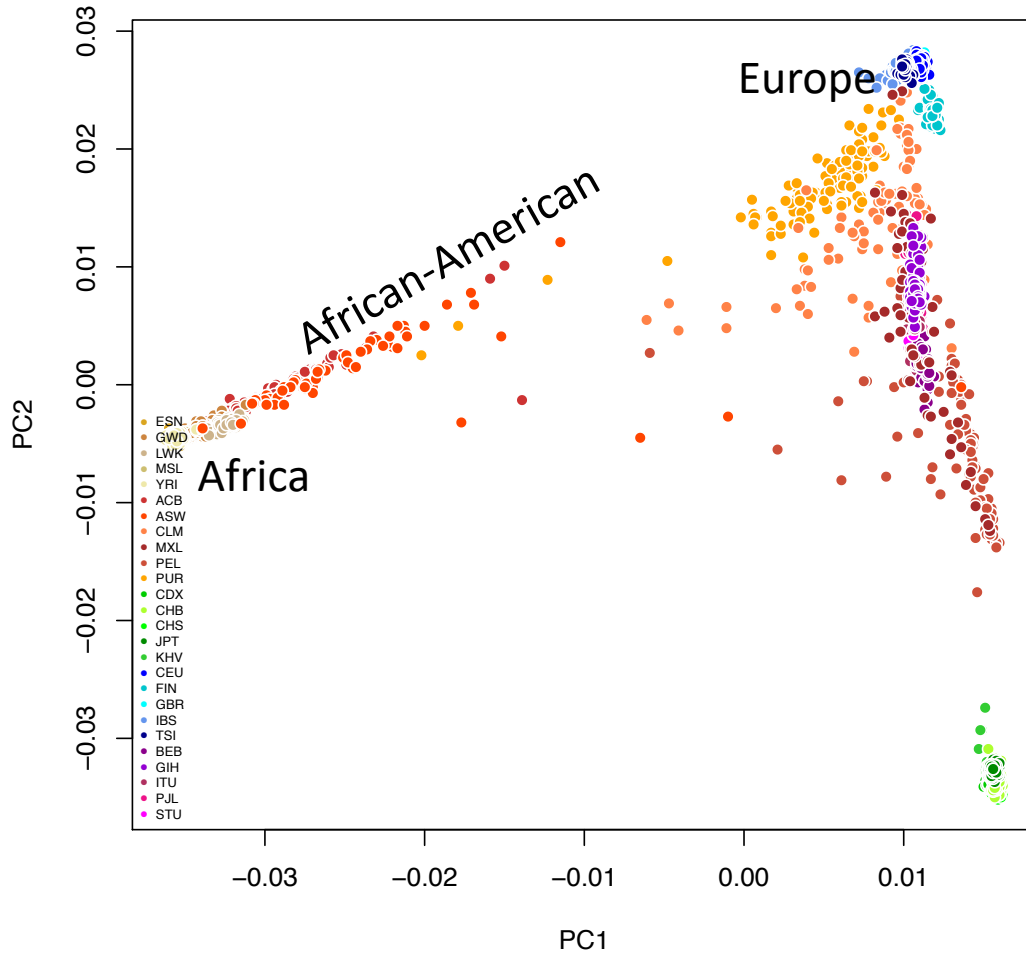
What causes these patterns?

1. Populations **splits** separate populations

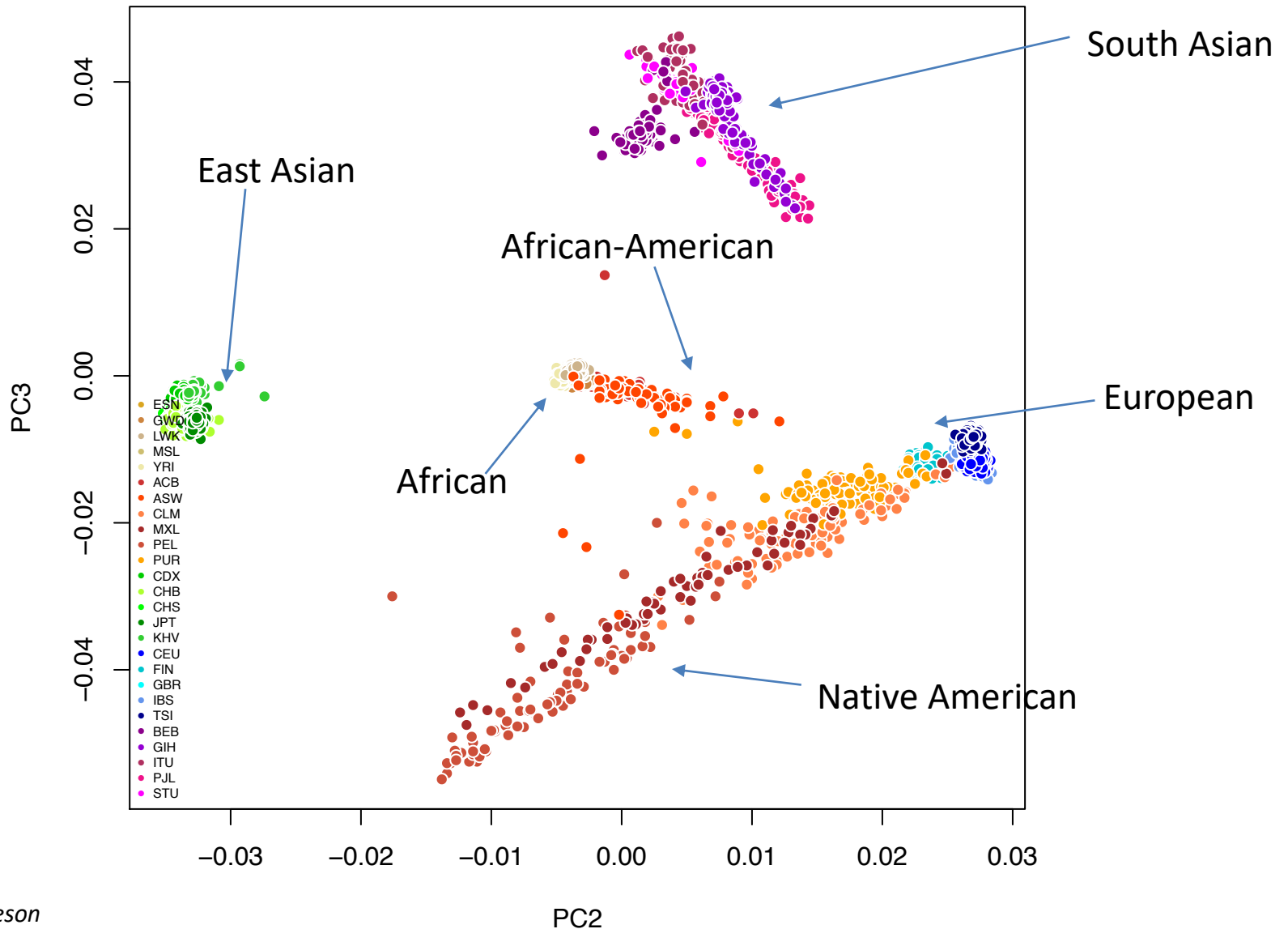


What causes these patterns?

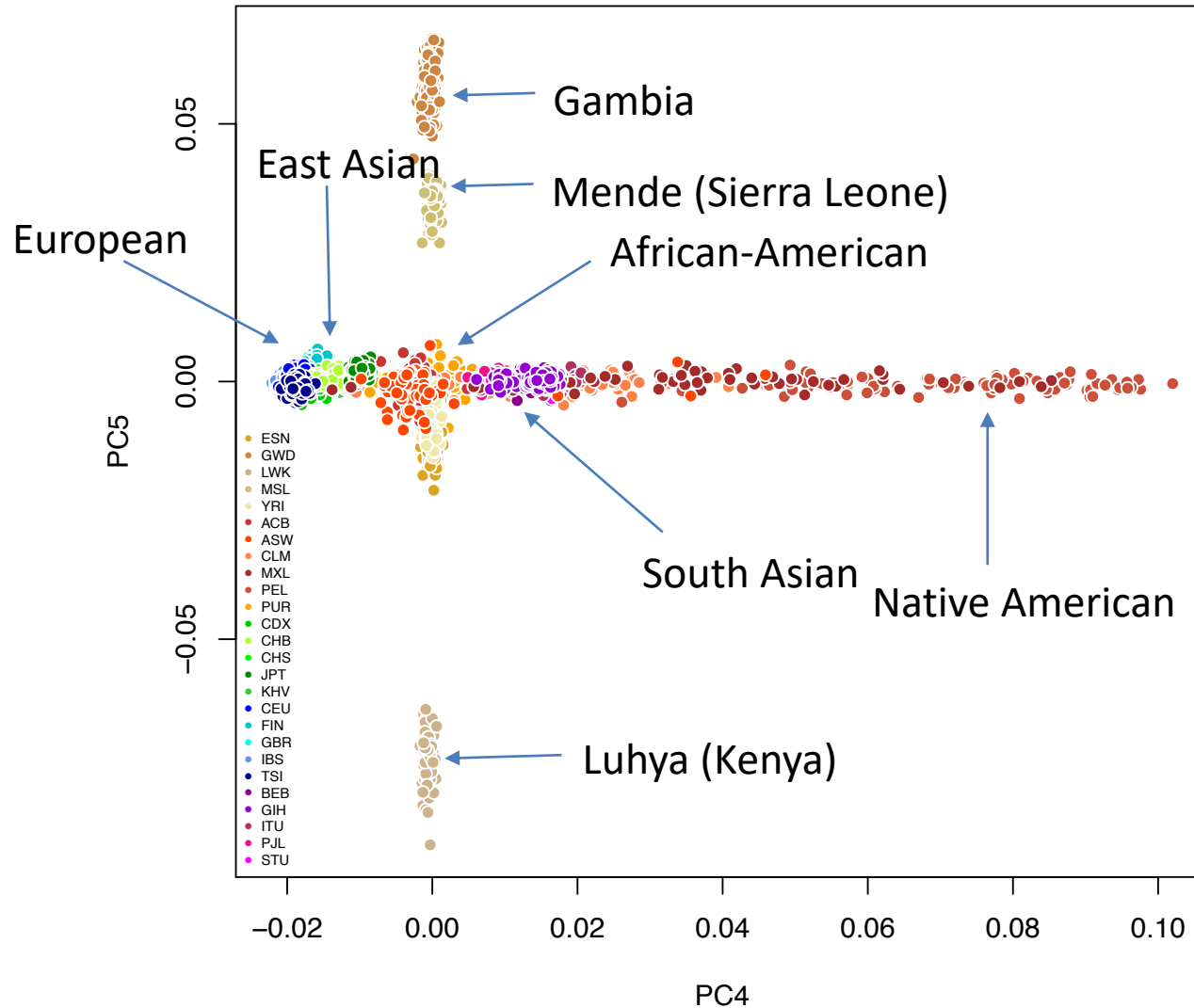
2. Admixture merges populations




Global population structure



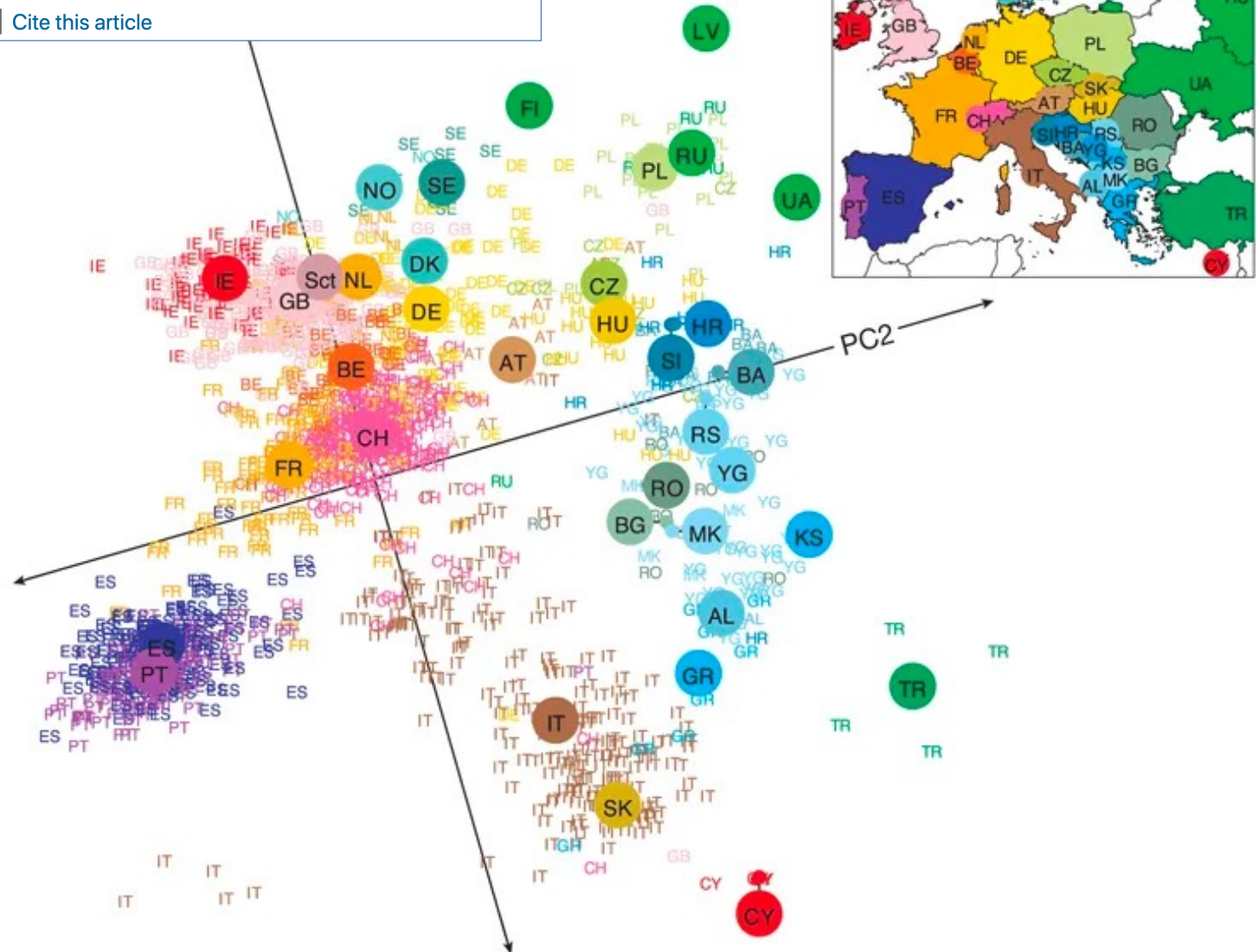
Global population structure



Genes mirror geography within Europe

John Novembre , Toby Johnson, Katarzyna Bryc, Zoltán Kutalik, Adam R. Boyko, Adam Auton, Amit Indap, Karen S. King, Sven Bergmann, Matthew R. Nelson, Matthew Stephens & Carlos D. Bustamante





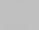




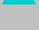



Nature **456**, 98–101(2008) | [Cite this article](#)



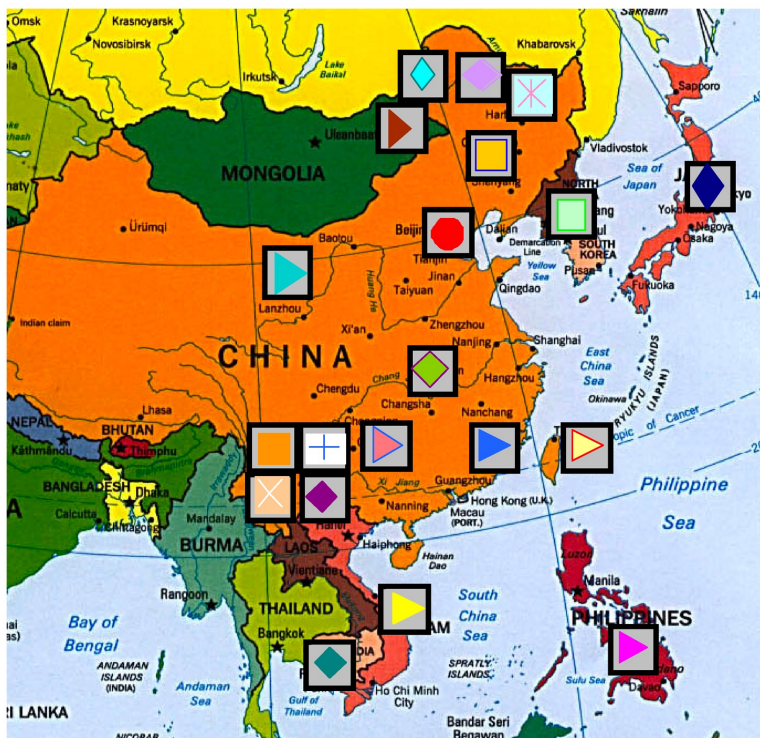
Analysis of East Asia Genetic Substructure Using Genome-Wide SNP Arrays

Chao Tian, Roman Kosoy, Annette Lee, Michael Ransom, John W. Belmont, Peter K. Gregersen, Michael F. Seldin 

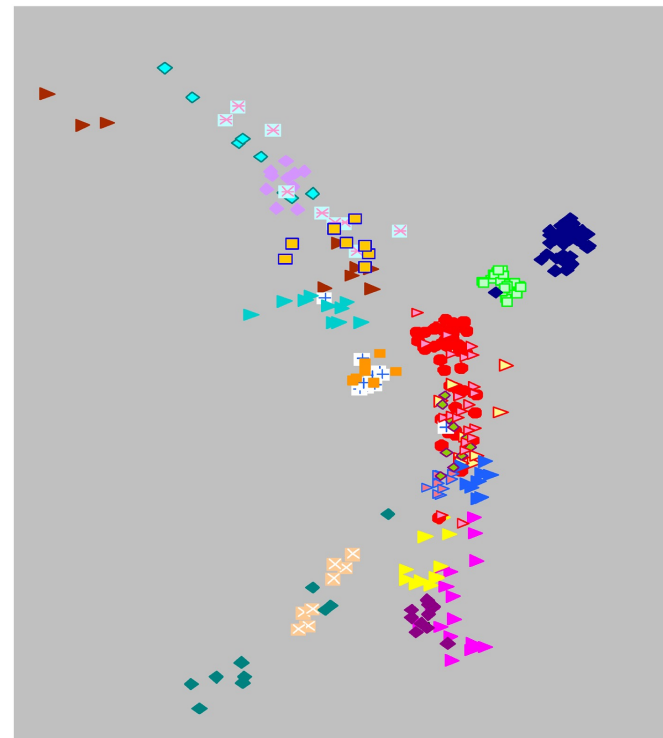
Published: December 5, 2008 • <https://doi.org/10.1371/journal.pone.0003862>

-  FIL
-  VIET
-  LAHU
-  DAI
-  CAMB
-  CHB
-  MGL
-  ORQ
-  DAUR
-  KOR
-  TWN
-  YI
-  HEZ
-  MIAO
-  NAXI
-  SHE
-  TU
-  TUJ
-  XIBO
-  CHA
-  JPT
-  YAK

C

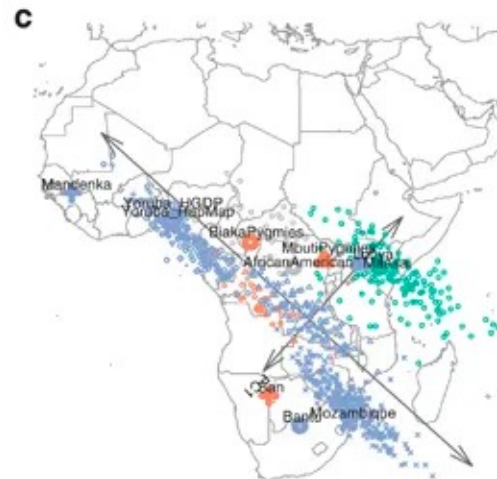
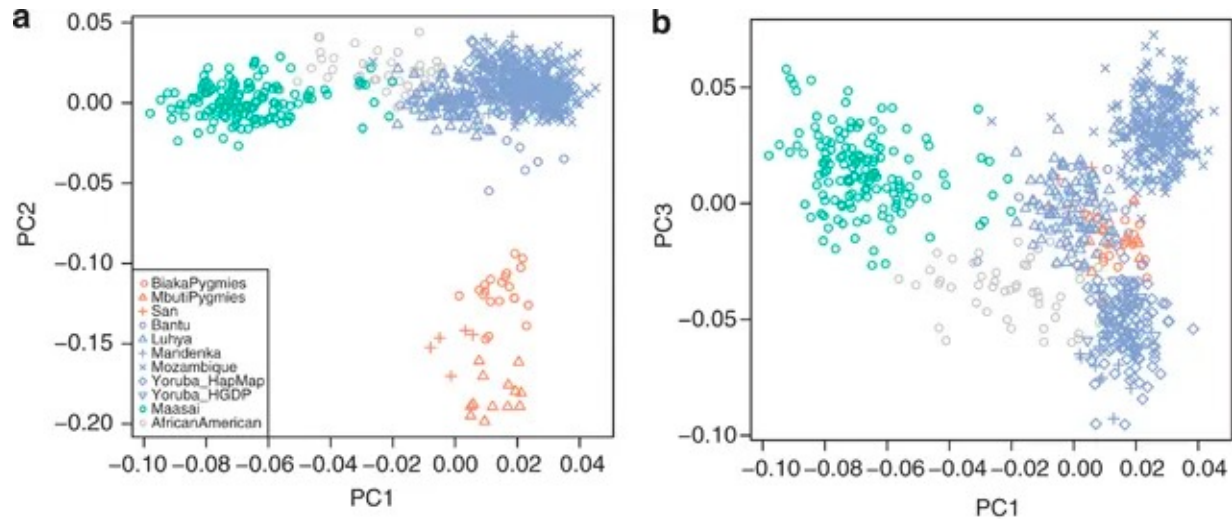


D



A genomic analysis identifies a novel component in the genetic structure of sub-Saharan African populations

Martin Sikora, Hafid Laayouni, Francesc Calafell, David Comas & Jaume Bertranpetit 



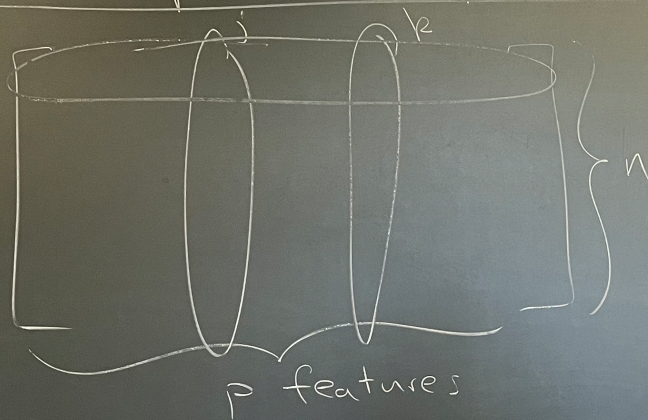
Outline for November 2

- Discuss final project
- Dimensionality reduction
- **PCA for data visualization**

Step 1

Principal Component

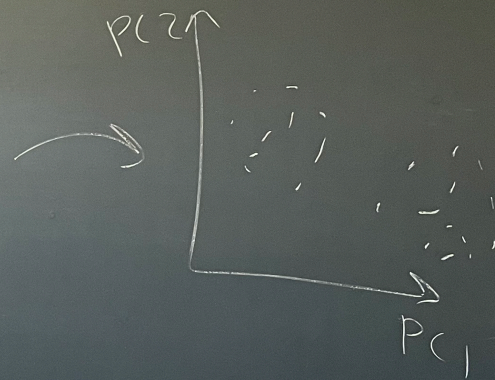
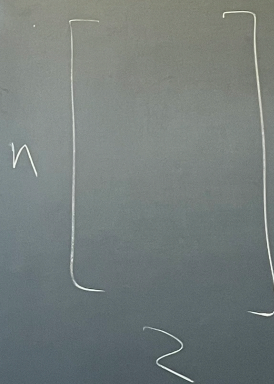
$X = X_{orig}$



$$p \gg n$$

Analysis

goal: create $n \times 2$ matrix for visualization

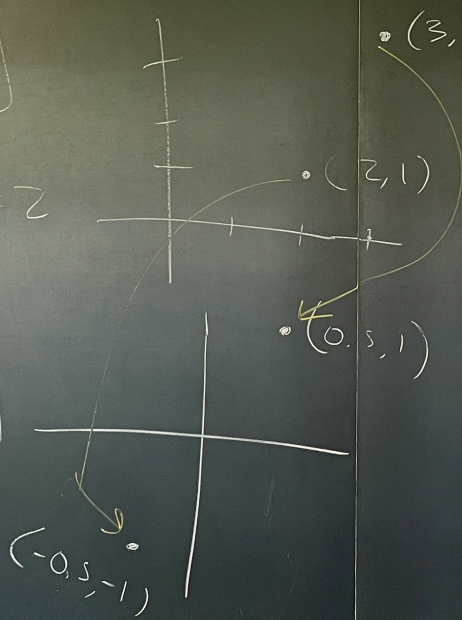


Step 2 subtract off column-wise mean

$$X_{\text{orig}} = \begin{bmatrix} 2 & 1 \\ 3 & 3 \end{bmatrix}$$

\downarrow \downarrow
 $\bar{x}_1 = 2.5$ $\bar{x}_2 = 2$

$$X = \begin{bmatrix} -0.5 & -1 \\ 0.5 & 1 \end{bmatrix}$$



Step 3

Compute covariance matrix A

$$\text{cov}(f, g) = \frac{1}{n-1} \sum_{i=1}^n (f_i - \bar{f})(g_i - \bar{g})$$

$$\text{cov}(f, f) = \text{var}(f) = \frac{1}{n-1} \sum_{i=1}^n (f_i - \bar{f})^2$$

2 features

runtime $\Rightarrow O(np^2)$

$$A = \begin{bmatrix} \text{cov}(f, f) & \text{cov}(f, g) \\ \text{cov}(g, f) & \text{cov}(g, g) \end{bmatrix}$$

Square!

Symmetric!

$p \times p$

Step 4

compute eigenvalues & eigenvectors of A

→ sort by eigenvalue high → low

$$A\vec{v} = \lambda\vec{v}$$

eigenvalue

eigenvector

$$\det(A - \lambda I) = 0$$

Solve for λ

$$\begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

$$\begin{bmatrix} \lambda & 0 \\ 0 & \lambda \end{bmatrix}$$

plug back in!

Step 5

$$W = \begin{bmatrix} \vec{v}_1 & \vec{v}_2 & \dots & \vec{v}_r \end{bmatrix}_{p \times r}$$

first eigenvector

usually $r=2$

★

$$\Rightarrow T_{n \times r} = X_{n \times p} W_{p \times r}$$

transformed data!

Step 6

$$PC_2 = \vec{v}_2$$

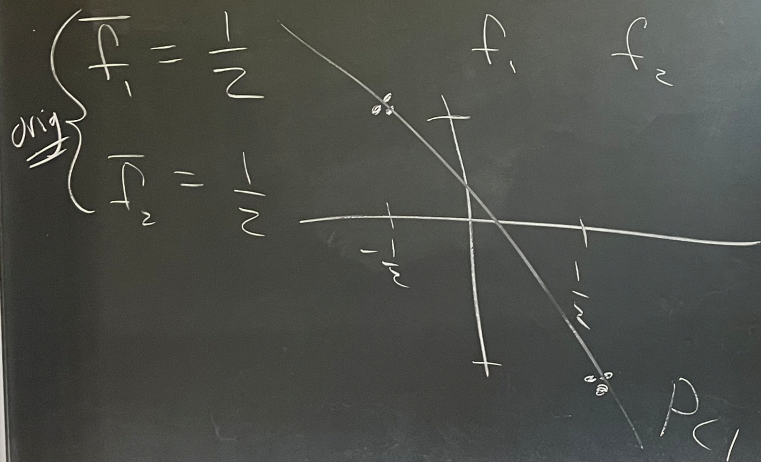
$\vec{v}_2 = PC_1$

Handout 16

Handout 16

Step 1
4
2

$$X = \begin{bmatrix} -1/2 & 1/2 \\ -1/2 & 1/2 \\ -1/2 & 1/2 \\ 1/2 & -1/2 \\ 1/2 & -1/2 \\ 1/2 & -1/2 \end{bmatrix}$$



Step 3

$$A = \begin{bmatrix} \text{var}(f_1) & \text{cov}(f_1, f_2) \\ \text{cov}(f_2, f_1) & \text{var}(f_2) \end{bmatrix}$$

$$\begin{aligned} \bar{f}_1 &= 0 \\ \bar{f}_2 &= 0 \end{aligned} \quad \text{cov}(f_1, f_2) = \frac{1}{6-1} \left(-\frac{1}{2} \cdot \frac{1}{2} \right) \cdot 6$$
$$= -\frac{3}{10}$$

$$\Rightarrow A = \begin{bmatrix} 3/10 & -3/10 \\ -3/10 & 3/10 \end{bmatrix}$$

Step 4

$$\det(A - \lambda I) = 0$$

$$\begin{bmatrix} 3/10 & -3/10 \\ -3/10 & 3/10 \end{bmatrix} - \lambda \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} = 0$$

$\frac{1}{2}) = 6$

$$\det \begin{pmatrix} 3/10 - \lambda & -3/10 \\ -3/10 & 3/10 - \lambda \end{pmatrix} = 0$$

$$\left(\frac{3}{10} - \lambda\right)^2 - \left(\frac{3}{10}\right)^2 = 0$$

$$\cancel{\left(\frac{3}{10}\right)^2} - 2 \cdot \frac{3}{10} \lambda + \lambda^2 - \cancel{\left(\frac{3}{10}\right)^2} = 0$$

$$\det \begin{bmatrix} a & b \\ c & d \end{bmatrix} = ad - bc$$

$$\lambda^2 - \frac{3}{5} \lambda = 0$$

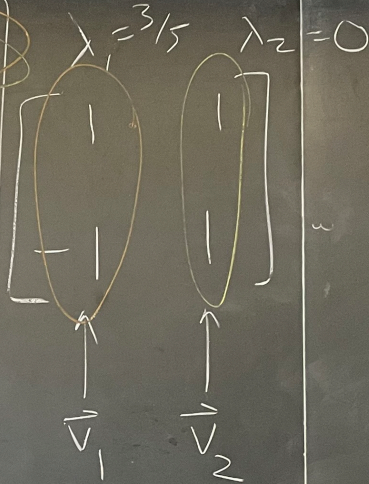
$$\lambda \left(\lambda - \frac{3}{5}\right) = 0 \Rightarrow$$

$$A\vec{v} = \lambda \vec{v}$$

$$\begin{array}{l} \lambda_1 = \frac{3}{5} \\ \lambda_2 = 0 \end{array}$$

$$T_2 = XW_2 =$$

$$\begin{bmatrix} -1/2 & 1/2 \\ -1/2 & 1/2 \\ -1/2 & 1/2 \\ 1/2 & -1/2 \\ 1/2 & -1/2 \\ 1/2 & -1/2 \end{bmatrix}$$



$$= \begin{bmatrix} -1 & 0 \\ -1 & 0 \\ -1 & 0 \\ -1 & 0 \\ -1 & 0 \\ -1 & 0 \end{bmatrix}$$

