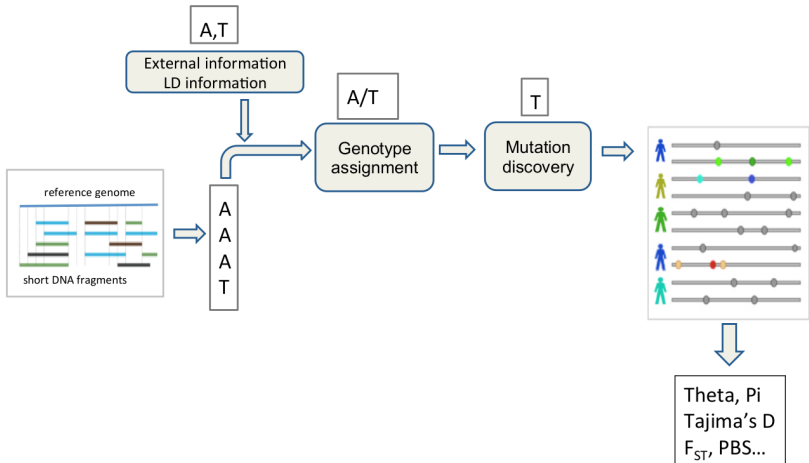
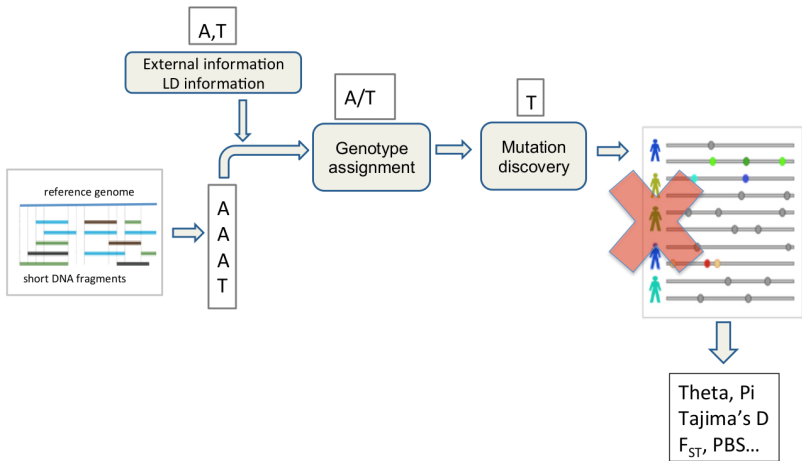


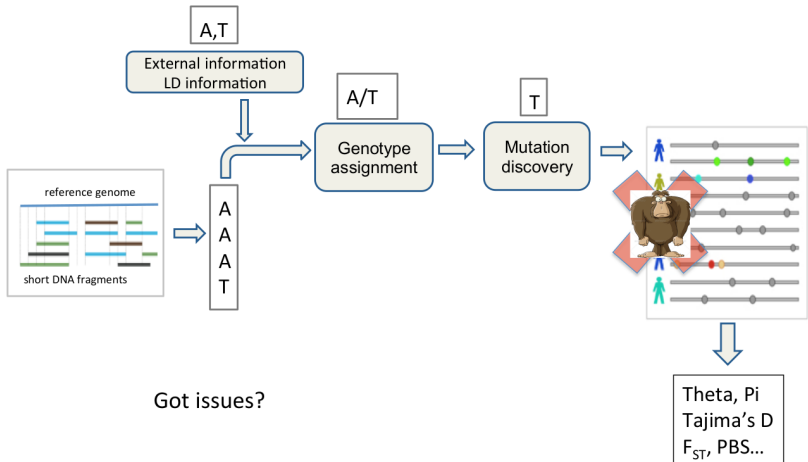
# *Population structure and admixture analysis*

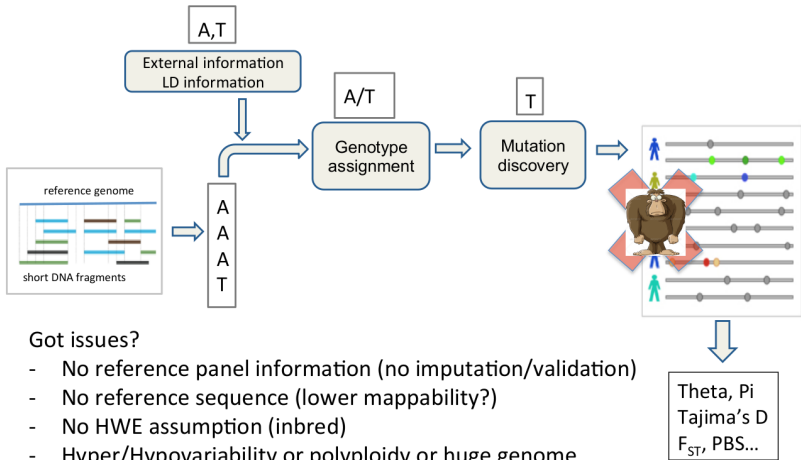
Matteo Fumagalli

---



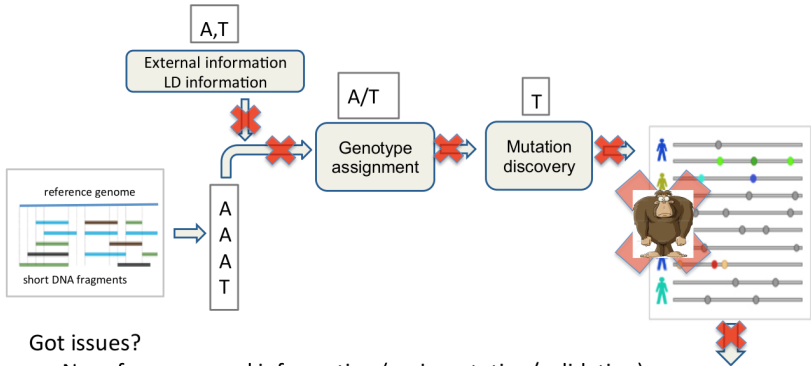






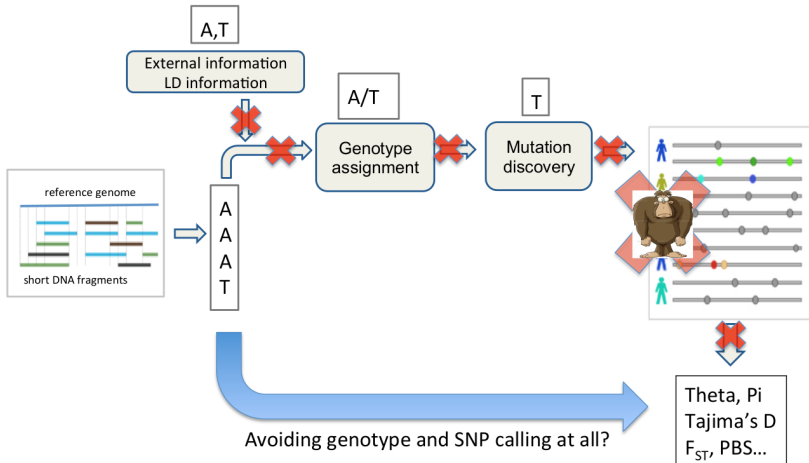
Got issues?

- No reference panel information (no imputation/validation)
- No reference sequence (lower mappability?)
- No HWE assumption (inbred)
- Hyper/Hypovariability or polyploidy or huge genome
- No money (?)
- ...



Got issues?

- No reference panel information (no imputation/validation)
- No reference sequence (lower mappability?)
- No HWE assumption (inbred)
- Hyper/Hypovariability or polyploidy or huge genome
- No money (?)
- **Your inferences will be wrong!**



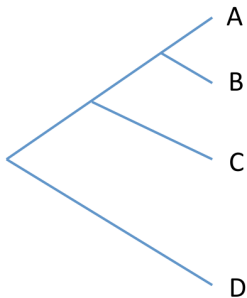
## Intended Learning Outcomes

By the end of this session you will be able to

- understand the theory underlying distance and covariance matrices
- appreciate how to extend such theory to low-coverage data
- acknowledge the process of inferring population structure and admixture from sequencing data
- implement a pipeline in ANGSD to perform the aforementioned analyses

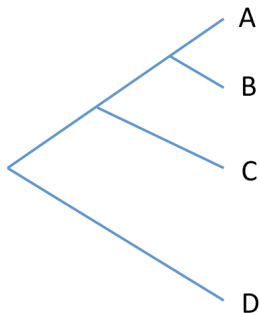


## Genetic distances



Genotype 1	Genotype 2	Distance
aa	aa	0
aa	aA	1
aa	AA	2
aA	aa	1
aA	aA	0
aA	AA	2
...	...	...

## Genetic distances



Genotypes are {aa, aA, AA} as {0, 1, 2}

For individuals  $i$  and  $j$  and  $N$  sites:

$$d(i, j) = -\log \left( 1 - \frac{1}{N} \sum_{s=1}^N \frac{|g(i, s) - g(j, s)|}{2} \right)$$

genotype of  $i$  at site  $s$

e.g.  $G(i=A, s=1)=0$  and  $G(j=B, s=1)=1$  then  $d(i, j)=1$

## Genetic distances from known genotypes

Genotypes are {aa, aA, AA} as {0, 1, 2}

For individuals i and j and N sites:

$$d(i,j) = -\log \left( 1 - \frac{1}{N} \sum_{s=1}^N \frac{|g(i,s) - g(j,s)|}{2} \right)$$

$$d(i,j) = 1 * 1.00 = 1.00/2$$

B

A

	0	1	2
0	0	1	0
1	0	0	0
2	0	0	0

## Expected genotype

### Expected value

The expected value of a discrete random variable is the probability-weighted average of all possible outcomes of the experiment.

It is equivalent to the average value when the experiment is performed many times.

Let's go back to the whiteboard!

## Expected genotype

### Expected value

The expected value of a discrete random variable is the probability-weighted average of all possible outcomes of the experiment.

It is equivalent to the average value when the experiment is performed many times.

$$E[X|D] = \sum_{i=1}^N x_i p(X = x_i|D)$$

## Genetic distances from (un)known genotypes

Genotypes are {aa, aA, AA} as {0, 1, 2}

For individuals i and j and N sites:

$$d(i, j) = -\log \left( 1 - \frac{1}{N} \sum_{s=1}^N \frac{|g(i, s) - g(j, s)|}{2} \right)$$

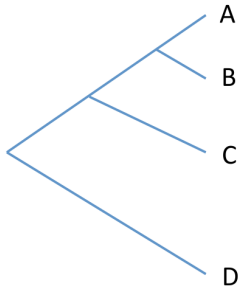
$$E[d(i, j)] = 0 \cdot 0.30 + 1 \cdot 0.50 + 2 \cdot 0.10 + 1 \cdot 0.10 + \dots = 0.80/2$$

B

	0	1	2
0	0.30	0.50	0.10
1	0.10	0	0
2	0	0	0

A

## Genetic distances from unknown genotypes



Genotypes are {aa, aA, AA} as {0, 1, 2}

For individuals i and j and N sites:

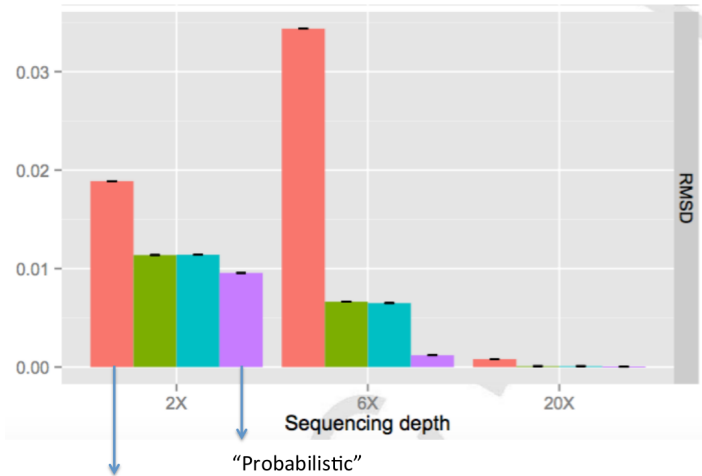
$$d(i, j) = -\log \left( 1 - \frac{1}{N} \sum_{s=1}^N \frac{|g(i, s) - g(j, s)|}{2} \right)$$

Iterate across all possible genotypes

Genotypes probability

$$d(i, j) = -\log \left( 1 - \frac{1}{N} \sum_{s=1}^N \sum_{g(i,s)=0}^2 \sum_{g(j,s)=0}^2 \frac{|g(i, s) - g(j, s)|}{2} * P(g(i, s), g(j, s)) \right)$$

## Genetic distances from unknown genotypes



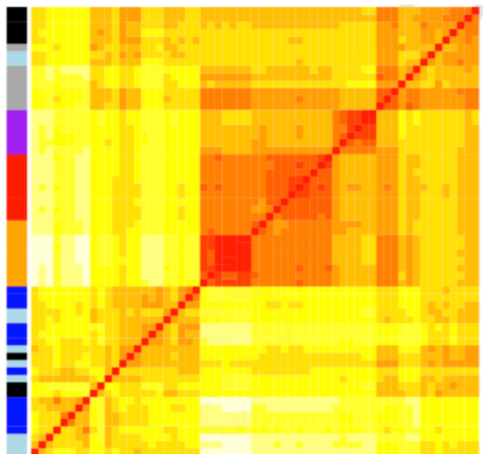
Genotype calling (no prior)

Vieira et al. BJLS 2016



## Clustering from unknown genotypes

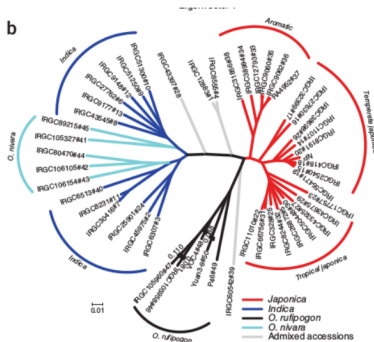
— nivara  
— rufipogon  
— Chinese rufipogon  
— Indica  
— aromatic  
— tropical japonica  
— temperate japonica



Original data: ~2M SNPs  
Here: 5.4M SNPs at 2X

# Population structure

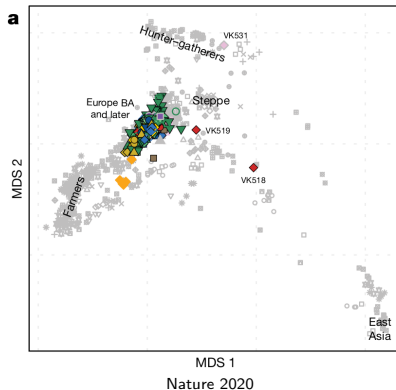
Graphical representation of genetic distance in form of a tree\*.



\* this should not be considered as a proper phylogenetic tree

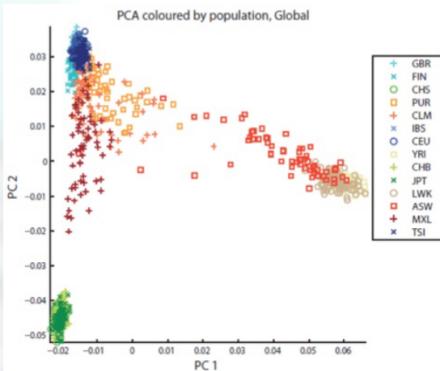
# Population genomics of the Viking world

<https://doi.org/10.1038/s41586-020-2688-8> Ashot Margaryan<sup>1,2,3,†</sup>, Daniel J. Lawson<sup>4,5,†</sup>, Martin Sikora<sup>1,†</sup>, Fernando Racimo<sup>1,†</sup>,

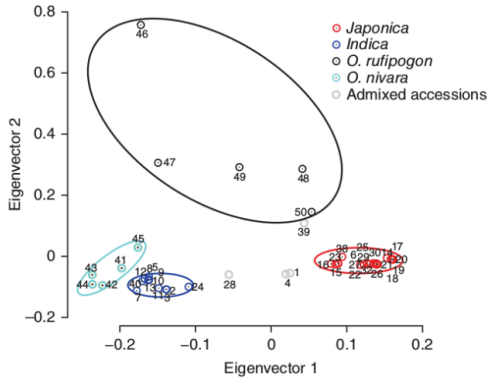


Principal Component Analysis (PCA) is a data reduction method for:

- visualisation
- correction for population stratification
- information on population history and differentiation?



- ❑ An eigenvector decomposition of the covariance matrix is computed. Eigenvectors are then plotted.
- ❑ PCA is a **descriptive** analysis of your dataset, for clustering individuals (and identify population assignment).



$p$  mutations

$$\mathbf{X} = \begin{bmatrix} x_{11} & x_{12} & \dots & x_{1j} & \dots & x_{1p} \\ x_{21} & x_{22} & \dots & x_{2j} & \dots & x_{2p} \\ \vdots & \vdots & \ddots & \vdots & \ddots & \vdots \\ x_{i1} & x_{i2} & \dots & x_{ij} & \dots & x_{ip} \\ \vdots & \vdots & \ddots & \vdots & \ddots & \vdots \\ x_{n1} & x_{n2} & \dots & x_{nj} & \dots & x_{np} \end{bmatrix} \quad n \text{ samples}$$

where  $\mathbf{X}$  is standardised to mean=0 and sd=1 for each column

## V. SOLVING PCA USING EIGENVECTOR DECOMPOSITION

We derive our first algebraic solution to PCA based on an important property of eigenvector decomposition. Once again, the data set is  $\mathbf{X}$ , an  $m \times n$  matrix, where  $m$  is the number of measurement types and  $n$  is the number of samples. The goal is summarized as follows.

Find some orthonormal matrix  $\mathbf{P}$  in  $\mathbf{Y} = \mathbf{P}\mathbf{X}$  such that  $\mathbf{C}_Y \equiv \frac{1}{n}\mathbf{Y}\mathbf{Y}^T$  is a diagonal matrix. The rows of  $\mathbf{P}$  are the *principal components* of  $\mathbf{X}$ .

$$\mathbf{C}_Y = \mathbf{P}\mathbf{C}_X\mathbf{P}^T$$

One interpretation of  $\mathbf{X}$  is the following. Each *row* of  $\mathbf{X}$  corresponds to all measurements of a particular type. Each *column* of  $\mathbf{X}$  corresponds to a set of measurements from one particular trial (this is  $\vec{X}$  from section 3.1). We now arrive at a definition for the *covariance matrix*  $\mathbf{C}_X$ .

$$\mathbf{C}_X \equiv \frac{1}{n}\mathbf{X}\mathbf{X}^T.$$

## Covariance matrix

Genotype (0,1,2)      Allele frequency

$$\text{cov}(i, j) = \frac{1}{(m-1)} \frac{\sum_{s=1}^m (G_s^{(i)} - 2\hat{p}_s)(G_s^{(j)} - 2\hat{p}_s)}{\sqrt{\hat{p}_s(1-\hat{p}_s)}}$$



## Covariance matrix

Genotype (0,1,2)  $\nwarrow$  Allele frequency  $\nearrow$

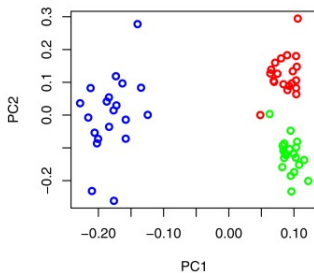
$$\text{cov}(i, j) = \frac{1}{(m-1)} \frac{\sum_{s=1}^m (G_s^{(i)} - 2\hat{p}_s)(G_s^{(j)} - 2\hat{p}_s)}{\sqrt{\hat{p}_s(1-\hat{p}_s)}}$$

Iterate across all genotypes  $\nwarrow$  Weight by their probability  $\nearrow$

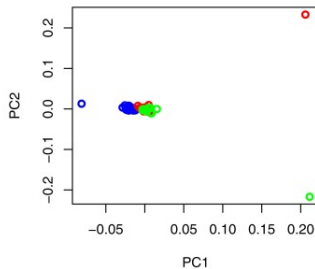
$$\text{cov}\hat{v}_{(i,j)} := \frac{1}{(\sum_{s=1}^m P_{\text{var},s}) - 1} \frac{\sum_{s=1}^m \left( \sum_{G_s^{(i)}=0}^2 \sum_{G_s^{(j)}=0}^2 (G_s^{(i)} - 2\hat{p}_s)(G_s^{(j)} - 2\hat{p}_s) P(G_s^{(i)}|X_s^{(i)}) P(G_s^{(j)}|X_s^{(j)}) \right) P_{\text{var},s}}{\sqrt{\hat{p}_s(1-\hat{p}_s)}}$$

Probability of the site being variable  
(to avoid SNP calling)

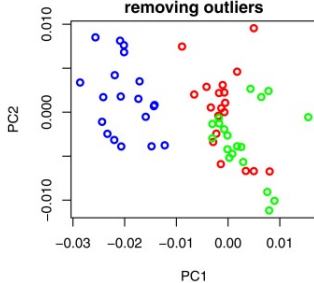
**True genotypes**



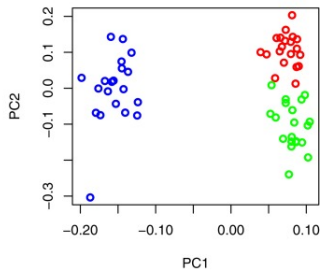
**Calling genotypes  
from posterior probabilities**

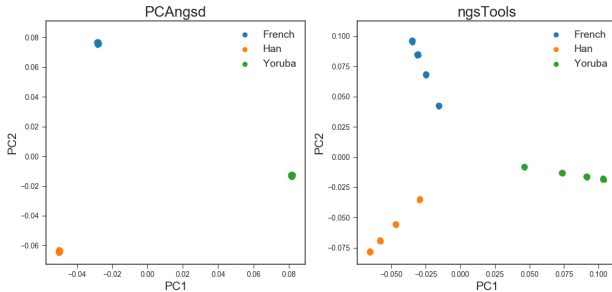


**Calling genotypes  
from posterior probabilities  
removing outliers**



**Without calling genotypes**





Jonas Meisner

$$c_{ij} = \frac{1}{m} \sum_{s=1}^m \frac{\sum_{g_i=0}^2 \sum_{g_j=0}^2 (g_i - 2\hat{p}_s)(g_j - 2\hat{p}_s) P(G_{is} = g_i, G_{js} = g_j \mid X_{is}, X_{js}, \hat{p}_s)}{2\hat{p}_s(1 - \hat{p}_s)}. \quad (3)$$

ngsTools splits up the joint posterior probability,  $P(G_{is}, G_{js} \mid X_{is}, X_{js}, \hat{p}_s)$ , into  $P(G_{is} \mid X_{is}, \hat{p}_s)P(G_{js} \mid X_{js}, \hat{p}_s)$  for  $i \neq j$  by assuming conditional independence between individuals given the estimated population allele frequencies.

$$c_{ij} = \frac{1}{m} \sum_{s=1}^m \frac{\sum_{g_i=0}^2 \sum_{g_j=0}^2 (g_i - 2\hat{p}_s)(g_j - 2\hat{p}_s) P(G_{is} = g_i, G_{js} = g_j \mid X_{is}, X_{js}, \hat{p}_s)}{2\hat{p}_s(1 - \hat{p}_s)}. \quad (3)$$

ngsTools splits up the joint posterior probability,  $P(G_{is}, G_{js} \mid X_{is}, X_{js}, \hat{p}_s)$ , into  $P(G_{is} \mid X_{is}, \hat{p}_s)P(G_{js} \mid X_{js}, \hat{p}_s)$  for  $i \neq j$  by assuming conditional independence between individuals given the estimated population allele frequencies.

The problem with this approach is that the assumption of conditional independence between individuals given the population allele frequency is only valid when there is no population structure. Here we propose a novel approach of estimating the covariance matrix using iteratively estimated individual allele frequencies to update the prior information of the posterior genotype probability. Thereby we condition on the individual allele frequencies as in the clustering-based approaches such as [Pritchard \*et al.\* \(2000\)](#); [Tang \*et al.\* \(2005\)](#); [Alexander \*et al.\* \(2009\)](#); [Skotte \*et al.\* \(2013\)](#).

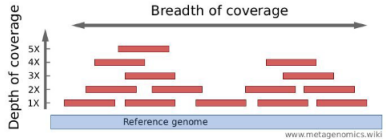
## Introduction to PCAngsd

- Principal Component Analysis of Next-Generation Sequencing Data (PCAngsd)
- Multithreaded Python framework
- Infers population structure using PCA
- Structured populations
- Low and medium sequencing depth
- Published in Meisner and Albechtsen, Genetics, 2018

- DNA data
  - $n$  **diploid** samples
  - $m$  SNPs
- Diallelic ( $G = \{0, 1, 2\}$ )
- Genotype matrix

$$\# \text{ of SNPs} \left\{ \begin{array}{cccccccccccc} 0 & 1 & 2 & \cdots & 0 & 0 & 0 & 0 & \cdots & 0 \\ 2 & 0 & 0 & \cdots & 1 & 0 & 0 & 2 & \cdots & 0 \\ 0 & 0 & 0 & \cdots & 0 & 0 & 2 & 0 & \cdots & 2 \\ 0 & 0 & 0 & \cdots & 0 & 0 & 0 & 0 & \cdots & 1 \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & 0 & 2 & \cdots & 0 & 2 & 0 & 0 & \cdots & 0 \end{array} \right.$$

- Sequencing data
  - Low (<5X) and medium (<15X) sequencing depth
  - Genotype uncertainty



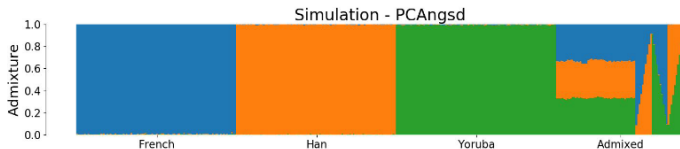
- Genotype likelihoods
  - Retain information of sequencing process
  - Take uncertainty into account

$$P(X_{is} | G = g)$$

Jonas Meisner



- Population allele frequencies
  - Average in discrete populations



- PCA
  - Dimension reduction method
  - Construct axes of genetic variation
- Individual allele frequencies
  - Infer underlying sampling parameter

$$g_{is} \sim \text{Binomial}(2, \pi_{is})$$

Jonas Meisner

- Posterior genotype probability

$$P(G = g \mid X_{is}, p_s) = \frac{P(X_{is} \mid G = g)P(G = g \mid p_s)}{\sum_{g'=0}^2 P(X_{is} \mid G = g')P(G = g' \mid p_s)}$$

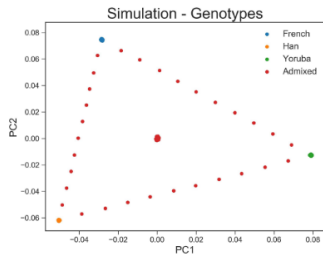
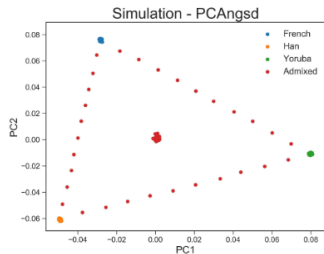
- Low-rank SVD reconstruction

$$\mathbb{E}[\mathbf{G} \mid \mathbf{X}] \approx \mathbf{U}_{1:D} \mathbf{\Delta}_{1:D} \mathbf{V}_{1:D}$$

$$\pi_{is} = \frac{1}{2} \mathbf{U}_{[i,1:D]} \mathbf{\Delta}_{1:D} \mathbf{V}_{[s,1:D]}$$

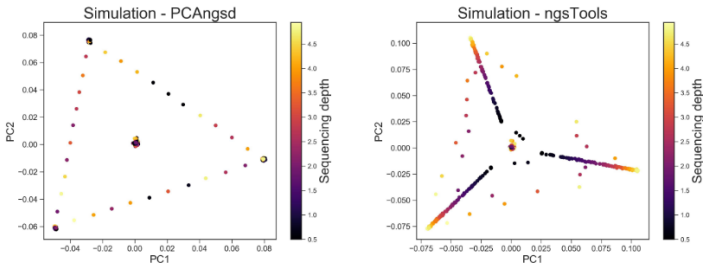
- Update prior information and iterate!

## Inferring population structure (0.5 - 5X)



Jonas Meisner

Existing methods are biased by sequencing depth!



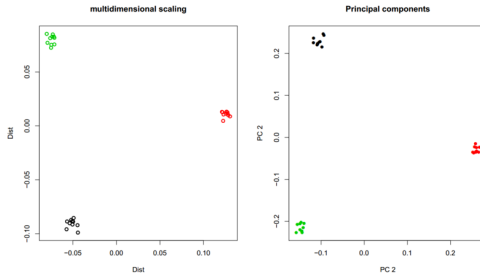
Jonas Meisner

PCAngsd is the current state-of-the-art for PCA from low-depth data

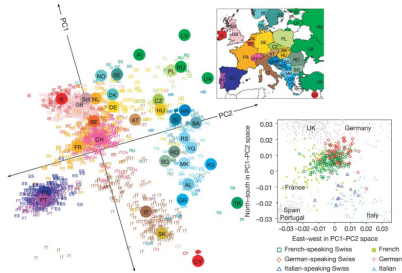
## PCA/MDS

with random sampling of read for IBS (for MDS) or covariance (for PCA) matrix.

It works even with very low depth data (1X), but it requires low error rate and known polymorphic sites.



## Beyond PCA?



Novembre et al. 2008

PCA represents genetic drift as allele frequency differences between populations.

There is a good theory behind it but possibly still ongoing issues on interpretability and quantification of admixture/gene flow.

# Admixture analyses

For admixed populations: assign ancestry proportion to each individual from multiple (ancestral) populations.

$$f_{A.1} = 0.6$$

AA: observed genotype

$$f_{A.2} = 0.2$$

$$\Pr(AA \mid \text{pop}=1) = ?$$

$$\Pr(AA \mid \text{pop}=2) = ?$$

# Admixture analyses

For admixed populations: assign ancestry proportion to each individual from multiple (ancestral) populations.


$$f_{A.1} = 0.6$$

AA: observed genotype

$$f_{A.2} = 0.2$$

$$\Pr(AA \mid \text{pop}=1) = ?$$

$$\Pr(AA \mid \text{pop}=2) = ?$$



Assuming HWE

$$\begin{aligned} P(AA) &= f^2 \\ P(AG) &= 2 * f * (1-f) \\ P(GG) &= (1-f)^2 \end{aligned}$$



# Admixture analyses

For admixed populations: assign ancestry proportion to each individual from multiple (ancestral) populations.

$$f_{A.1} = 0.6$$

$$f_{A.2} = 0.2$$

$$\Pr(AA \mid \text{pop}=1) = f_{A.1}^2 / (f_{A.1}^2 + f_{A.2}^2) = 0.90$$

$$\Pr(AA \mid \text{pop}=2) = f_{A.2}^2 / (f_{A.1}^2 + f_{A.2}^2) = 0.10$$

$$\Pr(AG \mid \text{pop}=1) = ?$$

$$\Pr(AG \mid \text{pop}=2) = ?$$

} Assuming HWE  
P(AA)= $f^2$   
P(AG)= $2*f*(1-f)$   
P(GG)=( $1-f$ )<sup>2</sup>

# Admixture analyses

For admixed populations: assign ancestry proportion to each individual from multiple (ancestral) populations.

$$f_{A.1} = 0.6$$

$$f_{A.2} = 0.2$$

$$\Pr(AA \mid \text{pop}=1) = f_{A.1}^2 / (f_{A.1}^2 + f_{A.2}^2) = 0.90$$

$$\Pr(AA \mid \text{pop}=2) = f_{A.2}^2 / (f_{A.1}^2 + f_{A.2}^2) = 0.10$$

$$\Pr(AG \mid \text{pop}=1) = 2 * f_{A.1} * (1 - f_{A.1}) / \dots = 0.60$$

$$\Pr(AG \mid \text{pop}=2) = \dots = 0.40$$

Assuming HWE  
 $P(AA)=f^2$   
 $P(AG)=2*f*(1-f)$   
 $P(GG)=(1-f)^2$

Bayesian (STRUCTURE) or Maximum Likelihood (ADMIXTURE) approaches.

Genomes of  $n$  samples,  $p$  mutations stored in  $n \times p$  matrix  $\mathbf{X}$

$x_{ij}$  contains the number of non-reference mutations for individual  $i$  at position (“SNP”)  $j$ , so that  $x_{ij} \in \{0,1,2\}$

Assume there are  $k = 1 \dots c$  ancestral populations. The frequency  $\mu_{jk}$  of mutation  $j$  in population  $k$  is unknown

Goal: infer the fraction of ancestry  $\alpha_{ik}$  for individual  $i$  and all values of  $k$ , subject to  $\sum_{k=1}^K \alpha_{ik} = 1$

We need to find  $\alpha_{ik}$  and  $\mu_{jk}$  that maximize  $P(\text{DATA} \mid \alpha_{ik}, \mu_{jk} \forall i, j, k)$ , which is proportional to

$$\prod_{i=1}^n \prod_{j=1}^p [\sum_{k=1}^c \alpha_{ik} \mu_{jk}]^{x_{ij}} [\sum_{k=1}^c \alpha_{ik} (1 - \mu_{jk})]^{2-x_{ij}}$$

STRUCTURE: Pritchard et al. Genetics, 2000

LDA: Blei et al. JMLR 2003, LDA [notes](#)

Adaptation of example by A. Price

Pier Palamara

Figure 2 displays ancestry plots for 1000 individuals from 10 populations across  $K=2$  to  $K=7$ . The populations are color-coded: TRJ (red), TEJ (blue), ARO (yellow), AUS (green), IND (purple), *O. nivara* (dark blue), and *O. rufipogon* (light blue). The plots show the proportion of ancestry from each population for each individual. The plots indicate that as  $K$  increases, the number of distinct ancestral components increases, and the proportion of ancestry from each population becomes more complex.

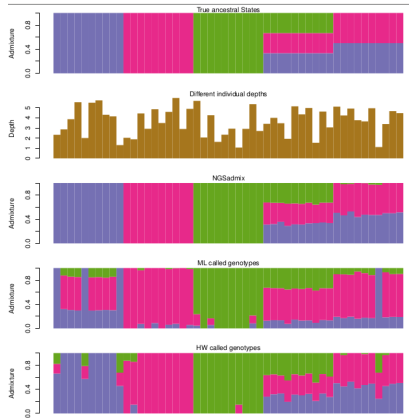
When the genotypes are observed, assuming that sites are independent, the likelihood is written as

$$p(G|Q, F) = \prod_{j=1}^M \prod_{i=1}^N p(G_{ij}|Q, F) = \prod_{j=1}^M \prod_{i=1}^N p(G_{ij}|h^{ij}). \quad (5)$$

If the sites are not independent, then this is a composite likelihood that will still have consistent estimates. This likelihood corresponds to the likelihood used in [Tang \*et al.\* \(2005\)](#) and [Alexander \*et al.\* \(2009\)](#) and will be used when dealing with called genotypes.

When using NGS data, the genotypes are not observed and we instead work with genotype likelihoods. The above likelihood is extended by summing over all possible genotypes:

$$\begin{aligned} p(X|Q, F) &= \prod_{j=1}^M \prod_{i=1}^N p(X_{ij}|Q, F) = \prod_{j=1}^M \prod_{i=1}^N p(X_{ij}|h^{ij}) \\ &= \prod_{j=1}^M \prod_{i=1}^N \sum_{G_{ij} \in \{0,1,2\}} p(X_{ij}|G_{ij}) p(G_{ij}|h^{ij}). \end{aligned} \quad (6)$$



in practice, PCAngsd extends its framework to admixture analysis

## Intended Learning Outcomes

At the end of this session you are be able to

- understand the theory underlying distance and covariance matrices
- appreciate how to extended such theory to low-coverage data
- acknowledge the process of inferring population structure and admixture from sequencing data
- implement a pipeline in ANGSD to perform the aforementioned analyses