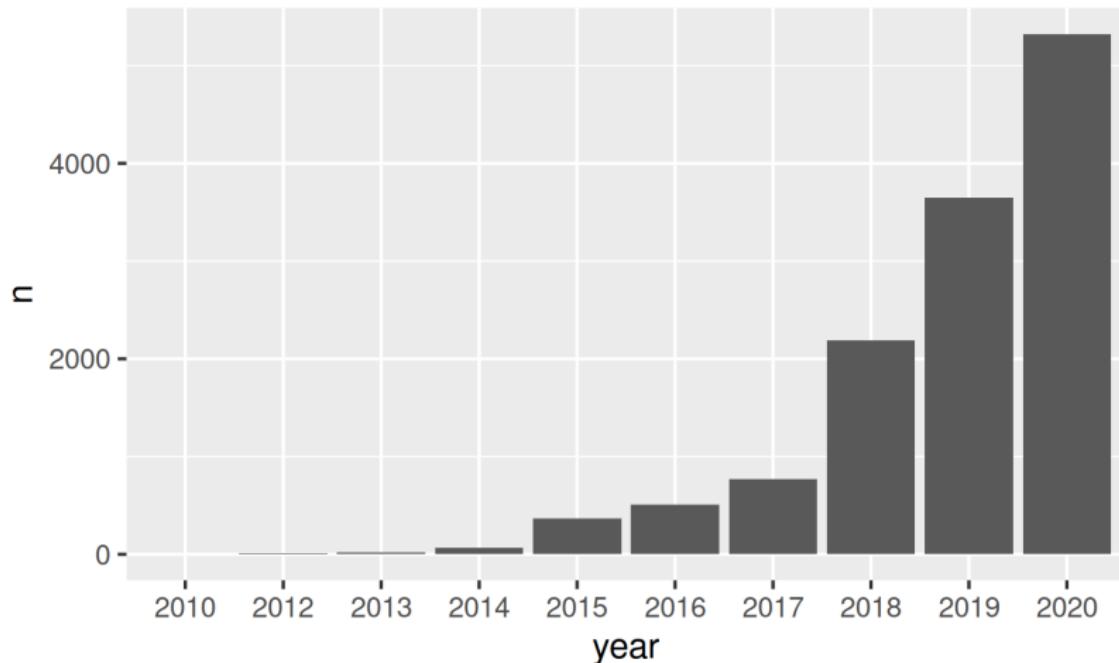


# F-statistics and PCA

Benjamin Peter

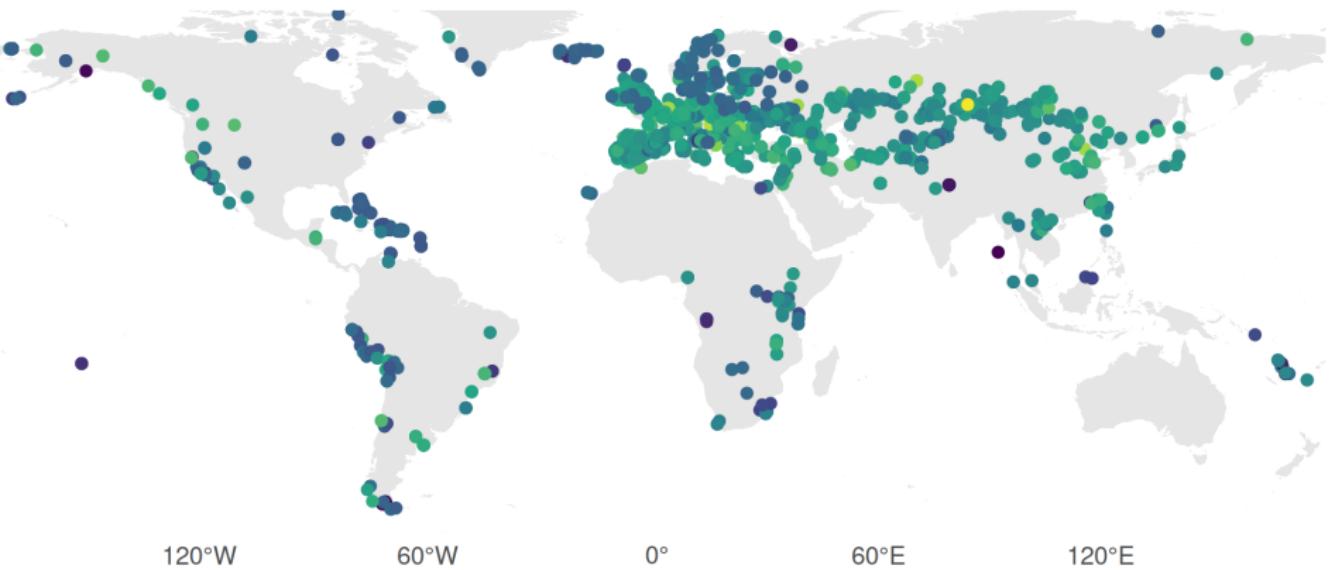
April 22, 2021

# Population structure and ancient DNA



<https://reich.hms.harvard.edu/>

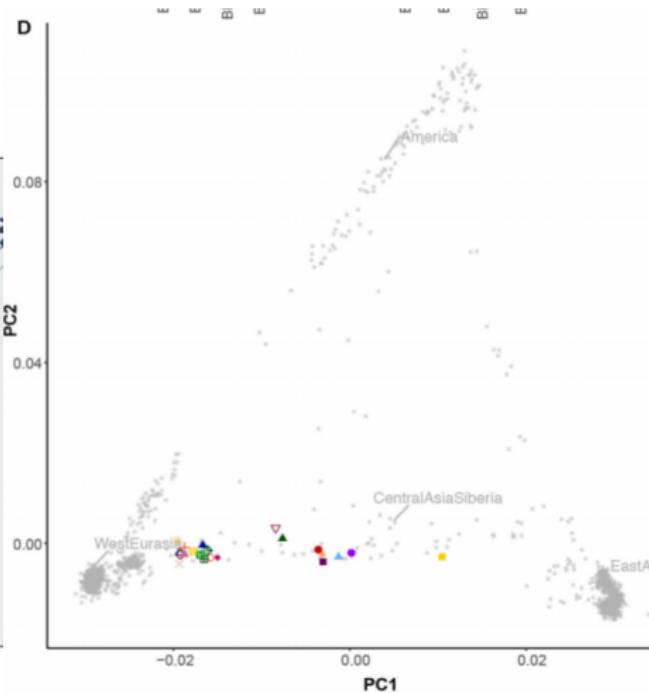
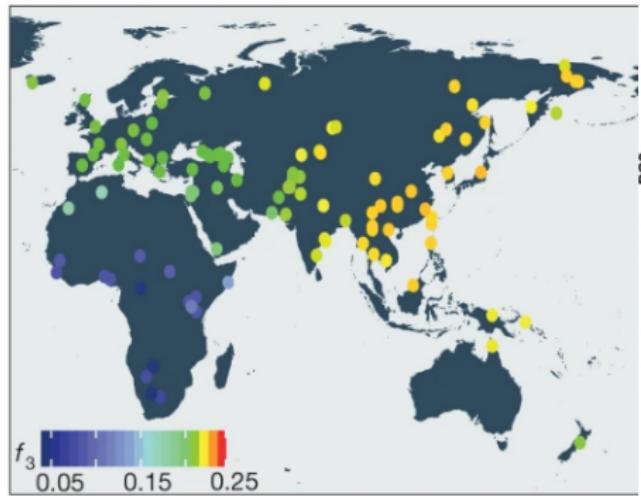
# Population structure and ancient DNA



<https://reich.hms.harvard.edu/>

# PCA and $F$ -statistics

$f_3(\text{Mbuti}; \text{IUP Bacho Kiro}, X)$



Hajdinjak et al. 2021

# Goals of this talk

- Technical & Conceptual Background
- Establish conceptual links between frameworks
  - ① How can we interpret PCA in context of  $F$ -stats?
  - ② How can we interpret  $F$ -stats in the context of PCA?
- (Use established links to improve data interpretation)

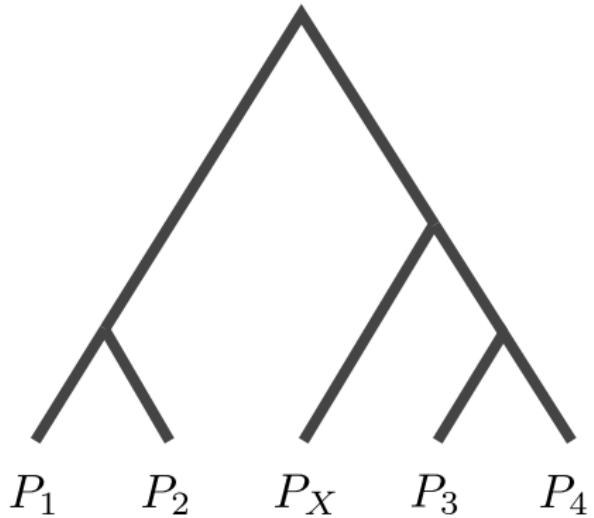
# Goals of this talk

- Technical & Conceptual Background
- Establish conceptual links between frameworks
  - ① How can we interpret PCA in context of  $F$ -stats?
  - ② How can we interpret  $F$ -stats in the context of PCA?
- (Use established links to improve data interpretation)

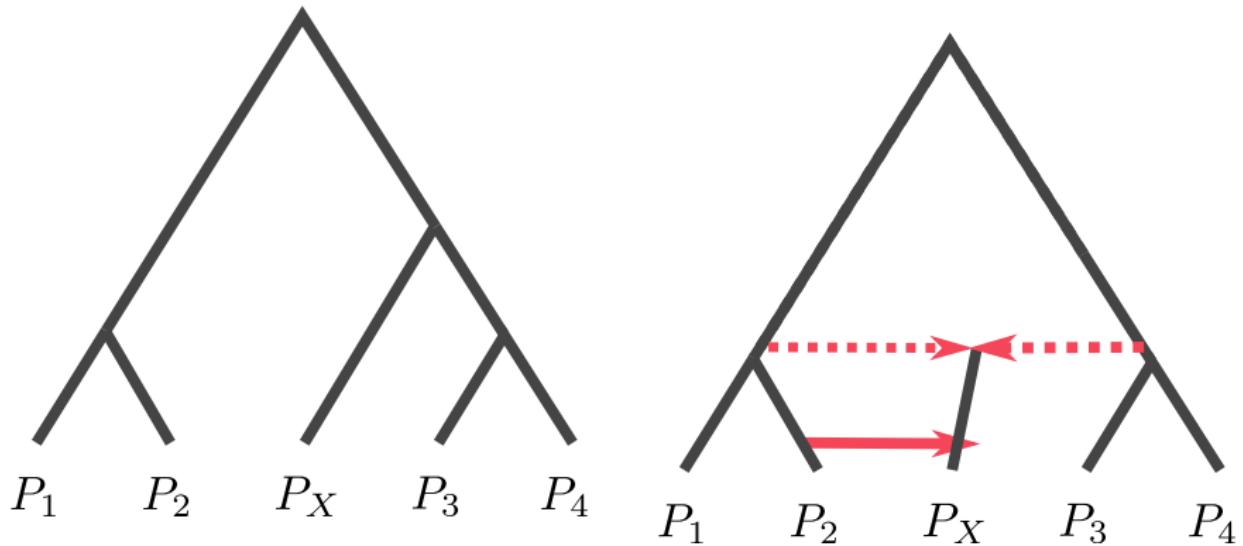
## Focus on intuition

Some details in terms of estimation, normalization, missing data will be glossed over

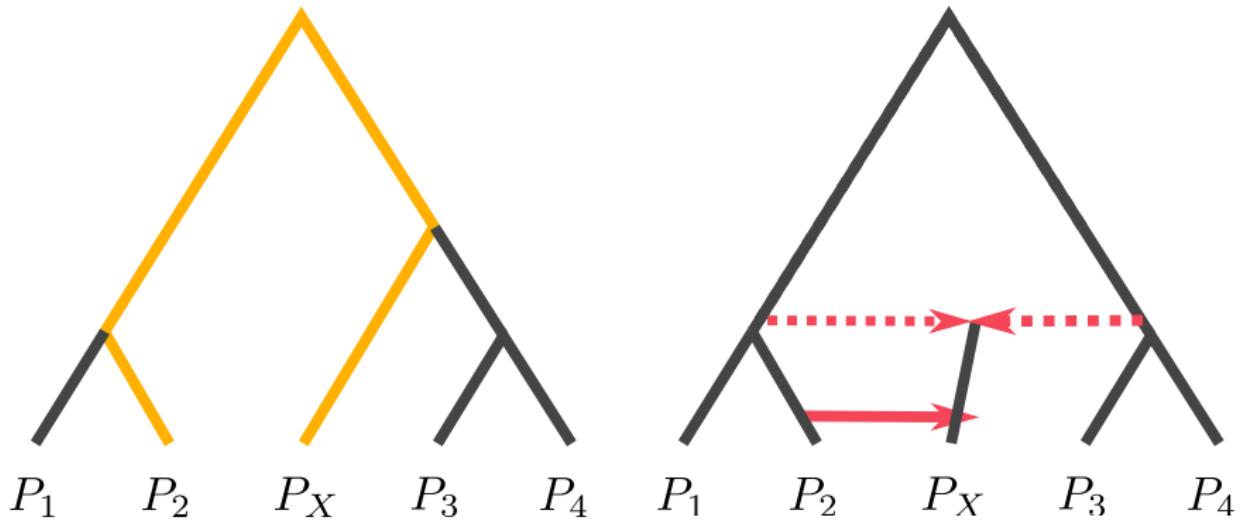
# Trees vs. Admixture Graphs



# Trees vs. Admixture Graphs



# Trees vs. Admixture Graphs

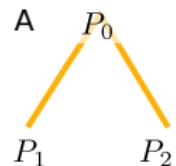


# $F$ -statistics

## Definition

$$F_2(X_1, X_2) = \sum_l (X_{il} - X_{jl})^2 - H_1 - H_2$$

Branch  
length

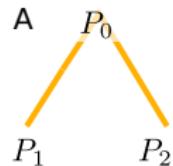


# $F$ -statistics

## Definition

$$F_2(X_1, X_2) = \sum_l (X_{il} - X_{jl})^2 - H_1 - H_2$$

Branch length



$$2F_3(X_x; X_1, X_2) = F_2(X_x, X_1) + F_2(X_x, X_2) - F_2(X_1, X_2)$$

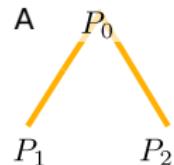


# $F$ -statistics

## Definition

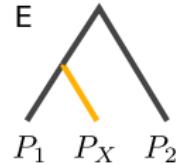
$$F_2(X_1, X_2) = \sum_l (X_{il} - X_{jl})^2 - H_1 - H_2$$

Branch length



$$2F_3(X_x; X_1, X_2) = F_2(X_x, X_1) + F_2(X_x, X_2) - F_2(X_1, X_2)$$

$$F_3(X_x; X_1, X_2) = \sum_l (X_{xl} - X_{1l})(X_{xl} - X_{2l}) - H_X$$

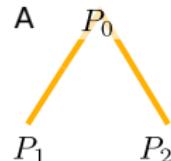


# $F$ -statistics

## Definition

$$F_2(X_1, X_2) = \sum_l (X_{il} - X_{jl})^2 - H_1 - H_2$$

Branch length

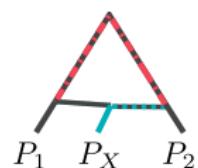


$$2F_3(X_x; X_1, X_2) = F_2(X_x, X_1) + F_2(X_x, X_2) - F_2(X_1, X_2)$$

$$F_3(X_x; X_1, X_2) = \sum_l (X_{xl} - X_{1l})(X_{xl} - X_{2l}) - H_X$$



“Admixture”- $F_3$ -statistic: If data is generated by a tree-like relationship,  $F_3(P_X; P_1, P_2) \geq 0$



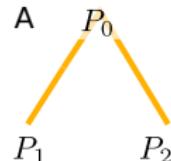
Patterson et al. 2012; Peter 2016

# $F$ -statistics

## Definition

$$F_2(X_1, X_2) = \sum_l (X_{il} - X_{jl})^2 - H_1 - H_2$$

Branch length



$$2F_3(X_x; X_1, X_2) = F_2(X_x, X_1) + F_2(X_x, X_2) - F_2(X_1, X_2)$$

$$F_3(X_x; X_1, X_2) = \sum_l (X_{xl} - X_{1l})(X_{xl} - X_{2l}) - H_X$$



“Outgroup”- $F_3$ -statistic: Most similar pops have highest  $F_3(P_2; P_X, P_1)$

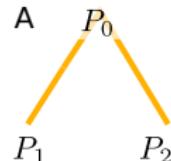


# $F$ -statistics

## Definition

$$F_2(X_1, X_2) = \sum_l (X_{il} - X_{jl})^2 - H_1 - H_2$$

Branch length



$$2F_3(X_x; X_1, X_2) = F_2(X_x, X_1) + F_2(X_x, X_2) - F_2(X_1, X_2)$$

$$F_3(X_x; X_1, X_2) = \sum_l (X_{xl} - X_{1l})(X_{xl} - X_{2l}) - H_X$$



$$F_4^{(B)}(X_1; X_2; X_3, X_4) = \sum_l (X_{1l} - X_{3l})(X_{2l} - X_{4l})$$



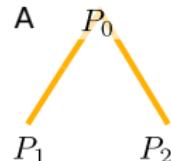
Patterson et al. 2012; Peter 2016

# $F$ -statistics

## Definition

$$F_2(X_1, X_2) = \sum_l (X_{il} - X_{jl})^2 - H_1 - H_2$$

Branch length



$$2F_3(X_x; X_1, X_2) = F_2(X_x, X_1) + F_2(X_x, X_2) - F_2(X_1, X_2)$$

$$F_3(X_x; X_1, X_2) = \sum_l (X_{xl} - X_{1l})(X_{xl} - X_{2l}) - H_X$$

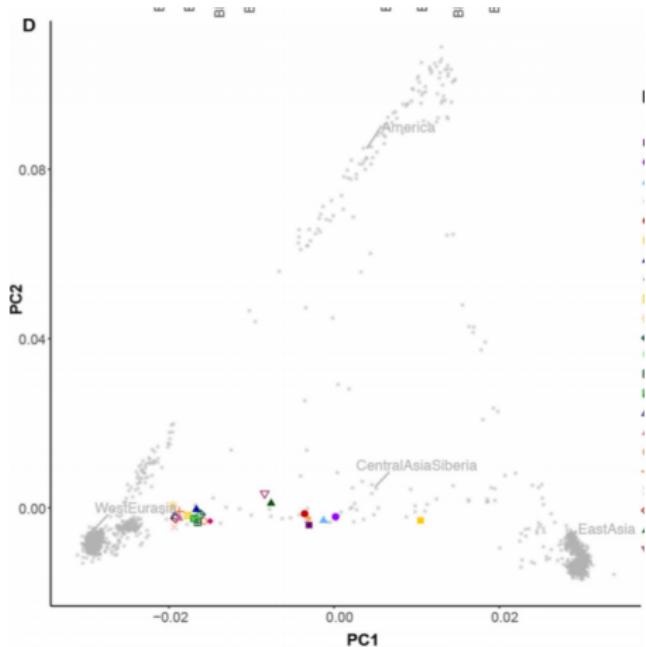


$$F_4^{(T)}(X_1; X_2; X_3, X_4) == \sum_l (X_{1l} - X_{2l})(X_{3l} - X_{4l})$$



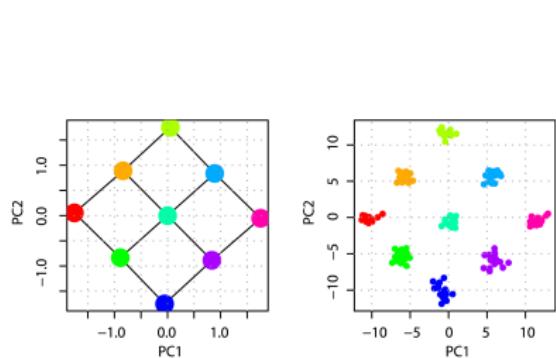
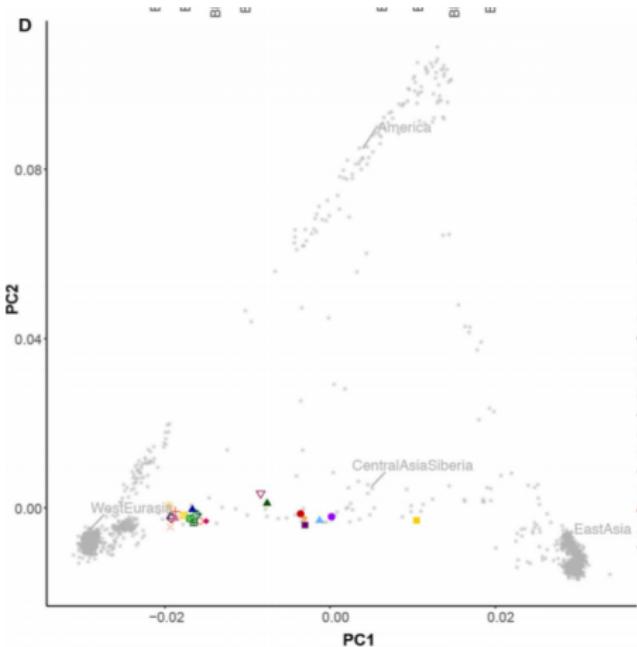
Patterson et al. 2012; Peter 2016

# Principal Component Analysis



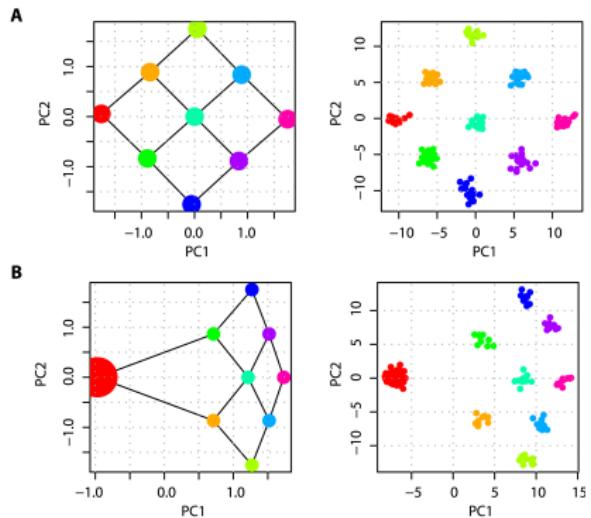
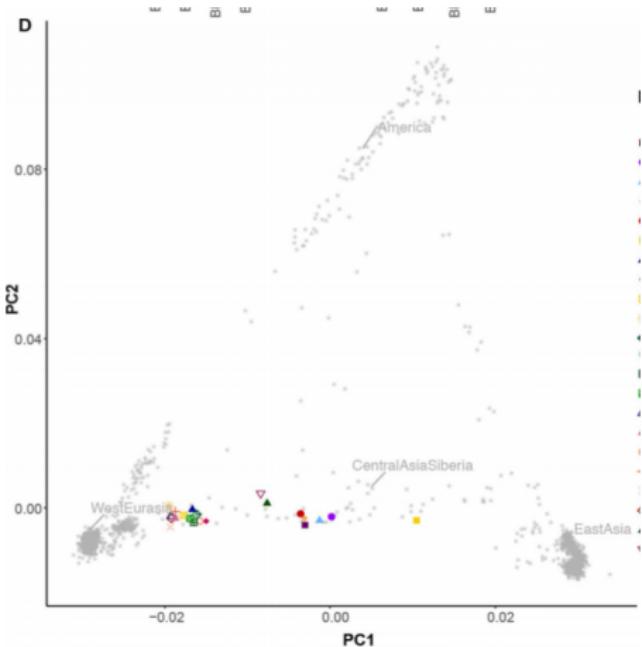
McVean, 2009

# Principal Component Analysis



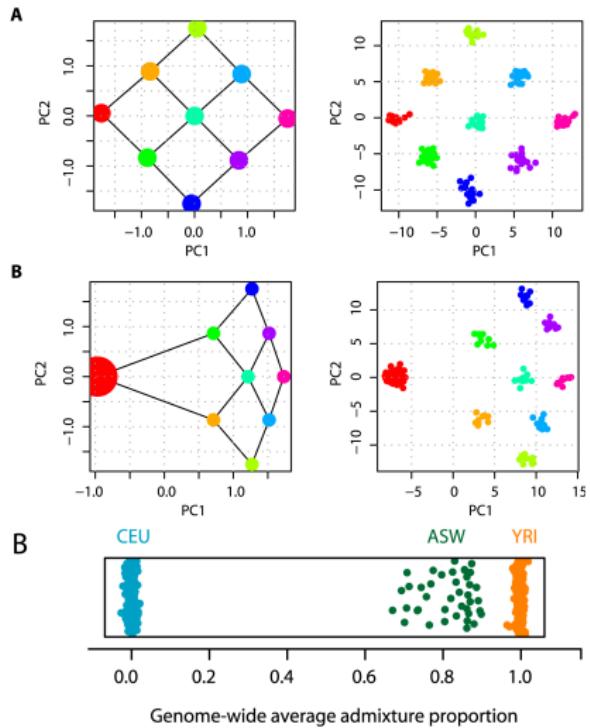
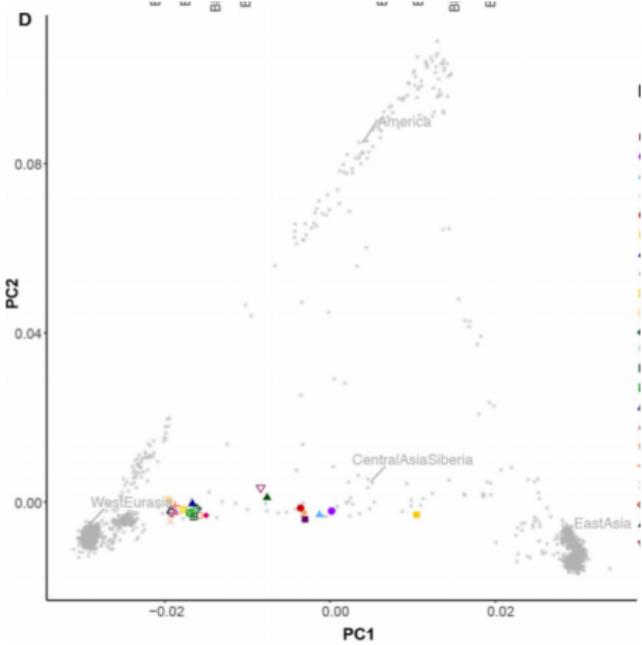
McVean, 2009

# Principal Component Analysis



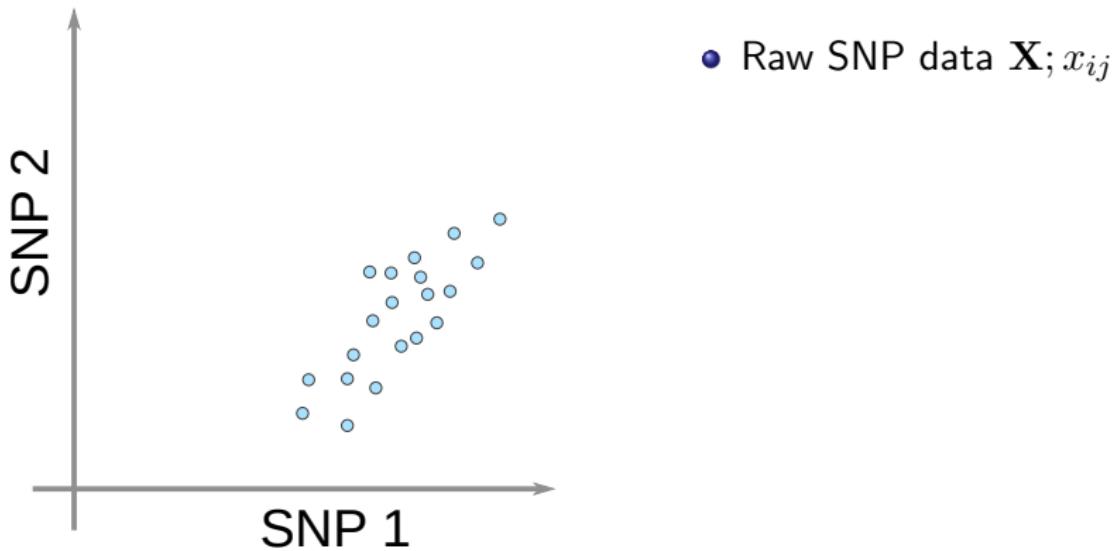
McVean, 2009

# Principal Component Analysis

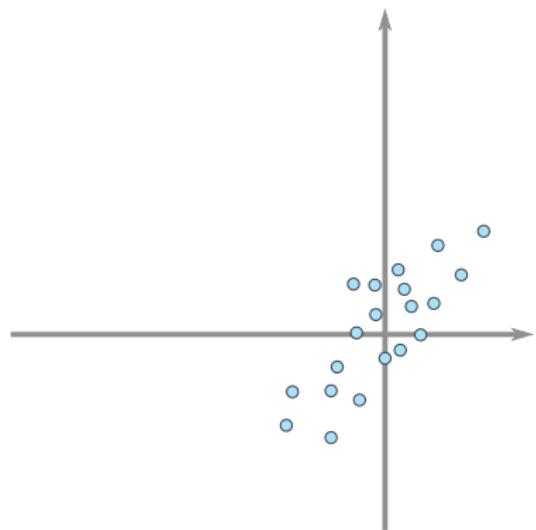


McVean, 2009

# Principal Component Analysis



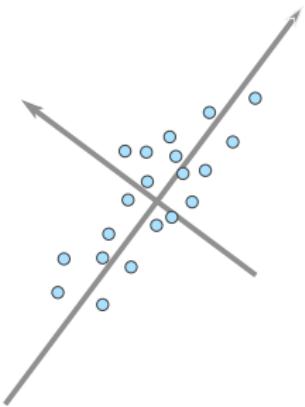
# Principal Component Analysis



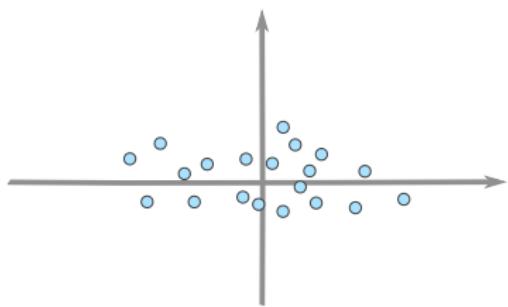
- Raw SNP data  $\mathbf{X}; x_{ij}$
- Centering  
 $\mathbf{Y} = \mathbf{C}\mathbf{X}; y_{ij} = x_{ij} - \mu_j$

# Principal Component Analysis

- Raw SNP data  $\mathbf{X}; x_{ij}$
- Centering  
 $\mathbf{Y} = \mathbf{C}\mathbf{X}; y_{ij} = x_{ij} - \mu_j$
- Rotation  $\mathbf{Y} = \underbrace{\mathbf{P}}_{\text{PCs}} \underbrace{\mathbf{L}}_{\text{Rotation}}$



# Principal Component Analysis



- Raw SNP data  $\mathbf{X}; x_{ij}$
- Centering  
 $\mathbf{Y} = \mathbf{C}\mathbf{X}; y_{ij} = x_{ij} - \mu_j$
- Rotation  $\mathbf{Y} = \mathbf{P}\mathbf{L}$

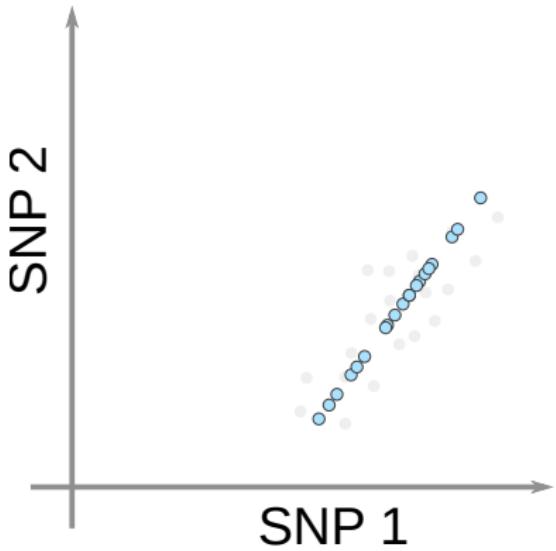
# Principal Component Analysis



- Raw SNP data  $\mathbf{X}; x_{ij}$
- Centering  
 $\mathbf{Y} = \mathbf{C}\mathbf{X}; y_{ij} = x_{ij} - \mu_j$
- Rotation  $\mathbf{Y} = \mathbf{P}\mathbf{L}$

- Truncation  $\hat{\mathbf{P}} = \begin{pmatrix} \mathbf{p}_1 \\ \mathbf{p}_2 \\ \vdots \\ \mathbf{p}_n \end{pmatrix}$

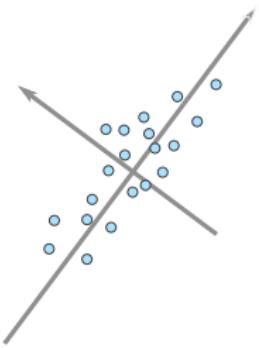
# Principal Component Analysis



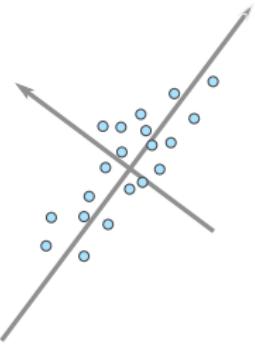
- Raw SNP data  $\mathbf{X}; x_{ij}$
- Centering  $\mathbf{Y} = \mathbf{C}\mathbf{X}; y_{ij} = x_{ij} - \mu_j$
- Rotation  $\mathbf{Y} = \mathbf{PL}$
- Truncation  $\hat{\mathbf{P}} = \begin{pmatrix} \mathbf{p}_1 \\ \mathbf{p}_2 \\ \vdots \\ \mathbf{p}_n \end{pmatrix}$
- Approximation  $\hat{\mathbf{Y}} = \hat{\mathbf{P}}\hat{\mathbf{L}}$

# How to find PCs

- Singular Value Decomposition:  
 $\mathbf{Y} = (\mathbf{U}\mathbf{D})\mathbf{L} = \mathbf{P}\mathbf{L}$

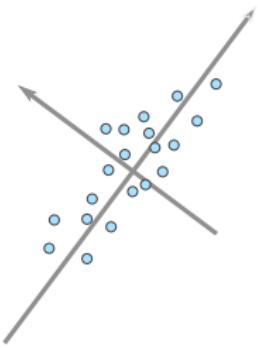


# How to find PCs



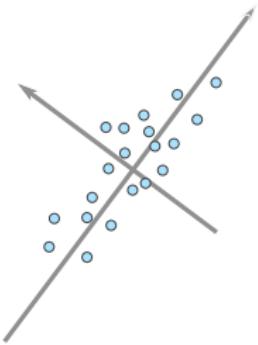
- Singular Value Decomposition:  
 $\mathbf{Y} = (\mathbf{U}\mathbf{D})\mathbf{L} = \mathbf{P}\mathbf{L}$
- Eigendecomposition of  $\mathbf{Y}\mathbf{Y}^T$ :  
 $\mathbf{Y}\mathbf{Y}^T = \mathbf{U}\mathbf{D}^2\mathbf{U}^T = \mathbf{P}\mathbf{P}^T$

# How to find PCs



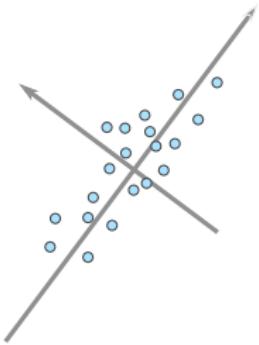
- Singular Value Decomposition:  
 $\mathbf{Y} = (\mathbf{UD})\mathbf{L} = \mathbf{PL}$
- Eigendecomposition of  $\mathbf{YY}^T$ :  
 $\mathbf{YY}^T = \mathbf{UD}^2\mathbf{U}^T = \mathbf{PP}^T$
- $(\mathbf{YY}^T)_{ij}$

# How to find PCs



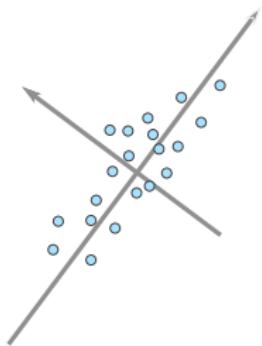
- Singular Value Decomposition:  
 $\mathbf{Y} = (\mathbf{UD})\mathbf{L} = \mathbf{PL}$
- Eigendecomposition of  $\mathbf{YY}^T$ :  
 $\mathbf{YY}^T = \mathbf{UD}^2\mathbf{U}^T = \mathbf{PP}^T$
- $(\mathbf{YY}^T)_{ij} = \sum_l (x_{il} - \mu_l)(x_{jl} - \mu_l)$

# How to find PCs



- Singular Value Decomposition:  
 $\mathbf{Y} = (\mathbf{UD})\mathbf{L} = \mathbf{PL}$
- Eigendecomposition of  $\mathbf{YY}^T$ :  
 $\mathbf{YY}^T = \mathbf{UD}^2\mathbf{U}^T = \mathbf{PP}^T$
- $(\mathbf{YY}^T)_{ij} = \sum_l (x_{il} - \mu_l)(x_{jl} - \mu_l)$
- $F_3(X_x; X_1, X_2) = \sum_l (X_{1l} - X_{xl})(X_{2l} - X_{xl})$

# How to find PCs

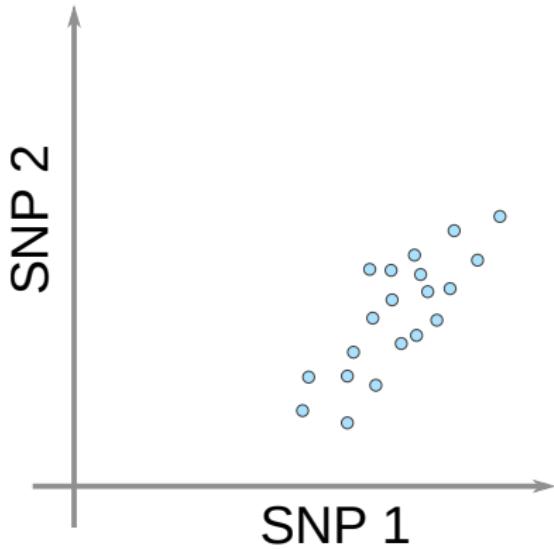


- Singular Value Decomposition:  
 $\mathbf{Y} = (\mathbf{UD})\mathbf{L} = \mathbf{PL}$
- Eigendecomposition of  $\mathbf{YY}^T$ :  
 $\mathbf{YY}^T = \mathbf{UD}^2\mathbf{U}^T = \mathbf{PP}^T$
- $(\mathbf{YY}^T)_{ij} = \sum_l (x_{il} - \mu_l)(x_{jl} - \mu_l)$
- $F_3(X_x; X_1, X_2) = \sum_l (X_{1l} - X_{xl})(X_{2l} - X_{xl})$
- $(\mathbf{YY}^T)_{ij} = F_3(\boldsymbol{\mu}; \mathbf{X}_i, \mathbf{X}_j)$

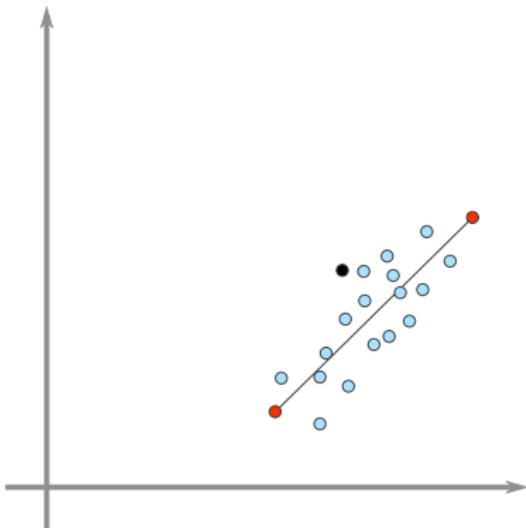
## Observation

PCA is equivalent to outgroup- $F_3$ -analysis with sample mean as outgroup

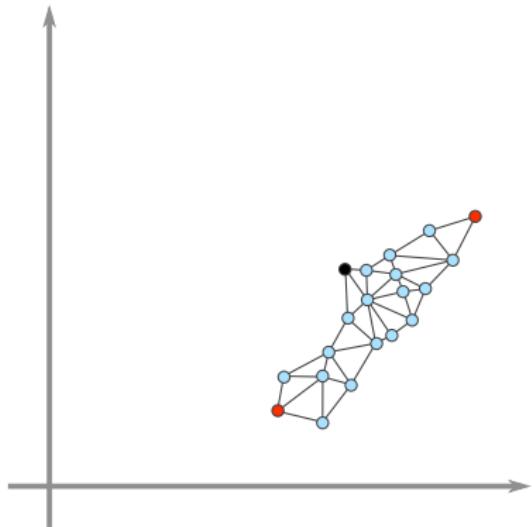
# (metric) Multi-Dimensional Scaling (MDS)



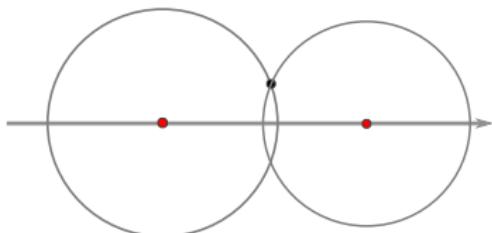
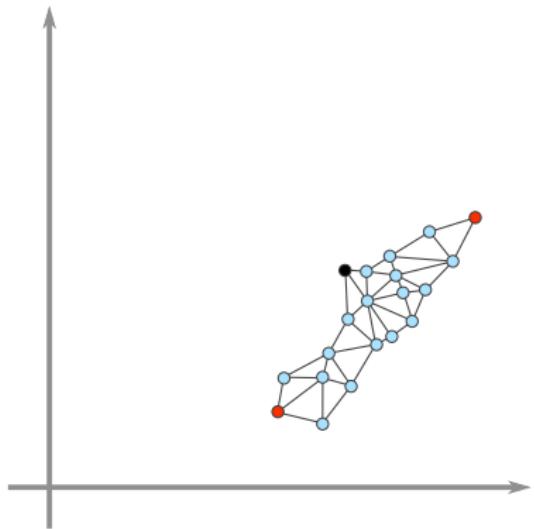
# (metric) Multi-Dimensional Scaling (MDS)



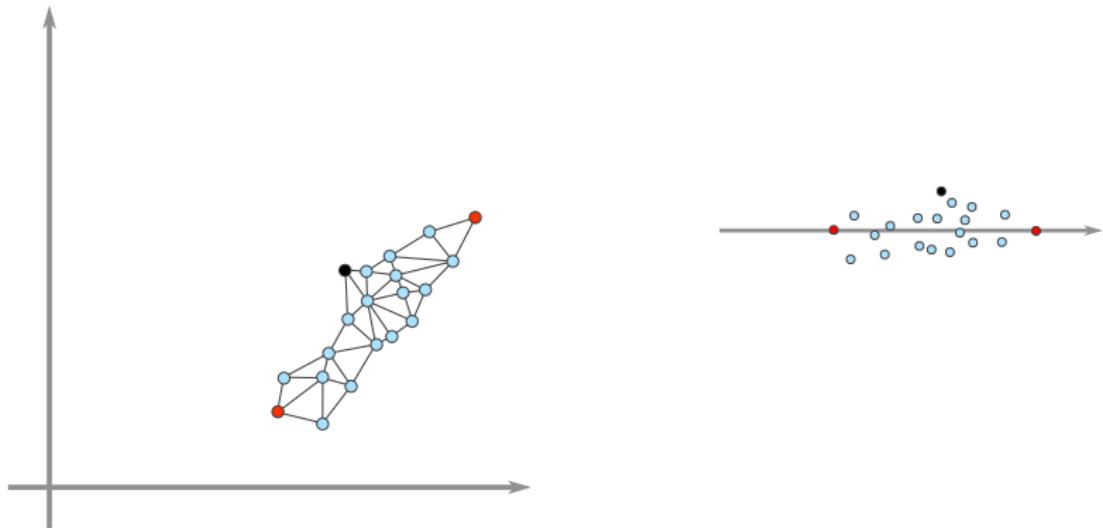
# (metric) Multi-Dimensional Scaling (MDS)



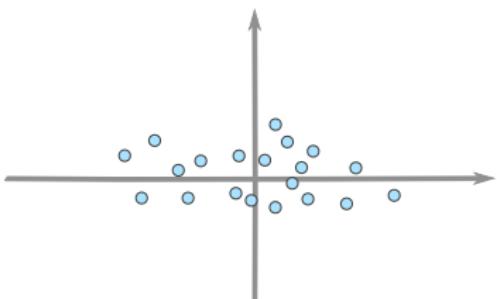
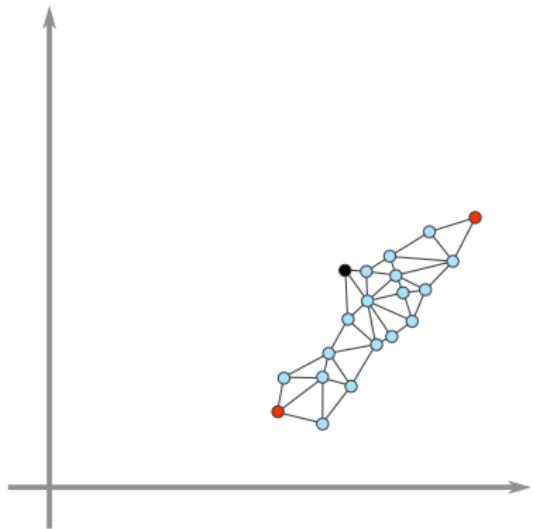
# (metric) Multi-Dimensional Scaling (MDS)



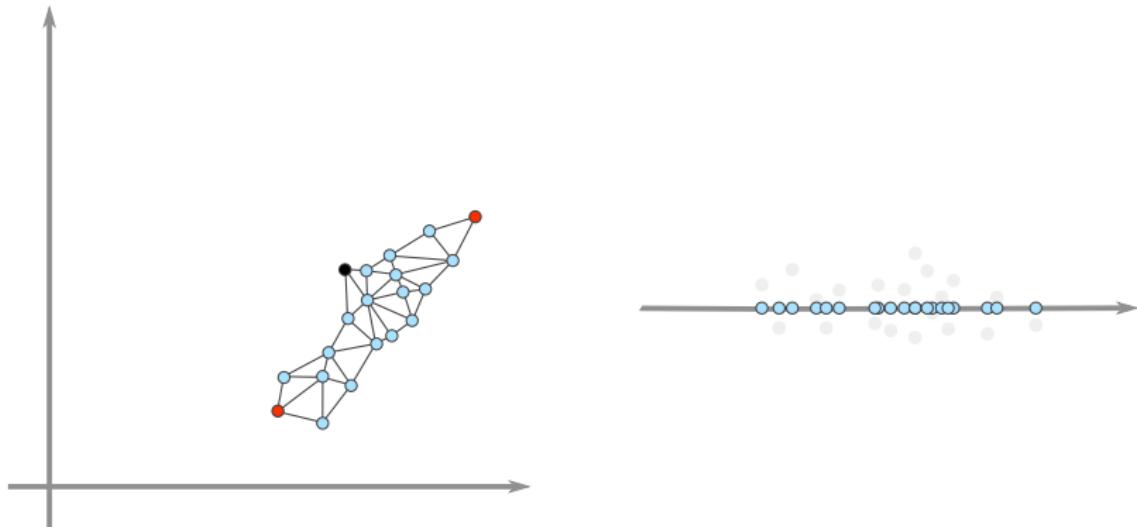
# (metric) Multi-Dimensional Scaling (MDS)



# (metric) Multi-Dimensional Scaling (MDS)



# (metric) Multi-Dimensional Scaling (MDS)



# PCA is MDS on $F_2$

- PCA is decomposition of Covariance matrix:  $\mathbf{Y}\mathbf{Y}^T$

# PCA is MDS on $\mathbf{F}_2$

- PCA is decomposition of Covariance matrix:  $\mathbf{Y}\mathbf{Y}^T$
- Consider  $\mathbf{F}_2$ ;  $f_{ij} = F_2(X_i, X_j) = \sum_l (X_{li}^2 - X_{lj}^2)$

# PCA is MDS on $\mathbf{F}_2$

- PCA is decomposition of Covariance matrix:  $\mathbf{Y}\mathbf{Y}^T$
- Consider  $\mathbf{F}_2$ ;  $f_{ij} = F_2(X_i, X_j) = \sum_l (X_{li}^2 - X_{lj}^2)$
- Consider  $\mathbf{F}_2$ ;  $f_{ij} = \sum_l X_{li}^2 + \sum_l X_{lj}^2 - 2 \sum_l X_{li}X_{lj}$

# PCA is MDS on $\mathbf{F}_2$

- PCA is decomposition of Covariance matrix:  $\mathbf{Y}\mathbf{Y}^T$
- Consider  $\mathbf{F}_2$ ;  $f_{ij} = F_2(X_i, X_j) = \sum_l (X_{li}^2 - X_{lj}^2)$
- Consider  $\mathbf{F}_2$ ;  $f_{ij} = \sum_l X_{li}^2 + \sum_l X_{lj}^2 - 2 \sum_l X_{li}X_{lj}$
- MDS is Eigendecomposition of  $-\frac{1}{2}\mathbf{C}\mathbf{F}_2\mathbf{C}$

# PCA is MDS on $\mathbf{F}_2$

- PCA is decomposition of Covariance matrix:  $\mathbf{Y}\mathbf{Y}^T$
- Consider  $\mathbf{F}_2$ ;  $f_{ij} = F_2(X_i, X_j) = \sum_l (X_{li}^2 - X_{lj}^2)$
- Consider  $\mathbf{F}_2$ ;  $f_{ij} = \sum_l X_{li}^2 + \sum_l X_{lj}^2 - 2 \sum_l X_{li}X_{lj}$
- MDS is Eigendecomposition of  $-\frac{1}{2}\mathbf{C}\mathbf{F}_2\mathbf{C}$
- $\mathbf{C}\mathbf{F}_2\mathbf{C} = \mathbf{C}\mathbf{X}_i^2\mathbf{C} + \mathbf{C}\mathbf{X}_j^2\mathbf{C} - 2\mathbf{C}\mathbf{X}\mathbf{X}^T\mathbf{C}$

# PCA is MDS on $\mathbf{F}_2$

- PCA is decomposition of Covariance matrix:  $\mathbf{Y}\mathbf{Y}^T$
- Consider  $\mathbf{F}_2$ ;  $f_{ij} = F_2(X_i, X_j) = \sum_l (X_{li}^2 - X_{lj}^2)$
- Consider  $\mathbf{F}_2$ ;  $f_{ij} = \sum_l X_{li}^2 + \sum_l X_{lj}^2 - 2 \sum_l X_{li}X_{lj}$
- MDS is Eigendecomposition of  $-\frac{1}{2}\mathbf{C}\mathbf{F}_2\mathbf{C}$
- $\mathbf{C}\mathbf{F}_2\mathbf{C} = \underbrace{\mathbf{C}\mathbf{X}_i^2\mathbf{C}}_0 + \underbrace{\mathbf{C}\mathbf{X}_j^2\mathbf{C}}_0 - 2\underbrace{\mathbf{C}\mathbf{X}\mathbf{X}^T\mathbf{C}}_{\mathbf{Y}\mathbf{Y}^T}$

# PCA is MDS on $\mathbf{F}_2$

- PCA is decomposition of Covariance matrix:  $\mathbf{Y}\mathbf{Y}^T$
- Consider  $\mathbf{F}_2$ ;  $f_{ij} = F_2(X_i, X_j) = \sum_l (X_{li}^2 - X_{lj}^2)$
- Consider  $\mathbf{F}_2$ ;  $f_{ij} = \sum_l X_{li}^2 + \sum_l X_{lj}^2 - 2 \sum_l X_{li}X_{lj}$
- MDS is Eigendecomposition of  $-\frac{1}{2}\mathbf{C}\mathbf{F}_2\mathbf{C}$
- $\mathbf{C}\mathbf{F}_2\mathbf{C} = \underbrace{\mathbf{C}\mathbf{X}_i^2\mathbf{C}}_0 + \underbrace{\mathbf{C}\mathbf{X}_j^2\mathbf{C}}_0 - 2\underbrace{\mathbf{C}\mathbf{X}\mathbf{X}^T\mathbf{C}}_{\mathbf{Y}\mathbf{Y}^T}$

## Observation

PCA is equivalent to MDS on  $\mathbf{F}_2$

# PCA is MDS on OutgroupF<sub>3</sub>

- PCA is decomposition of Covariance matrix:  $\mathbf{Y}\mathbf{Y}^T$

# PCA is MDS on Outgroup F<sub>3</sub>

- PCA is decomposition of Covariance matrix:  $\mathbf{Y}\mathbf{Y}^T$
- $\mathbf{F}_3(O); f_{ij} = F_3(O; X_i, X_j)$
- $f_{ij} = \sum_l [o_l^2 - o_l X_{li} - o_l X_{lj} + X_{li}X_{lj}]$

# PCA is MDS on Outgroup $F_3$

- PCA is decomposition of Covariance matrix:  $\mathbf{Y}\mathbf{Y}^T$
- $\mathbf{F}_3(O); f_{ij} = F_3(O; X_i, X_j)$
- $f_{ij} = \sum_l [o_l^2 - o_l X_{li} - o_l X_{lj} + X_{li}X_{lj}]$
- $\mathbf{C}\mathbf{F}_3\mathbf{C} = \underbrace{\mathbf{CO}^2\mathbf{C}}_0 - \underbrace{\mathbf{COX}_i\mathbf{C}}_0 - \underbrace{\mathbf{COX}_j\mathbf{C}}_0 + \underbrace{\mathbf{CXX}^T\mathbf{C}}_{\mathbf{YY}^T}$

# PCA is MDS on Outgroup $F_3$

- PCA is decomposition of Covariance matrix:  $\mathbf{Y}\mathbf{Y}^T$
- $\mathbf{F}_3(O); f_{ij} = F_3(O; X_i, X_j)$
- $f_{ij} = \sum_l [o_l^2 - o_l X_{li} - o_l X_{lj} + X_{li}X_{lj}]$
- $\mathbf{C}\mathbf{F}_3\mathbf{C} = \underbrace{\mathbf{CO}^2\mathbf{C}}_0 - \underbrace{\mathbf{COX}_i\mathbf{C}}_0 - \underbrace{\mathbf{COX}_j\mathbf{C}}_0 + \underbrace{\mathbf{CXX}^T\mathbf{C}}_{\mathbf{YY}^T}$

# PCA is MDS on Outgroup $F_3$

- PCA is decomposition of Covariance matrix:  $\mathbf{Y}\mathbf{Y}^T$
- $\mathbf{F}_3(O); f_{ij} = F_3(O; X_i, X_j)$
- $f_{ij} = \sum_l [o_l^2 - o_l X_{li} - o_l X_{li} + X_{li}X_{lj}]$
- $\mathbf{C}\mathbf{F}_3\mathbf{C} = \underbrace{\mathbf{CO}^2\mathbf{C}}_0 - \underbrace{\mathbf{COX}_i\mathbf{C}}_0 - \underbrace{\mathbf{COX}_j\mathbf{C}}_0 + \underbrace{\mathbf{CXX}^T\mathbf{C}}_{\mathbf{YY}^T}$

## Observation

Decomposition of *any* centered  $F_3$ -matrix is equivalent to PCA.

# 0-diagonal MDS

In aDNA applications (e.g. Fu et al 2016), MDS is performed on

$$\mathbf{M} = \begin{pmatrix} 0 & 1 - f_3(O; X_1, X_2) & \dots & 1 - f_3(O; X_1, X_n) \\ 1 - f_3(O; X_1, X_2) & 0 & \dots & 1 - f_3(O; X_2, X_n) \\ \vdots & \vdots & \ddots & 1 - f_3(O; X_i, X_n) \\ 1 - f_3(O; X_1, X_n) & 1 - f_3(O; X_2, X_n) & \dots & 0 \end{pmatrix}$$

It can be shown that

$$\mathbf{CMC} = \mathbf{CF}_2 \mathbf{C} + \mathbf{COC}$$

with  $o_{ii} = f_2(O, X_i) - 1$ ;  $o_{ij} = 0$

## How can we interpret PCA in context of $F$ -stats?

- PCA is decomposition of  $\mathbf{F}_3(\boldsymbol{\mu}, X_1, X_2)$ -matrix
- PCA is a function of (unnormalized)- $\mathbf{F}_2$ -matrix
  - same information content
  - differences in analysis choices (individual/population, normalization, estimate) need to be justified

## How can we interpret PCA in context of $F$ -stats?

- PCA is decomposition of  $\mathbf{F}_3(\mu, X_1, X_2)$ -matrix
- PCA is a function of (unnormalized)- $\mathbf{F}_2$ -matrix
  - same information content
  - differences in analysis choices (individual/population, normalization, estimate) need to be justified
- setting diagonal allows some leeway

## How can we interpret PCA in context of $F$ -stats?

- PCA is decomposition of  $\mathbf{F}_3(\mu, X_1, X_2)$ -matrix
- PCA is a function of (unnormalized)- $\mathbf{F}_2$ -matrix
  - same information content
  - differences in analysis choices (individual/population, normalization, estimate) need to be justified
- setting diagonal allows some leeway
- *Conjecture:* PCA is analogous to midpoint-rooted phylogenetic analyses; Outgroup- $F_3$  allow other rootings.

## $F$ -statistics on PCA-plot

- Recall that PCA is just translation + rotation

## $F$ -statistics on PCA-plot

- Recall that PCA is just translation + rotation
- Distances (such as  $F_2$ ) are invariant to translation + rotation

## $F$ -statistics on PCA-plot

- Recall that PCA is just translation + rotation
- Distances (such as  $F_2$ ) are invariant to translation + rotation
- 

$$F_2(X_1, X_2) = \sum_{\text{loci}} (x_{1l} - x_{2l})^2$$

## $F$ -statistics on PCA-plot

- Recall that PCA is just translation + rotation
- Distances (such as  $F_2$ ) are invariant to translation + rotation
- 

$$F_2(X_1, X_2) = \sum_{\text{loci}} (x_{1l} - x_{2l})^2$$

- 

$$F_2(X_1, X_2) = \sum_{\text{PCs}} (x_{1p} - x_{2p})^2$$

# $F$ -statistics on PCA-plot

- Recall that PCA is just translation + rotation
- Distances (such as  $F_2$ ) are invariant to translation + rotation
- 

$$F_2(X_1, X_2) = \sum_{\text{loci}} (x_{1l} - x_{2l})^2$$

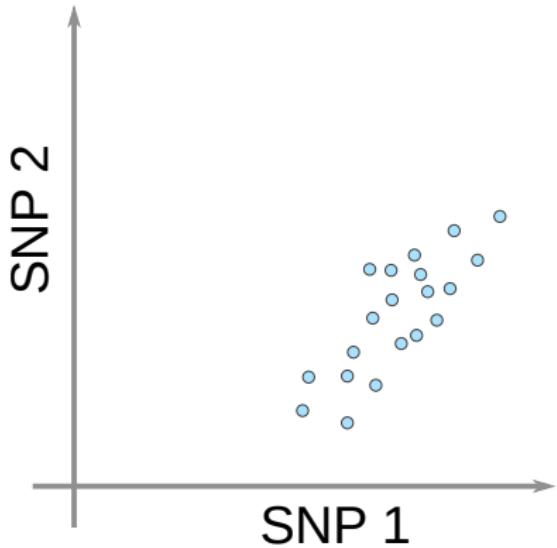
- 

$$F_2(X_1, X_2) = \sum_{\text{PCs}} (x_{1p} - x_{2p})^2$$

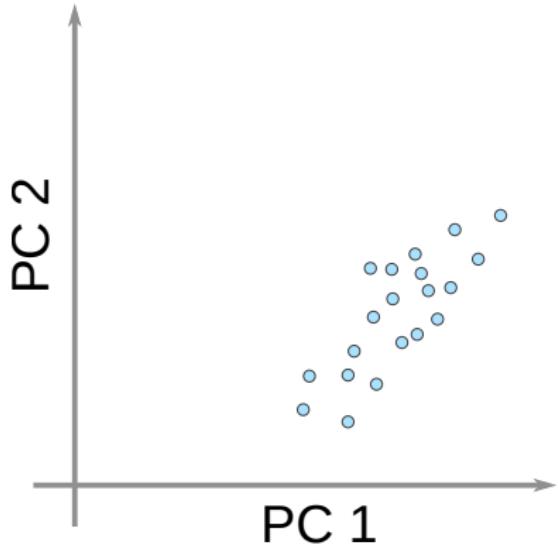
## Observation

$F_2$  can be decomposed in contributions of different principal components

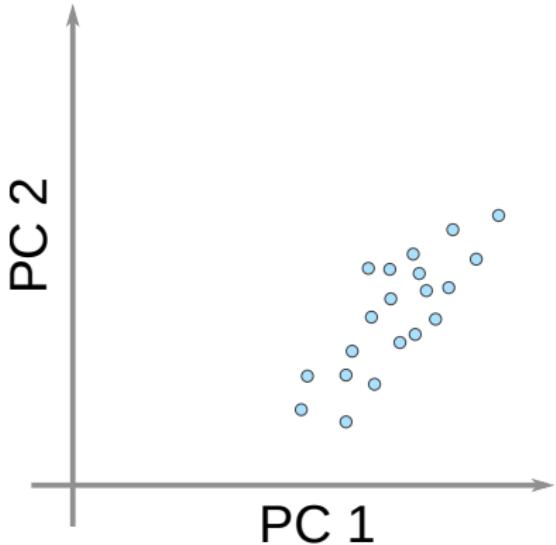
# $F$ -statistics on PCA-plot



# $F$ -statistics on PCA-plot

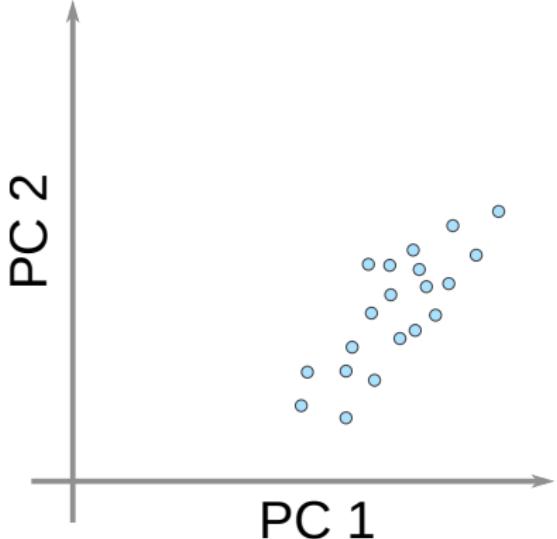


# $F$ -statistics on PCA-plot



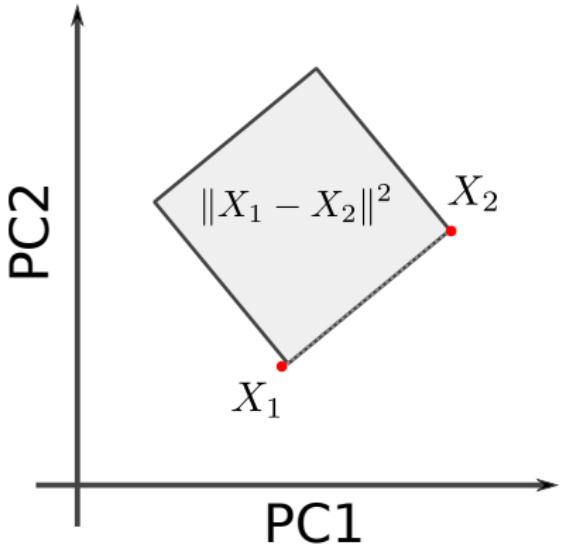
- $F$ -statistics have a geometrical representation on PCA-plot
- Exact only if we use *all* PCs

# $F$ -statistics on PCA-plot



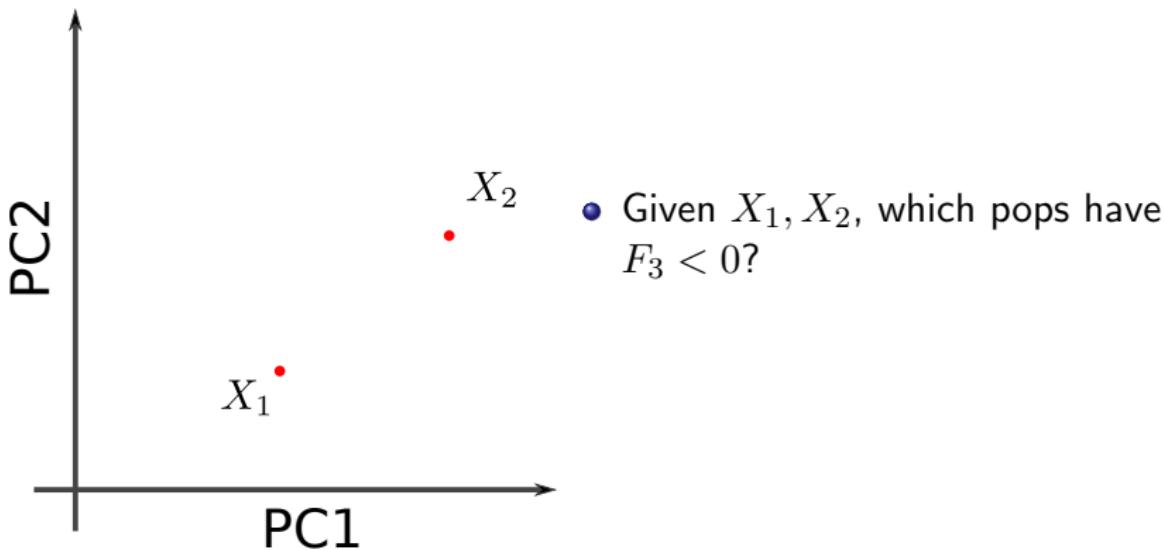
- $F$ -statistics have a geometrical representation on PCA-plot
- Exact only if we use *all* PCs
- Good approximation for 2D-plot if first 2 PCs capture relevant population structure

# $F_2$ -statistic on PCA-plot

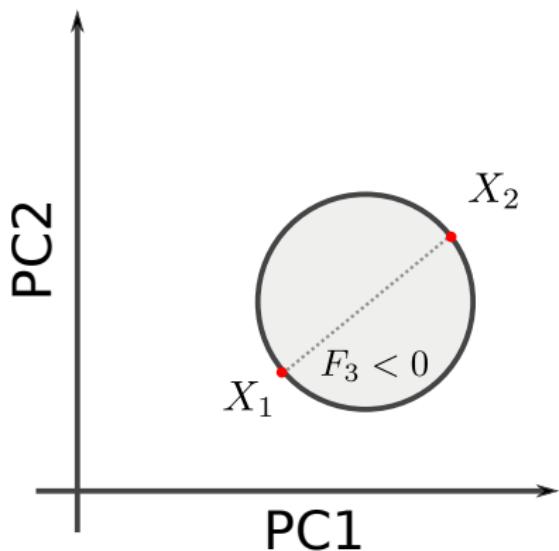


- $F_2(X_1, X_2) = \sum_l (X_{1l} - X_{2l})^2$
- $F_2(X_1, X_2) = \|X_1, X_2\|^2$

# Admixed populations ( $F_3$ ) on PCA-plot

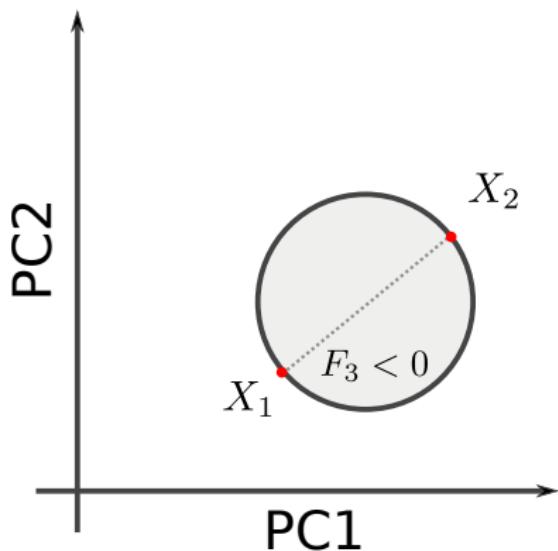


# Admixed populations ( $F_3$ ) on PCA-plot



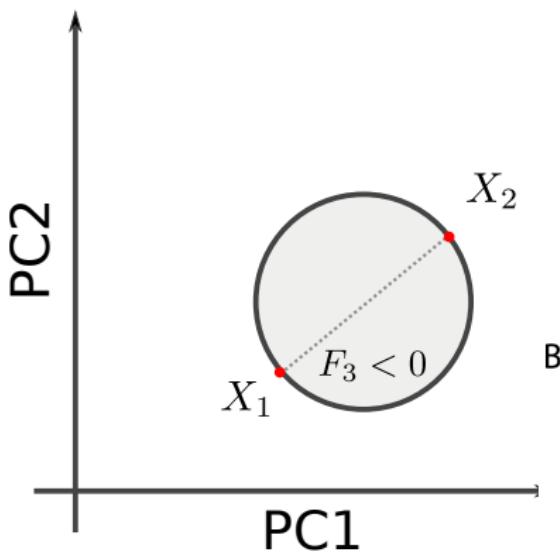
- Given  $X_1, X_2$ , which pops have  $F_3 < 0$ ?
- $F_3(Y; X_1, X_2) = 0$  is a circle!

# Admixed populations ( $F_3$ ) on PCA-plot

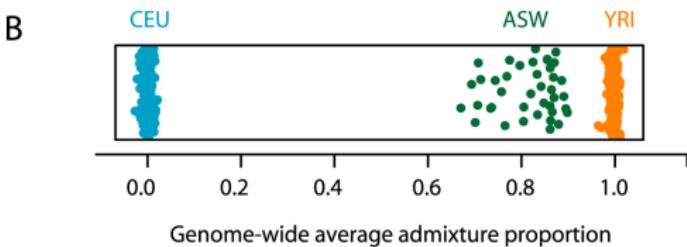


- Given  $X_1, X_2$ , which pops have  $F_3 < 0$ ?
- $F_3(Y; X_1, X_2) = 0$  is a circle!

# Admixed populations ( $F_3$ ) on PCA-plot

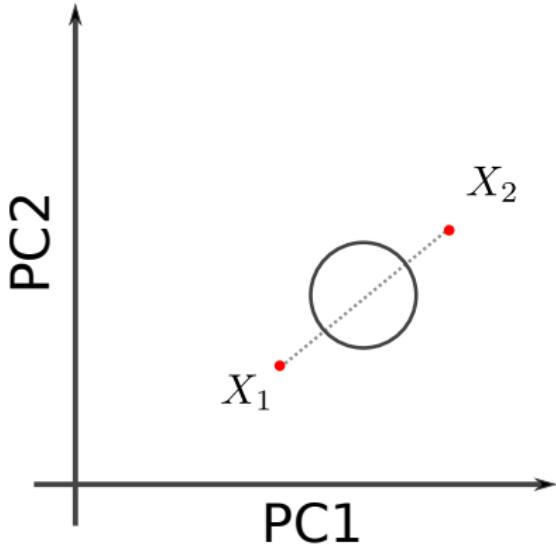


- Given  $X_1, X_2$ , which pops have  $F_3 < 0$ ?
- $F_3(Y; X_1, X_2) = 0$  is a circle!



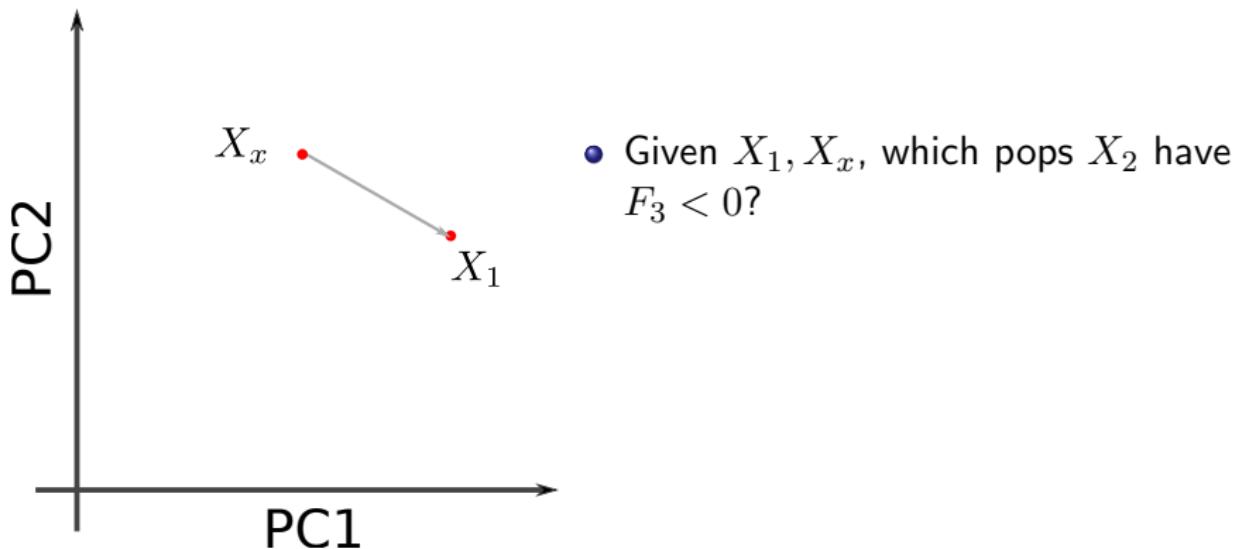
McVean 2009

# Admixed populations ( $F_3$ ) on PCA-plot

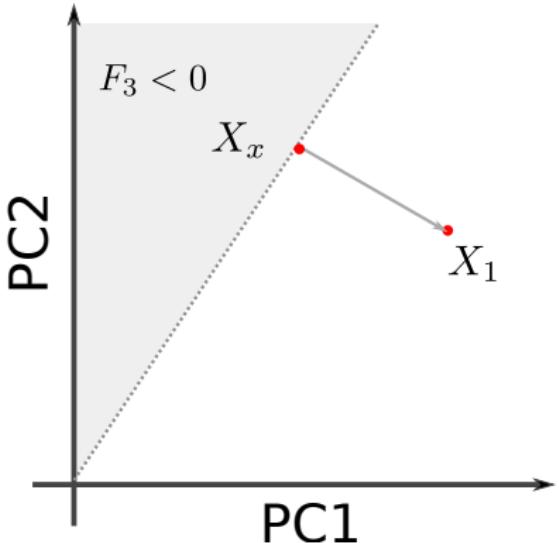


- Given  $X_1, X_2$ , which pops have  $F_3 < 0$ ?
- $F_3(Y; X_1, X_2) = 0$  is a circle!
- $F_3(Y; X_1, X_2) = k < 0$  is smaller circle

# Admixture $F_3$ -stats on PCA-plot

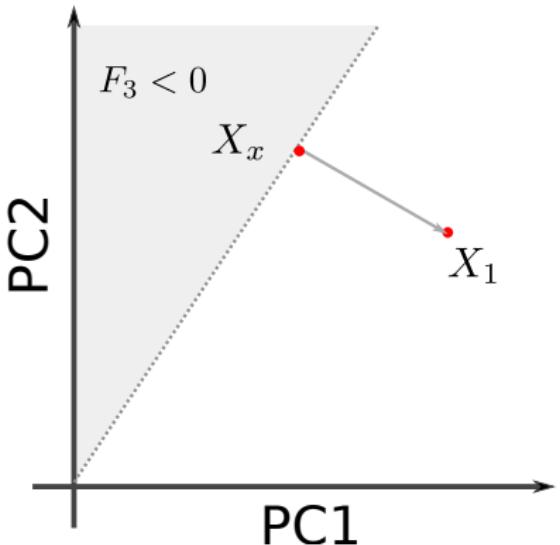


# Admixture $F_3$ -stats on PCA-plot



- Given  $X_1, X_x$ , which pops  $X_2$  have  $F_3 < 0$ ?
- $F_3$  is 0 if  $(X_x; X_1), (X_x; X_2)$  form a right angle!

# Admixture $F_3$ -stats on PCA-plot

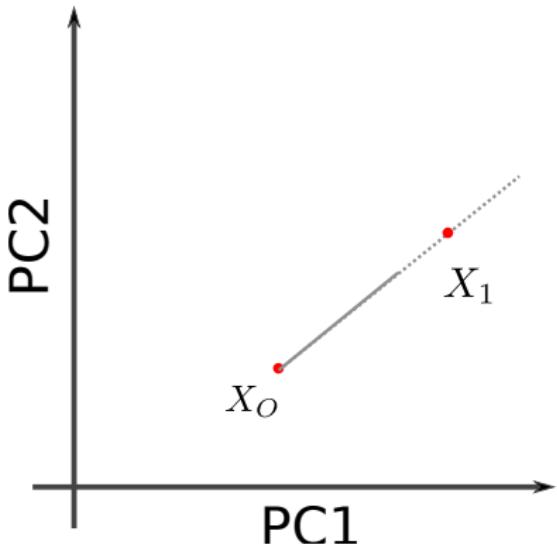


- Given  $X_1, X_x$ , which pops  $X_2$  have  $F_3 < 0$ ?
- $F_3$  is 0 if  $(X_x; X_1), (X_x; X_2)$  form a right angle!
- Inner (dot) product:  
$$F_3(X_x; X_1, X_2) = \langle X_x - X_1, X_x - X_2 \rangle$$

# Outgroup $F_3$ -stats on PCA-plot

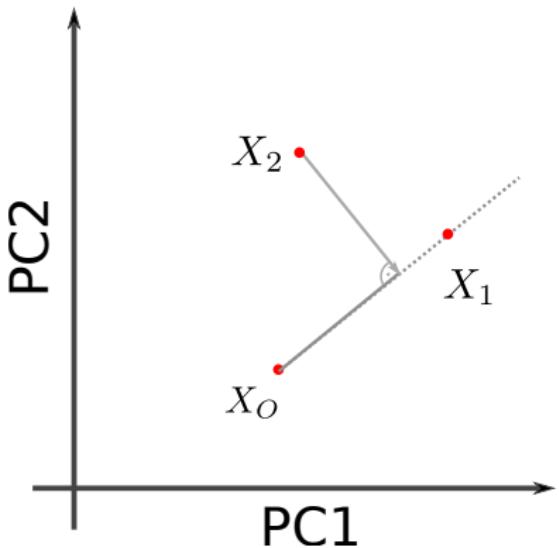
- Given  $X_O, X_1$ , which pop  $X_2$  has highest  $F_3$ ?

# Outgroup $F_3$ -stats on PCA-plot



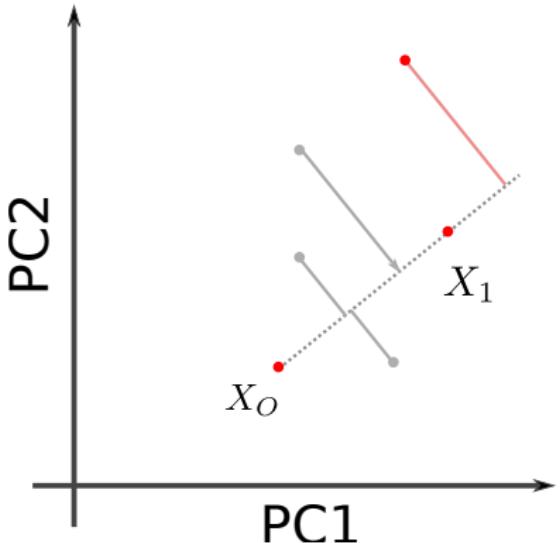
- Given  $X_O, X_1$ , which pop  $X_2$  has highest  $F_3$ ?
- As  $F_3 = \langle X_O - X_1, X_O - X_2 \rangle$ , project  $\overline{X_O X_2}$  on  $\overline{X_O X_1}$

# Outgroup $F_3$ -stats on PCA-plot



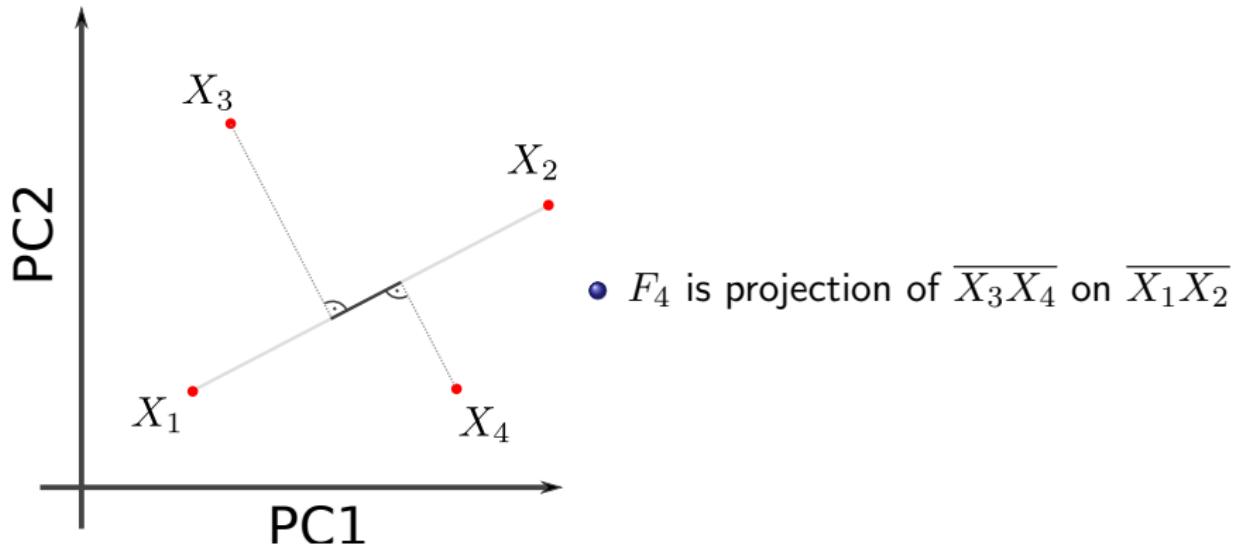
- Given  $X_O, X_1$ , which pop  $X_2$  has highest  $F_3$ ?
- As  $F_3 = \langle X_O - X_1, X_O - X_2 \rangle$ , project  $\overline{X_OX_2}$  on  $\overline{X_OX_1}$
- $F_3$  is proportional to segment in dir  $X_0 \rightarrow X_1$

# Outgroup $F_3$ -stats on PCA-plot

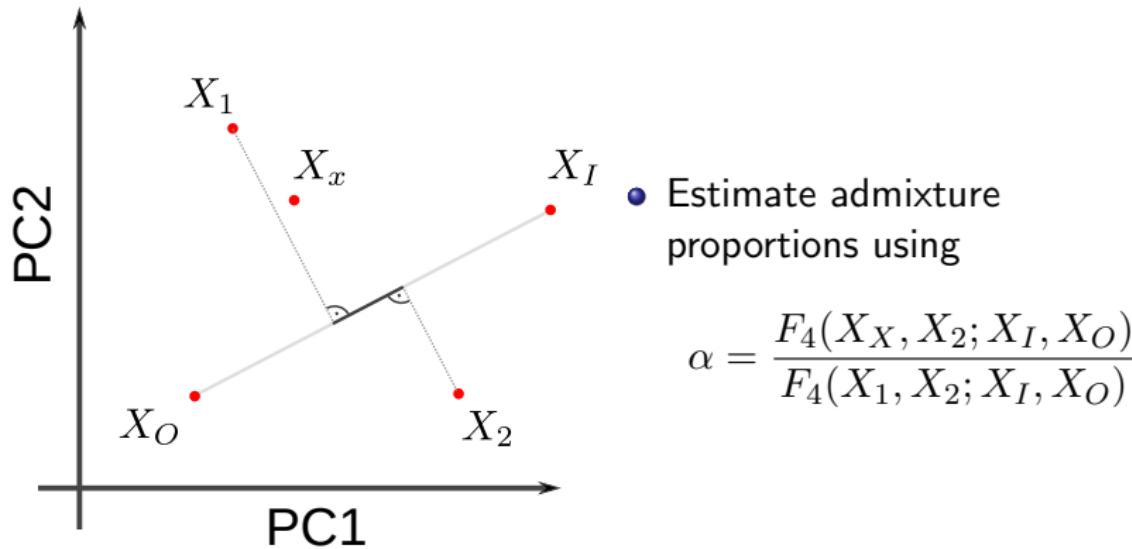


- Given  $X_O, X_1$ , which pop  $X_2$  has highest  $F_3$ ?
- As  $F_3 = \langle X_O - X_1, X_O - X_2 \rangle$ , project  $\overrightarrow{X_O X_2}$  on  $\overrightarrow{X_O X_1}$
- $F_3$  is proportional to segment in dir  $X_0 \rightarrow X_1$
- Pop furthest on axis has highest  $F_3$ .

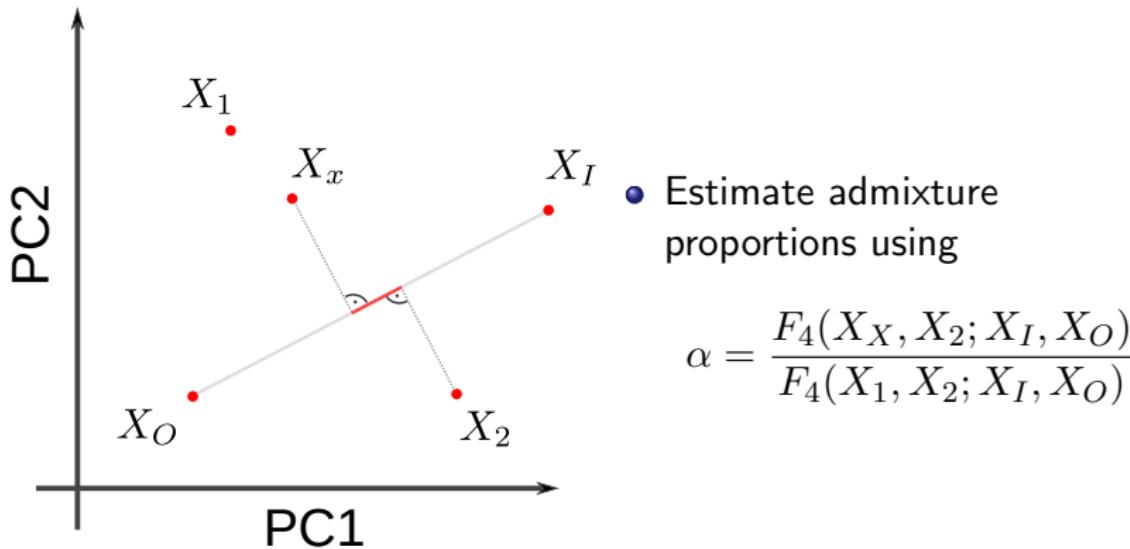
# $F_4$ -stats on PCA-plot



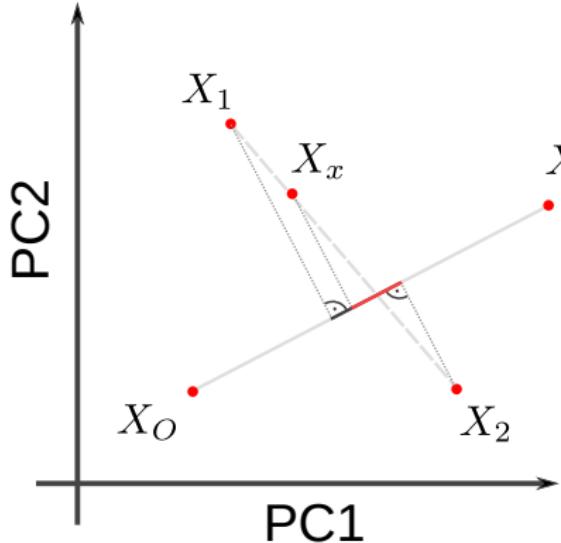
## $F_4$ -ratio on PCA plot



# $F_4$ -ratio on PCA plot



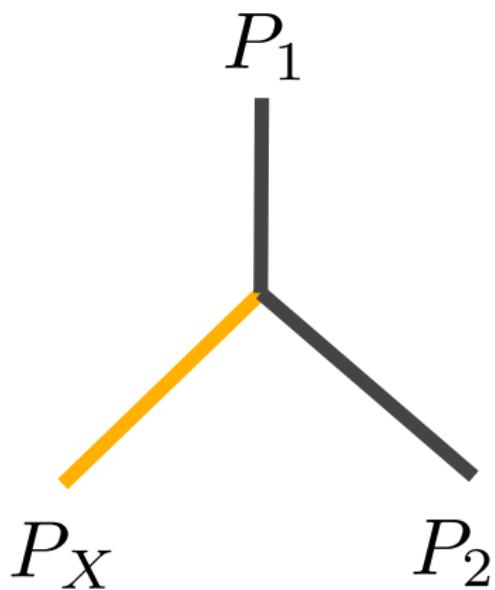
# $F_4$ -ratio on PCA plot



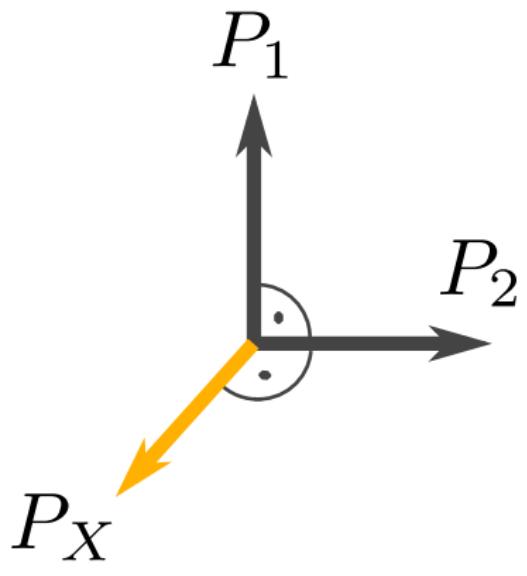
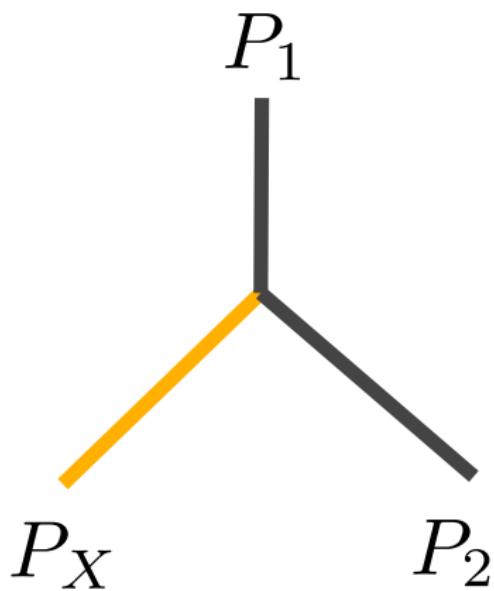
- Estimate admixture proportions using

$$\alpha = \frac{F_4(X_X, X_2; X_I, X_O)}{F_4(X_1, X_2; X_I, X_O)}$$

# Where does Orthogonality come from?



# Where does Orthogonality come from?



# How can we interpret $F$ -stats in the context of PCA?

- $F$ -stats have geometric interpretation

# How can we interpret $F$ -stats in the context of PCA?

- $F$ -stats have geometric interpretation
- exact in  $n$ -dimensional, approximate in 2D

# How can we interpret $F$ -stats in the context of PCA?

- $F$ -stats have geometric interpretation
- exact in  $n$ -dimensional, approximate in 2D
- $F_3$  is a (hyper)-circle

# How can we interpret $F$ -stats in the context of PCA?

- $F$ -stats have geometric interpretation
- exact in  $n$ -dimensional, approximate in 2D
- $F_3$  is a (hyper)-circle
- $F_4$  established orthogonality

## How can we interpret $F$ -stats in the context of PCA?

- $F$ -stats have geometric interpretation
- exact in  $n$ -dimensional, approximate in 2D
- $F_3$  is a (hyper)-circle
- $F_4$  established orthogonality
- Orthogonality stems from independence of drift in distinct parts of history

# Applications

- ① Better link  $F$ -stats and PCA results
  - use Dimensions / Orthogonality for useful data representations

# Applications

- ① Better link  $F$ -stats and PCA results
  - use Dimensions / Orthogonality for useful data representations
- ② Diagnostic plots
  - Easy checks whether assumptions for  $f_4$ -ratio, `qpadm`, etc. are satisfied

# Applications

- ① Better link  $F$ -stats and PCA results
  - use Dimensions / Orthogonality for useful data representations
- ② Diagnostic plots
  - Easy checks whether assumptions for  $f_4$ -ratio, qpAdm, etc. are satisfied
- ③ Distinguish admixture events

# Applications

- ① Better link  $F$ -stats and PCA results
  - use Dimensions / Orthogonality for useful data representations
- ② Diagnostic plots
  - Easy checks whether assumptions for  $f_4$ -ratio, qpAdm, etc. are satisfied
- ③ Distinguish admixture events
  - same  $F_3$  value may arise from distinct admixture events, PCs may point to differences

# Applications

- ① Better link  $F$ -stats and PCA results
  - use Dimensions / Orthogonality for useful data representations
- ② Diagnostic plots
  - Easy checks whether assumptions for  $f_4$ -ratio, qpAdm, etc. are satisfied
- ③ Distinguish admixture events
  - same  $F_3$  value may arise from distinct admixture events, PCs may point to differences
- ④ Understand discrepancies

# Applications

- ① Better link  $F$ -stats and PCA results
  - use Dimensions / Orthogonality for useful data representations
- ② Diagnostic plots
  - Easy checks whether assumptions for  $f_4$ -ratio, qpAdm, etc. are satisfied
- ③ Distinguish admixture events
  - same  $F_3$  value may arise from distinct admixture events, PCs may point to differences
- ④ Understand discrepancies
  - most likely due to data artifacts / higher PCs

# Applications

- ➊ Better link  $F$ -stats and PCA results
  - use Dimensions / Orthogonality for useful data representations
- ➋ Diagnostic plots
  - Easy checks whether assumptions for  $f_4$ -ratio, qpAdm, etc. are satisfied
- ➌ Distinguish admixture events
  - same  $F_3$  value may arise from distinct admixture events, PCs may point to differences
- ➍ Understand discrepancies
  - most likely due to data artifacts / higher PCs
- ➎ Standardize normalization
  - $F_2^{(\text{PCA})} = \frac{1}{\sigma} \sum (X_i - \bar{X})^2$
  - $F_2^{(F\text{-stats})} = \sum (X_i - \bar{X})^2$

# Applications

- ① Better link  $F$ -stats and PCA results
  - use Dimensions / Orthogonality for useful data representations
- ② Diagnostic plots
  - Easy checks whether assumptions for  $f_4$ -ratio, qpAdm, etc. are satisfied
- ③ Distinguish admixture events
  - same  $F_3$  value may arise from distinct admixture events, PCs may point to differences
- ④ Understand discrepancies
  - most likely due to data artifacts / higher PCs
- ⑤ Standardize normalization
  - $F_2^{(\text{PCA})} = \frac{1}{\sigma} \sum (X_i - \bar{X})^2$
  - $F_2^{(F\text{-stats})} = \sum (X_i - \bar{X})^2$
- ⑥ Better out-of-sample predictions
  - qpGraph and other tools fail with large samples

# Thank You!

