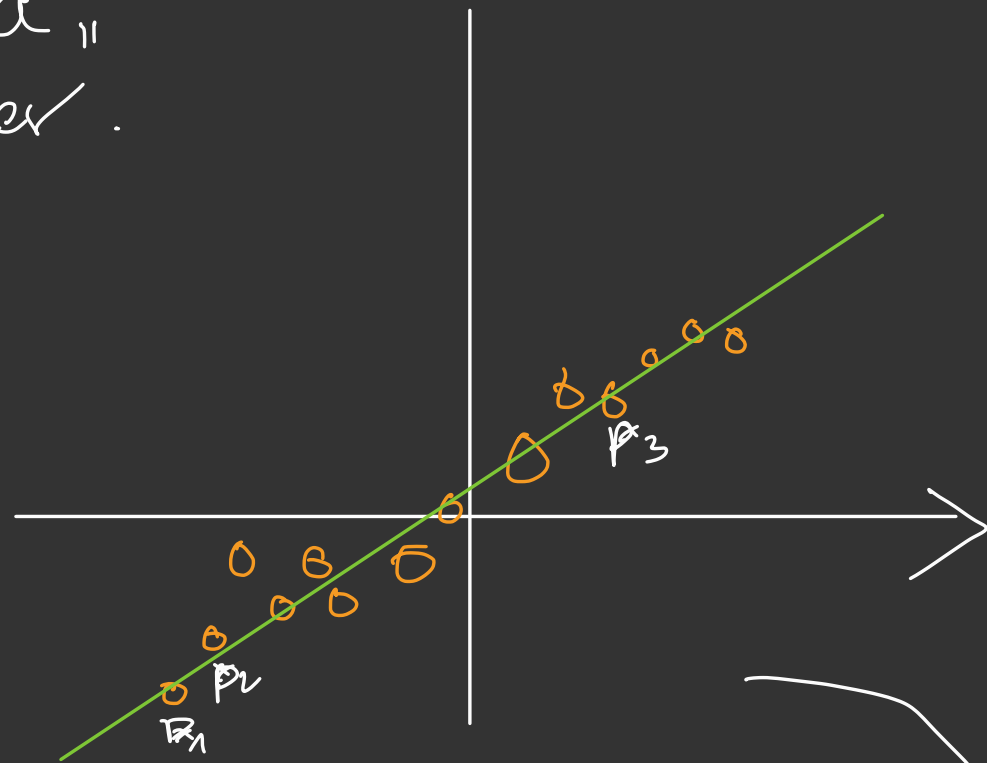


Unsupervised Learning: "Understand data" without teacher.

Principal Component Analysis (PCA)

Example: 2D dataset $S = (x_i, y_i)_{i=1}^n$

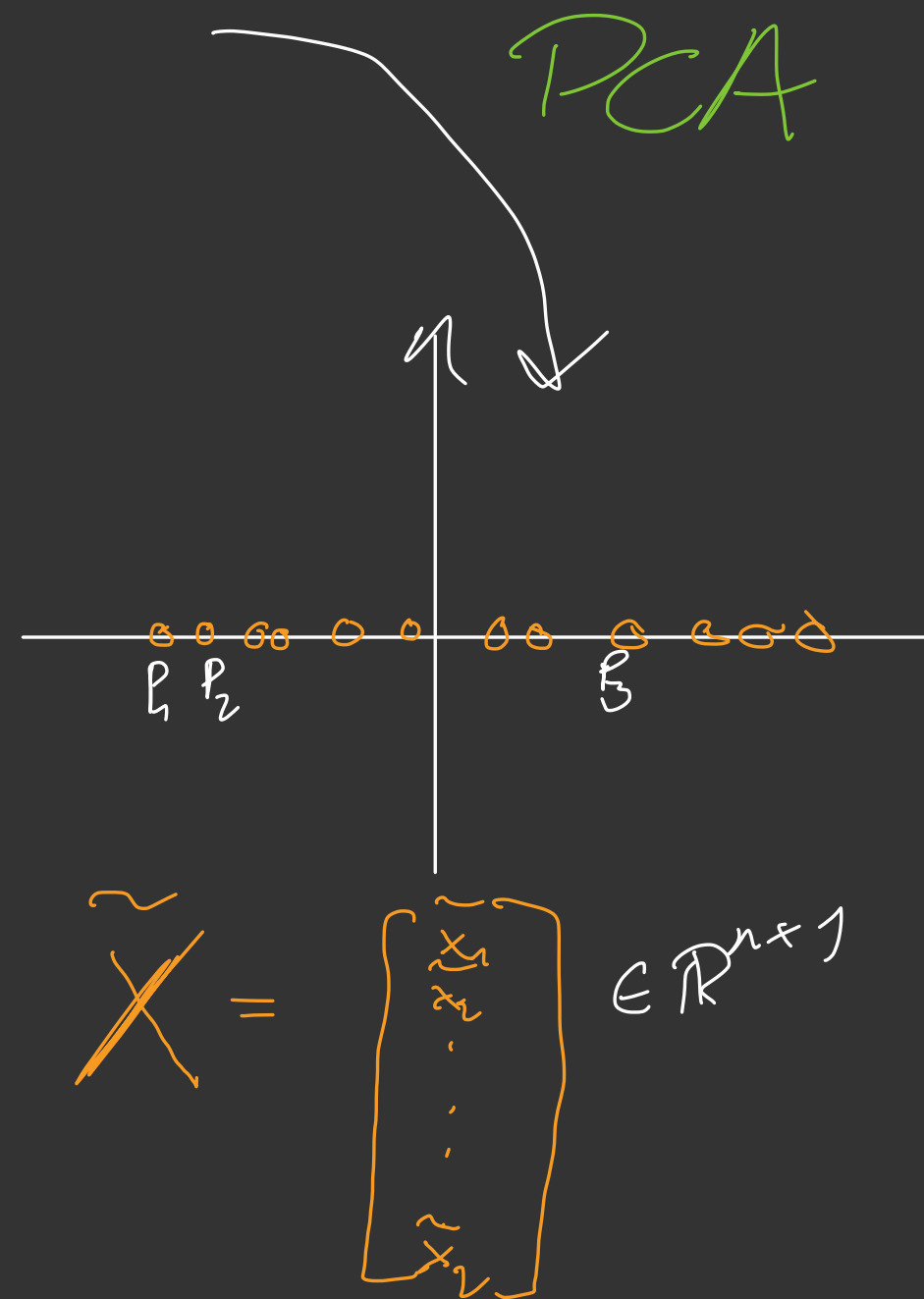
$$X = \begin{bmatrix} x_1 & y_1 \\ x_2 & y_2 \\ \vdots & \vdots \\ x_n & y_n \end{bmatrix} \in \mathbb{R}^{n \times 2}$$



Looks like "almost" a line, can we (as an approximation) represent S in 1D?

Idea: Dimension reduction for

- ▷ interpretability
- ▷ downstream computational savings



Human genetics

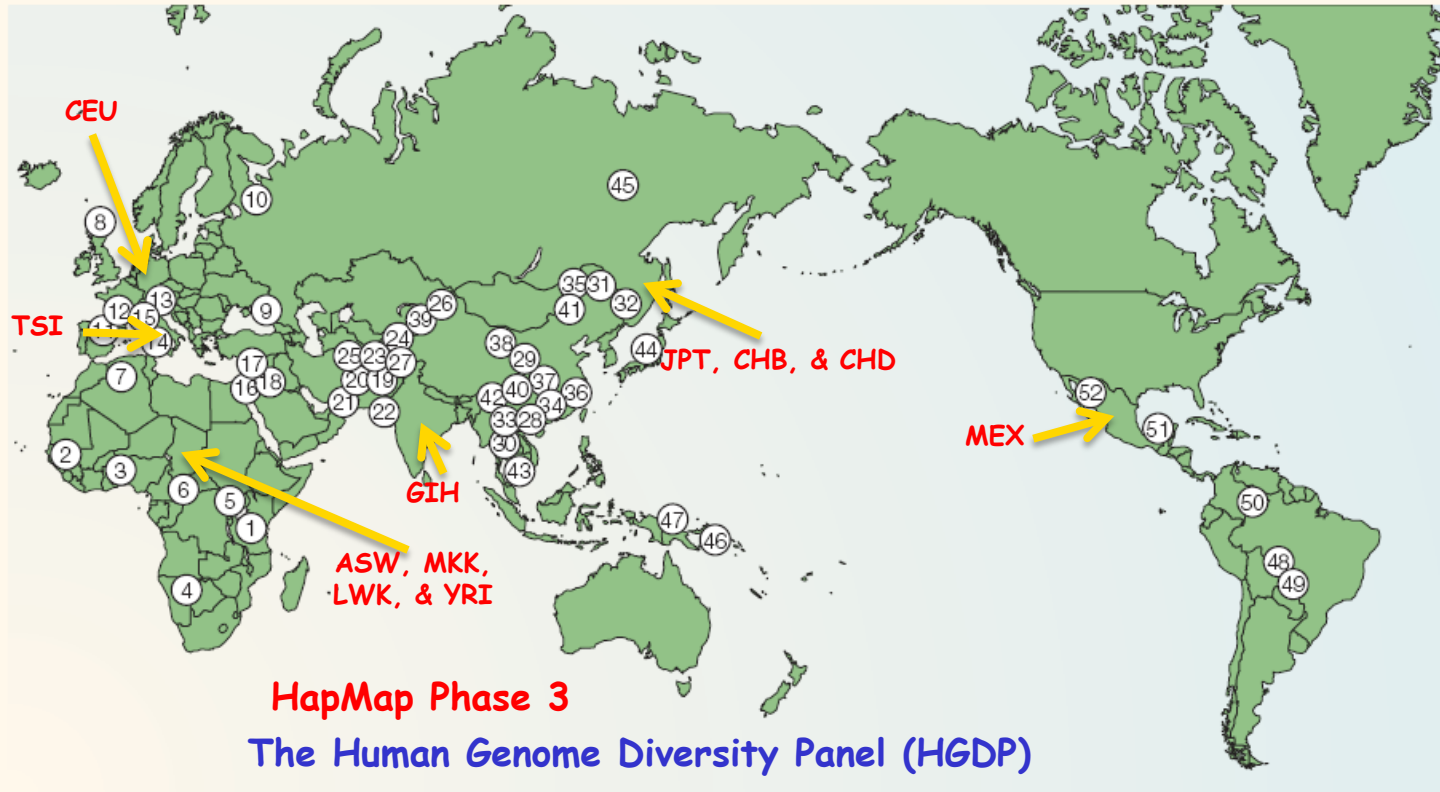
Single Nucleotide Polymorphisms: the most common type of genetic variation in the genome across different individuals.

They are **known** locations at the human genome where **two** alternate nucleotide bases (**alleles**) are observed (out of A, C, G, T).

SNPs



Matrices including thousands of individuals and hundreds of thousands if SNPs are available.



HapMap Phase 3
The Human Genome Diversity Panel (HGDP)

HGDP data

- 1,033 samples
- 7 geographic regions
- 52 populations

HapMap Phase 3 data

- 1,207 samples
- 11 populations

We will apply SVD/PCA on the (joint) HGDP and HapMap Phase 3 data.

Africans	Europeans	Western Asians	Eastern Asians	Oceanians
1 Bantu	8 Orcadian	16 Bedouin	28 Han (S. China)	46 Melanesian
2 Mandenka	9 Adygei	17 Druze	29 Han (N. China)	47 Papuan
3 Yoruba	10 Russian	18 Palestinian	30 Dai	
4 San	11 Basque		31 Daur	
5 Mbuti pygmy	12 French		32 Hezhen	
6 Biaka	13 North Italian		33 Lahu	
7 Mozabite	14 Sardinian		34 Miao	
	15 Tuscan		35 Oroqen	
		Central and Southern Asians	36 She	
		19 Balochi	37 Tujia	
		20 Brahui	38 Tu	
		21 Makrani	39 Xibo	
		22 Sindhi	40 Yi	
		23 Pathan	41 Mongola	
		24 Burusho	42 Naxi	
		25 Hazara	43 Cambodian	
		26 Uygur	44 Japanese	
		27 Kalash	45 Yakut	
				Native Americans
				48 Karitiana
				49 Surui
				50 Colombian
				51 Maya
				52 Pima

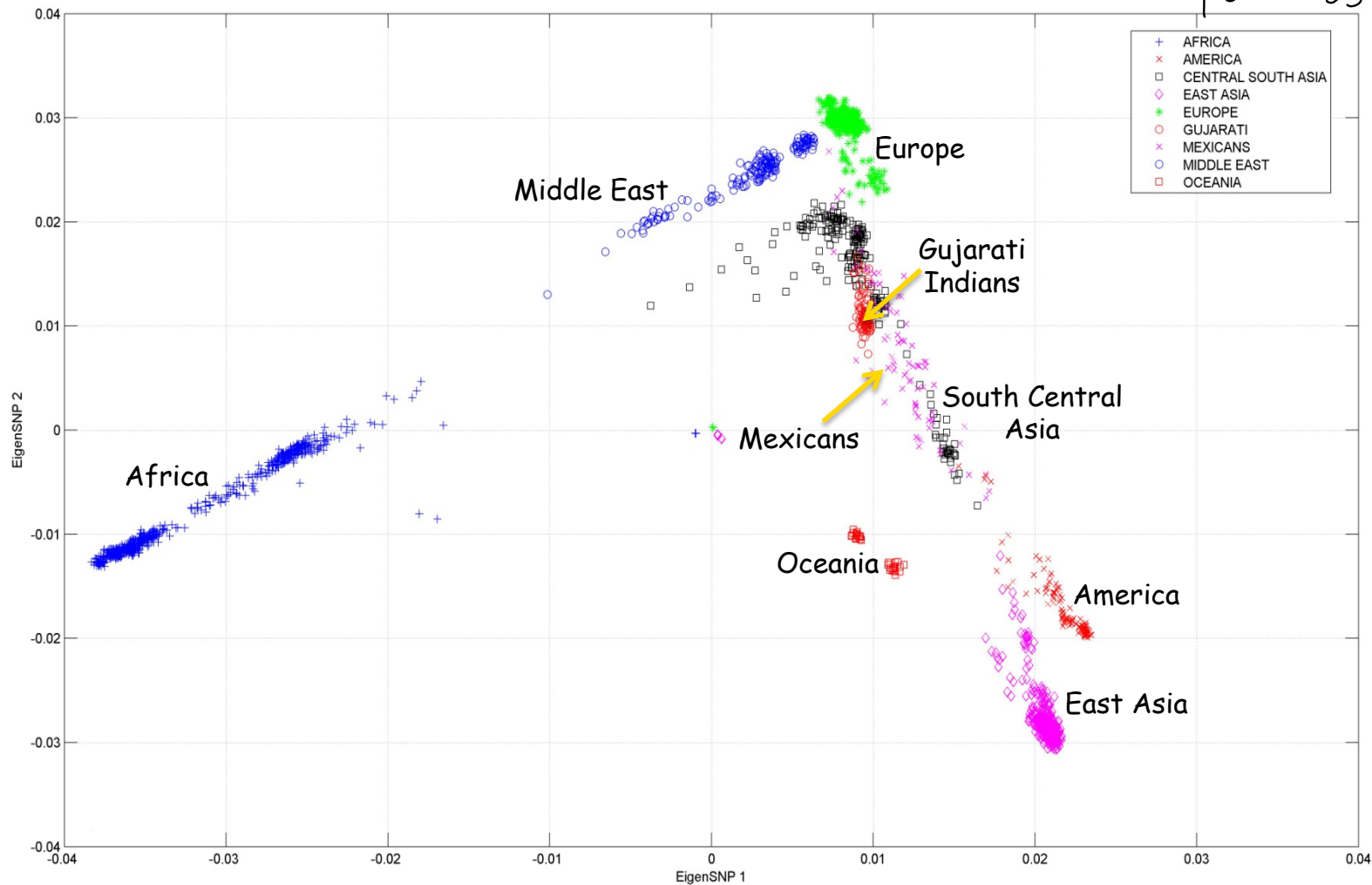
Cavalli-Sforza (2005) *Nat Genet Rev*
 Rosenberg et al. (2002) *Science*
 Li et al. (2008) *Science*
 The International HapMap Consortium (2003, 2005, 2007) *Nature*

Matrix dimensions:

2,240 subjects (rows)
 447,143 SNPs (columns)

Dense matrix:

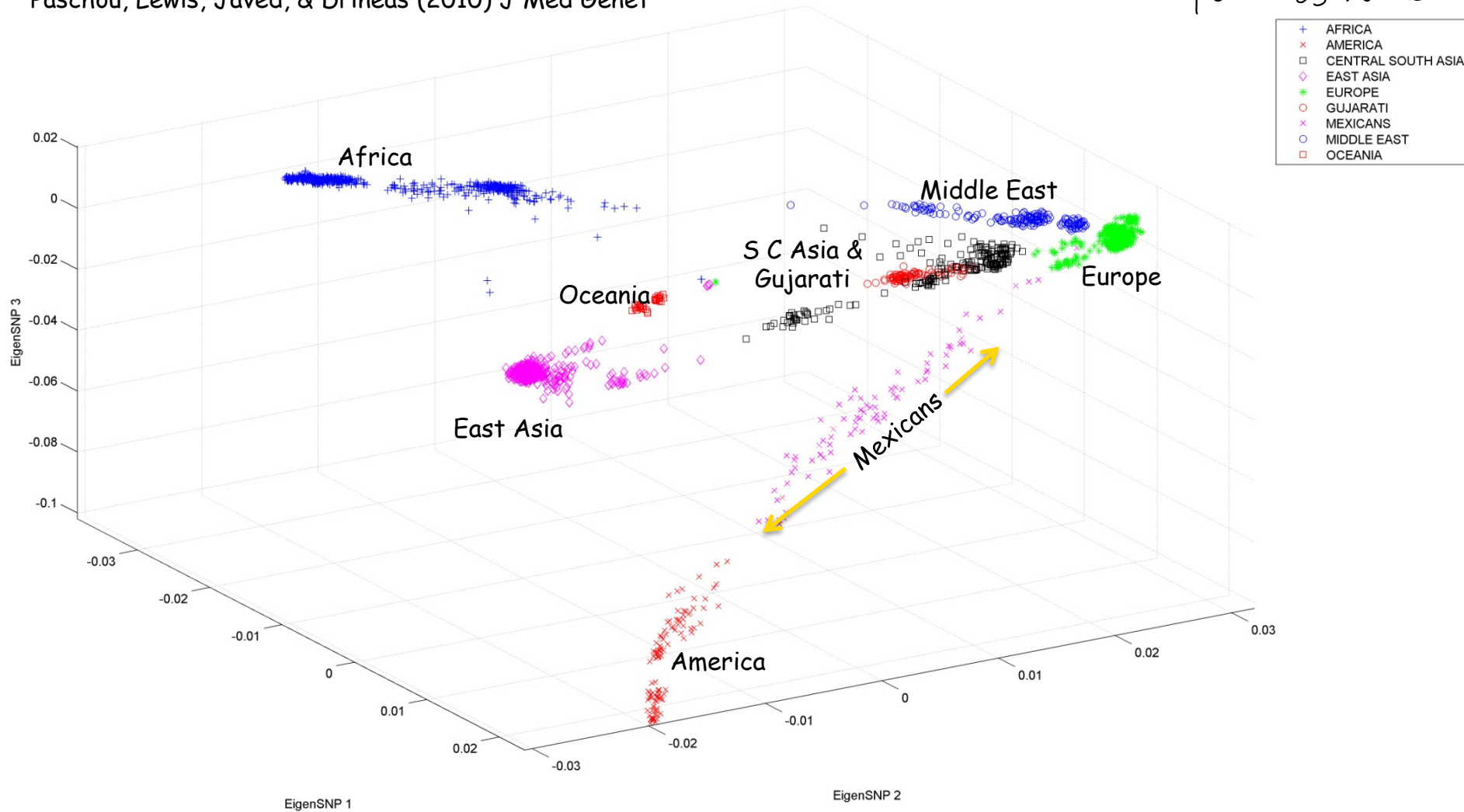
over one billion entries



- Top two Principal Components (PCs or eigenSNPs)

(Lin and Altman (2005) *Am J Hum Genet*)

- The figure renders visual support to the "out-of-Africa" hypothesis.
- Mexican population seems out of place: we move to the top three PCs.



Not altogether satisfactory: the principal components are linear combinations of all SNPs, and - of course - can not be assayed!

Can we find **actual SNPs** that capture the information in the singular vectors?

Formally: **spanning the same subspace.**

Setting of PCA: Let $X = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \in \mathbb{R}^{n \times k}$ Ex: n : nr. of subjects
 k : nr. of SKPs.

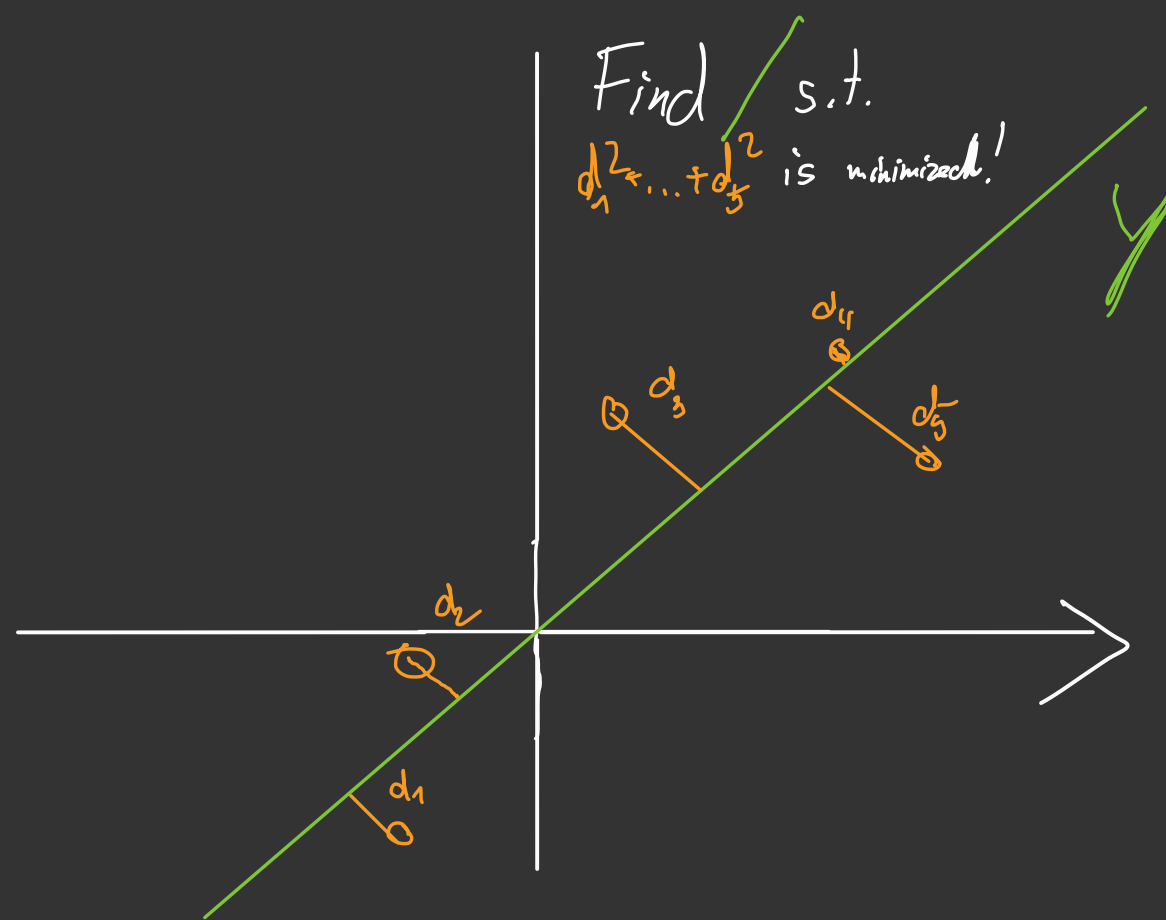
be a matrix of n data points with k features.

Assume without loss of generality that X is centered, i.e. $\frac{1}{n}(1, 1, \dots, 1)^T X = (0, \dots, 0)$
 (average value of each feature = 0).

Goal: For $r \ll k$, find r -dimensional subspace $Y \subset \mathbb{R}^k$, $\dim(Y) = r$,
 with orthonormal basis $V = \begin{bmatrix} | & & | \\ v_1 & \dots & v_r \\ | & & | \end{bmatrix} \in \mathbb{R}^{k \times r}$ such that

$$V = \underset{\substack{\tilde{V} \in \mathbb{R}^{k \times r} \\ \tilde{V}^T \tilde{V} = I_r}}{\operatorname{argmin}} \frac{1}{n} \|X - X \tilde{V} \tilde{V}^T\|_F^2 \quad (*)$$

"minimizes sum of squares of distances between points and projected points"



Properties of PCA: ▷ Obtain (approximate) low-dimensional representations

$$\tilde{Z} := X V = \begin{bmatrix} \tilde{z}_1 \\ \vdots \\ \tilde{z}_n \end{bmatrix} \in \mathbb{R}^{n \times \tau}$$

of the data points $\{x_1, \dots, x_n\}$.

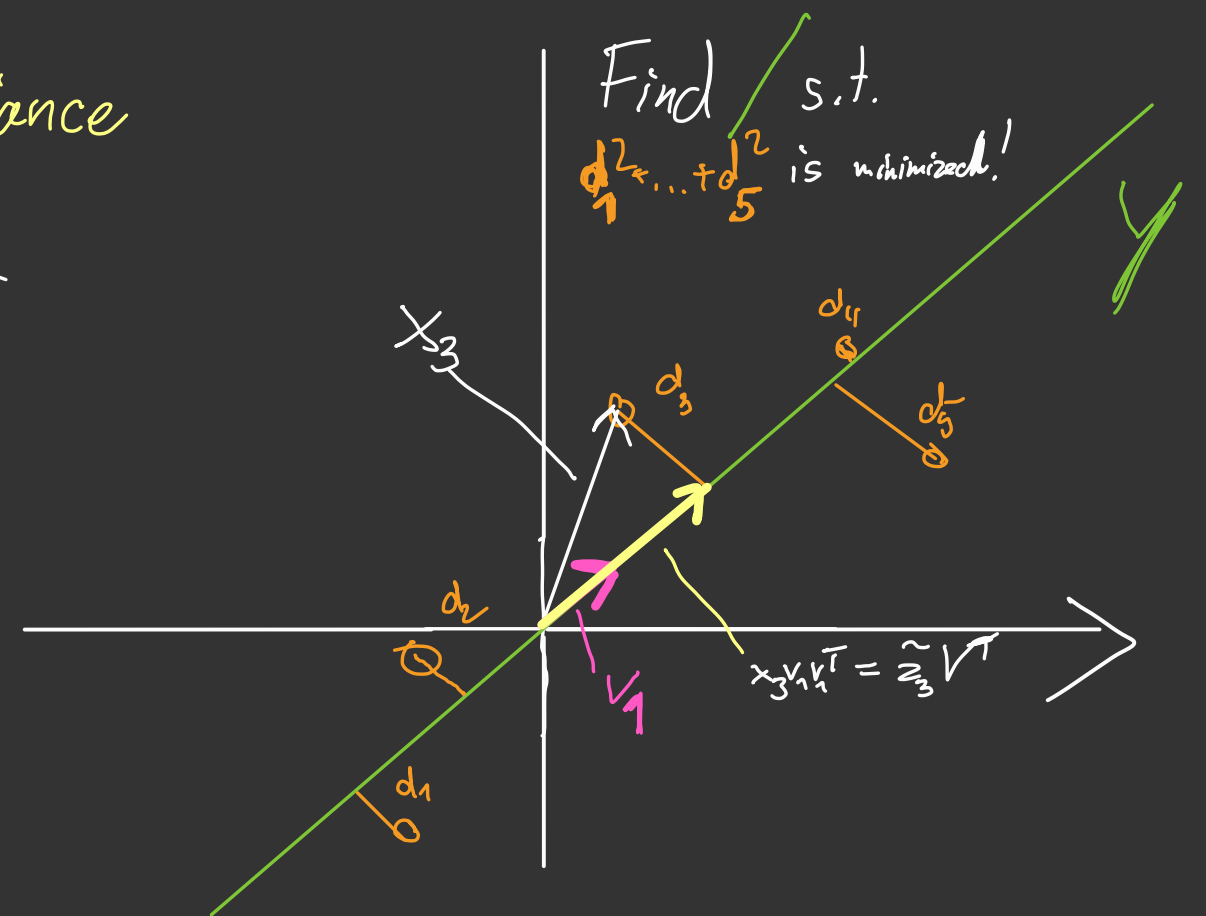
$$E_x: \begin{matrix} n & k \\ 2240 \times 447143 \\ \hline n & \tau \\ 2240 \times 3 \end{matrix}$$

▷ Memory savings: $\mathcal{O}(n\tau)$ instead of $\mathcal{O}(nk)$ parameters

▷ $V = \begin{bmatrix} | & & | \\ v_1 & \dots & v_\tau \\ | & & | \end{bmatrix}$: Matrix w/ eigenvectors corresponding

to τ largest eigenvalues of empirical covariance matrix

$$M = \frac{1}{n-1} \sum_{i=1}^n x_i x_i^T = \frac{1}{n-1} X^T X \in \mathbb{R}^{k \times k}$$



▷ PCA find ^(orthogonal) directions of maximal variance: (Interpretation favored in statistics)

Recall projection perspective of PCA:

$$(*) \quad V = \underset{\tilde{V} \in \mathbb{R}^{n \times k}: \tilde{V} \tilde{V}^T = I_k}{\operatorname{argmin}} \frac{1}{n} \|X - X \tilde{V} \tilde{V}^T\|_F^2$$

$$\cdot \frac{1}{n} \|X - X V V^T\|_F^2 = \frac{1}{n} \sum_{i=1}^n \|x_i^T - x_i^T V V^T\|_2^2, \text{ and}$$

$$\|x_i^T - x_i^T V V^T\|_2^2 = x_i^T x_i - 2 x_i^T V V^T x_i + \underbrace{x_i^T V V^T V V^T x_i}_{= x_i^T V V^T x_i} = x_i^T x_i - x_i^T V V^T x_i$$

Since $x_i^T x_i$ does not depend on V , rewrite (*) s.t.

$$V = \underset{\tilde{V}: \tilde{V} \tilde{V}^T = I_k}{\operatorname{argmax}} \underbrace{\frac{1}{(n-1)} \sum_{i=1}^n x_i^T \tilde{V} \tilde{V}^T x_i}_{\text{cyclicality of trace}} = \frac{1}{(n-1)} \sum_{i=1}^n \operatorname{tr}(x_i^T \tilde{V} \tilde{V}^T x_i)$$

$$\frac{1}{(n-1)} \sum_{i=1}^n \operatorname{tr}(\tilde{V}^T x_i x_i^T \tilde{V}) = \operatorname{tr}\left\{\tilde{V}^T \left(\frac{1}{n-1} \sum_{i=1}^n x_i x_i^T\right) \tilde{V}\right\}$$

⇒ \square connection to first k eigenvectors of M . = $\operatorname{tr}(\tilde{V}^T M \tilde{V})$

• Variance of dataset projected into first PC $\left(\left\{ x_i^T v_1 \right\}_{i=1}^n \right)$ if $r=1$

$$\text{Var}_x (x^T v_1) = \mathbb{E}_x [(x^T v_1)^2] = \sum_{i=1}^n \frac{1}{n} \|x_i^T v_1\|^2 = v_1^T \frac{1}{n} \sum_{i=1}^n x_i x_i^T v_1 = v_1^T M v_1$$

\Rightarrow PCA: Find Y that

- minimizes ^{sum of} squared distances $\sum d_i^2$
- maximizes variance of $\{x_i\}$ in Y .

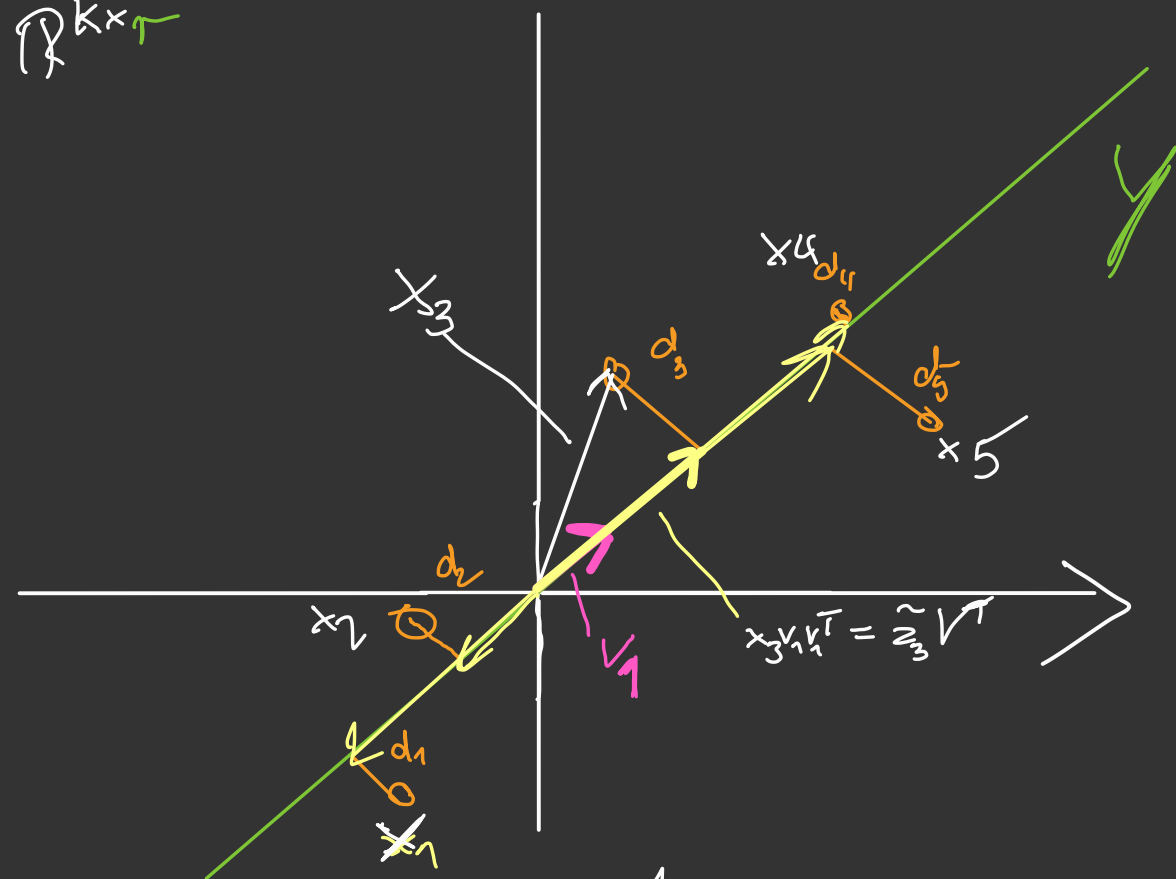
Terminology: \triangleright Columns v_i of $V = \begin{bmatrix} | & & | \\ v_1 & \dots & v_r \\ | & & | \end{bmatrix} \in \mathbb{R}^{k \times r}$

"Principal directions / axes"

\triangleright Columns Xv_i of $XV \in \mathbb{R}^{n \times r}$:
"Principal components" / "Scores"

$\triangleright \lambda_i$: "Variance explained by i -th PC"

$\triangleright \sqrt{\lambda_i} v_i$: " i -th loading"



Recall: $M := \frac{1}{n-1} X^T X$

$= \underset{\substack{\uparrow \\ \text{eigendecomposition}}}{[V \ V_{\perp}]} \begin{bmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_r \\ & & & \mathbb{I} \end{bmatrix} \begin{bmatrix} V^T \\ V_{\perp}^T \end{bmatrix}$

Terminology: \triangleright Columns v_i of $V = \begin{bmatrix} | & & | \\ v_1 & \dots & v_r \\ | & & | \end{bmatrix} \in \mathbb{R}^{k \times r}$

"Principal directions / axes"

$\triangleright \lambda_i$: "Variance explained by i -th PC"

\triangleright Columns Xv_i of $XV \in \mathbb{R}^{n \times r}$:

$\triangleright \sqrt{\lambda_i} v_i$: " i -th loading"

"Principal components" / "Scores": $Xv_i = \sqrt{\lambda_i} u_i$

• Instead of eigendecomposition of $M = \frac{1}{n-1} X^T X$, we can compute partial singular value decomposition of X s.t. $X \approx U \Sigma V^T$
 In this case, $\lambda_i = \frac{1}{n-1} \Sigma_{ii}^2$. $\in \mathbb{R}^{r \times k}$

\triangle : \triangleright All this holds for centered data $X = \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} \in \mathbb{R}^{n \times k}$.

\triangleright If $\mu := \frac{1}{n} \mathbf{1}^T X \neq 0$, center data first:

$X \xrightarrow{\sim} X - \frac{1}{n} \mathbf{1} \mathbf{1}^T X =: \tilde{X}$,
 proceed with \tilde{X} .

Clustering

Task: Find clusters/subgroups in a dataset $S = \{x_1, \dots, x_n\} \in \mathbb{R}^p$.

- ▷ Samples within subgroup similar/homogeneous
- ▷ Samples in different subgroups "distant"/heterogeneous from each other.

Ex:

- Customers of a company \rightarrow grouping for targeted marketing.
- Biology: Find groups of genes.

Q:

- What notion of similarity?
- Precise definition?

Note: Fundamentally different from classification problems!

"Unsupervised Learning":
No subset of data with "correct" classification/grouping available.

K-means clustering: [Steinhilber '56, Lloyd '57]

Given n points $x_1, \dots, x_n \in \mathbb{R}^D$, find K centroids $c_1, \dots, c_K \in \mathbb{R}^D$ and a partition $\Gamma_1 \cup \Gamma_2 \cup \dots \cup \Gamma_K = \{1, \dots, n\}$ ($\Gamma_i \cap \Gamma_j = \emptyset$ $\forall i \neq j$) such that

$$F(\{c_i\}_{i=1}^K, \{\Gamma_i\}_{i=1}^K) = \sum_{j=1}^K \sum_{i \in \Gamma_j} d(x_i, c_j) \quad \text{"k-means objective"}$$

is minimized, where $d: \mathbb{R}^D \times \mathbb{R}^D \rightarrow \mathbb{R}$ is a "distance" function

Ex: 1. $d(x, y) = \|x - y\|_2^2$ "squared Euclidean" ← by default
2. $d(x, y) = \|x - y\|_1$ "k-median"

Observation: K-means with is NP-hard for $k \geq 2$. [Drineas et al. '04].

Lloyd's algorithm (often called "k-means"):

Input: $\{x_i\}_{i=1}^n \in \mathbb{R}^p$, desired nr. of clusters k

1. Initialize $c_1, \dots, c_k \in \mathbb{R}^p$

Repeat until convergence:

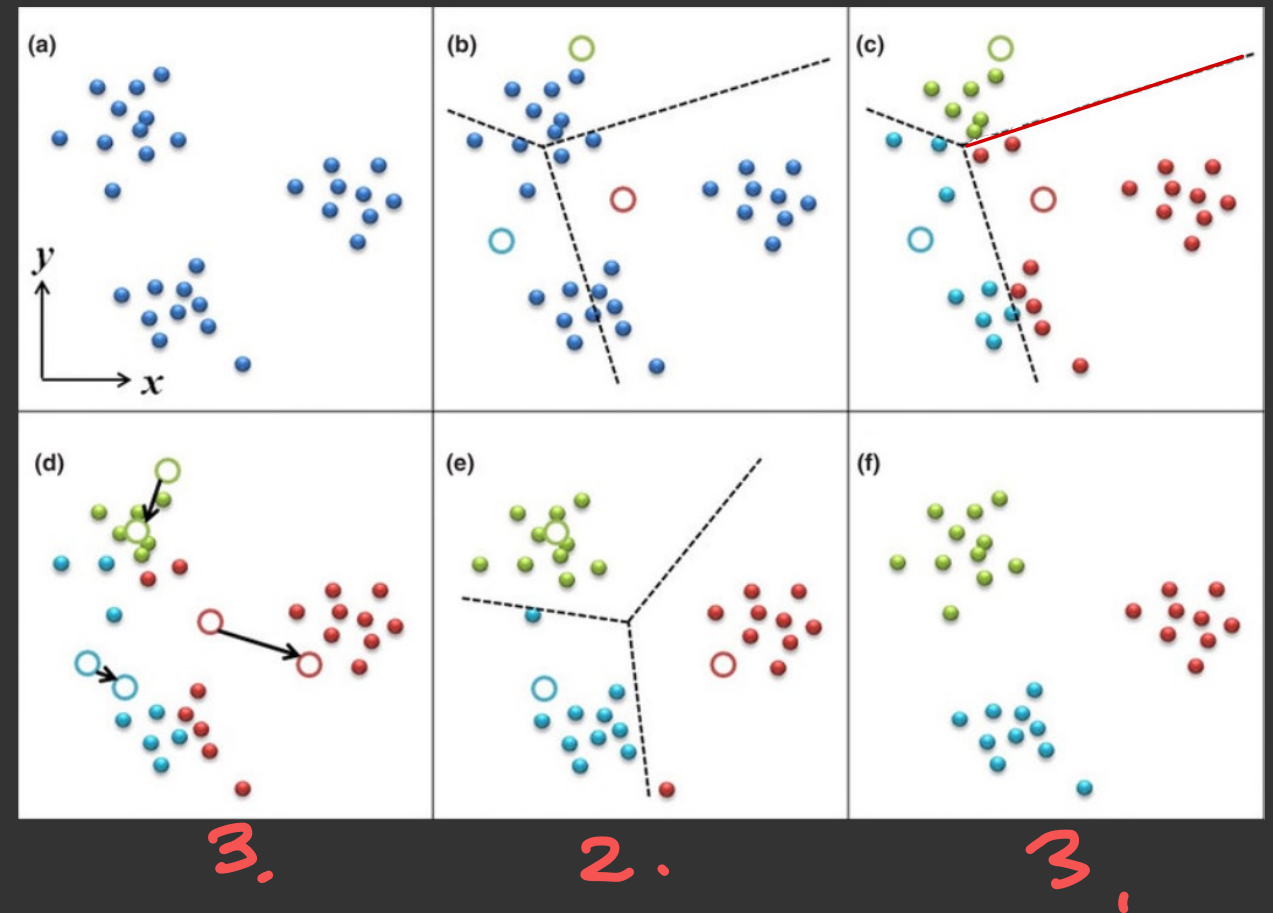
2. $\forall i=1, \dots, n$: Assign $x_i \in \mathcal{C}_j$ if c_j is closest centroid to x_i among $\{c_c\}_{c=1}^k$

3. Update $\forall j=1, \dots, k$:

$$c_j = \underset{c \in \mathbb{R}^k}{\operatorname{argmin}} \sum_{i \in \mathcal{C}_j} d(x_i, c) = \frac{1}{|\mathcal{C}_j|} \sum_{i \in \mathcal{C}_j} x_i$$

if $d(\cdot, \cdot) = \|\cdot - \cdot\|_2^2$

a) randomly among the $\{x_i\}_{i=1}^n$
 b) k-means++ : fancier



▷ Finds local optimum of (*)

▷ Works well for "convex" clusters



To consider:

▷ How to choose nr. of clusters k ?

▷ Which distance to choose (geometry of underlying space)?

▷ Initialization: If prior knowledge available, \rightarrow might be better than random.

▷ Needs a lot of pairwise distances. If $n \gg 10^5$ or so, slow \rightarrow "Minibatch KMeans"

Other clustering methods:

- Spectral Clustering: Based on Laplacian of similarity graph.

- Hierarchical Clustering: Creates trees