Task1 : Clustering    Task2 : feature extration

# IEMS 304 Lecture 8: Unsupervised Learning

Yiping Lu

yiping.lu@northwestern.edu

*Industrial Engineering & Management Sciences*

*Northwestern University*

NORTHWESTERN
UNIVERSITY

Clustering

$\{ (x_i) \}_{i=1}^{n} \longrightarrow$ Classification.
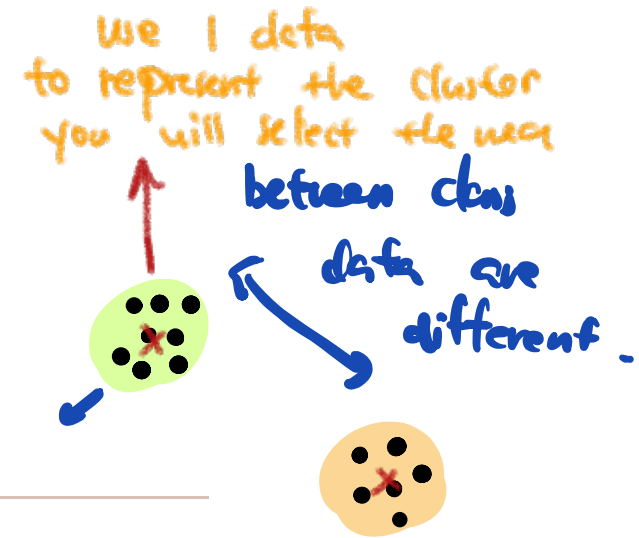
Supervised Classification. $\{ (x_i, z_i) \}_{i=1}^{n}$
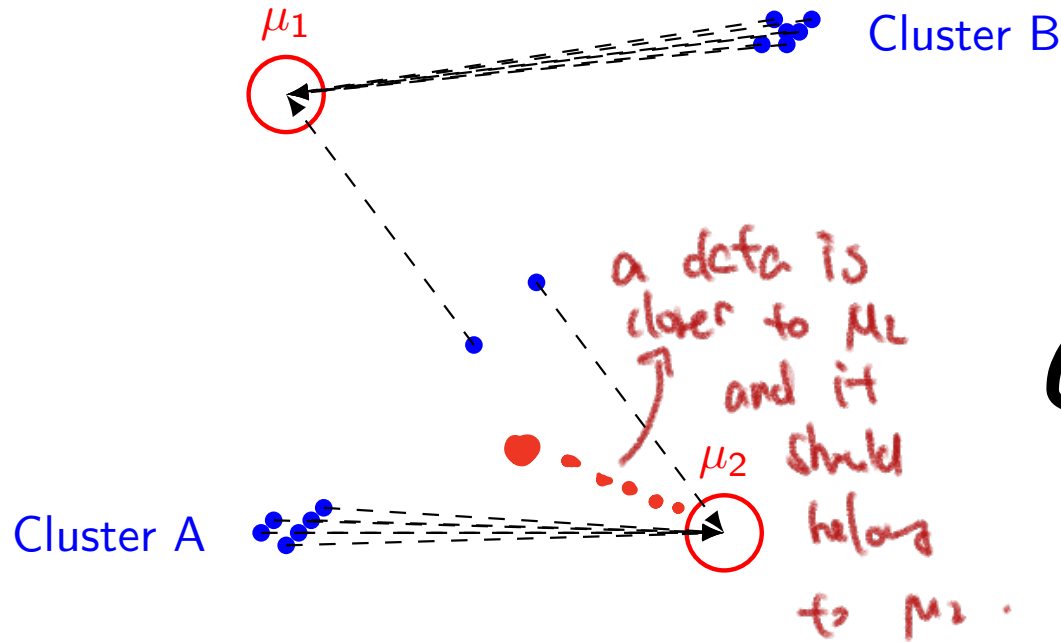
Classification Result.

we 1 data
to represent the cluster
you will select the mean

between class
data are
different.

k-means

in class
are similar

k is the number of
clues.

$\mu_1$

Cluster B

*a data is closer to $M_2$ and it should belong to $M_2$.*

$\mu_2$

Cluster A

① *randomly initialize the k-means.*

② *once you know the k-means, you can do, classification*

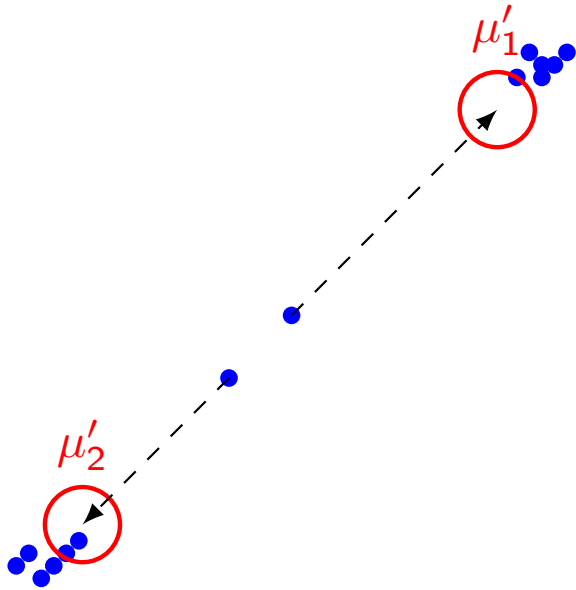*iterative .*

**Assignment Summary (Iteration 1):**

- $\mu_1 = (1, 4.5)$ gets: all Cluster B points (6 pts) + ambiguous point $(2.5, 2.5)$ [total 7 pts].
- $\mu_2 = (4.5, 1)$ gets: all Cluster A points (6 pts) + ambiguous point $(3, 3)$ [total 7 pts].

**Updated centroids (computed as the mean):**

$\mu_1' = \left(\frac{30+2.5}{7}, \frac{30+2.5}{7}\right) \approx (4.643, 4.643)$

$\mu_2' = \left(\frac{6.3+3}{7}, \frac{6.3+3}{7}\right) \approx (1.329, 1.329)$

③ update the cluster mean.

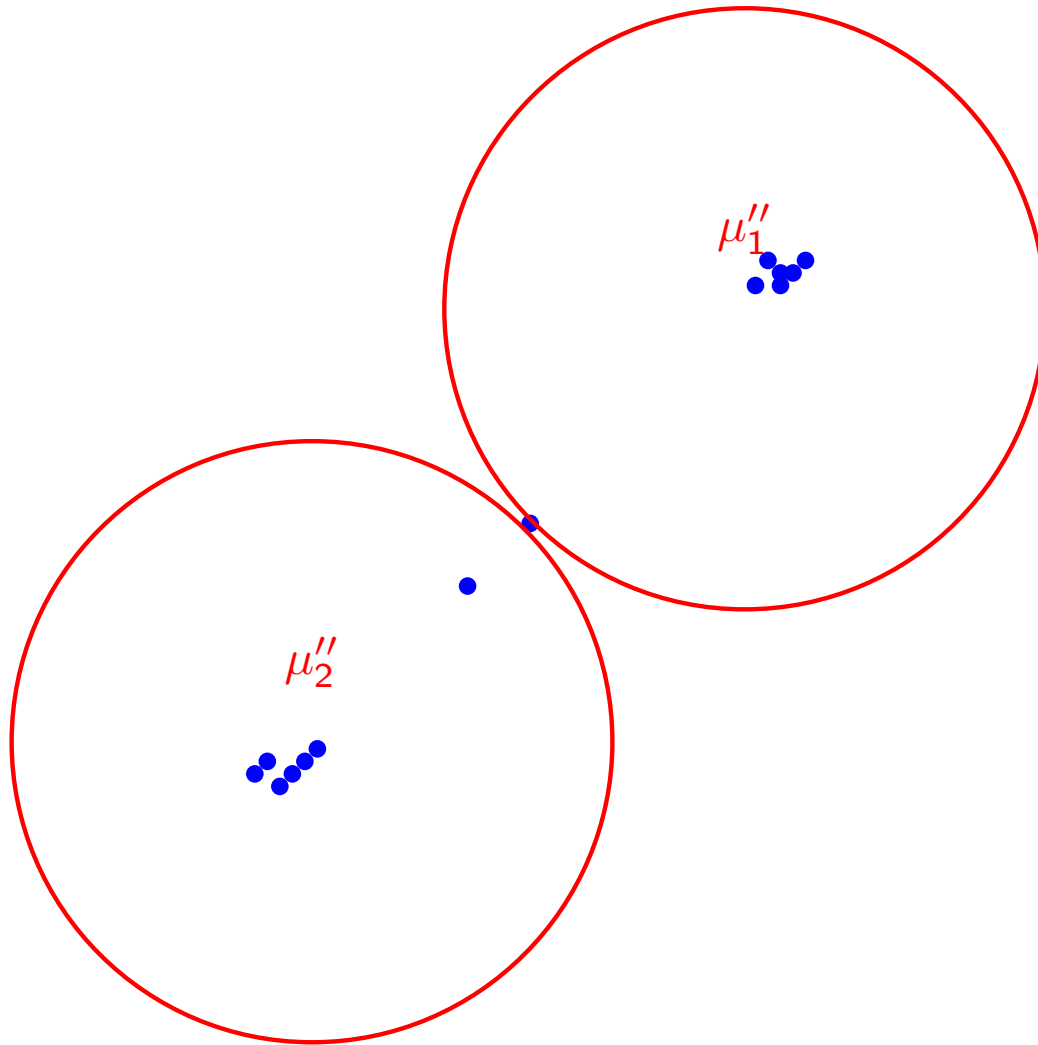cluster mean ⟷ clarification result.

**Reassignment (Iteration 2):**

- $(2.5, 2.5)$ switches from $\mu_1$ to $\mu_2'$ (closer to $(1.329, 1.329)$).
- $(3, 3)$ switches from $\mu_2$ to $\mu_1'$ (closer to $(4.643, 4.643)$).

**New centroids:**

$\mu_1'' = \left(\frac{30+3}{7}, \frac{30+3}{7}\right) = \left(\frac{33}{7}, \frac{33}{7}\right) \approx (4.714, 4.714)$
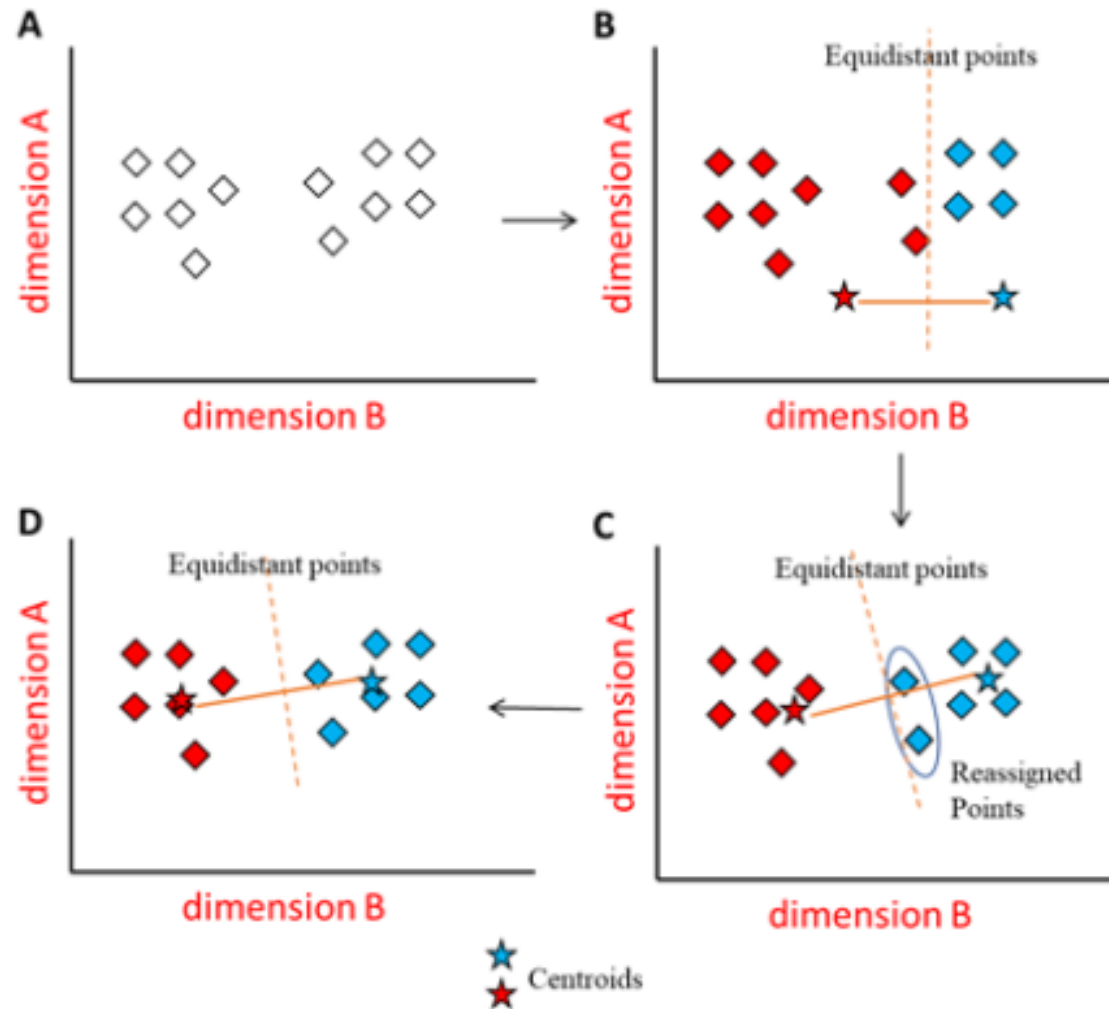
$\mu_2'' = \left(\frac{6.3+2.5}{7}, \frac{6.3+2.5}{7}\right) = \left(\frac{8.8}{7}, \frac{8.8}{7}\right) \approx (1.257, 1.257)$

3

**Convergence:** With centroids $\mu_1'' \approx (4.714, 4.714)$ and $\mu_2'' \approx (1.257, 1.257)$, all data points are now correctly grouped according to their true clusters.

$k-$means aims to minimize the total within cluster (square) distance

$$\min_{\{C_j\},\{\mu_j\}} \sum_{j=1}^{k} \sum_{x \in C_j} \|x - \mu_j\|^2$$

*sum over all the data in cluster $C_j$*

*How you cluster your data*

*the mean of the cluster*

$k-$means as alternating direction optimization algorithm

*fix $M_j$ and update $C_j$*

☐ **Assignment:** Assign each $x$ to its nearest $\mu_j$ (minimizes distance).

*fix $C_j$ and update $M_j$*

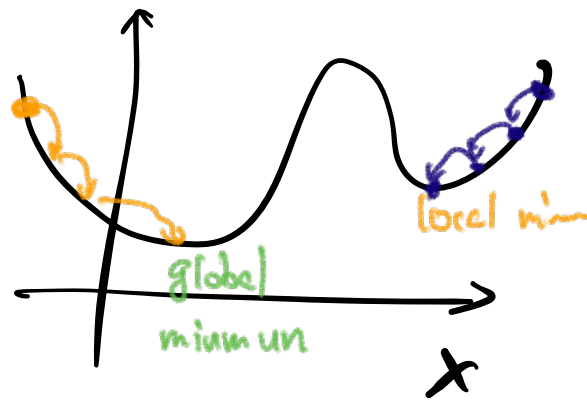☐ **Update:** Recompute $\mu_j$ as the mean of $C_j$ (minimizes variance).

① why we select the mean as the respresenfation

the mean $\mu$ is minicmise $\sum_{x_i \text{ in the cluster}} \|x_i - \mu\|^2$

Convex

linear Regression

Non-Convex

k- means

Neural Networks

global
minimun

local min

k-means ++ .
with high Prob.
goes to global minimum .

we can get.

we use

# Wrong *k* can be Problematic

Gaussian Mixture Data

Elbow Method for Selecting k

# Spectral Clustering

The faliure here is because
we just we mean to represent a cluter
⟹ we assume the cluters are conver

**K-means**

**Spectral clustering**

We first represent data as a weighted graph $G(V, E)$ with weights $w_{ij}$.

Consider the Dirichlet form,

$$\frac{1}{2} \sum_{i,j} w_{ij} (f(i) - f(j))^2 = f^T L f, \quad \text{(Why?)}$$

where L is the graph Laplacian defined as $L = D - W$ (where $D$ is the degree matrix).

*minimizing Dirichlet form.*

*e.g. Connect your data with its k-nearest neighbor*

$W_{AB}$

$W_{CD}$ 0.2

*w = h(~ distance between data)*

*Some function*

*assign a weight to the edge,*

*if two data are less similar then the weight mean smaller.*

B

F

A

C

D

E

What would happen if we minimizing this form?

# Quadratic Function as a Quadratic Form

put all variable as a vector.

$$v^T A v = \begin{pmatrix} x & y \end{pmatrix} \begin{pmatrix} 3 & 2 \\ 2 & 2 \end{pmatrix} \begin{pmatrix} x \\ y \end{pmatrix} = 3x^2 + 2xy + 2xy + 2y^2 = 3x^2 + 4xy + 2y^2.$$

all the coefficient

$$(x.y) \begin{pmatrix} 3x + 2y \\ 2x + 2y \end{pmatrix} = x(3x + 2y) + y(2x + 2y)$$

$$(x \cdot y) \begin{pmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{pmatrix} \begin{pmatrix} x \\ y \end{pmatrix} = a_{11}x^2 + (a_{12} + a_{21})xy + a_{22}y^2$$

fix $a_{12} = a_{21}$   which means $A$ is symmetric

$$A^T = A.$$

11

Consider the Dirichlet form:

$$\frac{1}{2} \sum_{i,j} w_{ij} \left( f(i) - f(j) \right)^2 = \frac{1}{2} \sum_{i,j} w_{ij} \left[ f(i)^2 - 2f(i)f(j) + f(j)^2 \right].$$

☐ terms involving $f(i)^2$:

$$\frac{1}{2} \left( \sum_{i,j} w_{ij} f(i)^2 + \sum_{i,j} w_{ij} f(j)^2 \right)$$

$$= \sum_{i} f(i)^2 \sum_{j} w_{ij} = \sum_{i} d_i f(i)^2.$$

☐ The cross term simplifies to:

$$- \sum_{i,j} w_{ij} f(i) f(j).$$

$$\frac{1}{2} \sum_{i,j} w_{ij} \left( f(i) - f(j) \right)^2 = \sum_{i} d_i f(i)^2 - \sum_{i,j} w_{ij} f(i) f(j).$$

At the same time,

$$f^T L f = \sum_{i} d_i f(i)^2 - \sum_{i,j} w_{ij} f(i) f(j), \text{ where } L = D - W,$$

12

# Understanding the Dirichlet Form

## Definition

The Dirichlet form on a graph is defined as:
$$\frac{1}{2} \sum_{i,j} w_{ij}(f(i) - f(j))^2 = f^T L f.$$

- It sums the squared differences of the function values $f(i)$ over every edge, weighted by $w_{ij}$.

- A small value of $f^T L f$ indicates that neighboring nodes (with high similarity $w_{ij}$) have similar function values.

- Minimizing the Dirichlet form under constraints leads to smooth functions on the graph, thus revealing inherent cluster structure.

Handwritten annotations on graph:
- $W_{AB}$, $W_{AC}$ edge labels
- $0.2$ on dashed edge between C and D
- $d_C = W_{AC} + W_{BC} + W_{CD}$

Handwritten matrix (top right):

|   | A | B | C | D |
|---|---|---|---|---|
| A |   | $W_{AB}$ | $W_{AC}$ |   |
| B | $W_{AB}$ |   | $W_{BC}$ |   |
| C | $W_{AC}$ | $W_{BC}$ |   | $W_{CD}$ |
| D |   |   | $W_{CD}$ |   |

$0.2$

**Step 1: Define the Matrices**

- **Weighted Adjacency Matrix** $W$:
  For each edge $(i,j)$, $w(i,j) = 1$
  except for the edge between $C$ and
  $D$ where $w(C, D) = 0.2$.

- **Degree Matrix** $D$: Diagonal with
  $d_A = 2, d_B = 2, d_C = 2.2, d_D = 2.2, d_E = 2, d_F = 2$

Handwritten: $\begin{bmatrix} d_A & d_B & d_C & d_D & d_E & d_F \end{bmatrix}$

**Step 2: Compute the Graph Laplacian**

$$L = D - W$$

$$= \begin{pmatrix} 2 & -1 & -1 & 0 & 0 & 0 \\ -1 & 2 & -1 & 0 & 0 & 0 \\ -1 & -1 & 2.2 & -0.2 & 0 & 0 \\ 0 & 0 & -0.2 & 2.2 & -1 & -1 \\ 0 & 0 & 0 & -1 & 2 & -1 \\ 0 & 0 & 0 & -1 & -1 & 2 \end{pmatrix}.$$

14

What is the smallest eigenvalue/eigenvectors of the graph laplacian?

What would happen if we have *l*-connected component

$$\min \quad \cancel{\text{n}} \enclose{circle}{f^\top L f} \qquad \text{s.t. } f^\top \mathbb{1} = 0, \|f\|_2 = 1$$

*Closed form!!! is the smallest eigenvector!*

Then run a $k-$means on the spectral clustering representation $f$. (homework)

$$f_1 + f_2 + \cdots + f_n = 0$$

① **Case1** Do No classification.
$(-f(i) = 1$ to all the data. $f^\top L f = 0)$

② **Case2** If half the value of $f$, then the quadratic function will greater
( the solution can be $f(i) = 0$ for all data)

$$f_1^2 + f_2^2 + \cdots + f_n^2 = 1$$

16

graph laplacian → several smallest eigen vectors



Data points

epsilon–graph, epsilon=0.3

kNN graph, k = 5

Mutual kNN graph, k = 5

The procedure of spectral Clustering.

1. Construct a graph based on if two data are similar

2. Compute the matrix L.

3. Compute the eigenvectors (the smallest four eigenvalue)
   of L

4. the Cluster is easier now. you can run a
   k-means



eign1        eigen L

A    2    4           • (2,-1)
B    2    -1
C    2    -1
D    4    2           (-1, 2)
E    -1   2              •
F    -1   2
G    -1   -1          (-1, -1)
H    -1   -1
I    -1   -1              •

# Dimension Reduction

# Framework of Dimension Reduction

Data: $X \in \mathbb{R}^n$      $\overset{\phi(x)}{\text{feature}} \in \mathbb{R}^m$

(We may consider $m$ is much smaller than)

e.g. in the face dataset.

     $n$: the number of pixel of face image

     $m$: $m = 2$: pxes.

$\phi(x)$ should include most of the information in your data $X$. (compression).

$$\underset{\text{data}}{X} \xrightarrow[\substack{\text{compress} \\ \text{the information}}]{} \underset{\text{feature}}{\phi(x)} \xrightarrow{\text{reconstruction}} r(\phi(x))$$

low dimension

we aim: $r(\phi(x)) \approx X$      auto-encoder

encoder          decoder

objective function $X_1 \dots X_n$ .

$$\min_{r, \phi} \sum_{i=1}^{n} \| r(\phi(x_i)) - x_i \|^2$$

If $r, \phi$ are Neural Networks.

$\Rightarrow$ auto-encoder

If $r, \phi$ are linear function
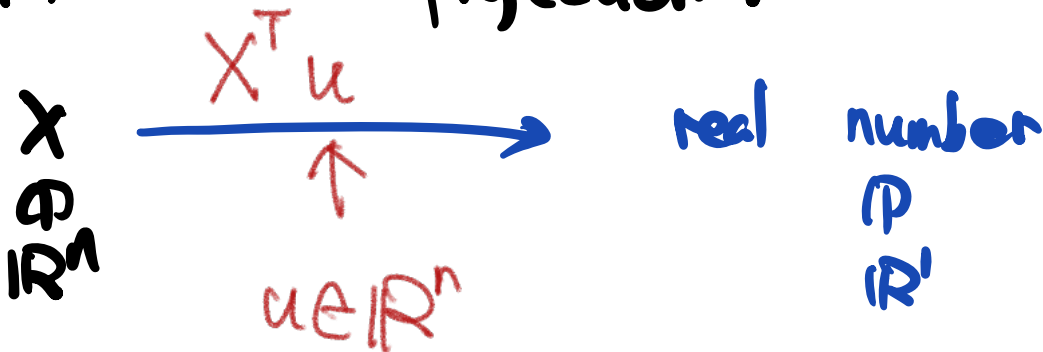
$\Rightarrow$ Principal Component Analysis

Close form !!
(PC4)

---

Review : Projection .

$X$ :   linear projection .

$$X \xrightarrow{\quad X^T u \quad} \text{real number}$$

$X$
$\in$
$\mathbb{R}^n$

$u \in \mathbb{R}^n$

$\mathbb{P}$
$\mathbb{R}^1$

$X^T u$ : is the inner/dot product

# Aim: X is the data.

## $U_1 \cdots U_m$ : (encoding)

$$X \rightarrow (X^T u_1, \; X^T u_2, \; \cdots, \; X^T u_m)$$

$X^T u_1$: scalar

$\mathbb{R}^m$

## Decoding $u_1' \cdots u_m' \in \mathbb{R}^n$

$$(a_1, \cdots, a_m) \rightarrow a_1 u_1' + a_2 u_2' + \cdots + a_m u_m'$$

$a_1$: scalar

$\mathbb{R}^m$

## Close-form Solution !!!!!!

$$U_i = U_i'$$

$U_1 \cdots U_m$ are the top
eigen vectors of $X X^T$ ($n \times n$ matrix)

covariance

$$\underline{X} = \begin{bmatrix} X_1 & \cdots & X_k \end{bmatrix} \in \mathbb{R}^{n \times k}$$

k-data: $X_i$, $X_i \in \mathbb{R}^n$

$\underline{XX^T}$ is a matrix of size $\mathbb{R}^{n \times n}$

---

Fact 1. $XX^T$ is a symmetric matrix

Fact 2. $U_1 \cdots U_m$ are orthogonal.

$$U_i^T U_j = 0$$

---

$$U = \begin{bmatrix} U_1 & \cdots & U_m \end{bmatrix}$$

Question! What is the projection of data matrix $X$ to the space spanned by $U$

$$\Rightarrow (U^TU)^{-1} U^T X = U^T X$$

$\underbrace{\quad}_{=I_m}$

Let's use the info that $U$ is

orthogonal $\begin{cases} u_i^T u_i = 1 \\ u_i^T u_j = 0 \end{cases}$

$$U^T U = \begin{bmatrix} u_1^T \\ u_2^T \\ \vdots \\ u_m \end{bmatrix} \begin{bmatrix} u_1 & u_2 & \cdots & u_m \end{bmatrix}$$
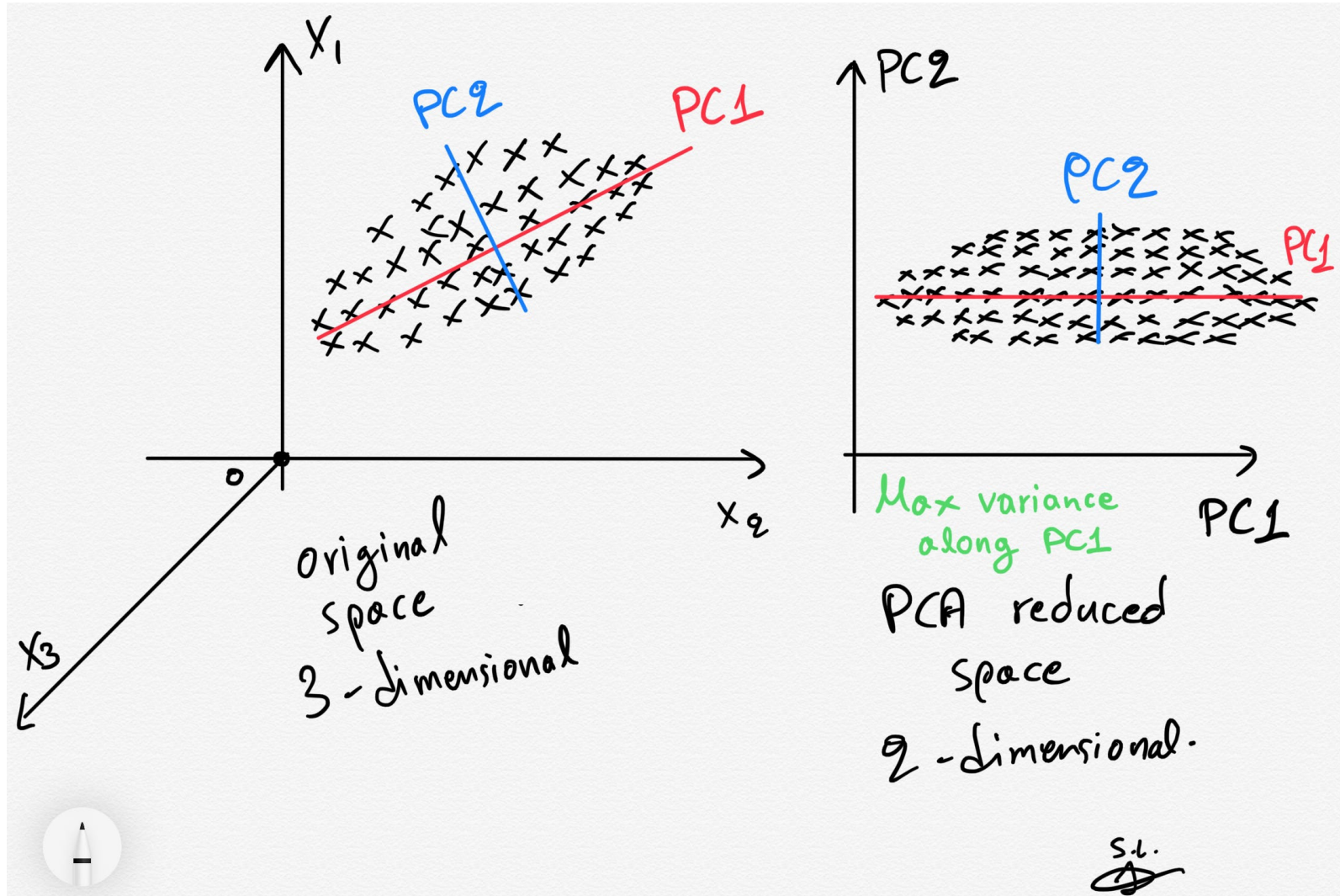
$$= \begin{bmatrix} 1 & 0 & \cdots & & 0 \\ 0 & 1 & \cdots & & 0 \\ \vdots & \vdots & \ddots & & \vdots \\ 0 & 0 & \cdots & & 1 \end{bmatrix} = I_m$$
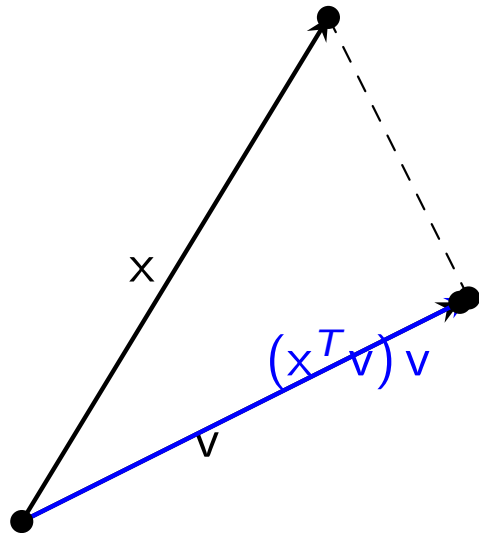
# PCA:

the projection that preserves most of the information of your data is the projection to the top eigenspaces of the covariance matrix,
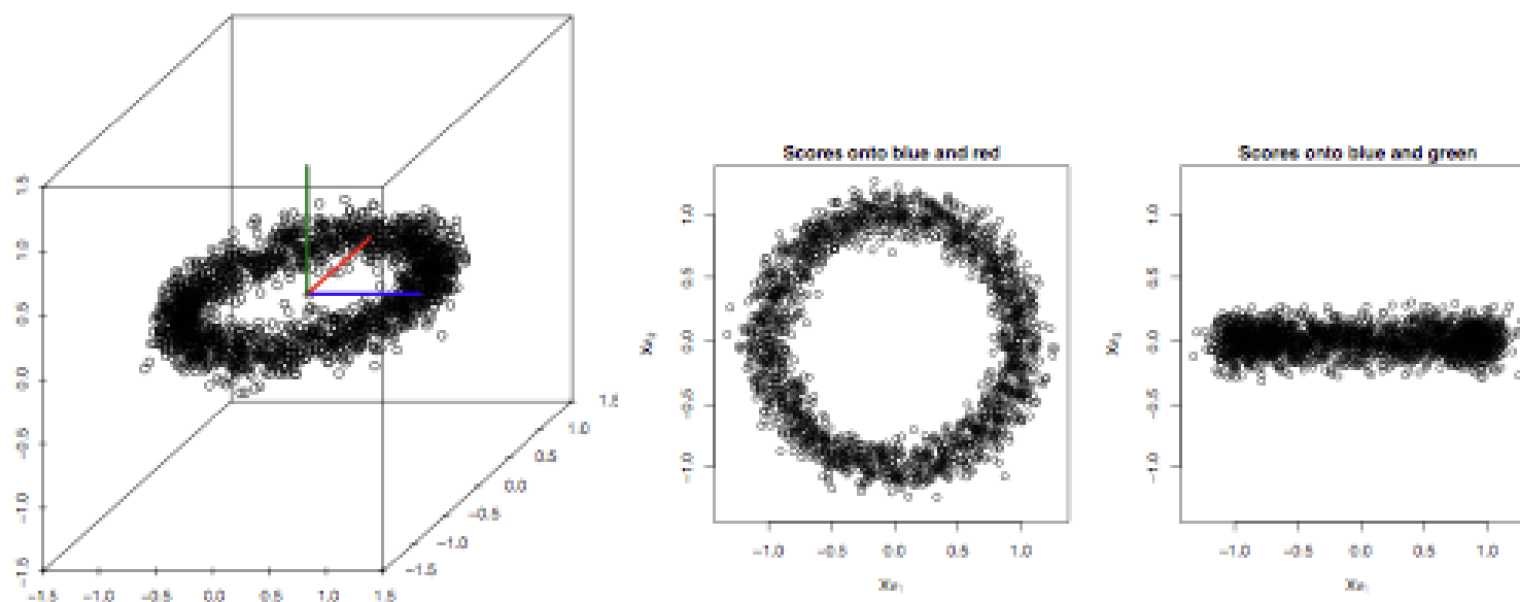
largest!

# number of basis $\Rightarrow$ bias-variance.

- $x^T v \in \mathbb{R}$ : score
- $(x^T v) v \in \mathbb{R}^p$ : projection

Example: $X \in \mathbb{R}^{2000 \times 3}$, and $v_1, v_2, v_3 \in \mathbb{R}^3$ are the unit vectors parallel to the coordinate axes



Not all linear projections are equal! What makes a good one?

We have $n$ $d$-dimensional data points $x_1, x_2, \ldots, x_n \in \mathbb{R}^d$ and a parameter $k \in \{1, 2, \ldots, d\}$. We assume that the data is centered, meaning that $\sum_{i=1}^{n} x_i = 0$. (How to do that?)

$\boxed{v^T \text{ matrix } v} \Rightarrow$ quadratic function of $v$

**AIM.** Find directions that maximize the information preserved

The output of the method is defined as $k$ orthonormal vectors $v_1, v_2, \ldots, v_k$ — the "top $k$ principal components" — that maximize the objective function :

$$\frac{1}{n} \sum_{i=1}^{n} \sum_{j=1}^{k} (x_i \cdot v_j)^2.$$

$\underbrace{(X^T v)^T (X^T v)}_{\text{Sum of square}} = v^T X X^T v$

---

**Question:** Why we want the principal components orthonormal?

(or I prefere the most varying direction of my data).

max the variance of the data after the projection

Let $A = [v_1, \cdots, v_k]$ where $v_1, \cdots, v_k$ are orthonormal. **<u>Remind.</u>** Least square solution: $A\beta \approx b$, then $\beta = (A^\top A)^{-1} A^\top b$ Then $A\beta = A(A^\top A)^{-1} A^\top b$

<u>**Review**</u>. Orthonormal means $A^\top A = I$

**<u>Check.</u>** Project $b$ to $\text{span}\{v_1, \cdots, v_k\}$ means

$$\langle v_1, b \rangle v_1 + \langle v_2, b \rangle v_2 + \cdots + \langle v_k, b \rangle v_k$$

# Matrix Formulation

**Matrix Formulation:** Define $V \in \mathbb{R}^{d \times k}$ with columns $v_1, \ldots, v_k$, representing the $k$ principal components.

The total variance captured when projecting the data onto the subspace spanned by $V$ is
$$\frac{1}{n}\|XV\|_F^2 = \mathrm{tr}\left( V^T \left( \frac{1}{n}X^T X \right) V \right) = \mathrm{tr}(V^T SV),$$
where $S = \frac{1}{n}X^T X$ is the covariance matrix.

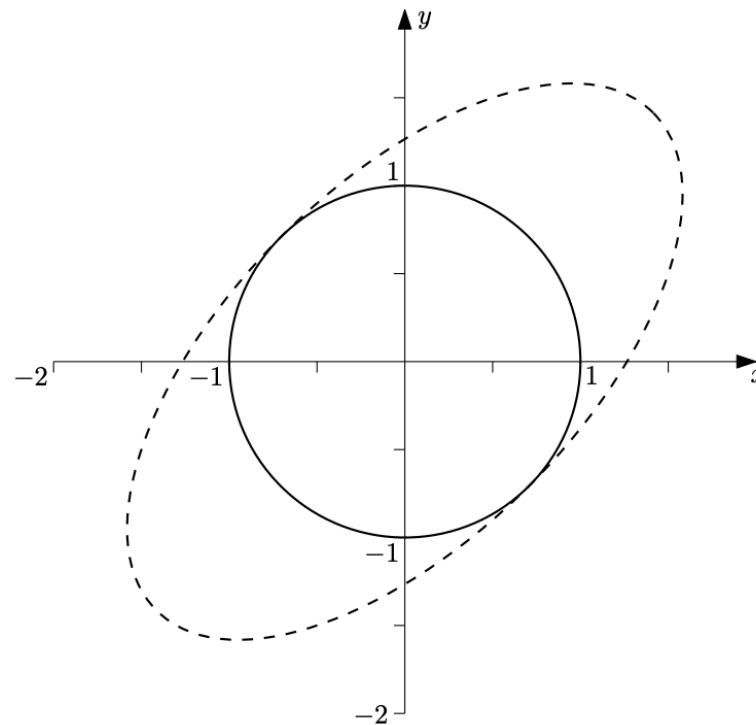Note that $\|A\|_F^2 = \mathrm{tr}(A^T A)$ For $A = XV$, we have:
$$\|XV\|_F^2 = \mathrm{tr}\left((XV)^T (XV)\right) = \mathrm{tr}(V^T X^\top XV). \qquad (\text{for } \mathrm{tr}(AB) = \mathrm{tr}(BA))$$

$$\max_{V \in \mathbb{R}^{d \times k}} \mathrm{tr}(V^T SV) \quad \text{subject to} \quad V^T V = I_k.$$

# Matrix Formulation

# Covariance Matrix: Rotation on Principal Component

$$\begin{pmatrix} \frac{3}{2} & -\frac{1}{2} \\ \frac{1}{2} & \frac{3}{2} \end{pmatrix} = \underbrace{\begin{pmatrix} \frac{1}{\sqrt{2}} & -\frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \end{pmatrix}}_{\text{rotate back } 45°} \cdot \underbrace{\begin{pmatrix} 2 & 0 \\ 0 & 1 \end{pmatrix}}_{\text{stretch}} \cdot \underbrace{\begin{pmatrix} \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \\ -\frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \end{pmatrix}}_{\text{rotate clockwise } 45°}$$
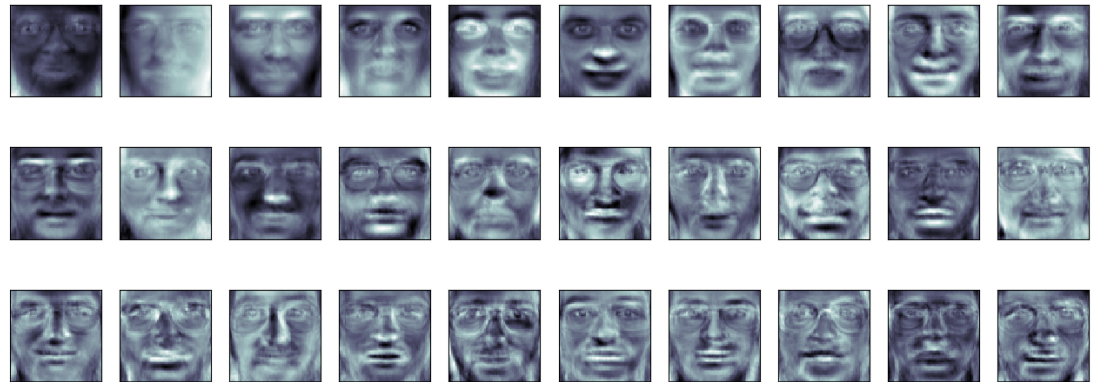
# PCA as Top Eigenvectors

PCA boils down to computing the $k$ eigenvectors of the covariance matrix $X^\top X$ that have the largest eigenvalues.
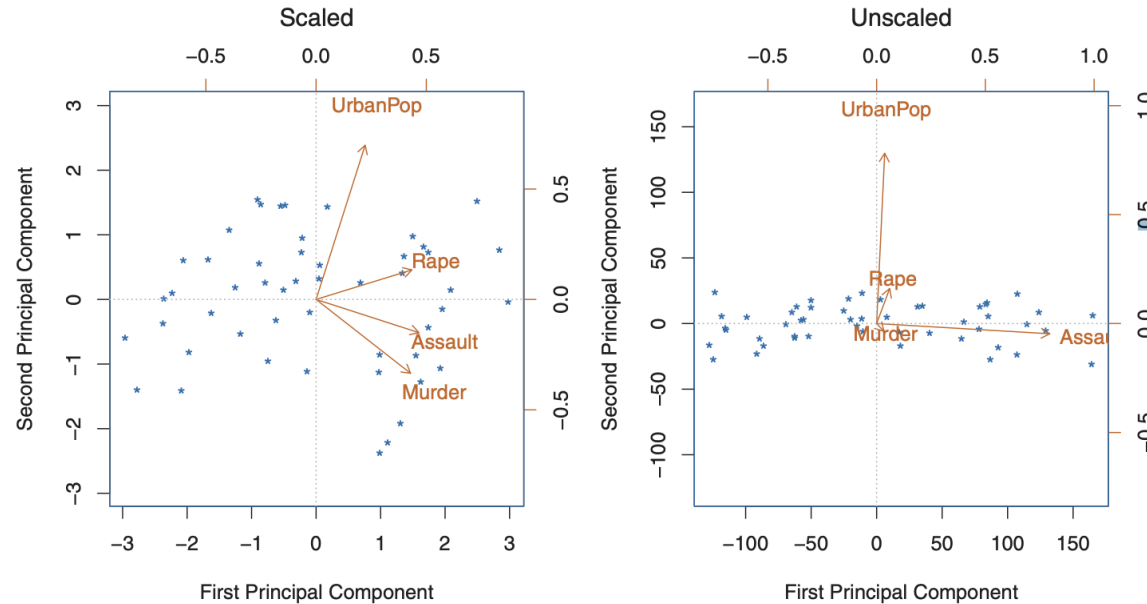
The components ("eigenfaces") are ordered by their importance from top-left to bottom-right. We see that the first few

components seem to primarily take care of lighting conditions; the remaining components pull out certain identifying

features: the nose, eyes, eyebrows, etc.

# Normalize Your Data



**Murder**, **Rape**, and **Assault** are reported as the number of occurrences per 100, 000 people, and **UrbanPop** is the percentage of the state's population that lives in an urban area. These four variables have variance 18.97, 87.73, 6945.16, and 209.5