

GOAL: UNSUPERVISED LEARNING

↳ REPRESENTATION LEARNING

- Given a set of "data points", "understand"
some thing "useful" about them.

DATA POINTS → VECTORS

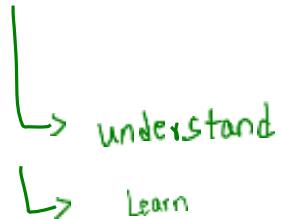
$$\begin{bmatrix} 180 \\ 75 \\ 37 \end{bmatrix} \rightarrow \begin{array}{l} \text{height} \\ \text{weight} \\ \text{age} \end{array} \in \mathbb{R}^3$$

understand ? useful ?

- RUNNING THEME

"COMPREHENSION is COMPRESSION"

- GEORGE CHAITIN



Problem:

Input: $\{x_1, x_2, \dots, x_n\}$ $x_i \in \mathbb{R}^d$ $d \leftarrow \# \text{ features}$

Output: Some "COMPRESSED" representation of the dataset.

Example:

$$\left\{ \begin{array}{c} x_1 \\ \left[\begin{array}{c} -5 \\ 2.5 \end{array} \right] \end{array}, \begin{array}{c} x_2 \\ \left[\begin{array}{c} -2 \\ 1 \end{array} \right] \end{array}, \begin{array}{c} x_3 \\ \left[\begin{array}{c} 4 \\ -2 \end{array} \right] \end{array}, \begin{array}{c} x_4 \\ \left[\begin{array}{c} 6 \\ -2.5 \end{array} \right] \end{array} \right\}$$

Question: How many numbers are needed to store this dataset? 8

Note:
Using this representation
we can "exactly"
reconstruct the dataset.

Representative

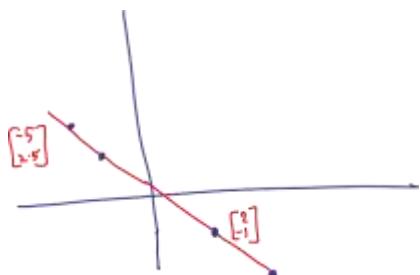
$$\begin{bmatrix} -2 \\ 1 \end{bmatrix}$$

Co-efficient

$$\left\{ 2.5, 1, -2, -2.5 \right\}$$

6

GEOMETRIC VIEW



R

$$\begin{bmatrix} 2 \\ 1 \end{bmatrix}$$

C

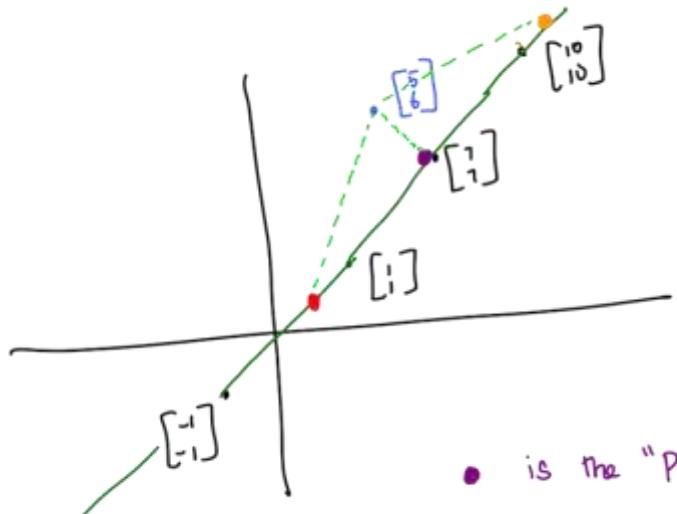
$$\left\{ 2.5, 1, -2, -2.5 \right\}$$

$$\begin{bmatrix} -2/\sqrt{5} \\ 1/\sqrt{5} \end{bmatrix}$$

$$\left\{ 2.5\sqrt{5}, \sqrt{5}, -2\sqrt{5}, -2.5\sqrt{5} \right\}$$

NOTE: R can be any point
on the line (other than $\begin{bmatrix} 0 \\ 0 \end{bmatrix}$)

original set of numbers 2^n
 Compressed representation $2 + n$



$$\text{REPR} \left\{ \begin{bmatrix} 1 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \end{bmatrix} \right\}$$

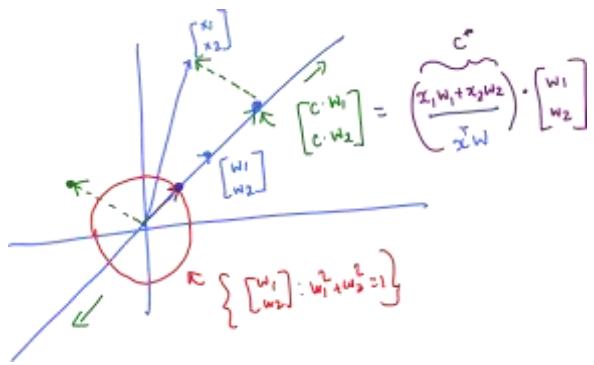
$$\text{Co-eff} \left\{ (1, 1), (1, 0), (-1, 1), (10, 10) \right\}$$

$$7 \times \begin{bmatrix} 1 \\ 0 \end{bmatrix} + 7 \begin{bmatrix} 0 \\ 1 \end{bmatrix} = \begin{bmatrix} 7 \\ 7 \end{bmatrix}$$

- is the "PROJECTION" of • on the green line.

$\text{length} \left(\begin{bmatrix} a \\ b \end{bmatrix} \right) = \sqrt{a^2 + b^2}$

 $\begin{bmatrix} a \\ b \end{bmatrix}$



$$\begin{bmatrix} c \cdot w_1 \\ c \cdot w_2 \end{bmatrix} + \begin{bmatrix} ? \\ ? \end{bmatrix} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \Rightarrow \text{ERROR VECTOR}$$

minimize c $(x_1 - c \cdot w_1)^2 + (x_2 - c \cdot w_2)^2$

$$f(c) = (x_1)^2 + c^2 (w_1)^2 - 2c x_1 w_1 + (x_2)^2 + c^2 (w_2)^2 - 2c x_2 w_2$$

$$f(c) = c^2 w_1^2 - 2x_1 w_1 + 2c^{w_2^2} - 2x_2 w_2$$

$$= 0$$

$$\Rightarrow C = \frac{x_1 w_1 + x_2 w_2}{w_1^2 + w_2^2}$$

INNER
PRODUCT /
DOT - PRODUCT

$$\left(\begin{bmatrix} x_1 \\ x_2 \end{bmatrix}, \begin{bmatrix} w_1 \\ w_2 \end{bmatrix} \right) = x_1 w_1 + x_2 w_2$$

$$\begin{matrix} x_1 & w_1 \\ \vdots & \vdots \\ x_d & w_d \end{matrix} = \underline{x_1 w_1 + \dots + x_d w_d} = \underline{\underline{x^T w}}$$

Summary

- "Comprehension is compression"

- Rep, co-efficients

- $\text{Proxy} \rightarrow (x^T w) \begin{bmatrix} w_1 \\ w_2 \end{bmatrix}$ (Projection)

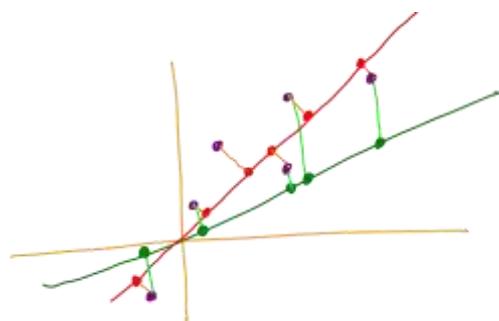
Question

- Given a line & a point, we know how to find proxy for the point along the line

- Who gives us the line ?

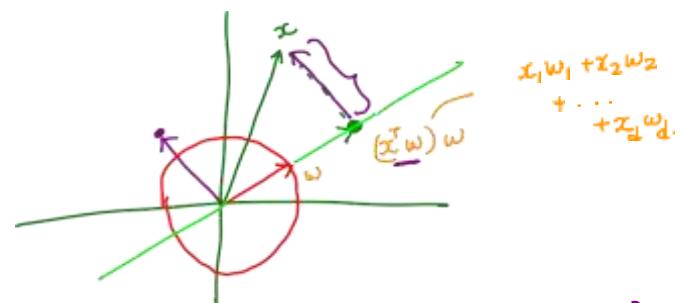
LECTURE 3

Goal: Develop a way to find a "compressed" representation of data
when points non-necessarily fall on a line
↳ RULE not an EXCEPTION



Goal : Find the line that has least "Error" \hookrightarrow Reconstruction

Dataset $\{x_1, x_2, \dots, x_n\} \quad x_i \in \mathbb{R}^d$



$$\begin{aligned} \text{ERROR}(\text{line}, \text{Dataset}) &= \sum_{i=1}^n \text{error}(\text{line}, x_i) \\ &= \sum_{i=1}^n \text{length}^2(x - (x^T w)w) \end{aligned}$$

$$(x^T w) \cdot \vec{w} + ? = \vec{x}$$

$$x - (x^T w)w$$

$$\text{error}(\omega, D) = \sum_{i=1}^n (x - (\underline{x}^\top \omega) \underline{\omega})^\top (x - (\underline{x}^\top \omega) \cdot \underline{\omega})$$

$$= \sum_{i=1}^n \|x - (\underline{x}^\top \omega) \cdot \underline{\omega}\|^2$$

Goal:

$$\min_{\omega: \underline{w}^\top \underline{w} = 1} \sum_{i=1}^n \|x - (\underline{x}^\top \omega) \cdot \underline{\omega}\|^2$$

$$f(\omega) = \frac{1}{n} \sum_{i=1}^n (\underline{x} - (\underline{x}^\top \underline{\omega}) \underline{\omega})^\top (\underline{x} - (\underline{x}^\top \underline{\omega}) \underline{\omega})$$

$$= \frac{1}{n} \sum_{i=1}^n (\underline{x}^\top \underline{x} - (\underline{x}^\top \underline{\omega})^2 - (\underline{x}^\top \underline{\omega})^2 + (\underline{x}^\top \underline{\omega})^2)$$

$$= \frac{1}{n} \sum_{i=1}^n (\underline{x}^\top \underline{x} - (\underline{x}^\top \underline{\omega})^2)$$

$$g(\omega) = -\frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i^\top \omega)^2$$

$$\begin{aligned}
 \max_{\omega: \|\omega\|=1} f_i(\omega) &= \frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i^\top \omega)^2 \\
 &= \frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i^\top \omega) \cdot (\mathbf{x}_i^\top \omega) \\
 &= \frac{1}{n} \sum_{i=1}^n (\omega^\top \mathbf{x}_i) \cdot (\mathbf{x}_i^\top \omega) \\
 &= \frac{1}{n} \sum_{i=1}^n \underbrace{\omega^\top}_{d \times 1} \underbrace{(\mathbf{x}_i \mathbf{x}_i^\top)}_{d \times d} \underbrace{\omega}_{d \times 1} \\
 &= \omega^\top \left[\frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top \right] \omega
 \end{aligned}$$

$$\begin{array}{c}
 \left[\begin{array}{c} x_1 \\ \vdots \\ x_d \end{array} \right]_{d \times 1} \quad \left[\begin{array}{c} \omega_1 \\ \vdots \\ \omega_d \end{array} \right]_{d \times 1} \\
 \left[\begin{array}{c} x_1 \dots x_d \end{array} \right]_{1 \times d}
 \end{array}$$

$$\begin{aligned}
 &f(a+b+c) \\
 &\omega^\top a + \omega^\top b + \omega^\top c \\
 &\underline{\frac{1}{n} \sum_{i=1}^n \omega^\top (\mathbf{x}_i^\top \omega)^2}
 \end{aligned}$$

$$\max_{\omega} \quad \omega^T C \omega$$

$\omega:$

$$\omega^T \omega = 1$$

-④

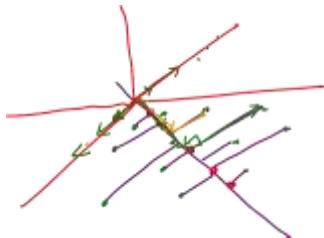
$$C = \frac{1}{n} \sum_{i=1}^n x_i x_i^T$$

Co-variance matrix

- what is the solution to ④?
- For the moment, let's say it's simple to solve ④.
So are we done?

$$w_1 = (x^T w)$$

$$w_2 = (x^T w_2)$$



$$x \in \mathbb{R}^d$$

↓

Find w

x be comes

$$(x^T w) w$$

↓

$$x - (x^T w) w$$

might
have
information.

Possible Algorithm

Input:
 $\{x_1, \dots, x_n\}$
 $x_i \in \mathbb{R}^d$

→ Find "best" line
 $w, \in \mathbb{R}^d$

→ Replace
 $x \rightarrow x - (x^T w) w$

→ Repeat to
 obtain w_2, \dots

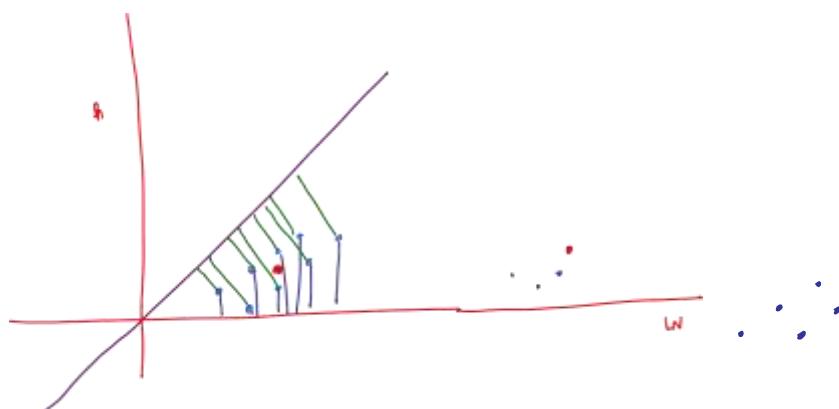
RECAP

$$\begin{aligned}
 & x \in \mathbb{R}^d \\
 & \downarrow \text{FIND } w \\
 & x \in \mathbb{R}^d \rightarrow (x^T w) \cdot w \in \mathbb{R}^d \\
 & \downarrow \text{RESIDUE / ERROR} \\
 & \underline{x - (x^T w) \cdot w} \\
 & \text{MIGHT NOT BE} \\
 & \text{ERROR BUT HAS} \\
 & \text{INFORMATION}
 \end{aligned}$$

- INPUT DATASET POSSIBLE ALGORITHM
- $\{x_1, \dots, x_n\}$ $x_i \in \mathbb{R}^d$
- FIND "BEST" LINE $w_i \in \mathbb{R}^d$
 - REPLACE $x_i \rightarrow \underline{x_i - (x_i^T w_i) w_i + i}$
 - REPEAT procedure to obtain w_2, w_3, \dots

QUESTIONS

- We look only for lines passing through ORIGIN . Is this OK?
- How to solve $\max_w w^T C w$?
w :
 $\|w\|^2 = 1$
- How many times to repeat procedure ?
- Where exactly is "COMPRESSION" happening ?
- What "REPRESENTATIONS" are we learning ?



$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$$

$$x'_i \rightarrow x_i - \bar{x} + \bar{b}$$

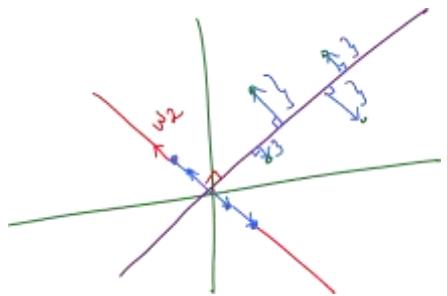
$$\frac{1}{n} \sum_{i=1}^n x'_i = 0$$

$$\frac{1}{n} \sum_{i=1}^n (x'_i - \bar{x}) = \underbrace{\frac{1}{n} \sum_{i=1}^n x_i}_{\bar{x}} - \underbrace{\frac{1}{n} \sum_{i=1}^n \bar{x}}_{\bar{x}} = 0$$

$$\left. \begin{array}{l}
 x_1' \rightarrow x_1 - (x_1^T w_1) w_1 \\
 x_2' \rightarrow x_2 - (x_2^T w_1) w_1 \\
 \vdots \\
 x_n' \rightarrow x_n - (x_n^T w_1) w_1
 \end{array} \right\} \quad \begin{array}{l}
 w_2 = \max_w \quad w^T C' w \\
 \|w\|^2 = 1
 \end{array}$$

$$C' = \frac{1}{n} \sum_{i=1}^n x_i' x_i'^T$$

QUESTION : What can we say about
 w_1 and w_2 ?



w_1

$$w_1^T w_2 = 0$$

By continuing the same procedure,
we get $\{w_1, w_2, \dots, w_R\}$ of

ORTHOGONAL VECTORS
(ORTHO NORMAL)

YES after Round 1

$$\left\{ \underbrace{x_1 - (x_1^T w_1) w_1}_{w_2}, \dots, \underbrace{x_n - (x_n^T w_1) w_1}_{w_2} \right\}.$$

$w_1^T w_2 = 0$.

RESIDUES after Round 2.

$$\left\{ x_1 - (x_1^T w_1) w_1 - \left((x_1 - (x_1^T w_1) w_1)^T w_1 \right) w_2, \dots \right\}$$

$$\left\{ x_1 - (x_1^T w_1) w_1 - \left(x_1^T w_2 - (x_1^T w_1) (w_1^T w_2) \right) w_2, \dots \right\}$$

$$\left\{ \underbrace{x_i - (x_i^T w_1) w_1 - (x_i^T w_2) w_2}_{\text{Residues after } d \text{ rounds}}, \dots \right\}$$

$x_i \in \mathbb{R}^d$

$$\text{Residues after } d \text{ rounds}$$

$$x_i \rightarrow x_i - (x_i^T w_1) w_1 - (x_i^T w_2) w_2 - \dots - (x_i^T w_d) w_d = 0.$$

$$x_i = \underbrace{(x_i^T w_1)}_{\in \mathbb{R}} w_1 + \underbrace{(x_i^T w_2)}_{\in \mathbb{R}^d} w_2 + \dots + \underbrace{(x_i^T w_d)}_{\in \mathbb{R}^d} w_d. \quad \forall i$$

$$x_i = \underbrace{\left(x_i^T \begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \right) \begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}}_{\text{wavy line}} + \underbrace{\left(x_i^T \begin{bmatrix} 0 \\ 1 \\ \vdots \\ 0 \end{bmatrix} \right) \begin{bmatrix} 0 \\ 1 \\ \vdots \\ 0 \end{bmatrix}}_{\text{solid line}} + \cdots + \underbrace{\left(x_i^T \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 1 \end{bmatrix} \right) \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 1 \end{bmatrix}}_{\text{solid line}}$$

What have we gained?

- If data lies in a "low" dimensional space, then residues become 0 much earlier than d rounds.

$$\text{Dataset} = \{x_1, \dots, x_n\} \quad x_i \in \mathbb{R}^{100} \quad +$$

$$x_i = \underbrace{(x_i^T w_1) w_1 + (x_i^T w_2) w_2 + (x_i^T w_3) w_3}_{\in \mathbb{R}^{100}} \quad + \leftarrow$$

Rep.

$$\{w_1, w_2, w_3\}$$

Co-BFF

$$x_i \rightarrow \begin{bmatrix} (x_i^T w_1) \\ (x_i^T w_2) \\ (x_i^T w_3) \end{bmatrix} \in \mathbb{R}^3$$

Original: $100 \times n$

Now: $\boxed{100 \times 3}_{w_1, w_2, w_3} + 3 \times n$

$$[a, b, c] \rightarrow [x^T w_1 \quad x^T w_2]$$

$$x = \begin{bmatrix} a \\ b \\ c \end{bmatrix}$$

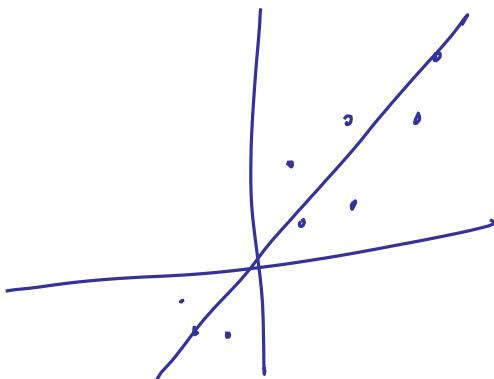
$$w_1 = \begin{bmatrix} w_{11} \\ w_{12} \\ w_{13} \end{bmatrix} \quad w_2 = \begin{bmatrix} w_{21} \\ w_{22} \\ w_{23} \end{bmatrix}$$

$$a \cdot w_{11} + b w_{12} + c w_{13}$$

$$a \cdot w_{21} + b w_{22} + c w_{23}$$

QUESTION

what if data approximately
low dimensional



Summary So Far

Algorithm:

Given data $\{x_1, x_2, \dots, x_n\}$

- > Center data so that mean is 0 (origin shift)
- > Use the “magic” box to find the best line

represented by w_1

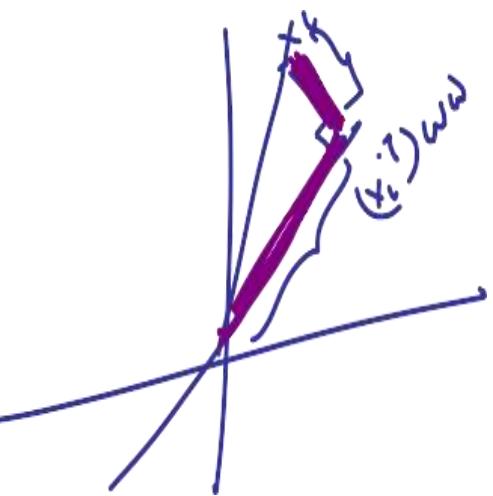
- > Create dataset with residues
- > Iterate the process to find w_2, w_3, \dots

Observations:

- All the lines obtained from the above process are orthogonal to each other
- The residues must necessarily become 0 after ‘d’ rounds
- We gain (compression wise) if residues becomes 0 much earlier than d rounds.

Question:

What if data lives in an approximately low dimensional space?



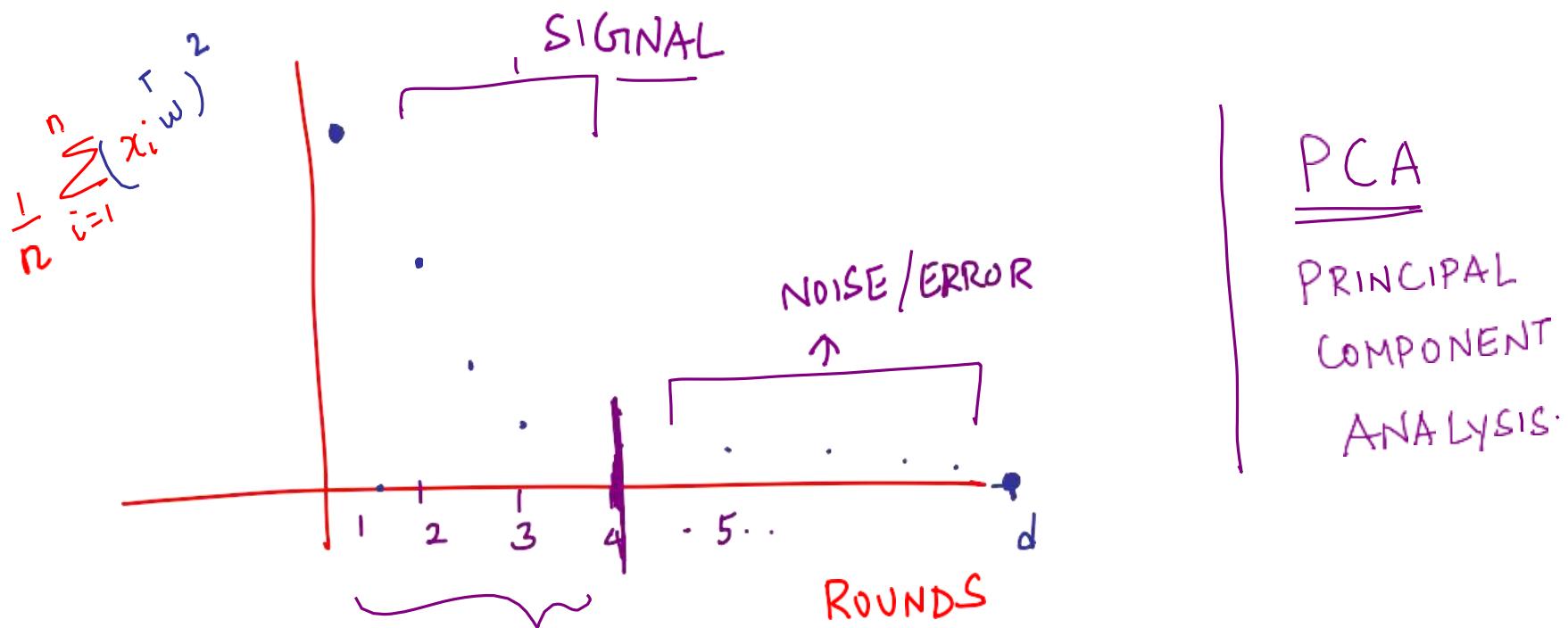
PYTHAGORES

$$+\lambda \|x_i\|^2 = \underbrace{\|x_i - (x_i^T w)w\|_2^2}_{+} + \underbrace{\|(x_i^T w)w\|^2}_{\sim}$$

$$\frac{1}{n} \sum_{i=1}^n \|x_i\|^2 = \frac{1}{n} \sum_{i=1}^n \|x_i - (x_i^T w)w\|_2^2 + \frac{1}{n} \sum_{i=1}^n (x_i^T w)^2$$

Smaller the value
better the fit

larger the value
better the fit



Question

$$\max_{\omega \in \mathbb{R}^d, \|\omega\|^2=1} \omega^T C \omega$$

$$C = \frac{1}{n} \sum_{i=1}^n x_i x_i^T$$

$x_i \in \mathbb{R}^d$
 $x_i^T \in \mathbb{R}^d$
 $C \in \mathbb{R}^{d \times d}$

→ CO-VARIANCE MATRIX

SOLUTION

ω_1 is eigenvector corresponding to
the "largest" eigenvalue of C

[HILBERT'S MIN-MAX
THEOREM]

In fact the eigenvectors $\{w_1, \dots, w_d\}$ of C form
 $\lambda_1, \dots, \lambda_d$
an orthonormal basis

$w_k \rightarrow$ Best line that we get in round k .

By definition

$$C w_i = \lambda_i w_i$$

$$w_i^T (C w_i) = w_i^T (\lambda_i w_i)$$

$$w_i^T C w_i = \lambda_i$$

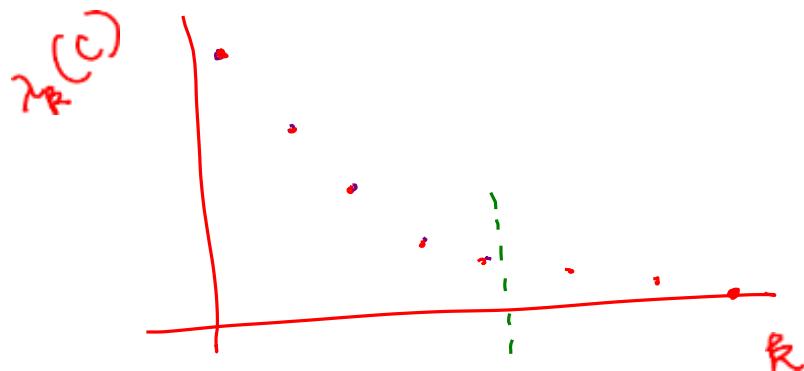
$$\omega_1^T \left(\frac{1}{n} \sum_{i=1}^n x_i x_i^T \right) \omega_1 = \lambda_1$$

$$\frac{1}{n} \sum_{i=1}^n (\omega_1^T x_i) (x_i^T \omega_1) = \lambda_1$$

\Rightarrow

$$\boxed{\lambda_1 = \frac{1}{n} \sum_{i=1}^n (\omega_1^T x_i)^2}$$

which is the term we used earlier



$$\left(\frac{\sum_{i=1}^k x_i}{d} \right) \geq 0.95$$

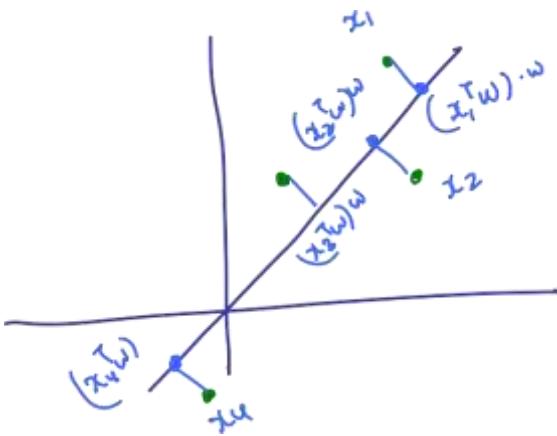
PCA

$$\underbrace{\{w_1, \dots, w_k\}}$$

PRINCIPAL
COMPONENTS

$$x_i \rightarrow \begin{bmatrix} x_i^T w_1 & \dots & x_i^T w_k \end{bmatrix} \in \mathbb{R}^d$$

A DIFFERENT VIEW OF PCA



$$x_i \rightarrow (x_i^T w) w$$

$$\left\{ (x_1^T w), \dots, (x_n^T w) \right\}$$

Average.

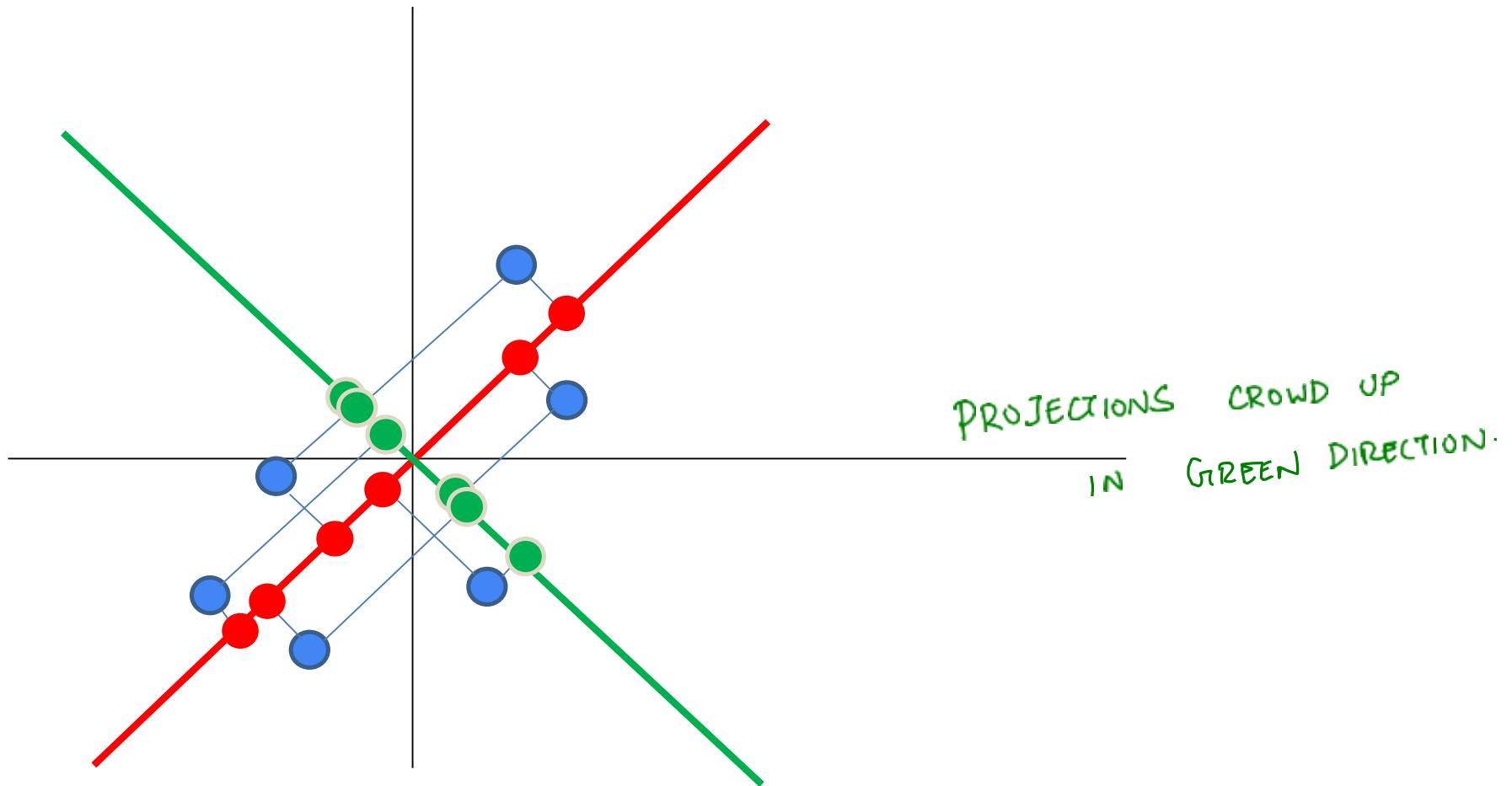
$$\frac{1}{n} \sum_{i=1}^n (x_i^\top w) = \left(\frac{1}{n} \sum_{i=1}^n x_i \right)^\top w$$

$\underbrace{\phantom{\sum_{i=1}^n}}_0^\top w = 0$

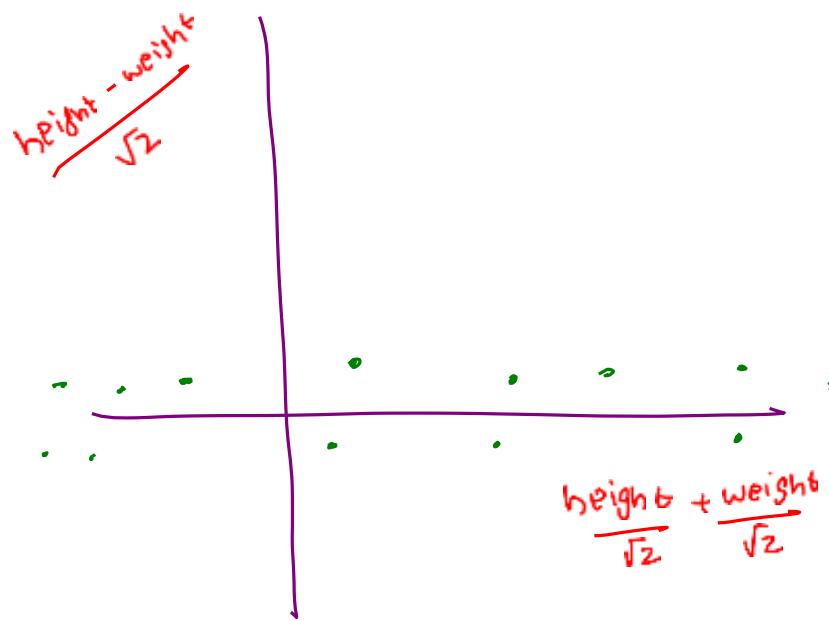
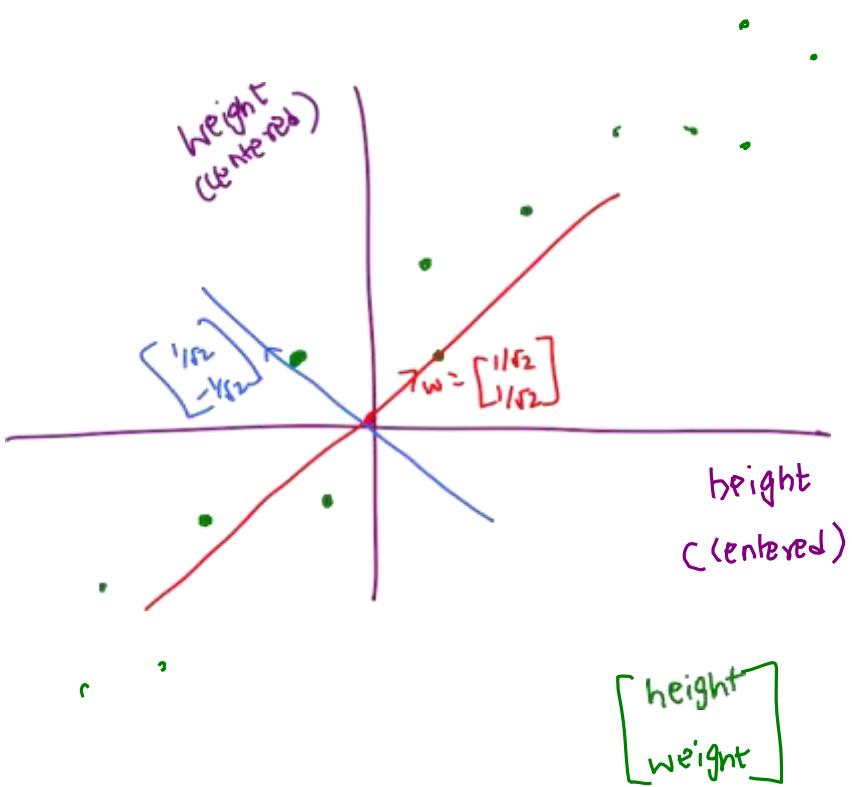
Variance

$$\frac{1}{n} \sum_{i=1}^n (x_i^\top w - 0)^2 = \boxed{\frac{1}{n} \sum_{i=1}^n (x_i^\top w)^2}$$





ONE MORE EXAMPLE

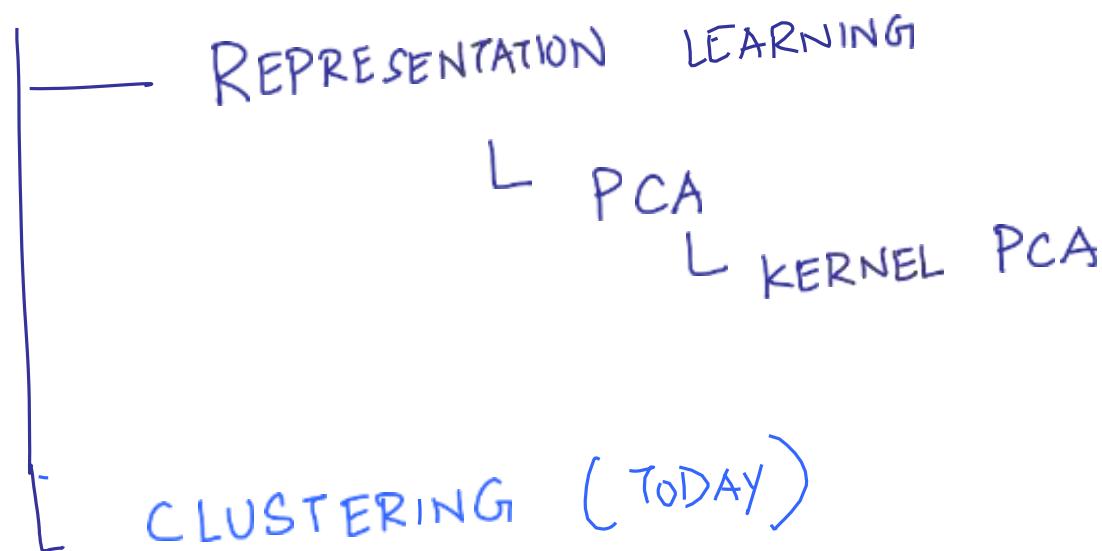


PCA finds combinations of features
that are "decorrelated"

So far

unsupervised

Learning



REAL WORLD EXAMPLE

faculty.washington.edu/sbrunton/me565/pdf/L29secure.pdf

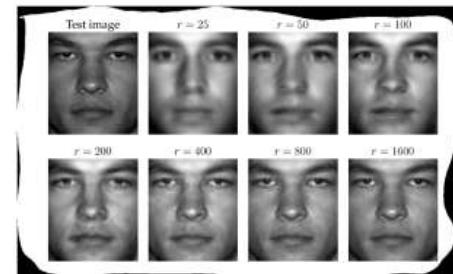
EXTENDED YALE FACE DATASET



EACH Datapoint $x_i \in \mathbb{R}^d$

$$d = \sim 32000 ; n = \sim 20000$$

PRINCIPLE COMPONENTS



TEST IMAGE 1



TEST IMAGE 2

Summary of PCA

- > Given a dataset, we wanted to find a compressed representation.
- > In other words, we wanted to find a subspace that best “represents” the data
- > Best representation => Error is minimum \Leftrightarrow variance is maximum
- > The principal directions can be chosen greedily ! (In practice using SVD)
- > Extremely useful for real world data such as images.

PCA

Time complexity \rightarrow Find Eigen vectors & Eigen values

$$C \in \mathbb{R}^{d \times d}$$

Typically $O(d^3)$

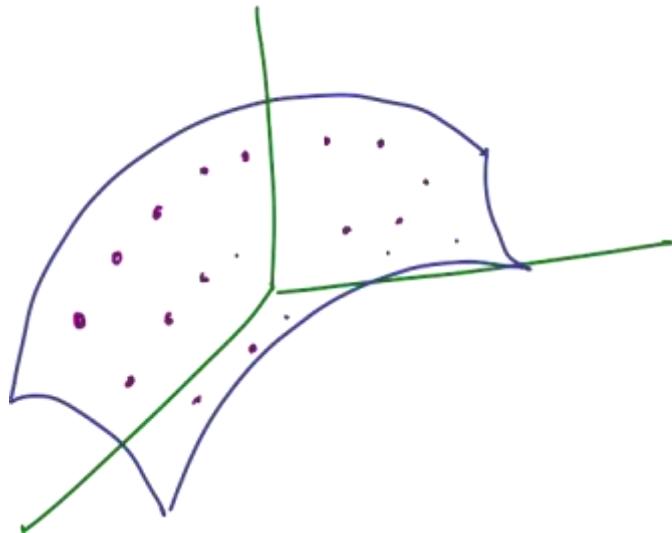
Issue 1: $d \gg n$

Features \downarrow number of data points \downarrow

Example: Face recognition

ISSUE 2

PCA finds "LINEAR" relationships
only



SURPRISING . . . SAME SOLUTION FOR
RESULT BOTH ISSUES !!

ISSUE 1: $d \gg n$ $O(d^3)$ for Eigen decomposition

$$X = \begin{bmatrix} | & | & | \\ x_1 & x_2 & \dots & x_n \\ | & | & & | \end{bmatrix} \in \mathbb{R}^{d \times n}$$

$$C = \frac{1}{n} \sum_{i=1}^n x_i x_i^T$$

$$XX^T = \begin{bmatrix} | & \dots & | \\ x_1 & \dots & x_n \\ | & & | \end{bmatrix}_{d \times n} \begin{bmatrix} -x_1 - \\ -x_2 - \\ \vdots \\ -x_n - \end{bmatrix}_{n \times d} = \sum_{i=1}^n x_i x_i^T \quad [\text{Show this}]$$

\Rightarrow

$$C = \frac{1}{n} \underline{X} \underline{X}^T$$

Let w_R be the eigenvector corresponding to
 R^{th} largest E.V of C (λ_R)

$$C \cdot w_R = \lambda_R w_R$$

$$\left(\frac{1}{n} \sum_{i=1}^n \underline{x}_i \underline{x}_i^T \right) w_R = \lambda_R w_R$$

$$\Rightarrow w_R = \frac{1}{n \lambda_R} \sum_{i=1}^n (\underline{x}_i^T w_R) \underline{x}_i$$

w_R is a
LINEAR
COMBINATION
OF
data points.

$$\begin{bmatrix} 1 & & & 1 \\ x_1 & \cdots & x_n \\ 1 & & & 1 \end{bmatrix} \begin{bmatrix} \alpha_{R1} \\ \vdots \\ \alpha_{Rn} \end{bmatrix} = \sum_{i=1}^n \alpha_{Ri} \cdot x_i$$

$$w_R = \sum_{i=1}^n \underbrace{\left(\frac{x_i^T w_R}{\|x_i\|_2} \right)}_{\text{approximate}} x_i$$

Assume for the moment we know $\alpha_R + R$
 [and not w_R].

Why is this useful?

$$x_{\text{test}} \in \mathbb{R}^d \rightarrow \begin{bmatrix} x_{\text{test}}^T w_1, \dots, x_{\text{test}}^T w_R \end{bmatrix}$$

$$x_{\text{test}}^T w_R = x_{\text{test}}^T \left(\sum_{i=1}^n \alpha_{Ri} x_i \right) = \sum_{i=1}^n \alpha_{Ri} \boxed{x_{\text{test}}^T x_i}$$

Can obtain compressed representation only by knowing alphas

Question: How to get α_B without computing w_B ?

Some algebra

$$\underline{w_B = X\alpha_B} \xrightarrow{\text{# } B} C w_B = \lambda_B w_B \quad [\text{by definition}]$$

$$\left(\frac{1}{n} XX^T\right) w_B = \lambda_B w_B$$

$$\frac{1}{n} XX^T (X\alpha_B) = \lambda_B (X\alpha_B)$$

$$XX^T (x\alpha_k) = n\lambda_k x\alpha_k$$

Premultiply by X^T on both sides

$$\underbrace{X^T (XX^T x\alpha_k)}_{= X^T (n\lambda_k x\alpha_k)}$$

$$(X^T X) (X^T x\alpha_k) = n\lambda_k (X^T x\alpha_k)$$

$$(X^T X)^2 \alpha_k = n\lambda_k (X^T x\alpha_k)$$

Let's call $X^T X$ as K

Want α_R

$$K^2 \alpha_R = \underbrace{n \lambda_R}_{} K \alpha_R.$$

s.t

$$x \alpha_R = \underline{w_R}$$

$$\underline{w_R^T w_R} = 1$$

$$(x \alpha_R)^T (x \alpha_R) = 1$$

$$\alpha_R^T (\underline{x^T x}) \alpha_R = 1$$

$$\underline{\alpha_R^T K \alpha_R} = 1$$

Observation:

Summary so far

PCA

ORIGINAL GOAL :

$$\left\{ \mathbf{x}_1, \dots, \mathbf{x}_n \right\} \xrightarrow{\mathbf{x}_i \in \mathbb{R}^d} \boxed{\mathbf{X} \mathbf{X}^T} \xrightarrow{\frac{nC}{\downarrow}} \left\{ \frac{\mathbf{w}_1}{\lambda_1}, \dots, \frac{\mathbf{w}_k}{\lambda_k} \right\} \quad \frac{n\lambda_1}{\lambda_1} \geq \dots \geq \frac{n\lambda_k}{\lambda_k}$$

OBSERVATION 1 :

$$\mathbf{w}_k = \mathbf{X} \alpha_k \quad \text{for some } \alpha_k \in \mathbb{R}^d$$

\downarrow
 $d \times n$

OBSERVATION 2 :

Suffices to find α_k satisfying

→ P1

$$K \alpha_k = (n \lambda_k) \alpha_k$$

(P2)

$$\alpha_k^T K \alpha_k = 1$$

$$K = \mathbf{X} \mathbf{X}^T \in \mathbb{R}^{n \times n}$$

OBSERVATION 3: $\underset{\text{**}}{nC = \frac{XX^T}{\cancel{n \times n}}} \rightarrow \left\{ \begin{array}{c} w_1 \in \mathbb{R}^d \\ \dots \\ w_l \end{array} \right\}$ $\frac{\|w_k\|^2}{n \lambda_k} = 1$

$$K = \frac{X^T X}{\cancel{n \times n}} \rightarrow \left\{ \begin{array}{c} B_1 \in \mathbb{R}^n \\ \dots \\ B_l \end{array} \right\} \frac{\|B_k\|^2}{n \lambda_k} = 1$$

Thus α_k that satisfies $(P1)$ and $(P2)$ can be obtained as

$$K \rightarrow \boxed{\quad} \rightarrow \left\{ \frac{B_1}{n \lambda_1}, \dots, \frac{B_l}{n \lambda_l} \right\}$$

$$\alpha_k := \frac{B_k}{\sqrt{n \lambda_k}}$$

[Verify $(P1)$ & $(P2)$
using observations
 $(1), (2) \& (3)$]

**** - The non-zero Eigenvalues of XX^T and X^TX are the same – proved using SVD of X**

PROCEDURE

Dataset - $\{x_1, \dots, x_n\}$ $x_i \in \mathbb{R}^d$ $[d \gg n]$

- Step 1: Calculate $K = \mathbf{x}^\top \mathbf{x}$
- Step 2: Eigen decompose K to get

$$\{\beta_1, \dots, \beta_d\}$$

$$\{n\lambda_1 \geq \dots \geq n\lambda_d\}$$

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

$$\|\beta_k\|^2 = 1$$
- Step 3: $\alpha_k = \frac{\beta_k}{\sqrt{n\lambda_k}}$ $\forall k$ $\alpha_k \in \mathbb{R}^n$
- Step 4: $w_k = \mathbf{x}\alpha_k$ $w_k \in \mathbb{R}^d$
 $\|w_k\|^2 = 1 ?$

Key advantages of this method



we need
only $\underline{\underline{X}}^T \underline{\underline{X}}$

i.e., only
"pairwise dot
products"

-

Step 2: Eigen decompose

K

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



to get

$$B_k \in \mathbb{R}^n$$
$$\|B_k\|^2 =$$

$O(n^3)$

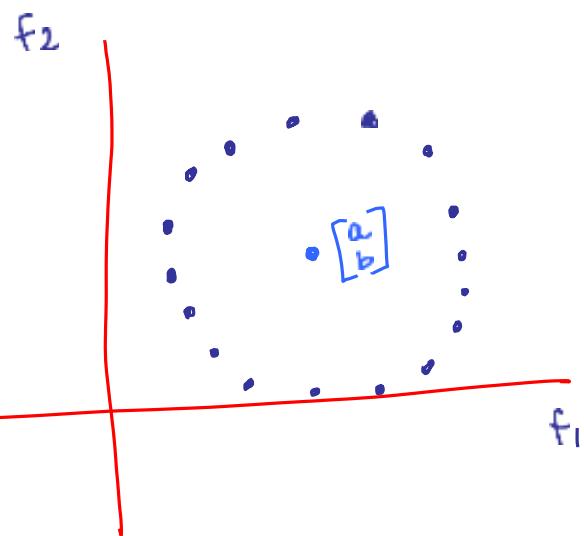


$$\{B_1, \dots, B_e\}$$

$$\{n\lambda_1 \geq \dots \geq n\lambda_e\}$$

ISSUE 2

→ Features could be
non-linearly related



Standard PCA would
pick 2 Eigen directions

$$(f_1 - a)^2 + (f_2 - b)^2 = r^2$$

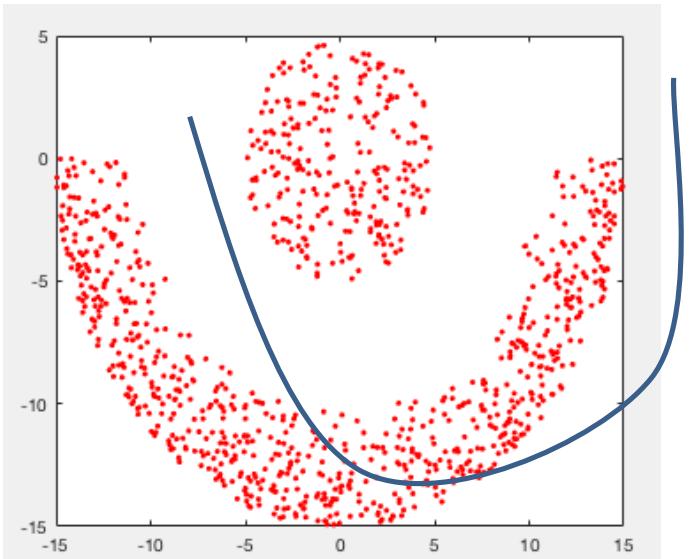
Idea: Transform features from
low dimension \mathbb{R}^d to \mathbb{R}^D
high dimension \mathbb{R}^D

$$x \rightarrow \phi(x)$$
$$\mathbb{R}^d \quad \mathbb{R}^D$$

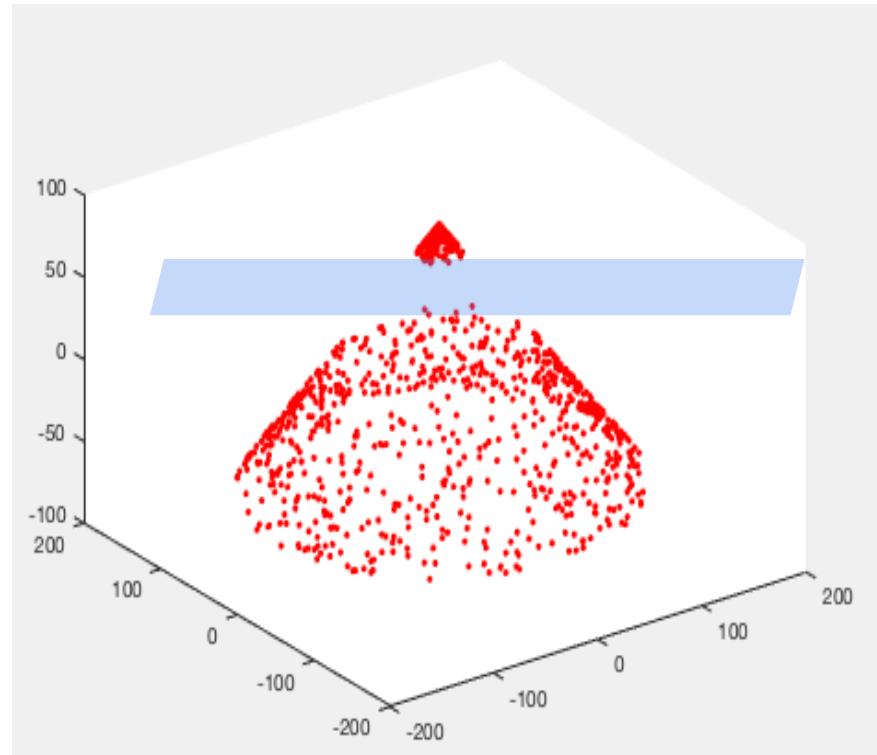
We already know how to handle
Case when $D \gg n$

[use $\phi(x)^\top \phi(x)$
instead of
 $\phi(x) \phi(x)^\top$]

ORIGINAL DATA



3D REPRESENTATION



Quadratic Map φ
to 6 Dimensions

PCA

w_1, w_2, w_3
Top 3 principal
components

Reconstruct
each point x as
 $[\phi(x)'w_1 \ \phi(x)'w_2 \ \phi(x)'w_3]$

KERNEL-PCA

Input: Dataset $\{x_1, \dots, x_n\}$ $x_i \in \mathbb{R}^d$
Kernel $R: \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$

Step 1: Compute $K \in \mathbb{R}^{n \times n}$

$$k_{ij} = R(x_i, x_j)$$

Step 2: Compute $\{\beta_1, \dots, \beta_l\}$ $\Rightarrow \alpha_k := \frac{\beta_k}{\sqrt{n\lambda_k}}$
 $\{\lambda_1, \dots, \lambda_l\}$
e.v of K

\times Step 3.

$$w_k = \boxed{\phi(x) \alpha_k}$$

"Defeats the purpose"

We cannot "reconstruct" eigenvectors of the covariance matrix

$$\phi(x_i)^T w_k = \phi(x_i)^T \left(\sum_{j=1}^n \alpha_{kj} \phi(x_j) \right)$$

$\sum_{j=1}^n \alpha_{kj} k(i,j)$

Typically enough for good downstream tasks.

What about Centering?

- Even if center the original data, it does not guarantee that the data will be centered after applying the transformation ϕ
- We want to work with the centered data w.r.t to $\{\phi(x_1), \phi(x_2), \dots, \phi(x_n)\}$
- Let $\mu = 1/n \sum \phi(x_\ell)$
- The (i,j) term in the pairwise dot product matrix in the high dimension will
- $$(\phi(x_i) - \mu)^T (\phi(x_j) - \mu) = \phi(x_i)^T \phi(x_j) - \phi(x_i)^T \mu - \phi(x_j)^T \mu + \mu^T \mu$$

$$= \phi(x_i)^T \phi(x_j) - \phi(x_i)^T (1/n \sum \phi(x_\ell)) - \phi(x_j)^T (1/n \sum \phi(x_\ell))$$

$$+ (1/n \sum \phi(x_\ell))^T (1/n \sum \phi(x_\ell))$$

$$= K(i,j) - 1/n \sum K(i, \ell) - 1/n \sum K(j, \ell) - 1/n^2 \sum \sum K(\ell, \ell')$$

Good news - Each term depends only on K

KERNEL-PCA

Input: Dataset $\{x_1, \dots, x_n\}$ $x \in \mathbb{R}^d$
 Kernel $R: \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$

Step 1: Compute $K \in \mathbb{R}^{n \times n}$

$$k_{ij} = R(x_i, x_j)$$

Step 1.5 Center Kernel using ①

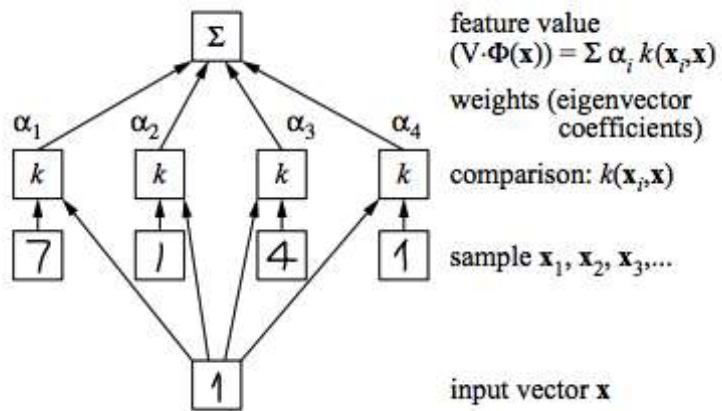
Step 2: Compute $\{\beta_1, \dots, \beta_l\}$ e.vect of K $\Rightarrow \alpha_R := \frac{\beta_k}{\sqrt{n\lambda_R}}$
 $\{\eta\lambda_1, \dots, \eta\lambda_l\}$ e.v of K

$$\phi(x_i)^T w_R$$

$$\phi(x_i)^T \left(\sum_{j=1}^n \alpha_{Rj} \phi(x_j) \right)$$

$$\boxed{\sum_{j=1}^n \alpha_{Rj} k(i,j)}$$

Typically enough good for downstream tasks



Nonlinear Component Analysis as a Kernel Eigenvalue Problem

Bernhard Schölkopf
Max-Planck-Institut für Neurologische Kybernetik, 7207 Tübingen, Germany

Alexander Smola
Klaus-Robert Müller
GMD First (Forschungszentrum Informationstechnik), 12489 Berlin, Germany

Application of Kernel-PCA to face recognition



400 images of 40 subjects with different expression/poses.

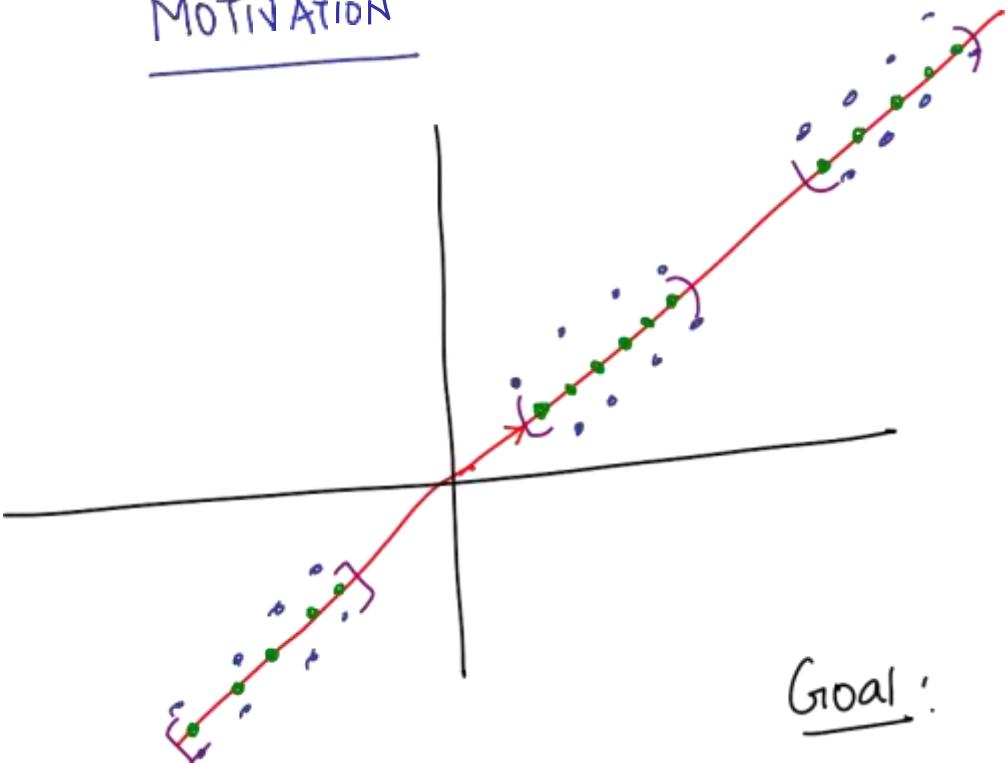
Each image is $(23 \times 28) = 644$ dimensional

- Apply Kernel PCA to obtain features
- Use features as input to a simple supervised learning algorithm
(which we will see later)

Results:

Method	Reduced space	Error Rate (%)
Eigenface	30	2.75
Kernel PCA, d=2	50	2.50
Kernel PCA, d=3	50	2.00
Kernel PCA, d=4	60	2.25
Kernel PCA, d=10	80	2.25

MOTIVATION

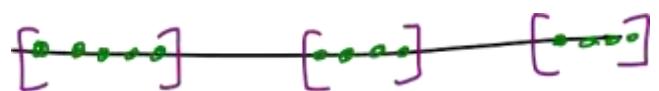


Goal:

Given $\{x_1, x_2, \dots, x_n\}$ $x_i \in \mathbb{R}^d$,

PARTITION the data into

K CLUSTERS

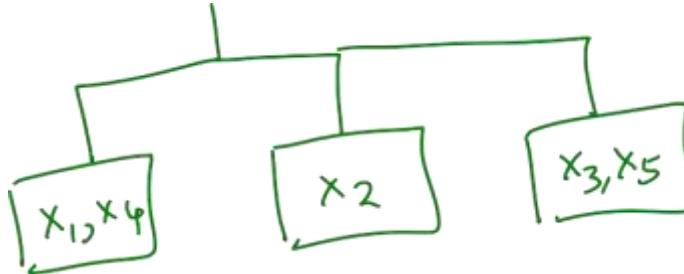


EXAMPLE

Say $K=3$

$$\{x_1, x_2, x_3, x_4, x_5\}$$

one way



another way



3^5 possibilities

$x_1, x_2, \dots, x_n \leftarrow$ DATA POINTS

$z_1, z_2, \dots, z_n \leftarrow$ CLUSTER INDICATOR

$$z_i \in \{1, \dots, k\}^{+i}$$

Given a cluster assignment, how
good is it

$$F(z_1, \dots, z_n) = \sum_{i=1}^n \|x_i - M_{z_i}\|^2$$

$\forall k \quad M_k = \frac{\sum_{i=1}^n x_i \underbrace{1(z_i = k)}_{\text{INDICATOR}}}{\sum_{i=1}^n 1(z_i = k)}$

INDICATOR

$$\begin{cases} 1(p) = 1 & \text{if } p \text{ is true} \\ 0 & \text{if } p \text{ is false} \end{cases}$$

Eg

$$M_1 = \frac{x_1(1) + x_2(0) + x_3(1) + x_4(0) + x_5(0)}{3} \quad k=2$$

x_1	x_2	x_3	x_4	x_5
1	2	1	1	2
z_1	z_2	z_3	z_4	z_5

Goal:

$$\min_{\{z_1, \dots, z_n\}} \sum_{i=1}^n \|x_i - \mu_{z_i}\|^2$$

NP - HARD

LLOYD'S ALGORITHM

?? Initialization

$$z_1^0, z_2^0, \dots, z_n^0 \in \{1, \dots, k\}$$

until Convergence.

μ_k^t Mean of cluster k
at iteration t

$$\mu_k^t = \frac{\sum_{i=1}^n x_i \mathbb{1}(z_i^t = k)}{\sum_{i=1}^n \mathbb{1}(z_i^t = k)}$$

REASSIGNMENT STEP

$$z_i^{t+1} = \arg \min_k \|x_i - \mu_k^t\|^2$$
$$\in \{1, \dots, k\}$$

FACT: Lloyd's Algorithm converges!

QUESTIONS



- CONVERGENCE ?
- INITIALIZATION
- CHOICE OF K.
- NATURE OF CLUSTERS

[sensitive to initialization]



LLOYD'S ALGORITHM CONVERGES

FACT - 1

Let $x_1, \dots, x_l \in \mathbb{R}^d$

$$\underline{v^*} = \underset{v \in \mathbb{R}^d}{\operatorname{argmin}} \sum_{i=1}^l \|x_i - v\|^2$$

$$f(v) = \sum_{i=1}^l \|x_i - v\|_2^2$$

Solution:

$$v^* = \frac{1}{l} \sum_{i=1}^l x_i$$

$$\nabla f(v) = 2 \sum_{i=1}^l (x_i - v)$$

[Exercise]

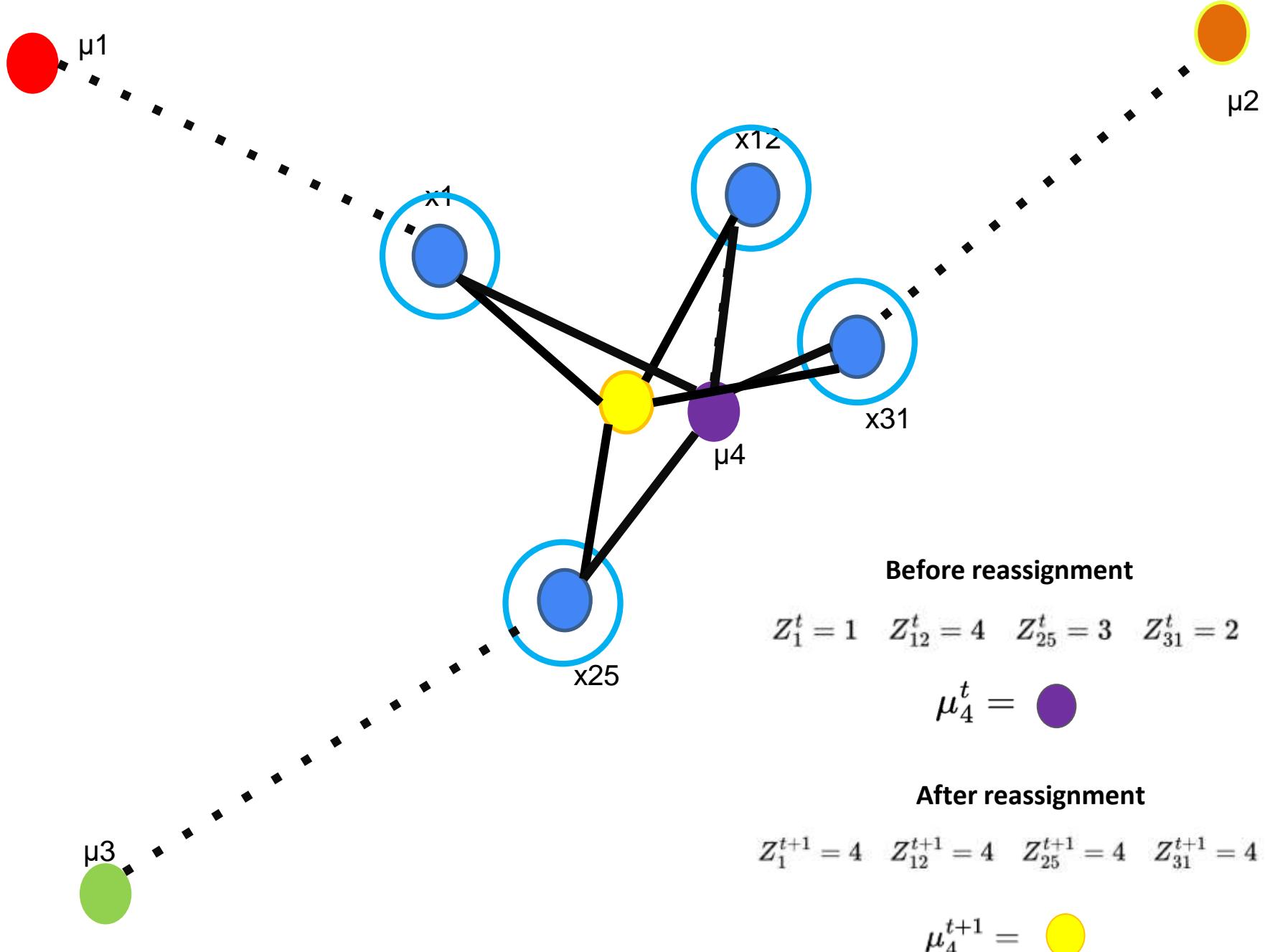
- Say we are at iteration t of Lloyd's

- $z_1^t, z_2^t, \dots, z_n^t \in \{1, \dots, k\}$

CURRENT ASSIGNMENTS

- $z_1^{t+1}, z_2^{t+1}, \dots, z_n^{t+1} \in \{1, \dots, k\}$

What can we say about the update?



By definition

$$\sum_{i=1}^n \|x_i - \mu_{z_i^{t+1}}^t\|^2 <$$

MEAN OF
Cluster
where
 x_i wants to
be

$$\sum_{i=1}^n \|x_i - \mu_{z_i^t}^t\|^2$$

Mean of
Current
Cluster
 x_i is
assigned to

$$\sum_{i=1}^n \|x_i - \mu_{z_i^{t+1}}^t\|^2 <$$

$$\sum_{i=1}^n \|x_i - \mu_{z_i^{t+1}}^t\|^2$$

using
(FACT 1)

If Algo does
not converge in
round t

\Rightarrow

$$\sum_{i=1}^n \|x_i - M_{z_i^{t+1}}\| < \sum_{i=1}^n \|x_i - M_{z_i^t}\|^2$$
$$F(z_1^{t+1}, \dots, z_n^{t+1}) < F(z_1^t, \dots, z_n^t)$$

\Rightarrow convergence?

- OBJECTIVE FUNCTION STRICTLY REDUCES IF
RE-ASSIGNMENT HAPPENS
- only "FINITE" possibilities (Partitions)
- Algorithm must converge.

NATURE OF CLUSTERS

Say $k=2$

- Lloyd's produces 2 clusters with means M_1, M_2
- What can we say about points assigned to Cluster 1?

Where are

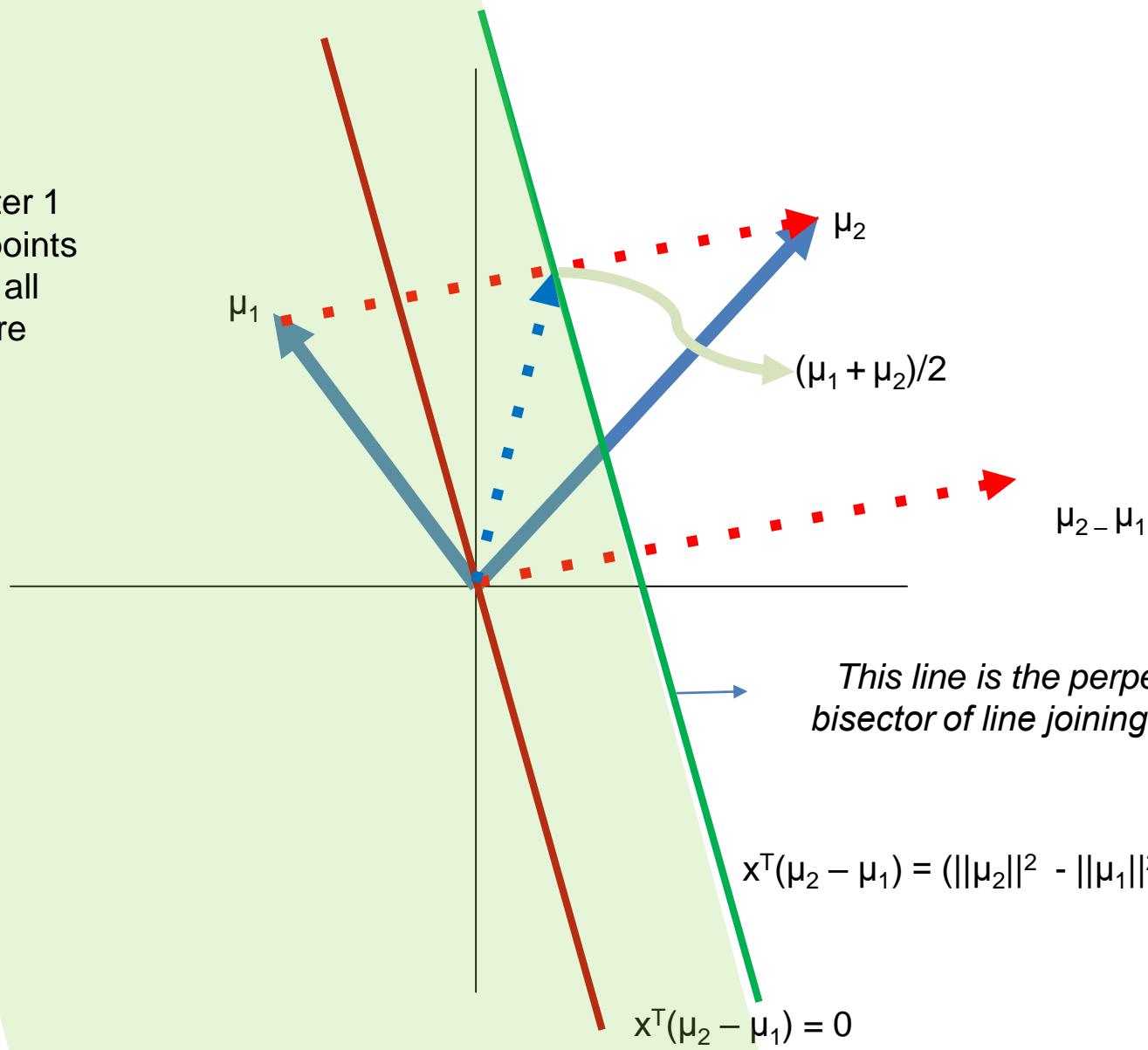
clusters?

$$\|x - \mu_1\|^2 \leq \|x - \mu_2\|^2$$

$$\begin{aligned} \|x\|^2 + \|\mu_1\|^2 - 2x^\top \mu_1 &\leq \|x\|^2 + \|\mu_2\|^2 \\ -2x^\top \mu_2 \\ x^\top (\mu_2 - \mu_1) &\leq (\|\mu_2\|^2 - \|\mu_1\|^2)/2 \end{aligned}$$

$$x^\top (\mu_2 - \mu_1) \leq (\|\mu_2\|^2 - \|\mu_1\|^2)/2$$

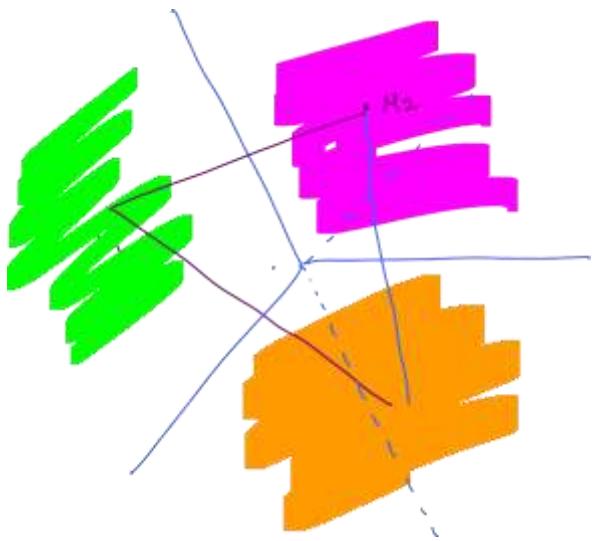
Cluster 1
data points
are all
here



This line is the perpendicular bisector of line joining μ_1 and μ_2

$$x^\top(\mu_2 - \mu_1) = (\|\mu_2\|^2 - \|\mu_1\|^2)/2$$

$$x^\top(\mu_2 - \mu_1) = 0$$



N / CELLS

So far

- CONVERGENCE ?
- INITIALIZATION
- CHOICE OF K.
- NATURE OF CLUSTERS

INITIALIZATION

→ RANDOM PARTITION

→ uniformly Sample k-means from dataset

K-MEANS ++

→ Choose first mean M_1° uniformly
at random from $\{x_1, \dots, x_n\}$

→ for $k = 2, \dots, K$
choose M_k° "probabilistically"
according to score $S(x) = \min_{j=1, \dots, k-1} \|x - M_j^{\circ}\|_2^2$

GUARANTEE

$$\mathbb{E} \left[\sum_{i=1}^n \|x_i - \mu_{z_i}\|^2 \right] \leq O(\log K) \left(\min_{z_1, \dots, z_n} \sum_{i=1}^n \|x_i - \mu_{z_i}\|^2 \right)$$

Over randomness
over algorithm's choices

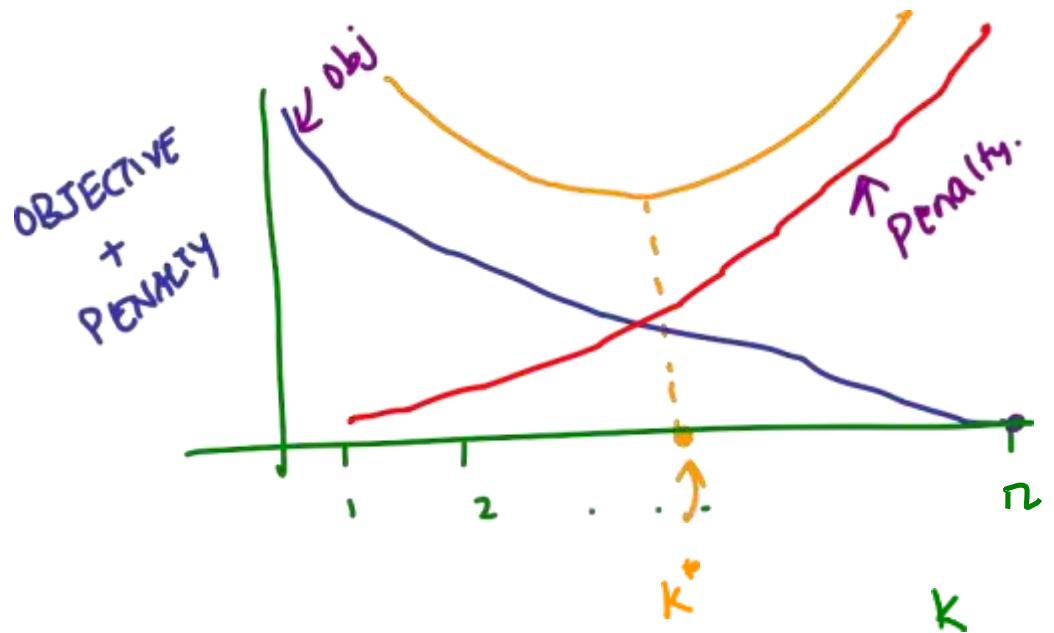
CHOICE OF K ?

$$\underbrace{F(z_1, \dots, z_n)}$$

Penalize large k.

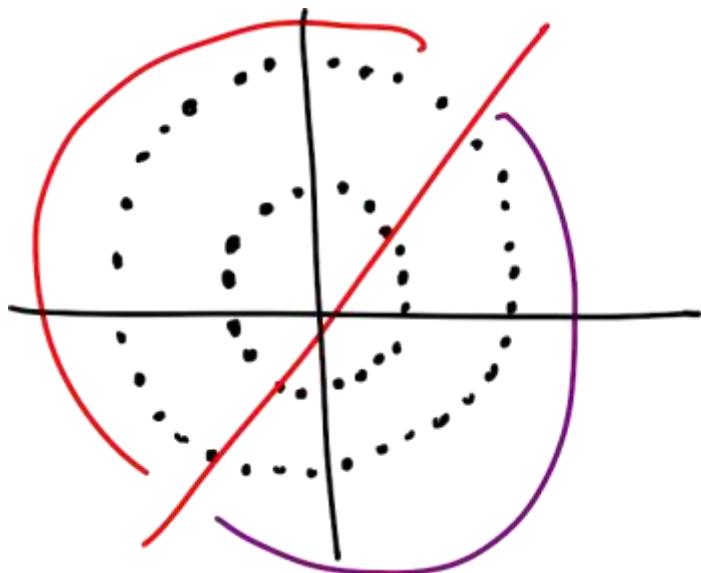
Find the k that has the

Smallest objective value + Penalty(k)



$AIC = 2K - 2 \underbrace{\log(L(\theta^*))}_{\text{Likelihood}}$
 - Akaike Information Criterion

$BIC = \frac{K \log(n)}{n} - \underbrace{2 \log(L(\theta^*))}_{\text{Likelihood}}$
 - Bayesian Information Criterion.



KERNELIZE

K-MEANS

ORIGINAL

$$\min_{z_1, \dots, z_n} F(z_1, \dots, z_n) = \sum_{i=1}^n \|x_i - \mu_{z_i}\|^2$$

NEW

$$\min_{z_1, \dots, z_n} \sum_{i=1}^n \|\phi(x_i) - \mu_{z_i}\|_2^2$$

NOT "KERNELIZED" YET!

$\neq k$

$\mu_k =$

$$\sum_{i=1}^n \phi(x_i) \mathbb{1}(z_i=k) \quad / \quad \sum_{i=1}^n \mathbb{1}(z_i=k)$$

G10 : Can we rewrite the new objective
using dot-products.

Consider

$$\begin{bmatrix} \phi(x_1) & \phi(x_2) & \dots & \phi(x_n) \\ | & | & & | \end{bmatrix} - \begin{bmatrix} 1 & 1 & \dots & 1 \\ M_3 & M_5 & \dots & M_1 M_3 \\ | & | & & | \end{bmatrix} M$$

Φ

$$M = \begin{bmatrix} 1 & | & M_1 & | & M_3 \\ | & M_5 & . & | & | \\ | & | & | & | & | \\ | & | & | & | & | \end{bmatrix}_{D \times N} = \begin{bmatrix} 1 \\ \frac{\phi(x_1) + \phi(x_3) + \phi(x_5)}{3} \\ | \\ = \end{bmatrix}_{\frac{1}{3} \times N} \rightarrow \begin{bmatrix} 1 \\ \frac{\phi(x_2) + \phi(x_4)}{2} \\ | \\ ? \end{bmatrix}_{\frac{1}{3} \times N}$$

$$= \begin{bmatrix} 1 & | & \phi(x_1) & | & \phi(x_2) & | & \phi(x_N) \\ | & | & | & | & | & | & | \end{bmatrix}_{D \times N} \cdot \begin{bmatrix} \frac{1}{3} & 0 & \frac{1}{3} \\ 0 & \frac{1}{2} & 0 \\ \frac{1}{3} & 0 & \frac{1}{3} \\ 0 & \frac{1}{2} & 0 \\ \frac{1}{3} & 0 & \frac{1}{3} \\ \vdots & \vdots & \vdots \\ 0 & 0 & 0 \end{bmatrix}_{N \times N}$$

$\overbrace{\quad}^{\Phi}$

$$\begin{bmatrix}
 l_3 & 0 & & \\
 0 & l_2 & & \\
 l_3 & 0 & \ddots & \\
 0 & l_2 & & \\
 l_3 & 0 & & \\
 \vdots & \vdots & & \\
 \end{bmatrix}_{N \times N} =
 \begin{bmatrix}
 0 & 0 & 0 & \cdots & 0 \\
 0 & 1 & 0 & \cdots & 0 \\
 0 & 0 & 1 & \cdots & 0 \\
 \vdots & & & & \\
 \end{bmatrix}_{N \times K} \begin{bmatrix}
 l_3 & & & & 0 \\
 l_2 & l_3 & & & \\
 0 & 0 & l_3 & & \\
 l_2 & l_2 & l_3 & l_3 & \\
 l_3 & & l_3 & & \\
 \vdots & & & & \\
 \end{bmatrix}_{K \times K} \begin{bmatrix}
 0 & 0 & 0 & \\
 0 & -1 & 0 & \\
 1 & 0 & -1 & \\
 \vdots & \vdots & \vdots & \\
 0 & 0 & 0 & \\
 \end{bmatrix}_{K \times N}$$

↗
 ASSIGNMENT MATRIX
 $Z \in \{0, 1\}^{N \times K}$
 ↘ COUNT MATRIX
 $L \in \mathbb{R}^{K \times K}$
 ↘
 Z^T

$$\bar{\Phi} - M = \bar{\Phi} - \bar{\Phi} Z L Z^T$$

$\downarrow_{D \times N}$

Want to
consider

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

$$= \text{trace} \left[(\phi - M)^T (\phi - M) \right]$$

$$= \text{trace} \left[(\phi - \phi Z L Z^T)^T (\phi - \phi Z L Z^T) \right]$$

$$= \text{trace} \left[\underbrace{\phi^T \phi}_{A} - \underbrace{\phi^T \phi Z L Z^T}_{B} - \underbrace{Z L Z^T \phi^T \phi}_{C} + \underbrace{Z L Z^T \phi^T \phi}_{A} \underbrace{Z L Z^T}_{B} \underbrace{\phi^T \phi}_{C} \right]$$

$$A = \begin{bmatrix} | & | & \dots & | \\ a_1 & a_2 & \dots & a_n \\ | & | & \dots & | \end{bmatrix}$$

$$A^T A = \begin{bmatrix} -a_1 & - \\ -a_2 & - \\ \vdots & \\ -a_n & - \end{bmatrix} \begin{bmatrix} | & \dots & | \\ a_1 & \dots & a_n \\ | & \dots & | \end{bmatrix}$$

$$= \begin{bmatrix} \|a_1\|^2 & & & \\ & \|a_2\|^2 & & \\ & & \ddots & \\ & & & \|a_n\|^2 \end{bmatrix}$$

$$\text{trace}(A^T A) = \sum_{i=1}^n \|a_i\|_2^2$$

Facts about trace

$$\cdot \text{trace}(A+B) = \text{trace}(A) + \text{trace}(B)$$

$$\cdot \text{tr}(c \cdot A) = c \cdot \text{tr}(A).$$

$$\left(|| - z|| - |\phi| \right) \sim_L |\phi^\top \phi|$$

$$(|\phi| - |\phi^\top z|)$$

Facts about trace

- $\text{trace}(A+B) = \text{trace}(A) + \text{trace}(B)$
- $\text{trace}(A \cdot B \cdot C) = \text{trace}(C \cdot A \cdot B) = \text{trace}(B \cdot C \cdot A)$.

$$\text{trace} \left(\phi^\top \phi - Z L Z^\top \cancel{\phi^\top \phi} - \phi^\top \phi Z L Z^\top + \underbrace{Z L \cancel{Z^\top} Z L^\top}_{L^{-1}} \phi^\top \phi \right) \quad (\text{verify})$$

~~$Z L Z^\top \cancel{\phi^\top \phi}$~~

$$\text{trace} \left(\phi^\top \phi - \phi^\top \phi Z L Z^\top \right)$$

Goal:

$$\min \text{trace} \left(\underline{\phi^T \phi} - \underline{\phi^T \phi Z L Z^T} \right)$$

Z s.t

it is a
valid
assignment
matrix

\equiv

$$\max_{\substack{Z \text{ s.t} \\ Z \text{ is} \\ \text{valid} \\ \text{assng} \\ \text{matrix}}} \text{trace} \left(\underline{\phi^T \phi Z L Z^T} \right)$$

Kernel.

$$\text{trace} \left(K Z L Z^T \right)$$

$$= \max_Z \text{trace} \left(\underbrace{K}_{A} \underbrace{Z L^{\frac{1}{2}}}_{B} \underbrace{L^{\frac{1}{2}} Z^T}_{C} \right)$$

$$= \max_Z \text{trace} \left(\boxed{(Z L^{\frac{1}{2}})^T} K (Z L^{\frac{1}{2}}) \right)$$

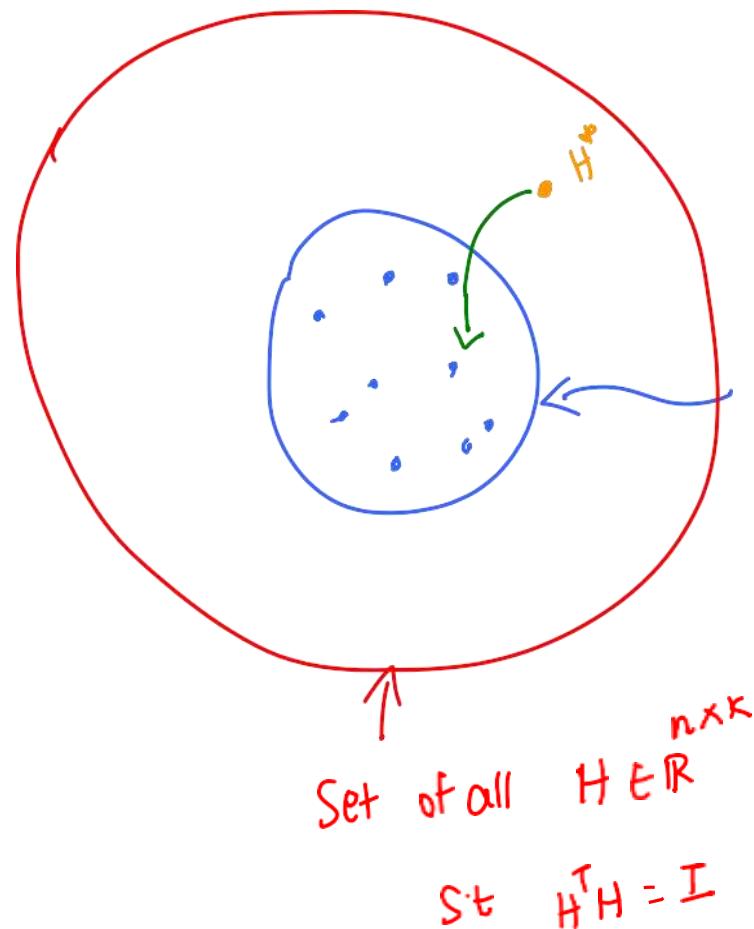
Z
 is a
 valid
 assignment
 matrix.

NP-HARD

Notice that

$$\begin{aligned}
 (Z L^{\frac{1}{2}})^T (Z L^{\frac{1}{2}}) &= L^{\frac{1}{2}} Z^T Z L^{\frac{1}{2}} \\
 &= L^{\frac{1}{2}} L^{-1} L^{\frac{1}{2}} = I
 \end{aligned}$$

Can we solve a relaxed version?



valid assignments
 $Z \in \{0,1\}^{n \times k}$ s.t.
 $Z L^2$ satisfies
 $(Z L^2)^T Z L^2 = I$

RELAXATION

$$\max_{\substack{H \in \mathbb{R} \\ n \times K}}$$

$$\text{trace}(H^T K H)$$

$$\text{st } H^T H = I$$

Solution: Top $\underset{E \in \mathbb{N}}{K}$ eigen vectors of K

kernel.



SPECTRAL CLUSTERING

→ Given $X = \{x_1, \dots, x_n\} \in \mathbb{R}^d$, kernel K

→ Compute $H^* = \begin{bmatrix} | & & | \\ x_1 & \cdot & x_k \\ | & & | \end{bmatrix} \in \mathbb{R}^{n \times k}$
 ↳ top k eigen vectors of K

→ Normalize rows of H^*

→ Run Lloyd's assuming each row is "new" data.

Why Normalize?

Observation:

- > H is a proxy for $ZL^{1/2}$
- > We need a way to go from H to Z
- > if we were lucky, $H = ZL^{1/2}$ for some Z .
- > How can we go from H to Z in that case?

Example:

$Z =$

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

$L =$

$$\begin{pmatrix} 0.5000 & 0 & 0 & 0 \\ 0 & 0.5000 & 0 & 0 \\ 0 & 0 & 0.5000 & 0 \\ 0 & 0 & 0 & 1.0000 \end{pmatrix}$$

$ZL^{1/2} =$

$$\begin{pmatrix} 0 & 0.7071 & 0 & 0 \\ 0.7071 & 0 & 0 & 0 \\ 0 & 0.7071 & 0 & 0 \\ 0 & 0 & 0.7071 & 0 \\ 0 & 0 & 0 & 1.0000 \\ 0.7071 & 0 & 0 & 0 \\ 0 & 0 & 0.7071 & 0 \end{pmatrix}$$

After row normalization

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

- In
ther
- H
- T
- N

150

100

50

0

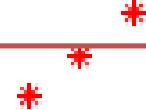
-50

-100

Kernel K = X^TX

20000	22000	24000
22000	24200	26400
24000	26400	28800
-20000	-22000	-24000
-22000	-24200	-26400
-24000	-26400	-28800

-20000	-22000	-24000
-22000	-24200	-26400
-24000	-26400	-28800
20000	22000	24000
22000	24200	26400
24000	26400	28800

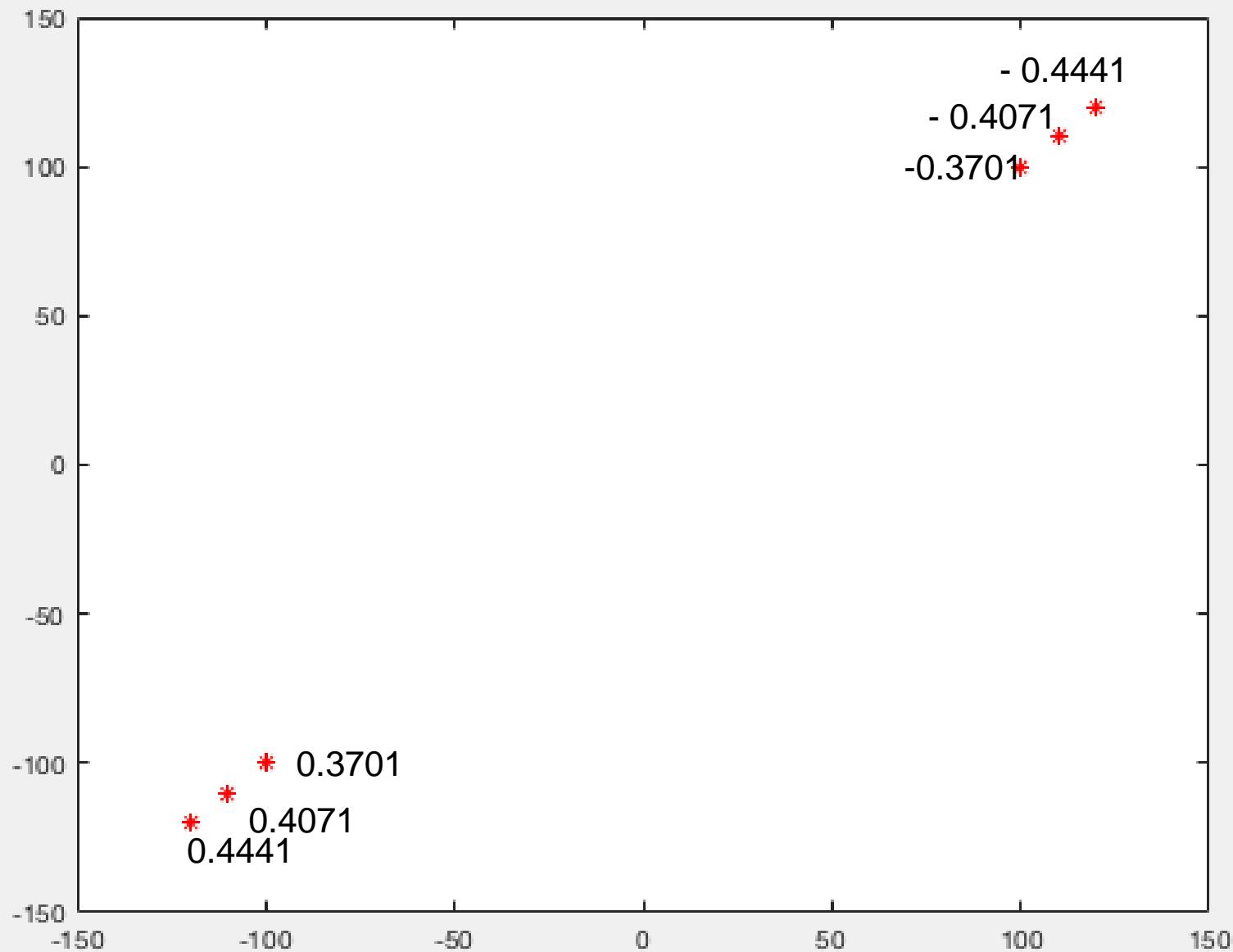


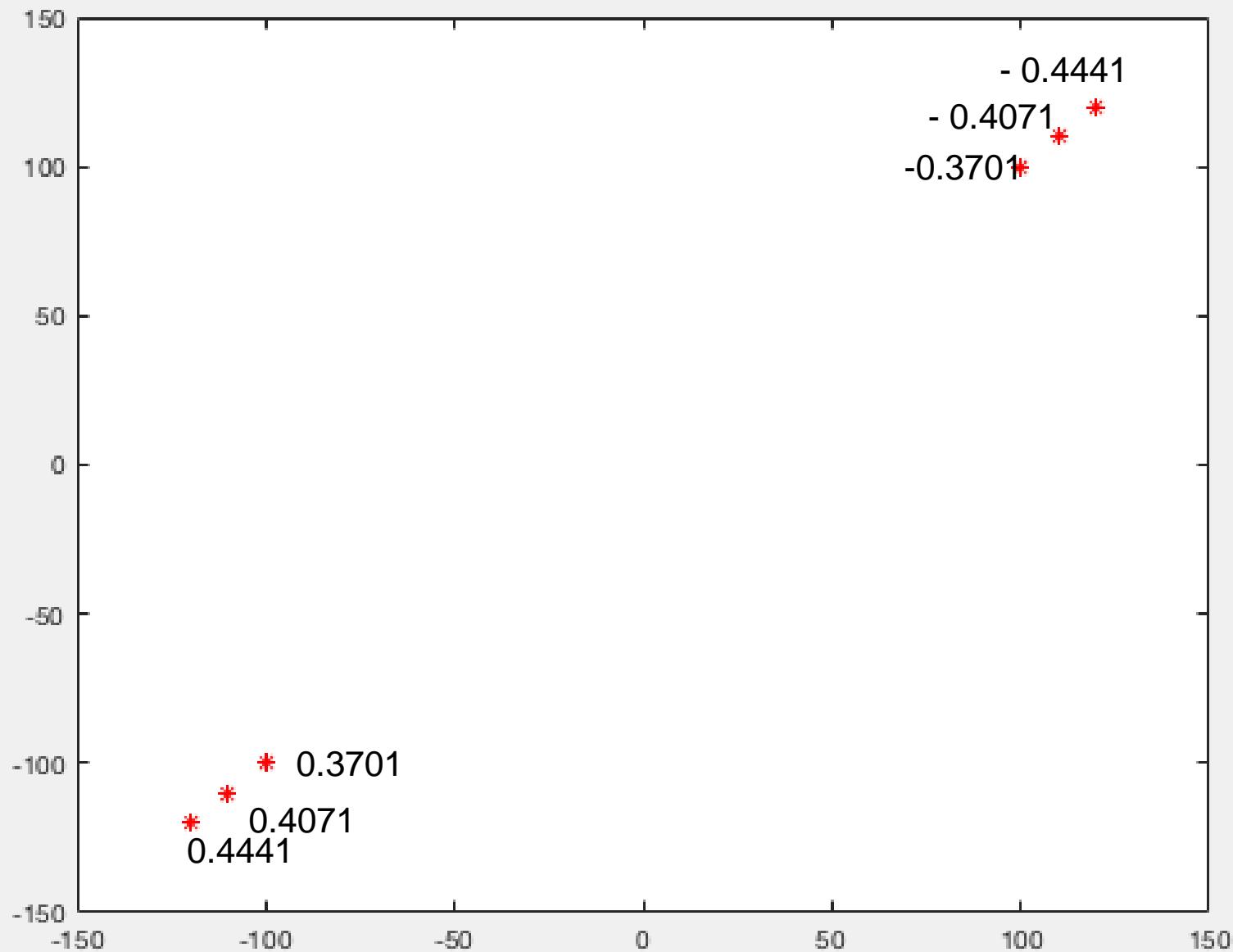
Observations:

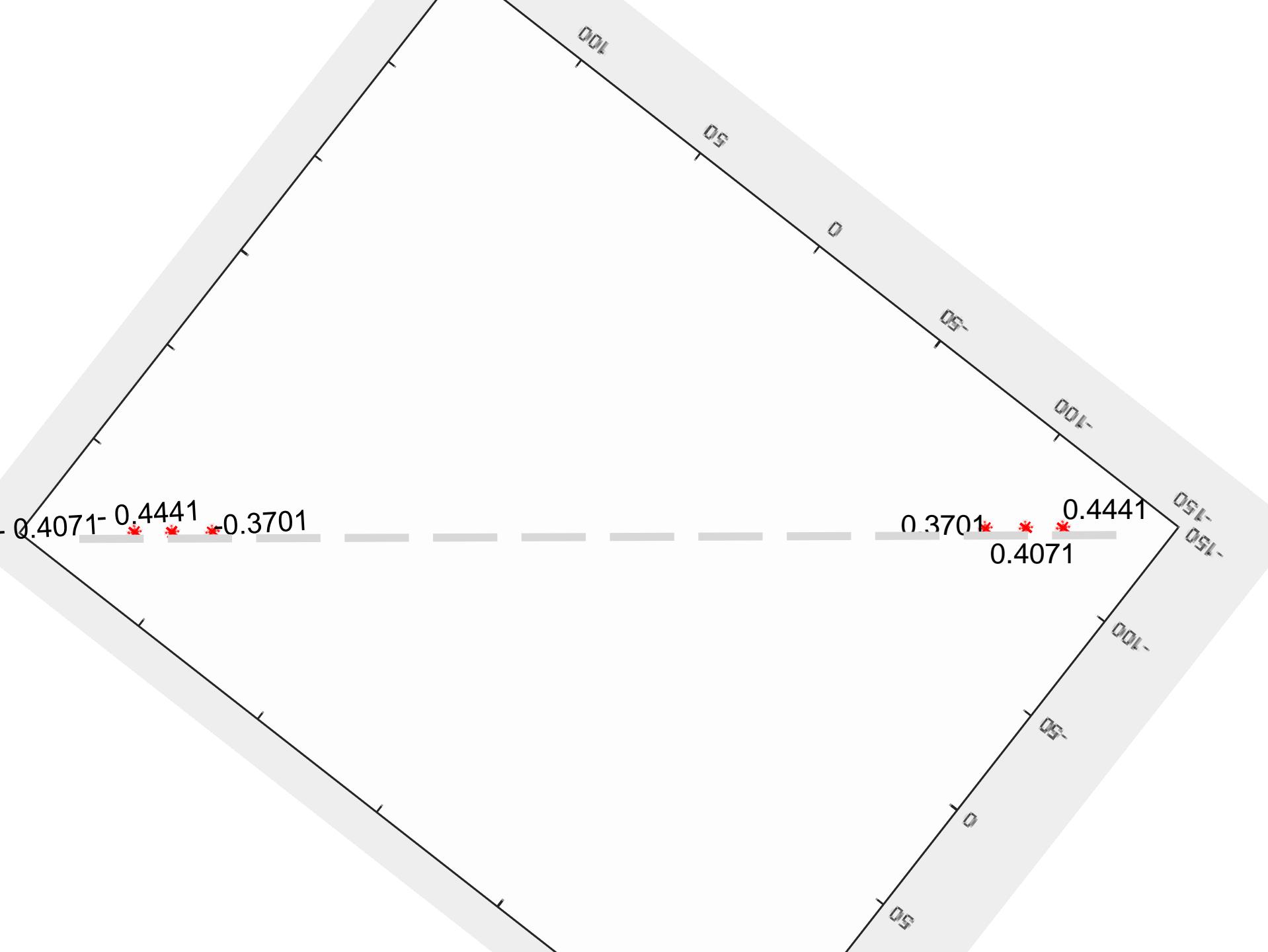
- K is a rank 1 matrix; K is real and symmetric
 $\Rightarrow K = \lambda(\alpha\alpha^T)$ for positive real number λ and vector α
- $\Rightarrow \lambda$ is the Eigen value and α is the corresponding Eigen vector

Thus, for all (i,j)

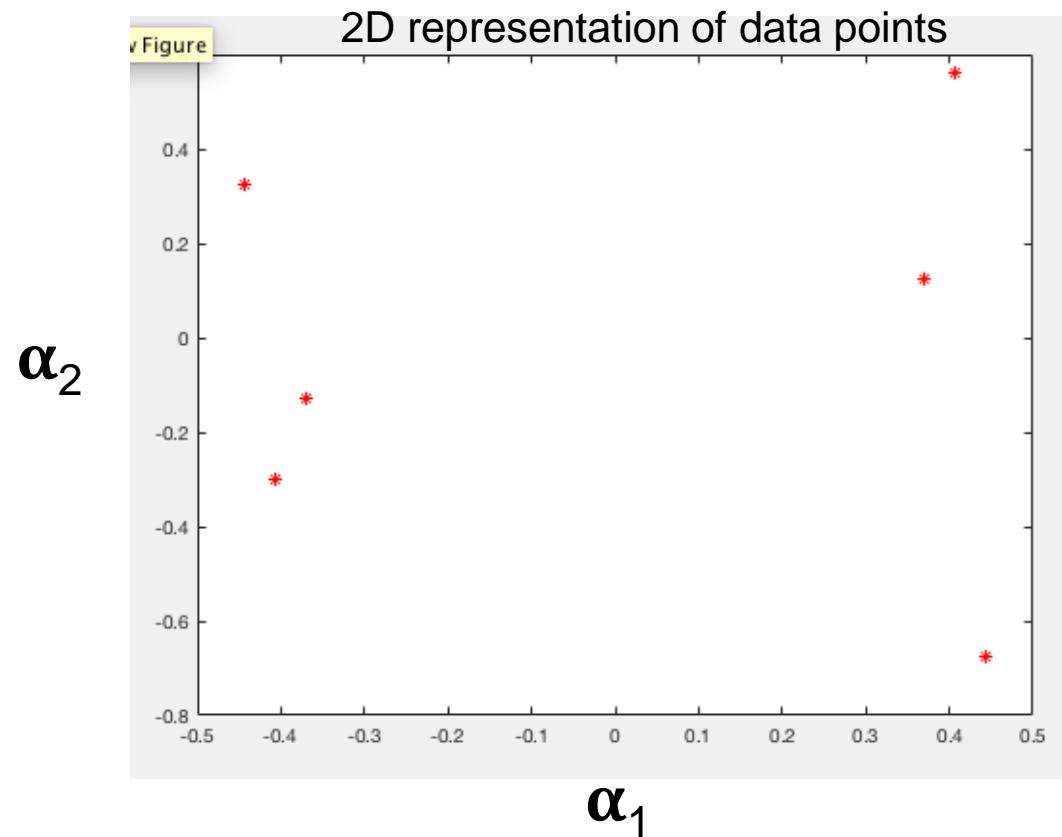
$$K_{ij} = \lambda \alpha_i \alpha_j$$





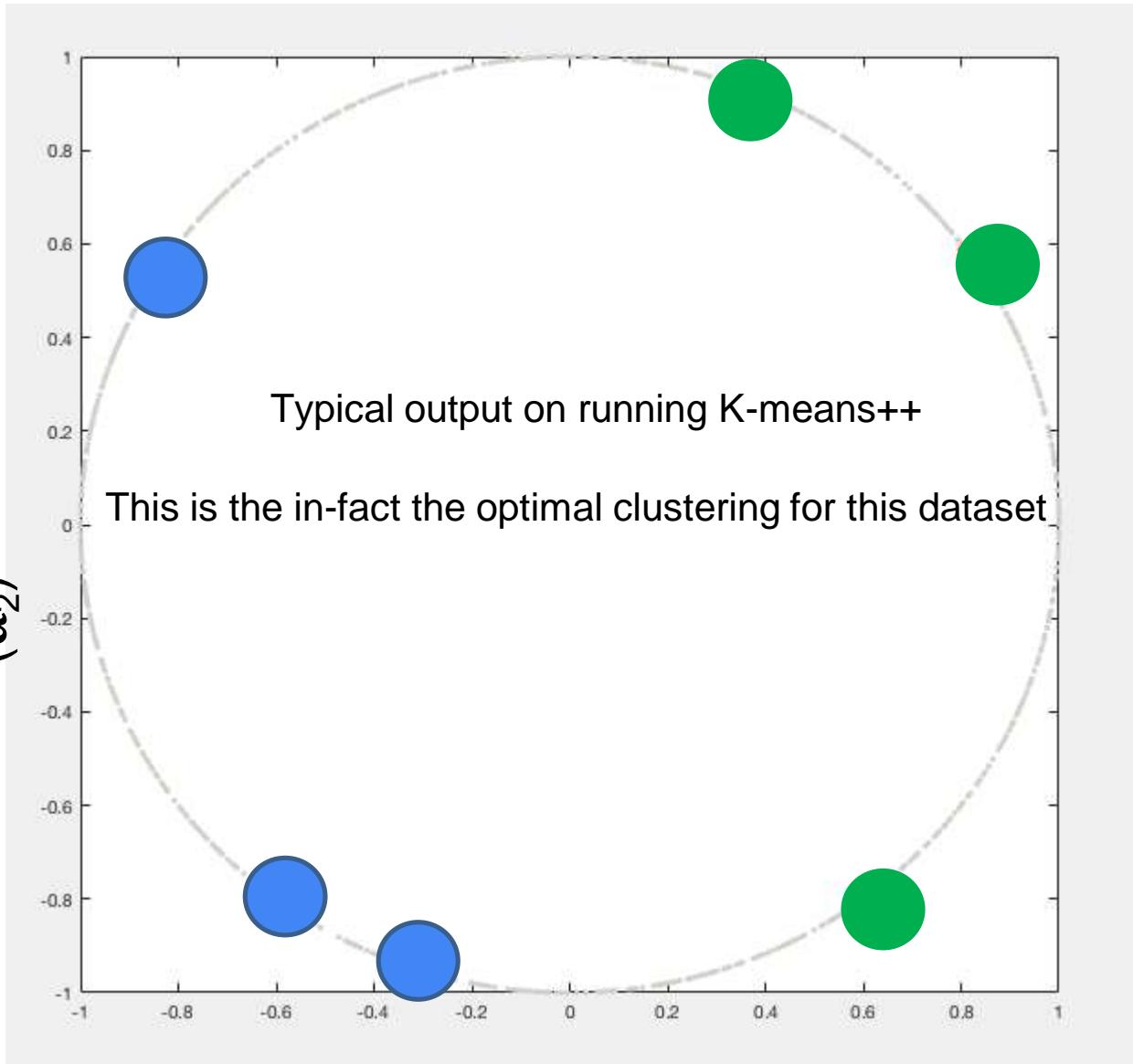


$$\begin{array}{c}
 \alpha_6 \\
 \left(\begin{array}{c} -0.1865 \\ -0.5742 \\ 0.5929 \\ 0.1441 \\ -0.4648 \\ 0.2171 \end{array} \right) \quad \left(\begin{array}{c} 0.4087 \\ -0.5908 \\ -0.2236 \\ -0.6512 \\ 0.0938 \\ 0.0320 \end{array} \right) \quad \left(\begin{array}{c} 0.7858 \\ -0.1177 \\ -0.0124 \\ 0.6067 \\ 0.0112 \\ 0.0187 \end{array} \right) \quad \left(\begin{array}{c} -0.1648 \\ -0.2281 \\ -0.5426 \\ 0.1849 \\ -0.5421 \\ -0.5462 \end{array} \right) \quad \left(\begin{array}{c} -0.1290 \\ -0.2992 \\ 0.3265 \\ 0.1261 \\ 0.5616 \\ -0.6752 \end{array} \right) \quad \left(\begin{array}{c} -0.3701 \\ -0.4071 \\ -0.4441 \\ 0.3701 \\ 0.4071 \\ 0.4441 \end{array} \right) \\
 \alpha_2 \\
 \alpha_1
 \end{array}$$



2D normalized representation of data points

$$\alpha_2 / \sqrt{(\alpha_1)^2 + (\alpha_2)^2}$$



$$\alpha_1 / \sqrt{(\alpha_1)^2 + (\alpha_2)^2}$$

Last point about Spectral Clustering

- If you search for spectral clustering on Google, you may find a graph-based algorithm
- Though **not covered in this course**, you are encouraged to read about it. In fact there are many variants that are called spectral clustering.
- If you think for a bit, you will realize it is nothing but kernel-K-means that we have seen with a specific choice of kernel (pseudo inverse of the Laplacian matrix associated with data)

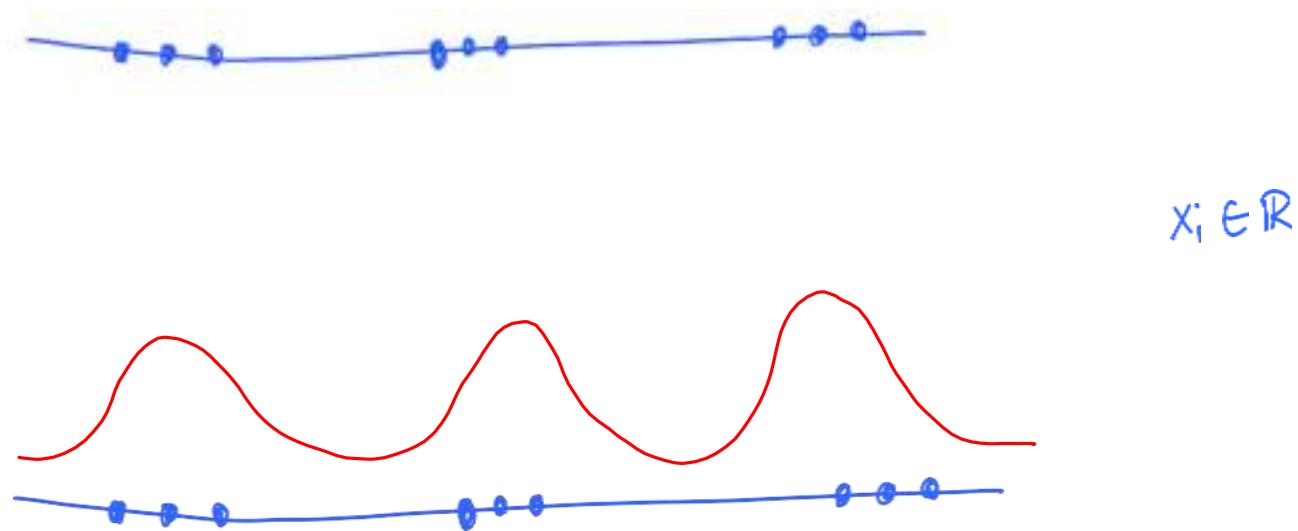
On Spectral Clustering: Analysis and an algorithm

Andrew Y. Ng
CS Division
U.C. Berkeley
ang@cs.berkeley.edu

Michael I. Jordan
CS Div. & Dept. of Stat.
U.C. Berkeley
jordan@cs.berkeley.edu

Yair Weiss
School of CS & Engr.
The Hebrew Univ.
yweiss@cs.huji.ac.il

Data



A new GENERATIVE STORY - MIXTURE OF GAUSSIANS

STEP 1: Pick which mixture a data point comes from

STEP 2: Generate a datapoint from that mixture

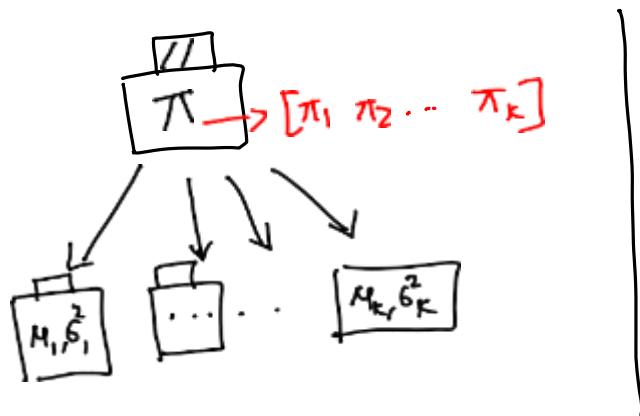
STEP 1: Generate a mixture component
among $\{1, \dots, k\}$ $z_i \in \{1, \dots, k\}$

$$P(z_i = \ell) = \pi_\ell$$

$$\left[\begin{array}{l} \sum_{\ell=1}^k \pi_\ell = 1 \\ 0 \leq \pi_\ell \leq 1 \end{array} \right]$$

STEP 2: Generate $x_i \sim N(\mu_{z_i}, \sigma_{z_i}^2)$

Gaussian Mixture Model



$\{x_1, \dots, x_n\} \rightarrow \text{OBSERVED}$

$\{z_1, \dots, z_n\} \rightarrow \text{UNOBSERVED/LATENT}$

LATENT VARIABLE MODEL

PARAMETERS?

$$\pi = [\pi_1, \dots, \pi_k]$$

$$\mu_k, \sigma^2_k \text{ + } k$$

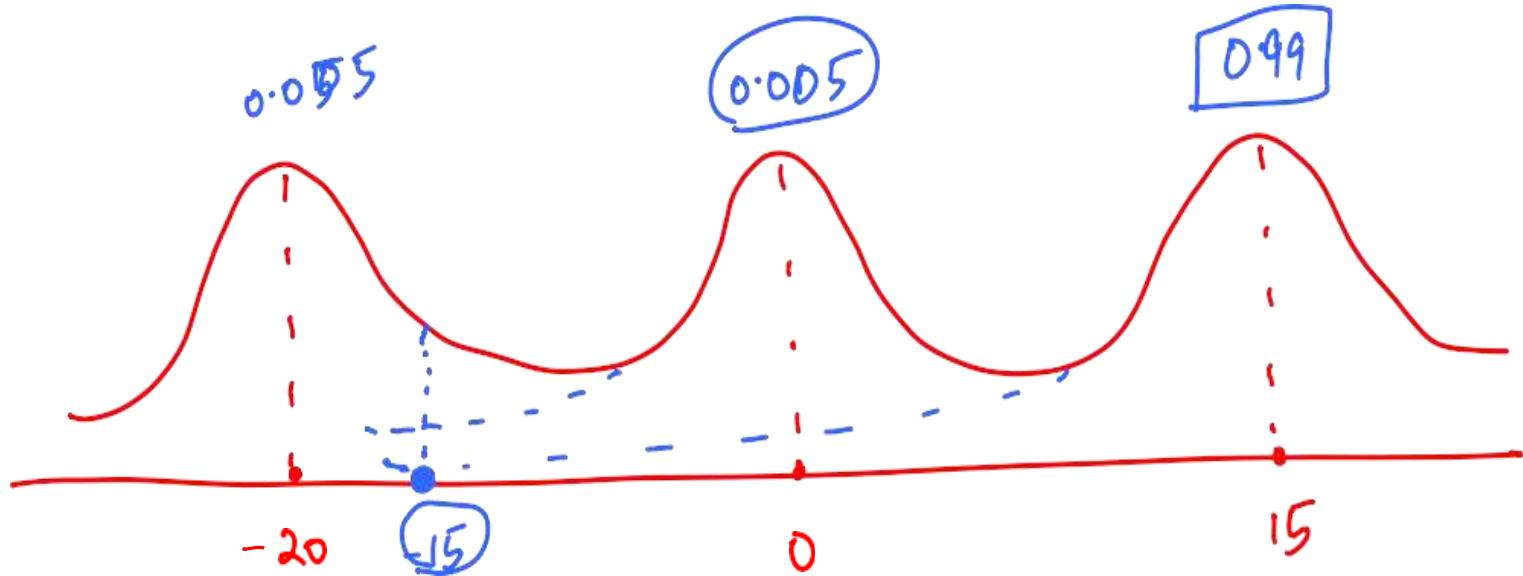
$$2k + k - 1$$

$$\text{Total: } 3k - 1$$

MAXIMUM LIKELIHOOD

$$L \left(\begin{matrix} \mu_1, \dots, \mu_K \\ \sigma_1^2, \dots, \sigma_K^2 \\ \pi_1, \dots, \pi_K \end{matrix} ; x_1, \dots, x_n \right) = \prod_{i=1}^n f_{\text{mix}} \left(x_i ; \begin{matrix} \mu_1, \dots, \mu_K \\ \sigma_1^2, \dots, \sigma_K^2 \\ \pi_1, \dots, \pi_K \end{matrix} \right)$$

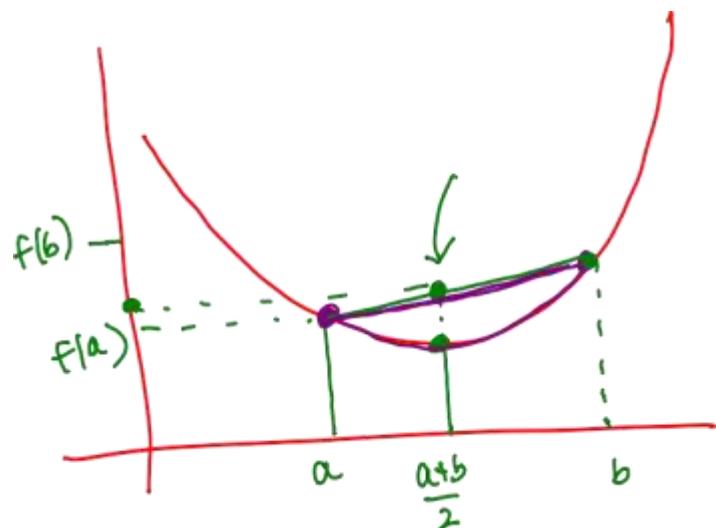
$$= \prod_{i=1}^n \left[\sum_{k=1}^K \pi_k f \left(x_i ; \mu_k, \sigma_k^2 \right) \right] \xrightarrow{\text{Gaussian density}}$$



$$\log L(\theta) = \sum_{i=1}^n \log \left(\sum_{k=1}^K \pi_k e^{-\frac{(x_i - \mu_k)^2}{2\sigma_k^2}} \cdot \frac{1}{\sqrt{2\pi}\sigma_k} \right)$$

- Not possible to solve analytically.
- Need some alternate ways to solve $\max_{\theta} \log L(\theta)$

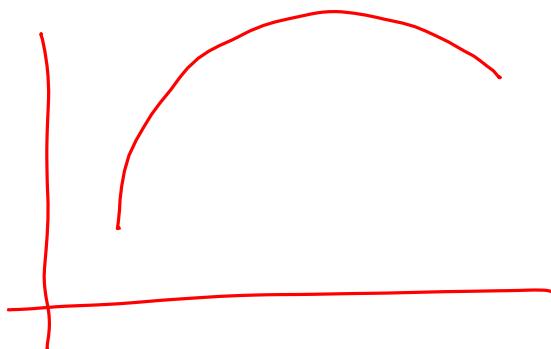
Quick Detour - Convex functions



$+ a, b$

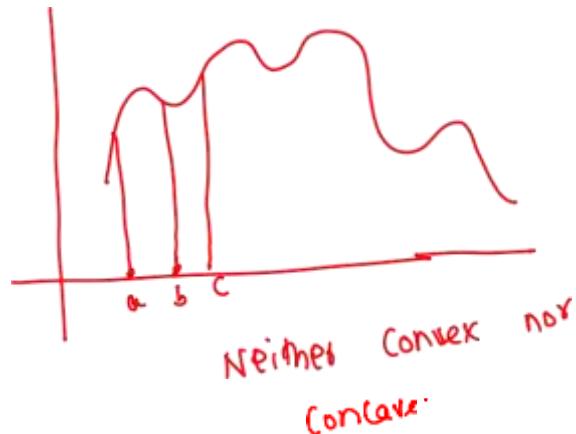
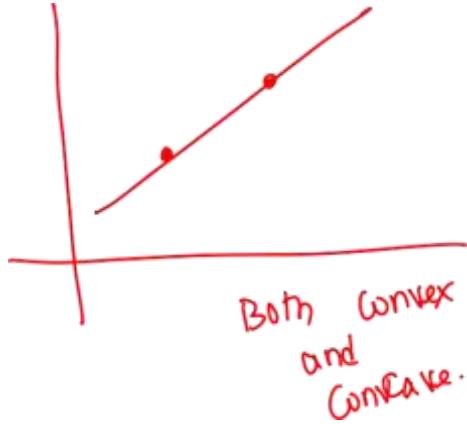
$$f\left(\frac{a+b}{2}\right) \leq \frac{f(a)+f(b)}{2}$$

\Rightarrow (convex function)



$$f\left(\frac{a+b}{2}\right) \geq \frac{f(a)+f(b)}{2}$$

\Rightarrow (concave function)



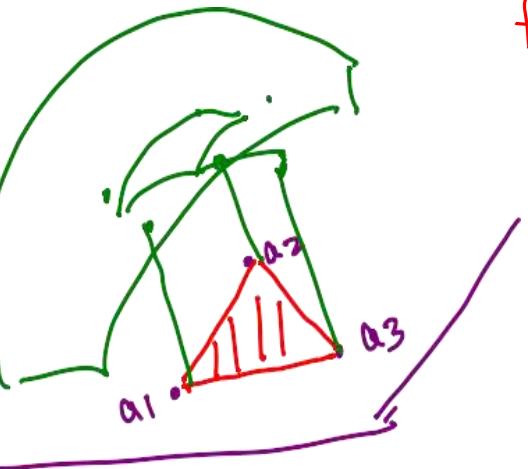
$$\forall a, b \quad f\left(\frac{1}{2}a + \frac{1}{2}b\right) \leq \frac{1}{2}f(a) + \frac{1}{2}f(b)$$

$$f(\lambda a + (1-\lambda)b) \leq \lambda f(a) + (1-\lambda)f(b) \quad \forall \lambda \in [0, 1]$$

f - CONCAVE

$$f(\lambda_1 a_1 + \lambda_2 a_2 + \dots + \lambda_k a_k) \geq \lambda_1 f(a_1) + \dots + \lambda_k f(a_k)$$

$\sum_{i=1}^k \lambda_i = 1 \quad 0 \leq \lambda_i \leq 1 \quad i$



JENSEN'S INEQUALITY

$$f\left(\sum_{k=1}^K \lambda_k a_k\right) \geq \sum_{k=1}^K \lambda_k f(a_k)$$

LOG IS A CONCAVE FUNCTION!

Recall

$$\log L(\theta) = \sum_{i=1}^n \log \left(\sum_{k=1}^K \left(\pi_k e^{-\frac{(x_i - \mu_k)^2}{2\sigma_k^2}} \frac{1}{\sqrt{2\pi}\sigma_k} \right) \right)$$

Introduce for every data-point $i \in \{\lambda_1^i, \dots, \lambda_K^i\}$ s.t

$$\sum_{k=1}^K \lambda_k^i = 1 \quad \text{and} \quad 0 \leq \lambda_k^i \leq 1 \quad \forall k.$$

$$\log L(\theta) = \sum_{i=1}^n \log \left(\sum_{k=1}^K \lambda_k^i \left(\pi_k e^{-\frac{(x_i - \mu_k)^2}{2\sigma_k^2}} \frac{1}{\sqrt{2\pi}\sigma_k} \right) \right)$$

$$\log L(\theta) \geq \text{modified-} \log L(\theta, \lambda)$$

$$= \sum_{i=1}^n \sum_{k=1}^K \lambda_k^i \log \left(\frac{\bar{\pi}_k e^{-\frac{(x_i - \mu_k)^2}{2\sigma_k^2}}}{\lambda_k} \right)$$

[JENSEN'S]

- Note the mod. $\log L(\theta, \lambda)$ gives a lower bound for $\log L(\theta)$ for any choice of λ

• What are we gaining?

KEY INSIGHT

- If we fix λ , it is "easy" to maximize wrt θ
- If we fix θ , maximize wrt λ is easy.

Fix λ and max over θ (mod- $\log L(\theta, \lambda)$)

$$\max_{\theta} \sum_{i=1}^n \sum_{k=1}^K \lambda_k^i \log \left(\frac{\pi_k e^{-\frac{(x_i - \mu_k)^2}{2\sigma_k^2}}}{\lambda_k^i} \right)$$

Take derivative w.r.t μ, σ^2 , we get

$$\hat{\mu}_R^{MML} = \frac{\sum_{i=1}^n \lambda_R^i x_i}{\sum_{i=1}^n \lambda_R^i}$$

$$\hat{\sigma}_R^{MML} = \frac{\sum_{i=1}^n \lambda_R^i (x_i - \hat{\mu}_R^{MML})^2}{\sum_{i=1}^n \lambda_R^i}$$

$$\max_{\pi_1, \dots, \pi_K} \text{mod_logL}(\theta, \lambda)$$

s.t. $\sum_{i=1}^K \pi_i = 1$

$$\hat{\pi}_R^{MML} = \frac{\sum_{i=1}^n \lambda_R^i}{n}$$

Fix θ and maximize w.r.t λ

max
 λ

$$\sum_{i=1}^n \left(\sum_{k=1}^K \lambda_k^i \log \left(\frac{\pi_k e^{-\frac{(x_i - \mu_k)^2}{2\sigma_k^2}}}{\lambda_k^i} \right) \right)$$

max for each i separately.

Solving this gives:

$$\hat{\lambda}_k^{MML} = \frac{\left(\frac{1}{\sqrt{2\pi}\sigma_k} e^{-\frac{(x_i - \mu_k)^2}{2\sigma_k^2}} \right) (\pi_k)}{\sum_{k=1}^K \frac{1}{\sqrt{2\pi}\sigma_k} e^{-\frac{(x_i - \mu_k)^2}{2\sigma_k^2}} \pi_k}$$

(Dempster et al)
1970's

EM - ALGORITHM

→ Initialization

$$\theta^0 = \{ \pi_1^0, \dots, \pi_k^0, \sigma_1^0, \dots, \sigma_k^0, \mu_1^0, \dots, \mu_k^0 \}$$

until convergence

$$(\|\theta^{t+1} - \theta^t\| \leq \epsilon)$$

Tolerance.

EXPECTATION
STEP

$$\lambda^{t+1} = \operatorname{argmax}_{\lambda} \text{modified-} \log L(\theta^t, \lambda)$$

MAXIMIZATION
STEP

$$\theta^{t+1} = \operatorname{argmax}_{\theta} \text{modified-} \log L(\theta, \lambda^{t+1})$$

end.

Expectation Maximization for GMM

- > Initialize all parameters (mean, variance, mixture proportions) that one has to estimate:
- > Until convergence,
 - > Assume the current parameters are correct and find the chance that each point came from every mixture
 - > Assume the chance that each point comes from every mixture is correct and find the best parameters.

(Dempster et al)
1970's

EM - ALGORITHM

→ Initialization

$$\theta^0 = \{ \pi_1^0, \dots, \pi_k^0, \sigma_1^0, \dots, \sigma_k^0 \}$$

Iteration

until convergence

$$(\|\theta^{t+1} - \theta^t\| \leq \epsilon)$$

Tolerance

EXPECTATION
STEP

$$\lambda^{t+1}$$

$$\underset{\lambda}{\operatorname{argmax}}$$

$$\text{modified-} \log L(\theta^t, \lambda)$$

MAXIMIZATION
STEP

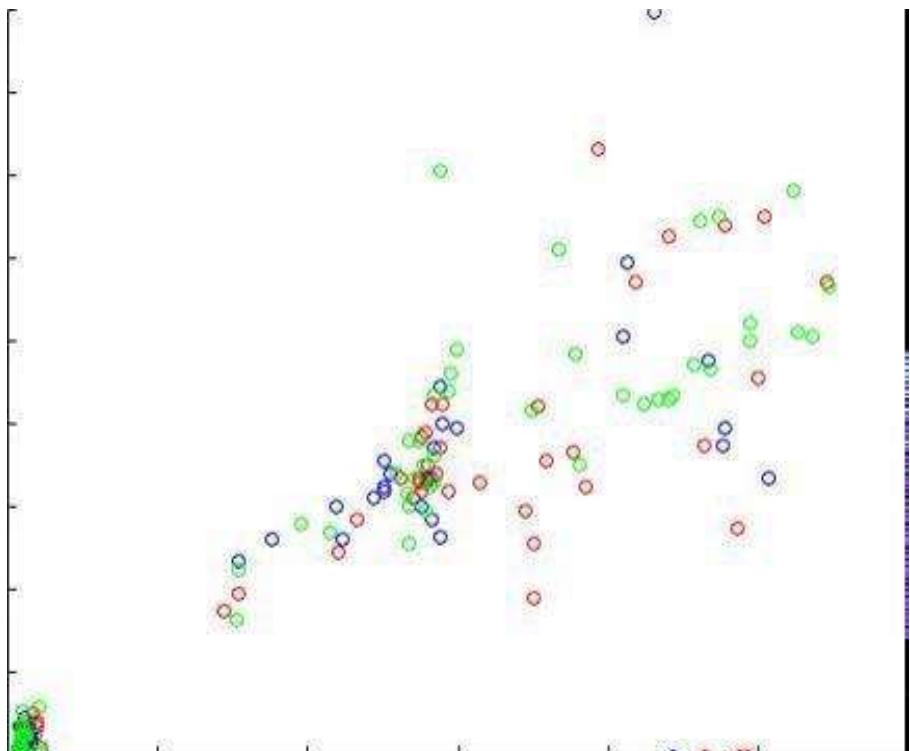
$$\theta^{t+1}$$

$$\underset{\theta}{\operatorname{argmax}}$$

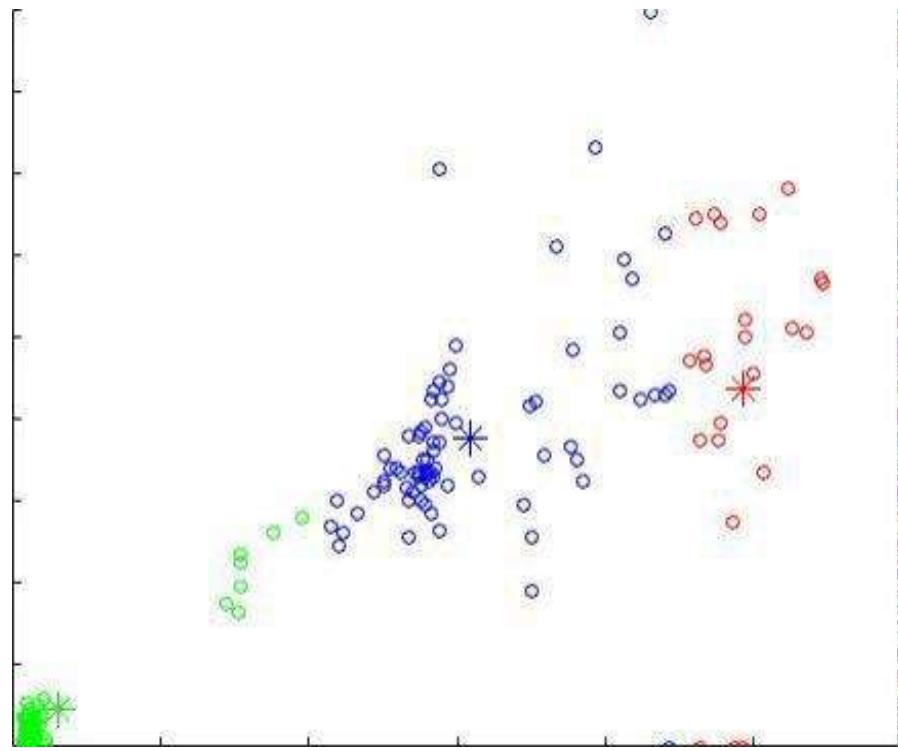
$$\text{modified-} \log L(\theta, \lambda^{t+1})$$

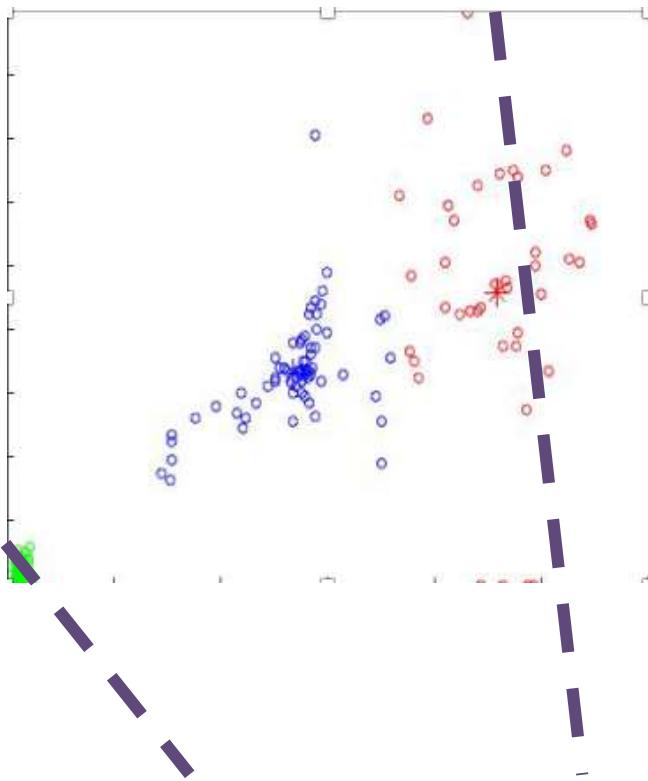
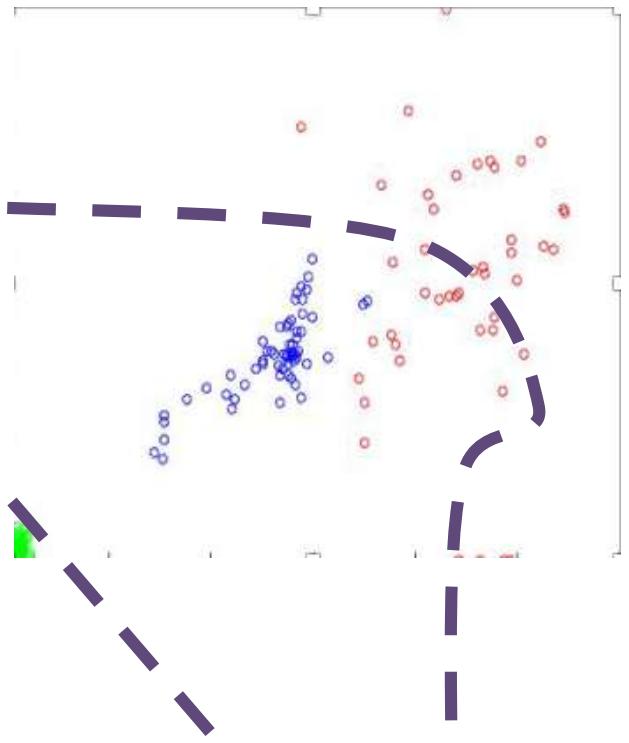
end.

EM-illustration



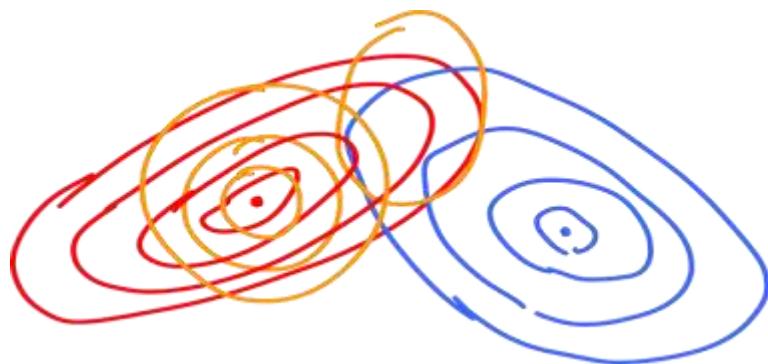
K means Illustration





A hard clustering is producing at every round by
assigning data to the cluster with highest probability

- EM produces "Soft Clustering" $\lambda_k^i = p(z_i | x_i)$
- EM takes into account variance of clusters as well



Why does EM work?

Recall

$$MML(\theta^t, \lambda | \theta^t) = \sum_{i=1}^n \sum_{k=1}^K \lambda_k^i \cdot \log(\pi_k^t \cdot f(x_i; \mu_k^t, (\sigma_k^2)^t) / \lambda_k^i)$$

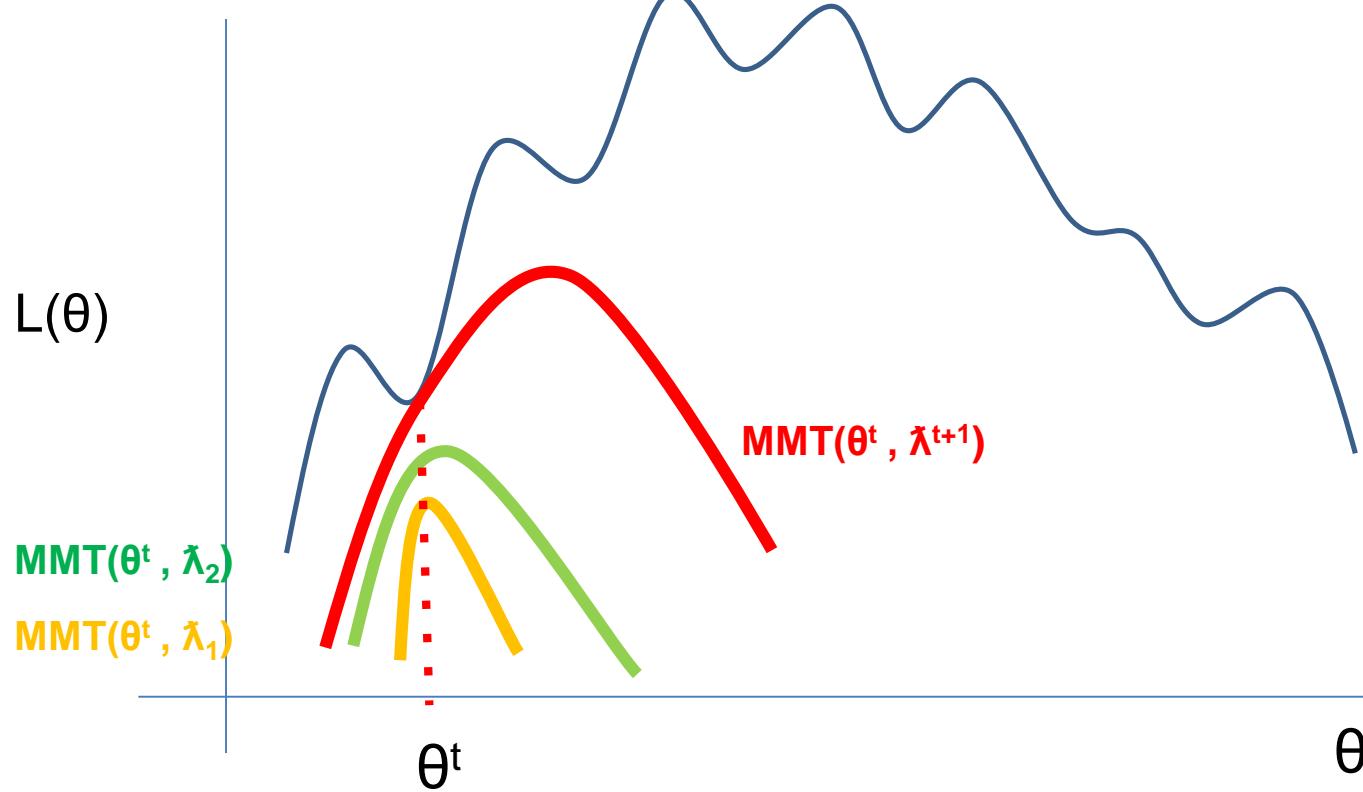
We know If we believed in the parameters Θ^t , then the best guess for λ is

$$\lambda_k^{t+1*} = \frac{f(x_i; \mu_k^t, \sigma_k^{t_2}) \cdot \pi_k^t}{f(x_i)}$$

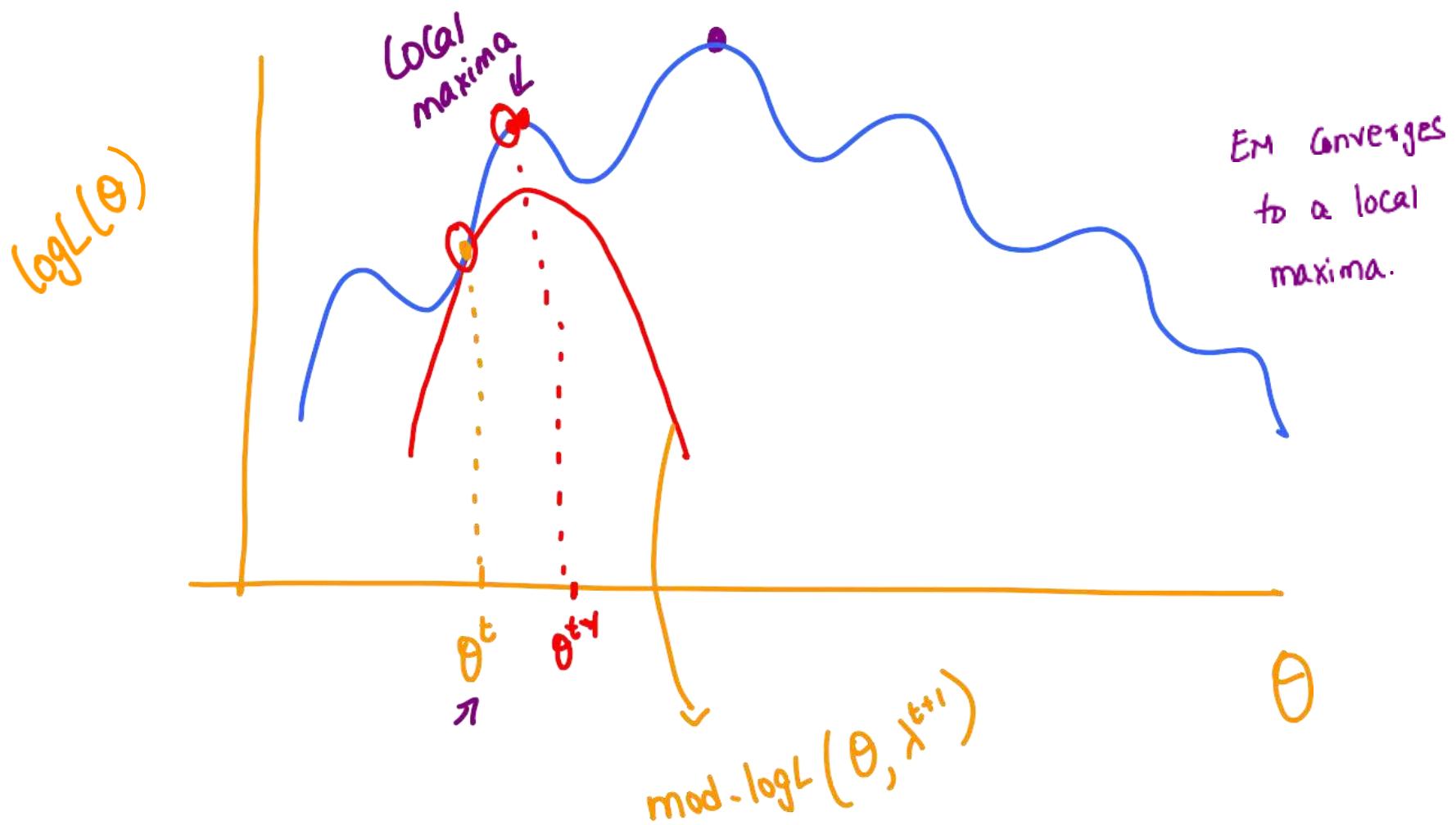
Substituting back the best values in MML, we can easily show that

$$MML(\theta^t, \lambda^{t+1}) = Loglikelihood(\theta^t)$$

So what?



$$MML(\lambda^{t+1}, \theta^t) = \text{Loglikelihood}(\theta^t)$$



$$\log L(\theta^t) = MML(\theta^t, \lambda^{t+1}) \leq MML(\theta^{t+1}, \lambda^{t+1}) \leq \log L(\theta^{t+1})$$

PREVIOUS
SLIDE

MAXIMIZATION

JENSEN

A few more points

Latent Variable Models are very popular.

EM is a general idea – not just for GMM

There is a probabilistic variant of PCA as well – called Probabilistic PCA

Summary

Unsupervised Learning

Representation learning
PCA, Kernel PCA

Clustering
K-means, Spectral Clustering

Estimation
Max-Likelihood, Bayesian
EM