

L 1.3 - Representation learning.

Goal: Given a set of "data points"
 "understand" something "useful"
 about it

Data points:- vectors in \mathbb{R}^d

$$\begin{bmatrix} \text{height} \\ \text{weight} \\ \text{age} \end{bmatrix} \in \mathbb{R}^3$$

Running Theme:- "Comprehension is compression"

Problem: Input = $\{x_1, x_2, \dots, x_m\}$
 $x_i \in \mathbb{R}^d$ $d \leftarrow$ no. of features

Output = Some compressed representation

Example:

$$\begin{array}{cccc} 7 & 2.5 & 0.5 & 0 \\ -16 & 5 & . & 0 \end{array}$$

→ 8 data points

→ A relationship

$$\begin{bmatrix} 1 \\ 2 \\ \vdots \\ n \end{bmatrix}; \{7, 2.5, 0.5, 0\} \}$$

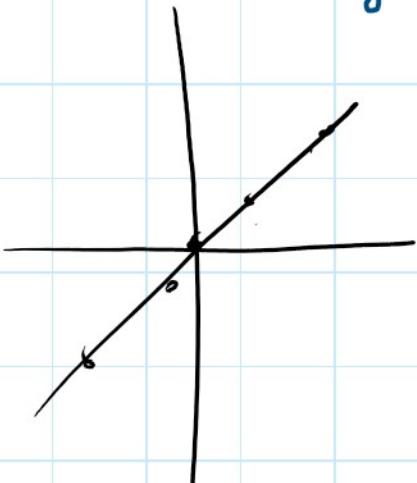
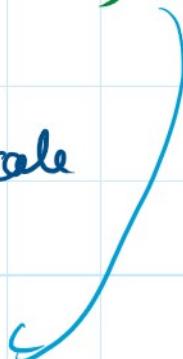
$|z_j, z^-, u^-, v^- \rangle$

→ 6 points

→ larger savings in scale

• Rep-I

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



Rep-II

$$\begin{bmatrix} 1/\sqrt{5} \\ 2/\sqrt{5} \end{bmatrix} \left\{ -7\sqrt{5}, 2\sqrt{5}, 0.5\sqrt{5}, 0 \right\}$$

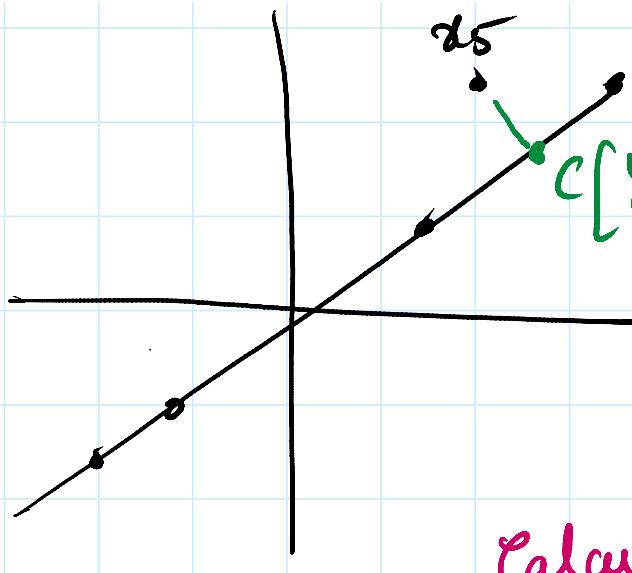
* Any vector along line can be chosen
except $\begin{bmatrix} 0 \\ 0 \end{bmatrix}$

Compression:

real number = d × n

compressed = d + n

* The above for perfectly fit. What if
there are other points



→ Something on the line act as proxy on line
→ Projection i.e dot product

Calculating projection

→ length of error vector

$$\begin{bmatrix} x_1 - c\omega_1 \\ x_2 - c\omega_2 \end{bmatrix}$$

$$\min_c (x_1 - (\omega_1))^2 + (x_2 - (\omega_2))^2$$

$$c^* = \left(\frac{x_1\omega_1 + x_2\omega_2}{\omega_1^2 + \omega_2^2} \right) \quad (\text{Scalar})$$

dot product

$$c[\omega_1, \omega_2] = \left(\frac{\mathbf{x}^\top \omega}{\|\omega\|} \right) \begin{bmatrix} \omega_1 \\ \omega_2 \end{bmatrix}$$

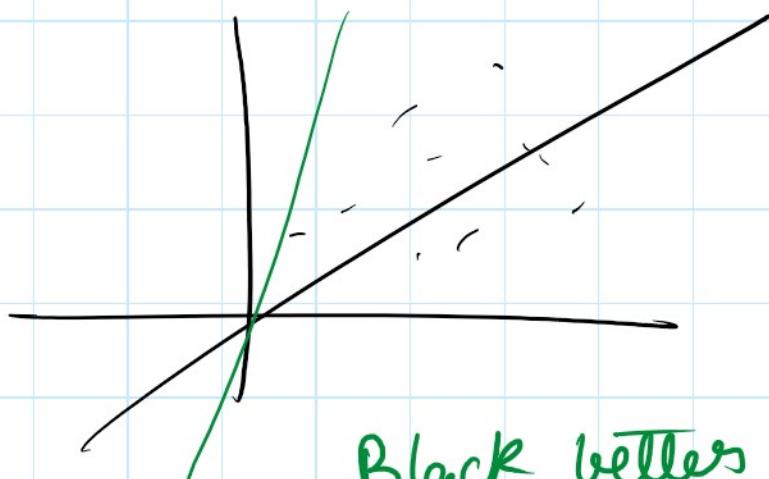
→ If choosing ω_1, ω_2 , where $\|\omega\|=1$,

$$c[\omega_1, \omega_2] = \mathbf{x}^\top \omega \begin{bmatrix} \omega_1 \\ \omega_2 \end{bmatrix}$$

2.1.4 - Representation learning p2

→ Real life lines not given, random dataset is given.

Goal: Develop way to find a "compressed" representation of data when data points not necessarily on line



Black better cuz less
reconstruction error.

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

ERROR (line, dataset) $\|\omega\| = 1$

$$= \sum_{i=1}^n \text{length}^2(x - (x^T \omega) \omega)$$

$$f(\omega) = \frac{1}{n} \sum_{i=1}^n \|x - (x^T \omega) \cdot \omega\|^2$$

$$= \frac{1}{n} \sum_{i=1}^n [x_i - (x_i^T \omega) \cdot \omega] [x_i - (x_i^T \omega) \cdot \omega]$$

$$= \frac{1}{n} \sum_{i=1}^n [x_i^T x_i - (x_i^T \omega)^2] = \cancel{(x_i^T \omega)^2} + \cancel{(x_i^T \omega)^2}$$

$$= \frac{1}{n} \sum_{i=1}^n [x^T x_i - (x_i^T \omega)^2]$$

↑
constant, irrelevant

$$g(\omega) = \frac{1}{n} \sum_{i=1}^n (x_i^T \omega)^2$$

$$= \frac{1}{n} \sum_{i=1}^n (\underbrace{\omega^T x_i}_{d \times 1}) (\underbrace{x_i^T \omega}_{1 \times d})$$

$$= \frac{1}{n} \sum_{i=1}^n \omega^T [x_i x_i^T] \omega$$

$$= \omega^T \left[\frac{1}{n} \sum_{i=1}^n (x_i x_i^T) \right] \omega$$

↓
co-variance matrix

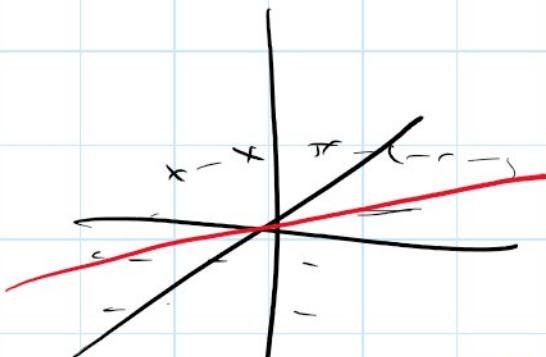
$$= \omega^T C \omega$$

- * ω is eigenvector corresponding to the max eigen value of C
- * $C \rightarrow$ covariance matrix
 - The above eigen vector is an optimization problem.

21.5- Representation learning p3

→ Is this compression satisfactory in all cases.

Eg:- 3-D where points on a plane.
line lies on the plane



- In 2-D cases, we are taking this as error

- In 3-D when points are on a plane, error also has some information

→ algorithm gives line but we also want to outcome that it is in 2-D.

to capture that it is in 2-D.

→ If we see the error vector, they lie along a single direction

Procedure

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

↓ Find ω

$$(x^T \omega) \cdot \omega$$

↓ error

$$\underline{x - (x^T \omega) \cdot \omega}$$

has some information
with it

Possible algorithm:

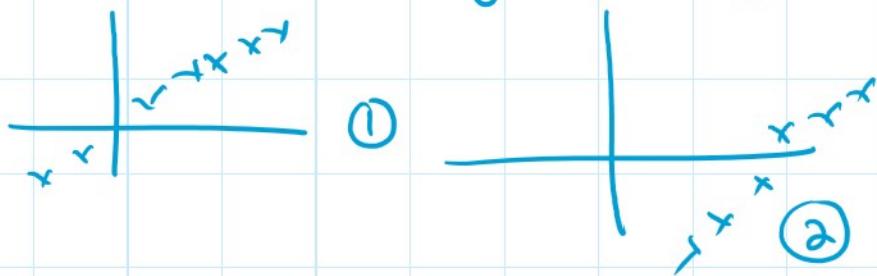
• Input: $\{x_1, \dots, x_n\}$ $x_i \in \mathbb{R}^d$

- $e = \frac{1}{n} \sum_{i=1}^n x_i$; $x_i = x_i - e$
- Find "best" line $\omega_1 \in \mathbb{R}^d$
- $x_i = x_i - (x_i^T \omega) \cdot \omega$
- Repeat to obtain ω_2

! Issue.

→ The above system works for
data about origin i.e. $y = mx$.

data about oxygen i.e $y = mx$



→ ② doesn't work.

→ so solution for that is to calculate mean & remove it

Questions for the above systems:-

→ How to solve $\max \omega^T C \omega$?

→ How many times to repeat?

→ Where exactly is compression happening?

→ What representations are we learning?

7.1.6- PCA - I

$$D = \{x_1, \dots, x_n\} \quad x_i \in \mathbb{R}^d$$

$$\omega_1 = \underset{\|\omega\|}{\operatorname{argmax}} \omega^T C \omega \quad \left| \begin{array}{l} \omega_2 = \operatorname{argmax} \omega^T C' \omega \\ \|\omega_2\|^2 = 1 \end{array} \right.$$

$$C = \frac{1}{n} \sum_{i=1}^n x_i x_i^T \quad \left| \begin{array}{l} C' = \frac{1}{n} \sum_{i=1}^n x_i' x_i'^T \end{array} \right.$$

$$x'_i = x_i - (x_i^T \omega_1) \omega_1$$

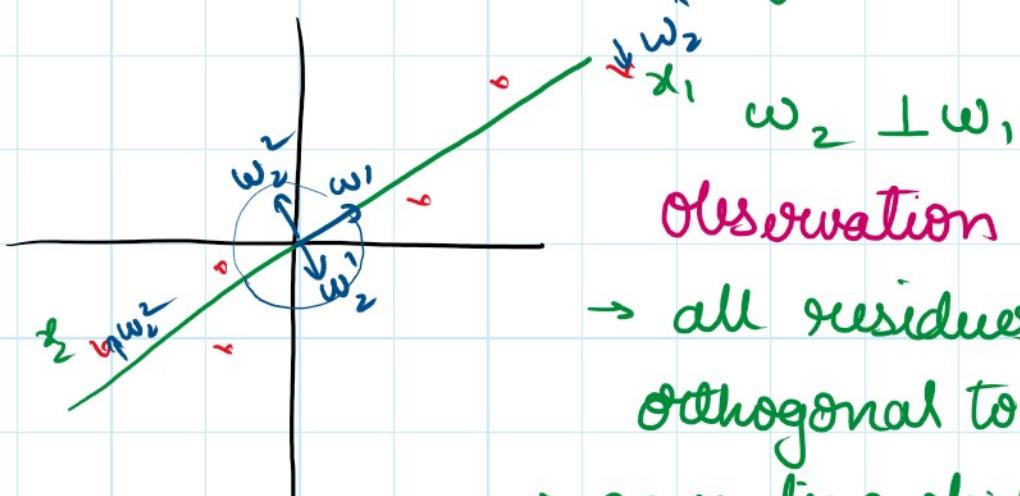
→ Best fit line
direction ...

→ Error line
direction ...

→ Best fit line
direction w_1

→ Error line
direction w_2

Question: what can we say about w_1 & w_2



Observation

→ all residues are
orthogonal to w_1

→ any line which
minimizes sum of errors
w.r.t. residues must
be orthogonal to w_1

$$w_2^T w_1 = 0$$

→ going into higher dimensions, we can
get residuals of residuals & so on

By continuing,

$\{w_1, \dots, w_d\} \rightarrow$ orthonormal
vector.

Properties:

- $\|w_k\|^2 = 1 + k$. [forall k]
- $w_i^T w_j = 0 \text{ if } i \neq j$

- $\omega_i^\top \omega_j = 0 \quad \forall i \neq j$

Residue after Round - I

$$\left\{ (x_i - (x_i^\top \omega_1) \omega_1), \dots, x_i \in \mathbb{R}^d \right. \\ \left. (x_n - (x_n^\top \omega_1) \omega_1) \right\}$$

- $\omega_2 \rightarrow$ Best fit line

- $\omega_1^\top \omega_2 = 0$

Residues after Round 2 .

$$\left\{ \left[x_i - (x_i^\top \omega_1) \omega_1 - \left[(x_i^\top \omega_1) \omega_1 \right] \omega_2 \right] \right.$$

$$\left. \left(x_i^\top \omega_2 - (x_i^\top \omega_1) \omega_1^\top \omega_2 \right) \omega_2 \right\}$$

$$\Rightarrow \left[x_i - (x_i^\top \omega_1) \omega_1 - (x_i^\top \omega_2) \omega_2 \right]$$

Residues after d rounds -

$$\forall i \quad x_i - ((x_i^\top \omega_1) \omega_1 + (x_i^\top \omega_2) \omega_2 + \dots + (x_i^\top \omega_d) \omega_d)$$

$$= x_i - \sum_{j=1}^d (x_i^\top \omega_j) \omega_j$$

$$x_i - \sum_{j=1}^d (x_i^\top \omega_j) \omega_j$$

If you run for d rounds, the residual will be 0 because the points will align

$$x_i - \left(\sum_{j=1}^d (x_i^\top \omega_j) \omega_j \right) = 0 \in \mathbb{R}^d$$

$$\Rightarrow x_i = \sum_{j=1}^d (x_i^\top \omega_j) \omega_j$$

What have we gained?

- If data lies in a lower dimensional subspace, residues become 0 much earlier than d rounds.

Example: 100d but residuals 0 after 3 rounds

$$D = \{x_1, \dots, x_n\} \subset \mathbb{R}^{100}$$

$$\forall i \quad x_i = (x_i^\top \omega_1) \omega_1 + (x_i^\top \omega_2) \omega_2 + (x_i^\top \omega_3) \omega_3$$

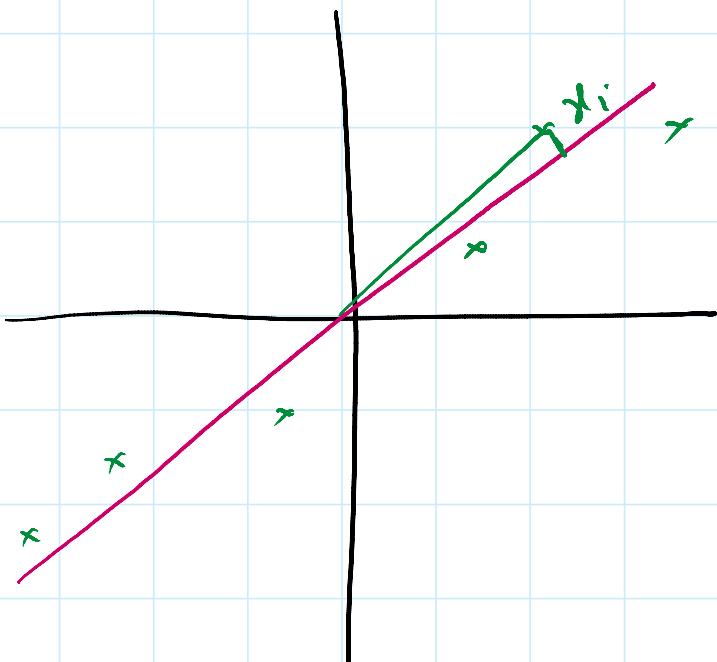
Rep
 $\{\omega_1, \omega_2, \omega_3\}$
 for full dataset
 original :- $100 \times n$ after :- $3 \times 100 + 3n$

Co-efficients
 $x_i \rightarrow [x_i^T \omega_1, x_i^T \omega_2, x_i^T \omega_3]$
 For each point

$$\rightarrow d \times n \Rightarrow d \times k + k \times n$$

What if data "approximately" fits in a k -dimensional space?

→ use a tolerance.



For any $w \in R^d$

For all $s.t \|w\| \leq 1$

$$\|x_i\|^2$$

$$= \|x_i - (x_i^T w) w\|^2 + \|(x_i^T w) w\|^2$$

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

↓
For any w
avg length
of dataset

↓
avg error, minimize

$$+ \frac{1}{n} \sum_{i=1}^n \|(\mathbf{x}_i^T w) w\|^2$$

↓
maximize ①

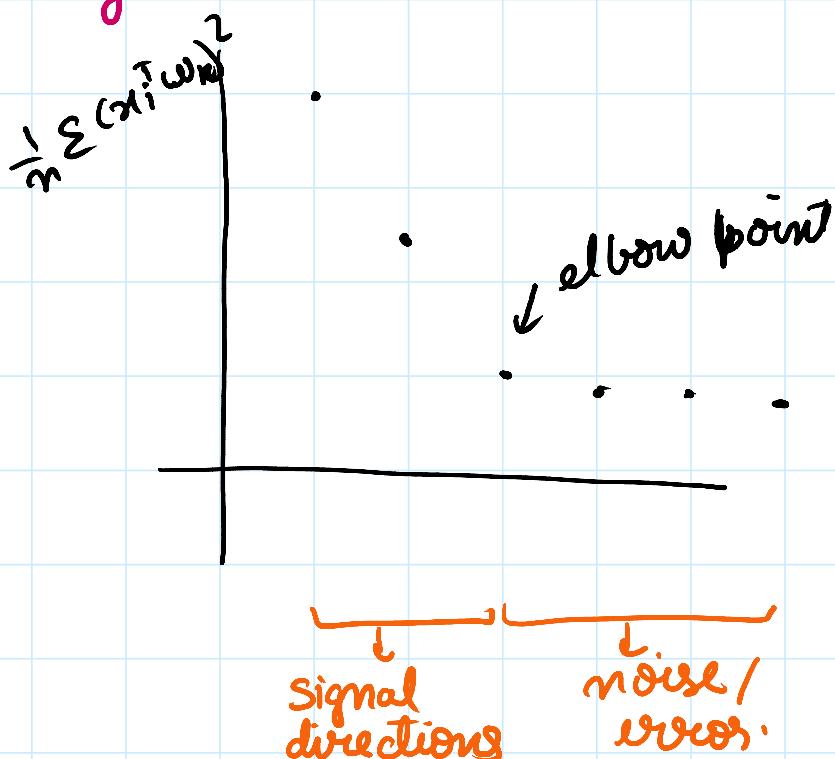
$$\textcircled{1} \rightarrow \frac{1}{n} \sum_{i=1}^n \|(\mathbf{x}_i^T w) w\|^2$$

$$= \|(\mathbf{x}_i^T w)^2 w^2\|$$

$$= \frac{1}{n} \sum_{i=1}^n \|(\mathbf{x}_i^T w)\|^2$$

* Larger the value of $\frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i^T w)^2$
the better the fit

graphing it: -



- We choose threshold.
 - Say 95% represented by k, we take only k

2.1.7 - PCA-II

Referring back to optimization problem

2.1.4

$$\text{argmax}_{\omega} \omega^T C \omega \quad C = \frac{1}{n} \sum_{i=1}^n x_i x_i^T$$

$$\omega \| \omega \|_2^2 = 1$$

$C \rightarrow$ covariance matrix

Soln:

- $\omega_1 \rightarrow$ eigenvectors corresponding to largest eigenvalue C
[Hilbert min-max problem]
- $\{\omega_1, \dots, \omega_d\}$ The eigenvectors of C forms orthonormal basis
- $\omega_k \rightarrow$ least line one can obtain in round k.

What do eigenvalues mean?

- we know

$$C \omega_i = \lambda_i \omega_i$$

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

$$\omega_i^T (\omega_i) = \omega_i^T (\lambda_i \omega_i)$$

$$\omega_i^T (\omega_i) = \lambda_i$$

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

□

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

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

← will
always
+ve.

* $\lambda_k(c)$ graph vs k gives same elbow graph.

* eigenvalues are non-negative

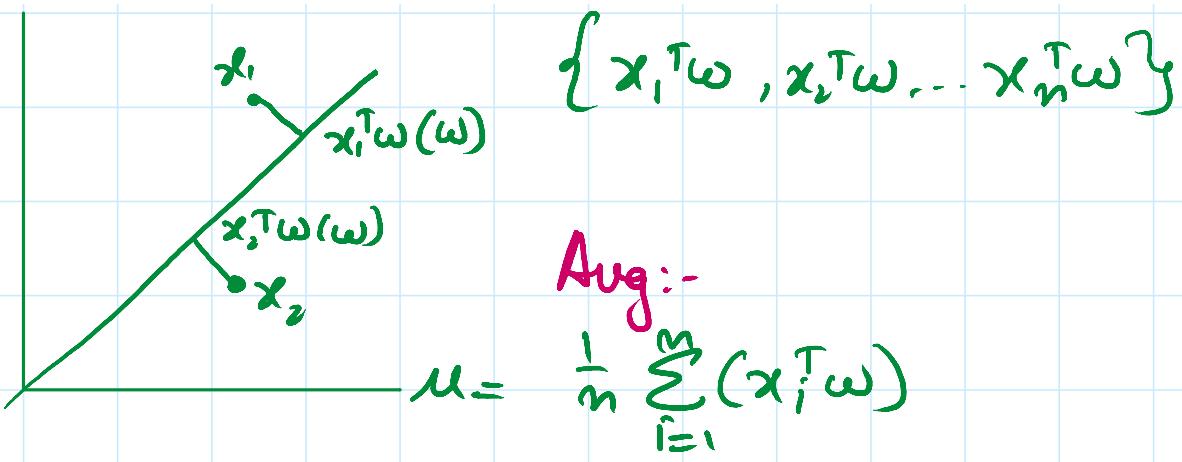
* Rule of Thumb - (For no. of dimensions)

$$\frac{\sum_{i=1}^k \lambda_i(c)}{\sum_{i=1}^d \lambda_i(c)} \geq \text{Threshold.}$$

[0.95 in practice]

Taking a collection:-

For a ω , taking collection of points



$$\mu = \left(\frac{1}{n} \sum_{i=1}^n x_i \right)^T \omega = (0^T \omega) \omega$$

↙

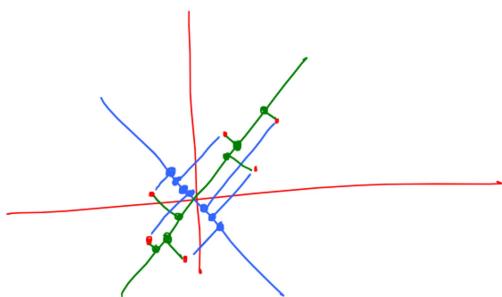
$\mu = 0$ [For centered dataset]

Variance:-

$$\frac{1}{n} \sum_{i=1}^n (x_i^T \omega - \text{mean})^2$$

$$\sigma^2 = \frac{1}{n} \sum_{i=1}^n (x_i^T \omega)^2$$

★ Error minimization on centered dataset \rightarrow variance maximization



Want directions where
Projections don't "crowd-up"
(i.e., Variance is not small.)

Eg:- Refer pdf.

basically that the components
are de correlated

PCA $\rightarrow \{\omega_1, \dots, \omega_k\}$

\rightarrow Means that

\rightarrow only k dimensions used to
represent d dimensions

\rightarrow PCA \rightarrow dimensionality reduction "

* PCA finds combination of features
that are de correlated.

\rightarrow Loosely speaking independent
of each other

Assignment: -

$$2) x \text{ on } L \quad x = (2, 5)$$
$$\omega = [1, 1]$$

$$\frac{x^T \omega}{\|\omega\|^2} \omega \quad [2 \ 5] \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

$$\frac{7}{2} = \begin{bmatrix} 3.5 \\ 3.5 \end{bmatrix}$$

$$\sqrt{2 \times 12.25} = \sqrt{24.5} \cdot$$

$$\sqrt{4.5} \cdot$$

$$d \times n$$

$$30,000$$

$$d \times k + k \times n \cdot$$

$$120 + 6000$$

$$30000$$

$$4120 \cdot$$

$$\frac{25880}{20000} = 86.25 \cdot$$

$$x_i = \begin{bmatrix} x \\ y \end{bmatrix}$$

$$x_i x_i^T = \begin{bmatrix} x \\ y \end{bmatrix} \begin{bmatrix} x & y \end{bmatrix}$$

$$\begin{bmatrix} x^2 & xy \\ xy & y^2 \end{bmatrix}$$

$$\overbrace{\begin{array}{cccc} 9 & 4 & 1 & 4 \\ 9 & 4 & 1 & 4 \end{array}}^{2 \times 2} \quad \begin{array}{c} 0 \\ -2 \\ -2 \\ - \end{array} \cdot$$

0 -2 -2

0 1 1

36
-4

-4
2