

Lecture 19] Summary of PCA

1) [optional] Normalize the data (usually per-variable separately)

↙ column means

2) Centering the data: $A_c = A - \bar{\mu}$

3) Compute covariance matrix Σ

$$\Sigma = \frac{1}{(N-1)} \cdot (A_c \cdot T) @ (A_c)$$

↳

shape: (M, M)

4) Compute eigen vectors and eigen values using Numpy: p ↳ shape: (M, M)

$$e\text{-Vals}, e\text{-Vecs} = np.linalg.eig(\Sigma)$$

↳ vector: len = $(M, 1)$

5) Sort e_vals $\xrightarrow{PC1}$ $\xrightarrow{PC2}$ $\xrightarrow{PC3}$ high \rightarrow low.

e.g. $\overbrace{\{3.2, 1.1, 0.03\}}^{e_vecs} = e_vals \quad M=3$

$\overbrace{e_vecs}$

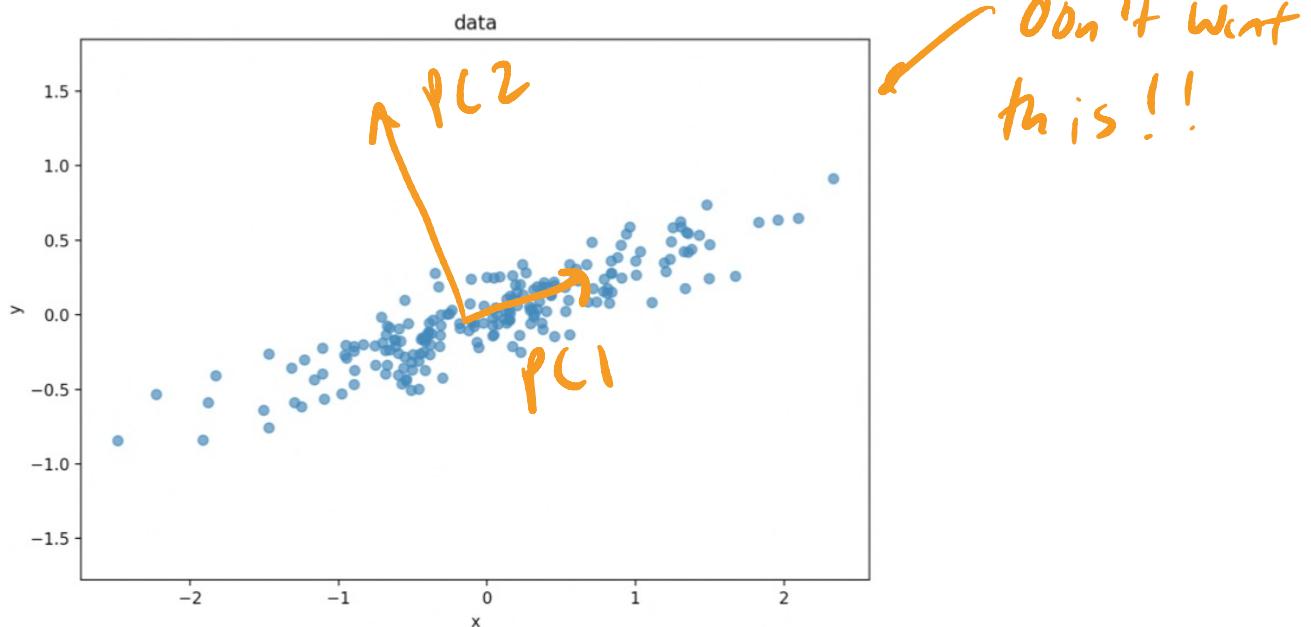
* Sort Columns of P in the same order

that you sorted the eigenvalues. $\xrightarrow{\text{after sorting according to } e_vals.}$

$$A = \begin{bmatrix} & x & y & z \\ & \{ & \{ & \{ \\ & \{ & \{ & \{ \end{bmatrix}$$

$$P = \begin{bmatrix} & PC1 & PC2 & PC3 \\ & \{ & \{ & \{ \\ & \{ & \{ & \{ \end{bmatrix}$$

problem: if don't sort e_vecs like e_vals



6 a) Compute prop-Vcr accounted for by each PC:

$$\text{total} = 3.2 + 1.1 + 0.03 = 4.33$$

$$\text{prop-Vcr} = \left[\frac{3.2}{4.33}, \frac{1.1}{4.33}, \frac{0.03}{4.33} \right]$$

6 b) Compute cumulative prop-Vcr accounted by by 1st K eigenvalues/PCs.

$$\text{cum-prop-Vcr} = \left[\underbrace{\frac{3.2}{4.33}}_{\text{PC1}}, \underbrace{\frac{(3.2+1.1)}{4.33}}_{\text{PC1+PC2}}, \underbrace{\frac{(3.2+1.1+0.03)}{4.33}}_{\text{PC1+PC2+PC3}} \right]$$

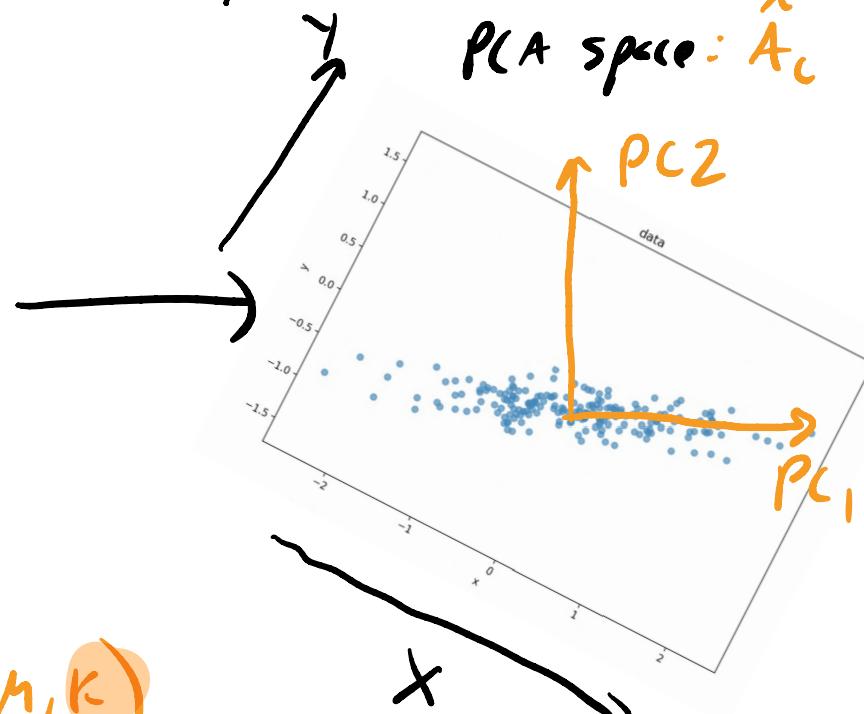
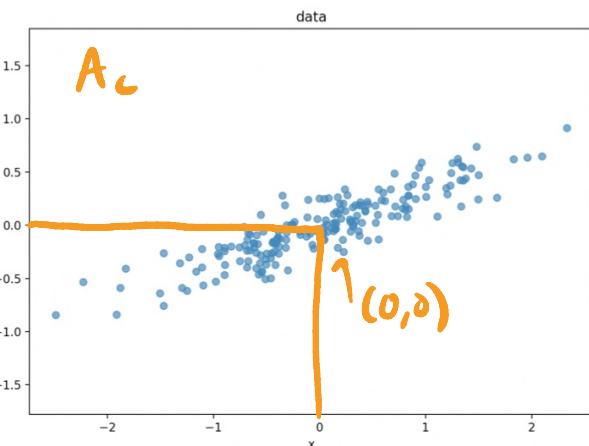
↓
Always = 1
(100%)

6 c) Threshold/cut off point: Want to keep

at least X% Vcriance:
(90%)

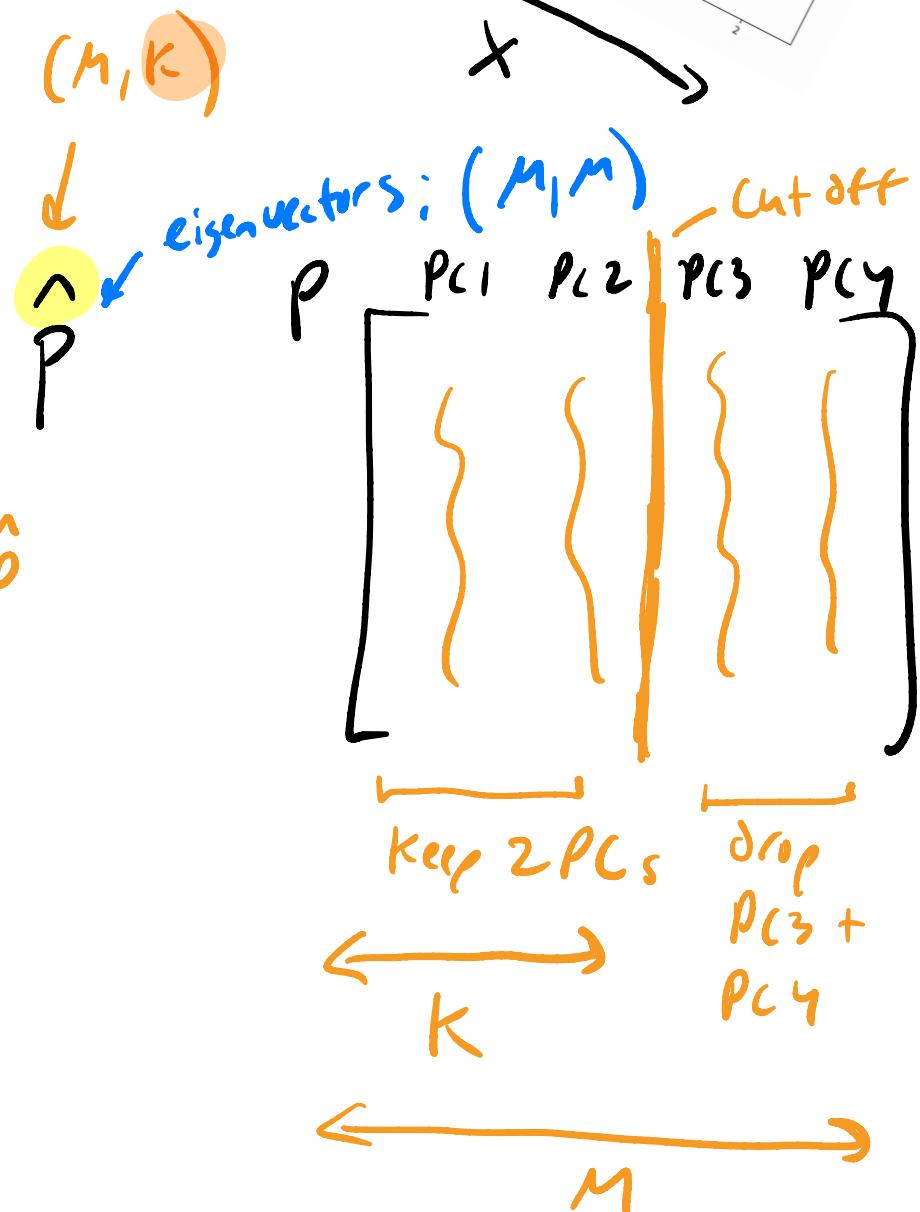
⇒ We ultimately keep K PCs

7) Project data from Data Space \rightarrow PCA Space
 Data Space PCA Space: \hat{A}_c



$$\hat{A}_c = A_c @ P$$

(N, K) (N, M)
 $\cap \cap \cap$ \downarrow
 $\hat{A}_c = A_c @ P$



Keep k PCs: \hat{P}

$$\hat{P} = P[:, :k]$$

Reduced Dimension

$$M \rightarrow k$$

⇒ Could be done

8) [optional] Reconstruction of data

P(A Space → data space)

if we did not Normalize

(K, M)

(M, K)

$$A_{\text{reconstruct}} = \underbrace{\hat{A}_c}_{(N, M)} @ \underbrace{\hat{P}^T}_{(N, M)} + \underbrace{\vec{\mu}}_{(1, M)}$$

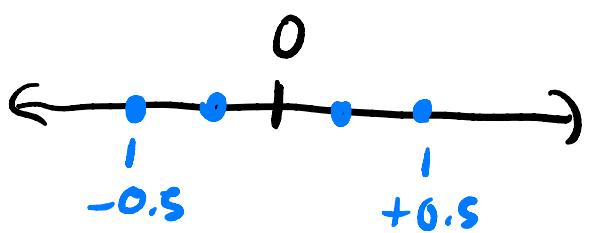
Keep N samples,
return to M data Vars
(original ones)

if we normalized: Original Data Range in each Var

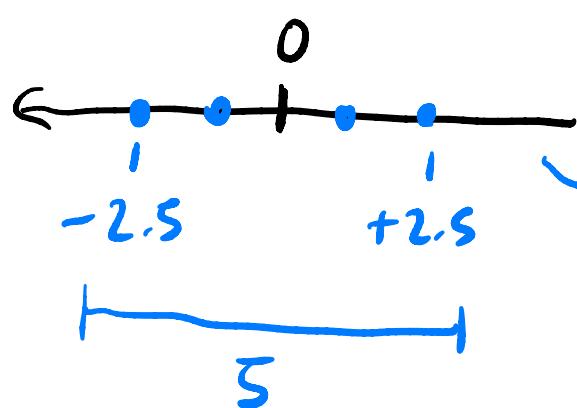
$$A_{\text{reconstruct}} = \vec{s} \cdot (\hat{A}_c @ \hat{P.T}) + \vec{\mu}$$



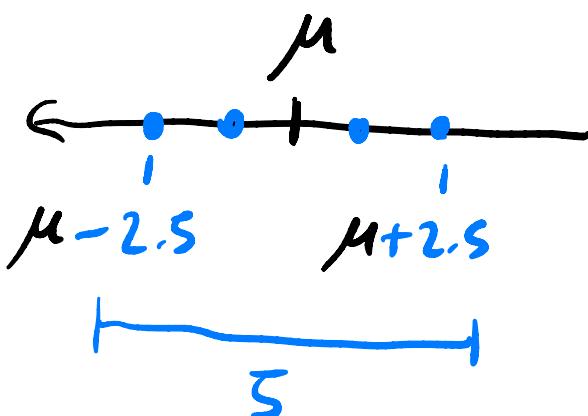
e.g. Original Range $s = 5$



multiply by s



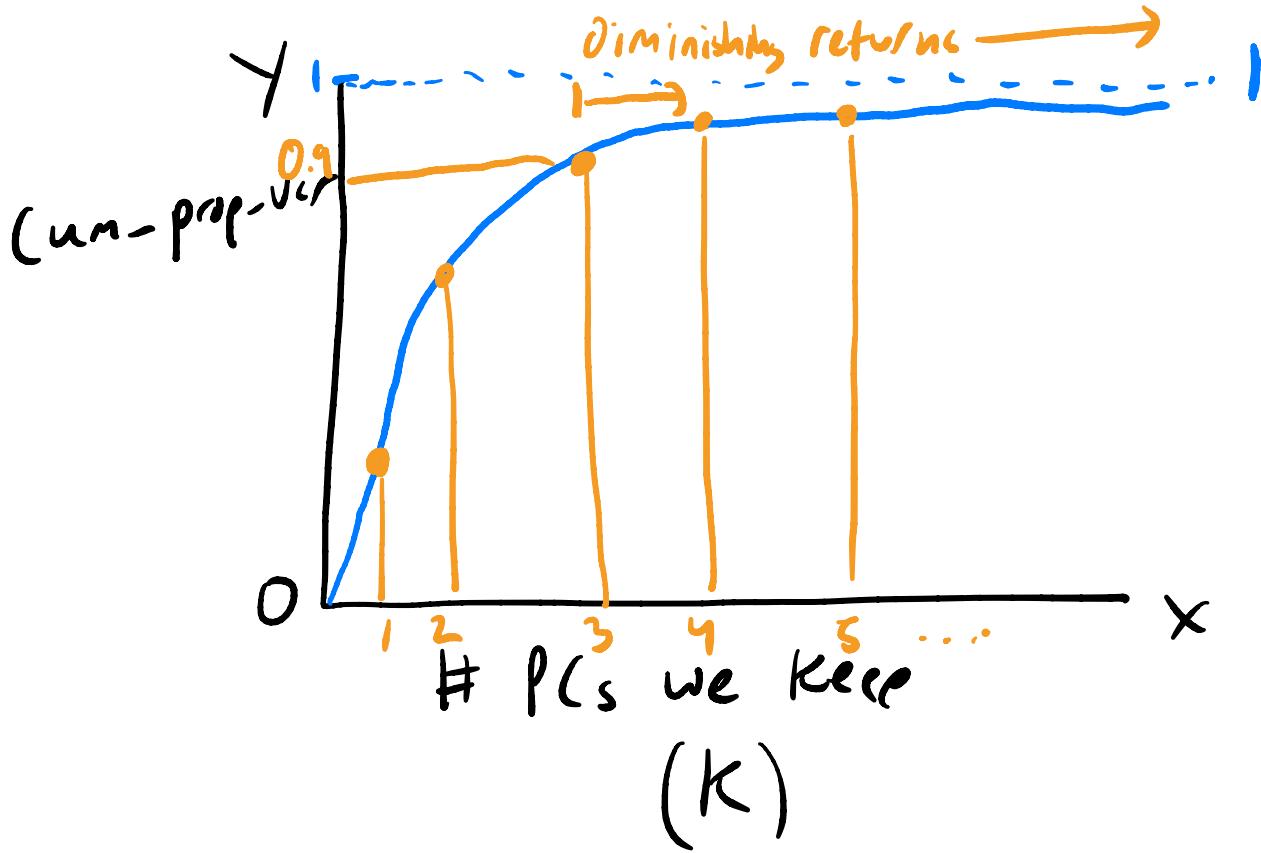
$+ \vec{\mu}$



 Elbow plot: useful if have large # variables
(e.g. $M=1000$)

Plot Cumulative prop Var (Y axis)

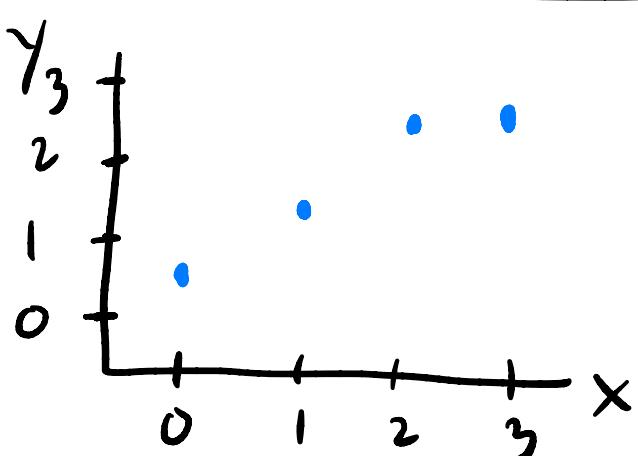
PC # on X axis



\Rightarrow 90% is a good cut off point: keep PC_1, PC_2, PC_3 ,
drop $PC_4 +$

Example:

X	Y
0	0.5
1	1.1
2	2.3
3	2.5



1) Normalize \rightarrow Skip.

2) Center data : $A_c = A - \vec{\mu}$

$$\vec{\mu} = (\mu_x, \mu_y) \quad \mu_x = \frac{0+1+2+3}{4} = 1.5$$

$$\mu_y = \frac{0.5+1.1+2.3+2.5}{4} = 1.6$$

X	Y
0	0.5
1	1.1
2	2.3
3	2.5

$$- [1.5, 1.6] \leftarrow \vec{\mu}$$

A_c

X_c	Y_c
-1.5	-1.1
-0.5	-0.5
0.5	0.7
1.5	0.9

3) Compute $\Sigma = \frac{1}{(N-1)} \cdot (A_c \cdot T) @ (A_c)$

$$\frac{1}{3} \begin{bmatrix} -1.5 & -0.5 & 0.5 & 1.5 \\ -1.1 & -0.5 & 0.7 & 0.9 \end{bmatrix} @ \begin{bmatrix} -1.5 & -1.1 \\ -0.5 & -0.5 \\ 0.5 & 0.7 \\ 1.5 & 0.9 \end{bmatrix}$$

$$= \begin{bmatrix} 1.67 & 1.2 \\ 1.2 & 0.92 \end{bmatrix} \Sigma$$

4) e-vals, e-vecs of Σ
 $\lambda_1 \quad \lambda_2$
eigen values : $[2.55, 0.04]$

eigen vectors : $\begin{bmatrix} 0.81 & -0.59 \\ 0.51 & 0.81 \end{bmatrix}$

$\overrightarrow{e_1} \quad \overrightarrow{e_2}$

$\overrightarrow{\lambda_1 \cdot e_1} \quad \overrightarrow{\lambda_2 \cdot e_2}$

5) Sort e-vecs/e-vals high → low according to e-vals