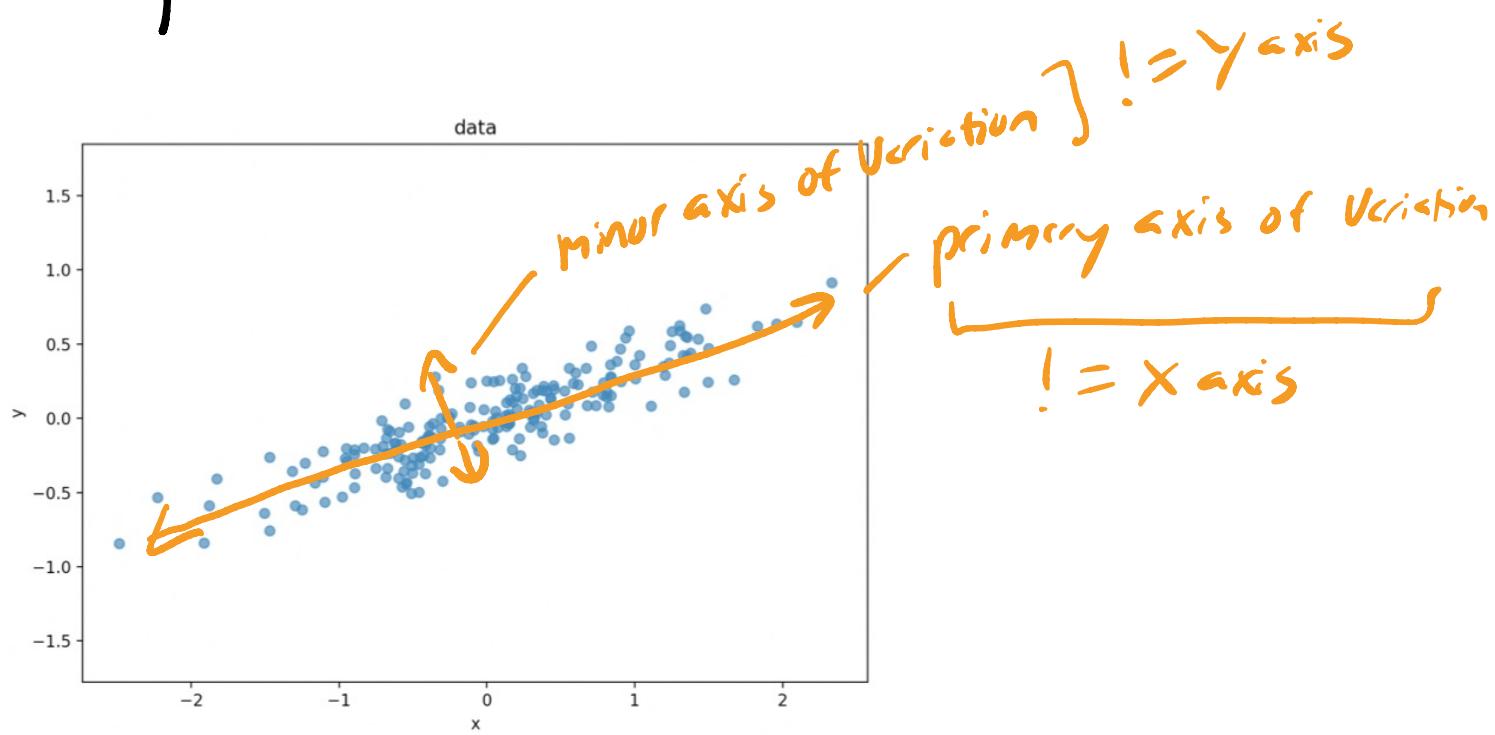
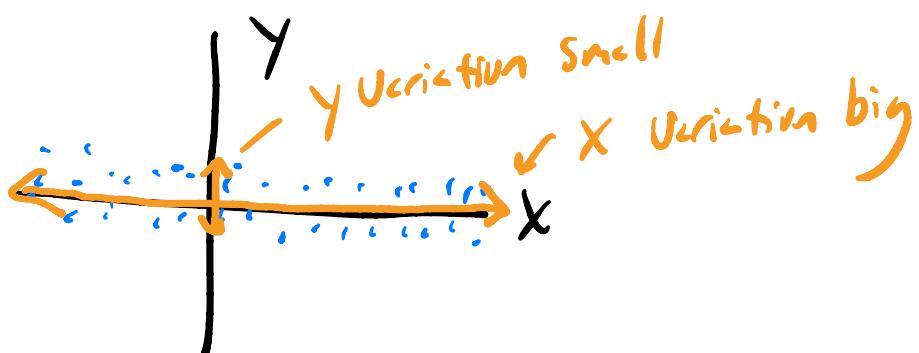
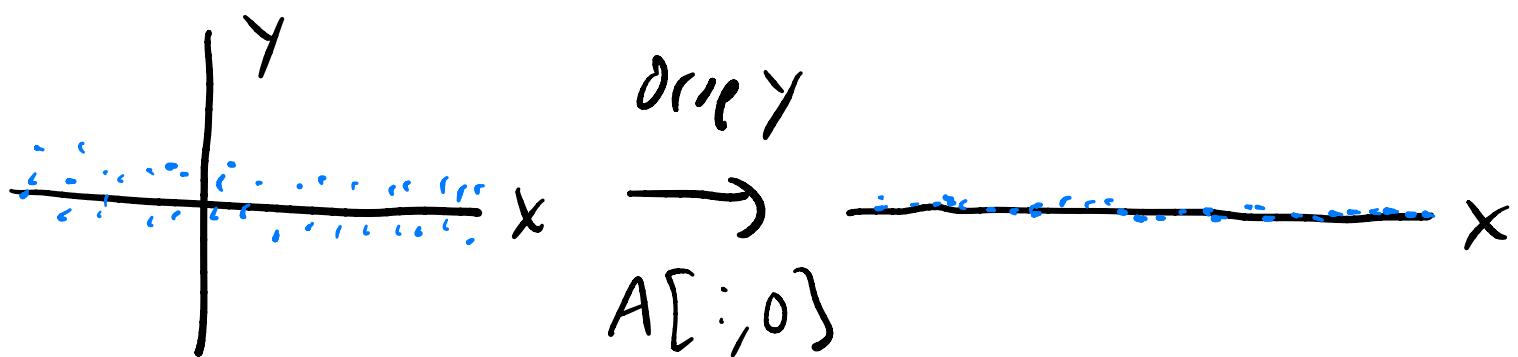
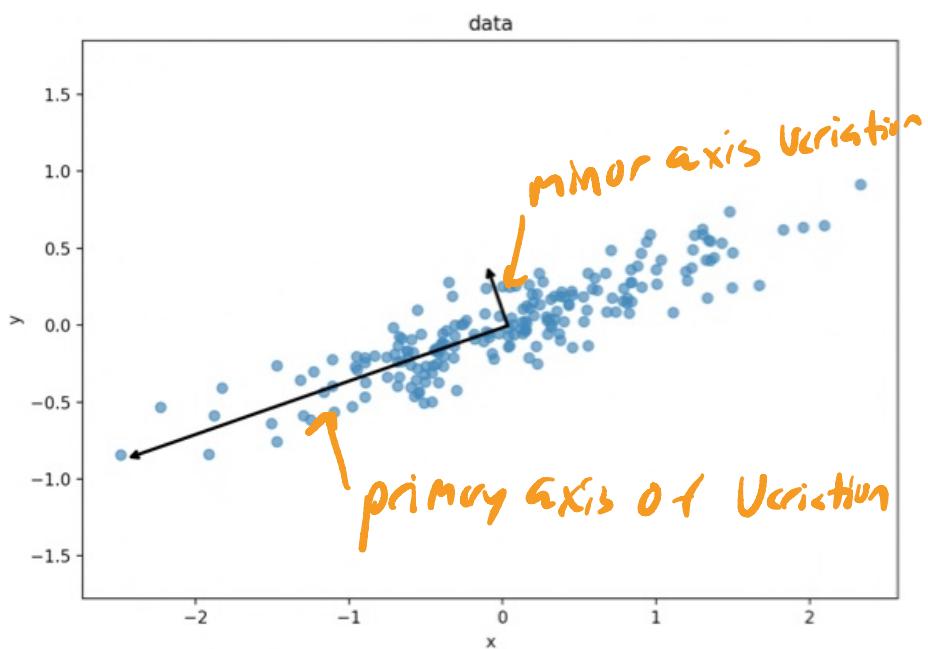
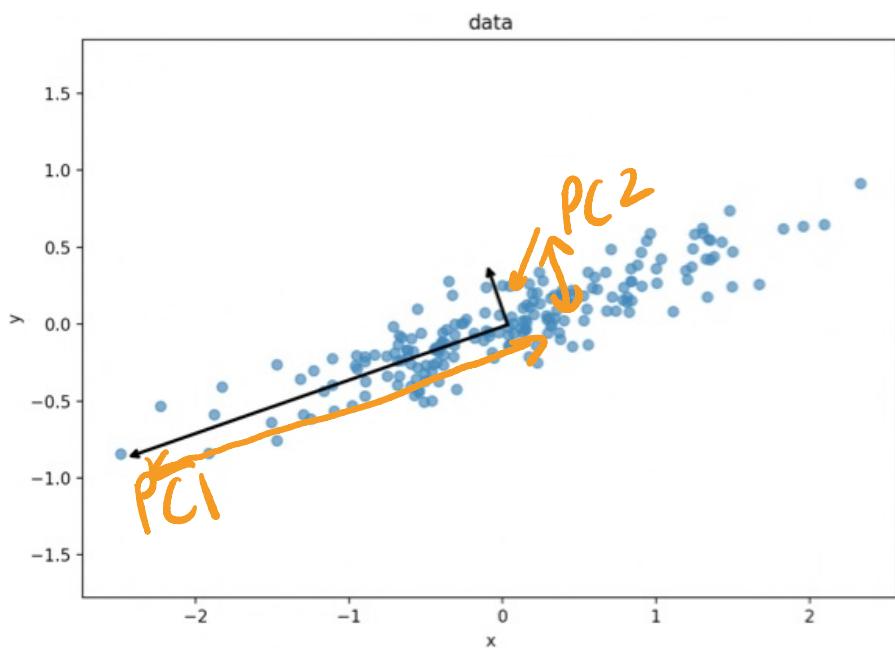


# Principal component analysis (PCA)





Principal Components (PCs): intrinsic axes of Variation in data

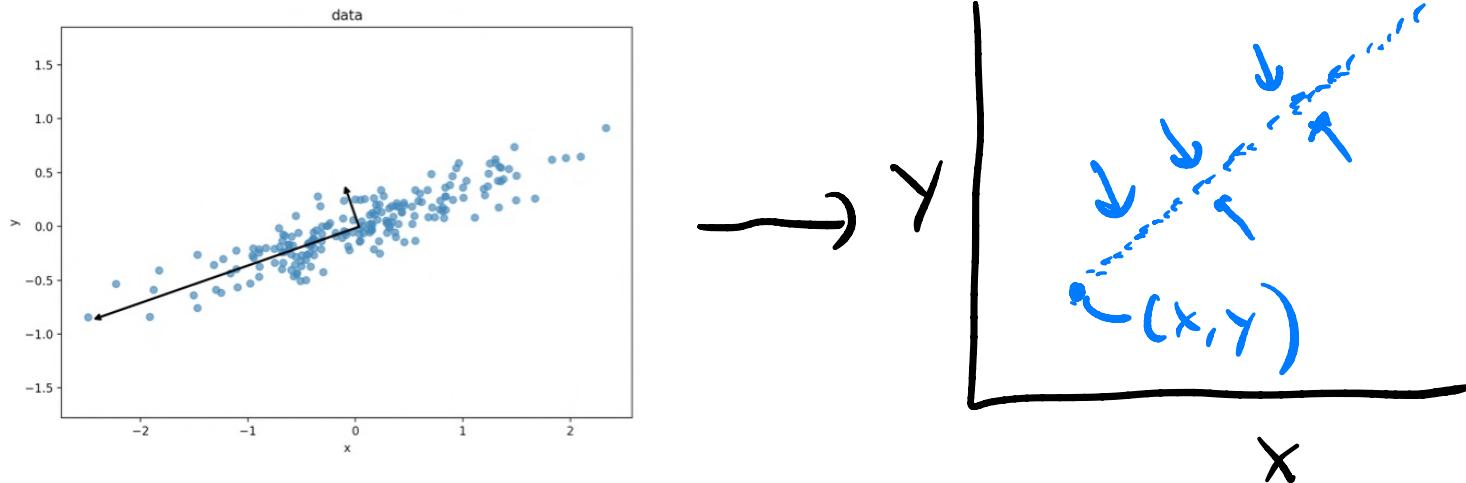


Smaller the PC number  
bigger the axis  
of Variation.

Goal of PCA

We want to drop PCs  
intrinsic axes of Variation  
that are insignificant / small

$\Rightarrow$  Do not want to drop data  
Variables



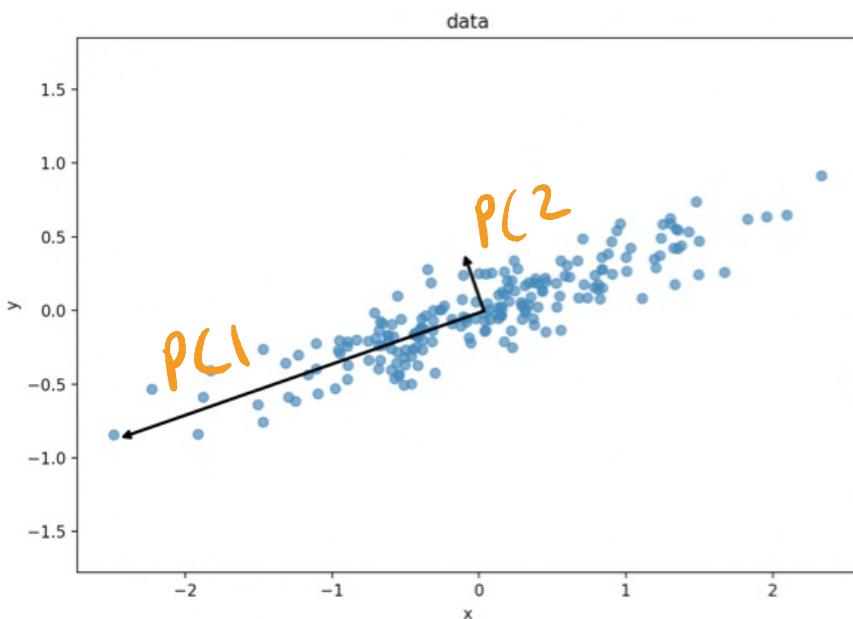
## Lecture 18 ] PCA (continued)

principal component

data variable

$$\text{PC1} ! = X$$

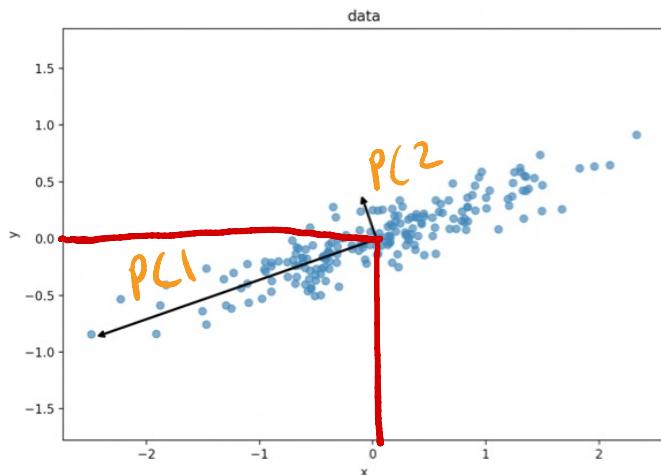
$$\text{PC2} ! = Y$$



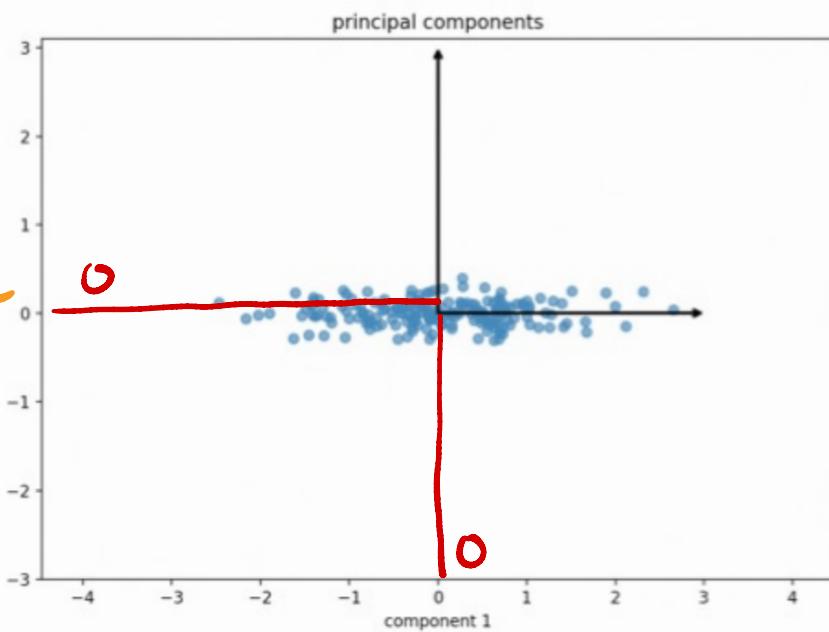
\* Want to rotate data so that PC1 becomes "x axis" and PC2 becomes "y axis"

Do that with a translation and multiplying by a rotation matrix.

$A_c$



Rotate



rotated data

PC1

rotation  
matrix

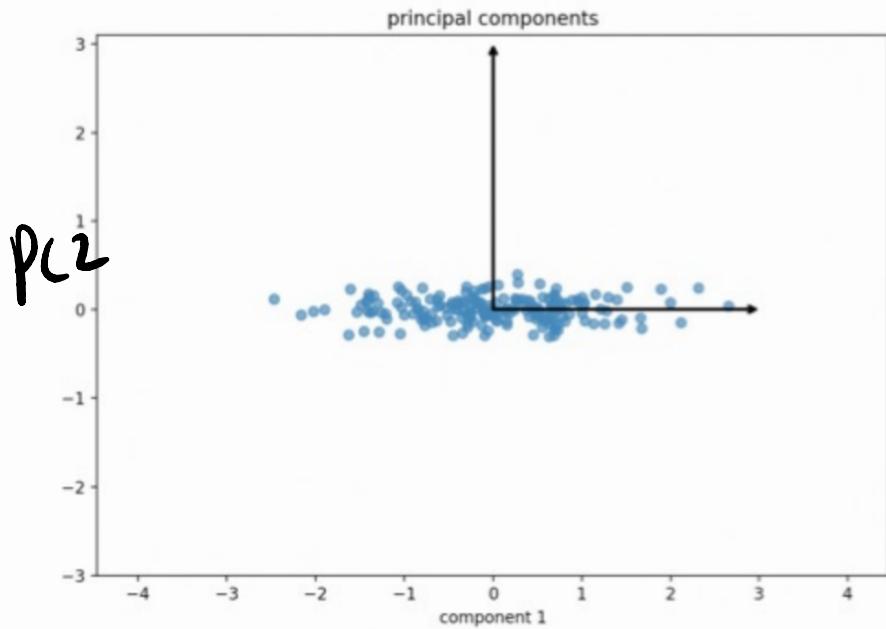
centered data matrix  
 $A_c = A - \bar{\mu}$   
 translation

$$\hat{A}_c = (P.T @ A_c).T$$

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

equal

$\hat{A}_c$ :



drop PC 2

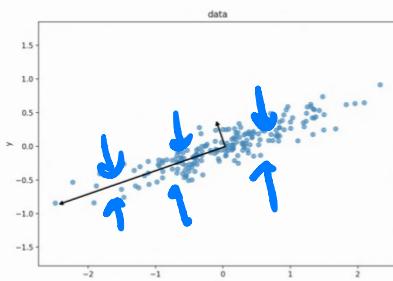
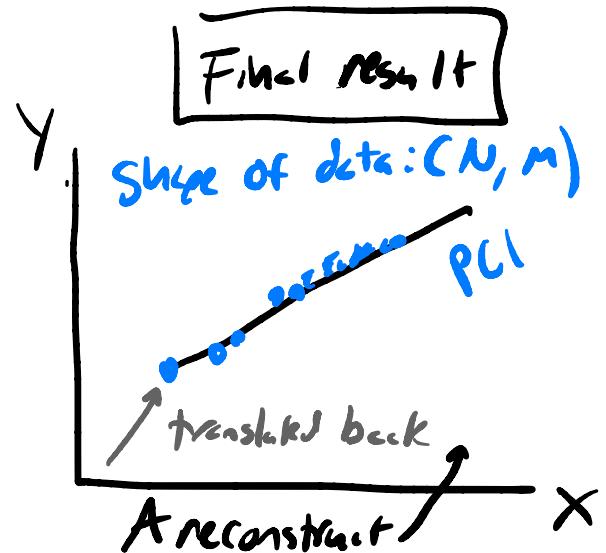
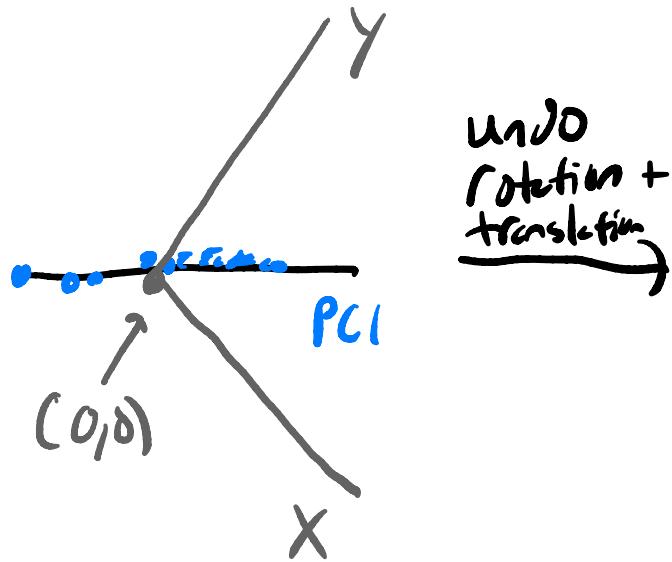
$\hat{A}_c[:, \delta]$

PC1

PC1

⇒ Can get back to original data variables

— reconstruct data — get our "axes" to  
be original data vars again



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

↑  
rotated data  
after PC2 removed

original  
data means

Same ↓

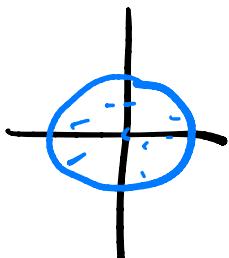
$$= \hat{A}_c @ \hat{P} . \bar{T} + \vec{\mu}$$

↑  
uncenter  
data

$A_c \rightarrow A$

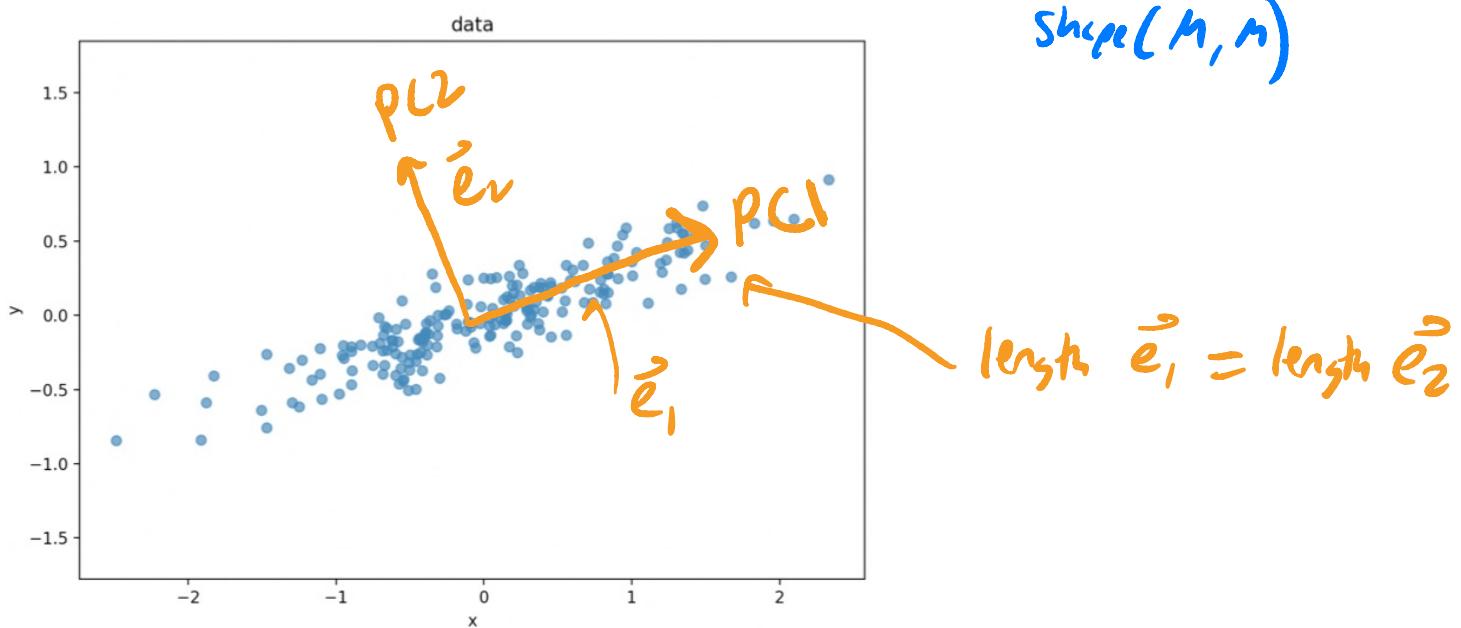
Qn What is the rotation matrix  $P$  ?

If we have  $\sum$  covariance matrix,  
shape here:  $(2, 2)$



$$\sum = \begin{bmatrix} 3 & 0 \\ 0 & 3 \end{bmatrix}$$

Shape  $(M, M)$



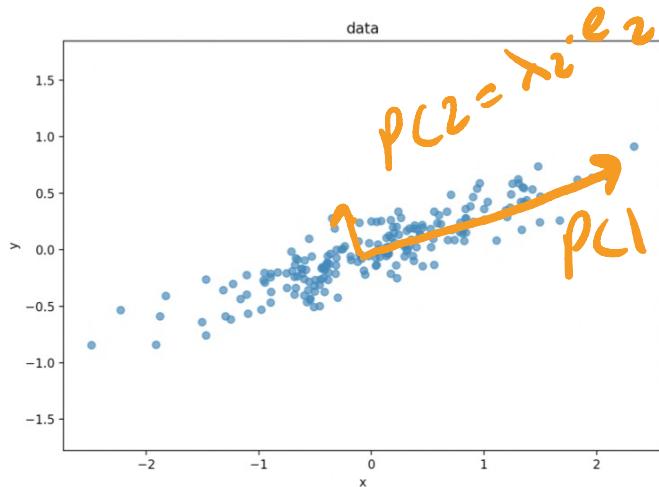
Eigenvectors of  $\Sigma$ : Vectors that tell us the direction of the PCs but not the amount of Variation in the PCs.

e.g. # eigenvectors = # data vars =  $n$   
 $\Rightarrow$  we have 2.

$$\vec{e}_1, \vec{e}_2$$

Eigenvalues of  $\Sigma$ : Scalars/floats — one number per eigenvector / data variable  
 $\lambda_1 = 2.5$        $\lambda_2 = 1.2$   
 $e\_vals = [\lambda_1, \lambda_2]$

tell us the amount of Variation in each PC direction



In numpy:

$$e\_vals, e\_vecs = \text{np.linalg.eig}(\Sigma)$$

P  
cols of P

$$P = [\vec{e}_1, \vec{e}_2, \vec{e}_3, \dots, \vec{e}_m]$$

## Strategy for removing insignificant PCs:

$$M = 4$$

not much variation  
in this PC → maybe safe to toss out

How much variation  
in this PC

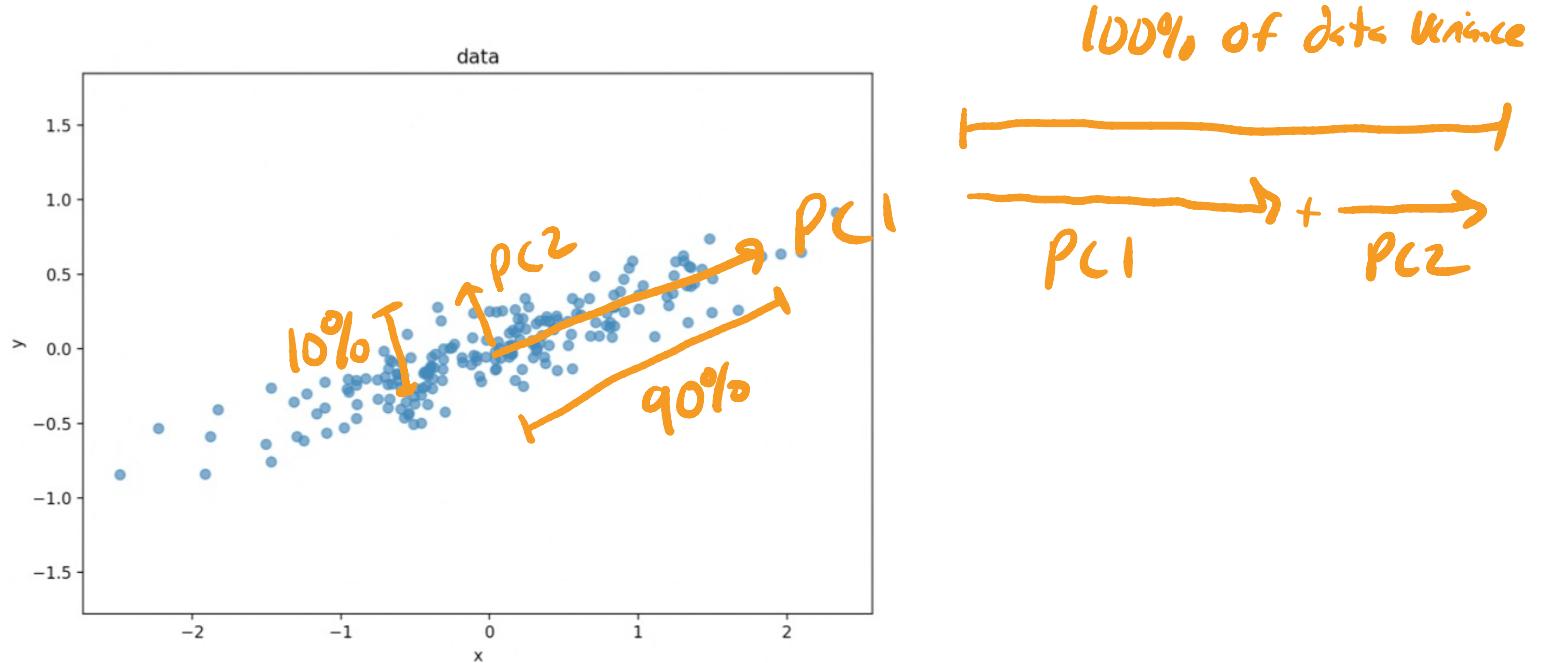
$$e\_vals = [0.1, 3, 0.7, 5]$$

→ Sort e\_vals high → low:  
  → maybe good to drop  
PC3 + PC4.

$$e\_vals = [5, 3, 0.7, 0.1]$$

↑  
PC1      ↑  
PC2      ↑  
PC3      ↑  
PC4

## Proportion Variance accounted for by each PC:



$$\text{total var} = \sum_{i=1}^m \lambda_i = 5 + 3 + 0.7 + 0.1 \\ = \underline{\underline{8.8}}$$

$$\text{prop-var} = \left[ \frac{5}{8.8}, \frac{3}{8.8}, \frac{0.7}{8.8}, \frac{0.1}{8.8} \right]$$

$\overbrace{\hspace{100px}}^{57\%}$ 
 $\overbrace{\hspace{100px}}^{34\%}$ 
 $\overbrace{\hspace{100px}}^{8\%}$ 
 $\overbrace{\hspace{100px}}^{1\%}$

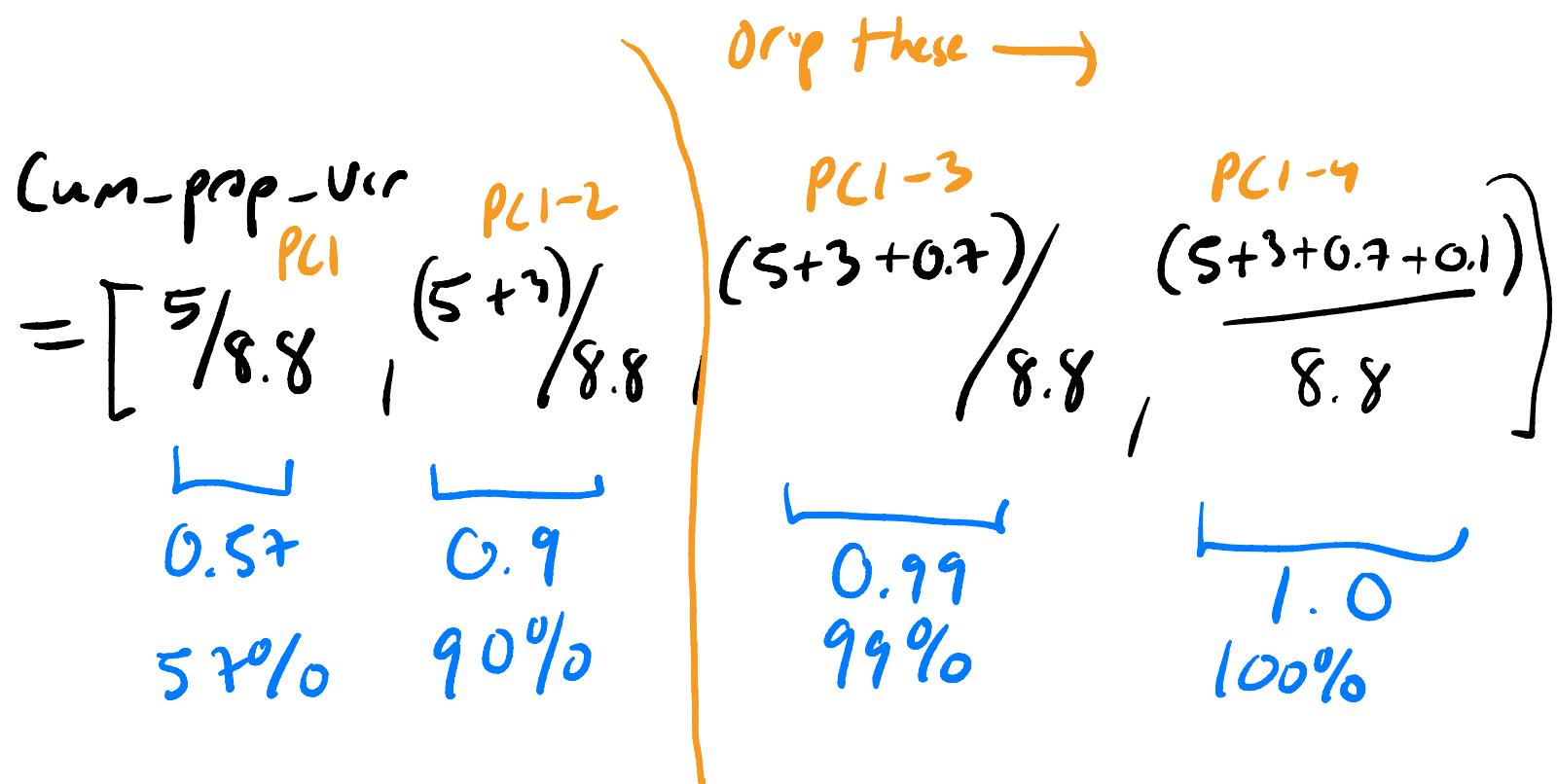
Of total Variance .

$\overbrace{\hspace{600px}}$   
Sums to 100% in total

Define Variance Cut off threshold :

Want to keep 90% of the total variance.

⇒ helpful to compute Cumulative Variance accounted for by PCs.



Goal: Keep at least 90% : keep PCI<sub>1</sub>, PCI<sub>2</sub>

Drop PCI<sub>3</sub>, PCI<sub>4</sub>