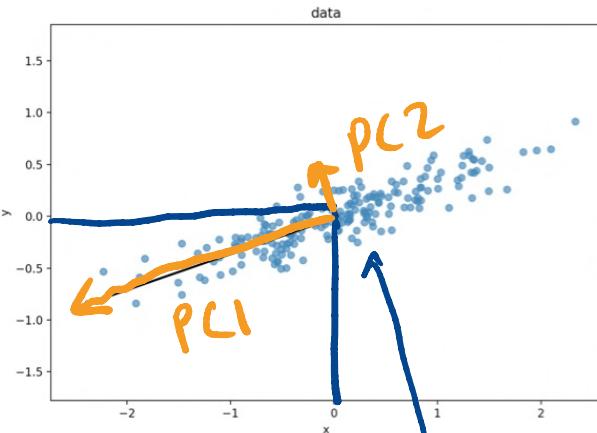


Lecture 19

A_C

data Space

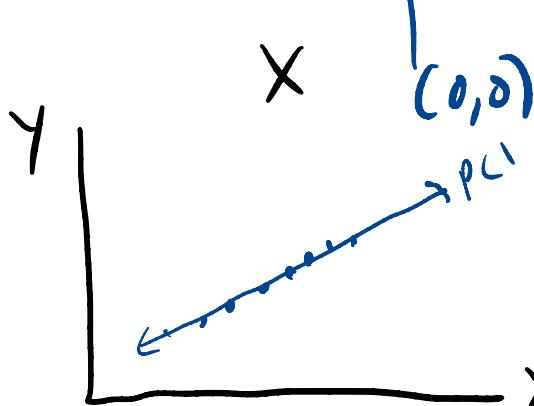
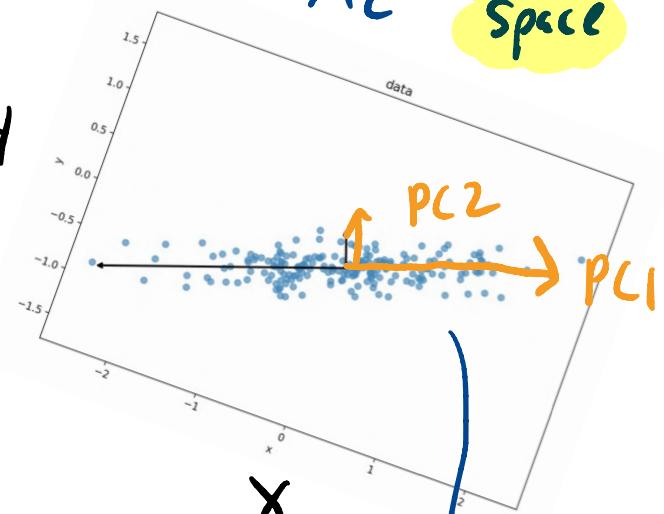


PCA: rotation



\hat{A}_C

PCA
Space



$undo R + T$

$reconstruction$

drop PC_2 :

$\hat{A}_C[:, 0]$

Data space \rightarrow PCA space:

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

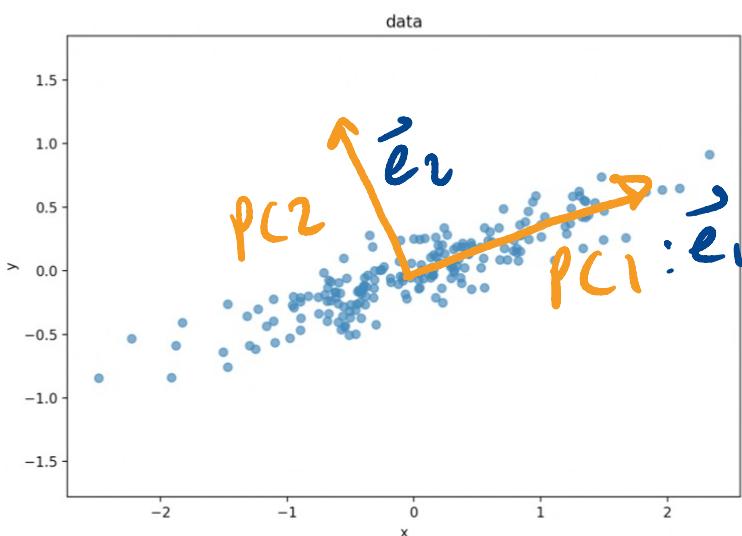
P : eigen vectors of Σ (covariance matrix)

of the centered data A_c

\downarrow eigenvectors
 \downarrow columns of P

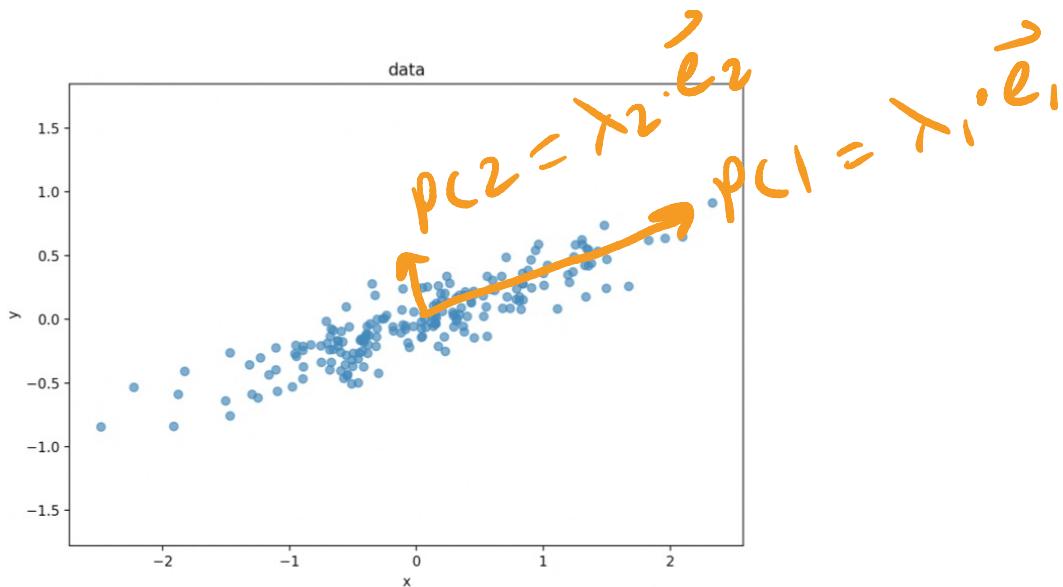
$$P = [\vec{e}_1, \vec{e}_2, \dots, \vec{e}_M]$$

\uparrow
all unit length (1)



Eigenvalues: $[\lambda_1, \lambda_2, \dots, \lambda_m]$

↓
Scalars



In numpy:

$e_vals, e_vecs = \text{np.linalg.eig}(\Sigma)$

Strategy for dropping insignificant PCs

$e_vals = [0.1, 3, 0.7, 5]$

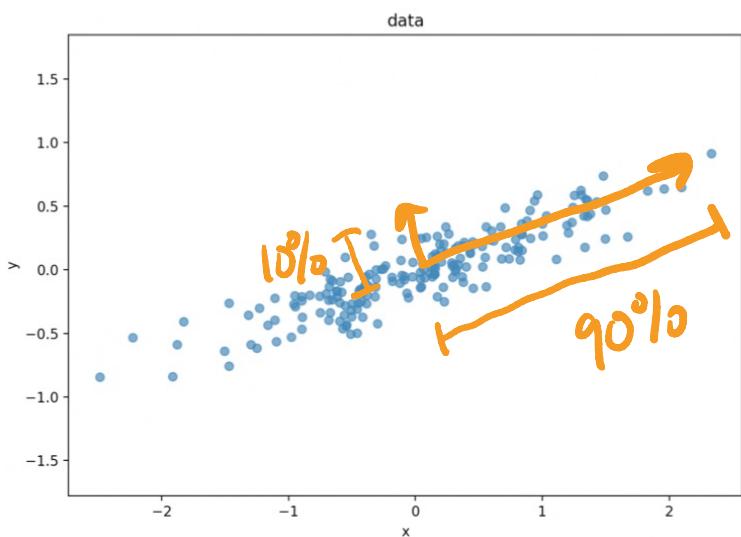
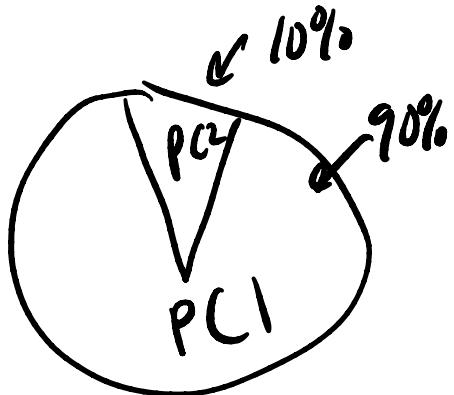
$M=4$
 Small Variance \rightarrow Maybe drop these PCs.

Sort e_vals list high-to-low:

$PC_1 \quad PC_2 \quad PC_3 \quad PC_4$
 $[5, 3, 0.7, 0.1]$

It's up to us to decide how many PCs to drop.

Proportion of total Variance accounted for by each PC:



$$\begin{matrix} \text{PC1} & \text{PC2} & \text{PC3} & \text{PC4} \\ [5, 3, 0.7, 0.1] \end{matrix}$$

$$\text{total Variance} = 5 + 3 + 0.7 + 0.1 = \underline{\underline{8.8}}$$

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

Underneath the fractions, blue brackets indicate the percentage of variance explained by each component:
PC1: 57%
PC2: 34%
PC3: 8%
PC4: 1%

Cumulative Variance accounted for by PCs:

E.g. Cutoff / threshold: Want to keep $\geq 90\%$ data variance.

Cum-prop-Var =

$$\begin{aligned} & \underbrace{\text{PC1}}_{[5/8.8,}, \quad \underbrace{\text{PC1+PC2}}_{(5+3)/8.8}, \quad \underbrace{\text{PC1+PC2+PC3}}_{(5+3+0.7)/8.8}, \quad \underbrace{\text{all PCs}}_{(5+3+0.7+0.1)/8.8} \\ & = [0.57, 0.9, \boxed{0.99, \underbrace{\text{PC3}}_{\substack{\text{keep} \leftarrow \\ \text{toss}}}}, 1.0] \end{aligned}$$

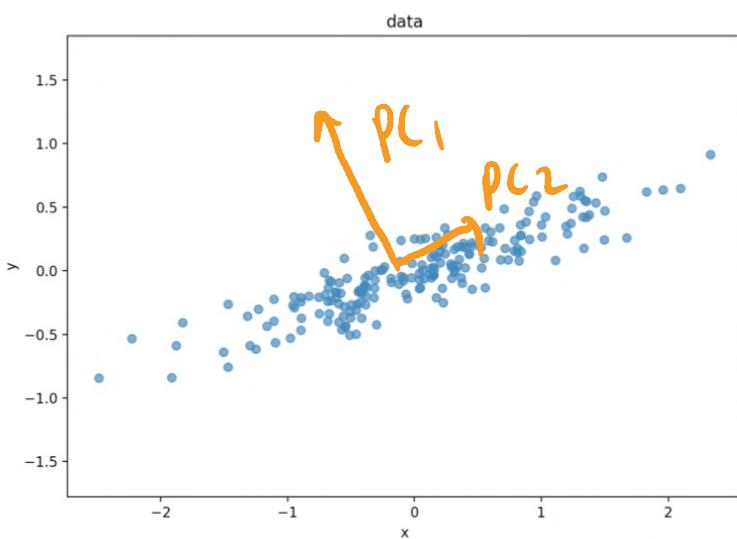
Summary of PCA

- 1) [optional] Normalize data.
 - per-variable : useful (maybe) if data have vastly different units (e.g. miles, cm)
 - global : same units, variables comparable. (e.g. miles, miles)
- 2) Center data : $A_c = A - \vec{\mu}$
- 3) $\Sigma = \frac{1}{(N-1)} A_c^T @ A_c$

P : Shape: (M, M)
- 4) Compute $\overbrace{e\text{-vecs}, e\text{-vals}}$ — via numpy
- 5) Sort high \rightarrow low for e-vals

* also sort **columns** of e-vecs (P)
in same way/order as the e-vals

problem (no sorting of PC's):



Don't want this

6a) Compute prop-var by each PC.

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

$$\Rightarrow \left[\frac{3.2}{4.33}, \frac{1.1}{4.33}, \frac{0.03}{4.33} \right]$$

6b) Compute cum-prop-var:

$$\left[\frac{3.2}{4.33}, \frac{(3.2+1.1)}{4.33}, \frac{(3.2+1.1+0.03)}{4.33} \right]$$

6c) Threshold: drop some # of PCs to preserve X% of info/var in data

keep K PCs (after dropping)

7) Project data to PCA Space:

$$\hat{A}_c = \underbrace{A_c}_{(N, M)} @ \begin{matrix} \hat{P} \\ \text{cols for top } k \text{ PCs} \end{matrix}$$

Shape: (N, M) (M, K)

$\underbrace{(N, K)}$
 $\underbrace{\text{Reduced}}$
 $\text{Dims } M \rightarrow K$

8) [optional] Reconstruct data. PCA Space \rightarrow data Space

if no normalization:

$$A_{\text{reconstruct}} = \underbrace{\hat{A}_c}_{(N, K)} @ \begin{matrix} \hat{P}^T \\ + \bar{\mu} \end{matrix}$$

$\underbrace{(K, M)}$
 $\underbrace{(N, M)}$

$\underbrace{\text{uncenter data}}_{(1, M)}$

if we normalized:

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

Range of original data

