

Summary of PCA

- 1) [optional] Normalize data if appropriate so that each data value in range $[0, 1]$.

Note: Can be done in vectorized way without matrix multiplication with one line of code.

- **Normalize per-Variable (separately)** may be appropriate if data variables have many different units (e.g. miles & cms)
⇒ You lose relative scales among vars
- **normalize globally (together)** may be appropriate if variables have similar units (e.g. inches)
⇒ Won't distort relative scales among vars

- 2) Center data by subtracting each variable by its mean $\vec{\mu}$

$$A_c = A - \vec{\mu}$$

- 3) Compute covariance matrix: $\Sigma = \frac{1}{N-1} A_c \cdot T \in A_c$

- 4) Compute eigenvectors P and eigenvalues via spectral decomposition of Σ — via Numpy.

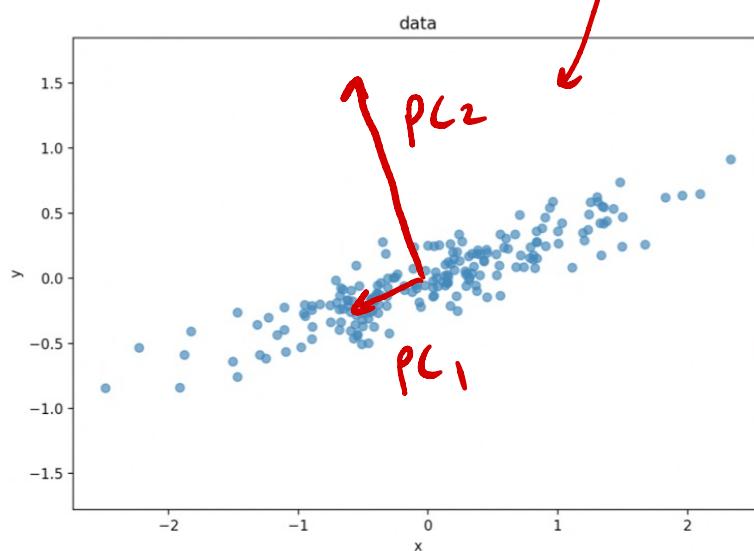
Shape: (M, M)

5) Sort e-vals high \rightarrow low

e.g. $[3.2, 1.1, 0.03]$

Sort P **(columns)** in the same order otherwise PC lengths & directions get Scrambled!

problem if we don't
sort P + evals
the same way!



6a) Compute proportion Variance accounted for by each PC

e.g. total = $3.2 + 1.1 + 0.03 = 4.33$

$\Rightarrow [3.2/4.33, 1.1/4.33, 0.03/4.33]$

6b) Compute cumulative Variance accounted for by PC_i :

e.g. $[3.2/4.33, (3.2+1.1)/4.33, (3.2+1.1+0.03)/4.33]$

6c) Threshold PCs based on desired amount of info/variance to keep. We keep K PCs.

e.g. $\geq 90\%$. If $PC1 + PC2 = 90\%$, toss out $PC3$.

7) project data into PCA Space [rotating data to align coordinate axes to PCs] \hat{P} is P with certain PCs / cols dropped

$$\hat{A}_c = \underbrace{\hat{A}_c}_{(N, K)} @ \underbrace{\hat{P}}_{(N, M)} \quad \text{keep top } k \text{ PCs}$$

8) [optional] Reconstruct data [rotating data after projecting away insignificant PCs to align coordinate axes with original data variables]

if you didn't normalize:

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

if you normalized:

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

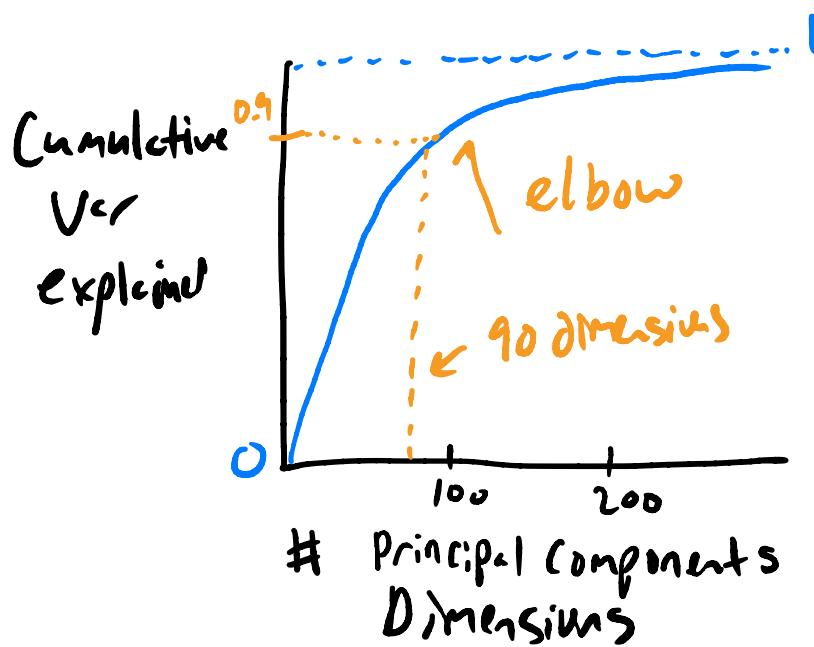
Original data range in each variable

Elbow plot

Step 6b was simple with 2D data, but generally get M values — Cumulative proportion variance explained by top K dimensions — $K \leq M$.

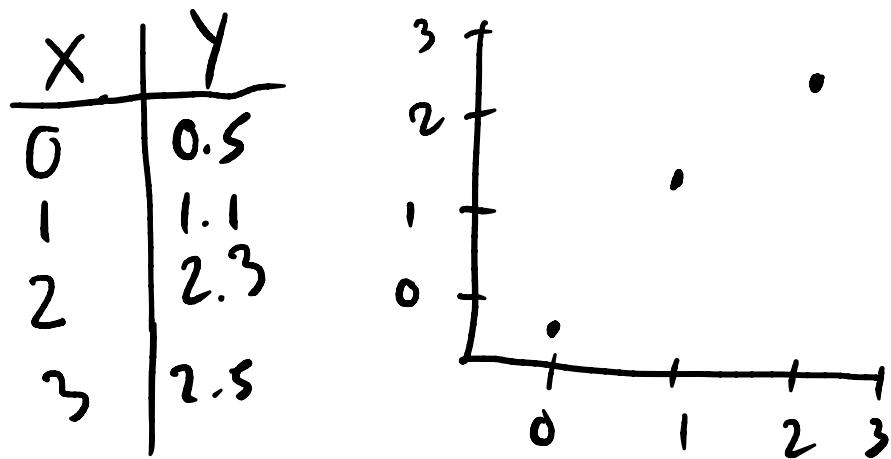
⇒ useful to plot to discover reasonable threshold on % Variance retained.

Elbow plot:



plot to determine point of "diminishing gains" visually
— cut off point of # PCs that explain good chunk of variance before plateauing

Example of PCA workflow by hand



1) Normalize (skip — assume data homogeneous)

2) Center data. $\mu_x = \frac{0+1+2+3}{4} = \frac{6}{4} = \frac{3}{2} = 1.5$

$$\vec{\mu} = (\mu_x, \mu_y) \quad \mu_y = \frac{0.5 + 1.1 + 2.3 + 2.5}{4} = \frac{6.4}{4} = 1.6$$

X_c	Y_c
$0 - 1.5$	$0.5 - 1.6$
$1 - 1.5$	$1.1 - 1.6$
$2 - 1.5$	$2.3 - 1.6$
$3 - 1.5$	$2.5 - 1.6$

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

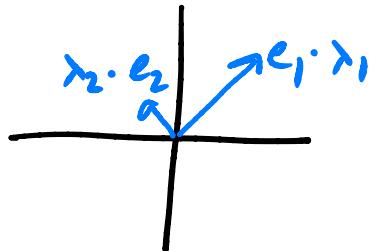
3) Compute covariance matrix : $\frac{A c^T A}{N-1}$

$$\frac{1}{3} \begin{bmatrix} -1.5 & -0.5 & 0.5 & 1.5 \\ -1.1 & -0.5 & 0.7 & 0.9 \end{bmatrix} \begin{bmatrix} -1.5 & -1.1 \\ -0.5 & -0.5 \\ 0.5 & 0.7 \\ 1.5 & 0.9 \end{bmatrix} = \begin{bmatrix} 1.67 & 1.2 \\ 1.2 & 0.92 \end{bmatrix}$$

4) Compute eigenvals/vecs of Σ .

Eigenvalues: $= \left[\frac{\lambda_1}{2.55}, \frac{\lambda_2}{0.04} \right]$

Eigenvectors: $P = \begin{bmatrix} 0.81 & -0.59 \\ 0.59 & 0.81 \end{bmatrix}$



Unit vectors (length = 1)

b/c Σ is Symmetric — See Spectral theorem coming soon

5) Sort eigenvalues high-to-low: [2.55, 0.04]

6) a) Compute proportion Variance accounted for

by each principal component (PC).

eigenvalues

$$\text{prop-var}_i = \frac{\lambda_i}{\sum_{j=1}^m \lambda_j} = \left[\frac{2.55}{2.55+0.04}, \frac{0.04}{2.55+0.04} \right] = [0.98, 0.02]$$

\vec{e}_1 accounts for 98% of variance