# PRINCIPAL COMPONENT ANALYSIS

$\downarrow$

Dimensionality reduction technique



$p = (x_1, x_2)$

$P \in \mathbb{R}^{N \times D}$   $K < D$

$\overline{P} \in \mathbb{R}^{N \times K}$
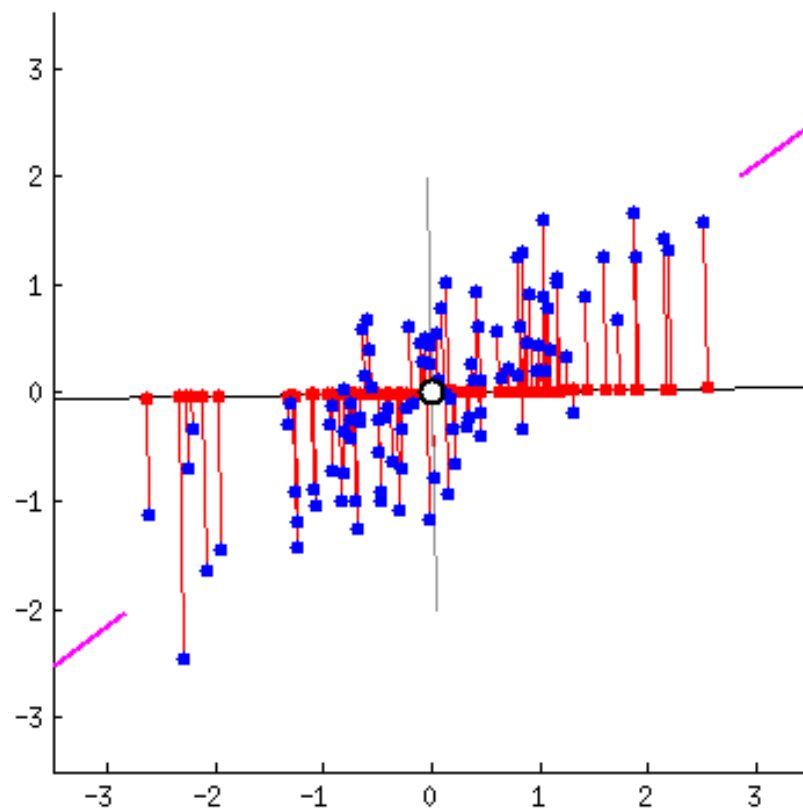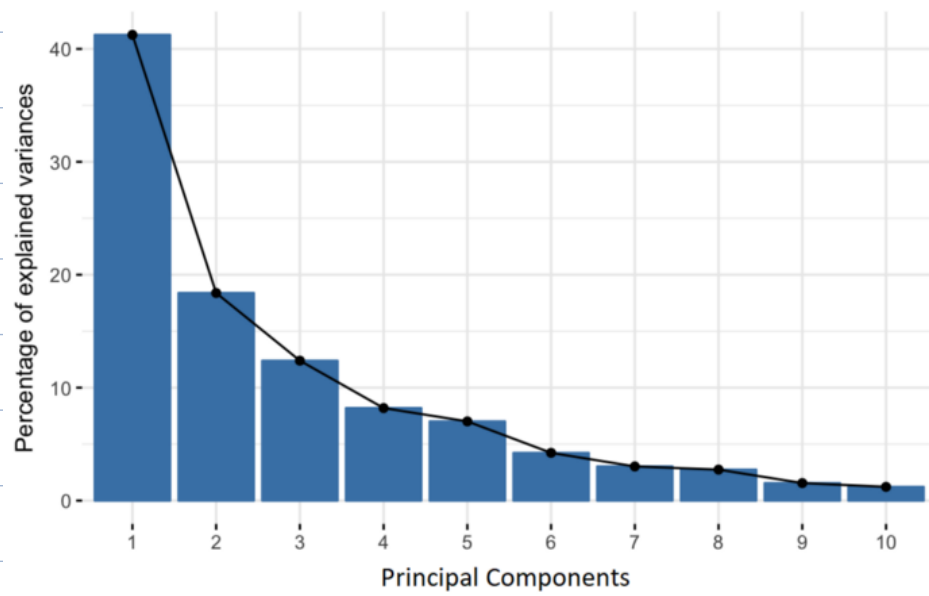
find the principle components i.e. directions along which the variance of the data is MAXIMAL

larger variance $\longrightarrow$ larger dispersion of data

$\downarrow$

larger information

# ALGORITHM

Consider a dataset with N data points

Each datapoint $\overline{p_i} \in \mathbb{R}^D$     $1 \leq i \leq N$     $(P_{N \times D})$

i.e   $\overline{p_i} = (x_1, x_2, x_3, \ldots x_D)$     $x_i \in \mathbb{R}$

① Standardize the data

$$x_i^1 = \frac{x_i - E[x_i]}{\sqrt{var(x_i)}}$$

② Compute the Covariance matrix $C_{D \times D}$

$$C = \begin{bmatrix} Cov(x_1, x_1) & \cdots & Cov(x_1, x_2) & \cdots & Cov(x_1, x_D) \\ \vdots & & \vdots & & \vdots \\ \vdots & & \vdots & & \vdots \\ Cov(x_D, x_1) & \cdots & Cov(x_D, x_2) & \cdots & Cov(x_D, x_D) \end{bmatrix}$$

$D \times D$

$\longrightarrow$ Symmetric

③ Compute the eigenvector and eigen values of the Cov matrix

$$C = V_{D \times D} \; \Sigma_{D \times D} \; V^T_{D \times D}$$

OR

$$Cv = \lambda v$$

$$\left( C - \lambda I \right) v = 0$$
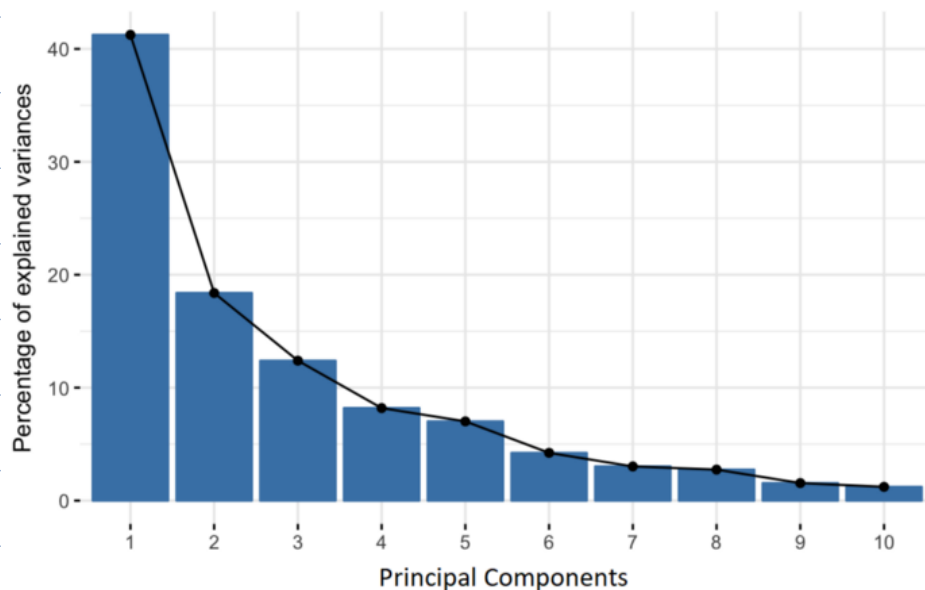
$$\boxed{\det \left( C - \lambda I \right) = 0}$$

np. linalg. eig

④ Select Top K Eigenvectors are the principal components

$$X = \begin{bmatrix} | & | & & | \\ v_1 & v_2 & \cdots & v_k \\ | & | & & | \end{bmatrix} \qquad v_i \in \mathbb{R}^D$$

$$D \times K$$

$(\lambda_1) \quad (\lambda_2) \quad\quad (\lambda_k)$

Such that $\lambda_1 > \lambda_2 > \lambda_3 \cdots > \lambda_k$



⑤ Recast the data along principal components

$$P_{N \times D} \, X_{D \times K} = \bar{P}_{N \times K}$$

$$\boxed{K \leq D}$$