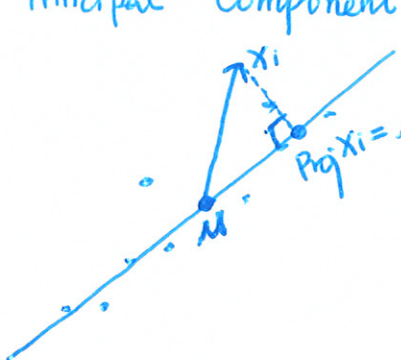


Principal Component Analysis

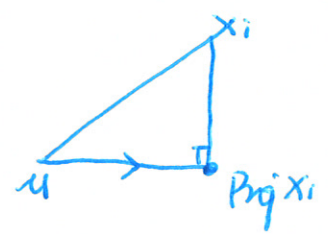


Approximate data $x_1, \dots, x_N \in \mathbb{R}^p$ by the plane

$$f(\lambda) = \mu + v_q \lambda$$

$$\mu \in \mathbb{R}^p \quad v_q \in \mathbb{R}^{p \times q}$$

$$\begin{aligned} \min_{\text{plane}} & \sum_{i=1}^N \|x_i - \text{Proj } x_i\|^2 \\ = \min_{\mu, v_q} & \sum_{i=1}^N \|x_i - \mu - v_q v_q^T (x_i - \mu)\|^2 \\ = \min_{\mu, v_q} & \sum_{i=1}^N \|x_i - \mu\|^2 - \|v_q v_q^T (x_i - \mu)\|^2. \end{aligned}$$



Step 1: minimize $\sum_{i=1}^N \|x_i - \mu\|^2 \Rightarrow \mu = \frac{1}{N} \sum_{i=1}^N x_i$ Mean.

Step 2. Find v_q to maximize $\sum_{i=1}^N \|v_q v_q^T (x_i - \mu)\|^2$.

$$\begin{aligned} \text{Notice: } \|v_q v_q^T (x_i - \mu)\|^2 &= \|v_q^T (x_i - \mu)\|^2 \\ &= \text{Trace} [v_q^T (x_i - \mu)(x_i - \mu)^T v_q] \end{aligned}$$

$$\begin{aligned} \max_{v_q} & \sum_{i=1}^N \text{Trace} [v_q^T (x_i - \mu)(x_i - \mu)^T v_q] \\ \Rightarrow \max_{v_q} & \text{Trace} \left[v_q^T \left(\sum_{i=1}^N (x_i - \mu)(x_i - \mu)^T \right) v_q \right] \end{aligned}$$

Let $A = \sum_{i=1}^N (x_i - \mu)(x_i - \mu)^T$ $A^T = A$.

max. $\text{Trace}(V_q^T A V_q) = \lambda_1 + \lambda_2 + \dots + \lambda_q$.

top q eigenvalues of A .

if $V_q = \begin{bmatrix} \vec{u}_1 & \vec{u}_2 & \dots & \vec{u}_q \end{bmatrix}$

the eigenvectors associated with $\lambda_1 \lambda_2 \dots \lambda_q$.

Equivalently: Let $X = \begin{bmatrix} x_1 - \mu & x_2 - \mu & \dots & x_n - \mu \end{bmatrix}$

SVD of X : $X = U D V^T$ → right singular vectors.

↙
left singular vectors

$V_q = \begin{bmatrix} \vec{u}_1 & \dots & \vec{u}_q \end{bmatrix}$

the left singular vectors associated with the largest q singular values.