

For a sample point $x^{(i)}$, we have

$$x^{(i)} = \underbrace{\left(x^{(i)T} u^{(1)}\right) u^{(1)} + \dots + \left(x^{(i)T} u^{(d)}\right) u^{(d)}}_{= \text{Proj}_S x^{(i)} \in S} + \underbrace{\dots + \left(x^{(i)T} u^{(p)}\right) u^{(p)}}_{= \text{Proj}_{S^\perp} x^{(i)} \in S^\perp}$$

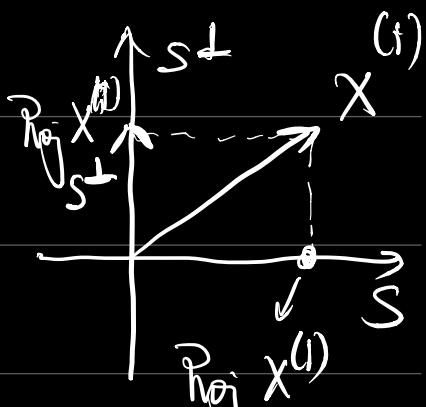
squared distance from $x^{(i)}$ to S

$$= \|\text{Proj}_{S^\perp} x^{(i)}\|^2$$

$$= \left\| \left(x^{(i)T} u^{(d+1)}\right) u^{(d+1)} + \dots + \left(x^{(i)T} u^{(p)}\right) u^{(p)} \right\|^2$$

$$\stackrel{\text{Pythagorean}}{=} \left| \left(x^{(i)T} u^{(d+1)}\right) \right|^2 + \dots + \left| u^{(i)T} u^{(p)} \right|^2$$

$$= \sum_{j=d+1}^p \left| x^{(i)T} u^{(j)} \right|^2$$



Therefore, the total squared distance from the sample points to S

$$= \sum_{i=1}^n \sum_{j=d+1}^p \left| x^{(i)T} u^{(j)} \right|^2 \quad (*)$$

(This is what is to be minimized)

We compute this in another way:

For $x^{(i)}$, Pythagorean theorem gives:

$$\|x^{(i)}\|^2 \stackrel{\text{Pythagorean}}{=} \|\text{Proj}_S x^{(i)}\|^2 + \|\text{Proj}_{S^\perp} x^{(i)}\|^2$$

$$= \|(x^{(i)\top} u^{(1)}) u^{(1)} + \dots + (x^{(i)\top} u^{(d)}) u^{(d)}\|^2$$

$$+ \|(x^{(i)\top} u^{(d+1)}) u^{(d+1)} + \dots + (x^{(i)\top} u^{(p)}) u^{(p)}\|^2$$

Pythagorean

$$\|x^{(i)\top} u^{(1)}\|^2 + \dots + \|x^{(i)\top} u^{(d)}\|^2$$

$$+ \|x^{(i)\top} u^{(d+1)}\|^2 + \dots + \|x^{(i)\top} u^{(p)}\|^2$$

$$= \underbrace{\sum_{j=1}^d \|x^{(i)\top} u^{(j)}\|^2}_{= \|\text{Proj}_S x^{(i)}\|^2} + \underbrace{\sum_{j=d+1}^p \|x^{(i)\top} u^{(j)}\|^2}_{= \|\text{Proj}_{S^\perp} x^{(i)}\|^2}$$

Insert $\sum_{j=d+1}^p \|x^{(i)\top} u^{(j)}\|^2 = \|x^{(i)}\|^2 - \sum_{j=1}^d \|x^{(i)\top} u^{(j)}\|^2$

into $(*)$:

the total squared distance =

$$\sum_{i=1}^n \left(\|x^{(i)}\|^2 - \sum_{j=1}^d |x^{(i)T} u^{(j)}|^2 \right)$$

$$= \underbrace{\sum_{i=1}^n \|x^{(i)}\|^2}_{\text{indep of } u^{(1)}, \dots, u^{(d)}} - \sum_{i=1}^n \sum_{j=1}^d |x^{(i)T} u^{(j)}|^2$$

It remains to maximize (by choosing $u^{(1)}, \dots, u^{(d)}$).

$$\sum_{i=1}^n \sum_{j=1}^d |x^{(i)T} u^{(j)}|^2 = \sum_{j=1}^d \left(\sum_{i=1}^n |x^{(i)T} u^{(j)}|^2 \right)$$

$$= \sum_{j=1}^d \left(\sum_{i=1}^n (x^{(i)T} u^{(j)})^T (x^{(i)T} u^{(j)}) \right) \text{scatter}$$

$$= \sum_{j=1}^d \left(\sum_{i=1}^n u^{(j)T} (x^{(i)} x^{(i)T}) u^{(j)} \right)$$

$$= \sum_{j=1}^d \left(u^{(j)T} \sum_{i=1}^n (x^{(i)} x^{(i)T}) u^{(j)} \right)$$

computed thus

$$\text{before} = \sum_{j=1}^d \left(u^{(j)T} X X^T u^{(j)} \right)$$

$$= u^{(1)T} (X X^T) u^{(1)} + \dots + u^{(d)T} (X X^T) u^{(d)}$$

This is the sum of d Rayleigh quotients.

To maximize $u^{(1)T} (X X^T) u^{(1)}$, we choose

$u^{(1)}$ to be a normalized eigenvector associated to the largest eigenvalue of $X X^T$.

To maximize $u^{(2)T} (X X^T) u^{(2)}$ with $u^{(2)} \perp u^{(1)}$
(because $\{u^{(1)}, \dots, u^{(d)}\}$ is an orthonormal basis of S)

we choose $u^{(2)}$ to be a normalized eigenvector associated to the second largest eigenvalue of

$$X X^T$$

:

Thm: The d -dimension subspace that best fits the sample points is spanned by the d eigenvectors associated to the d largest eigenvalues of XX^T .

3.3 PCA: Optimization Viewpoint:

Let $x^{(1)}, \dots, x^{(n)} \in \mathbb{R}^p$ be sample points and $X = [x_1^{(1)}, \dots, x_1^{(n)}]$ be the sample matrix.

Q: find a matrix $A \in \mathbb{R}^{p \times n}$ such that

$$\min_A \|X - A\|_F \quad \text{s.t. } \text{rank}(A) = d \quad (< p)$$

(i.e. find the optimal low-rank approximation of rank d)

3.3.1 Math Prep.

Recall: Given two matrices $A = (A_{ij})$, $B = (B_{ij})$

$\in \mathbb{R}^{m \times n}$, their Frobenius inner product is

$$(A, B)_F \stackrel{\text{def}}{=} \text{tr}(A^T B) = \text{tr}(B^T A) = \sum_{i=1}^m \sum_{j=1}^n A_{ij} B_{ij}$$

(because

$$A^T B = \begin{pmatrix} A_{11} & A_{21} & \cdots & A_{m1} \\ A_{12} & A_{22} & \cdots & A_{m2} \\ \vdots & \vdots & \ddots & \vdots \\ A_{1n} & A_{2n} & \cdots & A_{mn} \end{pmatrix} \begin{pmatrix} B_{11} & B_{12} & \cdots & B_{1n} \\ B_{21} & B_{22} & \cdots & B_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ B_{m1} & B_{m2} & \cdots & B_{mn} \end{pmatrix}$$

$$A_{11}B_{11} + A_{21}B_{21} + \cdots + A_{m1}B_{m1}$$

$$A_{12}B_{12} + A_{22}B_{22} + \cdots + A_{m2}B_{m2}$$

*

*

$$A_{1n}B_{1n} + A_{2n}B_{2n} + \cdots + A_{mn}B_{mn}$$

The induced Frobenius norm of A is

$$\|A\|_F^2 \stackrel{\text{def}}{=} \text{tr}(A^T A) = \sum_{i=1}^m \sum_{j=1}^n |A_{ij}|^2$$