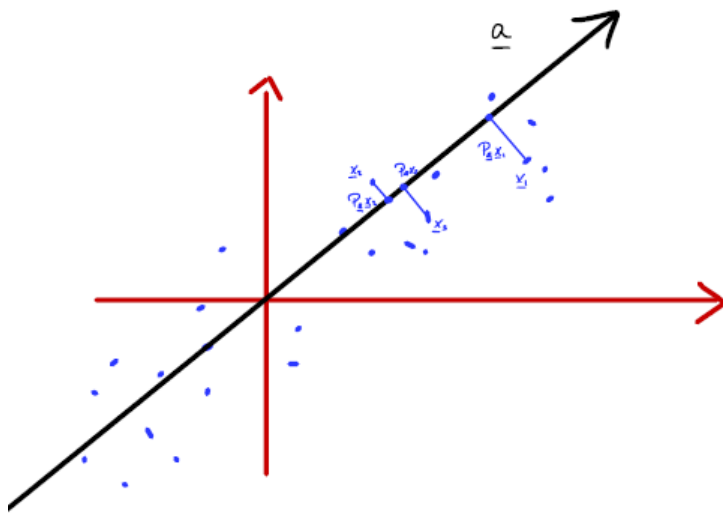


## Lecture 7. Introduction to Singular Value Decomposition

Observe  $\underline{x}_1, \underline{x}_2, \dots, \underline{x}_p \in \mathbb{R}^n$   $\underline{X} = [\underline{x}_1, \underline{x}_2, \dots, \underline{x}_p]$

Goal

find the 1D subspace that is closest to a set of points ("best" fits data)  
i.e., distance from  $\underline{x}_i$  to subspace to be as small as possible



$$d_i^2 = \|\underline{x}_i - \text{Proj}_{\underline{a}} \underline{x}_i\|_2^2$$

We want to minimize  $\sum_{i=1}^p d_i^2$

⊥

### Projection Matrices

If  $A \in \mathbb{R}^{n \times p}$  spans a subspace, then projection of  $\underline{X}$  onto  $\text{span}(\text{cols}(A)) = \text{Proj}_A \underline{X}$

If columns of  $A$  are linearly independent,

$$\text{Proj}_A \underline{X} = A(A^T A)^{-1} A^T \underline{X}$$

- $P_A = P_A^2 = P_A P_A = P_A^T = P_A^T P_A$
- If  $A = \underline{a}$ , then  $P_a = a(a^T a)^{-1} a^T = \frac{a a^T}{a^T a}$
- The orthogonal complement of a subspace is the orthogonal to the subspace.

Let B be a basis for orthogonal complement.

$$A^T B = 0$$

For any vector  $\underline{x} \in \mathbb{R}^n$ , it can be written as  $\underline{x} = P_A \underline{x} + P_B \underline{x} = (P_A + P_B) \underline{x}$

$$I \underline{x} = P_A \underline{x} + P_B \underline{x} = (P_A + P_B) \underline{x}$$

$$\rightarrow I = P_A + P_B$$

$$\rightarrow \mathbf{P}_B = \mathbf{I} - \mathbf{P}_A$$

$$d_i^2 = \|\underline{x}_i - P_{\underline{a}} \underline{x}_i\|_2^2$$

$$= \left\| \underline{x}_i - \frac{a a^T}{a^T a} \underline{x}_i \right\|_2^2$$

$$= \left\| \left( I - \frac{a a^T}{a^T a} \right) \underline{x}_i \right\|_2^2$$

Projection on  $\underline{a}$  (orthogonal complement)

$$= \underline{x}_i^T \left( I - \frac{a a^T}{a^T a} \right)^T \left( I - \frac{a a^T}{a^T a} \right) \underline{x}_i$$

$$= \underline{x}_i^T \left( I - \frac{a a^T}{a^T a} \right) \underline{x}_i$$

$$= \underline{x}_i^T \underline{x}_i - \frac{\underline{x}_i^T a a^T \underline{x}_i}{a^T a}$$

We want to minimize

$$\sum_{i=1}^p d_i^2 = \sum_{i=1}^p \underline{x}_i^T \underline{x}_i - \frac{\underline{x}_i^T \underline{a} \underline{a}^T \underline{x}_i}{\underline{a}^T \underline{a}}$$

$\underline{x}_i^T \underline{x}_i$  does not depend on  $\underline{a}$ .

$$\hat{\underline{a}} = \underset{\underline{a}}{\operatorname{argmin}} \sum_{i=1}^p d_i^2$$

$$\hat{\underline{a}} = \underset{\underline{a}}{\operatorname{argmin}} \sum_{i=1}^p -\frac{\underline{x}_i^T \underline{a} \underline{a}^T \underline{x}_i}{\underline{a}^T \underline{a}}$$

(Minimum of negative is the maximum of non-negative)

$$\hat{\underline{a}} = \underset{\underline{a}}{\operatorname{argmax}} \sum_{i=1}^p \frac{\underline{x}_i^T \underline{a} \underline{a}^T \underline{x}_i}{\underline{a}^T \underline{a}}$$

$$\hat{\underline{a}} = \underset{\underline{a}}{\operatorname{argmax}} \sum_{i=1}^p \frac{\underline{a}^T \underline{x}_i \underline{x}_i^T \underline{a}}{\underline{a}^T \underline{a}}$$

( $\underline{x}_i^T \underline{a}$ ,  $\underline{a}^T \underline{x}_i$  are scalars), ( $\underline{x}_i^T \underline{a} = \underline{a}^T \underline{x}_i$ )

$$\hat{\underline{a}} = \underset{\underline{a}}{\operatorname{argmax}} \frac{\underline{a}^T \mathbf{X} \mathbf{X}^T \underline{a}}{\underline{a}^T \underline{a}}$$

The vector  $\hat{\underline{a}}$  that achieves the maximum is called the **1<sup>st</sup> left singular vector** of X.

The value of  $\frac{\hat{\underline{a}}^T \mathbf{X} \mathbf{X}^T \hat{\underline{a}}}{\hat{\underline{a}}^T \hat{\underline{a}}} = \sigma_1^2$  is the **squared 1<sup>st</sup> singular value** of X.

## The Singular Value Decomposition (SVD)

Consider a matrix  $X \in \mathbb{R}^{n \times p}$ . There exist matrices  $U$ ,  $\Sigma$ ,  $V$  such that

$$X = U \Sigma V^T$$

$$X = U \Sigma V^T$$

$U \in \mathbb{R}^{n \times n}$  is orthogonal ( $U^T U = U U^T = I$ ), called left singular vectors  
 $V \in \mathbb{R}^{p \times p}$  is orthogonal ( $V^T V = V V^T = I$ ), called right singular vectors  
 $\Sigma \in \mathbb{R}^{n \times p}$  is diagonal; diagonal elements called singular values

The columns of  $U$  are called the "left singular vectors".

$U$  is an orthogonal matrix ( $U^T U = U U^T = I$ ).

The columns of  $U$  give an orthonormal basis for the columns of  $X$ .

The columns of  $V$  are called the "right singular vectors".

$V$  is an orthogonal matrix ( $V^T V = V V^T = I$ ).

The columns of  $V^T$  (rows of  $V$ ) are the basis coefficients (weights on the column of  $U$ ) need to represent each column of  $X$ .

$\Sigma$  is diagonal with non-negative diagonal elements.

$$\Sigma = \begin{bmatrix} \sigma_1 & & & \\ & \sigma_2 & & \\ & & \ddots & \\ & & & \sigma_n \end{bmatrix} \quad n=p$$

$$\left[ \begin{array}{ccc} \sigma_1 & & \\ & \ddots & \\ & & \sigma_p \\ \hline & & \end{array} \right] \quad \left. \begin{array}{l} p \times p \\ (n-p) \times p \end{array} \right\} \quad n > p$$

$$\left[ \begin{array}{cc} \sigma_1 & \\ & \ddots \\ & & \sigma_n \\ \hline & & \end{array} \right] \quad \left. \begin{array}{l} n \times n \\ n \times (p-n) \end{array} \right\} \quad n < p$$

$$\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_p$$

Let  $U = [u_1, u_2, \dots, u_n]$

$u_1$  is the best 1D subspace fit to  $x_i$ 's (all data).

$\tilde{x}_i^{(1)} = x_i - Proj_{u_1} x_i$  projection  $x_i$  onto  $u_1$

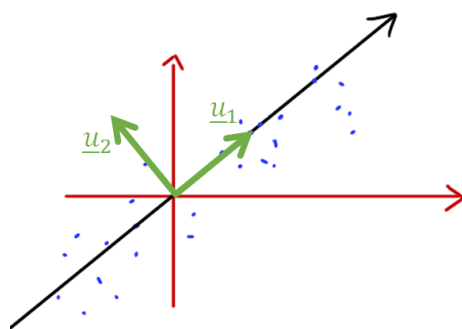
$u_2$  is the best 1D subspace fit to  $\tilde{x}_i^{(1)}$ 's.

$\tilde{x}_i^{(2)} = x_i - Proj_{u_1} x_i - Proj_{u_2} x_i$ .

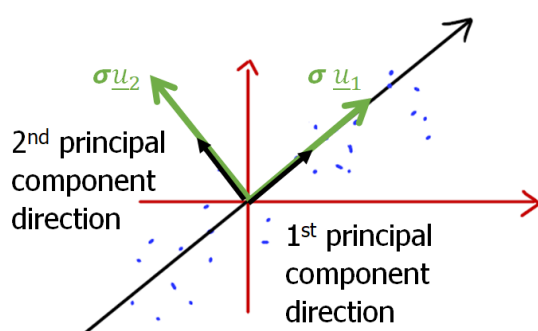
$u_3$  is the best 1D subspace fit to  $\tilde{x}_i^{(2)}$ 's.

...

$[u_1, u_2, \dots, u_k]$  is best k-dim subspace fit to  $x_i$ 's.



## Principal Component Analysis



If  $X = U\Sigma V^T$ , then left singular vectors of  $X$  are called **Principal Component Directions**.

The 1<sup>st</sup> principal component directions = the 1<sup>st</sup> left singular vector of a matrix

The 2<sup>nd</sup> principal component directions = the 2<sup>nd</sup> left singular vector of a matrix