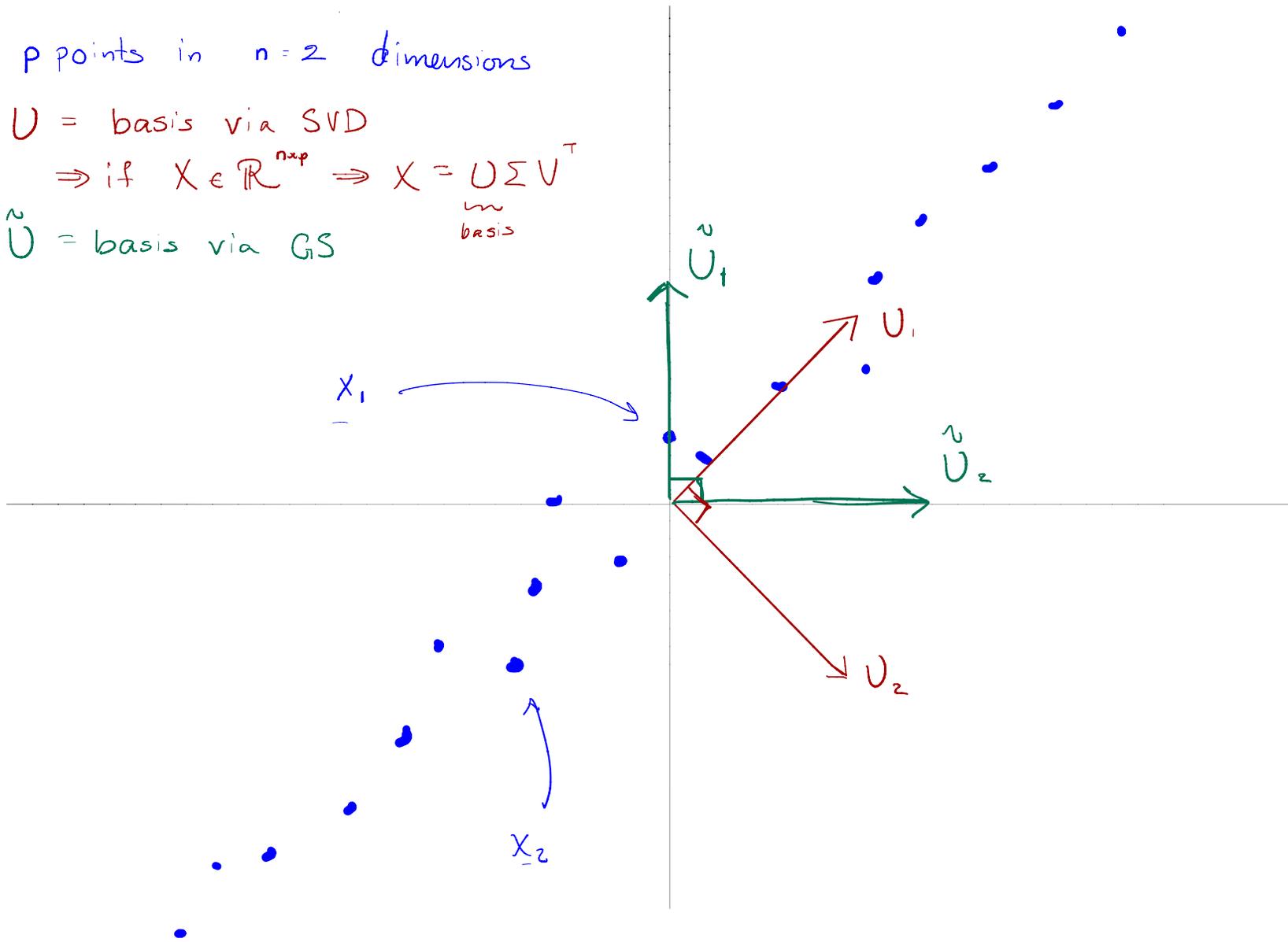# Lecture 7 : Introduction to the Singular Value Decomposition

p points in n=2 dimensions

U = basis via SVD
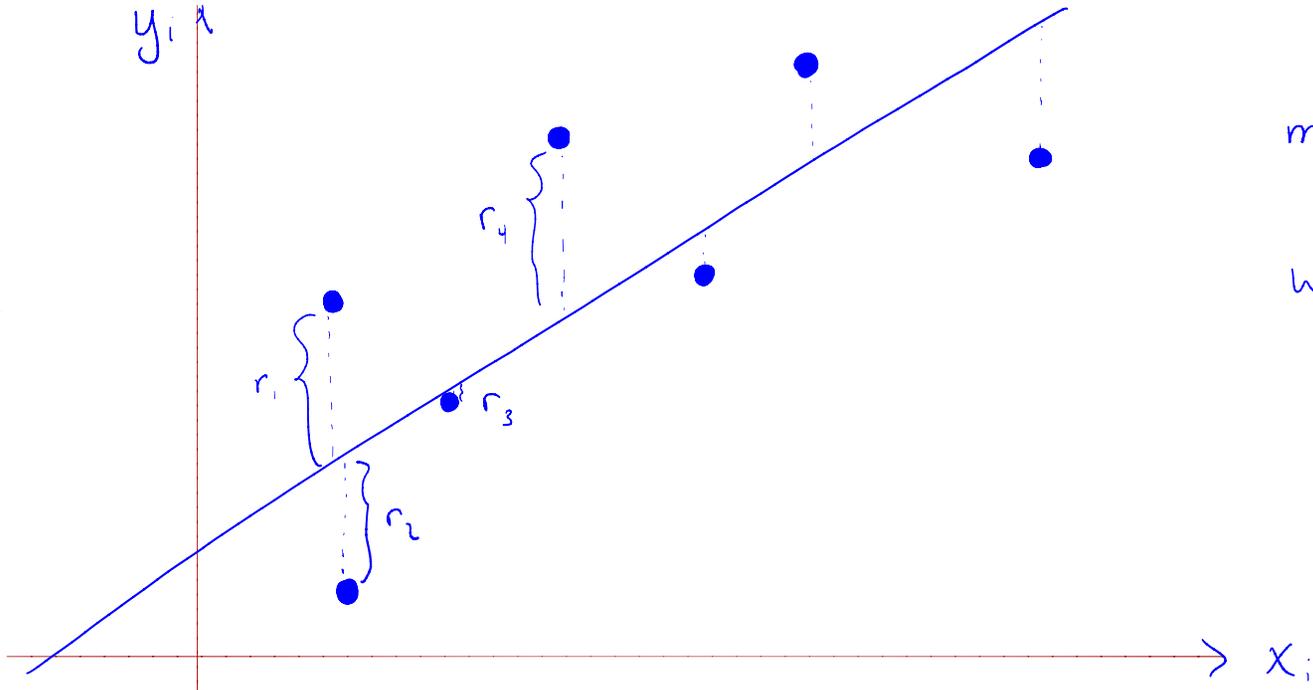  $\Rightarrow$ if $X \in \mathbb{R}^{n \times p} \Rightarrow X = U\Sigma V^T$
                                                      $\underbrace{\quad}_{basis}$
$\tilde{U}$ = basis via GS

$\tilde{U}_1$

$U_1$

$X_1$

$\tilde{U}_2$

$U_2$

$X_2$

$U_1$ is the 1d subspace that is closest to all the $X_i$'s (ie. best 1d subspace fit)

$x_i \in \mathbb{R}$

**L.S.**

$y_i$

$r_1$

$r_2$

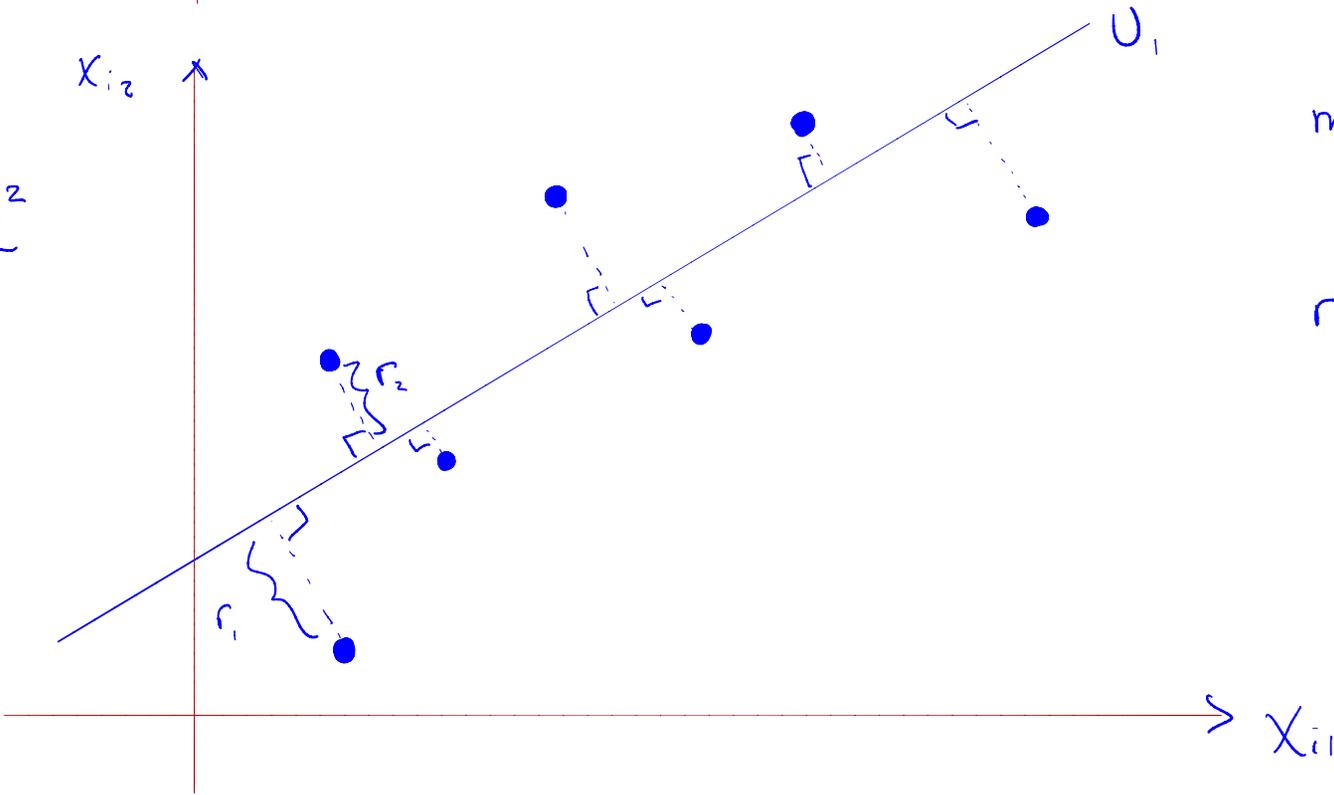$r_3$

$r_4$

$x_i$

minimize $\sum_i r_i^2$

write

$y_i \approx a x_i + b$

$\underline{w}$

$r_i = y_i - (a x_i + b)$

---

$x_i \in \mathbb{R}^2$

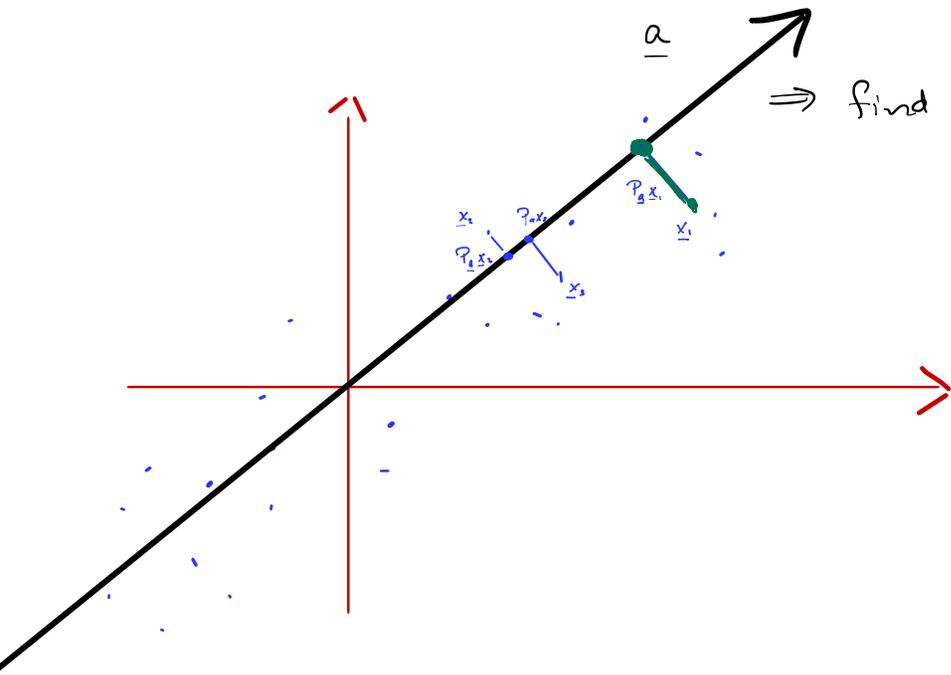**SVD**

$x_{i2}$

$U_1$

$r_1$

$r_2$

$x_{i1}$

minimize $\sum_i \|r_i\|_2^2$

$r_i = \underline{x}_i - P_{U_1} \underline{x}_i$

# Introduction to the Singular Value Decomposition (SVD)

Ex. find a 1d-subspace (line through origin) that is closest to a set of points $\underline{x}_1, \underline{x}_2, ..., \underline{x}_n \in \mathbb{R}^P$



$\Rightarrow$ find $\underline{a}$ to minimize sum of squared distances

## Projection Matrices

$$P_A = A(A^TA)^{-1}A^T \Rightarrow P_a = \underline{a}(\underline{a}^T\underline{a})^{-1}\underline{a}^T$$

$$P_A^2 = A(A^TA)^{-1}A^TA(A^TA)^{-1}A^T$$

$$= A(A^TA)^{-1}A^T = P_A$$

given a subspace $S$ spanned by columns of $A$, the orthogonal complement of $S$ (aka orthogonal complement of $A$) is the set of all vectors $b$ orthogonal to all columns of $A$ (i.e. orthogonal to every vector in $S$)

Let $A \in \mathbb{R}^{p \times r}$, $B \in \mathbb{R}^{p \times (p-r)}$ be orthogonal complements

$\Rightarrow A^TB = 0$ and any $x \in \mathbb{R}^P$ can be written as $\underline{x} = A\underline{u} + B\underline{v}$ for some $\underline{u} \in \mathbb{R}^r, \underline{v} \in \mathbb{R}^{p-r}$
$\underbrace{\qquad}_{P_A\underline{x}} \quad \underbrace{\qquad}_{P_B\underline{x}}$

$\Rightarrow P_A + P_B = I$

$\Rightarrow I - P_A = P_A + P_B - P_A = P_B$

$\Rightarrow I - P_A$ is also a Projection matrix!

distance from $\underline{x}_i$ to line $\underline{a}$ :

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

$$= \| \underline{x}_i - \underline{a}(\underline{a}^T \underline{a})^{-1} \underline{a}^T \underline{x}_i \|_2^2$$

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

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

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

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

$$= \underline{x}_i^T \underline{x}_i - \frac{(\underline{a}^T \underline{x}_i)^2}{\underline{a}^T\underline{a}}$$

Want to minimize $\sum_{i=1}^{P} d_i^2 = \sum_{i=1}^{P} \left( \underline{x}_i^T \underline{x}_i - \frac{(\underline{a}^T \underline{x}_i)^2}{\underline{a}^T \underline{a}} \right)$

<span style="color:red">constant with respect to $\underline{a}$</span>

$\Rightarrow \arg\min_{\underline{a}} \sum_{i=1}^{P} d_i^2(\underline{a}) = \arg\max_{\underline{a}} \sum_{i=1}^{P} \frac{\underline{a}^T \underline{x}_i \underline{x}_i^T \underline{a}}{\underline{a}^T \underline{a}} = \arg\max_{\underline{a}\,:\,\underline{a}^T\underline{a}=1} \sum_{i=1}^{P} \underline{a}^T \underline{x}_i \underline{x}_i^T \underline{a}$

<span style="color:red">Let $X = [\underline{x}_1, \underline{x}_2, \ldots, \underline{x}_P] \in \mathbb{R}^{n \times P}$</span>

$\Rightarrow \underline{U}_1 = \arg\max_{\underline{a}\,:\,\underline{a}^T\underline{a}=1} \sum_{i}^{P} \underline{a}^T \underline{x}_i \underline{x}_i^T \underline{a} = \arg\max_{\underline{a}\,:\,\underline{a}^T\underline{a}=1} \underline{a}^T X X^T \underline{a}$

Value of $\underline{a}$ that achieves maximum is $1^{st}$ left singular vector of $X$, denoted $\underline{U}_1$

also, $\sigma_1 = \max_{\underline{a}\,:\,\underline{a}^T\underline{a}=1} \| X^T \underline{a} \|_2 = \| X^T \underline{U}_1 \|_2$

is called the $1^{st}$ singular value of $X$

(also called the "operator norm" of $X$: $\| X \|_{op} = \| X \|_2$)

not like 2-norm of vector!

The bigger $\sigma_1$ is, the smaller $\sum_i d_i^2$ is

$\Rightarrow$ the better the $\underline{x}_i$'s are aligned with a 1-d subspace

# The Singular Value Decomposition

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

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

$U \in \mathbb{R}^{n \times n}$ 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

$n = p$        $n > p$        $n < p$

$$\Sigma = \begin{bmatrix} \sigma_1 & & & \\ & \sigma_2 & & \bigcirc \\ & & \ddots & \\ \bigcirc & & & \sigma_n \end{bmatrix} \quad \left.\begin{bmatrix} \sigma_1 & & \bigcirc \\ & \ddots & \\ \bigcirc & & \sigma_p \\ \hline & \bigcirc & \end{bmatrix}\right\} \begin{matrix} p \times p \\ \\ (n-p) \times p \end{matrix} \quad \underbrace{\begin{bmatrix} \sigma_1 & & \bigcirc \\ & \ddots & \\ \bigcirc & & \sigma_n \end{bmatrix}}_{n \times n} \Big| \underbrace{\bigcirc}_{n \times (p-n)}$$

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

The singular values weigh (scale the length) of the corresponding singular vectors.

- The number of non-zero singular vectors is the RANK of $X$.

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

- $U$ gives orthobasis for all of $\mathbb{R}^n$.

- 1st $r$ columns of $U$ give basis of best $r$-dim subspace fit to columns of $X$

- $\sigma_i$'s indicate how important each subspace dimension is to representing/approximating data

- 1st $r$ columns of $V$ give coordinates/locations of each $\underline{x}_i$ within the subspace spanned by $\underline{U}_1, \dots, \underline{U}_r$

- $\underline{U}_1$ = best 1d subspace fit to all data

$\tilde{\underline{X}}_i^{(1)} = \underline{X}_i - P_{U_1}\underline{X}_i = i^{th}$ residual $\forall i$

$U_2$ = best 1d subspace fit to all $\tilde{\underline{X}}_i^{(1)}$

$\vdots$

$\tilde{\underline{X}}_i^{(k)} = \underline{X}_i - P_{[U_1, \dots, U_k]}\underline{X}_i = \tilde{\underline{X}}_i^{(k-1)} - P_{U_k}\tilde{\underline{X}}_i^{(k-1)} \quad \forall i$

$U_{k+1}$ = best 1d subspace fit to all $\tilde{\underline{X}}_i^{(k)}$

Singular values $\sigma_1, \sigma_2$, etc. indicate how spread out points are in the subspace.

Recall Gram-Schmidt:

$U_1 = X_1 / \|X_1\|_2$

$\tilde{\underline{X}}_2 = \underline{X}_2 - P_{U_1}\underline{X}_1$

$U_2 = \tilde{\underline{X}}_2 / \|\tilde{X}_2\|_2$

order of points matters!

basis vectors less interpretable