

Lecture 8 :

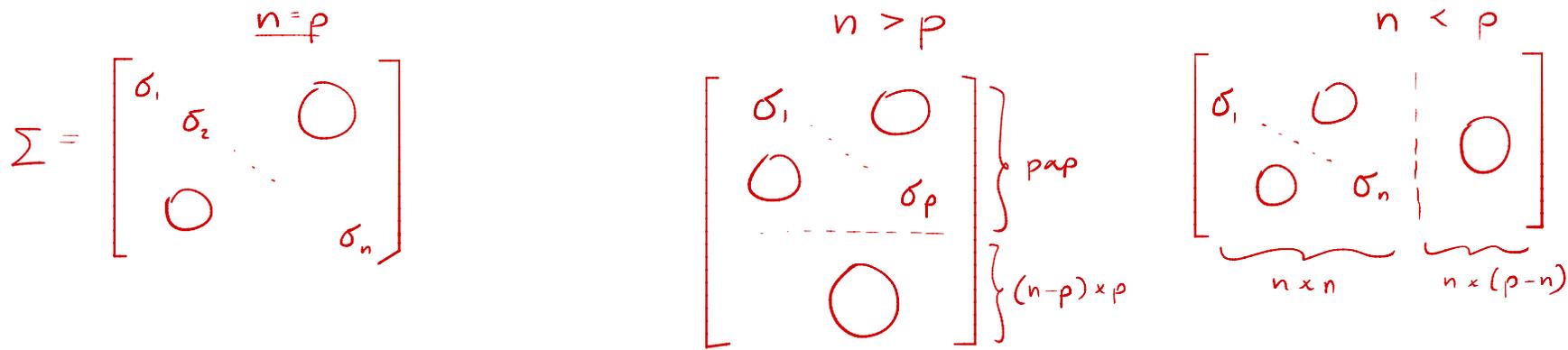
Singular Value Decomposition

The Singular Value Decomposition

Consider a matrix $X \in \mathbb{R}^{n \times p}$. There always exists matrices U, Σ, 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 - basis for cols of X
 $V \in \mathbb{R}^{p \times p}$ 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 - basis for rows of X



Let $r = \min(n, p)$. Then $X = \sum_{i=1}^r u_i \sigma_i v_i^T$, $u_i = i^{\text{th}}$ col of U , $v_i = i^{\text{th}}$ column of V
 = sum of rank-1 matrices

If X is square and has $\sigma_j = 0$ for any j , then X is not invertible, a "singular"

If X is square and not singular, then $X^{-1} = V \Sigma^{-1} U^T$

The SVD of $X^T = (U \Sigma V^T)^T = V \Sigma^T U^T$

\Rightarrow columns of U are basis for the columns of X
 and columns of V are basis for the rows of X

a note on multiplying by diagonal matrices:

$$\begin{bmatrix} A_1 & A_2 & A_3 \end{bmatrix} \begin{bmatrix} \sigma_1 & & \\ & \sigma_2 & \\ & & \sigma_3 \end{bmatrix} = \begin{bmatrix} \sigma_1 A_1 & \sigma_2 A_2 & \sigma_3 A_3 \end{bmatrix}$$

diagonal matrix on right \Rightarrow reweigh columns

$$\begin{bmatrix} \sigma_1 & & \\ & \sigma_2 & \\ & & \sigma_3 \end{bmatrix} \begin{bmatrix} -a_1^- \\ -a_2^- \\ -a_3^- \end{bmatrix} = \begin{bmatrix} -\sigma_1 a_1^- \\ -\sigma_2 a_2^- \\ -\sigma_3 a_3^- \end{bmatrix}$$

diagonal matrix on left \Rightarrow reweigh rows

- U gives orthonormal basis for all of \mathbb{R}^p .
- 1st r columns of U give basis of best r -dim subspace fit to columns of X
- 1st r columns of V give coordinates/locations of each \underline{x}_i within the subspace spanned by $\underline{U}_1, \dots, \underline{U}_r$
- σ_i 's indicate how important each subspace dimension is to representing/approximating data

recall $\sigma_i = \|X^T \underline{U}_i\|_2 = \left[\sum_1 (U_i^T \underline{x}_i)^2 \right]^{1/2}$ where $X = \begin{bmatrix} -x_1^T- \\ -x_2^T- \\ \vdots \\ -x_p^T- \end{bmatrix}$

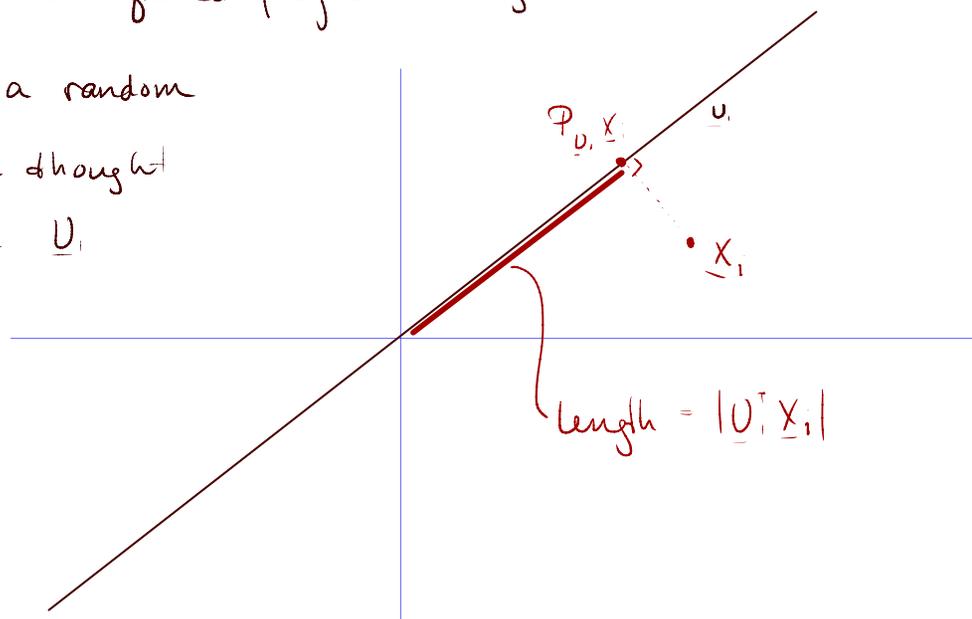
also recall $\mathcal{P}_{\underline{U}_i} \underline{x}_i = \underline{U}_i \underline{U}_i^T \underline{x}_i$ and that,

since \underline{U}_i is an orthonormal basis, it is length-preserving: $\|\underline{U}_i \underline{v}\|_2 = \|\underline{v}\|_2$ for any \underline{v} .

This means the length² of the projection of \underline{x}_i onto \underline{U}_i is $\|\underline{U}_i \underline{U}_i^T \underline{x}_i\|_2^2 = (U_i^T \underline{x}_i)^2$

In other words, $\sigma_i^2 = \sum_1 (U_i^T \underline{x}_i)^2$ is the sum of squared projection lengths

If we think of the \underline{x}_i 's as realizations of a random variable with mean = 0, then σ_i^2 can also be thought of as the variance of the \underline{x}_i 's in the \underline{U}_i direction.



The "Economy SVD"

Let $X \in \mathbb{R}^{n \times p}$, with rank $r \ll \min(n, p)$

Then $X = U \Sigma V^T$

$U \in \mathbb{R}^{n \times n}$ $\Sigma \in \mathbb{R}^{n \times p}$ $V^T \in \mathbb{R}^{p \times p}$

$$X = U \Sigma V^T = \tilde{U} \tilde{\Sigma} \tilde{V}^T$$

↑ this is the economy SVD

Netflix example

$n \approx 5k$ movies

$p \approx 100m$ customers

\Rightarrow storage = $5k \cdot 100m \cdot 4 \text{ bytes} \approx 2 \text{ TB} \approx$ okay to store, difficult to use in learning algorithms.

Now let's say we find a rank- r approximation to X using the subspace approximation theorem, and use its economy SVD. if $r=10$,

\tilde{U} takes $5k \cdot 10 \cdot 4 \text{ bytes} = 200 \text{ kB}$, \tilde{V} takes $100m \cdot 10 \cdot 4 \text{ bytes} = 4 \text{ GB}$, $\tilde{\Sigma}$ takes $10 \cdot 4 \text{ bytes}$

MUCH SMALLER!

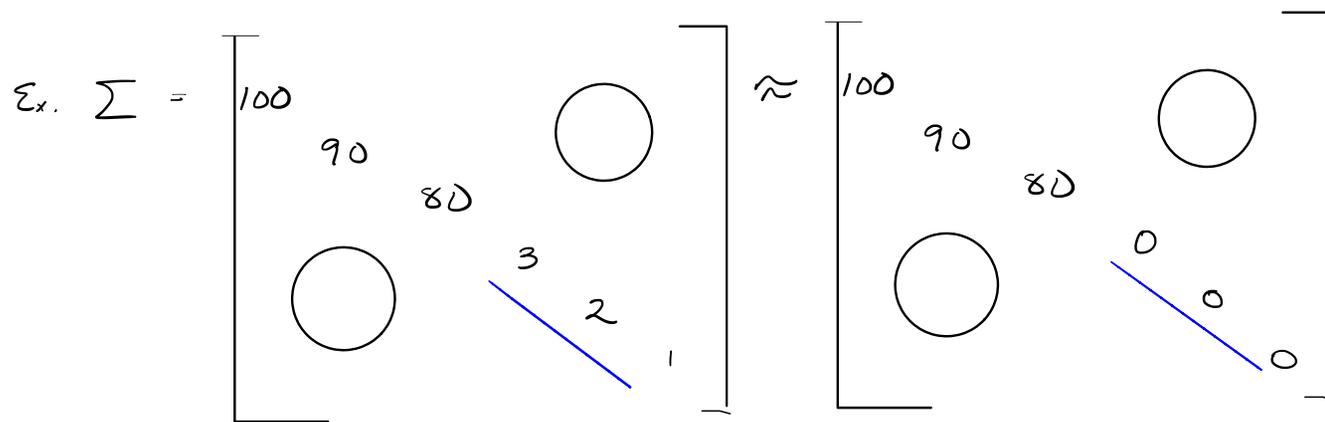
Subspace Approximation Theorem

If $X \in \mathbb{R}^{p \times n}$ has rank $r > k$, then

$$\operatorname{argmin}_{Z: \operatorname{rank}(Z)=k} \|X - Z\|_F^2 = U_k \Sigma_k V_k^T$$

$\underbrace{\quad}_{\substack{\text{1st } k \text{ cols} \\ \text{of } U}} \quad \underbrace{\quad}_{\substack{\text{1st } k \times k \\ \text{block of} \\ \Sigma}} \quad \underbrace{\quad}_{\substack{\text{1st } k \\ \text{cols of} \\ V}}$

and $\|X - X_k\|_F^2 = \sum_{i=k+1}^r \sigma_i^2$



Frobenius matrix norm

$$\|A\|_F = \left(\sum_{i,j} A_{ij}^2 \right)^{1/2}$$

if $A = [A_1, A_2, \dots, A_p]$,

$$\|A\|_F^2 = \sum_{i=1}^p \|A_i\|_2^2$$

$\Rightarrow X$ is "almost" low rank
 error of rank-3 approximation
 is $3^2 + 2^2 + 1^2 = 14$

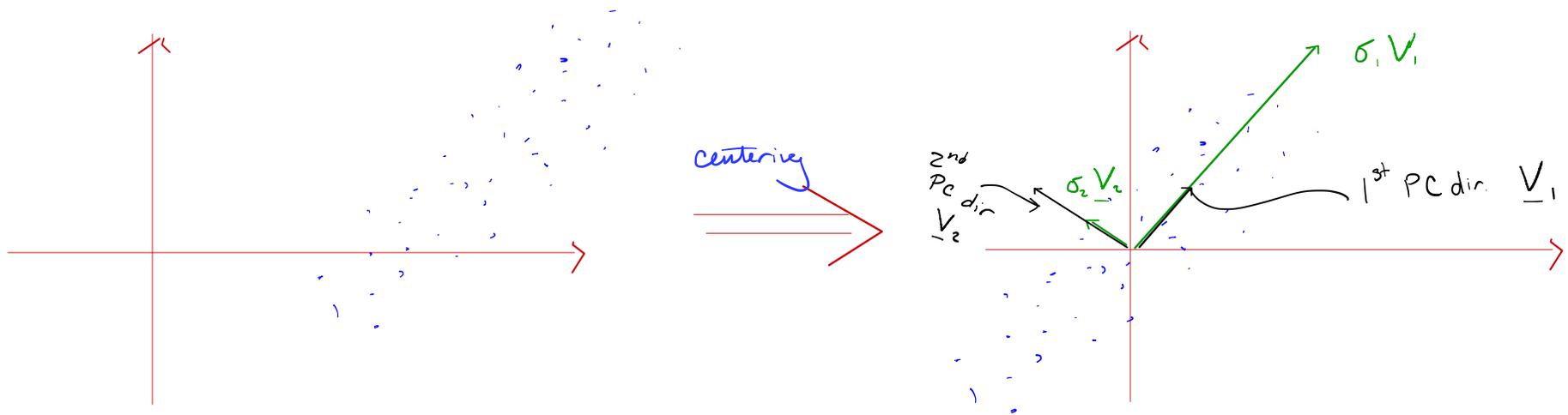
Principal Components Analysis

Let $X \in \mathbb{R}^{n \times p}$ be a data matrix with rows "centered" to have average value of 0.

(n points in p dimensions)

1st "center" data - is $X = \begin{bmatrix} | & | & \dots & | \\ X_1 & X_2 & \dots & X_p \\ | & | & \dots & | \end{bmatrix}$, for $i = 1, 2, \dots, p$, $\underline{X}_i \leftarrow X_i - \text{mean}(X_i)$
 $= X_i - \left(\frac{\sum_{j=1}^n X_{ij}}{n} \right) \underline{1}_{n \times 1}$

$X = U \Sigma V^T \Rightarrow V$ is basis matrix for \mathbb{R}^p



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

Equivalently. let $C = X^T X = (U \Sigma V^T)^T (U \Sigma V^T) = V \Sigma^T U^T U \Sigma V^T = V \Sigma^2 V^T$

called "eigenvalue decomposition"

right singular vectors of $X =$
eigenvectors of $X^T X$

Dimensionality Reduction

We have n points $\underline{x}_i \in \mathbb{R}^p$, $i=1, \dots, n$. Dimensionality reduction means defining new points $\underline{z}_i \in \mathbb{R}^k$, $i=1, \dots, n$ for $k < p$ that preserve important properties of the \underline{x}_i 's.

(e.g. $\|\underline{x}_i - \underline{x}_j\| \approx \|\underline{z}_i - \underline{z}_j\|$)

$$\text{Let } X = \begin{bmatrix} -\underline{x}_1^T \\ -\underline{x}_2^T \\ \vdots \\ -\underline{x}_n^T \end{bmatrix} = U \Sigma V^T \in \mathbb{R}^{n \times p} \quad \Rightarrow \quad X^T = \begin{bmatrix} \underline{x}_1 & \underline{x}_2 & \dots & \underline{x}_n \end{bmatrix} = V \Sigma U^T \in \mathbb{R}^{p \times n}$$

Consider $X_k^T = V_k \Sigma_k U_k^T$; let i^{th} column of X_k^T be $\underline{\tilde{x}}_i$

by subspace approximation theorem, we know $\|\underline{x}_i - \underline{\tilde{x}}_i\|_2^2 \leq \sum_{j=k+1}^r \sigma_j^2$ - if σ_j 's small for $j > k$, then \underline{x}_i is close to $\underline{\tilde{x}}_i$!

but $\underline{\tilde{x}}_i$'s are all in a k -dimensional subspace, so we can represent them

in terms of their coordinates in the subspace!

$$\Rightarrow \underline{z}_i = V_k^T \underline{\tilde{x}}_i = i^{\text{th}} \text{ col of } (\Sigma_k U_k^T) \in \mathbb{R}^k$$

Ex: $\underline{x}_i \in \mathcal{S} = \{\underline{x} \in \mathbb{R}^3 : x_3 = 0\}$ = horizontal plane

instead of representing each \underline{x}_i using 3d coordinates, we only need to represent where it is in the plane

Example: $X = \begin{bmatrix} 1 & -1 & -1 & 1 \\ -1 & 1 & -1 & 1 \\ 1 & -1 & -1 & 1 \\ -1 & 1 & -1 & 1 \\ 1 & -1 & 0 & 0 \end{bmatrix} \Rightarrow \underline{x}_i = (i^{\text{th}} \text{ row of } X)^T = i^{\text{th}} \text{ column of } X^T, \text{ for } i=1, \dots, 5$

↘ 4-dimensional

(e.g. $\underline{x}_1 = \begin{bmatrix} 1 \\ -1 \\ -1 \\ 1 \end{bmatrix}$)

$$= \begin{bmatrix} 1/\sqrt{5} & 1/2 \\ -1/\sqrt{5} & 1/2 \\ 1/\sqrt{5} & 1/2 \\ -1/\sqrt{5} & 1/2 \\ 1/\sqrt{5} & 0 \end{bmatrix} \begin{bmatrix} \sqrt{5}\sqrt{2} & 0 \\ 0 & 2\sqrt{2} \end{bmatrix} \begin{bmatrix} 1/\sqrt{2} & -1/\sqrt{2} & 0 & 0 \\ 0 & 0 & -1/\sqrt{2} & 1/\sqrt{2} \end{bmatrix}$$

} economy SVD

$$= \begin{bmatrix} \sqrt{5} & 0 \\ 0 & 2 \end{bmatrix} \begin{bmatrix} \sqrt{2} & 0 \\ 0 & \sqrt{2} \end{bmatrix}$$

$$= \begin{bmatrix} 1 & 1 \\ -1 & 1 \\ 1 & 1 \\ -1 & 1 \\ 1 & 0 \end{bmatrix} \begin{bmatrix} 1 & -1 & 0 & 0 \\ 0 & 0 & -1 & 1 \end{bmatrix}$$

$$\Rightarrow X^T = \underbrace{\left(\frac{1}{\sqrt{2}} \begin{bmatrix} 1 & 0 \\ -1 & 0 \\ 0 & -1 \\ 0 & 1 \end{bmatrix} \right)}_{\tilde{V}_k} \underbrace{\left(\begin{bmatrix} 1 & -1 & 1 & -1 & 1 \\ 1 & 1 & 1 & 1 & 0 \end{bmatrix} \cdot \sqrt{2} \right)}_{\sum_k \tilde{U}_k^T} \Rightarrow \underline{x}_i = \text{weighted sum of } \frac{1}{\sqrt{2}} \begin{bmatrix} 1 \\ -1 \\ 0 \\ 0 \end{bmatrix} \text{ \& } \frac{1}{\sqrt{2}} \begin{bmatrix} 0 \\ 0 \\ -1 \\ 1 \end{bmatrix}$$

$$\underline{z}_i = V_k^T \tilde{\underline{x}}_i \Rightarrow \underline{z}_1 = V_k^T \underline{x}_1 = \frac{1}{\sqrt{2}} \begin{bmatrix} 1 & -1 & 0 & 0 \\ 0 & 0 & -1 & 1 \end{bmatrix} \begin{bmatrix} 1 \\ -1 \\ -1 \\ 1 \end{bmatrix} = \frac{1}{\sqrt{2}} \begin{bmatrix} 2 \\ 2 \end{bmatrix} = \begin{bmatrix} \sqrt{2} \\ \sqrt{2} \end{bmatrix}$$

Let $\underline{z}_1 = \sqrt{2} \begin{bmatrix} 1 \\ 1 \end{bmatrix}$, $\underline{z}_2 = \sqrt{2} \begin{bmatrix} -1 \\ 1 \end{bmatrix}$, etc. each $\underline{z}_i \in \mathbb{R}^2 \Rightarrow$ reduced-dimensional representation of \underline{x}_i 's

$$X = \begin{bmatrix} 1.1 & 2.0 & 3.4 & 4.05 \\ 2.01 & 4.2 & 6.1 & 8.05 \\ 3.2 & 6.0 & 9.05 & 12 \\ 4 & 8.1 & 12 & 16 \\ 5 & 10 & 15 & 20 \end{bmatrix} = \underbrace{\begin{bmatrix} 1 \\ 2 \\ 3 \\ 4 \\ 5 \end{bmatrix}}_{\begin{bmatrix} 1 & 2 & 3 & 4 \\ 2 & 4 & 6 & 8 \\ 3 & 6 & 9 & 12 \\ 4 & 8 & 12 & 16 \\ 5 & 10 & 15 & 20 \end{bmatrix}} + \text{"errors"} \approx \underbrace{\begin{bmatrix} 1 \\ 2 \\ 3 \\ 4 \\ 5 \end{bmatrix}}_{\tilde{X}} \begin{bmatrix} 1 & 2 & 3 & 4 \end{bmatrix}$$

so $\tilde{X} \approx X$
 can find \tilde{X} by
 thresholding small
 singular values.

\tilde{X} is the same size as $X \Rightarrow$ no dimensionality reduction has occurred yet

(let's think of $X \in \mathbb{R}^{n \times p}$ - each row is a point x_i^T , $x_i \in \mathbb{R}^p$, $i=1, \dots, n$)

to perform dim red., want n new points $z_i \in \mathbb{R}^r$ for $r < p$, $i=1, \dots, n$

To do this, just use $k=1$, $V_1 = \begin{bmatrix} 1 \\ 2 \\ 3 \\ 4 \\ 5 \end{bmatrix} \frac{1}{\sqrt{55}} \Rightarrow \begin{aligned} z_1 &= V_1^T x_1 = 7.51 \approx \sqrt{55} \\ z_2 &= V_1^T x_2 = 14.94 \approx 2\sqrt{55} \end{aligned}$

More generally:

Step 1: find $\tilde{X} \Leftarrow$ see how entries of X are related / tied together

Step 2: exploit those ties to get simple representation of each x_i - set small singular values to zero,
 let $z_i = V_k^T x_i$ for $i=1, \dots, n$.