# Singular Value Decomposition
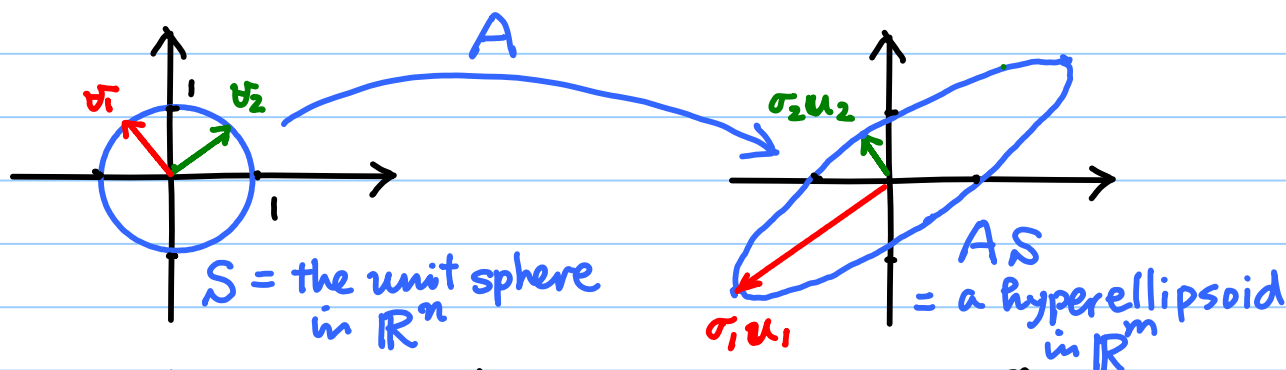
- SVD is a matrix factorization that is useful for many applications, e.g., search engines, LS problems, tomographic image reconstruction, ...

- SVD can be a concepual tool in linear algebra
  $\Rightarrow$ via SVD, we can check :
  - a given matrix is near singular
  - rank of the matrix
  - etc.

- $\exists$ a numerically stable algorithm to compute the SVD of a given matrix ( it's expensive though ...)
  In fact , one of the hottest topics in numerical linear algebra is how to compute a good approximation to the SVD of a _large_ matrix _fast_!

☆  A Geometric Observation

Let  $A \in \mathbb{R}^{m \times n}$, and consider how A maps an input vector in $\mathbb{R}^n$ to an output vector in $\mathbb{R}^m$.

"The image of the unit sphere under any $m \times n$ matrix is a hyperellipsoid"

A

$\vec{v}_1$  $\vec{v}_2$

$\sigma_2 u_2$

$\sigma_1 u_1$

S = the unit sphere in $\mathbb{R}^n$

AS = a hyperellipsoid in $\mathbb{R}^m$

Let $\{v_1, \cdots, v_n\}$ be an ONB of $\mathbb{R}^n$

ONB = ortho-normal basis

Let $\{u_1, \cdots, u_m\}$ be an ONB of $\mathbb{R}^m$

Let $\{\sigma_1, \cdots, \sigma_m\}$ be a set of $m$ scalars with $\sigma_i \geq 0$, $i = 1, \cdots m$.

Then, $\sigma_i u_i$ is the $i$th principal semiaxis with length $\sigma_i$ in $\mathbb{R}^m$.

Now, if rank$(A) = r$, then exactly $r$ of $\{\sigma_1, \cdots, \sigma_m\}$ are nonzero, and exactly $m - r$ of $\sigma_i$'s are zero.

So, if $\underline{m \geq n}$, then $\underline{\text{rank}(A) \leq n}$. i.e., at most $n$ of $\sigma_i$'s    ← full rank if $= n$ are nonzero.

For simplicity, let's assume $m \geq n$ and rank$(A) = n$ for the time being.

<u>Def.</u> The <span style="color:red">singular values</span> of A $\underset{\text{def}}{\Longleftrightarrow}$ The lengths of the $n$ principal semiaxes of the hyperellipsoid AS

Our convention: $\sigma_1 \geq \sigma_2 \geq \cdots \geq \sigma_n \gneq 0$

<u>Def.</u> The $n$ **left singular vectors** of $A$
$\stackrel{\text{def}}{\Longleftrightarrow}$ $\{u_1, \cdots, u_n\}$ : the unit vectors
in $\mathbb{R}^m$ along the principal semiaxes of $AS$.
So, $\sigma_i u_i$ is the $i$th largest principal
semiaxis of $AS$.

<u>Def.</u> The $n$ **right singular vectors** of $A$
$\stackrel{\text{def}}{\Longleftrightarrow}$ $\{v_1, \cdots, v_n\} \in S$ : the preimages
of the principal semiaxes of $AS$, i.e.,
$$A v_i = \sigma_i u_i \quad i = 1, \cdots, n.$$

☆ <u>Reduced SVD</u>

$$m \begin{bmatrix} A \end{bmatrix} \begin{bmatrix} v_1 & \cdots & v_n \end{bmatrix}^n = {}^m \begin{bmatrix} u_1, & \cdots & u_n \end{bmatrix} \begin{bmatrix} \sigma_1 & & 0 \\ & \ddots & \\ 0 & & \sigma_n \end{bmatrix}^n$$

$$\underbrace{\phantom{xxxxx}}_{\text{``}V} \qquad \underbrace{\phantom{xxxxx}}_{\text{``}\hat{U}} \qquad \underbrace{\phantom{xxxxx}}_{\text{``}\hat{\Sigma}}$$

$$\Rightarrow \quad A V = \hat{U} \hat{\Sigma}$$
$$\uparrow \quad \uparrow \qquad \uparrow \quad \uparrow$$
$$m \times n \quad n \times n \quad m \times n \quad n \times n$$

Since $V$ is an orthogonal matrix,

$$\boxed{A = \hat{U} \hat{\Sigma} V^T}$$

The **reduced** SVD of $A$.

$m \geq n$

$A \qquad \hat{U} \quad \hat{\Sigma} \quad V^T$

$m < n$

$A \qquad U \quad \hat{\Sigma} \quad \hat{V}^T$

## ✷ Full SVD

Note $\hat{U} \in \mathbb{R}^{m \times n}$ in the reduced SVD with $m \geq n$.

$\Rightarrow$ The column vectors of $\hat{U}$ do not form an ONB of $\mathbb{R}^m$ unless $m = n$.

$\Rightarrow$ Remedy: adjoin $m - n$ ON vectors to $\hat{U}$ to form an orthogonal matrix $U$. Then $\hat{\Sigma}$ must be changed to $\Sigma \in \mathbb{R}^{m \times n}$

$$\boxed{A = U \Sigma V^T}$$  The **full** SVD of A



$$m \geq n \qquad\qquad m < n$$
$$A \qquad U \quad \Sigma \quad V^T \qquad A \quad U \quad \Sigma \quad V^T$$

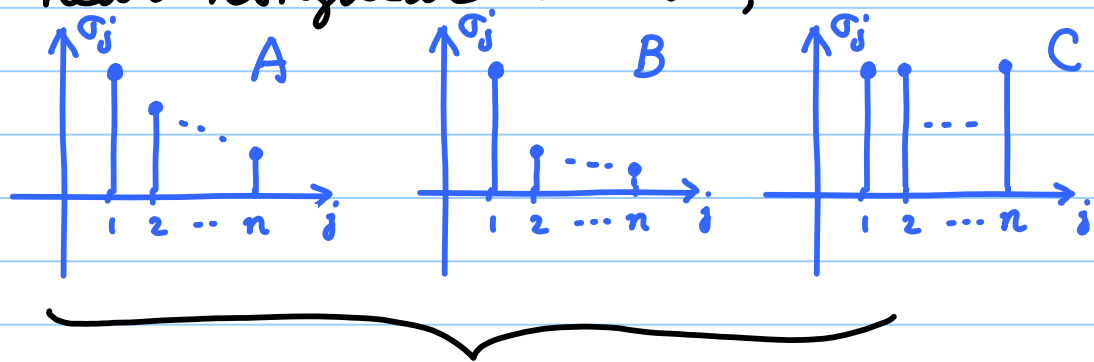For non-full rank matrices, i.e., $\text{rank}(A) = r < \min(m, n)$, $\exists$ only $r$ positive singular values.

So,

$$\Sigma = \begin{bmatrix} \sigma_1 & & & \\ & \ddots & & \\ & & \sigma_r & \\ & & & 0 \ddots 0 \\ \hline & & O & \end{bmatrix} \quad \text{or} \quad \begin{bmatrix} \sigma_1 & & & & \\ & \ddots & & & O \\ & & \sigma_r & & \\ & & & 0 \ddots 0 & \end{bmatrix}$$

$$m \geq n \qquad\qquad\qquad m \leq n$$

Let's consider $m = n$ and full rank case. Theoretically, it's invertible, non singular.

However, we can gain more info by checking the distribution of the singular values of $A$ $\Rightarrow$ We can see whether $A$ is near singular or not, etc.



Out of these three scenarios, which matrix do you think behaves best numerically?
$\Rightarrow$ C.

## ☆ Pseudoinverse via SVD

$$A^+ = V \Sigma^+ U^T$$

where

$$\Sigma^+ := \begin{bmatrix} \frac{1}{\sigma_1} & & & \\ & \ddots & & \\ & & \frac{1}{\sigma_r} & 0 \\ & & & \ddots & 0 \\ \hline & O & & \end{bmatrix}_{m \geq n} \quad \text{or} \quad \begin{bmatrix} \frac{1}{\sigma_1} & & & & \\ & \ddots & & & O \\ & & \frac{1}{\sigma_r} & 0 & \\ & & & \ddots & 0 \end{bmatrix}_{m \leq n}$$

Check: $AA^+ = U \Sigma V^T V \Sigma^+ U^T$

$$= U \Sigma \Sigma^+ U^T$$

$$= U \begin{bmatrix} \begin{smallmatrix} 1 & & \\ & \ddots & \\ & & 1 & \\ & & & 0 \\ & & & & \ddots & \\ & & & & & 0 \end{smallmatrix} \\ \hline O \end{bmatrix} U^T$$

$$= [u_1 \cdots u_r \ 0 \cdots 0] \begin{bmatrix} u_1^T \\ \vdots \\ u_m^T \end{bmatrix}$$

$$= \hat{U} \hat{U}^T$$

$\longrightarrow$ reduced version.

Similarly, $A^+ A = \hat{V} \hat{V}^T$

## The Moore - Penrose Conditions

For a given matrix $A \in \mathbb{R}^{m \times n}$, if $X \in \mathbb{R}^{n \times m}$ satisfies the following:

$\begin{cases} \text{(1)} \ AXA = A \\ \text{(2)} \ XAX = X \\ \text{(3)} \ (AX)^T = AX \\ \text{(4)} \ (XA)^T = XA \end{cases}$

then X is called the pseudoinverse (or the Moore-Penrose inverse) of A and written as $A^+$

$\exists$ many applications using $A^+$!

Note: If $\| AX - I_m \|_F \to min$
then $X = A^+$.

# S V D

★ <u>Formal Definition</u>

Let $A \in \mathbb{R}^{m \times n}$

Then SVD of $A$ is a factorization

full SVD → $A = U \Sigma V^T$

where $U \in \mathbb{R}^{m \times m}$ orthogonal
$\Sigma \in \mathbb{R}^{m \times n}$ diagonal
$V \in \mathbb{R}^{n \times n}$ orthogonal

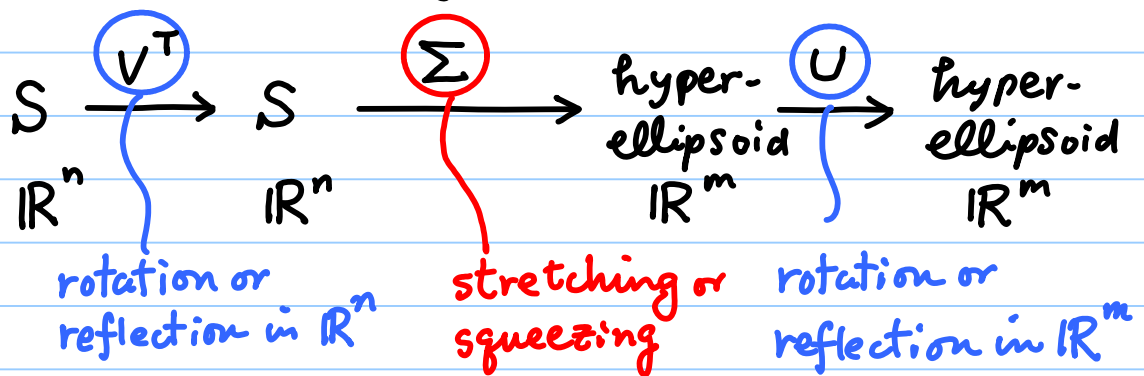$\text{diag}(\Sigma) = [\sigma_1, \sigma_2, \cdots, \sigma_p]^T$

$\sigma_1 \geq \sigma_2 \geq \cdots \geq \sigma_p \geq 0$.

$p = \min(m, n)$

$\text{rank}(A) = r \leq p$.

$A$ & $\Sigma$ are the same shape.

Geometrically,

$$S \xrightarrow{V^T} S \xrightarrow{\Sigma} \text{hyper-ellipsoid} \xrightarrow{U} \text{hyper-ellipsoid}$$

$\mathbb{R}^n \qquad \mathbb{R}^n \qquad \mathbb{R}^m \qquad \mathbb{R}^m$

rotation or reflection in $\mathbb{R}^n$ — stretching or squeezing — rotation or reflection in $\mathbb{R}^m$

So if we prove every $A \in \mathbb{R}^{m \times n}$ has an SVD, then we shall have proved that $A$ maps the unit sphere in $\mathbb{R}^n$ to a hyperellipsoid in $\mathbb{R}^m$.

**✷ Existence & Uniqueness of SVD**

→ We can get peace of mind
if we know that $\exists!$ SVD for any
given matrix.

___

<u>Thm</u>  Every matrix $A \in \mathbb{R}^{m \times n}$ has
an SVD. Furthermore, the singular
values $\{\sigma_j\}$ are <u>uniquely</u> determined.
If A is square and $\sigma_j$'s are distinct,
then singular vectors $\{u_j\}$, $\{v_j\}$
are <u>uniquely</u> determined up to signs
(i.e., $\pm 1$ factor).

___

(Proof : Existence)

Let's check the largest action of A
first, then do induction.

← definition

$$\text{Set } \sigma_1 = \| A \|_2 = \sup_{v \in S} \| A v \|_2$$

*This is often called "compactness" argument.*

∴ Because we are dealing with vectors
in $\mathbb{R}^n$ (i.e., finite dimensional space),
and $\| A \cdot \|_2$ is a continuous fcn,
$\exists v_1 \in S \subset \mathbb{R}^n$ s.t. $\| A v_1 \|_2 = \sigma_1$  is
attained.
Now set $\tilde{u}_1 = A v_1 \in \mathbb{R}^m$,
and consider orthogonal matrices
$V_1 = [v_1 \; v_2 \cdots v_n] \in \mathbb{R}^{n \times n}$,

$$U_1 = [u_1 \ u_2 \ \cdots \ u_m] \in \mathbb{R}^{m \times m}$$

where $\quad u_1 = \frac{1}{\sigma_1} \tilde{u}_1$

Note $\quad \|u_1\| = \frac{1}{\sigma_1} \|\tilde{u}_1\| = \frac{1}{\sigma_1} \|Av_1\|$

$$= \frac{1}{\sigma_1} \cdot \sigma_1 = 1 \quad \checkmark$$

Then, $U_1^T A V_1 = \begin{bmatrix} u_1^T \\ \vdots \\ u_m^T \end{bmatrix} A [v_1 \ \cdots \ v_n]$

$$= \begin{bmatrix} u_1^T \\ \vdots \\ u_m^T \end{bmatrix} [Av_1 \ \cdots Av_n]$$

<span style="color:red">$\tilde{u}_1 = \sigma_1 u_1$</span>

$$= \begin{bmatrix} \sigma_1 & w^T \\ 0 & \\ \vdots & B \\ 0 & \end{bmatrix}$$

<span style="color:red">$u_j^T u_1 = 0$ for $j \geq 2$.</span>

<span style="color:green">let's call $= \Sigma_1$</span> $\quad {}^{1 \times n-1}$

where $w^T = [u_1^T A v_2, \cdots, u_1^T A v_n] \in \mathbb{R}^{1 \times n-1}$

$$B = \begin{bmatrix} u_2^T A v_2 & \cdots & u_2^T A v_n \\ \vdots & & \vdots \\ u_m^T A v_2 & \cdots & u_m^T A v_n \end{bmatrix} \in \mathbb{R}^{m-1 \times n-1}$$

$$\left\| \underbrace{\begin{bmatrix} \sigma_1 & w^T \\ 0 & B \end{bmatrix}}_{\color{green}{``\Sigma_1}} \begin{bmatrix} \sigma_1 \\ w \end{bmatrix} \right\|_2 \geq \sigma_1^2 + w^T w$$

$$= \sqrt{\sigma_1^2 + \|w\|^2} \left\| \begin{bmatrix} \sigma_1 \\ w \end{bmatrix} \right\|_2$$

$$\Rightarrow \|\Sigma_1\|_2 \geq \sqrt{\sigma_1^2 + \|w\|^2} \quad —① $$

Since $U_1, V_1$ are orthogonal,

$$\|\Sigma_1\|_2 = \|A\|_2 = \sigma_1 \quad —②$$

From ① & ②, we can conclude
that $w = 0$, i.e.,

$$U_1^T A V_1 = \left[\begin{array}{c|c} \sigma_1 & 0 \\ \hline 0 & B \end{array}\right]$$

Hence if $m=1$ or $n=1$, we are done!
In general case, we can use the
induction hypothesis:
Suppose an SVD exists for any $m-1 \times n-1$
matrix. Then the above matrix $B$
has its SVD: $B = U_2 \Sigma_2 V_2^T$
Then $A = U_1 \underbrace{\left[\begin{array}{c|c} 1 & 0 \\ \hline 0 & U_2 \end{array}\right]}_{U} \underbrace{\left[\begin{array}{c|c} \sigma_1 & 0 \\ \hline 0 & \Sigma_2 \end{array}\right]}_{\Sigma} \underbrace{\left[\begin{array}{c|c} 1 & 0 \\ \hline 0 & V_2 \end{array}\right]^T}_{V^T} V_1^T$

This is an SVD of $A$!  ///

(Proof: Uniqueness)
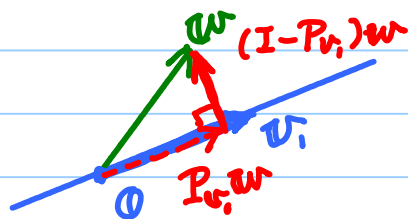Let $v_1 \in S \subset \mathbb{R}^n$ s.t.
$$\| A \|_2 = \| \tilde{u}_1 \|_2 = \| A v_1 \|_2 = \sigma_1$$
Suppose $\exists w \in S$, s.t., $w \neq v_1$,
$w$ is linearly independent from $v_1$,
and $\| A w \|_2 = \sigma_1$.
Let's define a unit vector $v_2 \in S$ by
$$v_2 := \frac{(I - P_{v_1}) w}{\| (I - P_{v_1}) w \|_2}$$

$v_2 \perp v_1$

Since $\|A\|_2 = \sigma_1$, by definition
$$\|A v_2\|_2 \leq \sigma_1 \quad ---- (a)$$
We now claim $\|A v_2\|_2 = \sigma_1$.
why? Because $w = P_{v_1} w + (I - P_{v_1}) w$
$$= c v_1 + s v_2$$

where $c, s$: constants satisfying $c^2 + s^2 = 1 \cdots (b)$
$$\sigma_1^2 = \|A w\|_2^2 = \|c A v_1 + s A v_2\|_2^2$$
$$= c^2 \|A v_1\|_2^2 + 2cs \underbrace{(A v_1)^T A v_2}_{= 0} + s^2 \|A v_2\|_2^2$$
$$= c^2 \sigma_1^2 + s^2 \|A v_2\|_2^2 \underset{(a)}{\leq} c^2 \sigma_1^2 + s^2 \sigma_1^2 \underset{(b)}{=} \sigma_1^2$$

This means that the inequality above
must be an equality, and hence $\|A v_2\|_2 = \sigma_1$ //

Hence, what we have proved is :
if $v_1$ is not unique, then the corresp.
singular value $\sigma_1$ is not simple
(i.e., has some multiplicity).
After determining $\sigma_1, u_1, v_1,$
we can use the induction argument.
In particular, for $A$ : square, $\{\sigma_j\}$ are
distinct (no multiple singular values),
then it's clear that $\{u_j\}, \{v_j\}$
are uniquely determined up to signs.

# More about SVD!

$*$ <u>"A Change of Bases" viewpoint</u>

$$A = U \Sigma V^T \in \mathbb{R}^{m \times n}$$

Pick any $x \in \mathbb{R}^n$ and consider

$$\tilde{x} = V^T x$$

Then $\tilde{x}$ is the expansion coefficient of $x$ w.r.t. the ONB $\{v_1, \cdots, v_n\}$

<u>why?</u> You should know this by now. But, just in case,

$$\tilde{x} = V^T x \iff x = V \tilde{x}$$

$$= \tilde{x}_1 v_1 + \cdots + \tilde{x}_n v_n$$

<span style="color:red">linear comb. of $\{v_1, \cdots, v_n\}$.</span> //

Now, let $b = A x \in \mathbb{R}^m$

Expand $b$ w.r.t. the ONB $\{u_1, \cdots, u_m\}$

$$\hat{b} = U^T b = U^T A x = U^T A V \tilde{x}$$

$$= \underbrace{U^T U}_{\color{red}{= I_m}} \Sigma \underbrace{V^T V}_{\color{red}{\text{"}I_n}} \tilde{x} = \Sigma \tilde{x}$$

Now, we know that $\Sigma$ is diagonal!

This again shows that

<span style="color:red">"$\Sigma$ represents the essence of $A$ in a much clearer manner!"</span>

# ✱ SVD vs Eigenvalue Decomposition

Let $A \in \mathbb{R}^{m \times m}$ be diagonalizable, i.e., $\exists$ the eigenvalue decomposition:

$$A = X \wedge X^{-1}$$

Note: where $X = [x_1 \cdots x_m] \in \mathbb{C}^{m \times m}$

Even if $A \in \mathbb{R}^{m \times m}$, satisfying

$\wedge = \text{diag}(\lambda_1, \cdots, \lambda_m) \in \mathbb{C}^{m \times m}$

its eigval's & eigvec's may be complex-valued!

$A x_j = \lambda_j x_j$, $j = 1, \cdots, m$

$$\Updownarrow$$

$$A X = X \wedge$$

Note that the eigenvectors $\{x_1, \cdots, x_m\}$ form a basis of $\mathbb{C}^m$, but not necessarily orthonormal in general

Ex. $A = \begin{bmatrix} 0 & -1 \\ 1 & 0 \end{bmatrix}$. unless $A^* = A$ (unitary)

Here $A^* := (\overline{a_{ji}}) = \overline{A}^T$

conjugate transposition of $A \in \mathbb{C}^{m \times m}$

"unitarity" is a generalization of "symmetry".

With the eigenvalue decomposition,
$b = A x$ can be simplified as

$\tilde{b} = \wedge \tilde{x}$ via $\begin{cases} \tilde{b} = X^{-1} b \\ \tilde{x} = X^{-1} x \end{cases}$

diagonal

change of bases again!

So, we can summarize as follows:
- SVD: Use two different ONB's $U, V$ and work for <u>any</u> matrix.
- EIG: Use one basis (not ONB in general) and work <u>only for square matrices</u>.

⭐ <u>Matrix Properties via SVD</u>
Let $A \in \mathbb{R}^{m \times n}$,
$$p := \min(m, n)$$
$$r := \# \text{ nonzero singular values}$$
$$\leq p.$$
<u>Thm</u>   $\text{rank}(A) = r$.

(Proof) Let $A = U \Sigma V^T$.
Since $U, V$ are orthogonal matrices, they are of full rank.
Hence,   $\text{rank}(A) = \text{rank}(\Sigma)$
$$= \# \text{ nonzero diagonal entries}$$

<span style="color:green">Recall $\langle u_1, \cdots, u_r \rangle$<br>$:= \text{span}\{u_1, \cdots, u_r\}$ ➘</span>
$$= r \quad /\!/\!/$$

<u>Thm</u>   $\text{range}(A) = \langle u_1, \cdots, u_r \rangle$
$$\text{null}(A) = \langle v_{r+1}, \cdots, v_n \rangle$$

(Proof) Since $\Sigma \in \mathbb{R}^{m \times n}$ is diagonal with only $r$ nonzero entries,

$$\text{range}(\Sigma) = \langle e_1, \cdots, e_r \rangle \subset \mathbb{R}^m$$

$\Longleftrightarrow \text{range}(A) = \langle u_1, \cdots, u_r \rangle \subset \mathbb{R}^m.$ ✓

On the other hand, it is clear that for any vector $X \in \mathbb{R}^n$ s.t.

$$X = [\underbrace{0, 0, \cdots, 0}_{r}, x_{r+1}, \cdots, x_n]^T,$$

$$\Sigma X = \begin{bmatrix} \sigma_1 & & & \\ & \ddots & & \\ & & \sigma_r & \\ & & & 0 \\ & & & & \ddots \\ & & & & & 0 \end{bmatrix} \begin{bmatrix} 0 \\ \vdots \\ 0 \\ x_{r+1} \\ \vdots \\ x_n \end{bmatrix} = \mathbb{O}.$$

So, $\text{null}(\Sigma) = \langle e_{r+1}, \cdots, e_n \rangle \subset \mathbb{R}^n$

Then, for such $X$, we have

$$A V X = U \Sigma V^T V X$$
$$= U \Sigma X = \mathbb{O}$$

i.e., any member of $\text{null}(A)$ should be of the form $V X$, $X \in \text{null}(\Sigma)$

i.e., $\text{null}(A) = \langle v_{r+1}, \cdots, v_n \rangle \subset \mathbb{R}^n$ ///

Thm $\quad \| A \|_2 = \sigma_1, \quad \| A \|_F = \sqrt{\sigma_1^2 + \cdots + \sigma_r^2}$

(Proof) Since $U, V$ are orthogonal,

$$\| A \|_2 = \| \Sigma \|_2 = \max_{1 \leq j \leq r} \{ |\sigma_j| \} = \sigma_1 \checkmark$$

The Frobenius norm is also invariant w.r.t. rotations (ortho. matrix multiplications)

Hence, $\|A\|_F = \|\Sigma\|_F = \sqrt{\sigma_1^2 + \cdots + \sigma_r^2}$ ///

**Thm** The nonzero singular values of $A$ are the square roots of the nonzero eigenvalues of $A^TA$ or $AA^T$.

(Proof) $A^TA = (U\Sigma V^T)^T(U\Sigma V^T)$
$$= V\Sigma^T\Sigma V^T$$
$$\Longleftrightarrow (A^TA)V = V(\underset{\sim\sim\sim}{\Sigma^T\Sigma})$$

$$\text{``} \operatorname{diag}(\sigma_1^2, \cdots, \sigma_r^2, 0, \cdots, 0)$$
$$\in \mathbb{R}^{n \times n}$$

So, the col's of $V$ are the eigenvectors of $A^TA$ and their nonzero eigval's are $\sigma_1^2, \cdots, \sigma_r^2$. You can show similarly that the col's of $U$ are the eigenvec's of $AA^T$, and their nonzero eigval's are $\sigma_1^2, \cdots, \sigma_r^2$. ///

**Thm** $A^T = A \Rightarrow \sigma_i(A) = |\lambda_i(A)|$

(Proof) HW #3 Prob 3 says:

<span style="color:red">Any symmetric matrix has only real-valued eigenvalues and the eigenvec's form an ONB.</span>

So, $A = Q\Lambda Q^T$, $Q$: ortho, $\Lambda$: diag
$$= Q|\Lambda|\operatorname{sgn}(\Lambda)Q^T$$

where $|\Lambda| := \begin{bmatrix} |\lambda_1| & & 0 \\ & \ddots & \\ 0 & & |\lambda_m| \end{bmatrix}$

$\text{sgn}(\Lambda) := \begin{bmatrix} \text{sgn}(\lambda_1) & & 0 \\ & \ddots & \\ 0 & & \text{sgn}(\lambda_m) \end{bmatrix}$

Now, it's clear that $Q\,\text{sgn}(\Lambda)$ is orthogonal if $Q$ is orthogonal.

<u>why?</u>

$$(Q\,\text{sgn}(\Lambda))(Q\,\text{sgn}(\Lambda))^T$$
$$= Q\,\text{sgn}(\Lambda)\,\text{sgn}(\Lambda)\,Q^T$$
$$= Q\,Q^T = I_m$$

So, $A = \underbrace{Q}_{U}\,\underbrace{|\Lambda|}_{\Sigma}\,\underbrace{(Q\,\text{sgn}(\Lambda))^T}_{V^T}$ ///

<u>Thm</u> For $A \in \mathbb{R}^{m \times m}$,
$$|\det(A)| = \prod_{i=1}^{m} \sigma_i = \sigma_1 \cdot \sigma_2 \cdots \cdot \sigma_m$$

(Proof) We'll use the following facts.

- $\det(AB) = \det(A) \cdot \det(B)$.
- $\det(A^T) = \det(A)$
- $\det(\text{diag}(a_1, \cdots, a_m)) = \prod_{i=1}^{m} a_i$
- For any $Q$: orthogonal, $|\det(Q)| = 1$.
  <u>why?</u> $\det(Q^TQ) = \det(Q^T) \cdot \det(Q) = (\det(Q))^2$
  $= \det(I) = 1$, so, $|\det(Q)| = 1$ ✓

Then, $|\det(A)| = |\det(U\Sigma V^T)| = |\det(\Sigma)|$
$$= \prod \sigma_i$$ ///

# Low Rank Approximations

Recall <u>Outer product</u> in Lecture 3.

Let $u \in \mathbb{R}^m = \mathbb{R}^{m \times 1}$,
$v \in \mathbb{R}^n = \mathbb{R}^{n \times 1}$.

Then, the outer product between $u$ and $v$ is :

$$u v^T = \begin{bmatrix} u_1 \\ \vdots \\ u_m \end{bmatrix} \begin{bmatrix} v_1 \cdots v_n \end{bmatrix} = \begin{bmatrix} u_1 v_1 & \cdots & u_1 v_n \\ \vdots & & \vdots \\ u_m v_1 & \cdots & u_m v_n \end{bmatrix}$$

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

This matrix has **rank 1** because

$$u v^T = \begin{bmatrix} v_1 u, & \cdots, & v_n u \end{bmatrix}$$

i.e., each column is just a scalar multiple of the same vector $u$.

Now SVD can be viewed as a sum of rank 1 matrices :

<u>Thm</u>  $A = \sum_{j=1}^{r} \sigma_j u_j v_j^T$ , $r = \text{rank}(A)$

(Proof) just obvious!

$$\begin{bmatrix} u_1 \cdots u_m \end{bmatrix} \begin{bmatrix} \sigma_1 & & & 0 \\ & \ddots & \sigma_{r_0} & \\ 0 & & & \ddots_0 \end{bmatrix} \begin{bmatrix} v_1^T \\ \vdots \\ v_n^T \end{bmatrix}$$ ///

Among all possible $m \times n$ matrices of rank $k$ $(k \le r)$,

$$\sum_{j=1}^{k} \sigma_j u_j v_j^T \text{ is the } \textbf{best approximation}$$

of $A$ in the following sense :

**Thm** For any $k$ with $0 \le k \le r$,

let $A_k := \sum_{j=1}^{k} \sigma_j u_j v_j^T$

If $k = p = \min(m, n)$, then define $\sigma_{k+1} = 0$. Then,

$$\| A - A_k \|_2 = \inf_{\substack{B \in \mathbb{R}^{m \times n} \\ \text{rank}(B) \le k}} \| A - B \|_2 = \sigma_{k+1}$$

(Proof)

$$\| A - A_k \|_2 = \left\| \sum_{j=k+1}^{p} \sigma_j u_j v_j^T \right\|_2$$

$$= \left\| U \begin{bmatrix} 0 & & & & O \\ & \ddots & O & & \\ & & \sigma_{k+1} & & \\ & O & & \ddots & \\ & & & & \sigma_p \\ & & & & & O \end{bmatrix} V^T \right\|_2$$

$$= \left\| \begin{bmatrix} 0 & & & & O \\ & \ddots & O & & \\ & & \sigma_{k+1} & & \\ & O & & \ddots & \\ & & & & \sigma_p \\ & & & & & O \end{bmatrix} \right\|_2 \qquad \textcolor{red}{\text{since } U, V : \text{orthogonal!}}$$

$$= \sigma_{k+1} \quad \text{by definition of the matrix norm}$$

<span style="color:green">Prove this ⇒ as an exercise!</span> **Note:** If $D = \text{diag}(d_1, \cdots, d_m) = \begin{bmatrix} d_1 & & O \\ & \ddots & \\ O & & d_m \end{bmatrix}$

then $\| D \|_p = \max_{1 \le j \le m} |d_j| \quad \forall p \ge 1$

Now, let $B \in \mathbb{R}^{m \times n}$ be any rank $k$ matrix. Then $\underline{\dim(\text{null}(B)) = n - k}$ why? Because of the following Thm:

<span style="color:red">For any $A \in \mathbb{R}^{m \times n}$, $\text{rank}(A) + \dim(\text{null}(A)) = n$</span>

Let $W := \text{null}(B) \cap \langle v_1, \cdots, v_{k+1} \rangle$

We know $W \neq \{0\}$ because

  $\dim(\text{null}(B)) = n - k$

  $\dim(\langle v_1, \cdots, v_{k+1} \rangle) = k+1$

  So, if these two do not intersect, $\mathbb{R}^n$'s

  dimension would become $n - k + k + 1 = n+1$

  This cannot happen! #

So let $h \in W$, $h \neq 0$.

We can always normalize $h$, so

can assume $\|h\|_2 = 1$.

Then,

$$\|A - B\|_2 \geq \|(A-B)h\|_2 \quad \text{by def.}$$

$$\overset{\textcircled{=}}{=} \|Ah\|_2 \quad \boxed{\text{since } h \in \text{null}(B)}$$

$$= \|U \Sigma V^T h\|_2$$

$\boxed{\begin{array}{l} \text{Since } h \in \langle v_1, \cdots v_{k+1} \rangle \\ V^T h = \begin{bmatrix} * \\ \vdots \\ * \\ 0 \\ \vdots \\ 0 \end{bmatrix} \begin{array}{l} \}k+1 \\ \\ \}n-k-1 \end{array} \end{array}}$ $$= \|\Sigma V^T h\|_2 \quad \begin{array}{l} \text{since } U: \\ \text{ortho.} \end{array}$$

$$\overset{\textcircled{\geq}}{\geq} \sigma_{k+1} \|V^T h\|_2$$

$$= \sigma_{k+1} \|h\|_2 = \sigma_{k+1} \quad /\!/\!/$$

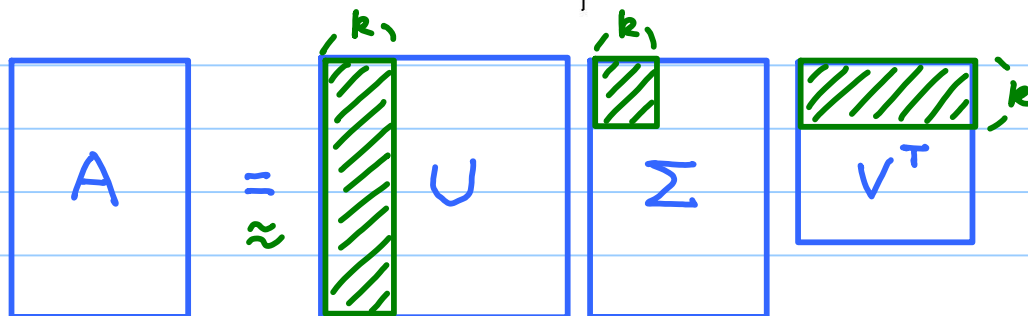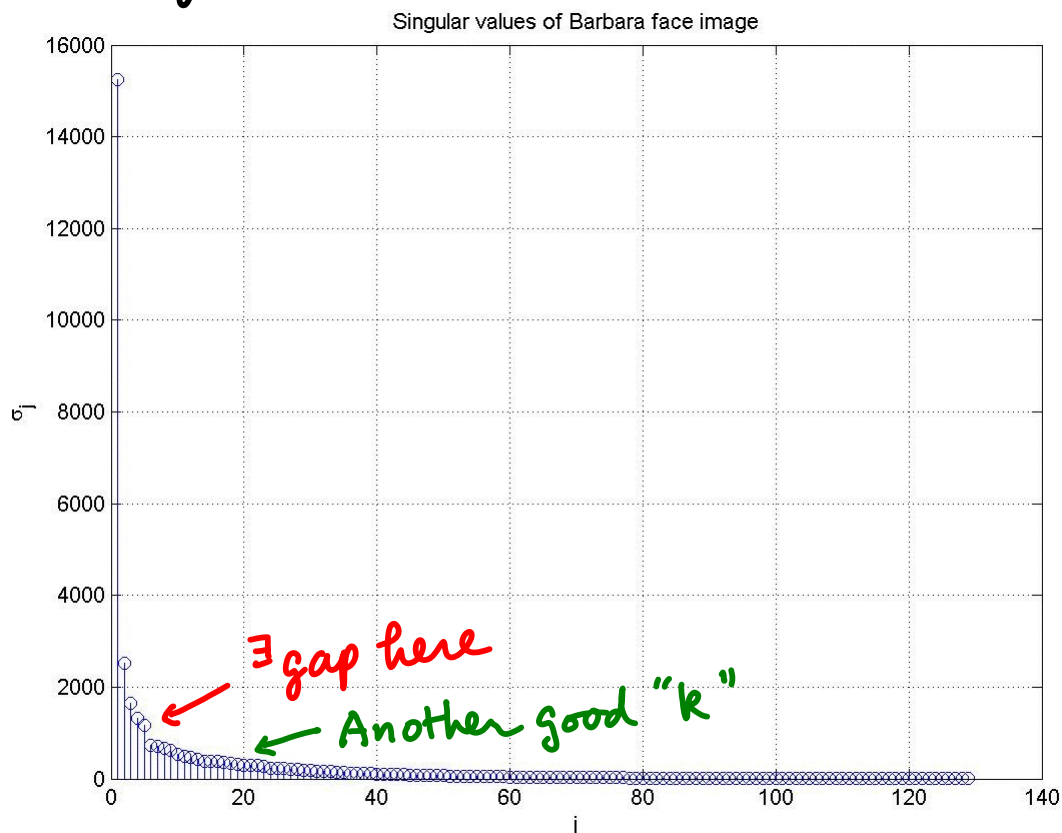**Thm** For any $k$ with $0 \leq k \leq r$,

$$\|A - A_k\|_F = \inf_{\substack{B \in \mathbb{R}^{m \times n} \\ \text{rank}(B) \leq k}} \|A - B\|_F$$

$$= \sqrt{\sigma_{k+1}^2 + \cdots + \sigma_r^2}$$

(Proof) Exercise!

So, for a given matrix, say, A how to determine a good "k" so that we can efficiently (i.e., compress) A without losing too much info of A?

$\Rightarrow$ Check the distribution of the singular values!

**Singular values of Barbara face image**



∃ gap here

← Another good "k"

rank k approximation of A only uses ///// portions!

## ☆ Condition Number and SVD

Recall the condition number for a square nonsingular matrix $A$:

$$\kappa(A) = \text{cond}(A) := \|A\|_2 \cdot \|A^{-1}\|_2$$

$\kappa(A)$: small $\Rightarrow$ $A$: well-conditioned.

$\kappa(A)$: large $\Rightarrow$ $A$: ill-conditioned,

$$\text{lose} \approx \log_{10}\kappa(A) \text{ digits}$$
$$\text{to solve } Ax = b.$$

If $A$: singular, $\kappa(A) = +\infty$.

Using SVD of $A$, we can nicely compute $\kappa(A)$ as follows.

$\|A\|_2 = \sigma_1$  $\rightarrow$ by definition

$\|A^{-1}\|_2 = 1/\sigma_m$  <u>why?</u> $A^{-1} = (U\Sigma V^T)^{-1} = V\Sigma^{-1}U^T$

$$= V \, \text{diag}(1/\sigma_1, \cdots, \underbrace{1/\sigma_m}_{\color{red}\text{largest}}) \, U^T$$

So, $\kappa(A) = \sigma_1/\sigma_m$

We can <span style="color:red">generalize</span> the definition of the <span style="color:red">condition number</span> for a rectangular matrix $A \in \mathbb{R}^{m \times n}$ using the pseudo-inverse $A^{\dagger}$ and SVDs as

$$\color{red}\kappa(A) := \|A\|_2 \cdot \|A^{\dagger}\|_2$$

$$\color{red}= \sigma_1/\sigma_r \qquad \begin{array}{l} r = \text{rank}(A) \\ \leq \min(m, n) \end{array}$$

# SVD and Least Squares Problems

## ⭐ LS via SVD

Recall the LS solution via
QR factorization:

$\begin{cases} \text{(1) Compute reduced QR of } A. \\ \text{(2) Compute } y = \hat{Q}^T b. \\ \text{(3) Solve } \hat{R} x = y \quad - (*) \end{cases}$

If $A$: full rank, then $\hat{R}_{ii} \neq 0$, $1 \leq i \leq n$,
and the triangular system (*) has a
unique LS solution.

Now using the reduced SVD of $A$,
i.e., $A = \hat{U} \hat{\Sigma} V^T$, we can also solve
the normal eqn:

$$A^T A x = A^T b$$
$$\Leftrightarrow (\hat{U} \hat{\Sigma} V^T)^T (\hat{U} \hat{\Sigma} V^T) x = (\hat{U} \hat{\Sigma} V^T)^T b$$
$$\Leftrightarrow V \hat{\Sigma}^T \hat{U}^T \hat{U} \hat{\Sigma} V^T x = V \hat{\Sigma} \hat{U}^T b$$
$$\Leftrightarrow V \hat{\Sigma}^T \hat{\Sigma} V^T x = V \hat{\Sigma}^T \hat{U}^T b$$
$$\Leftrightarrow \hat{\Sigma}^T \hat{\Sigma} V^T x = \hat{\Sigma}^T \hat{U}^T b \quad \text{since } V: \text{ortho.}$$
$$\Leftrightarrow \textcolor{red}{\hat{\Sigma} V^T x = \hat{U}^T b} \quad \begin{array}{l} \text{if } A: \text{ full rank,} \\ \text{i.e., } \sigma_j > 0, \; 1 \leq j \leq n \end{array}$$

This can be solved easily.

$\begin{cases} \text{(1) Compute reduced SVD of } A. \\ \text{(2) Compute } y = \hat{U}^T b. \\ \text{(3) Solve } \hat{\Sigma} w = y. \quad - (**) \\ \text{(4) Set } x = V w. \end{cases}$

<u>Note</u>: (**) is a diagonal system,
easier to solve than (*) !!

## ★ Pseudo inverse and SVD

Recall that if $A \in \mathbb{R}^{m \times n}$ is full rank,

$\underline{m > n :} \quad A^\dagger = (A^T A)^{-1} A^T$

$\underline{m = n :} \quad A^\dagger = A^{-1}$

$\underline{m < n :} \quad A^\dagger = A^T (A A^T)^{-1}$

However, we can define the pseudo inv. using SVD even if $A$ is not full rank!

$$A = U \Sigma V^T, \qquad \Sigma = \begin{bmatrix} \sigma_1 & 0 & 0 \\ 0 & \sigma_r & \\ 0 & 0 \end{bmatrix}$$

Define

$$A^\dagger := V \Sigma^\dagger U^T, \qquad \Sigma^\dagger := \begin{bmatrix} 1/\sigma_1 & 0 & 0 \\ 0 & 1/\sigma_r & \\ 0 & 0 \end{bmatrix}$$

As we discussed before, $A^\dagger$ satisfies the following <span style="color:red">Moore-Penrose conditions</span>:

(i) $A X A = A$; (ii) $X A X = X$

(iii) $(A X)^T = A X$; (iv) $(X A)^T = X A$.

Such $X$ is uniquely determined and $X = A^\dagger$ !!

## ✲ Pseudoinverse & Orthogonal Projectors

**Thm** $AA^+$ is an ortho. proj. onto range$(A)$

and $AA^+ = U_r U_r^T$

$A^+A$ is an ortho. proj. onto range$(A^T)$

and $A^+A = V_r V_r^T$

where $U_r \in \mathbb{R}^{m \times r}$, $V_r \in \mathbb{R}^{n \times r}$ consist
of the first $r$ columns of $U, V$, respectively.
$r = \text{rank}(A)$.

**(Proof)** Let $P_A := AA^+$, $P_{A^T} := A^+A$.

Now, $P_A = U \Sigma V^T V \Sigma^+ U^T$

$$= U \Sigma \Sigma^+ U^T = U \left[\begin{array}{c|c} I_r & 0 \\ \hline 0 & 0 \end{array}\right] U^T$$

$$= U_r U_r^T \quad \checkmark$$

$$P_A^2 = U_r \underbrace{U_r^T U_r}_{= I_r} U_r^T = U_r U_r^T = P_A \checkmark$$

so it's a proj.!

$$P_A^T = (U_r U_r^T)^T = (U_r^T)^T U_r^T = U_r U_r^T = P_A \checkmark$$

So it's an ortho. proj.!

Finally, it's also clear that
$P_A$ maps onto range$(A)$ since
range$(A) = \langle u_1, \dots, u_r \rangle$. $\checkmark$
You can do similarly for $P_{A^T}$ ///

<u>Note</u>: Consider any $X \in$ range$(A)$.
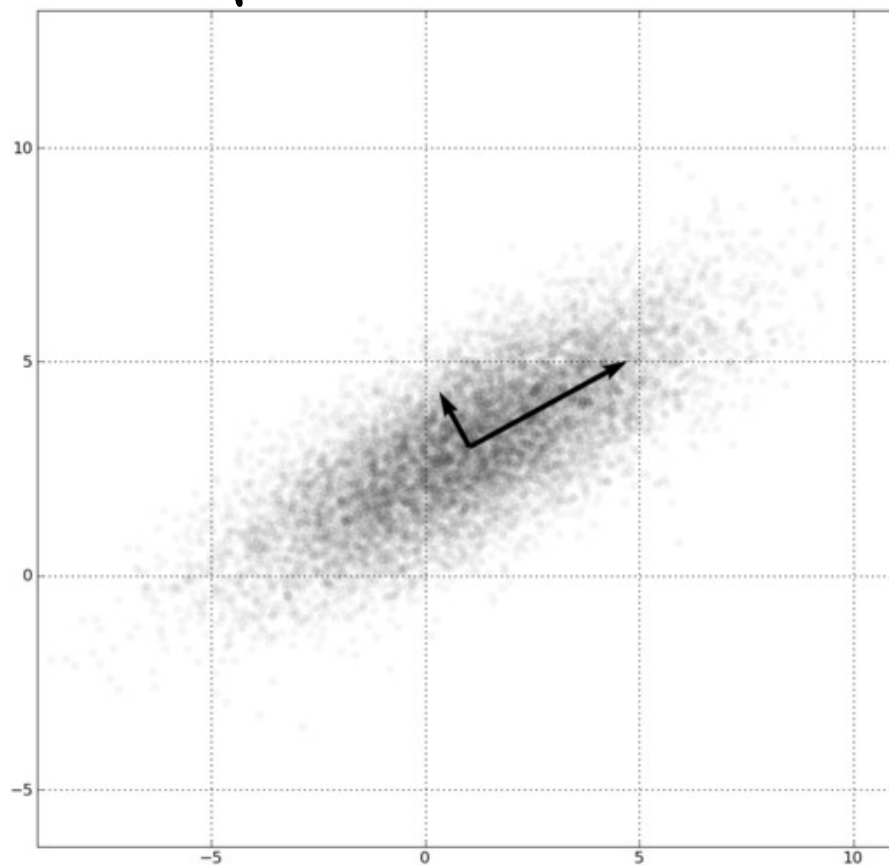Then $\exists y \in \mathbb{R}^n$ s.t. $X = Ay$.
Now $P_A X = AA^+ X = \underline{AA^+A} y$

$$= Ay = X. \quad \text{"}A\text{ via}$$

Moore-Penrose (i)

# ☆ Principal Component Analysis (PCA)

(a.k.a. Karhunen-Loève Transform)
is a data analysis technique that
uses an orthogonal transformation to
convert a set of observations of possibly
correlated variables into a set of
linearly uncorrelated variables called
"principal components."

2D example (from Wikipedia)



One can understand PCA using
SVD! But before doing so, we need
a bit of statistics.

Suppose we are given a set of vectors (observations)

often these → are viewed as $n$ realizations of some stochastic process.

$$X_1, X_2, \cdots, X_n$$

and each $X_j \in \mathbb{R}^d$.   $d$: could be huge (ex. a face image database).

Let $X := [X_1 \; X_2 \cdots X_n] \in \mathbb{R}^{d \times n}$

You know the mean (or average) of this data set

$$\overline{X} := \frac{1}{n} \sum_{j=1}^{n} X_j$$

And define the **centered** data matrix

$$\tilde{X} := [X_1 - \overline{X} \quad X_2 - \overline{X} \quad \cdots \quad X_n - \overline{X}]$$

<u>Note</u>:   $\tilde{X} = X\left(I_n - \frac{1}{n} \mathbb{1}_n \mathbb{1}_n^T\right)$

↳ Good exercise!

Now the **sample covariance matrix** $S$ is defined as

$$S := \frac{1}{n} \tilde{X} \tilde{X}^T \quad \in \mathbb{R}^{d \times d}$$

$S_{ij}$ indicates the **covariance** or **mutual correlation** between the $i$th and $j$th entries of data vectors.

**PCA is nothing but an eigenvalue decomposition of $S$, i.e.,**

$$S = \Phi \Lambda \Phi^T, \quad \Lambda = \text{diag}(\lambda_1, \cdots, \lambda_d)$$

Let's sort $\lambda_i$'s as $\lambda_1 \geq \lambda_2 \geq \cdots \geq \lambda_d$
Because $S^T = S$, and $S = \frac{1}{n} \widetilde{X} \widetilde{X}^T$,
we can show that $\lambda_i \geq 0$, $1 \leq i \leq d$.
$$\Phi = [\phi_1 \cdots \phi_d] \in \mathbb{R}^{d \times d}$$
is a matrix containing the eigenvectors.
Also thanks to $S^T = S$, $\Phi$ is an
orthogonal matrix whose columns
form an ONB of $\mathbb{R}^d$.
The change of the bases from
$[e_1 \cdots e_d]$ to $[\phi_1 \cdots \phi_d]$
is achieved simply by $\Phi^T \widetilde{X}$.

$\phi_j^T \widetilde{X}$ is called <span style="color:red">the $j$th principal components</span> of $X$.

PCA was known for a long time,
e.g., since the time of Pearson (1901)
and Hotelling (1933).
Those days, the measurement dimension
$d$ was much smaller than the number
of samples $n$, i.e. $d \ll n$
This is called the "classical" setting.
Ex. 5 exam scores of 2000 students
$\quad d = 5$, $n = 2000$.
Due to the advent of computers and
sensor technology, now we often have
$d \gg n$, the "neo-classical" setting.
$\quad$ Ex. The face database: $d = 128^2$, $n = 143$.

# PCA & SVD

Recall the **centered data matrix**
$$\widehat{X} := [\widetilde{x}_1 \cdots \widetilde{x}_n] \in \mathbb{R}^{d \times n}$$
$$\widetilde{x}_j := x_j - \overline{x}, \quad \overline{x} := \frac{1}{n}\sum_{j=1}^{n} x_i,$$

and the **sample covariance matrix**
$$S := \frac{1}{n} \widetilde{X}\widetilde{X}^T$$

Then, **PCA** is nothing but **the eigendecomposition of S**
$$S = \Phi \Lambda \Phi^T, \quad \Lambda = \text{diag}(\lambda_1, \cdots, \lambda_d)$$

$\lambda_1 \geq \lambda_2 \geq \cdots \geq \lambda_d \geq 0.$
$\Phi := [\phi_1 \cdots \phi_d] \in \mathbb{R}^{d \times d}$ is
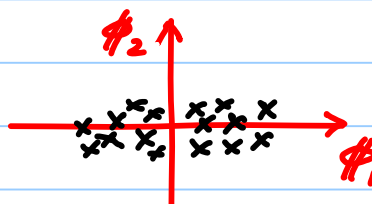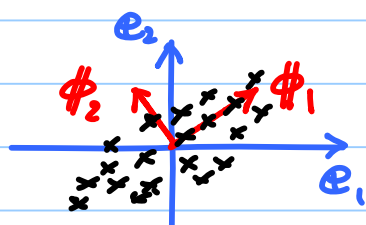an ortho. matrix, and $\{\phi_1, \cdots, \phi_d\}$
form an ONB of $\mathbb{R}^d$.
$\phi_j^T \widetilde{X}$ is said to be **the $j$ th**

**principal components** of $\widetilde{X}$.
These are nothing but the expansion
coefficients of $\widetilde{X}$ w.r.t. the ONB
vector $\phi_j$.

If $\widehat{X}$ forms a    then $\phi_j^T \widetilde{X}$ are the
"cigar" shape,    coordinate values of $\widetilde{X}$
            under the rotated axes

- Hence viewing the given dataset under the principal axes $\#_1, \#_2, \cdots,$ provides us better interpretations of the data than viewing them under the original axes $e_1, e_2, \cdots$.

- PCA is also often used as a tool to do <span style="color:red">dimension reduction</span> and <span style="color:red">feature extraction</span> by keeping only <span style="color:blue">top $k$</span> PCA coordinates where $k \ll d$, i.e.,

$$\Phi_k := [\, \#_1 \cdots \#_k \,] \in \mathbb{R}^{d \times k}$$

$$\mathbb{R}^d \ni \tilde{x}_j \longmapsto \Phi_k^T \tilde{x}_j \in \mathbb{R}^k$$

<span style="color:blue">top $k$</span> PCA coordinates or <span style="color:blue">top $k$</span> Principal components of $x_j$.

Note that using these <span style="color:blue">top $k$</span> principal components, we can approximate the original data $x_j$ by

$$x_j \approx \bar{x} + \Phi_k \Phi_k^T \tilde{x}_j$$

Of course the approximation gets better and better as $k$ increases. In fact, if $k = d$, then $x_j$ is recovered exactly (within machine $\varepsilon$).

Now we'll face the problem when we compute the eigendecomposition of $S = \Phi \Lambda \Phi^T$:

(1) If $d$ is large, we cannot compute this eigendecomposition because we cannot hold $\Phi \in \mathbb{R}^{d \times d}$ in computer memory, and its computational cost is $O(d^3)$, i.e., too expensive to compute.

(2) Fortunately, we often do not need all $d$ eigenvectors, most likely, only first $k$ eigenvectors $k \ll d$.

(3) Moreover if $d > n$, then rank$(S) = n-1$ if $x_j$'s are linearly indep. So, after the first $n-1$ eigenvectors are useless!

Why?    $S = \frac{1}{n} \tilde{X} \tilde{X}^T = \frac{1}{n} \{ \underbrace{\tilde{x}_1 \tilde{x}_1^T}_{\text{rank 1}} + \cdots + \underbrace{\tilde{x}_n \tilde{x}_n^T}_{\text{rank 1}} \}$

So looks like rank$(S) = n$.
But since $\tilde{x}_1 + \cdots + \tilde{x}_n = 0$ because the mean $\bar{X}$ is subtracted from each data vector $x_j$ (i.e., $\tilde{x}_j = x_j - \bar{X}$)
Hence, $S$ loses 1 rank.
So, rank$(S) = n-1$.

Now, let's consider the reduced
SVD of $\tilde{X}$ :
$$\tilde{X} = \hat{U} \hat{\Sigma} V^{T}$$

$d \geq n$         $d < n$



$\tilde{X} \quad \hat{U} \quad \hat{\Sigma} \quad V^{T}$      $\tilde{X} \quad U \quad \hat{\Sigma} \quad \hat{V}^{T}$

Just consider the "neo-classical" setting,
i.e., $d \geq n$ (e.g., the face image database)

Then consider the sample covariance
matrix $S$ using the above SVD:
$$S = \frac{1}{n} \tilde{X} \tilde{X}^{T} = \frac{1}{n} \hat{U} \hat{\Sigma} V^{T} V \hat{\Sigma}^{T} \hat{U}^{T}$$

$$= \frac{1}{n} \hat{U} \hat{\Sigma} \hat{\Sigma}^{T} \hat{U}^{T} = \frac{1}{n} \hat{U} \hat{\Sigma}^{2} \hat{U}^{T}$$

Now $\hat{\Sigma} = \text{diag}(\sigma_1, \cdots, \sigma_{n-1}, \underline{0})$
if $X_1, \cdots, X_n$ are linearly indep.
So, $\hat{\Sigma}^{2} = \text{diag}(\sigma_1^2, \cdots, \sigma_{n-1}^2, 0)$.

Finally, $S$ can be written as
$$S = \hat{U} \left( \frac{1}{n} \hat{\Sigma}^{2} \right) \hat{U}^{T}$$
$\uparrow$    $= \text{diag}(\sigma_1^2/n, \cdots, \sigma_{n-1}^2/n, 0)$

columns are orthonormal.

Comparing this with the eigendecomposition

$S = \Phi \Lambda \Phi^T$, we can conclude that

$$\begin{cases} \Phi(:,1:n) = \hat{U} \\ \Lambda(1:n,1:n) = \frac{1}{n}\hat{\Sigma}^2 = diag(\sigma_1^2/n, \cdots, \sigma_{n-1}^2/n, 0) \end{cases}$$

In fact, only the $1:n-1$ portion is useful since $\sigma_n = 0$.

Hence, we should <span style="color:red">use the reduced SVD of</span> <span style="color:blue">$\tilde{X}$</span> <span style="color:red">(not $S$) for computing PCA!!</span> Do not use the eigendecomposition of $S$ unless $d$ is small.

<u>Note</u>: $\tilde{X}V = \hat{U}\hat{\Sigma}V^TV = \hat{U}\hat{\Sigma}$

$\qquad\qquad\qquad = [\sigma_1 u_1, \cdots \sigma_{n-1}u_{n-1}, 0]$

$= [\tilde{X}v_1, \cdots, \tilde{X}v_n]$

So, $\quad u_j = \frac{1}{\sigma_j}\tilde{X}v_j$, $\quad j=1,\cdots,n-1$.

In other words, <span style="color:red">each principal axis $u_j$ is just a linear combination of the (centered) input vectors $\tilde{X}_1, \cdots, \tilde{X}_n$!</span>

Now let's do MATLAB experiments using the face image database consisting of 143 faces each of which has $128 \times 128 = 16384$ pixels, i.e., $d = 16384$, $n = 143$.