# SVD and Least Squares Problems

**★ LS via SVD**

Recall the LS solution via QR factorization:

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

If $A$: full rank, then $\hat{R}_{ii} \neq 0$, $1 \leq i \leq n$, and the <u>triangular system</u> (*) 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:

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

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 <u>diagonal system</u>, easier to solve than (*) !!

## ★ Pseudo inverse and SVD

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

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

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

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

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

$$A = U \Sigma V^T, \qquad \Sigma = \begin{array}{c} \overbrace{\quad}^{r} \overbrace{\quad}^{n-r} \\ \left[ \begin{array}{cc|c} \sigma_1 \; \begin{smallmatrix} \ddots \end{smallmatrix} \; 0 & 0 \\ 0 \quad \sigma_r & \\ \hline 0 & 0 \end{array} \right] \begin{array}{c} \} r \\ \\ \} m-r \end{array} \end{array}$$

Define

$$A^\dagger := V \Sigma^\dagger U^T, \qquad \Sigma^\dagger := \begin{array}{c} \overbrace{\quad}^{r} \overbrace{\quad}^{m-r} \\ \left[ \begin{array}{cc|c} 1/\sigma_1 \; \begin{smallmatrix} \ddots \end{smallmatrix} \; 0 & 0 \\ 0 \quad 1/\sigma_r & \\ \hline 0 & 0 \end{array} \right] \begin{array}{c} \} r \\ \\ \} n-r \end{array} \end{array}$$

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** $AA^{\dagger}$ is an ortho. proj. onto range$(A)$

and $AA^{\dagger} = U_r U_r^T$

$A^{\dagger}A$ is an ortho. proj. onto range$(A^T)$

and $A^{\dagger}A = V_r V_r^T$

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 $P_A := AA^{\dagger}$, $P_{A^T} := A^{\dagger}A$.

Now, $P_A = U\Sigma V^T V \Sigma^{\dagger} U^T$

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

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

$$P_A^2 = U_r \underbrace{U_r^T U_r}_{= I_r} U_r^T = U_r U_r^T = P_A \quad \checkmark \quad \text{so it's a proj.!}$$

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

So it's an ortho. proj.!

Finally, it's also clear that $P_A$ maps onto range$(A)$ since range$(A) = \langle u_1, \cdots, u_r \rangle$. $\checkmark$

You can do similarly for $P_{A^T}$ ///

Note: Consider any $X \in$ range$(A)$.
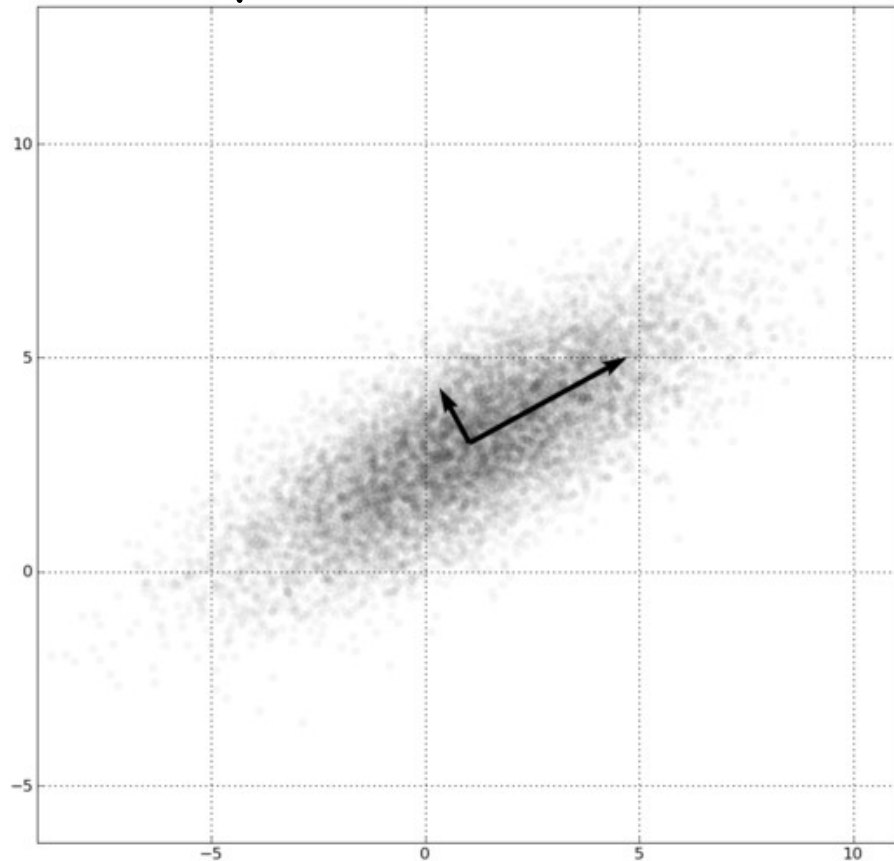
Then $\exists y \in \mathbb{R}^n$ s.t. $X = Ay$.

Now $P_A X = AA^{\dagger} X = AA^{\dagger}Ay$

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

Moore-Penrose (i)

☆ <u>Principal Component Analysis (PCA)</u>
(a.k.a.   Karhunen-Loève Transform)
is a data analysis technique that
uses an <u>orthogonal transformation</u> to
convert a set of observations of possibly
<u>correlated</u> variables into a set of
<u>linearly uncorrelated</u> 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)

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

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

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

and define the **centered** data matrix

$$\tilde{X} := [X_1 - \bar{X} \;\; X_2 - \bar{X} \cdots X_n - \bar{X}]$$

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

<span style="color:green">↳ Good exercise!</span>

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

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

$S_{ij}$ 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 $S = \frac{1}{n} \tilde{X} \tilde{X}^T$,
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 $\Phi^T \tilde{X}$.

$\phi_j^T \tilde{X}$ 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
    $d = 5$, $n = 2000$.
Due to the advent of computers and
sensor technology, now we often have
$d \gg n$, the "neo-classical" setting.
    Ex. The face database: $d = 128^2$, $n = 143$.