# DIMENSIONALITY REDUCTION

D input features

N data points

Data matrix

dimensionality reduction

$D'$

N

$D'=4$

feature selection

$x = \{(x_i)\}_{i=1}^{N}$ where $x_i \in \mathbb{R}^D$

we will select a subset of $\{1, 2, \ldots, D\}$
size of K

$K \ll D$

feature extraction

$x_i \in \mathbb{R}^D \Rightarrow z_i \in \mathbb{R}^K$
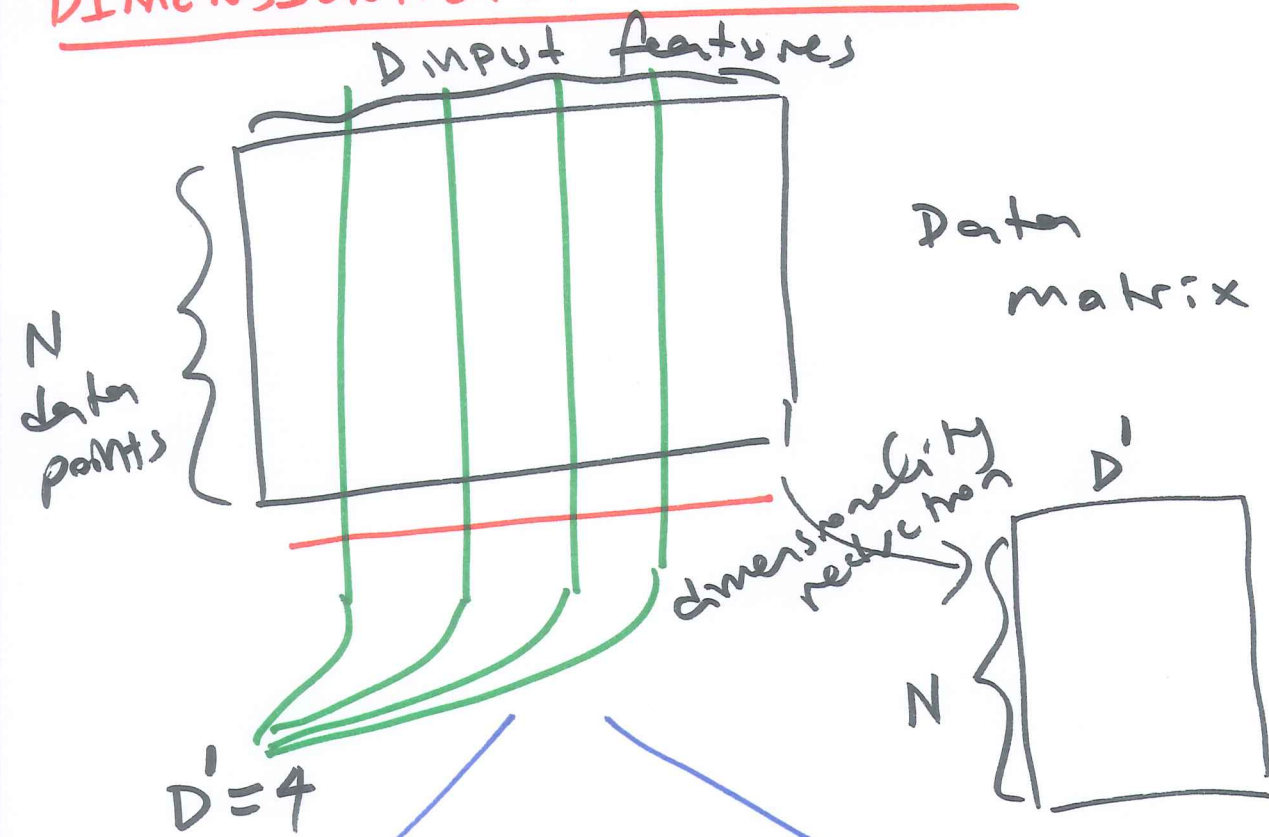
eliminate

[ X ]  NxD
[ x' ]

[ X ] #4  DxK
$\Rightarrow$ [ Z ]
① [ z' ] NxK

Reasons

1) To reduce computational complexity.

2) To reduce storage complexity.

3) To reduce data acquisition cost

4) To increase robustness

5) To increase interpretability.

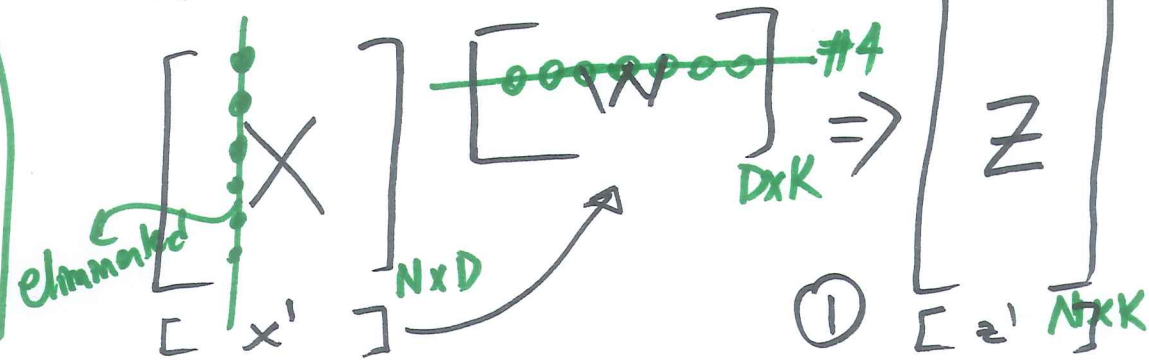6) To enable visualization (2D or 3D)

...

# SUBSET SELECTION

$$F : \{1, 2, \ldots, D\} \rightarrow \hat{F} : \{\ ?\ \}$$

$\hookrightarrow$ # of possible subsets : $\underline{2^D - 1 - 1}$

## Forward selection:

$\hat{F} = \emptyset$

At each iteration, find the best
new feature to be <u>added</u> to $\hat{F}$ $\Big\}$ $d^* = \arg\min\limits_{d} Error(\hat{F} \cup d)$    union

Add $d^*$ to $\hat{F}$ if

$$Error(\hat{F} \cup d^*) < Error(\hat{F})$$

$t=1 \Rightarrow$ ① 2  3  4  5  6

$t=2 \Rightarrow$ 1-2 , 1-3, 1-4 , 1-5, 1-6

$t=3 \Rightarrow$ 1-4-2, 1-4-3, 1-4-5, 1-4-6

$t=4 \Rightarrow$ 1-4-5-2, 1-4-5-3, 1-4-5,6

$\hookrightarrow$ best solution

## Backward selection:

$\hat{F} = F$

At each iteration, find the best
feature to be removed from $\hat{F}$ $\Big\}$ $d^* = \arg\min\limits_{d} Error(\hat{F} / d)$

Remove $d^*$ from $\hat{F}$ if

$$Error(\hat{F} / d) < Error(\hat{F})$$

$\hookrightarrow$ set difference

$t=1 \Rightarrow 2,3,4,5,6 / 13456 /$

②

# PRINCIPAL COMPONENT ANALYSIS (PCA)

PCA is a feature extraction algorithm

$$Z = W^T \cdot X$$

$X \in \mathbb{R}^D \qquad Z \in \mathbb{R}^k \qquad W \in \mathbb{R}^{D \times K}$

$K \times 1 \qquad K \times D \qquad D \times 1$

$$\underline{\underline{K < D}}$$

$$\left[ \begin{array}{cccc} \downarrow & \downarrow & \cdots & \downarrow \end{array} \right]_{D \times K}$$

$w_1 \; w_2 \quad w_k \qquad \underline{\underline{K=1}}$

We would like to find the direction that maximizes the variance.

$$VAR(Z_1) = VAR(\underline{\underline{W_1^T}} \cdot X)$$

$$= W_1^T VAR(X) \cdot W_1 \longrightarrow \text{covariance matrix}$$

$$= W_1^T \cdot \underset{X}{\Sigma} \cdot W_1$$

$$\underset{1 \times D}{} \quad \underset{D \times D}{} \quad \underset{D \times 1}{} = \underline{\underline{1 \times 1}}$$

$$\boxed{VAR(\underline{a X}) = a^2 VAR(X)}$$

maximize $VAR(Z_1) = \underset{1 \times D}{W_1^T} \cdot \underset{X}{\Sigma} \cdot W_1 \longrightarrow \text{decision variables.}$

$\qquad\qquad\qquad\qquad \hookrightarrow \text{constant}$

$$\left[ \text{subject to} : \|W_1\|^2 = 1 . \right] \alpha$$

$$W^{\#} \Rightarrow 2W^{\#}$$

$$(2W^{\#})^T \cdot \Sigma_X (2W^{\#}) = 4 \left[ (W^{\#}) \Sigma_X W^{\#} \right]$$

③

$L_p:$ $w_1^T . \Sigma_x w_1 - \alpha \left( \|w_1\|^2 - 1 \right)$        $\|w\|^2 = w_1^T . w_1$

$= w_1^T . \Sigma_x w_1 - \alpha \left( w_1^T . w_1 - 1 \right)$

$$\frac{\partial L_p}{\partial \boxed{w_1}}_{D \times 1} = 2 . \underbrace{\Sigma_x}_{D \times D} . \underbrace{w_1}_{D \times 1} - 2 . \underbrace{\alpha . w_1}_{D \times 1} = 0$$

$\not{2} . \Sigma_x . w_1 = \not{2} . \alpha . w_1$

$$\boxed{Ax = \lambda x}$$        $$\boxed{\Sigma_x . w_1 = \alpha . \underline{w_1}}$$

$w_1 \Rightarrow$ is the first eigenvector of $\Sigma_x$.

$\alpha_1, \alpha_2, \alpha_3 \ldots , \alpha_D$   $\Rightarrow$   $\boxed{\alpha_1} \geqslant \boxed{\alpha_2} \geqslant \alpha_3 \geqslant \ldots \geqslant \alpha_D$

$\begin{bmatrix} z_1 \\ z_2 \end{bmatrix}$        $z_2 = w_2^T . x .$    maximize $VAR(z_2) = w_2^T . \Sigma_x w_2$

subject to: $\left. \begin{array}{l} w_2^T . w_2 = 1 \\ w_2^T . w_1 = 0 \end{array} \right\} \begin{array}{l} \alpha \\ \beta \end{array}$

$w_2^T . \Sigma_x . w_2 - \alpha \left( w_2^T . w_2 - 1 \right) - \beta w_2^T . w_1$

$w_1^T \left[ 2 . \Sigma_x . w_2 - 2 . \alpha . w_2 - \beta . w_1 \right] = 0$        $\Sigma_x . w_2 - \alpha . w_2 = 0$

$$\boxed{\Sigma_x . w_2 = \alpha . w_2}$$

$w_2$ is the second eigenvector

④

Step 1: Calculate $\Sigma_x$

Step 2: Find first K eigenvectors (W)

Projection Step: $z_i = W^T \cdot (x_i - m) \quad \forall i$

$\swarrow$ centering.

$m$ = sample mean

$$m = \frac{\sum_{i=1}^{N} x_i}{N}$$

K = ?

$\partial_1, \partial_2, \ldots, \partial_D \geq 0$

POVE = proportion of variance explained.

$$POVE(K) = \frac{\partial_1 + \partial_2 + \ldots + \partial_K}{\partial_1 + \partial_2 + \ldots + \partial_D}$$



95% variance.

$K \geq 4$