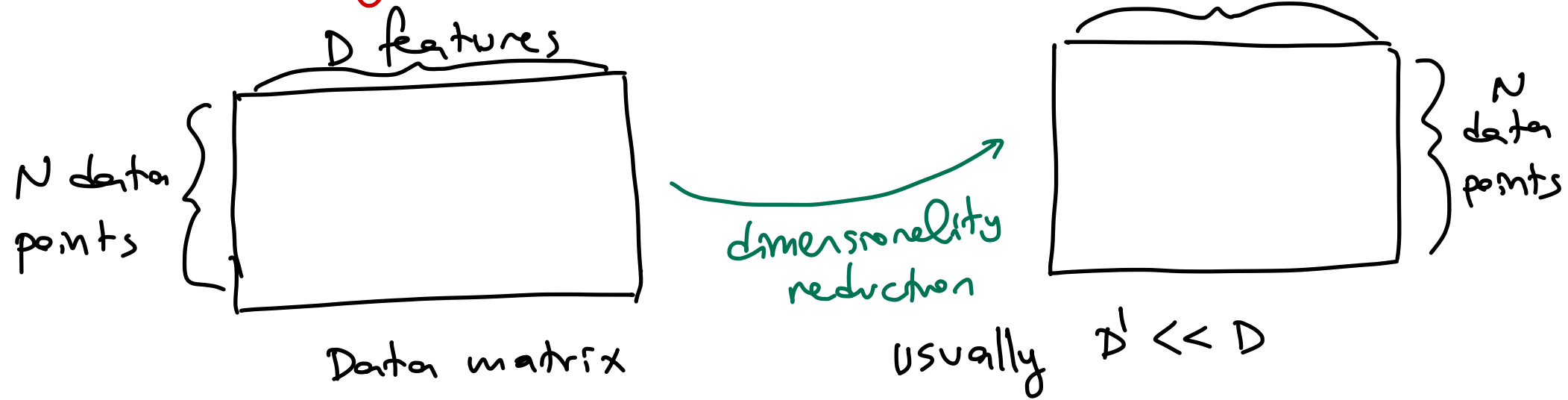


Dimensionality Reduction



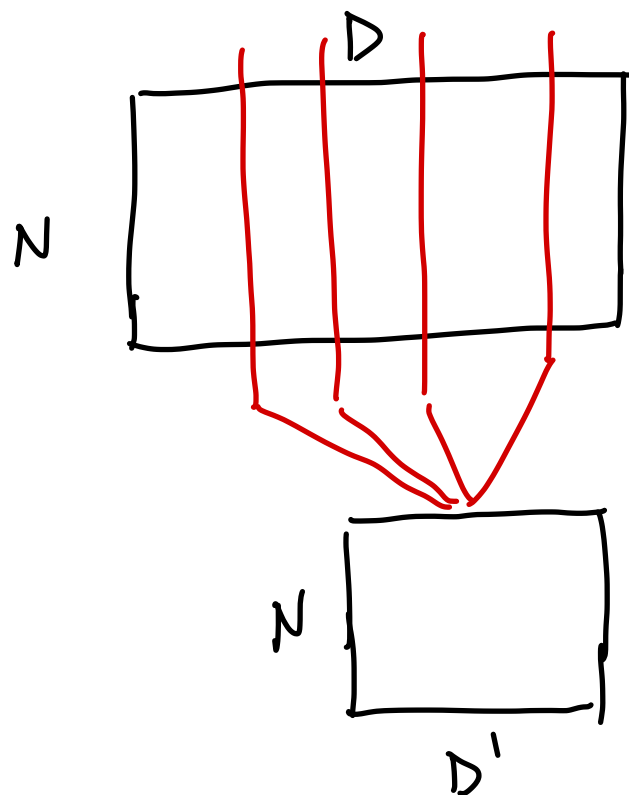
Reasons

- 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 ($D'=2$ or $D'=3$)
- ...

Feature Selection

$$\mathcal{X} = \{x_i\}_{i=1}^N \text{ where } x_i \in \mathbb{R}^D$$

we will select a subset
of $\{1, 2, \dots, D\}$.



$$F = \{1, 2, 3, \dots, D\}$$

$$F' = \{ \}$$

of possible subsets of $F = 2^D - 1 - 1$
 full set \rightarrow empty set

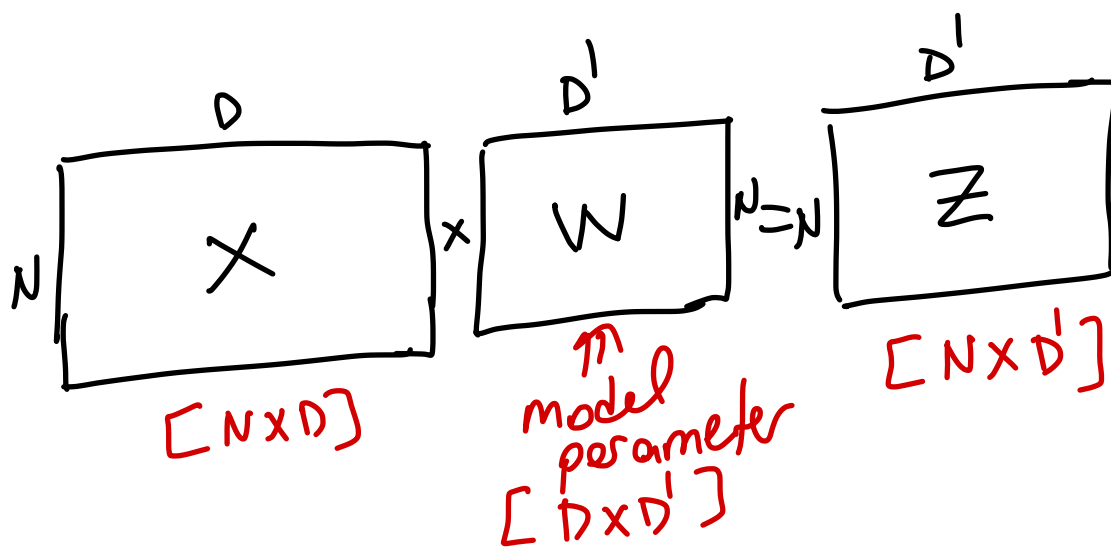
Feature Extraction

$$\mathcal{X} = \{x_i\}_{i=1}^N \text{ where } x_i \in \mathbb{R}^D$$

$$x_i \in \mathbb{R}^D \longrightarrow z_i \in \mathbb{R}^{D'}$$

$$z_i = W^T \cdot x_i$$

$[D' \times 1]$ $[D' \times D]$ $[D \times 1]$



① Forward Selection

- $F' = \emptyset$
- At each iteration, find the best new feature to be added to F'

$$d^* = \arg \min_d \text{Error}(F' \cup d)$$

→ separate set of data points other than the training data.
- Add d^* to F' if $\text{Error}(F' \cup d) < \text{Error}(F')$

$F' = \emptyset$

$t=1 \Rightarrow$ ① 2 3 4 5 6 $\Rightarrow F' = \{1\}$
 $\{1,2\}$ $\{1,3\}$ $\{1,4\}$ $\{1,5\}$ $\{1,6\} \Rightarrow F' = \{1,4\}$

$t=2 \Rightarrow$ — $\{1,4,2\}$ $\{1,4,3\}$ — $\{1,4,5\}$ $\{1,4,6\} \Rightarrow F' = \{1,4,5\}$
 if $\text{Error} \{1,4\} < \text{Error} \{1\} \Rightarrow \text{YES}$
 if $\text{Error} \{1,4,5\} < \text{Error} \{1,4\} \Rightarrow \text{YES}$

$t=3 \Rightarrow$ — $\{1,4,5,2\}$ $\{1,4,5,3\}$ — — $\{1,4,5,6\} \Rightarrow \underline{\text{STOP}}$
 if $\text{Error} \{1,4,5,2\} < \text{Error} \{1,4,5\} \Rightarrow \underline{\text{NO}}$

of ML models that we trained = $6 + 5 + 4 + 3 = 18$ out of 62

② Backward Elimination

- $F' = F$
- At each iteration, find the best feature to be removed from F'

$$d^* = \arg \min_d \text{Error}(F' / d)$$

→ set difference.

- Remove d^* from F' if $\text{Error}(F' / d) < \text{Error}(F')$

$t=1 \Rightarrow$

$\{2, 3, 4, 5, 6\}$ <u>remove 1</u>	$\{1, 3, 4, 5, 6\}$ <u>remove 2</u>	$\{1, 2, 4, 5, 6\}$ <u>remove 3</u>	$\{1, 2, 3, 5, 6\}$ <u>remove 4</u>
	$\{1, 2, 3, 4, 6\}$ <u>remove 5</u>	$\{1, 2, 3, 4, 5\}$ <u>remove 6</u>	$\Rightarrow F' = \{1, 3, 4, 5, 6\}$

if $\text{Error}\{1, 3, 4, 5, 6\} < \text{Error}\{1, 2, 3, 4, 5, 6\} \Rightarrow \text{YES}$

$t=2 \Rightarrow$

$\{3, 4, 5, 6\}$	$\{1, 4, 5, 6\}$	$\{1, 3, 5, 6\}$	$\{1, 3, 4, 6\}$	$\{1, 3, 4, 5\}$
------------------	------------------	------------------	------------------	------------------

if $\text{Error}\{1, 4, 5, 6\} < \text{Error}\{1, 3, 4, 5, 6\} \Rightarrow \text{NO} \Rightarrow \underline{\text{STOP}}$

Return $F' = \{1, 3, 4, 5, 6\}$

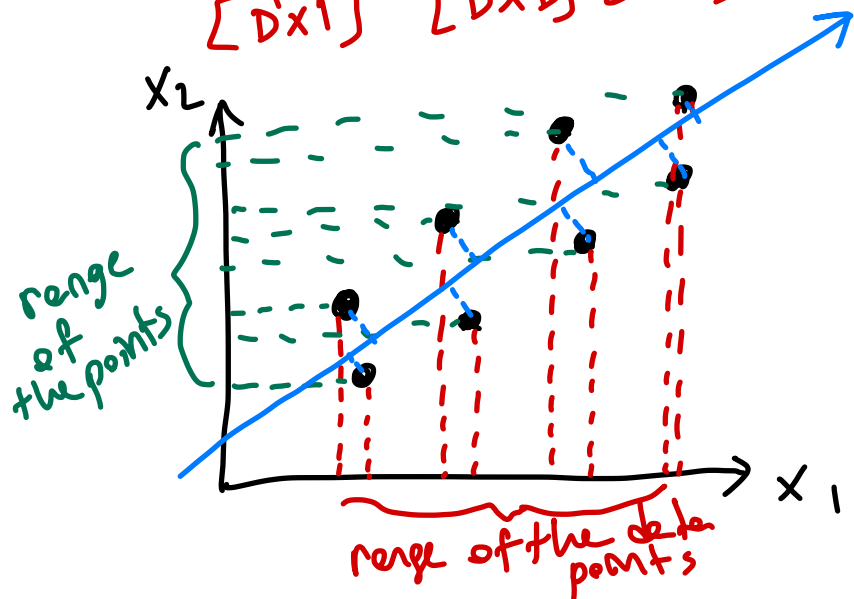
Principal Component Analysis (PCA)

- PCA is a feature extraction algorithm

$$x \in \mathbb{R}^D \quad z \in \mathbb{R}^{D'} \quad W \in \mathbb{R}^{D \times D'}$$

$$z = W^T \cdot x$$

$$[D' \times 1] \quad [D' \times D] \quad [D \times 1]$$



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

$$\text{VAR}(z) = \text{VAR}(W^T \cdot x) \quad D' = 1$$

$$= W^T \cdot \text{VAR}(x) \cdot W$$

$$= W^T \cdot \sum_{D \times D} x \cdot W_{D \times 1}$$

$$\text{VAR}(aX) = a^2 \text{VAR}(X)$$

$$X = \begin{bmatrix} x_{11} & x_{12} \\ x_{21} & x_{22} \\ \vdots & \vdots \\ x_{81} & x_{82} \end{bmatrix} \begin{bmatrix} 1 \\ 0 \end{bmatrix} = \begin{bmatrix} x_{11} \\ x_{21} \\ \vdots \\ x_{81} \end{bmatrix}$$

$$X = \begin{bmatrix} x_{11} & x_{12} \\ x_{21} & x_{22} \\ \vdots & \vdots \\ x_{81} & x_{82} \end{bmatrix} \begin{bmatrix} 0 \\ 1 \end{bmatrix} = \begin{bmatrix} x_{12} \\ x_{22} \\ \vdots \\ x_{82} \end{bmatrix}$$

$$\text{maximize } \text{VAR}(z) = w^T \cdot \Sigma_x \cdot w$$

assume w^* is the optimum solution

$$\tilde{w} = 2 \cdot w^*$$

$$\tilde{w}^T \cdot \Sigma_x \cdot \tilde{w} = (2 \cdot w^*)^T \cdot \Sigma_x (2 \cdot w^*)$$

$$> w^{*T} \Sigma_x w^*$$

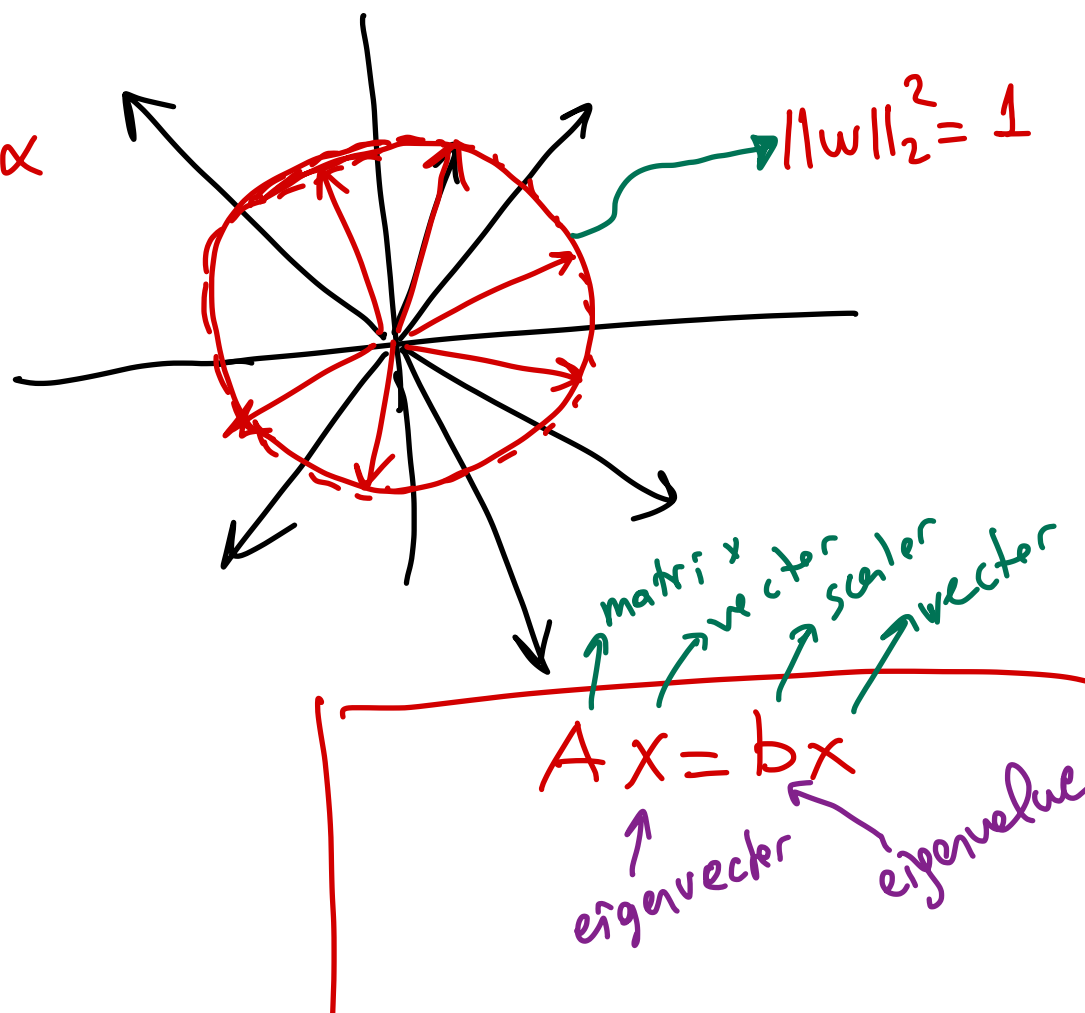
$$\begin{aligned} &\text{maximize } w^T \cdot \Sigma_x \cdot w \\ &\text{subject to: } \|w\|_2^2 = 1 \end{aligned} \quad \alpha$$

$$\begin{aligned} L_p &= w^T \cdot \Sigma_x \cdot w - \alpha \cdot (\|w\|_2^2 - 1) \\ &= w^T \cdot \Sigma_x \cdot w - \alpha \cdot (w^T \cdot w - 1) \end{aligned}$$

$$\frac{\partial L_p}{\partial w} = 2 \cdot \Sigma_x \cdot w - 2 \cdot \alpha \cdot w = 0$$

$$\boxed{\Sigma_x \cdot w = \alpha \cdot w}$$

$P=2 \Rightarrow 2 \times 2 \quad 2 \times 1 \quad 1 \times 1 \quad 2 \times 1$



D eigenvalues

$$\alpha_1, \alpha_2, \dots, \alpha_D \Rightarrow \alpha_1 \geq \alpha_2 \geq \dots \geq \alpha_D$$

w^* \Rightarrow the eigenvector that corresponds to the largest eigenvalue
[the first eigenvector]

Exercise: If $D' = 2 \Rightarrow$ we need to pick the first two eigenvectors.

$$W = \begin{bmatrix} | & | \\ w_1 & w_2 \\ | & | \end{bmatrix}$$

first eigenvector

second eigenvector

PCA Algorithm:

Step 1: Calculate Σ_x .

Step 2: Find first D' eigenvectors.

$$W = \begin{bmatrix} | & | & \dots & | \\ w_1 & w_2 & \dots & w_{D'} \\ | & | & \dots & | \end{bmatrix}^{D \times D'}$$

Projection Step: $z_i = W^T \cdot (x_i - \hat{\mu}) \quad \forall i$
centering

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

sample mean