

# # PCA steps

$$\underline{x} = [\underline{x}_1, \underline{x}_2, \dots, \underline{x}_N], \quad \underline{x}_i \in \mathbb{R}^D, D = 784$$

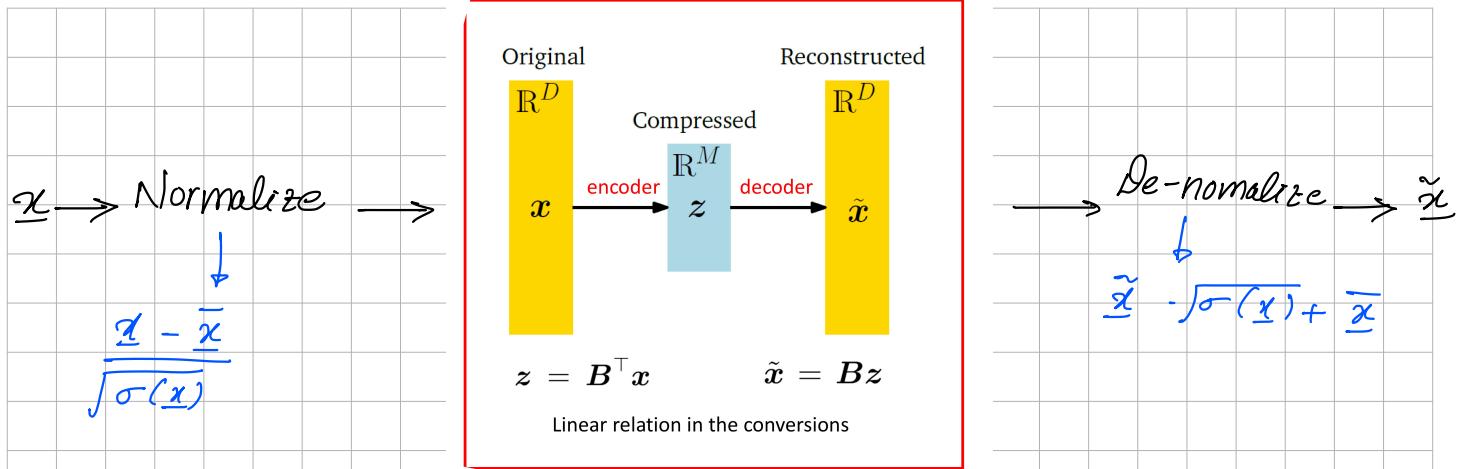
$$\underline{S} = \frac{1}{N} \sum_{n=1}^N \underline{x}_n \underline{x}_n^T$$

$\begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix}_{D \times 1} \cdot \begin{bmatrix} \vdots & \vdots & \vdots & \vdots \\ \cdots & \cdots & \cdots & \cdots \\ \vdots & \vdots & \vdots & \vdots \end{bmatrix}_{1 \times D} \Rightarrow \begin{bmatrix} \vdots & \vdots & \vdots & \vdots \\ \vdots & \ddots & \ddots & \vdots \\ \vdots & \vdots & \ddots & \vdots \\ \vdots & \vdots & \vdots & \ddots \end{bmatrix}_{D \times D}$

$$\begin{bmatrix} \begin{bmatrix} 1 \\ \underline{x}_1 \\ 1 \end{bmatrix} \begin{bmatrix} 1 \\ \underline{x}_2 \\ 1 \end{bmatrix} \cdots \begin{bmatrix} 1 \\ \underline{x}_N \\ 1 \end{bmatrix} \end{bmatrix}_{D \times N} = \begin{bmatrix} [\underline{x}_1] \\ [\underline{x}_2] \\ \vdots \\ [\underline{x}_N] \end{bmatrix}_{N \times D} \quad D \times D$$

$$\Rightarrow \begin{bmatrix} 1 \\ \underline{x}_1 \\ 1 \end{bmatrix}_{D \times 1} [ -\underline{x}_1 - ]_{1 \times D} + \begin{bmatrix} 1 \\ \underline{x}_2 \\ 1 \end{bmatrix}_{D \times 1} [ -\underline{x}_2 - ]_{1 \times D} + \cdots + \begin{bmatrix} 1 \\ \underline{x}_N \\ 1 \end{bmatrix}_{D \times 1} [ -\underline{x}_N - ]_{1 \times D}$$

$$\Rightarrow \begin{bmatrix} \vdots & \vdots & \vdots \\ \vdots & \ddots & \vdots \\ \vdots & \vdots & \ddots \end{bmatrix}_{D \times D} + \begin{bmatrix} \vdots & \vdots & \vdots \\ \vdots & \ddots & \vdots \\ \vdots & \vdots & \ddots \end{bmatrix}_{D \times D} + \cdots + \begin{bmatrix} \vdots & \vdots & \vdots \\ \vdots & \ddots & \vdots \\ \vdots & \vdots & \ddots \end{bmatrix}_{D \times D} \Rightarrow \sum_{n=1}^N \underline{x}_n \underline{x}_n^T$$



$\rightarrow [-\lambda_1 -], [-\underline{v}_1 -] = \text{eigendecomposition } (\underline{\underline{\Sigma}})$

$\rightarrow$  rearrange eigenvalues in descending order

$$[\lambda_1, \lambda_2, \dots, \lambda_D] \quad \lambda_1 > \lambda_2 > \dots > \lambda_D$$

$\rightarrow$  rearrange eigenvectors in the same order

$$[\underline{v}_1, \underline{v}_2, \dots, \underline{v}_D], \quad \underline{v}_D \in \mathbb{R}^D$$

$\rightarrow$  Collect first  $M$  eigenvectors. ( $M$  = no. of Principal Components  
= reduced dimension)  
 $\underline{\underline{B}} = [\underline{v}_1, \underline{v}_2, \dots, \underline{v}_M], \quad M \leq D$

$$\underline{x} \rightarrow \underline{\underline{B}}^\top \underline{x} \rightarrow \underline{z} \rightarrow \underline{\underline{B}} \underline{z} \rightarrow \tilde{\underline{x}}$$

encode    decode

$$\begin{bmatrix} \underline{v}_1 \\ \underline{v}_2 \\ \vdots \\ \underline{v}_M \end{bmatrix}_{M \times D} \begin{bmatrix} \underline{x}_1 \\ \underline{x}_2 \\ \vdots \\ \underline{x}_N \end{bmatrix}_{D \times N} = \begin{bmatrix} z_1 \\ z_2 \\ \vdots \\ z_M \end{bmatrix}_{M \times N}$$

$$\begin{bmatrix} \underline{v}_1 \\ \vdots \\ \underline{v}_M \end{bmatrix}_{D \times M} \begin{bmatrix} z_1 \\ z_2 \\ \vdots \\ z_N \end{bmatrix}_{M \times N} = \begin{bmatrix} \tilde{x}_1 \\ \tilde{x}_2 \\ \vdots \\ \tilde{x}_N \end{bmatrix}_{D \times N}$$

$$\underline{x} \rightarrow \underline{\underline{B}}^\top \underline{x} \rightarrow \underline{z} \rightarrow \underline{\underline{B}} \underline{z} \rightarrow \tilde{\underline{x}}$$

$$\begin{bmatrix} \underline{v}_1 \\ \vdots \\ \underline{v}_M \end{bmatrix}_{M \times D} \begin{bmatrix} \underline{x}_1 \\ \vdots \\ \underline{x}_N \end{bmatrix}_{D \times N} = \begin{bmatrix} z_1 \\ \vdots \\ z_M \end{bmatrix}_{M \times N}$$

$$\begin{bmatrix} \underline{v}_1 \\ \vdots \\ \underline{v}_M \end{bmatrix}_{D \times M} \begin{bmatrix} z_1 \\ \vdots \\ z_N \end{bmatrix}_{M \times N} = \begin{bmatrix} \tilde{x}_1 \\ \vdots \\ \tilde{x}_N \end{bmatrix}_{D \times N}$$

→ Recall linear regression

$$h_{\underline{w}, b}(\underline{x}) = \underline{w} \cdot \underline{x} + b, \quad \underline{w} \in \mathbb{R}^d, \quad b \in \mathbb{R}$$

→ Used for regression, i.e. to predict values.

How to use it for classification?

Sigmoid function

$$\text{sig}(x) = \frac{1}{1+e^{-x}} = \frac{e^x}{1+e^x}$$

$$\text{if } x \in \mathbb{R}, \quad \text{sig}(x) \in (0, 1)$$

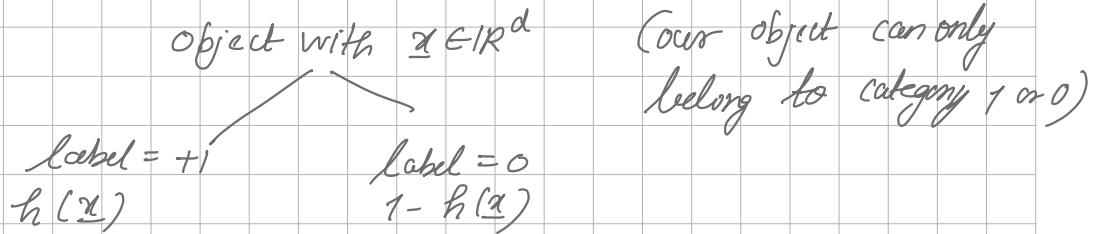
Thus, we can convert a predicted value from regression to a probability value.

→ Logistic Regression

$$h_{\underline{w}, b}(\underline{x}) = \text{sig}(\underline{w} \cdot \underline{x} + b), \quad \underline{w} \in \mathbb{R}^d, \quad b \in \mathbb{R}$$

$$h(\underline{x}) = \text{IP}(\text{Object with features } \underline{x} \in \mathbb{R}^d \text{ has label } +1)$$

$$\Rightarrow \text{IP}(\text{Object with features } \underline{x} \in \mathbb{R}^d \text{ has label } 0) = 1 - h(\underline{x})$$



→ Likelihood

$$\text{IP}(Y=y_i | X=\underline{x})$$

$$\text{IP}(Y=1 | X=\underline{x}) = h_{\underline{w}, b}(\underline{x})$$

$$\Rightarrow \text{IP}(Y=0 | X=\underline{x}) = 1 - h_{\underline{w}, b}(\underline{x})$$

$$= \text{IP}(Y=1 | X=\underline{x}) \cdot \text{IP}(Y=0 | X=\underline{x})$$

$$= h_{\underline{w}, b}(\underline{x}) \cdot (1 - h_{\underline{w}, b}(\underline{x}))$$

Example: Coin tossing  $\rightarrow y \in \{H, T\}$ ,  $IP(Y=H | X=\text{fair coin}) = \frac{1}{2}$   
 for two toss,  $IP(\text{getting } H \text{ atleast once})$

$$= IP(H, H) + IP(H, T) + IP(T, H)$$

$$= \frac{1}{2} \cdot \frac{1}{2} + \frac{1}{2} \cdot \frac{1}{2} + \frac{1}{2} \cdot \frac{1}{2} = \frac{3}{4}$$

### Maximum Likelihood Estimation

Given: sample with  $m$  no. of datapairs  $(x_i, y_i) \in \mathbb{R}^d \times \{1, 0\}$

Goal: Find the weight vector  $w_s$  and bias  $b_s$  that best explains the data

OR: Find the pair  $(w_s, b_s)$  which yields the longest likelihood

for observing all the data pairs  $(x_i, y_i)$

Assuming that our data pairs are i.i.d, we can say,

$$\mathcal{L}_{w,b}(s) = IP_{w,b}(y_1=y_1, \dots, y_m=y_m | X_1=x_1, \dots, X_m=x_m)$$

Probability of finding the correct labels  $y_1, \dots, y_m$  given all the samples  $x_1, \dots, x_m$

Since our samples are i.i.d, and  $P(A \cap B) = P(A) \cdot P(B)$

for independent A and B,

$$\mathcal{L}_{w,b}(s) = \prod_{i=1}^m IP_{w,b}(y=y_i | x=x_i)$$

$$\mathcal{L} = \prod_{i=1}^m h_{w,b}(x_i) (1 - h_{w,b}(x_i)) \stackrel{!}{\longrightarrow} \max$$

$$h_{\underline{w}, b}(\underline{x}) = \frac{1}{1 + \exp(-\underline{w} \cdot \underline{x} + b)} \quad \text{let } \underline{w} \cdot \underline{x} + b = t;$$

$$h(\underline{x}) = \frac{1}{1 + e^{-t}}, \quad 1 - h(\underline{x}) = 1 - \frac{1}{1 + e^{-t}} = \frac{e^{-t}}{1 + e^{-t}} = \frac{1}{1 + e^t}$$

$$\therefore 1 - h_{\underline{w}, b}(\underline{x}) = \frac{1}{1 + \exp(\underline{w} \cdot \underline{x} + b)}$$

### Log-likelihood

$$L = \prod_{i=1}^m h_{\underline{w}, b}(x_i) (1 - h_{\underline{w}, b}(x_i))$$

$$\ln(L) = \sum_{i=1}^m [y_i \ln(h_{\underline{w}, b}(x_i)) + (1-y_i)(1 - \ln(h_{\underline{w}, b}(x_i)))]$$

!  
≡ max

$$\text{Aim: } (\underline{w}, b) = \underset{\underline{w}, b}{\operatorname{argmax}}(\ln(L))$$

OR equivalently:

$$(\underline{w}, b) = \underset{\underline{w}, b}{\operatorname{argmin}}(-\ln(L)) \quad !!!$$

Many optimization  
algorithms to minimize

$$\underline{w}_s, b_s = \underset{\underline{w}, b}{\operatorname{argmin}} \left[ \sum_{i=1}^m -y_i \ln(h_{\underline{w}, b}(x_i)) - (1-y_i)(1 - \ln(h_{\underline{w}, b}(x_i))) \right]$$

$L(\underline{w}, b)$

Thus, we can find  $\underline{w}_s, b_s$  as follows:

$$\frac{\partial}{\partial \underline{w}} \mathcal{L}(\underline{w}, b) = \sum_{i=1}^m (h_{\underline{w}, b}(x_i) - y_i) \cdot 1$$

$$\frac{\partial}{\partial b} \mathcal{L}(\underline{w}, b) = \sum_{i=1}^m (h_{\underline{w}, b}(x_i) - y_i) \cdot 1$$

Update rule by gradient descent:

$$\underline{w}^{(k+1)} = \underline{w}^{(k)} - \alpha \frac{\partial}{\partial \underline{w}} \mathcal{L}(\underline{w}, b)$$

$$b^{(k+1)} = b^{(k)} - \alpha \frac{\partial}{\partial b} \mathcal{L}(\underline{w}, b)$$