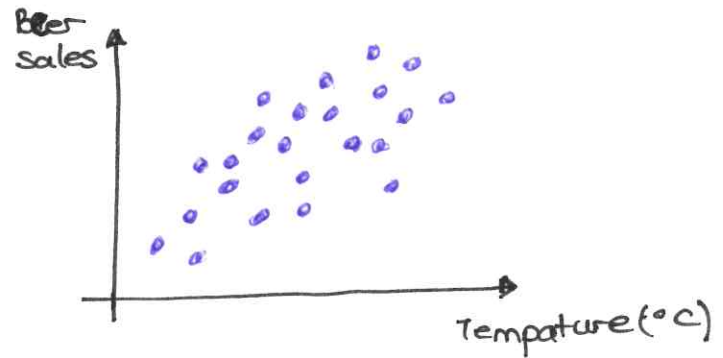
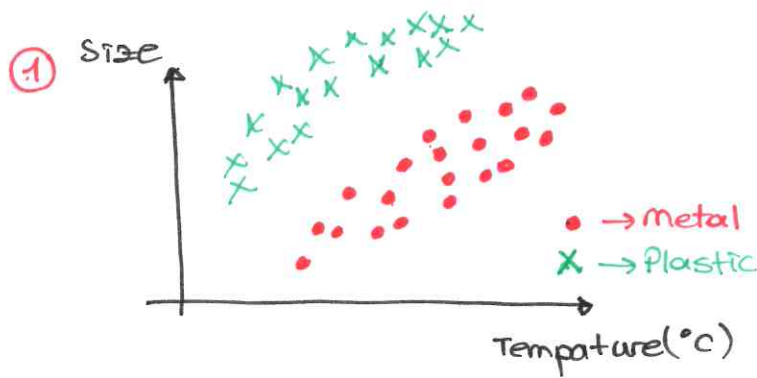


Exercise - 2 / Question - 1② Classification

⇒ We have training set $X = \{x_1, \dots, x_N\}$, i.e. we have N data points.

⇒ we know target values $T = \{t_1, \dots, t_N\}$ for each data point in training set X .

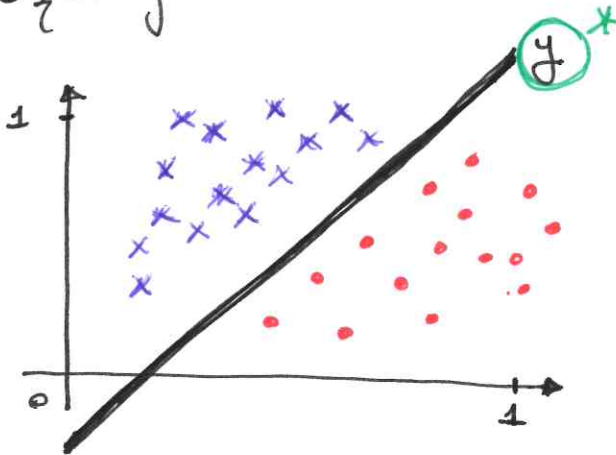
Goal: for a given new data x , classify these new data to one of K classes C_k .

2-class

K-class

$$t \in \{0, 1\}$$

$$t_n = (0, 1, 0, 0)^T$$



⇒ In this question, we have 2-class problem

⇒ Define a hyperplane y such that:

$$y(x) = w^T x + w_0$$

Goal: Finding the hyperplane that separates the data points into two classes (blue or red).

⇒ Let's consider a single data point:

$$\begin{aligned}
 x &= \begin{pmatrix} x_1 \\ \vdots \\ x_D \end{pmatrix} \quad w = \begin{pmatrix} w_1 \\ \vdots \\ w_D \end{pmatrix} \Rightarrow y(x) = w^T x + w_0 \\
 \text{"D dimensional vector"} \quad \text{"D dimensional vector"} & \quad \quad \quad = \sum_{d=1}^D w_d x_d + w_0
 \end{aligned}$$

$$(w_1 \dots w_D) \begin{pmatrix} x_1 \\ \vdots \\ x_D \end{pmatrix} + w_0$$

$$\underbrace{1 \times D \times D \times 1}_{1 \times 1}$$

$$\tilde{x} = \begin{pmatrix} x_0 = 1 \\ x_1 \\ \vdots \\ x_D \end{pmatrix} \quad \tilde{w} = \begin{pmatrix} w_0 \\ w_1 \\ \vdots \\ w_D \end{pmatrix}$$

$$\begin{aligned}
 y(x) &= w^T x + w_0 \\
 &= \sum_{d=1}^D w_d x_d + w_0 \\
 &= \sum_{d=0}^D w_d \tilde{x}_d \\
 &= \tilde{w}^T \tilde{x}
 \end{aligned}$$

! In that way, this equation becomes the multiplication of two vectors of size $(D+1) \times 1$.

⇒ Now, let's generalize it to K-class classification problem:

$$y_k(x) = w_k^T x + w_{k0} \quad k = 1, \dots, K$$

⇒ Turn this equation to vector notation:

$$① \quad y_k(x) = w_k^T x + w_{k0} \quad k = 1, \dots, K$$

$$② \quad y(x) = \tilde{w}^T \tilde{x} = [\tilde{w}_1, \tilde{w}_2, \dots, \tilde{w}_K]^T \tilde{x}$$

We know that,
 \tilde{w}_k is a vector of
 size $(D+1) \times 1$

$$③ \quad \tilde{W}^T = [\tilde{w}_1, \tilde{w}_2, \dots, \tilde{w}_K]^T$$

$$= \begin{pmatrix} \begin{matrix} w_{10} \\ w_{11} \\ \vdots \\ w_{1D} \end{matrix} & \begin{matrix} w_{20} \\ w_{21} \\ \vdots \\ w_{2D} \end{matrix} & \dots & \begin{matrix} w_{K0} \\ w_{K1} \\ \vdots \\ w_{KD} \end{matrix} \end{pmatrix}^T$$

$(D+1) \times 1 \quad (D+1) \times 1 \quad (D+1) \times 1$

when, we take transpose of this.

$$= \begin{pmatrix} w_{10} & w_{11} & \dots & w_{1D} \\ w_{20} & w_{21} & \dots & w_{2D} \\ \vdots & \vdots & \ddots & \vdots \\ w_{K0} & w_{K1} & \dots & w_{KD} \end{pmatrix}$$

$K \times (D+1)$

$$④ \quad y(x) = \tilde{W}^T \cdot \tilde{x}$$

$\downarrow \quad \downarrow$
 $K \times (D+1) \quad (D+1) \times 1$
 $= K \times 1$

As the output's
 a vector of $K \times 1$;
 we can directly compare
 it with target values.

⇒ Until now, we have talked about the equations for a single data point \tilde{x} .

⇒ Let's write them for the entire dataset with N data points:

(1)

$$Y(\tilde{X}) = \tilde{X} \tilde{W}$$

The entire dataset with N data points

A little explanation:

why it is written like $\tilde{X} \tilde{W}$ and why not $\tilde{W}^T \tilde{X}$

Actually, when we take transpose of $(\tilde{W}^T \tilde{X})$

$$(\tilde{W}^T \tilde{X})^T = \tilde{X}^T \tilde{W}$$

In the notation, we assume that \tilde{X}^T is \tilde{X} .

(2)

$$\tilde{W} = [\tilde{w}_1, \dots, \tilde{w}_K]$$

$(D+1) \times 1$ $(D+1) \times 1$

$$\Rightarrow (D+1) \times K$$

(3)

$$\tilde{X} = \begin{bmatrix} x_1^T \\ \vdots \\ x_N^T \end{bmatrix}$$

$$= \begin{bmatrix} x_{10} & x_{20} & \dots & x_{N0} \\ \vdots & \vdots & \ddots & \vdots \\ x_{1D} & x_{2D} & \dots & x_{ND} \end{bmatrix}^T$$

$$x_1 \rightarrow (D+1)$$

$$x_1^T \rightarrow 1 \times (D+1)$$

$$= \begin{bmatrix} x_{10} & \dots & x_{1D} \\ x_{20} & \dots & x_{2D} \\ \vdots & \ddots & \vdots \\ x_{N0} & \dots & x_{ND} \end{bmatrix}$$

$$\Rightarrow N \times (D+1)$$

④

Page 5

$$T = \begin{bmatrix} t_1^T \\ \vdots \\ t_N^T \end{bmatrix}$$

where \tilde{t} is the target value for a single data point \tilde{x} which is expressed as $\tilde{t} = [t_1, \dots, t_K]^T$

$(K \times 1)$

$$T = \begin{bmatrix} t_1^T \\ \vdots \\ t_N^T \end{bmatrix}$$

$$= \begin{bmatrix} t_{11} & \dots & t_{N1} \\ t_{12} & \dots & t_{N2} \\ \vdots & & \\ t_{1K} & \dots & t_{NK} \end{bmatrix}^T$$

for a single data point

$t_1 \rightarrow K \times 1$

$t_1^T \rightarrow 1 \times K$

$$= \begin{bmatrix} t_{11} & t_{12} & \dots & t_{1K} \\ t_{21} & t_{22} & \dots & t_{2K} \\ \vdots & \vdots & & \vdots \\ t_{N1} & t_{N2} & \dots & t_{NK} \end{bmatrix}$$

when we take transpose

$\Rightarrow N \times K$

⑤

$$Y(\tilde{X}) = \tilde{X} \cdot \tilde{W}$$

$$N \times (D+1) \quad (D+1) \times K$$

$$N \times K$$

$$Y(\tilde{X}) = \tilde{X} \cdot \tilde{W} \Rightarrow N \times K$$

⇒ As the output vector T has the size of $N \times K$ we can directly compare $\tilde{X}\tilde{W}$ with T :

$$y(\tilde{X}) - T = \underbrace{\tilde{X}\tilde{W}}_{(N \times K) - (N \times K)} - T$$

$\Rightarrow (N \times K)$ *

Goal: Choose \tilde{W} such that this is minimal!

⇒ We directly try to minimize the sum of squares error:

$$\begin{aligned} \textcircled{1} \quad E_D(\tilde{W}) &= \frac{1}{2} \sum_{k=1}^K \sum_{n=1}^N (y_k(x_n; w) - t_{kn})^2 \\ &= \frac{1}{2} \sum_{k=1}^K \sum_{n=1}^N (w_k^T x_n - t_{kn})^2 \\ &\quad \hookrightarrow \frac{1}{2} (\tilde{X}\tilde{W} - T)^2 \end{aligned}$$

$$\textcircled{2} \quad E_D(\tilde{W}) = \frac{1}{2} \text{Tr} \left((\tilde{X}\tilde{W} - T)^T (\tilde{X}\tilde{W} - T) \right)$$

Use:

$$\sum_{i,j} a_{ij}^2 = \text{Tr} \{ A^T A \}$$

where $a_{ij} \rightarrow w_k x_n - t_{kn}$
and A is $(\tilde{X}\tilde{W} - T)$

⇒ Take the derivative of $E_D(\tilde{w})$ for minimization

①

$$\frac{\partial E_D(\tilde{w})}{\partial \tilde{w}} = \frac{1}{2} \frac{\partial}{\partial \tilde{w}} \text{Tr} \left((\tilde{x}\tilde{w} - T)^T (\tilde{x}\tilde{w} - T) \right)$$

use chain rule

$$\frac{\partial \mathcal{L}}{\partial X} = \frac{\partial \mathcal{L}}{\partial Y} \frac{\partial Y}{\partial X}$$

where $\mathcal{L} \rightarrow \text{Tr} \left((\tilde{x}\tilde{w} - T)^T (\tilde{x}\tilde{w} - T) \right)$

$X \rightarrow \tilde{w}$

$$\frac{\partial \mathcal{L}}{\partial X} = \frac{\partial \text{Tr} \left((\tilde{x}\tilde{w} - T)^T (\tilde{x}\tilde{w} - T) \right)}{\partial \tilde{w}}$$

so, Y ??

Y will be $(\tilde{x}\tilde{w} - T)^T (\tilde{x}\tilde{w} - T)$ so
we get that:

②

$$\frac{\partial E_D(\tilde{w})}{\partial \tilde{w}} = \frac{1}{2} \frac{\partial \text{Tr} \left((\tilde{x}\tilde{w} - T)^T (\tilde{x}\tilde{w} - T) \right)}{\partial \left((\tilde{x}\tilde{w} - T)^T (\tilde{x}\tilde{w} - T) \right)} \cdot \frac{\partial \left((\tilde{x}\tilde{w} - T)^T (\tilde{x}\tilde{w} - T) \right)}{\partial \tilde{w}}$$

use:

$$\frac{\partial}{\partial A} \text{Tr} \{A\} = I$$

for us, $A \rightarrow (\tilde{x}\tilde{w} - T)^T (\tilde{x}\tilde{w} - T)$

③

$$\frac{\partial E_D(\tilde{w})}{\partial \tilde{w}} = \frac{1}{2} \cdot I \cdot \frac{\partial}{\partial \tilde{w}} \left((\tilde{x}\tilde{w} - T)^T (\tilde{x}\tilde{w} - T) \right)$$

④

$$\frac{\partial E_D(\tilde{w})}{\partial \tilde{w}} = \frac{1}{2} \cdot \cancel{2} \cdot \tilde{X}^T (\tilde{X} \tilde{w} - T)$$

$$= \tilde{X}^T (\tilde{X} \tilde{w} - T)$$

⇒ In order to minimize the sum-of-squares error, set the derivative $\frac{\partial E_D(\tilde{w})}{\partial \tilde{w}}$ to 0.

⑤

$$\frac{\partial E_D(\tilde{w})}{\partial \tilde{w}} = \tilde{X}^T (\tilde{X} \tilde{w} - T) \stackrel{!}{=} 0$$

Since \tilde{X}^T cannot be 0 as it is our data points, only $(\tilde{X} \tilde{w} - T)$ can be 0.

⑥

$$\tilde{X} \tilde{w} - T = 0$$

$$\tilde{X} \tilde{w} = T$$

$$\tilde{X}^T \cdot \tilde{X} \cdot \tilde{w} = \tilde{X}^T T$$

In order to leave \tilde{w} alone first take transpose of \tilde{X} on both side of the equation

⑦

$$\tilde{X}^T \tilde{X} \tilde{w} = \tilde{X}^T T$$

$$\underbrace{(\tilde{X}^T \tilde{X})^{-1}}_I (\tilde{X}^T \tilde{X}) \tilde{w} = \underbrace{(\tilde{X}^T \tilde{X})^{-1} \tilde{X}^T}_\tilde{X^+} \cdot T$$

$$I \tilde{w} = \tilde{X}^+ T$$

$$\tilde{w} = \tilde{X}^+ T$$

where \tilde{X}^+ is "pseudo-inverse"

Multiply both side of equation w/ $(\tilde{X}^T \tilde{X})^{-1}$