

EX.1

if $B = b_1$ and $A = a_1$ then [-] ①

if $B = b_1$ and $A = a_2$ and $C = c_1$ then [+] ②

if $B = b_1$ and $A = a_2$ and $C = c_2$ then [+] ③

if $B = b_1$ and $A = a_2$ and $C = c_3$ then [-] ④

if $B = b_1$ and $A = a_3$ then [+] ⑤

if $B = b_2$ and $C = c_1$ then [-] ⑥

if $B = b_2$ and $C = c_2$ then [-] ⑦

if $B = b_2$ and $C = c_3$ then [+] ⑧

1.2

An hypothesis h is consistent with a set of training examples D of target concept c if and only if $h(x) = c(x)$ for each training sample $(x, c(x))$ in D

$$\text{consistent}(x, c) \Leftrightarrow h(x) = c(x) \quad \forall x \in D$$

1.3

S1 is consistent because of ①

S2 is consistent because of ③

S3 is consistent because of ⑧

S4 is NOT consistent because of ⑦

EX.2

2.1

$h_{MAP} = \arg\max P(h|D) = \arg\max \frac{P(D|h)P(h)}{P(D)}$ if we assume $P(h_i) = P(h_j)$

we obtain $h_{ML} = \arg\max_{h \in H} P(D|h)$

2.2

Naive bayesian classifier uses conditional independence to approximate the solution more in detail

$$P(X_j|Y_1, Y_2) = P(X_j|Y_2)P(Y_1|Y_2) = P(X_j|Y_2)P(Y_1|Y_2)$$

assuming a target function $f: X \rightarrow V$ where each x is describable by attributes $\langle a_1, \dots, a_n \rangle$ we compute

$$\arg\max P(V_j|X, D) = \arg\max P(V_j | a_1, \dots, a_n, D)$$

with Bayesian assumption

$$(18) \quad P(a_1, \dots, a_n | V_j, D) = \prod_{i=1}^n P(a_i | v_i, D)$$

we obtain $V_{NB} = \max_{v_j \in V} P(v_j | D) \prod_i P(a_i | v_j, D)$

2.3

Naive bayes classifier is based on two assumption

① conditional independence

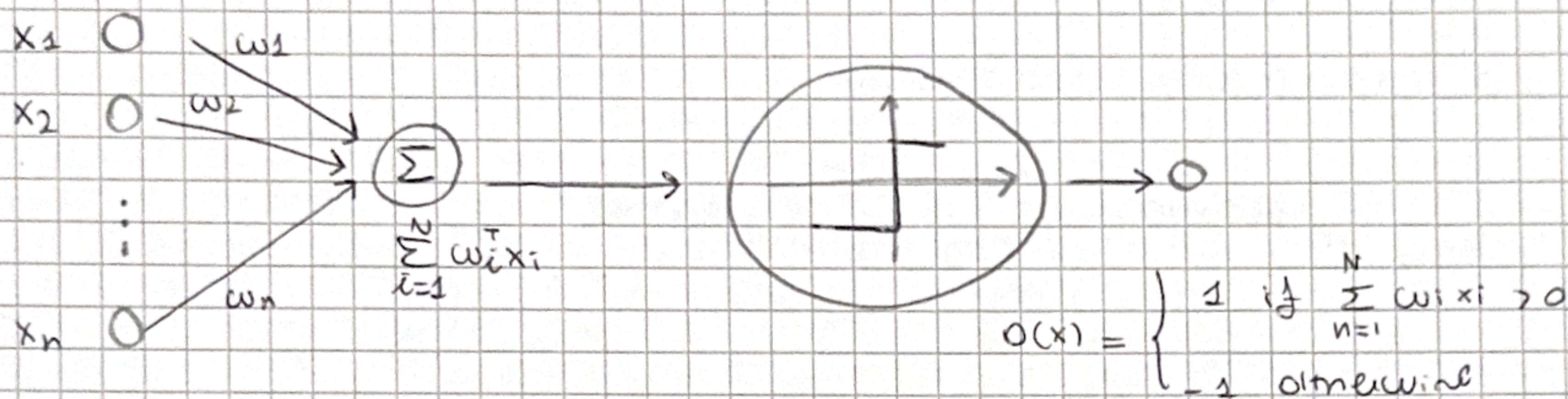
$$P(X, Y | Z) = P(X|Y, Z)P(Y|Z) = P(X|Z)P(Y|Z)$$

② Bayesian assumption

$$P(a_1, \dots, a_n | v_j, D) = \prod_i P(a_i | v_j, D)$$

EX.3

3.1



We want to minimize the squared error function

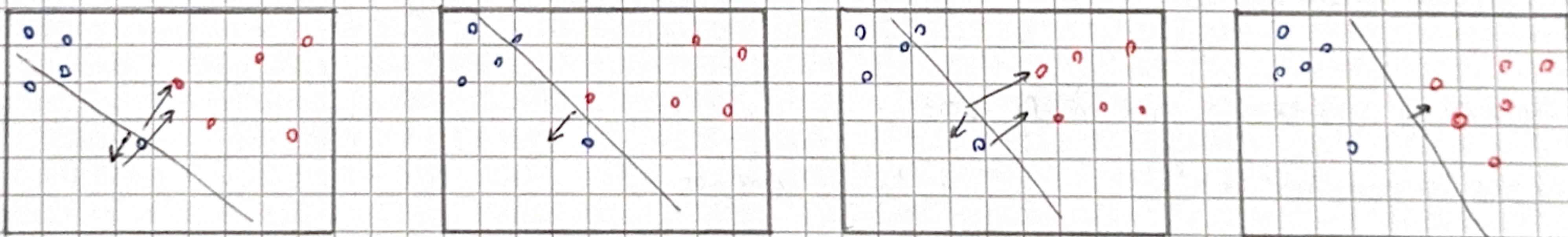
$$\text{to obtain } E(w) = \frac{1}{2} \sum_{n=1}^N (t_n - o_n)^2 = \frac{1}{2} \sum_{n=1}^N (t_n - w^T x_n)^2$$

$$\frac{\partial E}{\partial w} = \sum_{n=1}^N (t_n - w^T x_n) (-x_{i,n})$$

Algorithm will be with this update rule of weights

$$w_i \leftarrow w_i + \Delta w_i \quad \text{where } \Delta w_i = -n \cancel{\frac{\partial E}{\partial w}} = +n \sum_{n=1}^N (t_n - w^T x_n) (x_{i,n})$$

3.2



Ex 4.

4.1

Consider a linear model $y(x, w) = w^T x$ and a dataset $D = \{(x_n, t_n)\}_{n=1}^N$

$$E(w) = \frac{1}{2} \sum_{n=1}^N (w^T x_n - t_n)^2 + \lambda \|w\|^2$$

$$E(w) = (t - w^T x)^T (t - w^T x) + \lambda \|w\|^2$$

$$\text{where } X = \begin{bmatrix} x_1^T \\ \vdots \\ x_N^T \end{bmatrix} \text{ and } t = \begin{bmatrix} t_1 \\ \vdots \\ t_N \end{bmatrix}$$

optimal solution $\nabla E(w) = 0$

$$y(x, \omega^*) = \sum_{n=1}^N \alpha_n \omega_n^\top x \text{ where}$$

$$\omega_n = (x x^\top + \lambda I_N)^{-1} \text{ and } x x^\top = \text{gram matrix}$$

$$\begin{bmatrix} x_1^\top x_1 & \dots & x_1^\top x_M \\ \vdots & \ddots & \vdots \\ x_N^\top x_1 & \dots & x_N^\top x_M \end{bmatrix}$$

for kernelized version we use kernel trick where we substitute inner product $x^\top x$ with kernel function $k(x, x)$

so we will have $y(x, \omega) = \sum_{n=1}^N \alpha_n k(x, x_n)$ and gram matrix =

$$\begin{bmatrix} k(x_1, x_1) & \dots & k(x_1, x_M) \\ \vdots & \ddots & \vdots \\ k(x_N, x_1) & \dots & k(x_N, x_M) \end{bmatrix}$$

4.2

$f: \mathbb{R}^4 \rightarrow \mathbb{C}$ with $|f| = 3$ and D with 100 sample

$X \in \mathbb{R}^{100 \times 4}$

$X \in \mathbb{R}^{100 \times 4}$

$K \in \mathbb{R}^{100 \times 100}$

and $t \in \mathbb{R}^{100 \times 3}$

EX. 5

5.1

The role of backpropagation is to compute the gradient instead in not ~~as~~ a phase of learning but just compute the gradient in an efficient way different is for stochastic gradient descent that is an learning algorithm that allow to model to find the best values of weights ~~using~~ to solve the task

5.2

STOCHASTIC GRADIENT DESCENT

Require learning rate $\eta \geq 0$

Require initial values of $\theta^{(0)}$

$k \leftarrow 1$

while stopping criterion not meet do:

sample a subset $\{x^{(i)}, \dots, x^{(m)}\}$ from dataset

Compute the gradient estimate $g \leftarrow \frac{1}{m} \sum_i L(f(x^{(i)}, \theta^{(k)}), t^{(i)})$

update $\theta^{(k+1)} \leftarrow \theta^{(k)} - \eta g$

$k \leftarrow k + 1$

end while

$\eta \geq 0 \rightarrow$ always $\eta \ll 0$ can happen will never converge

$$P(C_3|X) = \frac{P(X|C_3)P(C_3)}{P(X)} = \frac{P(X|C_3)P(C_3)}{P(X|C_1)P(C_1) + P(X|C_2)P(C_2) + P(X|C_3)P(C_3)}$$

Considering a parametric model $P(X|C_i) = N(x; \mu_i; \Sigma)$

the number of parameters will be

Σ = size of covariance matrix $\frac{d(d+1)}{2}$ because is symmetric

μ = dimension of \mathbf{x} (number of dimensions)

now considering:

$$\ln \left(\frac{P(X|C_1)P(C_1)}{P(X|C_2)P(C_2) + P(X|C_3)P(C_3)} \right) = \ln \left(\frac{N(x; \mu_1; \Sigma)P(C_1)}{N(x; \mu_2; \Sigma)P(C_2) + N(x; \mu_3; \Sigma)P(C_3)} \right)$$

$$= w^T x + w_0$$

thus

$$P(C_1|X) = \sigma(w^T x + w_0)$$

now because $P(C_1) = \pi_1 \quad P(C_2) = \pi_2 \quad P(C_3) = 1 - (\pi_1 + \pi_2)$

we need to estimate $\pi, \mu_1, \mu_2, \Sigma$

EX-6

6.1

$$\text{#params} = ((100 \cdot 50) + 50) + ((50 \cdot 10) + 10) = 10060 \text{ trainable params}$$

\downarrow bias \downarrow bias

6.2

because we are dealing with ~~one~~ multiclass classification problem
outputs units \rightarrow softmax activation function

$$y_i = \text{softmax}(\alpha^{(i)}) = \frac{\exp(\alpha_i^{(i)})}{\sum_j \exp(\alpha_j^{(i)})}$$

6.3

Loss function categorical cross-entropy

$$J_i(\theta) = -\ln(\text{softmax}(\alpha^{(i)}))$$

$$\text{where } \alpha^{(i)} = w_i^T h + b_i$$

I choose this loss function because we are dealing with categorical multiclass classification problem