

2. ML general properties

ML EXAMS 23/03/2018

Exercise 1

I. Given a dataset D and a hypothesis space H : $h \in H$ describes training data if exist an hypothesis h' such that

$$\text{Error}(h) < \text{Error}(h') \quad \text{and} \quad \text{error}(h) > \text{error}(h')$$

s

s

D

D

II in decision tree is very probable to have overfitting in fact two different solutions for a problem and tree is deeper than others for one we have overfitting for reduce them.

1. stop growing when data split are not significant

2. grow full tree then post prune

for to do this we need to divide dataset in train and validation and use validation to understand if pruned

Question 2

1. Naive bayesian classifier uses conditional independence to approximate the algorithm

X is conditionally independent of V given Z

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

given a target function $f: X \rightarrow V$ where $V = \{v_1, v_2, \dots, v_n\}$

knowing that each sample can be explained as $\langle a_1, \dots, a_n \rangle$

$$P(V, | X, D) = P(V | a_1, \dots, a_m, D) = \frac{P(a_1, \dots, a_n | V, D) P(V | D)}{P(a_1, \dots, a_n | D)}$$

$$VNB = \underset{V \in V}{\operatorname{argmax}} P(a_1, \dots, a_n | V, D) P(V | D)$$

now with bayesian assumption $P(a_i, \dots, a_m | V, D) = \prod_{j=1}^m P(a_j | v_j, D)$

$$VNB = \underset{V \in V}{\operatorname{argmax}} P(V | D) \prod_i P(a_i | v_i | D)$$

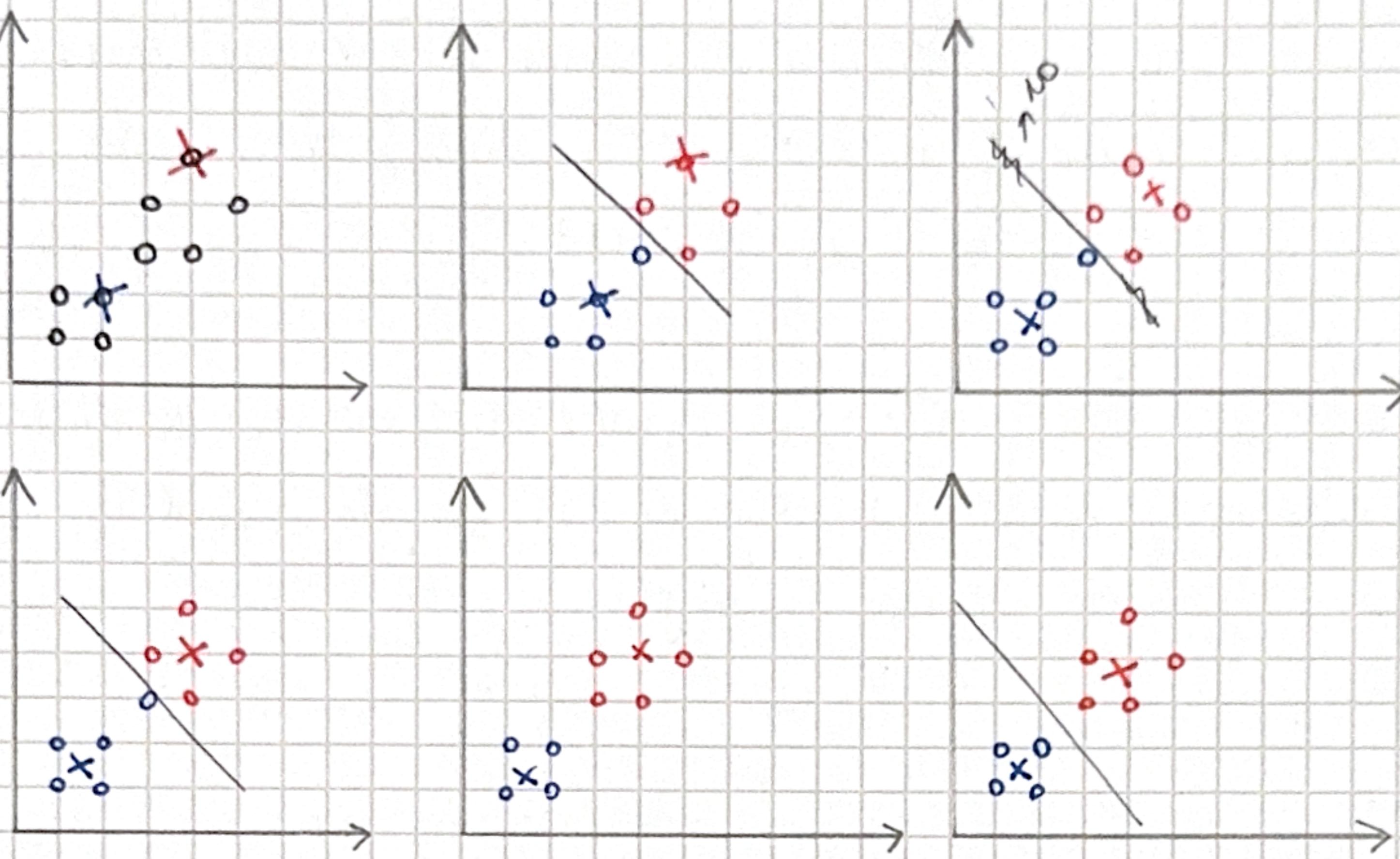
2. the target function is $f: X \rightarrow Y$ with $X = \{\text{title}, \text{x}, \text{abstact}, \text{x}, \text{author}, \text{x}, \text{title}\}$ and $Y = \{\text{ML}, \text{KR}, \text{RL}\}$ the dataset is $D = \{(d_i, y_i)\}_{i=1}^N$

$$VNB = \underset{V \in V}{\operatorname{argmax}} P(V | D) \prod_{i=1}^N P(d_i | V, D) \quad \hat{P}(f(d_i | V, D)) = \frac{\text{TF}_{ij} + 1}{\text{TF}_{ij} + 2}$$

Question 3

1. In a semi-supervised learning task we have a target function $f: X \rightarrow Y$ and a dataset $D = \{(x_n)\}_{n=1}^N$. It's possible to use we don't have labels the dataset is composed only by input values.

2.



Use K-MEANS to solve this problem

3. we can use an initially a point in classified as blue but some iterations is classified as red

Question 4.

1. classification problem $f: X \rightarrow C$ with dataset $D = \{(x_n, t_n)\}_{n=1}^N$

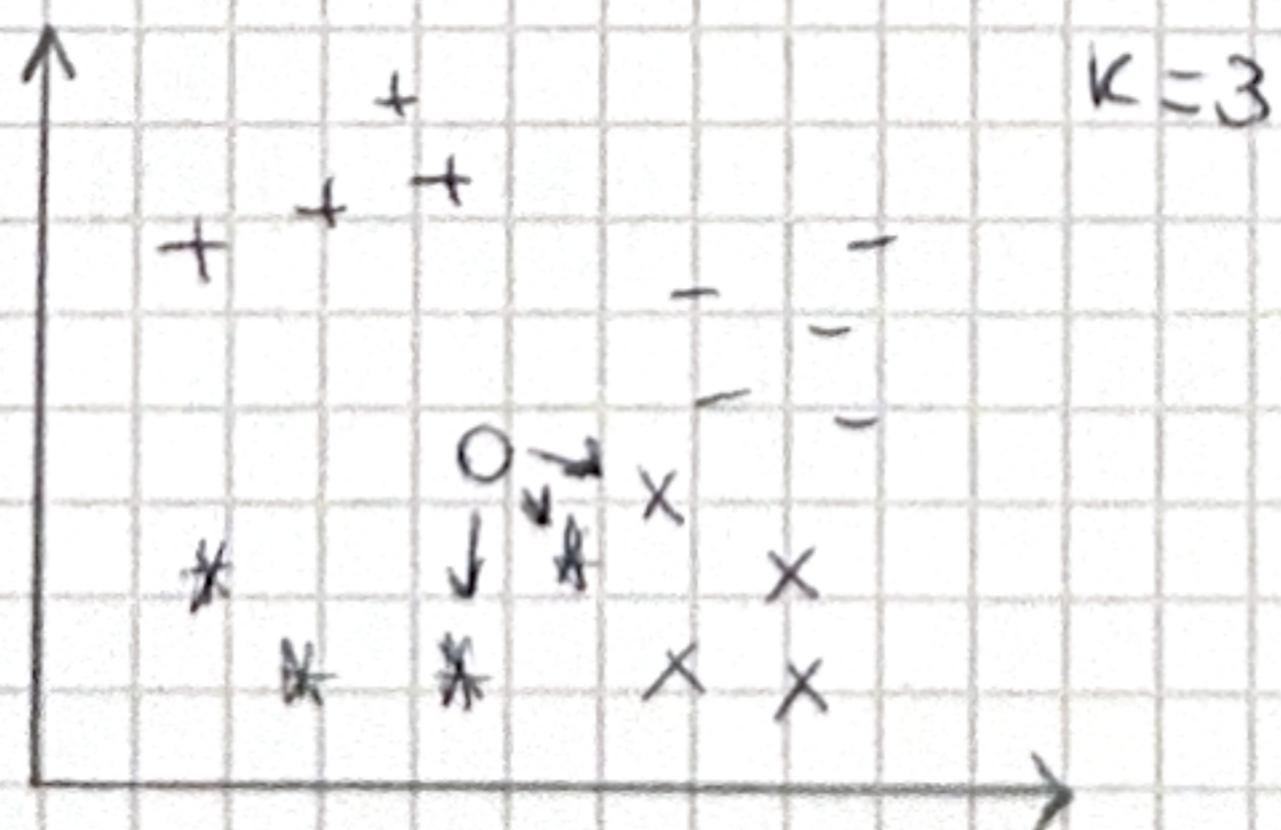
- ① find k nearest neighbours of new instance x
- ② assign to x the most common label among the majority of neighbours

likelihood of class c of new instance x is:

$$P(c|x, D) = \frac{1}{k} \sum_{x_n \in N_k(x, D)} \mathbb{I}(t_n=c)$$

where $\mathbb{I}(e) = \begin{cases} 1 & \text{if } e = \text{true} \\ 0 & \text{otherwise} \end{cases}$

2.



Question 5

1.

BACKPROPAGATION: is an algorithm used to compute the gradient. It is not a learning algorithm and is used in ANN. In this way we propagate the gradient through all the net.

FORWARD AND BACKWARD:

In the forward step we compute the value of the parameter

of the function (given x, t compute α, h and $J = L(t, y)$)

In the backward we compute the gradient (given x, t, α, h, y, J compute $\frac{\partial J}{\partial w_{ij}^{(k)}}$).

SGD: is a learning algorithm that is used to find the minimum or a maximum of a function is an iterative method and is a optimization strategy.

2.

Recall

FORWARD STEP

Require: depth l , weight matrices $W^{(i)}$

bias parameters $b^{(i)}$, x input

$h^{(0)} = x$ & t target value

for $k=1, \dots, l$ do

$$\alpha^{(k)} = b^{(k)} + W^{(k)} h^{(k-1)}$$

$$h^{(k)} = f(\alpha^{(k)})$$

end for

$$y = h^{(l)}$$

$$J = L(t, y)$$

BACKWARD STEP

$$g \leftarrow \nabla_y J = \nabla_y L(t, y)$$

for $k=l, \dots, 1$ do

$$g \leftarrow \nabla_{\alpha^{(k)}} J = g \odot f'(\alpha^{(k)})$$

$$\nabla_{b^{(k)}} J = g$$

$$\nabla_{w^{(k)}} J = g (h^{(k-1)})^T$$

$$g \leftarrow \nabla_{h^{(k-1)}} J = (W^{(k)})^T g$$

end for

Question 6

1. we have a target function $f: X \rightarrow Y$ with $x \in \mathbb{R}^N$ and $y \in \mathbb{R}$ our model

is $y(x, w) = \sum_{n=1}^N w_n \phi_n(x)$ that is linear in w and $D = \{(x_n, t_n)\}_{n=1}^N$

We know that $t = y(x, w) + \epsilon$ where ϵ is additive Gaussian noise

$P(\epsilon | D) = N(\epsilon | 0, \beta^{-1})$ assuming that we have

now we use $p(t | x, w, \beta) = N(t | y(x, w), \beta^{-1})$

now we use iid hypothesis the error function will be:

$$E_D(w) = \frac{1}{2} \sum_{n=1}^N (t_n - w^T \phi(x_n))^2 \text{ and } w^* = \arg \min (E_D(w))$$

2. we can use least squares $EJ(w) = \frac{1}{2} (t - \phi(w))^T (t - \phi(w))$ and
 $w^* = \phi^T t$

we can also use sequential algorithm $\hat{w} \leftarrow w + \eta [t_n - w^T \phi(w)] \phi(x_n)$

ML EXAMS 18/01/2018 A

Exercise 1

1.

if $C = c_1$ and $B = b_1$ then YES NO ①

if $C = c_1$ and $B = b_2$ then YES ②

if $C = c_2$ and $B = a_1$ then YES ③

if $C = c_2$ and $B = a_2$ and $B = b_1$ ④

if $C = c_2$ and $B = a_2$ and $B = b_1$ ⑤

if $C = c_2$ and $B = a_3$ then NO ⑥

if $C = c_3$ then NO ⑦

2.

in commitment with S_1 because of ①

in commitment with S_2 because of ④

NOT in commitment with S_3 because of ⑦

NOT in commitment with S_4 because of ⑤

Exercise 2

1. the maximum a posteriori hypothesis is $h_{MAP} = \arg\max_{h \in H} P(D|h) \frac{P(h)}{P(D)}$
if we assume $P(h_j) = P(h_i)$ we simplify and
we obtain $h_{ML} = \arg\max_{h \in H} P(D|h)$

2. Bayes is an optimal classifier that returns the most probable
class ~~class~~ prediction:

$$VNB = \arg\max_{v \in V} \sum_{h \in H} P(v|x,D) p(h|x,D) = \sum_{h \in H} P(v|x,h) P(h|D)$$

3. we can use it if we have analytical solution or if hypothesis space
is not too large otherwise it is not practical