

12/02/2018

①

[Ex-1]

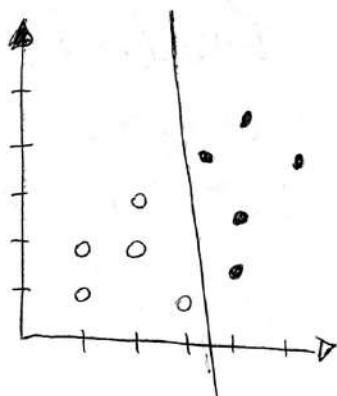
In Supervised learning we have dataset  $D = \{(x_n, t_n)\}_{n=1}^N$ ,  
 $x_n$  is input and  $t_n$  is target function expected to  
that particular input.

In supervised learn we have classification and regression.  
The main difference between them is that for learning  
function  $f: X \rightarrow Y$ , classification  $Y = \{C_1, \dots, C_K\}$  and  
in regression  $Y \subseteq \mathbb{R}$ .

In unsupervised learning we have dataset  $D = \{(x_n)\}$   
contains only information about input, value of  
target function is not available.

The aim is to define or find a structure of  
our input using some application such as clustering.

[Ex 2]



Perception method want to classify  
instances with this kind of weight to  
explore weight that we use for  
change our decision boundary  
~~and~~ in iteration when  
we find a point that is  
classified in a bad way

In particular we have func<sup>t</sup>  $\delta(x) = \begin{cases} 1 & \text{IF } x^T x > 0 \\ -1 & \text{otherwise} \end{cases}$

$$\theta = w^T x$$

Given a dataset  $D = \{(x_n, t_n)\}_{n=1}^N$

We want to minimize error function

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

so we need to compute derivative respect to  $w$

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

After this we need to update weight

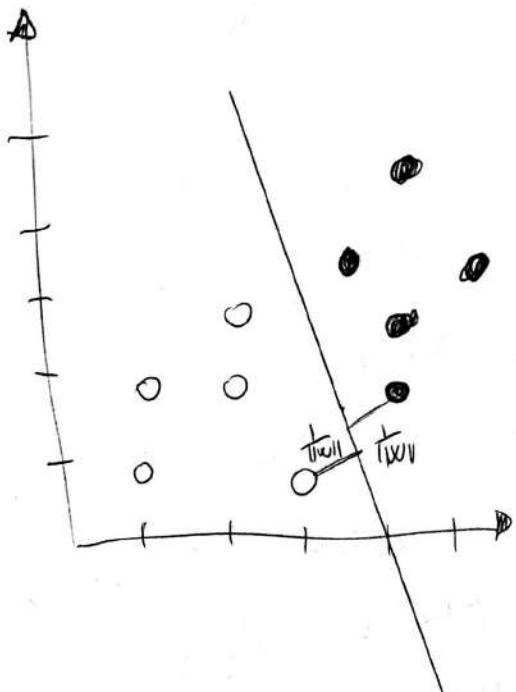
$$w \leftarrow w + \Delta w_i$$

$$\Delta w_i = -\eta \nabla E = \eta \sum_{n=1}^N (t_n - w^T x_n)$$

$\eta$  is learning rate

Perception start with  $w$  chosen in a random way  
and then ~~off~~ in each iteration update weights  
until this method converge, so we don't need to  
update weight anymore.

Method converge IF dataset is linearly separable  
and learning rate is small, otherwise IF dataset  
is not linearly separable perceptron will never  
converge



SVM we want to maximize the margin to provide a better accuracy.

Margin is smallest distance between closest points of decision boundary and the decision boundary itself. We compute  $w^*$  and  $w_0$  using lagrange multipliers. At the end this method converge if dataset is linearly separable. Otherwise we need to introduce slack variable to avoid problem of overfitting and try to converge else if there are some points that are misclassified.

~~We~~ I prefer SVM because ~~it~~ it's dataset is ~~not~~ almost linearly separable can use slack variable to maximize this margin while softly penalty points that are ~~misclassified~~ missclassified.

Ex 3

1)

IF  $C = c_1 \wedge B = b_1$  THEN No

IF  $C = c_1 \wedge B = b_2$  THEN Yes

IF  $C = c_2 \wedge A = a_1$  THEN Yes

IF  $C = c_2 \wedge A = a_2 \wedge B = b_1$  THEN Yes

IF  $C = c_2 \wedge A = a_2 \wedge B = b_2$  THEN No

IF  $C = c_3$  THEN No

2)

$$S = \{ S_1 = \langle a_1, b_1, c_1, \text{No} \rangle, S_2 = \langle a_2, b_2, c_2, \text{Yes} \rangle, S_3 = \langle a_1, b_2, c_3, \text{No} \rangle, S_4 = \langle a_3, b_1, c_1, \text{Yes} \rangle \}$$

For  $S_1$

for the first rule we can say that  $T$  is consistent because  $C = c_1 \wedge B = b_1$ , then no such as  $S_1$

For  $S_2$

IF  $A_1 = a_2 \wedge B = b_1 \wedge C = c_2 \Rightarrow \text{Yes}$ , is true  
with ~~four rules~~ four rule so also for this set  $T$  is consistent

For  $S_3$

IF  $C_3 \neq \text{No}$  No

Also for this set  $T$  is consistent

or S<sub>4</sub>

T is consistent for this set of samples because

$\text{IFC} = c_1 \wedge B = b_2$  we obtain YES such ~~as~~ as  
the second rule

EX 3.4

Boosting is an approach in we we consider to train different learners instead of just one and then collect their result, but in particular with Boosting we train different models in a sequential way this means that each classifier is trained on weighted data and that ~~the~~ weights depends on performance of the previous classifier.

If classifier give same error this means that we obtain weights with ~~very~~ higher values.

At the end predictions are based on weighted majority of votes.

3)

Adaboost can be explained also with sequential minimization of exponential error function

Consider error

$$\text{In } E = \sum_{n=1}^N \exp(-t_n f_n(x_n))$$

Otherwise

$$f_M(x_n) = \frac{1}{2} \sum_{m=1}^M \alpha_m Y_m(x_n) \quad t_n \in \{-1, +1\}$$

We want to minimize

Even function w.r.t respect to  $\alpha_m Y_m(x)$   
 $m = 1, \dots, M$

Instead of minimize over globally we ~~over~~  
assume

$$Y_1(x), \dots, Y_{M-1}(x)$$

$$\alpha_1, \dots, \alpha_M \quad \text{fixed}$$

~~fixed~~,

vary  $Y_M$  and  $\alpha_M$  explicit we obtain

$$E_{\text{fix}} = \sum_{n=1}^N \exp \left( -t_n f_{M-1}(x_n) - \frac{1}{2} \alpha_M \alpha_M Y_M(x_n) \right)$$

$$= \sum_{n=1}^N \exp \left[ -\frac{1}{2} \alpha_M \alpha_M Y_M(x_n) \right]$$

$$W_M^{(M)} = \exp \left( -t_M f_{M-1}(x_n) \right)$$

We know that from sequential optimality of  $E$

$$\alpha_M^{(M+1)} = \alpha_M^{(M)} \exp \left( \alpha_M I(Y_M(x_n) \neq t_n) \right)$$

$$\text{Sgn}(f_M(x)) = \text{Sgn} \left( \frac{1}{2} \sum_{m=1}^M \alpha_m Y_m(x) \right)$$

equivalent  $Y_M(x) = \text{sgn} \left( \frac{1}{2} \sum_{m=1}^M \alpha_m Y_m(x) \right)$

(5)

$$1) f^{(2)}(f^{(1)}(x, \theta^{(1)}), \theta_2^{(2)})$$

- 2) activation function ~~should~~ for input  
could be ~~value~~ because we need to  
predict a scalar  $x \in \mathbb{R}^3$  and for  
the activation function of the output we can have  
linear function because we don't work with  
binary classes or multiclass so we just  
need linear ~~as~~ function ~~and~~
- 3) Since we consider linear function as an  
activation function of the output we  
choose Mean squared error as a loss  
function for training the Network

$$J(\theta) = \frac{1}{2} \sum_{n=1}^N (t_n - f(x, \theta))^2$$

x6

1)

$$\mathcal{L} = (K + \lambda I)^{-1} t$$

Where  $K = X X^T = \begin{bmatrix} K(x_1, x_1) & \dots & K(x_1, x_n) \\ \vdots & \ddots & \vdots \\ K(x_n, x_1) & \dots & K(x_n, x_n) \end{bmatrix}$

this is Gram matrix

2) For Kernel trick we can say that if we find input vector  $x$  in form of inner product  $x^T x$  we can replace this product with kernel function

In case of linear regions we can say that

$$y(x) = \sum_{i=1}^n d_i K(x_i, x) \quad \text{with } \mathcal{L} = (K + \lambda I)^{-1} t$$

19/01/2018

Ok

(Ex) Input  $1242 \times 378 \times 3$

- First layer we have conv1, relu & pool1

After conv1 we obtain

$$W_{out1} = \frac{W_{in} - K + 2 \cdot P}{S} + 1 = \frac{1242 - 5 + 4}{1} + 1 = 1242$$

$$h_{out1} = \frac{W_{in} - K + 2 \cdot P}{S} + 1 = 378$$

the output after conv1 is  $1242 \times 378 \times 64$  features

- After ReLU dimension of outputs doesn't change  
and so we pass through pooling 1

The output after pooling should be  $W_{out} = \frac{W_{out1} - \text{Poolsize}}{\text{Stride}} + 1$

$$\Rightarrow W_{out} = 621 \quad h_{out} = 189$$

So After First layer we obtain

$$621 \times 189 \times 64$$

## Second Layer

We repeat some thing but in this case conv2 has a Kernel  $3 \times 3$

$$W_{out2} = 310 \quad h_{out2} = 94 \quad \Rightarrow 310 \times 94 \times 128$$

And after pool2 we have

$$78 \times 24 \times 128$$

2)

Nº parameters in ~~convolutional~~ conv layer would be  
 $(m \times n) + 1$  features

dimension of kernel

So For conv 1 we have 1664 parameters

ReLU and ~~pooling~~ don't change number of parameter

Nº of conv 2 are 1280

In total we obtain 2844 parameters for this CNN

EX 2

1) The dimensionality of old space is the dimension of images that in this case is  $15 \times 12 \times 12$  because we have this matrix in input and for the ~~assumed~~ intrinsic dimensionality we have 3 because it corresponds of 2D ~~rotation~~ translation and 1 rotation, so we have 3 degrees of freedom

2) We assume basis vectors are sorted by decreasing eigenvalues, e.g.: simple  $x_i$  can be transformed into the new space using projection but first we must ~~compute the mean~~  $\{x_1, \dots, x_6\}$  to compute mean of  $x_1, \dots, x_6$  so  $\bar{x} = \frac{1}{6} \sum_{n=1}^6 x_n$

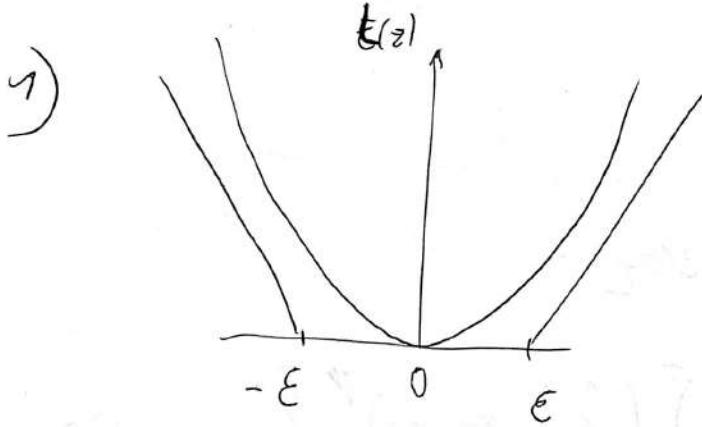
then  $[u_1 \bar{x}_1, \dots, u_6 \bar{x}_6]$ , we extract  $u_n$  and obtain

$$u_i^T S u_i = \lambda_i$$

When  $S$  is known and to have maximum variance we consider to take ~~high~~  $u_i$ , eigen vector corresponding to the largest eigen value  $\lambda_i$ . This is called principal components.

3) Marginally different of 3 because principal components of PCA are not latent variables

Ex 3



The problem of minimizing  $E$  is that it is not differentiable & it is difficult to solve.

2) this slack variable  $\epsilon^+$  and  $\epsilon^-$  have reduced overcon of problem of ~~non-differentiable~~ in particular with them we can compute a quadratic program that can be easily to solve and are short time and more the problem of ~~non-differentiable~~

Ex4

Least Squares is a linear classification method that has the aim to move prediction as close as possible to the target function

Given a dataset  $D = \{(x_n, t_n)\}$

$$y = \tilde{w}^T \tilde{x}$$

We want to compute  $\tilde{w}$

We consider

$$\tilde{x} = \begin{pmatrix} \tilde{x}_1 \\ \vdots \\ \tilde{x}_N \end{pmatrix} \quad + = \begin{pmatrix} t_1 \\ \vdots \\ t_n \end{pmatrix}$$

$$\tilde{w} = (\tilde{w}_1, \tilde{w}_K, \tilde{w}_K)$$

We want to minimize error

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

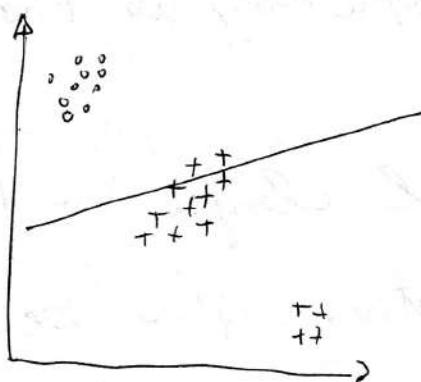
In a closed form we obtain

$$\tilde{w} = \cancel{\tilde{x}^T \tilde{x}} (\tilde{x} \tilde{x}^T)^{-1} \tilde{x}^T \cdot \cancel{+} = \underbrace{(\tilde{x}^T)^{-1} \tilde{x}^T}_{\text{pseudo inverse}} \cdot \cancel{+}$$

$$y = \tilde{w}^T \tilde{x} = \cancel{\tilde{w}^T} \tilde{x} = \cancel{\tilde{w}^T} (\tilde{x}^T)^{-1} \tilde{x}^T \cdot \cancel{+} \cdot \tilde{x}$$

This method have problem with outliers because the metric  $E(w)$  that is a sum of each

int ~~on~~ contribute to compute this error and  
 In particular with outliers we can have a decision  
 boundary that cannot divide perfectly ~~the~~  
~~the~~ instances that have different classification  
 value but we could have distance that ~~are~~ are  
~~not~~ not in the right side



The decision boundary  
 is NOT good because we  
 have outliers that  
 create problem so this  
 method is not robust to  
 outliers.

EX5

1) Maximum a posteriori hypothesis we find the  
 most probable hypothesis

$$h_{MAP} = \underset{h \in H}{\operatorname{argmax}} \frac{P(D|h)P(h)}{P(D)} = \underset{h \in H}{\operatorname{argmax}} P(D|h)P(h)$$

If  $p(h_i) = p(h_j)$

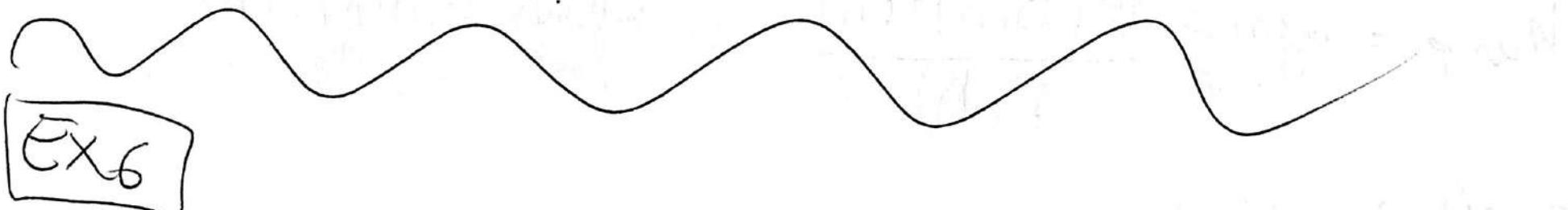
$$h_{MC} = \underset{h \in H}{\operatorname{argmax}} P(D|h)$$

2) Bayes Optimal classifier  $f: X \rightarrow V = \{v_1, \dots, v_k\}$

$$V_{BO} = \underset{v \in V}{\operatorname{argmax}} P(v_j | x, D) = \underset{v_j \in V}{\operatorname{argmax}} \sum_{h_i \in H} P(v_j | x, h_i) P(h_i | D)$$

Bayes Optimal classifier is used to compute most probable classification for new instance  $x$ , see what ~~we~~ we use prior probability of  $h_i$  to compute this most probable classification, we find the optimal value

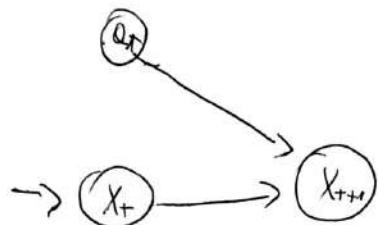
3) the problem of Bayes Optimal classifier is that it's not protocol for huge hypothesis space because we cannot know all prior probabilities of all hypothesis, we solve this prob by using Naive Bayes classifier that use conditional independence to ~~compute~~ estimate most probable classifier, see it more approximations.



Ex 6

Marks property said that current state contain all information to know and it does not depend of ~~precedent~~ precedent history.

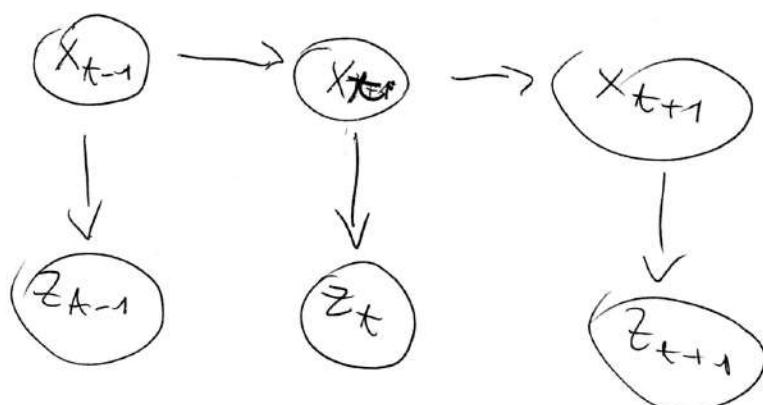
DP composed by  $\langle X, A, S, V \rangle$  where  $X$  is set of state  
 $A$  is set of actions,  $S$  is transition function,  $V$  is reward  
 Function. States are fully observable so we don't need  
 observation.



Next state depend only of current state, ~~so they~~ they  
 are input of next ts and past observe given the current state

~~Then~~ HMM is composed by  $\langle X, Z, \Pi \rangle$  where  $X$  is set  
 of states,  $Z$  is set of observation and  $\Pi^0$  is initial  
 distribution.

In this case state are non-observable so  
 we need observation but also in this case we have  
~~not~~ always markov property



23/03/2018

Ex 1

1) Given an hypothesis space  $H$ , we can say that  $h \in H$  overfits training that if chosen on  $h' \in H$  we have that

$$\text{error}_S(h) < \text{error}_S(h')$$

so

$$\text{error}_D(h) > \text{error}_D(h')$$

2)

In The decision tree we can have problems of overfitting because when we have a dataset that contains noisy data, or when training dataset is too small to produce representative sample of target function but we can choose to have two decision Tree  $T$  and  $T'$  obtained by different configuration of  $10^3$  with a dataset that contains noisy data and we notice that

$$\text{accuracy}(T') > \text{accuracy}(T)$$

We cannot say that  $T'$  is better than  $T$   
because certainly it has a problem with overfitting

For Decision tree To avoid overfitting we can stop  
To grow the tree, but it's is very difficult to  
estimate so we can use method of <sup>post</sup>Pruning  
Post pruning cannot to prune the tree after  
growing full tree.

There are two different method: Reduced-Error Pruning  
and Rule Post-Pruning

Reduced error pruning we split our dataset in  
Training and Validation set

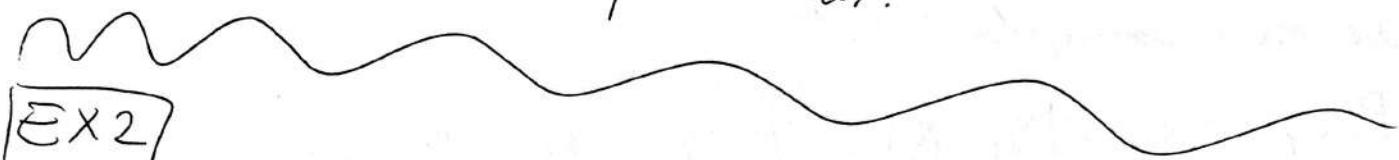
then until pruning is harmful we evaluate the  
Impact of on Validation set of pruning each possible  
Node so then we gradually remove the node that  
~~most~~ ~~the~~ improves ~~the~~ the validation ~~set~~ accuracy

This Pruning can give us ~~a~~ only optimal subset  
of tree and for a small training set it can  
give us bad result

Rule Post pruning is a pruning method that infer  
the tree allowing for overfitting

Then we convert decision ~~the~~ tree in set of rules

This is because in this way we can distinguish different contexts, & in fact each path of tree has different rule set. In this way we can prune paths independently. At the end, after pruning each path independently of each other, we set final rules.



**EX 2**

$$\text{f: } X \rightarrow V = \{v_1, \dots, v_k\} \quad X = \{x_1, \dots, x_m\} \notin D, \text{ dataset}$$

$$V_B = \arg \max_{V_j \in V} P(V_j | D) \prod_i P(x_i | V_j, D)$$

Naive Bayes classifier ~~is~~ is used for Hypothesis space that is very big so Bayes optimal classifier cannot know all prior probability of hypothesis  $h_i$ .

For this reason we use Naive Bayes classifier that use independence condition to do an approximation for new instance  $x$ .

$x$  is described ~~of~~ of ~~at~~ values of attribute

$x = \{x_1, \dots, x_n\}$  and so there are no more hypotheses

What we want to find is ~~the~~ the most probable classification and we do it using ~~very~~ ~~approximate~~ ~~approximation~~ approximation.

In particular

$$\begin{aligned} P(v_j | x, D) &= P(v_j | \alpha_1, \dots, \alpha_n, D) \\ &= P(v_j | D) (P(\alpha_1, \dots, \alpha_n | v_j, D)) \end{aligned}$$

$$\text{or max } P(v_j, x, D)$$

We make assumption that

$$P(\alpha_1, \dots, \alpha_n | v_j, D) = \prod_i P(\alpha_i | v_j, D)$$

$$\Rightarrow v_{NB} = \operatorname{argmax}_{v_j \in V} P(v_j | D) \prod_i P(\alpha_i | v_j, D)$$

②

3 stages of pages

We can use Naive Bayes with Multinomial distribution in which we consider that we cannot observe of  $w_i$  in documents

$$\text{f. } C = \{c_1, \dots, c_n\} \quad \text{Document } d_i \quad \text{Delete } \langle d_i, c_i \rangle$$
$$\underbrace{P(d_i | c_j, D)}_{\text{length}(d_i)} = \prod_{i=1}^{\text{length}(d_i)} P(\alpha_i = w_k | c_j, D)$$

We move thus conditional independence assumption

$$P(\alpha_i = w_k | c_j, D) = P(\alpha_m = w_k | v_j, D) \rightarrow k, m$$
$$\Rightarrow P(w_k | c_j, D)$$

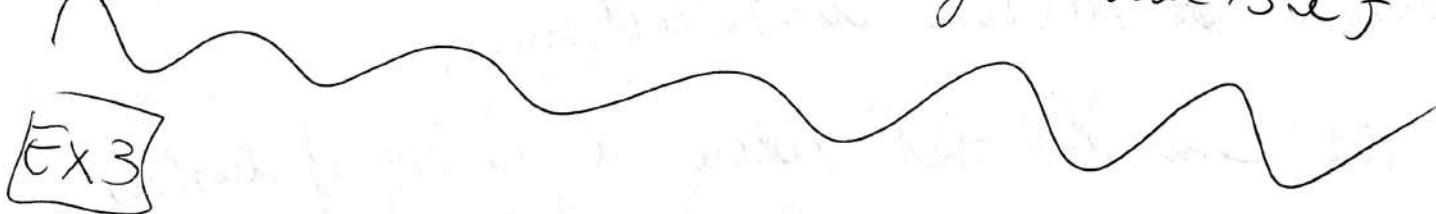
or multinomial distribution

$$P(d_i | c_j, D) = \text{Multinomial}(d_i, n, \theta)$$

$$P(w_i | c_j, D) = \frac{TF_{ij} + 2}{TF_j + |V|} = \text{approximate number}$$

Where  $TF_{ij}$  is total number of times that word  $w_i$  occurs in subset of  $V$  documents of  $D$  for which target value is  $c_j$

$TF_j$  is total number of word in subset of documents of  $D$  for which target value is  $c_j$



1) In Unsupervised learning problem we have dataset  $D = \{x_n\}$  where we have only information about inputs and the target value is not available. We want to find the structure of our inputs using some applications, methods such as clustering.

2) For example we have followers ~~as~~ on twitter  
~~and we collect them~~ but we don't know about  
outputs so we can analyse salary, or no. of  
comment we don't know a prior we find with  
some method such as clustering that  
through data find natural cluster if they exist,  
Data compression such as dimensionality reduction  
with PCA



3) We have  $\mathcal{X}$  followers of twitter,  $\{x_n\}_{n=1}^N$ , value  $K$   
can consider K-means so we want to compute  
 $K$  means of data  $\mu_1, \dots, \mu_k$  generate from  $K$   
Gaussian mixture distribution.

- ① We consider that value  $K$  is no. of clusters
- ② We ~~select~~ ~~set~~ best any initial partition that  
classifies the data into  $K$  clusters:  
We take first  $K$  ~~data~~ training samples as  
a single-element clusters. Then we assign  
each of remaining training samples  $\mathcal{N}-K$   
to a cluster with nearest centroid and we  
update this centroid after each assignment.
- ③ Take each sample in sequence and we check if

is in the cluster with nearest centroid. If sample is not in the cluster with nearest centroid, then this sample change ~~and~~ its cluster with cluster with closest centroid. After this we recompute centroids of two clusters involved.

④ We repeat ③ until there is a convergence.

There is a convergence where there is no more sample to switch.

We obtain a circular shape for our input, if it converges and it converge if there is a finite number of partition of ~~processes~~ into ~~processes~~  $K$  cluster and if distance from ~~any~~ sample onto ~~a~~ group of samples of centroid decreased.

EX 4

1) KNN is classification algorithm for non parametric model

We have Dataset  $D$ , target function  $f: X \rightarrow C$  and  $K$ . The algorithm is composed by two steps for instance  $x \notin D$

- ① Choose  $K$  nearest neighbors of instance  $x$
- ② Assign to  $x$  the most common value among  $K$  nearest neighbors

Likelihood Function

$$P(c|x, D, K) = \frac{1}{K} \sum_{i \in N_K(x, D)} I(x_i=c)$$

When  $N_k(x, \Delta)$  are Nearest Neighbors and  $I(e) = \{ \cdot \}$

2)

$K=3$

4 classes

$O_1$  has 3

2 + and 1 \*

so it is

classified as '+'

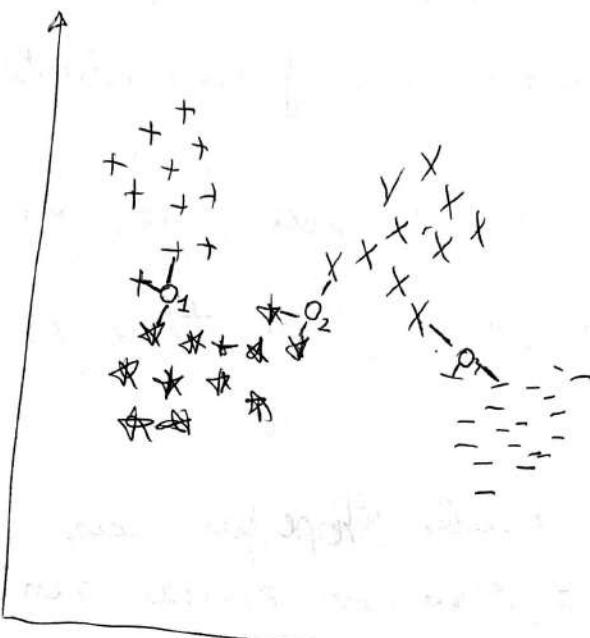
$O_2$  has 2 \*

and 1 x so it is classified as '\*'

$O_3$  has 2 -

and 1 x so it is classified as '-'

EX 5



1) Backpropagation is algorithm used only to compute gradient ~~of cost function~~ of cost function  $J(\theta)$ . In particular the approach is bottom-up, we start from the output layer and we propagate at each level the gradient.

Stochastic Gradient Descent is a learning algorithm that we use to improve the learning using the

gradient, in each call with it we update weights at each step.

Forward and backward pass : Forward pass calculate value from inputs to outputs and backward pass we perform the backpropagation that start from output and go to the inputs applies chain rule to compute the gradient

2)

### Forward pass

Require :  $l$  depth of network

Require :  $x$  inputs

Require :  $t$  target value

Require :  $W^i$  matrices of weights  $i = 1, \dots, l$

Require :  $b^{(i)}$  bias  $i = 1, \dots, l$

$$h^0 = x$$

for  $k = 1, \dots, l$

$$d^k = W^k \cdot h^{k-1} + b^k$$

$$h^k = f(d^k)$$

end for

$$y = \mathbf{f} h^{(l)}$$

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

Backward pass

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

for  $K = l \dots 1$

Propagate gradients to previous layer activations

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

$$\cancel{\nabla}_{\theta^{(K)}} J = g$$

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

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

end for

EX6

1) Linear regression we have learning function  $f: X \rightarrow Y$

where  $Y \subseteq \mathbb{R}$ , we want to find relationship

between inputs and output

In particular in linear regres we have

$$Y(X; w) = w_0 + w_1 x_1 + \dots + w_d x_d = w^T X$$

$$w = \begin{pmatrix} w_0 \\ \vdots \\ w_d \end{pmatrix} \quad X = \begin{pmatrix} 1 \\ \vdots \\ x_d \end{pmatrix}$$

w = weights  
x = inputs

We can consider linear regression with non function of input variable

$$y(x; w) = w^T \phi(x)$$

$$w = \begin{pmatrix} w_0 \\ w_1 \\ \vdots \\ w_{M-1} \end{pmatrix}$$

$$\phi(x) = \begin{pmatrix} \phi_0(x) \\ \vdots \\ \phi_{M-1}(x) \end{pmatrix}$$

We want to minimize loss function

$$\underset{w}{\text{min}} \sum_{n=1}^N (t_n - w^T \phi(x_n))^2$$

this means

$$E_\delta(w) = \frac{1}{2} (T - w\phi)^T (T - w\phi)$$

$$t = \begin{pmatrix} t_1 \\ \vdots \\ t_N \end{pmatrix}$$

$$\phi = \begin{pmatrix} \phi_0(x_1) & \cdots & \phi_{M-1}(x_1) \\ \vdots & & \vdots \\ \phi_0(x_N) & \cdots & \phi_{M-1}(x_N) \end{pmatrix}$$

$$W_{ML} = (\phi^T \phi)^{-1}$$

$$\phi^T T = \phi^+ t$$

For sequential learning we use it when the dataset is very huge so we use stochastic

$$w \leftarrow w + \eta \nabla E$$

$$\hat{W} \leftarrow W + \eta \sum_{n=1}^N (t_n - W^\top \phi(x_n)) \cdot \phi(x_n)$$

that converge for suitable value of learning rate  $\eta$ .

13/12/2012

Ex 1

1)  $F = \text{Furniture}$

$NR = \text{Nr Rooms}$

$NK = \text{New Kitchen}$

$f: F \times NR \times NK \rightarrow \{\text{Yes, No}\}$

with  $F = \{\text{Yes, No}\}$ ,  $NR = \{3, 4\}$   
 $NK = \{\text{Yes, No}\}$

2)

Variable are chosen using highest value of information gain where information gain is

$$\text{Gain}(S, A) = \text{Entropy}(S) - \sum_{\text{values}(A)} \frac{|S_v|}{|S|} \text{Entropy}(S_v)$$

that is expected Reducer of entropy

$$\text{Entropy}(S) = -p_+ \log_2 p_+ - p_0 \log_2 p_0$$

$p_+$  portion of positive samples

$p_0$  portion of negative samples

$$③ S_A = [3+, 2-] \quad \begin{matrix} \text{means 3 positive} \\ 2 negative \end{matrix} \quad \text{for acceptable}$$

$$S_F = \{ \text{Yes}, \text{No} \} \quad \text{Entropy}(S_A) = 0.969$$

$$S_{Fy} = [1+, 1-] \quad \text{Entropy}(S_{Fy}) = 1$$

$$S_{FN} = [2+, 1-] \quad \text{Entropy}(S_{FN}) = 0.918$$

$$\text{Gain}(S_F, F) = 0.969 - \frac{2}{5} \cdot 1 - \frac{3}{5} \cdot 0.918 = 0.0182$$

$$S_{NR3} = [1+, 2-] \quad \text{Entropy}(S_{NR3}) = 0.918$$

$$S_{NR4} = [2+; 0] \quad \text{Entropy}(S_{NR4}) = 0$$

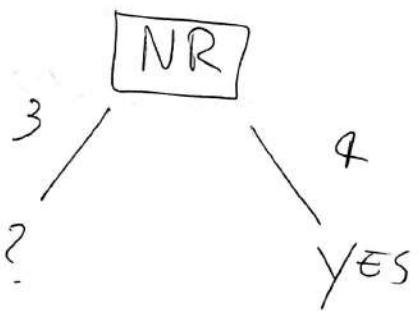
$$\text{Gain}(S, NR) = 0.969 - \frac{3}{5} \cdot 0.918 = 0.4182$$

$$S_{Nky} = [1; 0] \quad \text{Entropy}(S_{Nky}) = 0$$

$$S_{NKN} = [2+; 2-] \quad \text{Entropy}(S_{NKN}) = 1$$

$$\text{Gain}(S, NK) = 0.969 - \frac{4}{5} = 0.169$$

The highest gain we have with Attribute NR  
So the root is NR



Consider form of other in that condition so when

$$NR = 3 \quad S_3 = [1+ 2-] \quad \text{Entropy}(S_3) = 0.918$$

~~$$SFY = [0; 1-] \quad \text{Entropy}(SFY) = 0$$~~

~~$$SFN = [1+, 1-] \quad \text{Entropy}(SFN) = 1$$~~

$$\text{form}(S, P) = 0.918 - \frac{2}{3} = 0.251$$

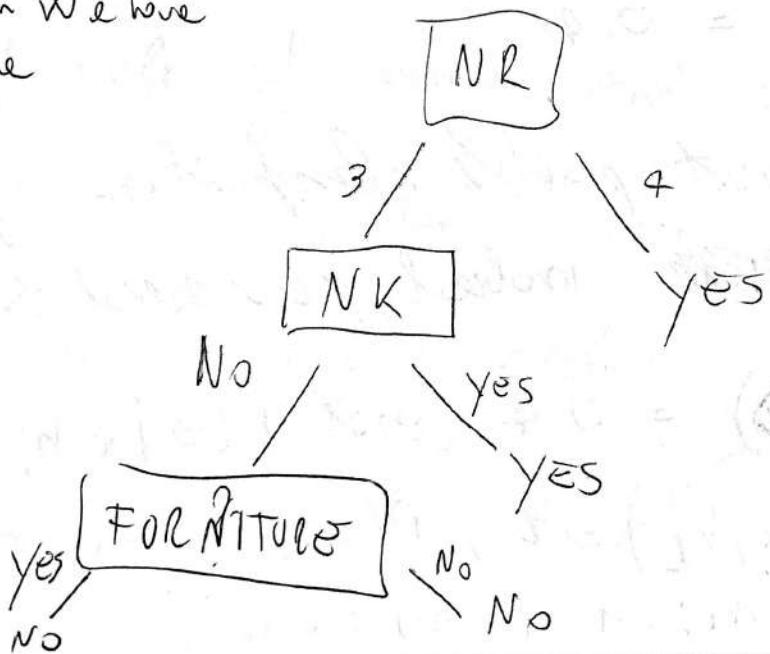
$$SNKY = [1+; 0] \quad \text{Entropy}(SNKY) = 0$$

$$SNKN = [0; 2+] \quad \text{Entropy}(SNKN) = 0$$

$$\text{form}(S, NK) = 0.918$$

We choose NK  
so then we have

Formiture



Ex 2

Given hypothesis space  $H$  give a data  $D$   
with  $h \in H$

$$P(h|D) = \frac{P(D|h)P(h)}{P(D)}$$

We compute the most probable hypothesis

If  $P(h_i) = P(h_j)$  with maximum likelihood

$$h_{MC} = \operatorname{argmax}_{h \in H} P(h|D) = \operatorname{argmax}_{h \in H} P(D|h)$$

We drop  $P(D)$  and  $P(h)$  because are constant

2) No, because if we consider 3 hypothesis

$$\begin{aligned} h_1 &= \oplus \\ h_2 &= \ominus \\ h_3 &= \ominus \end{aligned}$$

$$P(h_1|D) = 0.4 \quad P(h_2|D) = 0.3 \quad P(h_3|D) = 0.3$$

With maximum likelihood if we consider  $x$  a new instance we assign to it  $\oplus$  because

$$h_{MC} = \operatorname{argmax}_{h \in H} P(h|D) = 0.4$$

but this is NOT most probable classification of  
new instance ~~because~~ indeed you need to

$$\text{consider } P(\oplus|\oplus|x, D) = 0.4 \text{ and } P(\ominus|x, h_i) = 0.6$$

$$\text{If we consider } P(\oplus|x, h_1) = 1, P(\oplus|x, h_2) = 0, P(\oplus|x, h_3) = 0 \\ P(\ominus|x, h_1) = 0, P(\ominus|x, h_2) = 1, P(\ominus|x, h_3) = 1$$

We can say a instance  $x$  should be classified as + or not as - as we can see with maximum circled

Ex 3

1) We suppose to use a polynomial func

$$K(x'_i, x) = (\beta^T x + \gamma)^d \quad \text{where } d \geq \{2, \dots\}$$

In this case we can use a polynomial of degree 3 so  $d=3$

2)

We can use slack variable  $\varepsilon_n \geq 0$  when dataset is not perfectly separable and in particular for dataset  $D = \{(x_n, t_n)\}_{n=1}^N$  we should be under soft margin constraint in which  $t_n y(x_n) \geq 1 - \varepsilon_n$

If  $\varepsilon_n \geq 1$  means that this point is in the wrong side of boundary and in the margin

If  $0 < \varepsilon_n < 1$  point is inside the margin but in the right side.

### EX4

$$1) t = f(x, \theta)$$

We consider to use as activation function ~~linear~~

For output a sigmoid function because we consider  
~~linear~~ binary classification problem in which we  
 hope that it assume value 1 or 0 or  $\frac{1}{2}$   
 in output & result of sigmoid func is  
~~sigmoid Function~~ ~~function~~ there are  $\approx$ , work  
~~with binary or multiclass in this problem~~

2) For this use we choose ~~linear~~ <sup>sigmoid</sup> function  
 buts so we ~~don't~~ have problem of  
 saturation ~~but~~ only when we give the const  
~~these~~ ~~but~~ ~~don't~~ ~~saturate~~ on we

The act function that we can use is ~~softmax~~  
~~softmax over loss function~~ binary - complementary



### EX5

1)

$$y_{out} = \frac{w_{in} - 3 + 2p}{s} + 1$$

$$y_{out} = \frac{w_{in} - 3 + 2p}{s} + 1$$

$$r = w_{IN}$$

$$w_{IN} = \frac{w_{IN} - 3 + 2P}{S} + 1$$

$$h_{IN} = \frac{h_{IN} - 3 + 2P}{S} + 1$$

$$S \quad w_{IN} = w_{IN} - 3 + 2P + S$$

$$sh_{IN} = h_{IN} - 3 + 2P + S$$

$$S = \frac{w_{IN} - 3 + 2P}{w_{IN} - 1}$$

$$2P = S(h_{IN} - 1) + 3 - h_{IN}$$

$$P = \frac{S(h_{IN} - 1) + 3 - h_{IN}}{2}$$

$$= \frac{(w_{IN} - 3 + 2P)(h_{IN} - 1)}{2(w_{IN} - 1)}$$

$$P = \frac{(w_{IN} - 3)(h_{IN} - 1) + 3 - h_{IN} + 2P(h_{IN} - 1)}{2(w_{IN} - 1)}$$

$$P = \frac{(w_{IN} - 3)(h_{IN} - 1) + 3 - h_{IN}}{2(w_{IN} - 1)}$$

$$\frac{(w_{IN}h_{IN} - 3h_{IN} + 3 + h_{IN})}{2(w_{IN} - 1)} - 3h_{IN} - 2h_{IN} + 2h_{IN}$$

$$= \frac{2 \cdot \left( 1 - 2(h_{IN} - 1) \right)}{2(w_{IN} - 1)}$$

$$\frac{2(w_{IN} - 1) - 4(h_{IN} - 1)}{(w_{IN} - 1)}$$

$$= \frac{w_{IN}h_{IN} - 3h_{IN} + 3 - 3w_{IN} + 3 + w_{IN}h_{IN}}{2(w_{IN} - 1) - 4(h_{IN} - 1)} - 3(w_{IN} + h_{IN})$$

$$= \boxed{\begin{aligned} & \frac{-3(w_{IN} + h_{IN})}{2(w_{IN} - 1)} \\ & + \frac{2(w_{IN} - 2h_{IN} + 1)}{2(w_{IN} - 1) - 4(h_{IN} - 1)} \end{aligned}}$$

$$S = \frac{w_{IN} - 3}{w_{IN} - 1} + \frac{-3(w_{IN} + h_{IN})}{(w_{IN} - 1)(w_{IN} - 2h_{IN} + 1)}$$

2)

1) Hyperbolic Tangent

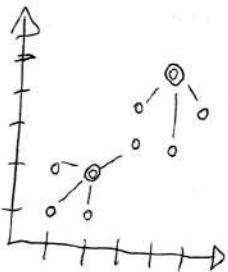
2) ReLU

3) Sigmoid

Ex 6

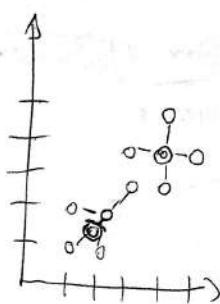
K = 2

With K-Means sequence and we correlate each sample in with nearest centroid we check if that sample is the cluster



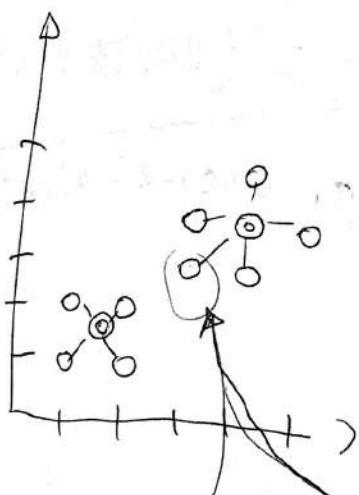
We show that

○ are centroid  
In the next  
step, we need  
to update  
our centers



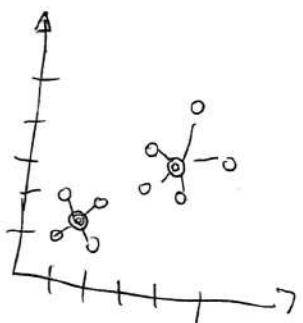
We update  
centroids

& Now we  
need to check  
if there is a  
sample that should  
be move from one  
to the other cluster



Now we move this sample to the other cluster  
~~on top~~ because it is cluster with nearest  
centroid

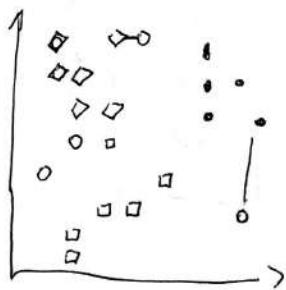
Goal to update ~~one~~ centroid of two clusters



We can say that there are no more points to switch between two cluster, so this ~~is~~ K-means converge

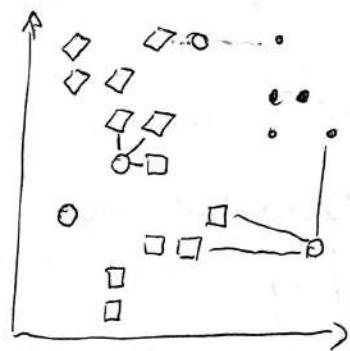
# EXAM TEST

EX1



■ K=1

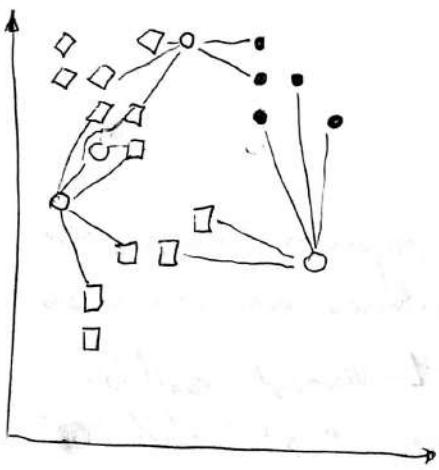
We choose 1-nearest neighbors to our new instances two of this 4 new instances we can classify because we take 1 1-nearest neighbor, in particular we assign "□" and "●" for two instances that we collected in the picture, the other two instances we cannot assign 1-nearest neighbors because we have more than 1 neighbor that are in the same distance respect to new instance, we cannot assign a classified value



■ K = 3

In this case we can classify two of four new instances, one with value "□" and because it is the most common value among 3-nearest neighbors and also for the other instance in picture we assign value "○".

For other two instances we have some problem than in the previous example there is problem that new instance is in the point that ~~more~~ two or more nearest neighbor are near to him at the same distance and you can choose best of them but you don't know which



■  $K=5$  we can compare values only on 3 of 4 max instances  
 In particular we can see that for 1 instance we compare value "□"  
~~another one we have 5~~ so this value should be "□" and then the last one has 3. 0 and 2 □ so we compare to it "●".  
 Only one instance cannot have a value because when we need to choose the last nearest neighbor there are two neighbors at the same distance, indeed it shouldn't be a problem for classifier because we have already 3 □ and 1 □ so it should be "□" but we don't know how to choose the last one nearest neighbor.

Ex2

layer

NN is convolutional ~~and not composed by~~  
 3 stage: convolution, detection and pooling

In convolutional stage we apply convolutional function

$$(I * K)(i, j) = \sum_{m=-M}^{M} I(i+m, j+m) K(m, m)$$

We can also apply some techniques such as sparse connectivity in which we consider that output depend of fewer inputs, we select a subset of input ~~parameters~~ that affect the output so in this way we can make less operations and increase the ~~efficiency~~ efficiency

We use another technique : parameter sharing. This consist that all units of that layer share parameter with each other and ~~this~~ with this we don't need to learn too much parameters for this layers and we can ~~use~~ re-use them many times.

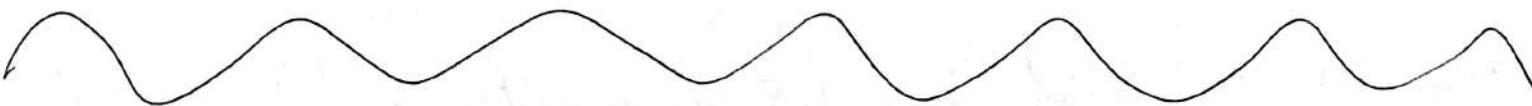
The last one thing in convolutional stage is the padding that it should be valid, so output contains only valid values, or some, so input layer is padded with zeros.

In the Detector stage we only decide which activation function we want to use for that layer, for example ReLu or hyperbolic tangent.

The last one stage is pooling stage where we modify the output, in particular pooling function replace ~~of~~ output of the net at certain location with a summary statistic of the nearby outputs

We have two different pooling: max pooling and average pooling  
Max pooling we consider maximum value of ~~neighbor~~ regions, instead in average pooling we consider ~~the~~ average value of ~~neighbor~~ region.

If we apply stride we reduce the size of output.



Ex 3)

Using PCA we reduce dimensionality of our dataset consisting of many parameters correlated with each other while retaining the variation present in the dataset.

We can do + Maximize variance or MINIMIZE error.

If we maximize variance we need to compute a dataset

$D = \{x_n\}_{n=1}^N$ , and our goal is to maximize variance after projection in some direction  $\mu_1$ .

Compute mean and variance

$$\bar{x} = \frac{1}{N} \sum_{n=1}^N x_n \quad \Rightarrow \text{projected mean } \mu_1^T x$$

$$\frac{1}{N} \sum_{n=1}^N (\mu_1^T x_n - \mu_1^T \bar{x})^2 = \mu_1^T S_{\mu_1} \mu_1 \quad & \begin{matrix} \text{projected} \\ \text{variance} \\ \text{when } S \text{ is variance} \end{matrix}$$

Maximize  $\mu_1^T S_{\mu_1} \mu_1$

We use Lagrange so we obtain

$$\max \mu_1^T S_{\mu_1} \mu_1 + \lambda_1 (1 - \mu_1^T \mu_1)$$

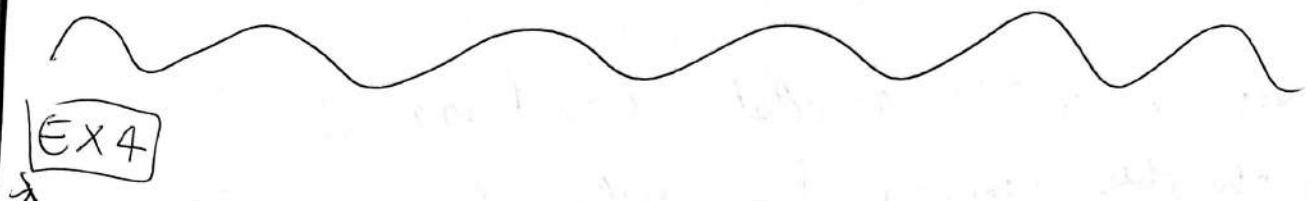
We set derivative to zero

~~$\mu_1$~~   $S_{\mu_1} \mu_1 + \lambda_1 \mu_1 = 0$

$$\mu_1^T S_{\mu_1} \mu_1 = \lambda_1$$

We can say that  $\mu_1$  is eigenvector and  $\lambda_1$  is eigenvalue

To maximize variance we need to find eigen vector  
 $M_1$  corresponding to the largest eigenvalue  $\lambda_1$ , this is  
principle component



EX 4

In MDP we consider to have

$$\langle X, A, S, r \rangle$$

Where  $X$  is set of states

In this case we have only 2 states, the initial state

$A$  is set of actions and in this case there are 3 possible paths

$S$  transition function is the probability that action  $a$  in state  $s$  at certain time  $t$  will lead state at time  $t+1$

$r$  is reward function, in this case reward ~~is~~ that we want to maximize

2) We don't know reward function so

What we need to use is the Q-learning function to find optimal hypothesis that maximize reward minimally

Time

training rule is

$$\hat{Q}_{\text{new}}(x, a) \leftarrow r + \gamma \max_a (\hat{Q}(x', a'))$$

where  $r$  is the immediate reward one  $x'$  is the state resulting from apply action  $a'$  in state  $x$

3)

Exploitation we try to select best action that maximize  $\hat{Q}(x, a)$  and exploration we select a random action  $a$  with low value of  $\hat{Q}(x, a)$

We prefer to use  $\epsilon$ -greedy strategy

We select random action with probability  $\epsilon$  and then select best action with probability  $1-\epsilon$

When  $\epsilon$  decrease over time

Exs

- 1) Confusion matrix is a performance measure used when we want to summarize result of mult. classes, in particular in the diagonal we have accuracy of each class and out of the diagonal we have the error of classification

$$2) f: \{0,1\}^4 \rightarrow \{\text{Low}, \text{Medium}, \text{High}\}$$

150 samples, 70 for low, 30 for medium, 50 for high

	Low	Medium	High	
Low	35	15	25	40
Medium	5	20	5	10
High	5	10	35	15
	10	25	30	90

3) Accuracy of <sup>65</sup> absorption matrix is given by

$$\frac{TP + TN}{TP + TN + FP + FN}$$

Where  $TP =$  true positive,  $TN =$  true negative,  $FP =$  false positive and  $FN =$  false negative

true positive are sample that are classified well for that class and true negative are value of the diagonal in which is not classified that particular class but the other class

For example for class low

$$TP = 35 \quad TN = 20 + 35$$

False negative are sample of that class ~~predicted~~ predicted some value in wrong way, for example for class low we predicted ~~not~~ ~~too~~ class incorrectly

Predict is of low to be High and is to be High

Finally False positive are true sample that are predicted in wrong way, for example

~~In the first column we have that goes mostly predicted to be low~~

In the first column we have 5 Medium that are wrongly predicted to be Low and 3 High that are wrongly predicted to be Low

$$\text{Accuracy} = \frac{90+90}{180+63+65} = 0.58 \Rightarrow 58\% \text{ accuracy}$$

Ex 6

1) Your answer is

$$E(Y) = \beta_0 + \beta_1 X$$

2) Yes the model used is in the form of linear regression and in particular to establish relation between input and output

$$y = \theta^T x \quad \theta = \begin{pmatrix} \theta_1 \\ \theta_2 \\ \theta_3 \\ \theta_4 \end{pmatrix} \quad x = \begin{pmatrix} 1 \\ x_1 \\ x_2 \\ x_3 \end{pmatrix}$$

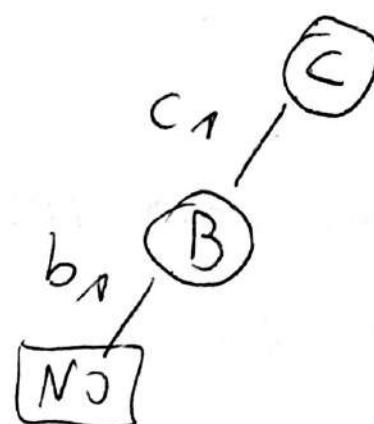
2) We can use the approach of Sequential learning in which we use the stochastic gradient

# ERASMIUS

## EX1

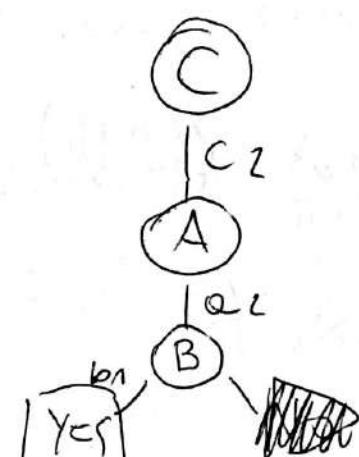
- 1) IF  $C = c_1 \wedge B = b_1$  THEN NO
- IF  $C = c_1 \wedge B = b_2$  THEN YES
- IF  $C = c_2 \wedge A = a_1$  THEN YES
- IF  $C = c_2 \wedge A = a_2 \wedge B = b_1$  THEN YES
- IF  $C = c_2 \wedge A = a_2 \wedge B = b_2$  THEN NO
- IF  $C = c_2 \wedge A = a_3$  THEN NO
- IF  $C = c_3$  THEN NO

- 2)  $\blacksquare S_1 = \langle a_1, b_1, c_1, \text{No} \rangle$



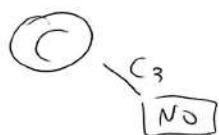
We can say that i draw part of tree given that value of set  $S_1$  and we can see that T is consistent with  $S_1$  because NO of the set is also result of tree for the values of attributes  $C, B, A$

- $\blacksquare S_2 = \langle a_2, b_1, c_2, \text{Yes} \rangle$



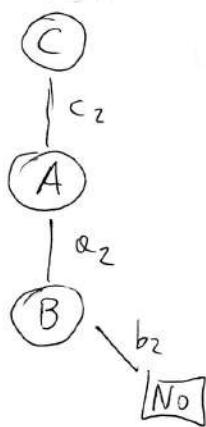
In this case T is ~~not~~ constant because with sequence  $c_2, a_2, b_1$  we obtain value Yes that is the same of set

■  $s_3 = \langle a_1, b_2, c_3, \text{Yes} \rangle$



In this case  $T$  is not consistent because we should have "No" value as we can see in the ~~the~~ tree but  $s_3$  set have value Yes so  $T$  is not consistent with  $s_3$

■  $s_4 = \langle a_2, b_2, c_2, \text{Yes} \rangle$



$T$  is not consistent with  $s_4$  for the same reason in which  $T$  is not consistent with  $s_3$

### EX 2

1) Maximum a posteriori is used to compute most probable hypothesis

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

We drop  $P(D)$  because it is constant and independent of  $h$

↳ IF  $p(h_i) = p(h_j)$

We have maximum likelihood

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

We drop  $p(h_i)$  because it is constant

Bayes optimal classifier is classifier that we use to determine the most probable classification for new instance

Given a dataset  $D$  and target function  $f: X \rightarrow V$ ,  $V = \{v_1, \dots, v_n\}$  and a new instance,  $x \notin D$

$$V_{OB} = \operatorname{argmax}_{v_j \in V} \sum_{h_i \in H} p(v_j | x, h_i) P(h_i | D)$$

~~if~~  $h_i$  is independent from  $D$  so

$$p(v_j | x, h_i, D) = p(v_j | x, h_i)$$

and also  $h_i$  is independent from  $x \notin D$

$$p(h_i | D, x) = p(h_i | D)$$

3) Bayes optimal classifier is optimal because it uses prior probability hypothesis to compute result but it's not good when we have huge hypothesis space because we cannot have prior probability of each hypothesis, so in this case we use approximation using Naive Bayes Classifier

[Ex 3]

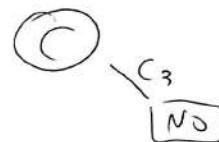
Perceptron is a model used for linear classification

In particular we have function  $\delta(x) = \begin{cases} 1 & \text{if } w^T x > 0 \\ -1 & \text{otherwise} \end{cases} \Rightarrow \delta = w^T x$

in which we have  $y = w^T x$  and dataset  $D = \{(x_n, t_n)\}$

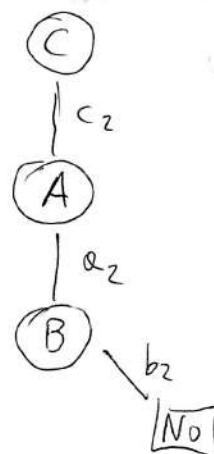
The aim is to find good separation to classify new instances. We start with value of  $w$  that is random and then set

■  $s_3 = \langle a_1, b_2, c_3, \text{Yes} \rangle$



In this case  $T$  is not consistent because we should have "No" value as we can see in the ~~the~~ tree but  $s_3$  set have value Yes so  $T$  is not consistent with  $s_3$

■  $s_4 = \langle a_2, b_2, c_2, \text{Yes} \rangle$



$T$  is not consistent with  $s_4$  for the same reason in which  $T$  is not consistent with  $s_3$

## Ex 2

1) Maximum a posteriori is used to compute most probable hypothesis

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

We drop  $P(D)$  because it is constant and independent of  $h$

↳ IF  $p(h_i) = p(h_j)$

We have maximum likelihood

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

We drop  $p(h_i)$  because it is constant

Bayes optimal classifier is classifier that we use to determine the most probable classification for new instance

Given a dataset  $D$  and target function  $f: X \rightarrow V$ ,  $V = \{v_1, \dots, v_n\}$   
and a new instance,  $x \notin D$

$$V_{OB} = \operatorname{argmax}_{v_j \in V} \sum_{h_i \in H} P(v_j | x, h_i) P(h_i | D)$$

~~if~~  $h_i$  is independent from  $D$  so

$$P(v_j | x, h_i, D) = P(v_j | x, h_i)$$

and also  $h_i$  is independent from  $x \notin D$

$$P(h_i | D, x) = P(h_i | D)$$

3) Bayes optimal classifier is optimal because it uses prior probability of hypothesis and to compute result but it is not practical when we have huge hypothesis space because we cannot have prior probability of each hypothesis, so in this case we use approximation using Naive Bayes Classifier

Ex 3

Perceptron is a model used for linear classification

In particular we have function  $\delta(x) = \begin{cases} 1 & \text{IF } w^T x > 0 \\ -1 & \text{otherwise} \end{cases} \Rightarrow \delta = w^T x$

in which we have  $y = w^T x$  and dataset  $D = \{(x_n, t_n)\}$

The aim is to find good separation to classify new instances.  
We start with value of  $w$  that is random and then at

each iteration this value change with stochastic gradient  
in particular we minimize error:

$$E(w) = \frac{1}{2} \sum_{n=1}^N (t_n - w^\top x_n)^2 = \frac{1}{2} \sum_{n=1}^N (t_n - w^\top x_n)^2$$

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

We use stochastic gradient descent

$$w \leftarrow w + \Delta w_i$$

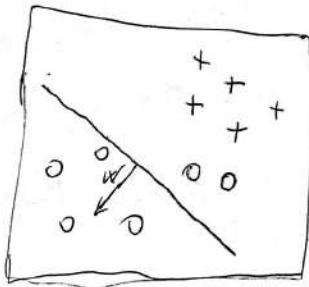
$$\Delta w_i = -\eta \nabla E \quad \text{where } \eta \text{ is learning rate}$$

$$= +\eta \sum_{n=1}^N (t_n - w^\top x_n) x_{i,n}$$

This means that ~~we~~ the value of  $w$  change at each iteration  
when we find instance in the dataset that is missclassified  
this cycle terminate ~~when~~ after some finite iteration if dataset  
is linearly separable and  $\eta$  is small otherwise it will never  
converge.

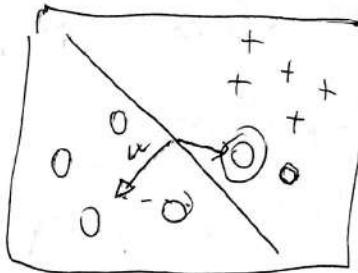
2)

a)

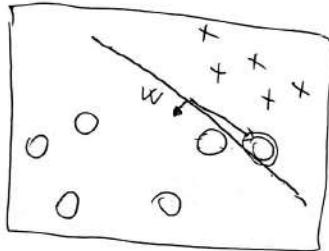


We initialize with decision boundary  
input space because we choose  $w$   
randomly, on the better ~~the~~ should be  
instances of ~~the~~ class '+' and one  
instance of class '-'

b)

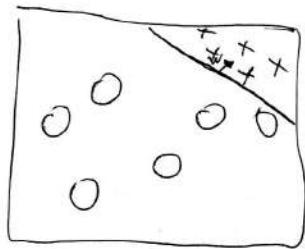


A class is missclassified, so we  
need to update our weight  
and this means that decision boundary  
change and also value of  $w$



As we can see there is another instance that is miss classified so we need to update weights again

d)



Now we have ~~one~~ that perception converge because there are no more instance missclassified and ~~the~~ dataset that we create is linearly separable, and we can do it in small

**Ex 4**

1)  $W_1$  We pass from 128 to 50 units so dimensions of  $W_1$  should be ~~128~~  $128 \times 50$  and dimension of  $W_2$  should be  $50 \times 10$

2) we have ReLU activation function for hidden layer so

$$h = g(x^T W_1 + b)$$

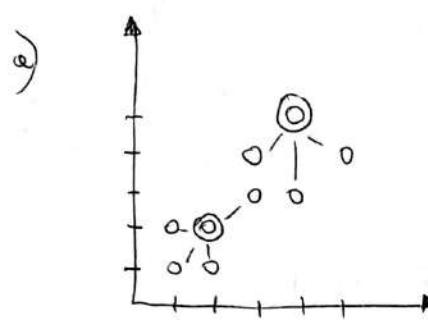
$$\text{where } g = \max(0, \cdot)$$

for output we have linear activation function so consider

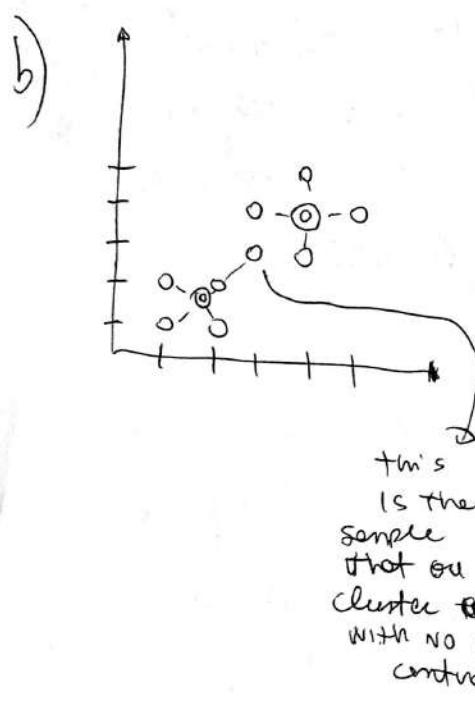
$$y = W_2^T h + c = W_2^T \cdot \max(0, x^T W_1) + b + c$$

# Ex 5

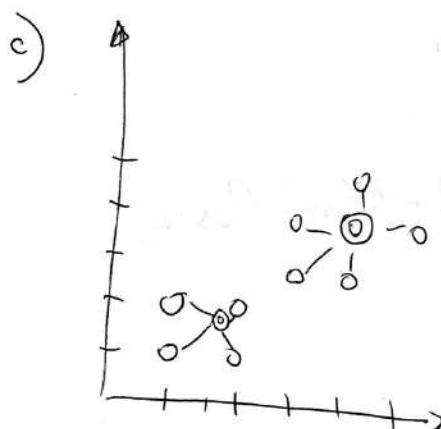
$K = 2$



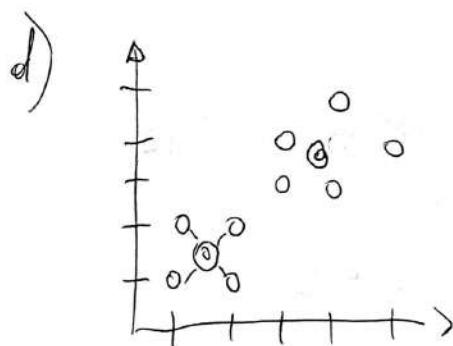
~~We have 2 sets samples in~~  
we understand that we want to update our centroid and it show in the picture sample belong to each cluster.



As we can see we update our centroid and then we take in sequence each sample to check if it stay in cluster with nearest centroid  
We see that there is a cluster that is not belong to the ~~00~~ cluster with nearest centroid for it so we need to switch cluster for this sample



Now we need to recompute centroids of the clusters involved



Now after update centroid we take each sample ~~one~~ in sequence and check again, we can see that all sample are in ~~one~~ cluster with nearest centroid so 2-Means converge

Autoencoder is a neural network with reduced size of hidden layer and learn to reconstruct

~~structure~~ structure of inputs by minimize sum-of-squared error.

In particular output and input have same size and hidden layers have a size minor of inputs and outputs

For example we have  $D$  for input and  $M$  for hidden

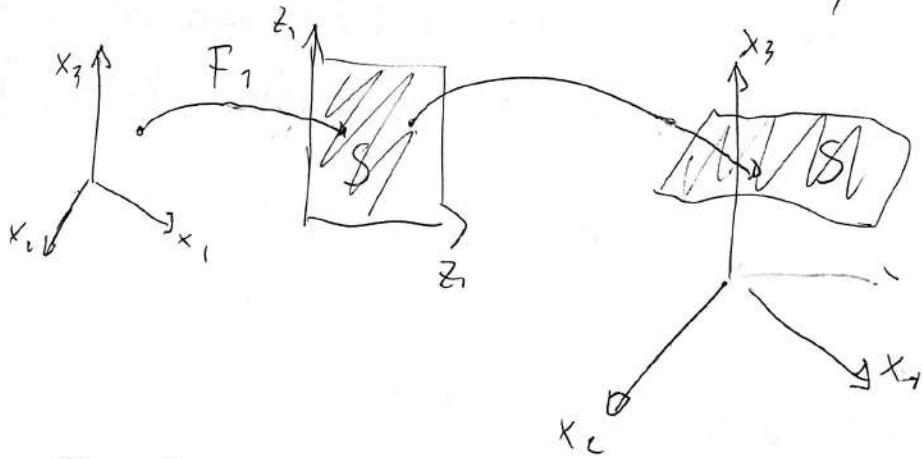
$$\begin{matrix} & F_1 & F_2 \\ \begin{matrix} x_{d0} \\ 0 & 0 & 0 & \dots & 0 \\ 0 & 0 & 0 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & \dots & 0 \\ x_d \end{matrix} & \xrightarrow{\quad} & \begin{matrix} x_1 \\ \vdots \\ x_M \end{matrix} \end{matrix}$$

$F_1$  convert dimensionality of  $D$  to  $M$  and  $F_2$  convert dimensionality of  $M$  to  $D$

These are non linear function for hidden layer

2)

We consider 3D - input , 2D - hidden , 3D - output



$F_1$  map from  $D$  to  $M$

$F_2$  map from  $M$  to  $D$

When  $M = 2$

$$D = 3$$

If  $F_2$  is non linear  
we obtain 1/V entries  
o non linear function

6  
Autoencoder is a neural network with reduced size of hidden layer and learn to reconstruct

~~keep~~ structure of inputs by minimize sum-of-squared error.

In particular output and input have same size and hidden layers have a size minor of inputs and outputs

For example we have  $D$  for input and  $M$  for hidden

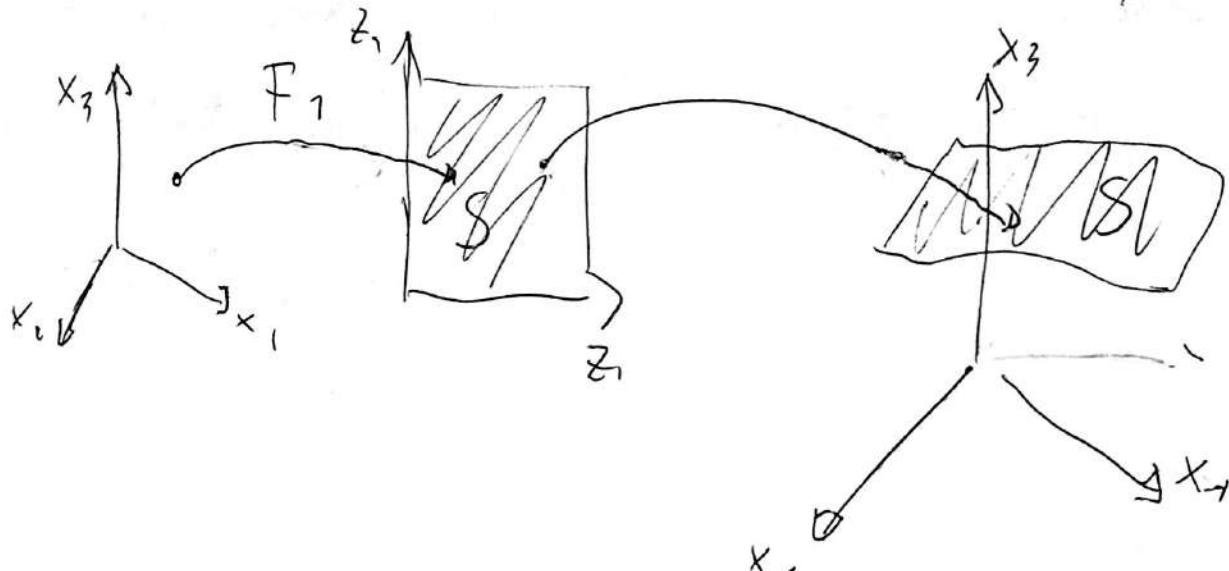
$$\begin{matrix} & F_1 & F_2 \\ \begin{matrix} x_{d0} \\ 0 & 0 & 0 & \dots & 0 \\ 0 & 0 & 0 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & \dots & 0 & \dots & x_1 \end{matrix} & \xrightarrow{\quad} & \begin{matrix} x_1 \\ \vdots \\ x_M \end{matrix} \end{matrix}$$

$F_1$  convert dimensionality of  $D$  to  $M$  and  $F_2$  convert dimensionality of  $M$  to  $D$

These are non linear function for hidden layer

2)

We consider 3D - input , 2D - hidden , 3D - output



$F_1$  map from  $D$  to  $M$

$F_2$  map from  $M$  to  $D$

When  $D = 3$

If  $F_2$  is non linear  
we obtain non linear  
and non linear function

EX1

Given on Hypothesis space  $H$ , we can say that  $h \in H$  overfit training data if exist  $h' \in H$  such that:

$$\text{error}_S(h) < \text{error}_S(h') \Rightarrow \text{error}_D(h) > \text{error}_D(h')$$

Where

$\text{error}_S(h)$  is sample error  
 $\text{error}_D(h)$  is true error

2)

In Decision tree we consider to have overfitting when for example we have dataset that contains noisy data and ~~the decision tree~~ we have  $T_1$  is and two decision tree  $T'$  and  $T$  made of different ID's.

We can have that  $T'$  has accuracy  $>$  than  $T$  but this doesn't mean that  $T'$  is better than  $T$  because  $T'$  overfitting causes of noisy data so there are different technique to avoid overfitting, in particular we can:

- 1) Stop growing tree, but it is not easy to do because it's difficult to estimate when to stop
- 2) Grow full tree then post-prune

There are two different techniques of post pruning

- 1) Reduced Error Pruning
- 2) Rule post Pruning

- Reduced Error pruning we split dataset in training and validation set, then we until pruning is harmful we evaluate ~~and~~ the impact on validation set of pruning each possible node then greedily remove node that most improves validation set accuracy
- This method give us optimal subtree and furthermore when tree size is small ~~is~~ it can give us best results
- Rule post Pruning we infer the tree allowing overlapping then we convert tree in set of rules this because in this way we can distinguish each path. In this way we can prune each path independently to each other at the end after pruning we order set of rules as we want.

3

$$D = \{(x_1, s_1), p_1\}, \dots, (x_N, s_N), p_N\}$$

1) likelihood function

$$P(P|W) = \prod_{n=1}^N y_n^{p_n} (1-y_n)^{1-p_n}$$

where  $y_n = \delta(w^T x_n)$

$w$  - vector of weights

$$x = \begin{pmatrix} (x_1, s_1) \\ \vdots \\ (x_N, s_N) \end{pmatrix}$$

2) We should learn  $w$ , weights  
in particular

$$w^* = \underset{w}{\operatorname{arg\min}} E(w)$$

3)  $E(w)$  is cross entropy

$$E(w) = -\ln P(P|W) = -\sum_{n=1}^N p_n \ln y_n + (1-p_n) \ln(1-y_n)$$

Exs

- 1) In stochastic and non deterministic actions, state result of execution of an action is not known before execution of the action but it is fully observable after execution.

2) MDP is a process that use decision policy and it is expressed by,

$$\langle X, A, S, r \rangle$$

X is set of states

A is set of actions  
and states

These are fully observable

S is transition function

r is reward function

Partial POMDP is combination of MDP and thus  
se POMDP has also observation respect to MDP with  
partial observation function, this means that ~~receives~~ receives states on

$$\langle X, A, Z, S, r, \delta \rangle$$

X set of states

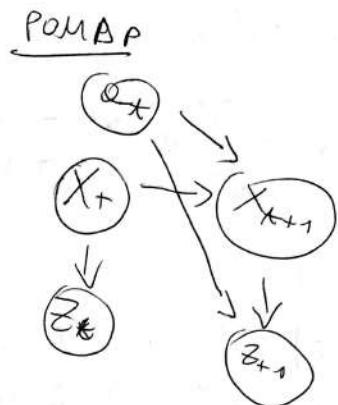
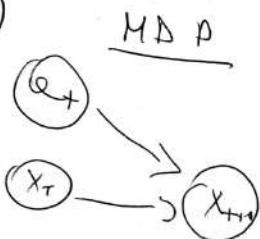
A set of actions

Z set of observation

S is transition function

r is reward function

$\delta$  is observe function



(Ex)

- 1) Boosting is a method wifh whom we want to use different classifier and combine their result.  
In particular Boosting use different classifier in sequence. Each learner is trained in sequence and it is trained on weighted data. Each weighted data depends on performance of the previous classifier. Points that are missclassified give vs ~~to~~ high weight.  
At the end predictions ~~of~~, ~~are~~ be final output is done with majority of votes.
- 2) AdaBoost is a kind of Boosting in particular the algorithm is const, given a dataset  $D = \{(x_n, t_n)\}_{n=1}^N, x_n \in X$ , We initializ  $w_n^{(1)} = \frac{1}{N}$ ,  $n=1, \dots, N$ ,  $t_n \in \{-1, +1\}$   
then for  $m = 1 \dots M$  do

We train weak learner  $y_m(x)$  by minimize weighted error function

$$T_m = \sum_{n=1}^N w_n^{(m)} I(y_m(x) \neq t_n)$$

With

$$I(e) = \begin{cases} 1 & \text{IF } e \text{ is True} \\ 0 & \text{IF } e \text{ is False} \end{cases}$$

then we evaluate  $\epsilon_m$  and  $d_m$

Update weighted coefficient

$$w_n^{(m+1)} = w_n^{(m)} \exp(d_m I(y_m(x) \neq t_n))$$

Output of classifier is

$$y_M(x) = \operatorname{sgn} \left( \sum_{m=1}^M d_m y_m(x) \right)$$

Ex

- 1) SVM we want to find MAXIMUM margin ~~margin~~ to provide better accuracy. In particular we find hyperplane with maximum margin separation on the training sets, in this way we classify our points in good way respect margin that is too small

1) Unsupervised learning is learning that has a dataset  
 $D = \{x_n\}_{n=1}^N$ . We have input but no value of target function corresponding to the input (is not available). In unsupervised learning we want to find a structure of input using different method such as clustering.

2)



3) We want to find  $k$ -means provided by  $K$  Gaussian distributions

Given a dataset  $D = \{x_n\}$  and  $K=2$ , such as previous problem

- 1) We consider  $K=2$  number of clusters
- 2) Now we need to do an initial partition that classifies data in  $K$  clusters

- 3) We take first  $K$  training samples ~~as~~ as initial clusters
- 4) Assign to ~~the~~ remaining  $N-K$  training sample to cluster with Nearest centroid.  
After each assignment we update centroid of ~~each~~ clusters
- 5) Take each sample in sequence and we check if that sample is in the cluster with Nearest centroid  
If sample is not in cluster with Nearest centroid  
we switch cluster for this sample and then  
We update centroids of two clusters involved
- 6) We repeat ③ UNTIL the algorithm converge,  
it converge when there are no more sample that  
need to switch cluster.



i)

Backpropagation is an algorithm that is used to compute gradient with respect to parameter  $\theta$  of loss function;  $\nabla_{\theta} J(\theta)$

We want to propagate this gradient from output to input after computing loss function.

Backpropagation is algorithm that is used to compute only gradient, it is NOT a training algorithm for neural network.

Forward and back pass are two steps of backpropagation algorithm

In forward pass we compute loss function step by step from input and get to the output, then in backpass we use loss function output to compute gradient and propagate this in each layer from output layer to input.

Stochastic Gradient Descent is a training algorithm that want to improve performance of neural network using gradient. ~~gradient~~

In particular is an algorithm that update parameters

$$\theta \leftarrow \theta^{(k+1)} \leftarrow \theta^{(k)} + \eta_j g$$

where  $g$  is gradient

## 2) Forward pass

Require:  $l$  depth of neural network

Require:  $W^{(i)}$  matrices of weights  $i = 1, \dots, l$

Require:  $b^{(i)}$  bias  $i = 1, \dots, l$

Require: input vector  $x$

Require: target vector  $t$

$$h^{(0)} = x$$

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

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

$$h^{(k)} = f(d^{(k)})$$

end for

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

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

where  $J$  is cost function; loss function

## Backprop

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

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

Propagate gradient to non-linear activations

$$g \leftarrow \nabla_{x^{(k)}} T = g \odot f'(x^k)$$

element wise product

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

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

Propagate gradients to the next layer and hold layer

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

End for

**Ex**

1) We have a binary classification problem so far

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

we can use for activation function in output sigmoid function

$y = \sigma(w^T h + b)$  and cost function should be

binary cross entropy

$$T(\theta) = \text{softmax}(1 - 2t) \cdot \lambda$$

$$\lambda = w^T h + b$$

2)

As we can see several junks we have that our units structure only when ~~the~~ they give us the correct ~~correct~~ answer.

Ex

1)

$$W_{out} = W_{in}$$

$$W_{out} = \frac{W_{in} - 3\text{Kernel} + 2p}{S} + 1$$

$$W_{out} = \frac{W_{in} - 3 + 2p}{S} + 1$$

$$\textcircled{1} \quad W_{in} \cdot S = W_{in} - 3 + 2p + 5$$

$$\textcircled{2} \quad S = \frac{W_{in} - 3 + 2p}{W_{in} - 1}$$

$$h_{out} = \frac{h_{in} - 3 + 2p}{S} + 1$$

$$h_{in} \cdot S = h_{in} - 3 + 2p + 5$$

$$S = \frac{h_{in} - 3 + 2s(h_{in}-1) + 6}{h_{in}-1}$$

$$2p = \cancel{h_{in}} s(h_{in}-1) + 3$$

$$p = \frac{s(h_{in}-1) + 3}{2}$$

$$S = \underbrace{W_{IN} - 3 + S_{IN} - S}_{W_{IN-1}} + 6$$

~~$$S_{(W_{IN}-1)} + S = W_{IN} + 6$$~~

~~$$S (W_{IN} - 1 + 1 - h_{IN}) = 3$$~~

~~$$S (W_{IN} - h_{IN}) = 3$$~~

$$S = \frac{3}{W_{IN} - h_{IN}}$$

$$P = \frac{3(W_{IN}-1) + 3}{2(W_{IN} - h_{IN})}$$

## Tex

Convolution layer is composed of 3 stages:

convolutional stage, detector stage and pooling stage.

In Convolutional stage we have convolution function between ~~Kernel~~ input and Kernel

$$(I * K)(i, j) = \sum_m \sum_m I(i+m, j+m) K(m, m)$$

We can apply different techniques.

④ Sparse connectivity and parameter sharing

With sparse connectivity we have output of layer depends only of few input in this way we can reduce number of parameter and increase statistically efficiency

With parameter sharing we have that all units of that layer share subset of parameter, this means that we don't learn more parameter and we can use <sup>some</sup> parameter ~~for other~~ with other functions.

Finally we have pooling that could be same or valid; If it some input is padded with zeros and output doesn't depend on kernel size

If pooling is valid output contains ~~as~~ valid values and output depends of kernel size.

After convolutional stage we have linear function that ~~in~~ in detector stage use some non linear function, like ReLU, to transfer those linear func in non linear.

In fully stage we modify the output of neural network, in particular ~~with~~ pooling function replace ~~at~~ at certain position the output with sum of statistically outputs nearby

we see two different Pooling form:

Max pooling: we take maximum of neighbors outputs

Average Pool: we take average of neighbors outputs

With stride we decrease size of output

Ex

PCA is technique reduce dimensionality of dataset  
With many parameter correlated with each other  
while keeping present variation.

PCA maximize variance or minimize error

When we maximize variance we need to consider points projected to perturb direction

We need to ~~choose~~ consider direction vector  $\mu_i$

We compute mean project  
and variance project

and after ~~choose~~ maximise of it, so ~~choose~~  
use covariance and then ~~choose~~ set it down to  
zero we obtain

$S\mu_i = \lambda \mu_i$  when  $S$  is variance, we know

$$\mu_i^T S \mu_i = \lambda$$

To MAXIMIZE variance we need to choose  $\mu_i$ ,  
(as eigen vector correspondingly to largest ~~large~~ eigen values  
thus is principal components

EX

→ K Fold cross validation is a method to ~~estimate~~ estimate accuracy of learning algorithm

We consider dataset  $D = \{(x_n, t_n)\}_{n=1}^N$

We split our dataset in  $K$   $S_i$  with  $|S_i| \geq 30$

for  $i = 1, \dots, K$

$S_i$  is test set and  $T_i$  is training set where

$$T_i = \cancel{D - S_i}$$

$$h_i = L(T_i)$$

We compute prediction using function  $L$

$$\delta_i = \text{error}_S(h_i)$$

sample error of  $h_i$

$$\text{error}_S(h) = \frac{1}{K} \sum_{i=1}^K \delta_i$$

$$\text{accuracy}_{S, D}(h) = 1 - \text{error}_{S, D}(h)$$

2) K-fold is a validation method used also when you need to compare which algorithm is better in terms of accuracy

We consider  $(A, L_B)$

We split our dataset  $D$  in  $K$   $S_i$   $|S_i| \geq 30$

for  $i = 1, \dots, K$

$S_i$  is test set

$T_i$  is training set =  ~~$\{D\}$~~   $\{D - S_i\}$

$$h_A = L_A(T_i)$$

$$h_B = L_B(T_i)$$

We compute prediction of A and B over some testing set

$$\delta_i = \text{error}_S(h_A) - \text{error}_S(h_B)$$

$$\bar{\delta} = \frac{1}{K} \sum_{i=1}^K \delta_i$$

If  $\bar{\delta} < 0$  we can say that  $L_A$  is better than  $L_B$  because  $\text{error}_S(h_A)$  is minor of  $\text{error}_S(h_B)$

Ex

One state markov process we have

$$(X_0, A, S, r)$$

where  $X_0$  is unique state

$A$  is set of actions

$S$  is transition function

$r$  is reward function

K-armed bandit problem is stochastic problem for one-state MDP

In particular we have  $K$  actions and

$v(e_i) = N(\mu_i, \sigma_i)$  gaussian distribution  
 K-armed Bandit use  $\epsilon$ -greedy strategy  
 In which we choose ~~action~~ ~~not~~ random action with probability  $\epsilon$  and then we select best action with probability  $1-\epsilon$  and  $\epsilon$  decreases over time

2) The强化 rule is based on Q-learning

$$Q_n(e_i) \leftarrow Q_{n-1}(e_i) + \alpha [r_t + \gamma Q_{n-1}(e_i)]$$

$$\alpha = \frac{1}{V_{n-1}(e_i)}$$

When  $V_{n-1}(e_i)$  no. of ~~action~~ exclude of  $e_i$   
upto time  $n-1$

Ex

1) Generative is a probabilistic model that ~~use~~ use Bayes rule to estimate

$P(C_i | x)$  for classification, instead discriminant probabilistic model use generative model

to estimate value of  $P(C | x)$ , in particular logistic regression

CNN

1) 2 layer

■ Output of CNN is  $W_{out}$  and  $h_{out}$

$$W_{out} = \frac{W_{in} - W_{kernel} + 2 \cdot padding}{\text{Stride}} + 1$$

$$h_{out} = \frac{h_{in} - h_{kernel} + 2 \cdot p}{\text{Stride}} + 1$$

- Output value is sum of output in input because activation function doesn't change dimension of output
- output of pool is

$$W_{out} = \frac{W_{in} - W_{pool}}{\text{Stride}} + 1$$

$$h_{out} = \frac{h_{in} - h_{pool}}{\text{Stride}} + 1$$

We consider first layer

We have often conv 1

$$W_{out} = \frac{1242 - 5 + 4}{4} + 1 = 1242$$

$$h_{out} = 378$$

$\Rightarrow$  output after conv 1 is  $1242 \times 378 \times 64$

where 64 are features

- after relu dimension doesn't change and after pool it will halve

$$W_{out} = \frac{1282 - 2}{2} + 1 = 621 \quad h_{out} = \frac{388 - 2}{2} = 183$$

Output of first layer is  $621 \times 183 \times 64$

thus should be input of second layer

- Conv 2 -  $3 \times 3$  AND 128 feature map = 0 stride = 2  
Output is given by  $621 \times 183 \times 64$

$$W_{out} = 310$$

$$h_{out} = 94$$

We obtain  $310 \times 94 \times 128$

- We can skip pool 2 because dimensions don't change after pool 2

$$W_{out} = 78$$

$$h_{out} = 24$$

Final output  $78 \times 24 \times 128$

- 2) N° of neurons depend only of convolution stage  
this because pool and ~~ReLU~~ Relu doesn't contribute of number of neurons

$N_1 = N^o$  neurons after first layer

$$= W_k \times h_k \cdot \text{Feature}_{in} \cdot \text{Feature}_{out} + \text{Feature}_{out}$$

$$= 5 \cdot 5 \cdot 3 \cdot 64 + 64 = 4864$$

$$N_2 = 3 \cdot 3 \cdot 6 \cdot 178 + 178 = 73856$$

Ex

- 1) the dimensionality of obespace is the input  
Size, so in this case  $12 \times 12 = 144$   
the intrinsic dimensionality of your obj is  
translation and rotation, since is the same  
image we have 2 translation and 1 rotation  
so intrinsic dimensity is 3

2)



$$\tilde{x}_n = \sum_{i=1}^M z_{ni} \cdot w \quad \text{where } z_{ni} = x_i^T w$$

means that  $z_{ni}$  depend on particular object part

- 3) Not necessary, usually ~~M~~ M is usually  
different of intrinsic dimension because PCA is  
linear so ~~non-linear~~ principal component of PCA are not  
latent variable

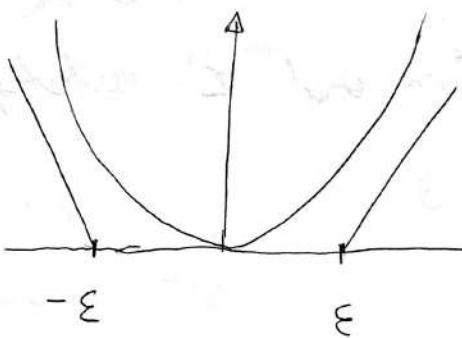
$$Q(x, \alpha) \leq \bar{V} + \gamma \max_{\alpha'} (Q(x', \alpha'))$$

$$Q_n(\alpha_i) \leq \max_{\alpha'} Q_{n-1}(\alpha_i) + \gamma [v(\alpha_i) +$$

$$\cdot L = \frac{1}{v_{n-1}(\alpha_i)} [Q_{n-1}(\alpha_i)]$$

**EX3**

1)



The problem to minimize  $J$  is that ~~as~~  $\epsilon$  is not differentiable

2)  $\epsilon_n^+$  and  $\epsilon_n^-$  slack variable with  
 $\epsilon_n^+ > 0 \quad \epsilon_n^- > 0$

The role of them is to ~~be~~ be minimum of  $J$  without probk of not differentiable met by  $\epsilon$  in particular we have

$$J(V, K) = C \sum_{i=1}^N (\epsilon_n^+ + \epsilon_n^-) + \frac{1}{2} \|w\|^2$$

15/10/2019

EX1

1) CONFUSION MATRIX is a performance measure used to ~~suspect~~ in a classification problem with many classes, in the diagonal are the true accuracy of each class. Furthermore we compute how many times an instance of class  $C_i$  is classified in class  $C_j$  outside the diagonal are false errors.

EXAMPLE

Predicted value

	$C_1$	$C_2$	$C_3$
$C_1$			
$C_2$			
$C_3$			

2)

3 classes generate 100 sample of each class

2)

	FN			
	C <sub>1</sub>	C <sub>2</sub>	C <sub>3</sub>	
C <sub>1</sub>	70	20	10	30
C <sub>2</sub>	15	75	10	25
C <sub>3</sub>	25	15	60	40
FP	40	35	20	205

	C <sub>1</sub>	C <sub>2</sub>	C <sub>3</sub>
C <sub>1</sub>	0.2	0.2	0.2
C <sub>2</sub>	0.15	0.75	0.1
C <sub>3</sub>	0.25	0.15	0.60

$$3) \frac{TP}{FP} = \frac{TN}{FN} = 205 \quad FN = 95 \\ FP = 95$$

$$Acc = \frac{TP+TN}{TP+TN+FN+FP} = \frac{410}{410+190} = 0.68$$

EX2

1) We perform regression method to find good relationship between inputs and outputs in particular in this case we can choose polynomial of degree 3

$$y(x_i) = w_0 + w_1 x_1 + w_2 x_2^2 + w_3 x_3^3$$

w<sub>0</sub>, w<sub>1</sub>, w<sub>2</sub>, w<sub>3</sub> are weights

x input variable

2) In regression we use regularization to avoid overfitting, in particular we add a

regularization term with error  $E_D$

$$E_D(w) \frac{1}{2} \sum_{n=1}^N (t_n - w^T x_n)^2$$

In particular we have

$$\text{Min } E_D(w) + \lambda E_W(w)$$

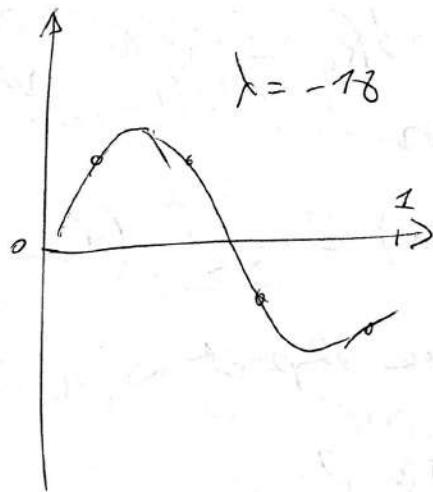
where  $\lambda$  is regularization factor,  $\lambda > 0$

And  $E_W(w)$  is regularization term

We can choose  $E_W(w) = \frac{1}{2} w^T w$  for our model

3)

the plot is



Ex 3

- 1) Perception is modeled using in linear classification  
in particular it has aim to classify perfectly inputs  
of object or updating weight in each iteration if  
it ~~is~~ Find some instances that ~~are~~ are classified  
in bay way

IT update weights using learning rate and gradient of error, but we can say that after this operation the decision boundary change until we don't find any instances that are not misclassified (IF it's possible) usually depends on dataset)

2) we have  $y = w^T x$ , but in particular in this case we have function  $\sigma(x) = \begin{cases} 1 & \text{if } w^T x > 0 \\ -1 & \text{otherwise} \end{cases}$

Given a dataset  $D = \{(x_n, t_n)\}_{n=1}^N$

minimise loss function

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

We compute derivative respect to  $w$  and we obtain

$$\frac{\partial E(w)}{\partial w_i} = - \sum_{n=1}^N (t_n - w^T x_n)^2 x_{in} = \bar{V}_E$$

$$w \leftarrow w + \Delta w_i$$

where

$$\Delta w_i = -\eta \bar{V}_E = +\eta \sum_{n=1}^N (t_n - w^T x_n)^2 x_{in}$$

$\eta$  is learning rate

found a decision border that classify our instances  
 in a much way there is no instance that is  
 misclassified, if dataset  $\Delta$  is linearly separable  
 and if learning rate  $\eta$  is sufficiently small.  
 If  $\eta$  is too small, the convergence will be  
 too slow.

Ex 4

- 1) In SVM we want to maximize margin to  
provide better accuracy.  
Margin is the smallest distance between closest  
points to the hyperplane (decision border) and  
the hyerplane itself.  
To find  $w^*$  and  $w_0^*$  we minimize that margin,  
thus means that we need to use Lagrange  
Method to compute  $w^*$  and  $w_0^*$ .  
 $w^*$  is compute using support vectors  
of the inst we have

$$y(y(x')) = \text{sign} \left( \sum_{k=1}^m w_k x_k^{*T} x' + w_0^* \right)$$

At the end we often that we have two closest  
 points  $x_k^{(1)}$  and  $x_k^{(2)}$  who distance should

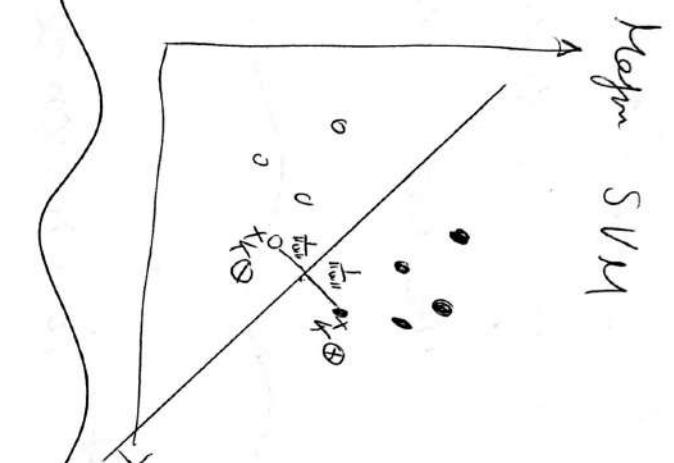
be the same  $\frac{1}{\sqrt{N} \cdot \text{H}}$

$$w^T x_k^\oplus + w_0^\oplus = -1$$

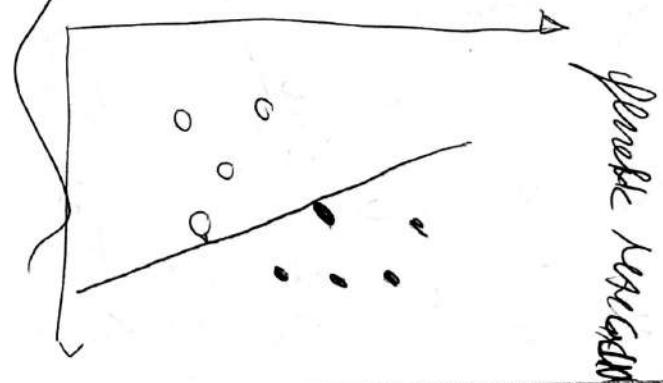
IN FACT  $w^T \cdot x_k^\oplus + w_0 = -1$  OR

2)

i)



i)



Ex 5

→ K Nearest Neighbor algorithm is algorithm  
very for non parametric model.

In particular we consider a shortest D ->  $\{x_i, t_i\}_{i=1}^n$   
and ~~Centroid~~ function  $f: X \rightarrow \mathbb{K}$  and  $K$ , and  $x \notin D$   
the ~~first~~  $\Rightarrow$  steps is to compute K nearest  
neighbors to instance  $x$  then we can find the  
K-nearest neighbor  
In particular k-th nearest function is

K-nearest neighbor

In particular k-th nearest function is

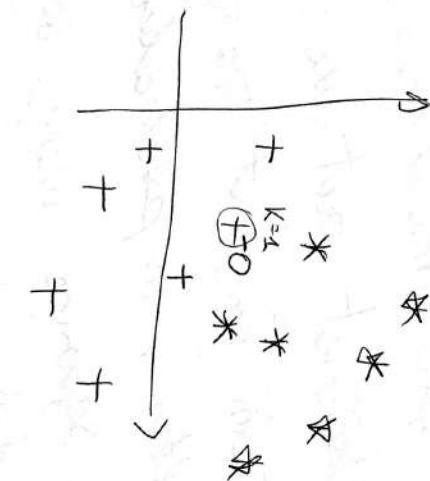
$$c(x, D, k) = \frac{1}{k} \sum_{N_k(x, \Delta)} T_{T_i=c}$$

$N_k(x, \Delta)$  are  $k$  nearest neighbors of  $x$

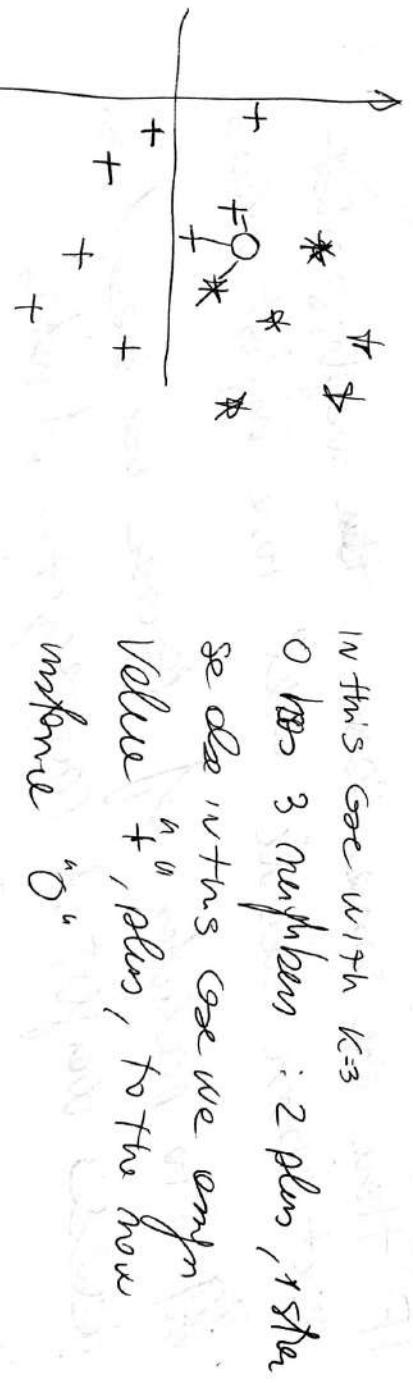
$$T(c) = \begin{cases} 1 & \text{if } c \text{ is true} \\ 0 & \text{otherwise} \end{cases}$$

2)  $\boxed{k=1}$

With  $k=1$  we can see that "+" is nearest neighbor of "0" so in this case the value of this new instance should be plus (+)

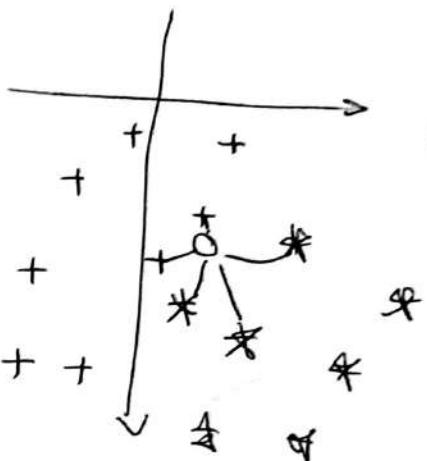


$\boxed{k=3}$



In this case with  $k=3$   
 0 has 3 neighbors : 2 plus, 1 star  
 so also in this case we assign  
 Value "+" plus, to the new  
 instance "0"

$\blacksquare$   $k = 5$



In this case neighbors are  
3 stars and 2 plus  
stars are major than plus  
so the new instance "no"  
will have value of stars(\*)

- 1) Boosting is based on concept that we want  
to use not a single complex model / learner to  
train a particular dataset but we want  
more learners that train on that training dataset  
these models thus a sequence weights so each  
learner train on weights obtained from performance  
of previous learner

If there are some error / weights that  
next learner decide will have highest values.

The predictor of learner are based on  
highest weighted majority of votes

2) We can use a weak learner, see classifier,

$f(x)$ , we then 1) minimize over function  
 then what we do is to evaluate it, so make  
 prediction and after this we update the slope  
 weight coefficient to obtain the output  
 of the first layer

$$Y_M(x) = \sum_{m=1}^M d_m Y_m(x)$$

$$d_m = \ln \left( \frac{1 - \epsilon_m}{\epsilon_m} \right)$$

$\epsilon_m$  is the evaluation

$$\epsilon_m = \frac{\sum_{n=1}^N w_n^{(m)} (f_m(x_n) \neq t_n)}{\sum_{n=1}^N w_n^{(m)}}$$

I forgot to say that first we initialize

$w^{(1)} = \frac{1}{N}$  and for  $m = 1, \dots, M$  we then the  
 weak layers

12/04/2019

[ext]

- 1) Given Hypothesis space  $H$ , given  $h \in H$  and with error on training set  $\text{error}(h)$ , where entire distribution  $D$  is  $\text{error}(h)$ . Consider to have another  $h' \in H$  in which we have errors( $h$ ) < errors( $h'$ )  $\Rightarrow \text{error}(h) > \text{error}(h')$ . This means that  $h$  overfit training data.

2)

We can have problem of overfitting in decision tree in particular when we have noise data in our dataset. ~~• Furthermore, also when training data cannot represent our when you create a tree that is too deep and too specific for a particular training set because also if you obtain an accuracy that is better than another decision tree that is more generic you risk to have a huge overfit.~~ We can avoid overfit using two different methods: we can stop tree or we build tree and then post-prune. Estimate to stop the tree is difficult because usually you don't know exactly when you have to stop. Post-prune is best idea.

• THERE ARE TWO TYPE OF POST-PRUNE: REDUCED-ERROR PRUNE AND ROCES POST-PRUNING

WITH REDUCED-ERROR PRUNING WE SPILT DATA IN TRAINING TEST AND VALIDATION TEST.

UNTIL PRUNE IS HARMFUL WE EVACUATE IMPACT ON VALIDATION SET OF PRUNING EACH POSSIBLE NODE AND THEN WE SELECT NODE ~~TO~~ TO REMOVE THAT IS THE MOST IMPROVES

#### VALIDATION SET ACCURACY

THIS REDUCED-ERROR PRUNING IS GOOD ONLY FOR SUBTREE OF OUR DECISION TREE AND FURTHERMORE IF WE HAVE A SMALL TRAINING SET IT CAN GIVE ~~US~~ BAD RESULTS.

ROCES-POST PRUNING IS A VARIANT OF REDUCED POST-PRUNING AND WE ~~DO~~ INFER THE DECISION TREE ACCORDING FOR OVERFITTING AND THEN WE CONVERT OUR TREE IN TO AN EQUIVALENT SET OF RULES, THIS BECAUSE ~~WE~~ EACH PATH OR THE TREE PRODUCE DIFFERENT RULE AND CONVERTING YOU CAN SEE THIS DISTINCTION, SO ~~WE~~ YOU CAN PRUNE IN ~~A~~ DIFFERENT WAY. FINALLY SORT FINAL RULES IN TO A SEQUENCE THAT YOU WANT TO USE.

Ex2

NAIVE BAYES CLASSIFIER IS A CLASSIFIER TO FIND MOST PROBABLE CLASSIFICATION FOR A NEW INSTANCE THAT IS NOT IN THE DATA SET, IT USE ~~SMALL~~ VACUES OF ATTRIBUTE, IN PARTICULAR INSTANCE IS DESCRIBED ~~WITH~~ WITH VACUES OF

ribute.

Given  $D = \{(x_n, t_n)\}_{n=1}^N$  dataset, target function  
 $f: X \rightarrow V, V = v_1, \dots, v_k$

$$P(V_j | x, D) = P(V_j | \alpha_1, \dots, \alpha_n, D)$$

We need to maximize this to compute naive bayes class  
of new instance  $x$

$$V_j = \operatorname{argmax}_{V_j \in V} P(V_j | D) \prod_i P(\alpha_i | V_j, D)$$

Naive Bayes respect to Bayes optimality but it do approximation don't compute "optimal value" but it do approximation and it doesn't need hypothesis and prior probability to find best probability classification for new instance, and this mean that Naive Bayes don't have problem of size or of course we lost optimality but approximation of Naive Bayes classifier are good enough.

2)

3 classes, so I need to use Naive Bayes with

multinomial

with probability  $P(C_i)$ , probability of that

particular class between  $\mu_1, \mu_2, \mu_3$

We must estimate  $P(C_j)$  and  $P(x_i | C_j)$

where ~~prob~~  $x_i$  is word  $i$  in sentence and  $P(x_i | C_j)$  is the probability to find word  $x_i$  given a class  $C_j$

$$\boxed{V_{NB}} = \arg \max_{C_j} P(C_j) \cdot \prod_{i=1}^{\text{len}(d)} P(x_i | C_j, D)$$

$d$  is document that we need to analyze

[ex3]

1) Given a dataset

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

$$D = \{(Q_1, S_1), P_1, \dots, (Q_n, S_n), P_n\}$$

$$P(+|w) = \prod_{i=1}^N y_i^{P_i} (1 - y_i)^{1-P_i}$$

with

$$y_i = \sigma(w^T x_i) = P(C_j | x_i)$$

posterior probability  
of class  $C_j$

We use cross-entropy function

$$E(w) = -\ln P(+|w) = -\sum_{i=1}^N P_i \ln y_i + (1 - P_i) \ln (1 - y_i)$$

$$w^* = \arg \min_w E(w)$$

In this way we can estimate if student pass or not exams, where  $w$  are weights.

We know iterative option to minimize error: Newton-Raphson Gradient error respect to  $w$

$$\nabla E = \sum_{i=1}^N (y_i - P_i) x_i$$

gradient descent step

$$w^{new} \leftarrow w^{old} - H^{-1} \nabla E$$

$H^{-1}$  is hessian matrix of  $E(w)$

$$H^{-1} = \nabla \nabla E = \sum_{i=1}^n y_i(1-y_i)x_i x_i^T = X^T R X$$

$$R = y_i(1-y_i)$$

$$\nabla E(w) = X(y - x) \quad \Rightarrow \quad w \leftarrow w - (X^T R Y)^{-1} X^T (y - x)$$

2) Parameter  $\theta$  that you have to learn on weight, so  $w$  and you update it using iterative method IRLS until there are termination condition

3) Error Function is cross entropy

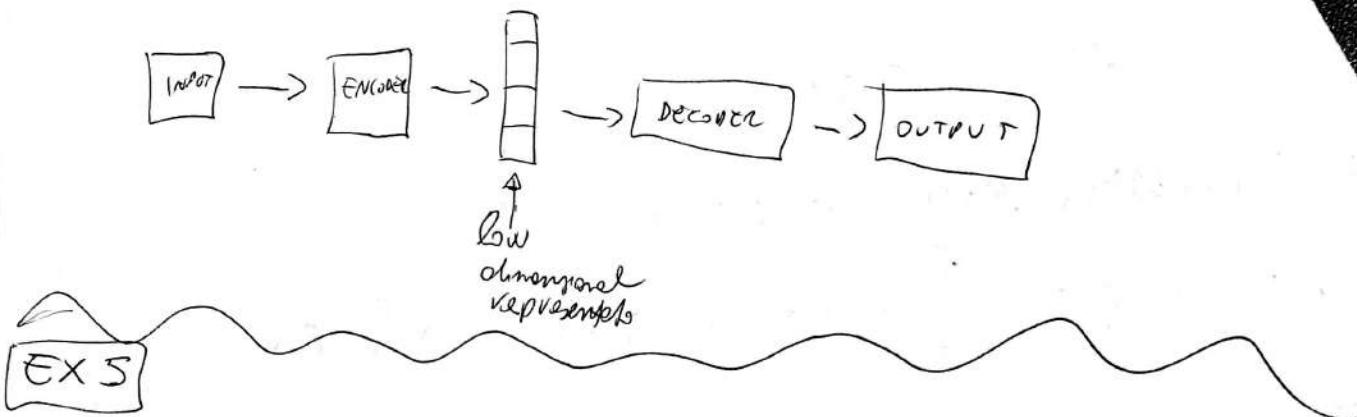
$$E(w) = -\ln p(t|w) = \sum_{i=1}^n p_i \ln y_i + (1-p_i) \ln (1-y_i)$$

[Ex 4]

1) Autoencoder are combination of encoder and decoder, take images in input and return a vector such as CNN, decoder takes input a vector and return images as it is composed by deconvolutional layers.

Autoencoder Trained is base on reconstruction loss and provides low dimensional representation.

2)



1)

STATES ARE FULLY OBSERVABLE BECAUSE IN PARTICULAR IN  
THE PRESENCE OF NON DETERMINISTIC OR STOCHASTIC ACTIONS,  
THE STATE RESULTING FROM EXECUTION OF ACTION ~~IS~~ IS  
NOT KNOWN BEFORE ~~EXECUTION~~ EXECUTION OF THE ACTION,  
BUT AFTER IS FULLY OBSERVABLE.

2)

IN MDP WE HAVE STATES THAT ARE FULLY OBSERVABLE SO  
WE DON'T NEED OBSERVATIONS & INDEED IN ~~POMDP~~ POMDP  
WE HAVE ~~STATE~~ COMBINATION OF NOT FULLY OBSERVABLE AND  
FULLY OBSERVABLE AND FURTHER  
MAY THEY HAVE DIFFERENCE OF STRUCTURE OF MODEL.

IN MDP WE HAVE

$$\langle X, A, \delta, r \rangle$$

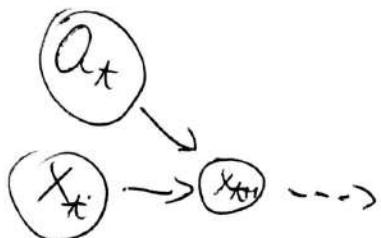
Where  $X$  IS SET OF STATES,  $A$  IS SET OF ACTIONS,  
 $\delta$  IS TRANSITION FUNCTION AND  $r$  IS REWARD FUNCTION  
IN POMDP WE HAVE

$$\langle X, A, z, \delta, r, o \rangle$$

$X$  IS SET OF STATES,  $A$  IS SET OF ACTIONS,  $z$  IS SET

OBSERVATION,  $\delta$  TRANSITION FUNCTION,  $V$  IS REWARD FUNCTION,  
 $\pi$  IS PROBABILITY DISTRIBUTION OVER OBSERVATION.

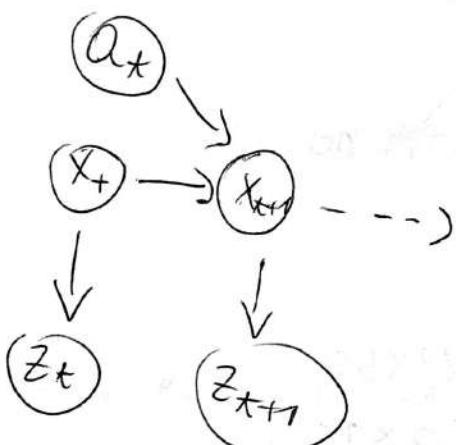
3) MDP



AS WE CAN SEE

$x_{t+1}$  DEPEND OF ONLY  
OF CURRENT STATE  
AND WE DON'T HAVE ANYTHING  
THAT ACTIONS AND STATES

POMDP



HERE WE HAVE THAT  
 $x_{t+1}$  DEPEND OF COURSE OF  
CURRENT STATE BUT ALSO  
ACTIONS AND OBSERVATION  
BECAUSE WE HAVE A  
COMBINATION OF MDP

HMM PROCESS IN WHICH  
ONE IS FULLY OBSERVABLE AND  
THE OTHER ONE NO.

EX6

two layer ANN

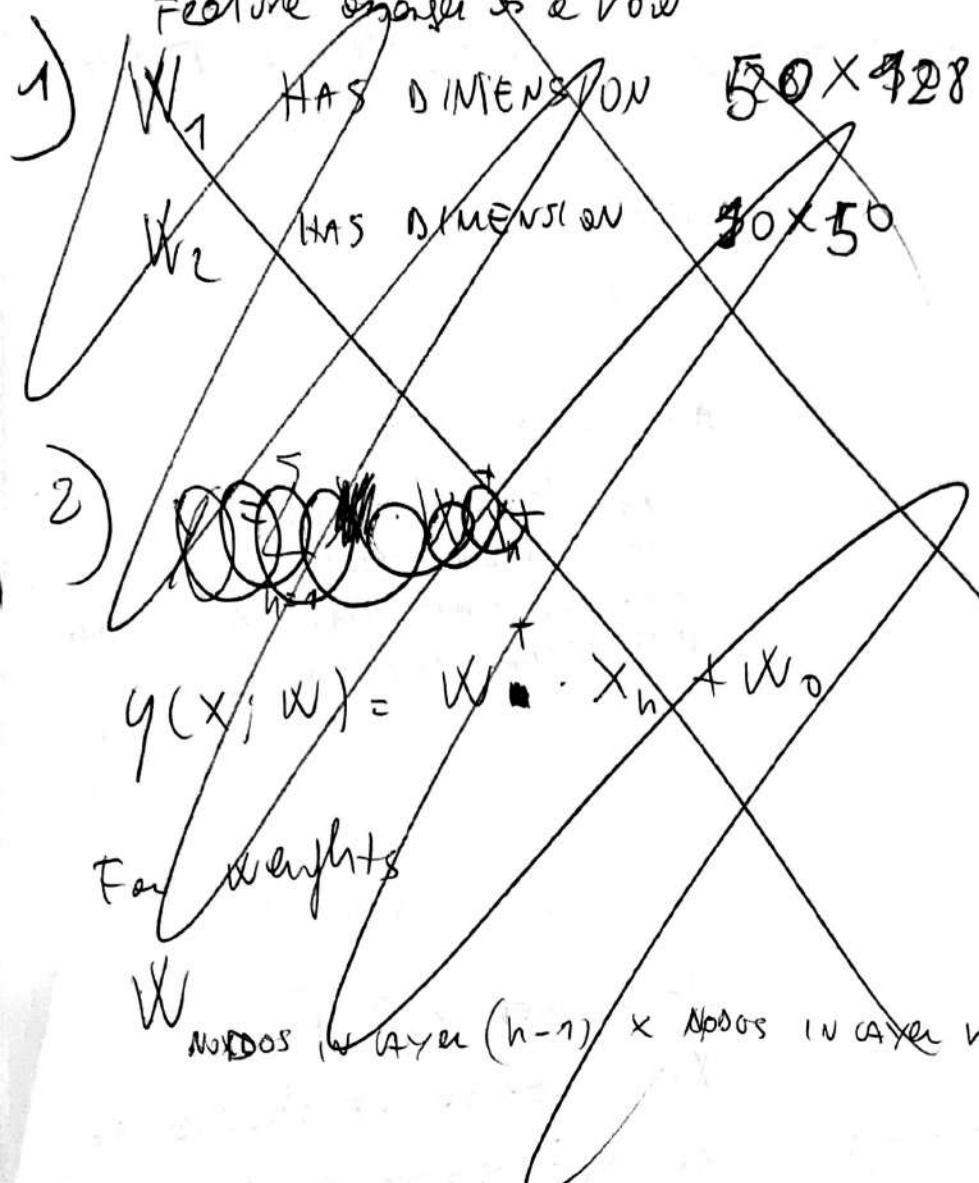
INPUT 128  
OUTPUT 10  
HIDDEN 50

④

WE DO NOT USE BIAS

WE DO NOT USE

Feature vector is a row



1) DIMENSION OF  $W_1 = 128 \times 50$

DIMENSION OF  $W_2 = 50 \times 10$

2) ~~Weights~~ we have seen the calculate function of hidden layer and linear for output

Forward pass dimension is ~~= NODES IN LAYER (N-1)  $\times$  NODES IN LAYER N~~

$$h_1 = g(W_1^T \cdot x + w_{01})$$

output for first layer

where  $g(d) = \max(0, d)$

$$y = W_2^T \cdot h_1 + w_{02}$$

output of out layer using a linear function

EXAM TEST 18/09/2019

EXERCISE 1

1)  $f: X \rightarrow Y$

is target function which

$X = F \times NR \times NK$

consists by attributes of furniture  $F$ , Nr room or  $NR$ ,  
New Kitchen or  $NK$

$Y = \{0, 1\} = \{Yes, No\}$

$F = \{Yes, No\}$

$NR = \{3, 4\}$

$NK = \{Yes, No\}$

No consideration NOT ACCEPTABLE

Yes acceptable

$D = \{(x_n, t_n, y_{n=1}^s, y)\}$  dataset

$x_n$  input of dataset based on attribute  $F, NR, NK$   
of that specific n house and  $t_n$  is target value,  
so ~~0 OR 1~~ 0 OR 1 associated to that  $x_n$

2)

You can choose attribute using information gain, higher is this information but it is so  
the attribute with the highest information gain is  
done during each step of build a decision

tree.

Furthermore, Information gain is how good that particular attribute separates well the training sets with respect to a target function

$$\text{Gain}(S, A) = \text{Entropy}(S) - \sum_{v \in \text{Values}(A)} \frac{|S_v|}{S} \text{Entropy}(S_v)$$

$$S_v = \{s \in S | A(s) = v\}$$

Entropy is measure of a collection of samples, in particular give a collection  $S$ , containing positive and negative sample of some target function

$$\text{Entropy}(S) = -P_+ \log_2 P_+ - P_- \log_2 P_-$$

$P_+$  = portion of positive samples

$P_-$  = portion of negative samples

$$\text{Values}(F_{out}) = \{Y, N\} \quad |S| = 5$$

$$S = [3+, 2-]$$

this means that in 3 times is acceptable so 3+ and 2 times is not acceptable so 2-

$$S_{F_{NO}} = [2+, 2-] \quad S_{F_{YES}} = [1+, 1-]$$

$$\text{Values}(NR) = \{3, 4\}$$

$$S_{NR3} = [1+, 2-] \quad S_{NR4} = [2, 0-]$$

$\{K\} = \{\text{Yes}, \text{No}\}$

$$= [1+, 0-]$$

$$S_{NK\text{No}} = [2+, 2-]$$

$$(S) = -\frac{3}{5} \log_2 \left(\frac{3}{5}\right) - \frac{2}{5} \log_2 \left(\frac{2}{5}\right) = 0.9710$$

$$(S, F) = 0.9710 - \frac{3}{5} \text{Entropy}(S_{F\text{No}}) - \frac{2}{5} \text{Entropy}(S_{F\text{Yes}})$$

$$= 0.9710 - 0.5509 - 0.4 = 0.0201$$

$$(S, NR) = +0.9710 - \frac{3}{5} \text{Entropy}(S_{NR_3}) - \frac{2}{5} \text{Entropy}(S_{NR_4})$$

$$-\frac{3}{5} \text{Entropy}(S_{NR_3}) - \frac{2}{5} \text{Entropy}(S_{NR_4})$$

$$= 0.9710 - 0.5460 = 0.425$$

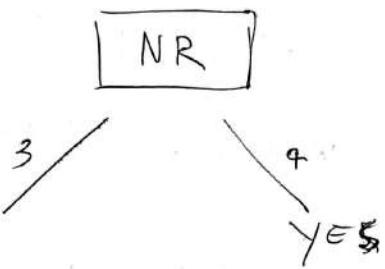
$$\text{gain}(S, NK) = 0.9710 - \frac{1}{5} \text{Entropy}(S_{N\text{Yes}}) - \frac{4}{5} \text{Entropy}(S_{N\text{No}})$$

$$= 0.9710 - 0 - 0.8 = 0.171.$$

3) 113

Example on not all positive or all negative and attribute's list is not empty so we need to find attribute to assign to the root of tree, so to choose attribute we want to see the ~~gains~~ information gain and choose attribute with

highest information gain, in this case it is NR  
so at first we have



Now when NR is equal to 3 so when  $A(S) = 3$   
we consider gain knowing 3  
 $|S_3| = 3$

$$\text{Gain}(S_3, F) = \text{Entropy}(S_3) - \frac{1}{3} \text{Entropy}(S_{3F_Y}) - \frac{2}{3} \text{Entropy}(S_{3F_N})$$

$$\text{Entropy}(S_3) = -\frac{1}{3} \log_2\left(\frac{1}{3}\right) - \frac{2}{3} \log_2\left(\frac{2}{3}\right) = 0.9183$$

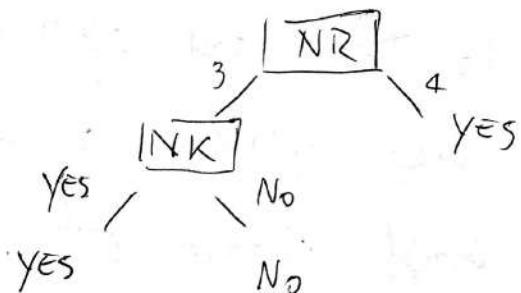
$$S_{3F_Y} = [0, 1+]$$

$$S_{3F_N} = [2, -]$$

$$\Rightarrow \text{Gain}(S_3, F) = 0.9183 - 0 - 0 = 0.9183$$

$$\begin{aligned} \text{Gain}(S_3, NK) &= 0.9183 - \frac{1}{3} \text{Entropy}(S_{3NK_Y}) - \frac{2}{3} \text{Entropy}(S_{3NK_N}) \\ &= 0.9183 - 0 - 0 = 0.9183 \end{aligned}$$

Doesn't matter what you choose because they have highest gain in both cases



1) We have a Dataset  $D$  and hypothesis space  $H$   
 we want to compute <sup>probability</sup> distribution over  $H$  given  $D$

$P(H|D)$   $\rightarrow$  with Bayes

$$P(h|D) = \frac{P(D|h) \cdot P(h)}{P(D)} \quad \text{where } h \in H$$

Where  $P(h)$  ~~posterior~~ prior probability of  $h$

$P(D)$  prior probability of  $D$

$P(D|h)$  is a ~~posterior~~ probability of  $D$  given  $h$ , posterior probability of  $h$

We want to find most probable hypothesis

and also we consider that  $P(h_i) = P(h_j)$  with  $i \neq j$

$$\Rightarrow \boxed{h_{MC} = \arg \max_{h \in H} P(h|D)}$$

$$= \arg \max_{h \in H} \frac{P(D|h) \cdot P(h)}{P(D)}$$

$$= \arg \max_{h \in H} P(D|h)$$

This because  $P(D)$  and  $P(h)$  are constant so we can not consider



2) No this statement is false because MAXIMUM LIKELIHOOD is NOT good method for ~~classification~~ as most probable classification of new instance this because it

For example we have

Given new instance  $x$  we have

$$h_1(x) = \Theta \quad P(D|h_1) = 0.4$$

$$h_2(x) = \Theta \quad P(D|h_2) = 0.03$$

$$h_3(x) = \Theta \quad P(D|h_3) = 0.3$$

If we use MAXIMUM LIKELIHOOD we choose that the most probable ~~classification~~ for  $x$  is  $\Theta$  because we do ~~approx~~ of  $P(D|h_i)$   $i=1,3$  and we obtain that is 0.4 so

0.4 compared to  $h_2$  that is positive but this is not true, indeed  $x$  should be classified  $\Theta$  thus

$$\sum_{h_i \in H} P(\Theta|x, h_i) P(h_i|D) = 0.6 \text{ and } \cancel{\text{this}}$$

$$\sum_{h_i \in H} P(\Theta|x, h_i) P(h_i|D) = 0.2$$

so the most probable classification for  $x$  is  $\Theta$  and we can do using Bayes optimal classifier in which we can find optimal classifier for new instance if we have

$$P(v_j|x, D) = \sum_{h_i \in H} P(v_j|x, h_i) P(h_i|D)$$

where probability of  $h_i = v_j$  is independent from  $D$  given  $h_i$  and ~~also~~ also doesn't depend of  $x \in D \Rightarrow P(h_i|x, D) = P(h_i|D)$

3

In linear model we have target function  $f: X \rightarrow Y$

where  $X \subseteq \mathbb{R}^d$  and  $Y = \{c_1, \dots, c_n\}$

We assume that our dataset  $D$  is linearly separable, this means that we can divide instance space in two regions in which differently classified instances are separated.

We consider for example two classes

$$y: X \rightarrow \{c_1, \dots, c_n\}$$

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

where  $w_0$  is bias and  $w$  is weight vector

We have Dataset  $D = \{(x_n, t_n)\}_{n=1}^N$  and least square method we need to find

$$y(x) = \tilde{w}^T \tilde{x}$$

where

$$\tilde{w} = (\tilde{w}_1, \dots, \tilde{w}_k)$$

$$\tilde{w}_k = \begin{pmatrix} w_0 \\ w_k \end{pmatrix}$$

$$\tilde{x} = \begin{pmatrix} 1 \\ x \end{pmatrix}$$

So our aim is to minimize sum-of-squares error

$$E(\tilde{w}) = \frac{1}{2} \text{Tr} \{ (\tilde{x} \tilde{w} - T)^T (\tilde{x} \tilde{w} - T) \}$$

we obtain

$$\tilde{w} = (\tilde{x}^T \tilde{x})^{-1} \tilde{x}^T \cdot T \Rightarrow y(x) = T^T ((\tilde{x}^T \tilde{x})^{-1} \tilde{x}^T)$$

where this is a closed-form solution

This method has some problem with outliers, in particular when there are outliers the decision boundary ~~cannot~~ could be in a position to divide perfectly two classes in different region. This because sum-of-square error is a sum so also outlier give a contribution and when we want to minimize it to find a  $y(x)$ , so then final decision boundary, we obtain wrong division of our instance space and some instances are on decision boundary and you can't classify or in particular we ~~are~~ misclassify.

Ex 4

$$1) \hat{A} = [\mathbf{X} \mathbf{X}^T + \lambda I]^{-1} +$$

$K = \mathbf{X} \mathbf{X}^T$  is a gram matrix

$$K = K(\mathbf{x}_m, \mathbf{x}_n) = \begin{bmatrix} K(\mathbf{x}_1, \mathbf{x}_1) & \dots & K(\mathbf{x}_1, \mathbf{x}_N) \\ \vdots & & \vdots \\ K(\mathbf{x}_N, \mathbf{x}_1) & \dots & K(\mathbf{x}_N, \mathbf{x}_N) \end{bmatrix}$$

(is a  $N \times N$  matrix)

2) linear model for regression

$y = w^T x$  and  $D = \{(x_n, t_n)\}_{n=1}^N$  we have that we want to minimize loss function

$$J(w) = \sum_{n=1}^N E(y_n, t_n) + \lambda \|w\|^2 \quad \text{where } y_n = w^T x_n$$

$$\text{So we consider } E(y_n, t_n) = (y_n - t_n)^2$$

$$\tilde{w} = (X^T X + \lambda I)^{-1} X^T x = X^T d$$

$$\Rightarrow y(x; w) = \sum_{n=1}^N \hat{d}_n \cancel{x_n} x_n^T x$$

We ~~can~~ apply Kernel trick, that consist to substitute inner product  $x_i^T x$  with Kernel function

$$\Rightarrow y(x; w) = \sum_{n=1}^N \hat{d}_n K(x_n, x)$$



**Ex 5**

- 1) We can perform regression based of some basis functions for example we know polynomial, Radial, Sigmoid function and we want to find one of them that is useful for this problem.

In this case we can use polynomial basis function of degree 3

$$y(x) = w_1 x + w_2 x^2 + w_3 x^3 + w_4 = x^T w$$

$$w = \begin{pmatrix} w_0 \\ w_1 \\ w_2 \\ w_3 \end{pmatrix}$$

- 2) for linear regression we can avoid overfitting using regularization. In particular in linear regression we want to maximize likelihood, so minimize the error

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

Regularization technique we have

$$\min E_D(w) + \lambda E_R(w)$$

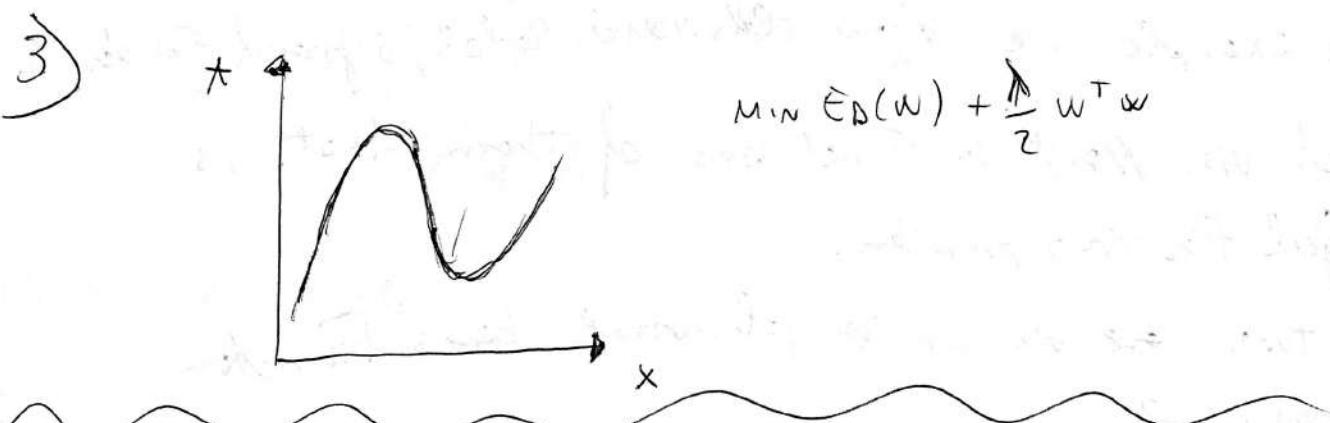
where  $\lambda$  is regularization factor  $\lambda > 0$

we can choose two different choice for  $E_R(w)$  that is regularization term

$$① E_R(w) = \frac{1}{2} w^T w$$

or

$$② E_R(w) = \sum_{j=0}^{M-1} |w_j|^q$$



[Ex6]

- 1) Classifier of k nearest neighbors is an algorithm used for non-parametric model, in particular we choose ~~k nearest neighbor~~ mode classification of a new instance  $x$  and so we have

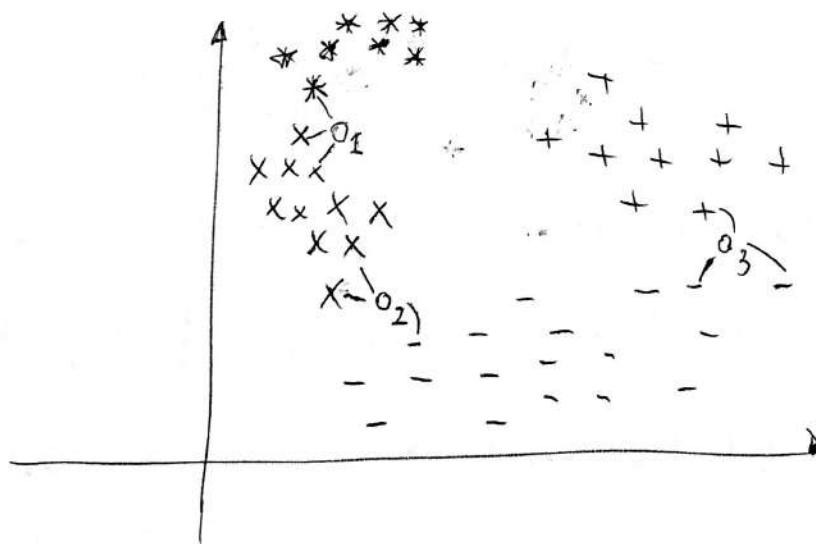
$\mathcal{D} = \{(x_n, t_n)\}_{n=1}^{\infty}$ , and  $f: X \rightarrow C$  Target Function.

First we choose  $k$  nearest neighbors of new instance  $x$ , then we assign to  $x$  the most common value / label ~~avg~~ among majority of its  $k$  nearest neighbors chosen.

2)

$$(*, x, +, -) \quad K = 3$$

Assume that our instances are "o" that we want to classify, in particular  $(o_1, o_2, o_3)$



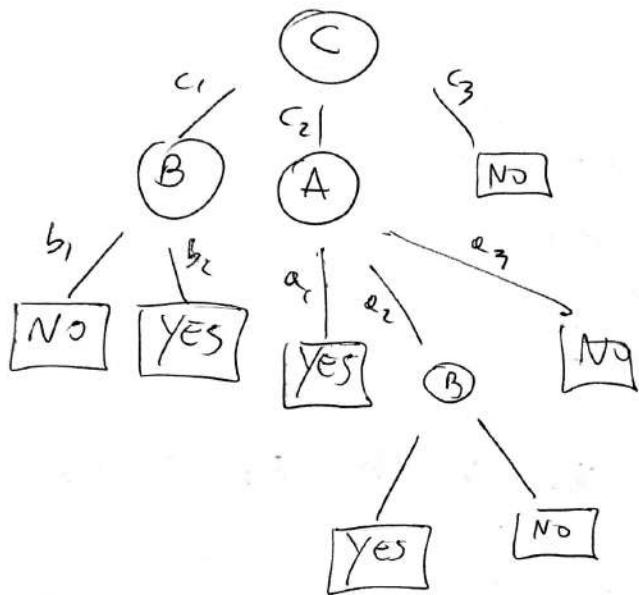
$o_1$  is classify as "x"

$o_2$  is classify as "x"

$o_3$  is classify as "-"

18/01/2019

Ex1



1)

IF  $C = c_3$  THEN No

IF  $(C = c_1) \wedge (B = b_2)$  THEN Yes

IF  $(C = c_1) \wedge (B = b_1)$  THEN No

IF  $(C = c_2) \wedge (A = a_1)$  THEN Yes

IF  $(C = c_2) \wedge (A = a_2) \wedge (B = b_1)$  THEN Yes

IF  $(C = c_2) \wedge (A = a_2) \wedge (B = b_2)$  THEN No

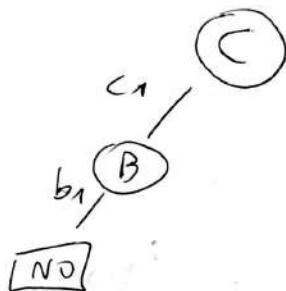
IF  $(C = c_2) \wedge (A = a_3)$  THEN No

2)  $S = \{S_1 = \langle \varrho_1, b_1, c_1, \text{No} \rangle, S_2 = \langle \varrho_2, b_1, c_2, \text{Yes} \rangle,$   
 $S_3 = \langle \varrho_1, b_2, c_3, \text{Yes} \rangle, S_4 = \langle \varrho_2, b_2, c_2, \text{Yes} \rangle\}$

We consider  $s_1 = \langle a_1, b_1, c_1, \text{No} \rangle$

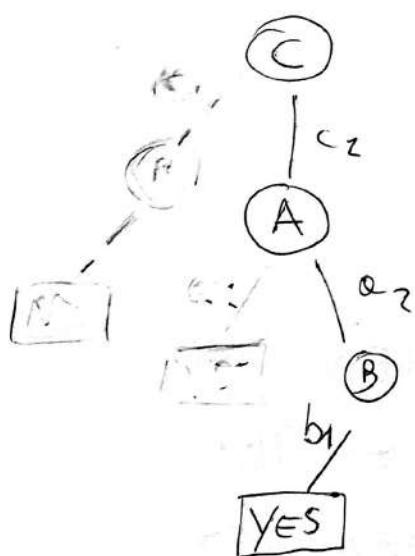
$(C=c_1) \wedge (B=b_1) \wedge (A=A_1)$  ~~is~~

We can see that



give us No value  $\infty$  with  $s_1 + t$  is consistent

For  $s_2 = \langle a_2, b_1, c_2, \text{Yes} \rangle$



IT'S CONSISTENT WITH THIS SET  $s_2$

For  $s_3 = \langle a_1, b_2, c_3, \text{Yes} \rangle$



IT IS NOT CONSISTENT because in  $s_3$  we have Yes as output but in our T we have No value at the end

$s_4 = \langle a_2, b_2, c_2, \text{Yes} \rangle$



Also in this case we have some problem because in set  $s_4$  it said to have output Yes but in our Tree we have No

=D We can conclude that T is consistent for  $s_1$  and  $s_2$  and Not consistent for  $s_3$  and  $s_4$

Ex 2

① We have data set  $D = \{(x_n, t_n)\}_{n=1}^N$ ,  $x_n$  is input and  $t_n$  value of target function associated to that particular  $x_n$ . We have hypothesis  $h \in H$  (Hypothesis space, see for bayes)

$$P(h|D) = \frac{P(D|h) \cdot P(h)}{P(D)}$$

Maximum a posteriori create most probable hypothesis so we need to maximize  $p(h|D)$

$$h_{MAP} = \arg \max_{h \in H} p(h|D) = \arg \max_{h \in H} p(D|h) \cdot P(h)$$

We drop  $P(D)$  because it is a constant and does not depend on  $h$ .

If we do consider  $P(h_i) = P(h_j)$  with  $i \neq j$

We compute MAXIMUM LIKELIHOOD

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

We drop  $P(h)$  because it is a constant for the particular hypothesis

2)

Bayes Optimal classifier is a classifier that finds most probable classification for a new instance

In particular, given a dataset  $D = \{(x_n, t_n)\}_{n=1}^N$ ,  
• target function  $f: X \rightarrow V$ ,  $V = \{v_1, \dots, v_k\}$  and  
 $x \notin D$ . We have that probability that  $x$  is classified for a particular  $v_j$  is

$$P(v_j|x, D) = \sum_{h_i \in H} P(v_j|x, h_i) P(h_i|D)$$

where

$P(v_j|x, h_i)$  means that probability that  $h_i = v_j$  is independent from  $D$  so  $P(v_j|x, h_i; D) = P(v_j|x, h_i)$

else  $h_i$  doesn't depend on  $x \notin D$

$$P(h_i | X, D) = P(h_i | D)$$

This means that Bayes Optimal classifier is that

$$V_{OB} = \operatorname{arg\max}_{h_i \in V} \sum p(v_j | x, h_i) P(h_i | D)$$

We maximize probability that new ~~new~~ instance is  
classified correctly

Bayes optimal classifier finds global value because  
no other classifier uses some hypothesis space and some  
prior knowledge

3) Bayes Optimal classifier has problem that for  
Hypothesis space that are very large is very difficult  
to know prior probability of all hypotheses  
belong to that hypothesis space, so is powerful for  
Hypothesis spaces that are small enough.

EX3

1) Goal of Linear Regression is to find a relation  
ship in which predictor  $y$  and input  $x$  are well  
defined, or a kind of curve that can represent

all point of class in a good way.  
Define a model  $y(x; w)$  with parameter  $w$  to approximate target function

$$y(x; w) = w_0 + w_1 x_1 + \dots + w_d x_d = w^T x$$

with

$$x = \begin{bmatrix} 1 \\ x_1 \\ \vdots \\ x_d \end{bmatrix}$$

$$w = \begin{bmatrix} w_0 \\ w_1 \\ \vdots \\ w_d \end{bmatrix}$$

2)

$$\mathcal{D} = \{(x_1, t_1), \dots, (x_n, t_n)\}$$

To find parameter  $w$

We need to maximize likelihood maximum likelihood but  
in particular we use feature space to represent linear regression

$$y(x; w) = \sum_{j=0}^{M-1} w_j \cdot \phi_j(x) = w^T \phi(x)$$

$$w = \begin{bmatrix} w_0 \\ \vdots \\ w_{M-1} \end{bmatrix}$$

$$\phi(x) = \begin{bmatrix} \phi_0(x) \\ \vdots \\ \phi_{M-1}(x) \end{bmatrix}$$

$$\phi_0(x) = 1$$

the target value  $t$  is given by  $y(x; w)$   
affected by noise  $\epsilon$  so

$$t = y(x; w) + \epsilon$$

common Gaussian noise

$$N(\epsilon | 0, \beta^{-1}) = P(\epsilon | \beta)$$

$\beta$  is precision

$$P(+ | x_i, w, \beta) = N(+ | \phi(x_i; w), \beta^{-1})$$

Likelihood function is

$$P(+_1, \dots, +_n) | (x_1, \dots, x_n), w, \beta = \prod_{i=1}^n N(+_i | w^\top \phi(x_i), \beta^{-1})$$

Maximize likelihood means minimize error so we search to minimize likelihood, we compute it at first

$$\ln (P(+_1, \dots, +_n) | (x_1, \dots, x_n), w, \beta)$$

$$\Rightarrow \text{We obtain } -\beta E_D - \frac{N}{2} \ln (e + \beta^{-1})$$

$$E_D = \frac{1}{2} \sum_{i=1}^N (+ - \phi^\top w)^2$$

We need to minimize error

$$\min E_D = \min \left( \frac{1}{2} \sum_{i=1}^N (+ - w^\top \phi(x_i))^2 \right)$$

$$= \min \left( \frac{1}{2} \sum_{i=1}^N (+ - \phi^\top w)^\top (+ - \phi^\top w) \right)$$

optimal condition is when  $\boxed{V_{D0} = 0 \Leftrightarrow \phi^\top \phi w = \phi^\top +}$

where  $\left[ \begin{array}{c} \phi_1(x_1) \\ \vdots \\ \phi_N(x_N) \end{array} \right]$

$$\phi = \begin{cases} & \\ & \\ & \end{cases}$$

$$W_{\text{MC}} = [(\phi^T \phi)^{-1} \phi^T \cdot t] = \phi^* t$$

for sequential learning - we can say that we use '+' when data set is sufficiently large so we obtain sequential learning using stochastic gradient descent

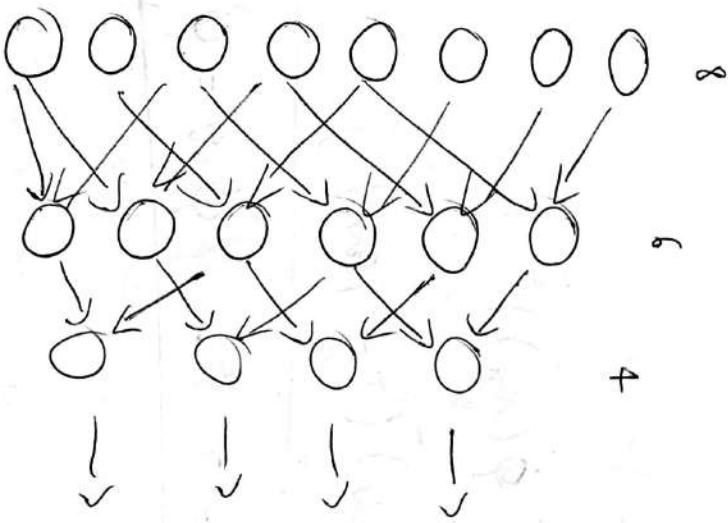
$$\hat{w} \leftarrow \hat{w} + \eta V_{\text{EM}}$$

$\eta$  is learning rate

$$\text{So } \hat{w} \leftarrow \hat{w} - \eta [\hat{x}_n - \hat{w}^T \cdot \phi(x_n)] \phi(x_n)$$

And ~~if~~ algorithm converge for small value of  $\eta$

$$\begin{aligned} \text{No parameter} &= (6 \times 8) + (6 \times 4) + (6 \times 1) \\ &= 82 \end{aligned}$$



①

Backpropagation is an off algorithm that only compute gradient, it is not a learning algorithm, for example a ~~gradient~~ SGD (stochastic gradient descent) with Momentum suffer of problem of local minima because we we compute the gradient and then velocity, given by Momentum, we can change direction and get stuck in local minima

3) No, Backpropagation is not affected by overfitting because is not a model or a learning algorithm, so it doesn't have problem of overfitting, it just compute gradient

(Ex 5)

1) We see 3 classes in figure and you can separate them with a line, curve line because they're one of projected in some separate area so they are not mixed together, this is reason we can separate them in 3 different region of our figure

2) We can choose different type of kernel but we can consider SVM so for classification is useful to choose linear kernel because in

SVM we have  $w = \sum_{n=1}^N d_n x_n$

$$\text{sign}(x_i \cdot d)) = \text{sign} \left( \sum_{n=1}^N d_n x_n^\top x_i + w_0 \right)$$

so we can replace kernel function if

we have inner product  $x^\top x$

$$\Rightarrow y(x_i \cdot d) = \text{sign} \left( \sum_{n=1}^N d_n k(x_n, x) + w_0 \right)$$

3) For SVM for classifier we have

And to find we we ned to use long way

$$L(a) = \sum_{k=1}^K \sum_{n=1}^N (Q_{kn} - \frac{1}{2}) \sum_{m=1}^N \sum_{l=1}^N Q_{kn} Q_{km} t_k t_m K(t_k t_m)$$

So we can compute with multiplies  $Q_{kn}$  of lagrange

$$W_k = \sum_{k=1}^K Q_{kn} t_k x_n$$

and constant support vector we can have

$$w_0 = \frac{1}{4} \sum_{k=1}^K \left( t_k - \sum_{j=1}^K Q_{kj} t_j \right)$$

x6

Perception model is a linear classification model  
 that ~~uses different classes~~ want to classify new  
 instance, we consider to have ~~function~~

$$\delta(x) = \begin{cases} 1 & \text{if } w^T x > 0 \\ 0 & \text{otherwise} \end{cases}$$

$$\delta = w_1 x_1 + \dots + w_d x_d = w^T x$$

Now we consider a data set  $D = \{(x_n, t_n)\}_{n=1}^N$  and  
 we want to minimize error where sum-of-squares  
 error is

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

Minimize means do derivative respect to  $w$  of  $E(w)$

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

After this the training rule of percept is to update  
 weights in such to find a decision boundary that  
 can divide instance space in two ~~different regions~~  
~~where there are classes with different classification values~~  
~~that are separated.~~

Regions where in one should be a class and ~~other~~  
 in the other other class with different classification value

If and only if dataset is linearly separable.  
 So updating weights, that is perceptron when yet find  
 a instance that is misclassified, ~~will~~  
~~the~~ model will converge (IFF Dataset is linearly  
 separable) in a finite number of steps

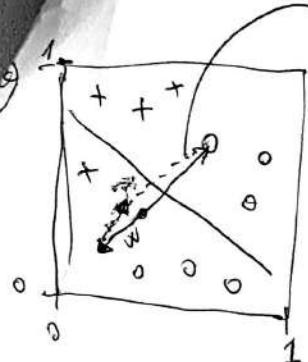
$$w_i^{new} \leftarrow w_i^{old} + \Delta w$$

$$\Delta w = -\eta \nabla E = -\eta \sum_{n=1}^N (t_n - \sigma_n)(x_n) \sum_{i=1}^m (t_n - w^T x_n) (x_{i,m}) \\ = \eta \sum_{n=1}^N (t_n - \text{sign}(w^T x_n)) x_{i,n}$$

$\eta$  is learning rate

Perceptron algorithm has aim to find a vector  
 that can perfectly classify ~~dataset~~ inputs in  
 our set and it guaranteed to find solution  
 in a finite number of step if Dataset is  
 linearly separable ~~but may have this case~~  
 but ~~not~~ often model converge if initial values  
 that we have to choose for  $w$  are good,  
 but you choose them arbitrary so you don't know  
 very well

2)

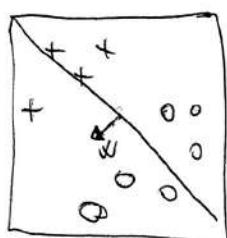


→ We FIND

a INSTANCE THAT IS MISCLASSIFY  
Because ON TOP we choose to classify  
INSTANCE OF CLASS "+" and on bottom  
INSTANCE OF CLASS "o"

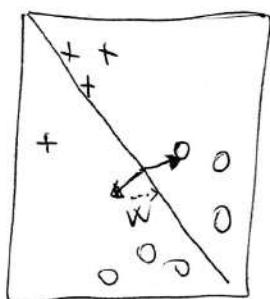
THIS MEANS THAT we need to update  
weight and w change

b)

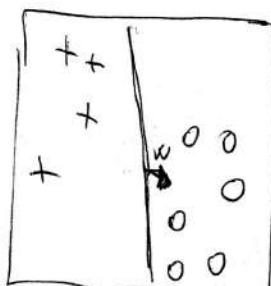


w is smaller than in the previous picture  
thus because arrow that go to MISCLASSIFY point  
we proud to w FIND length of ~~vector~~ vector  
w with update weights

c)



We find another point that is misclassify  
so we need to update weight and w  
change and also decision boundary



at the end we don't reach any  
points that is Misclassified so we  
converge, as because our offset set is  
linearly separable and we don't need  
to update weights anymore

2/02/2019

[Ex-1]

1) Reinforcement learning problem we have a dataset

$$D = \{x_1, Q_{x_1}, V_{x_1}, \dots, x_n, Q_n, V_n\}$$

We want to maximize reward so we need to find optimal policy to obtain this

In Reinforcement learning we have only inputs but the input is a collection of state, action and reward and the output is not available.

2) Q-Learning is an algorithm used to find optimal policy when transition function ( $\pi$ ) and reward function ( $R$ ) are not known, so we need to do an approximate with Q-Learning when training rule of this algorithm is

$$\hat{Q}(x, a) \leftarrow \bar{r} + \gamma \max(\hat{Q}(x', a'))$$

where  $\bar{r}$  is the immediate reward,  $\gamma$  is scaling factor and  $\hat{Q}$  is the approximation function of  $Q$ ,  $\pi^*$  is optimal policy and  $\hat{\pi}^* = \arg \max_{a \in A} \hat{Q}(x, a)$

The algorithm is:

$$Q_{(t+1)}(x, a) \leftarrow 0$$

then we observe current state  $x$

for  $t=1, \dots, T$  (until terminal condition) do choose action  $a$

execute action  $a$

observe resulting state  $x'$

collect immediate reward  $r$

update table entry of  $\hat{Q}(x, a)$

$$\hat{Q}_t(x, a) \leftarrow r + \gamma \max_{a'} \hat{Q}_{t+1}(x', a')$$

$$x \leftarrow x'$$

end for

$$\pi^* = \arg \max_{a \in A} \hat{Q}_T(x, a)$$

We need to execute action and collect reward to get value of  $\hat{Q}(x, a)$ , but we notice how this being about need transition function and reward function to compute then optimal policy

Ex2

CNN is a convolutional neural network, this means that we use convolutional function between ~~one~~ input and kernel. To avoid overfitting in CNN we use parameter sharing, a method ~~that~~ in which all unit of that layer share a subset of parameter. In this way ~~we~~ we don't need to learn more parameter and this also mean that a particular parameter can be used more than one time for a function. Another methods most common is Dropout, in which we eliminate ~~some~~ randomly units of network with a probability  $p$ .

Ex3

In SVM we want to maximize margin to provide better accuracy.

Margin is smallest distance between decision boundary and closest points of + (support vectors)

The main idea is to find  $w^*$  and  $w^*$  maximize this margin. In particular we have a dataset  $D = \{(x_n, t_n)\}_{n=1}^N$  ~~where  $t_n \in \{-1, +1\}$~~

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

$$\text{Margin is } = \min_{n=1, \dots, N} \frac{|y(x_n)|}{\|w\|}$$

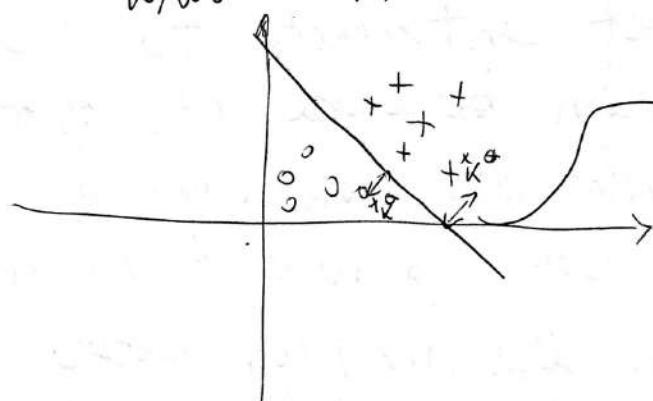
$$\begin{aligned} \text{We know that } & y(x_n) \geq 0 \quad \text{if } t_n = 1 \\ & y(x_n) \leq 0 \quad \text{if } t_n = -1 \end{aligned}$$

$$\Rightarrow t_n y(x_n) \geq 0$$

$$\text{So } |y(x)| = t_n y(x_n)$$

$$\text{Margin } \underset{\|w\|}{=} \min_{n=1, \dots, N} t_n y(x_n) = \frac{1}{\|w\|} \min_{n=1, \dots, N} t_n \cdot (w^T x_n + w_0)$$

$$w^*, w_0^* = \frac{1}{\|w\|} \operatorname{argmax}_{w, w_0} \left( \min_{n=1, \dots, N} (t_n (w^T x_n + w_0)) \right)$$

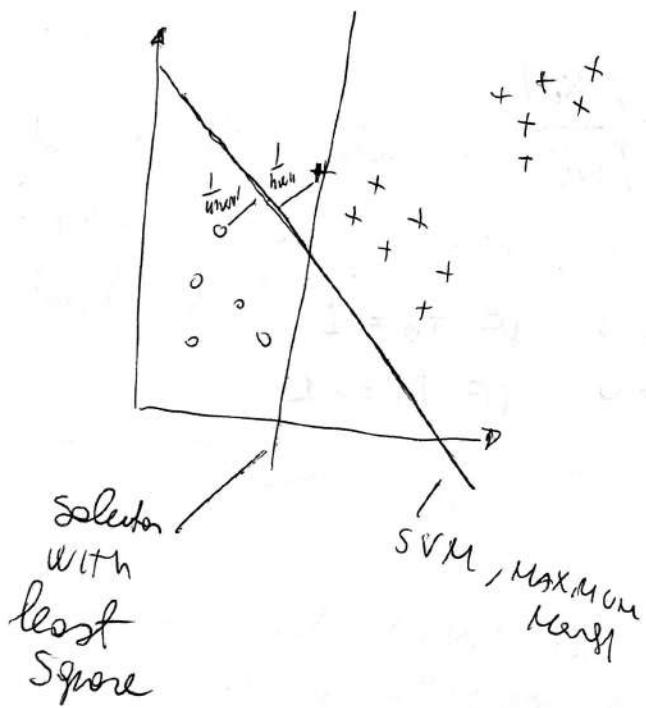


margin that we maximize  
perfect we notice that we have 2 closest point that are at some distance

in where one  $x_k^*$   $w^T x_k^* + w_0 = 1$   
 (point  
in the  
top.)

and other other point  $x_k^*$  when  $w^T x_k^* + w_0 = -1$   
 distance between these two points from Hypoplane (decision  
 boundary) is  $\frac{2}{\|w\|}$

2)



3)

We can notice how Least Square is not good in this problem because it is sensitive to outliers, thus because it's error is a function that sum ~~square~~, so every points in the dataset contribute to the solution but this create problem because it's ~~one~~ decision boundary is not in right position and don't classify better solution, instead SVM is robust to the outliers so it better solution also in this case.

Furthermore Maximum Margin solution is preferred because

is a method that creates a good separation for classification of instances in this way we don't have to worry about classification of instance in wrong way because we have a sufficient margin that can provide to classify better close instance near margin, this not happen in perceptron because we can have decision boundary that has margin that is very small and so there is a risk to classify instance not in good way.

#### Ex 4

1) Confusion matrix is a performance metric used for multiclass problem and it shows how many time instance of class  $C_i$  is classified class  $C_j$ .

In the diagonal we have the accuracy of each class and outside it we have errors

	$C_1$	$C_2$	$C_3$
$C_1$			
$C_2$			
$C_3$			

Confusion matrix of 3 classes

	$C_1$	$C_2$	$C_3$	
$C_1$	60	20	20	40
$C_2$	10	70	20	30
$C_3$	20	15	65	35
FP	30	35	40	105
Absolute Values	30	35	40	195

	$C_1$	$C_2$	$C_3$
$C_1$	0.6	0.2	0.2
$C_2$	0.1	0.7	0.2
$C_3$	0.25	0.15	0.65

Percentage Values

3)

$$\text{Accuracy} = \frac{TP + TN}{TP + TN + FP + FN} = 0.65 \Rightarrow 65\% \text{ accuracy}$$

$$TP = \text{True positive} = 195$$

$$TN = \text{True Negative} = 195$$

$$FP = \text{False Positive} = 105$$

$$FN = \text{False Negative} = 105$$



**Ex 5**

1) GMM

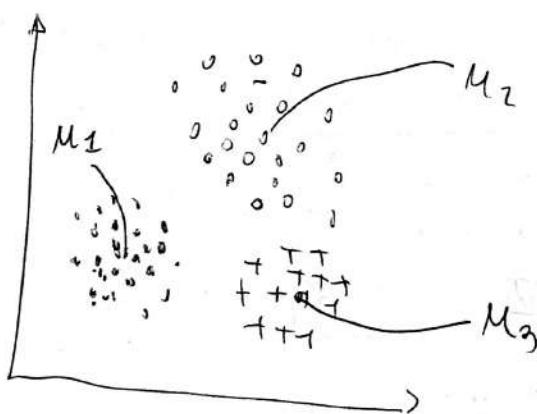
We have dataset  $D = \{x_n\}$ , we are in unsupervised learning and  $K$  is  $K$  different gaussian distribution

$$P(x) = \sum_{k=1}^K \pi_k N(x; \mu_k, \Sigma_k)$$

$\mu_k$  is mean,  $\Sigma_k$  is covariance

$$P(z_k = 1) = \pi_k \quad \text{when } z_k \text{ is a variable } \in \{0, 1\}$$

2)



Matrix covariance size  $2 \times 2$

MEAN has size 2

but we have 3 clusters so, we obtain to have

$(2 \times 2 + 2) \cdot 3 + 3$  because we have 3 cluster

$$(2 \times 2 + 2) \cdot 3 + 3 = 21$$

↓

MATRIX  
COVARIANCE  
+ MEAN of dim  
2

3 P.  
that we  
have to consider ("weights")



Ex 6

KNN is an algorithm for classification of non parametric model, in which we have  $f: X \rightarrow C$  and dataset

$$D = \{(x_n, t_n)\}_{n=1}^N$$

this algorithm has two steps:

- 1) Choose K nearest neighbors to new instance  $x \notin D$
- 2) Assign to  $x$  most common value among the majority of neighbors

likelihood function

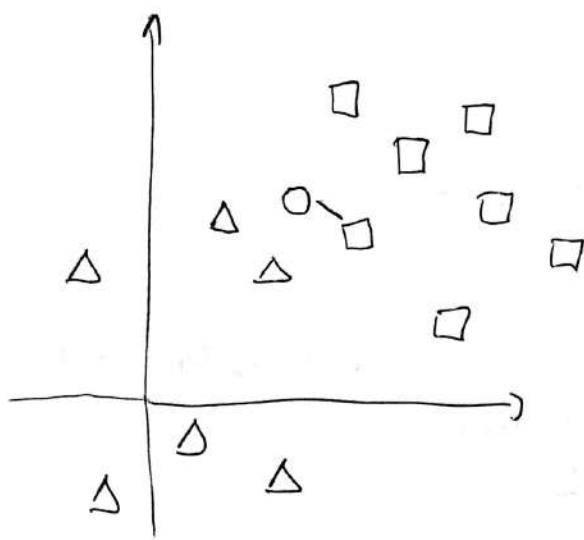
$$P(c|x, K, D) = \sum_{K \in N_K(x, D)} I(t_i = c)$$

where  $N_K(x, D)$  are nearest K neighbor

and

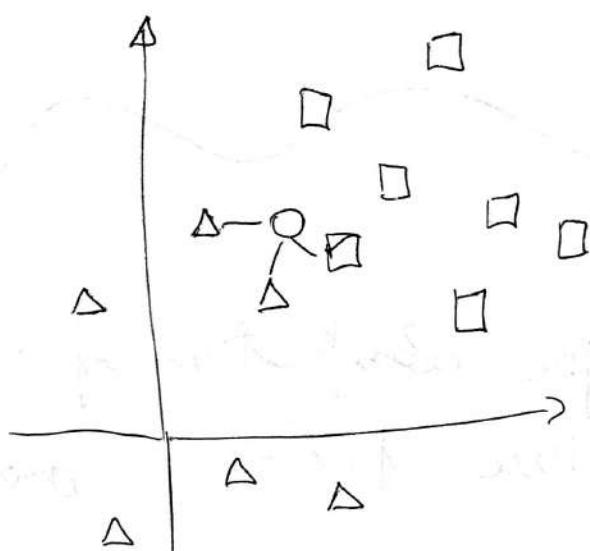
$$I(e) = \begin{cases} 1 & \text{IF } e \text{ is true} \\ 0 & \text{IF } e \text{ is false} \end{cases}$$

2)  $\boxed{K=1}$



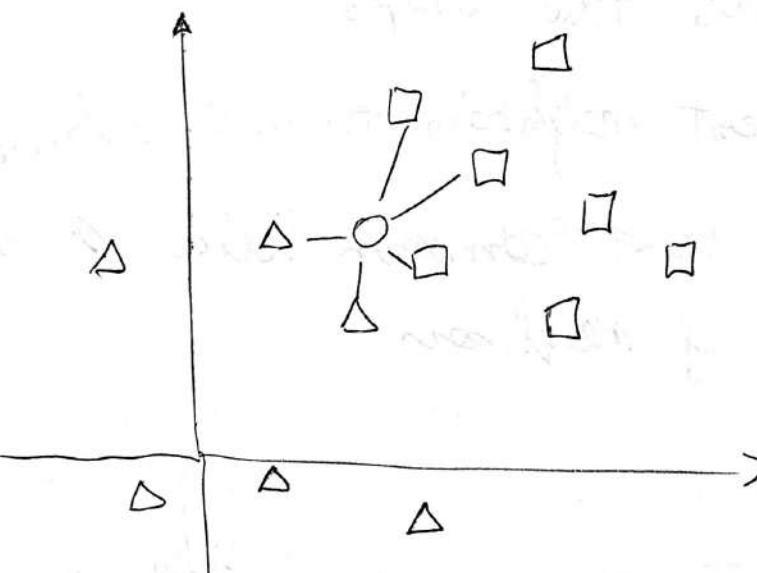
1 - nearest neighbor is "□" thus means that O should be classified as "□" For  $K=1$

$\boxed{K=3}$



IN this case  
3-nearest neighbors  
are 2 △ and 1 □  
So we can say  
that O is classified  
as "△" for  $K=3$

$\boxed{K=5}$



IN this case  
5-nearest neighbors  
are 3 "□" and  
2 "△" so  
the new instance  
'O' is classified  
as "□" for  $K=5$

# EXAM TEST 2

(EX1)

(a)

$$1) d = (X X^T + \lambda I)^{-1} t$$

$K = X X^T$  gram matrix

$$K = \begin{pmatrix} K(x_1, x_1) & \dots & K(x_N, x_1) \\ \vdots & \ddots & \vdots \\ K(x_1, x_N) & \dots & K(x_N, x_N) \end{pmatrix}$$

2) for linear regression

$$f(\theta) = \frac{1}{2} \sum_{n=1}^N \epsilon(y_n, t_n) + \lambda \|w\|^2$$

We want to minimize this so

We have  $\epsilon(y_n, t_n) = (y_n - t_n)^2$  and in closed form solution

We have

$$w = \sum_{n=1}^N d_n x_n^*$$

$$\text{So } y(x) = w^T x = \sum_{n=1}^N d_n x_n^T x$$

We can use ~~use~~ Kernel trick to obtain Vernerized version for linear regression

$$y(x) = \sum_{n=1}^N d_n K(x_n, x)$$

(Ex 2)

- 1) INPUT vector DIM 128  
output vector DIM 10  
HIDDEN layer DIM 50 UNIT

dimension of  $W_1$  is  $128 \times 50$  and dimension of  $W_2$  is  $50 \times 10$

2)

We have hidden layer that use ReLU function, this means that

$$h = g(x^T W_1 + c)$$

where  $g(l) = \max(0, l)$  and it is ReLU activation function  
for output we have linear function

$$y = W_2^T h + b$$

$$y = W_2^T \cdot \max(0, x^T W_1 + c) + b$$

EX 3

Autoencoder is a neural network that ~~has~~ reduced size of hidden layer, which learn to reconstruct by MINIMIZE sum-of-squared error

In particular input and output have some dimension that we consider  $D$ , and hidden layer  $M$  dimension with  $M < D$

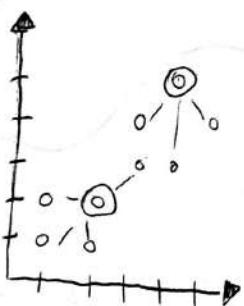
$X_D$	$M$	$Y_D$
0	0	0
0	0	0
.	0	.
;	0	;
;	0	;
$X_1$ , 0	0	1
		$\underbrace{y_s}$
		new linear

Input and output are linear  
and hidden layer ~~are~~ are  
non linear

This means that we have, if we consider 3D of input layer and 2D of hidden layer we have function  $F_1$  that maps  $D$  to  $M$  and function  $F_2$  that maps  $M$  to  $D$  and if  $F_1$  is not linear the output cannot be planar



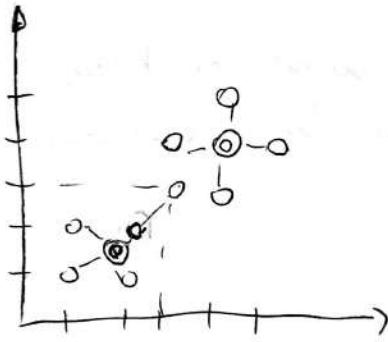
EX 4



In K-means we need to find

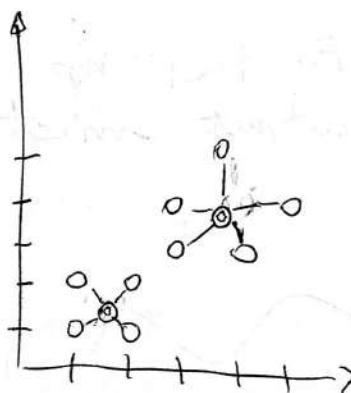
K means of K cluster, in this case we have 2 cluster

I draw in the picture samples of each cluster with specific centroids of the two clusters, now we can understand that we need to update our centroid because the cluster on the top of the picture doesn't have centroid in ~~good~~ right position and also centroid in the bottom



After computing centroids we take each sample of them in sequence and we check if they are in the cluster with nearest centroid.

We notice that sample in (3, 3) position doesn't stay in cluster with nearest centroid, so we need to switch this centroid with the other and update centroids of two clusters involved.



Now we can see that there is no more sample that is not in the cluster with nearest centroid.

K-Means converge, with  $k=2$

### Ex 5

1) KNN algorithm for classification is used for non parametric model

We have target function  $f: X \rightarrow C$  and dataset

$$D = \{(x_n, t_n)\}_{n=1}^N$$

This algorithm is composed by two steps

① Find K neighbors of new instance  $x$

② Assign to  $x$  most common value among the majority of neighbors

likelihood function

$$P(c|x, K, D) = \frac{1}{K} \sum_{N_K(x, D)} I(t_i = c)$$

Where  $N_K(x, D)$  is  $K$  nearest neighbor of  $x$

$$I(e) = \begin{cases} 1 & \text{IF } e \text{ is true} \\ 0 & \text{IF } e \text{ is false} \end{cases}$$

2)



We can say that  
We choose  $K=3$   
nearest neighbor of  
 $o_1, o_2, o_3$

For  $o_1$  we have 2 "+"  
and 1 "\*" so this  
means that we classify  
 $o_1$  as "+".

For  $o_2$  we have 2 "x" and 1 "+"  
so we classified  $o_2$  as "x"

The last one  $o_3$ , in this case we have 2 "-" and 1 "x"  
so  $o_3$  should be classified "-".



Ex 6

We need to consider that if this classifier are not trained yet  
we can use boostify so we train classifier in sequence  
~~this classifier~~ Each classifier is trained on weighted  
sets.

Weighted vote depends on performance of the classifier, so if point is misclassified we have high weight

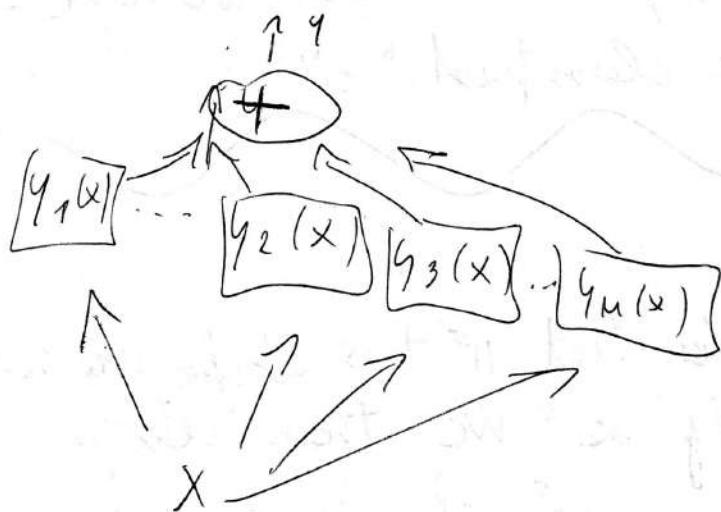
then the prediction of instance is based on the majority of votes

$$\textcircled{1} \quad \begin{array}{c} w_h^{(1)}(x) \\ \downarrow \\ y_1(x) \end{array} \quad \begin{array}{c} w_h^{(2)}(x) \\ \downarrow \\ y_2(x) \end{array} \quad \dots \quad \begin{array}{c} w_h^{(n)}(x) \\ \downarrow \\ y_n(x) \end{array}$$

$$y_m = \text{sign} \left( \sum_{m=1}^n w_m y_m(x) \right)$$

~~Method~~

Instead if our classifier are already trained we consider to use method for multiple learner such as voting in which the approach is to consider that each learner gives a result and all these results will be used to compute final output.



In mixture of expert these learner should be  
specific ~~so~~ because fitting model, which  
look at the input will choose learner  
that is responsible of final output  
But usually this depend, so we use these  
multiple learner method to understand and  
combine results of classifiers • Rather than  
considering only just one

18/01/2019

(B)

Ex-1

- 1) Yes data are separable because exist a line that can divide 3 classes in to 3 different part of a region.

Data sets is not linearly separable but we can separate these classes

- 2) We can use RBF function or polynomial of degree ~~to~~ three as Kernel function both cases should be right

$$K(x, x') = (\beta x^T x' + \gamma)^d \quad d = \{2, \dots\}$$

CASE  
OF  
POLYNOMIAL

$$K(x, x') = \exp(-\beta \|x - x'\|^2)$$

CASE OF  
RBF

- 3) If we apply SVM for classifier of multiple classes we use Lagrange to obtain value of  $w^*$  and  $w_0^*$  that we obtain maximize

Marginal

For K-classes

$$w^* = \sum_{n=1}^N \sum_{k=1}^K Q_{kn} + w_{kn} x_{kn}$$

$$SV = \{x_k \in X_d \mid t_k y(x_k) = 1\}$$

$$y(x) = \text{sign}(w^T x + w_0^*)$$

and

$$w_0^* = \frac{1}{|SV|} \sum_{x_i \in SV} \left( t_i - \sum_{x_j \in S} \cancel{Q_{kj}} + \cancel{t_{kj}} x_{kj}^T \cancel{x_{kj}} \right)$$

[EX 2]

Perceptron is a model for classification that want to find good separator to classify instance.

In particular perceptron works with function  $\delta(x)$ , where  $\delta(x) = \begin{cases} 1 & \text{if } w^T x > 0 \\ -1 & \text{otherwise} \end{cases}$  that we know

$y = w^T x$  so we want to find  $w$  to have good separation of our classes.

Given a dataset  $D = \{(x_n, t_n)\}_{n=1}^N$ , we want to minimize sum of squared error

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

Need to compute

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

thus the main concept of perceptron is to update weights of each feature with this training rule

$$W \leftarrow W + \Delta w_i$$

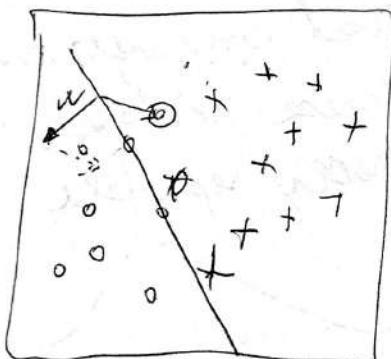
$$\Delta w_i = -\eta \nabla E = +\eta \sum_{n=1}^N (+_n - w^T x_n)$$

$\eta$  is learning rate

Perceptron update weight when it find an instance that is misclassified, so it update weights and the decision boundary changes, we move this until termination condition.

In particular Perceptron converge after a FINITE number of steps IF dataset is linearly separable otherwise it will never converge.

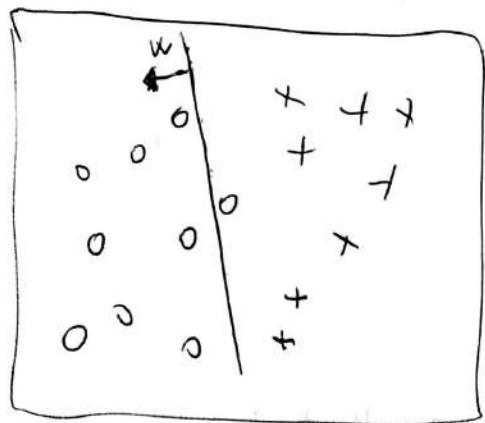
2)



We choose  $w$  randomly, then we see at each iteration if there is an instance that is misclassified. On the bottom we have class 'o' and on top class '+'. We can see that we find an instance of class 'o' over the decision boundary, we need to

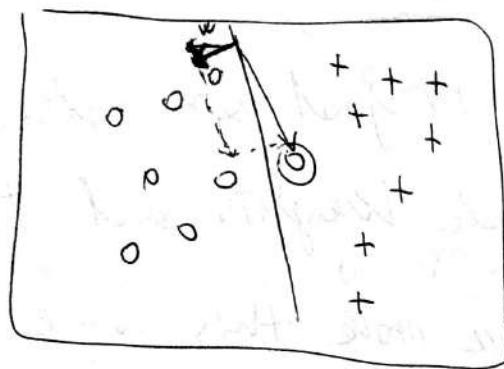
update our weights

b)



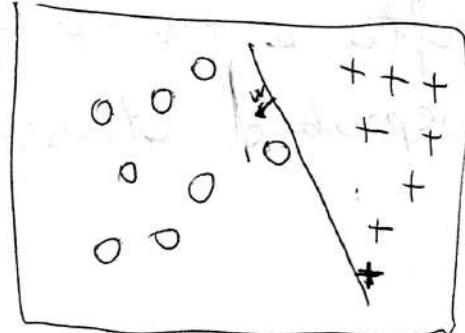
$w$  is smaller than before  
because  $\Delta w_i = -\eta \nabla E$  reduce  
value of  $w$

c)



We find another instance  
that is miss classified  
so we want to update  
weights again and thus  
decision boundary change

d)

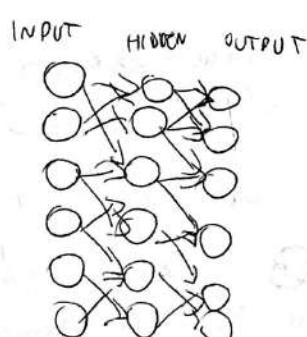


$w$  is reduced again  
but we can see now  
that there is no instance  
that is miss classify  
so method converge and this  
is because data set is  
linearly separable



6x3

1

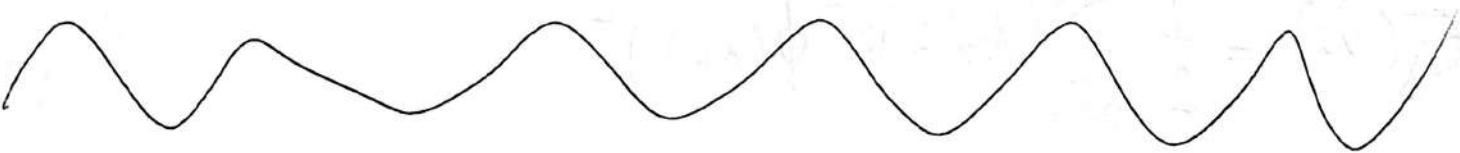


In regression problem we have output layer so activation function of output is linear and this means that loss function is mean squared error.

Furthermore for hidden layer we can choose ReLU function because it's similar of linear function so we prefer for this model

$$\text{No. of parameters } (6 \times 6) + (6 \times 6) + (6 + 6) = 36 + 36 + 12 \\ = 84$$

- 2) No because Backpropagation is algorithm that only compute the gradient is not an algorithm the is a training algorithm so it doesn't suffer of local MINIMA, stochastic gradient descent suffer of local MINIMA
- 3) No because Backpropagation is not a training algorithm so it doesn't suffer problem of over fitting, because it only propagate gradient value from loss function, computed in forward pass to the input in the backward pass.



### EX 4)

Linear regression describe linear relationship between input variable and output variable, in particular learning function  $f: X \rightarrow Y$ , where  $X \subseteq \mathbb{R}^d$  and  $Y = \mathbb{R}$ , given  $D = \{(x_n, y_n)\}_{n=1}^N$ .

$$y = w_0 + w_1 x_1 + w_2 x_2 + \dots + w_d x_d = w^\top x$$

$$w = \begin{pmatrix} w_0 \\ w_1 \\ \vdots \\ w_d \end{pmatrix} \quad x = \begin{pmatrix} 1 \\ x_1 \\ \vdots \\ x_d \end{pmatrix}$$

The goal is to ~~predict~~ define a model that can be used to predict the ~~yes~~ values of  $y$  depend of  $x$ .

2) In linear regression we can use obse known function  $\phi$

$$y = \sum_{i=0}^{m-1} w_i^\top \cdot \phi_i(x_n) = w^\top \phi(x)$$

$$w = \begin{pmatrix} w_0 \\ \vdots \\ w_{m-1} \end{pmatrix} \quad \phi(x) = \begin{pmatrix} \phi_0(x) \\ \vdots \\ \phi_{m-1}(x) \end{pmatrix} \quad \text{where } \phi_0(x) = 1$$

We want to minimize sum of squared error

$$E_g(w) = \frac{1}{2} \sum_{n=1}^N (t_n - w^\top \phi(x_n))^2$$

$$\min E_D(w)$$

but in particular  $t = \begin{pmatrix} t_1 \\ \vdots \\ t_n \end{pmatrix}$   $\phi = \begin{pmatrix} \phi_0(x_1) & \dots & \phi_{n-1}(x_1) \\ \vdots & & \vdots \\ \phi_0(x_N) & \dots & \phi_{n-1}(x_N) \end{pmatrix}$

$$E_D(w) = \frac{1}{2} (t - \phi w)^T (t - \phi w)$$

optimal solution when

$$\nabla E = 0 \Leftrightarrow \phi^T t = \phi^T \phi \cdot w$$

thus

$$w_M = (\phi^T \phi)^{-1} \phi^T \cdot t$$

IF we have dataset that is very huge we can think to do sequential learning we need to use stochastic gradient so

$$\hat{w} \leftarrow \hat{w} + \eta \nabla E$$

$$\hat{w} \leftarrow \hat{w} + \eta \sum_{n=1}^N (t_n - w^T \phi(x_n)) \phi(x_n)$$

Algorithm converge for small value of  $\eta$ , learning rate

Ex 5

$$h_{MAP} = \underset{h \in H}{\operatorname{argmax}} P(h | D) = \underset{h \in H}{\operatorname{argmax}} P(D|h) \cdot P(h)$$

We can drop  $P(D)$  because is constant and independent of  $h$

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

We use maximum a posteriori and maximum likelihood to compute the most probable hypotheses

- 2) Given a dataset  $D = \{(x_i, t_i)\}_{i=1}^N$  and function  $f: X \rightarrow V$ ,  $V = \{v_1, \dots, v_m\}$  and a new instance  $x \notin D$  that is described with values of attribute  $\{e_1, \dots, e_n\}$

$$v_{NB} = \arg \max_{v_j \in V} P(v_j | D) \prod_i P(e_i | v_j, D)$$

If we assume that

$$P(e_1, \dots, e_n | v_j, D) = \prod_i P(e_i | v_j, D)$$

Naive Bayes use ~~assumption~~ conditional independence

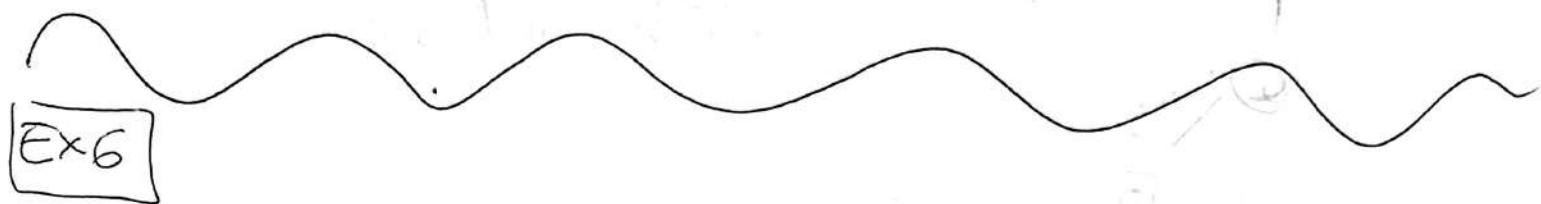
approximate  $P(e_i | v_j, D)$  ~~because~~ this because

we don't have prior probability of hypothesis and in this case we have values of attribute so we don't compute "optimal" value for most probable classification of new instance

but the approximation give us good result.

- 3) In Naive Bayes classifier we don't find optimal

such as Bayes Optimal Classifier but it is practical when we have huge hypothesis space because we cannot have all prior probabilities of hypotheses so we need to use approximation to consider most probable hypothesis for new instance



1) IF  $B = b_1 \wedge A = \alpha_1$  THEN -

IF  $B = b_1 \wedge A = \alpha_2 \wedge C = c_1$  THEN +

IF  $B = b_1 \wedge A = \alpha_2 \wedge C = c_2$  THEN +

IF  $B = b_1 \wedge A = \alpha_2 \wedge C = c_3$  THEN -

IF  $B = b_1 \wedge A = \alpha_3$  THEN +

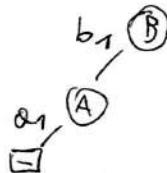
IF  $B = b_2 \wedge C = c_1$  THEN -

IF  $B = b_2 \wedge C = c_2$  THEN -

IF  $B = b_2 \wedge C = c_3$  THEN +

2)

$S_1 = \langle \alpha_1, b_1, c_1, - \rangle$

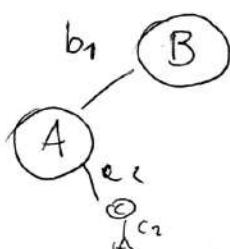


T is consistent because of First Rule we have that IF  $B = b_1 \wedge A = \alpha_1$  THEN -

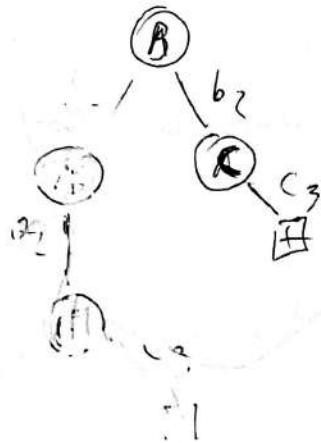
$S_2 = \langle \alpha_2, b_1, c_1, + \rangle$

T is consistent ~~with this set~~ with this ~~subset~~ of sample for third rule

ANSWER



$$\boxed{S_3} = \langle e_1, b_2, c_3, + \rangle$$



As we can see

+ is ~~not~~ consistent  
with set  $S_3$

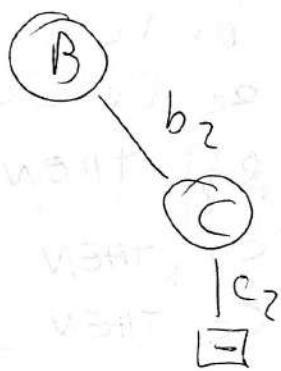
$$\boxed{S_4} = \langle e_2, b_2, c_2, + \rangle$$

T is NOT

consistent with

$S_4$  because we  
expect - but

~~the~~ set  $S_4$  has '+'



27/06/2019

OK

[Ex1]

i) We have  $L$ , is an ~~algorith~~ algorithm  
Split dataset  $D$  in  $K$  set  $S_i \quad i=1, \dots, K$   
 $|S_i| \geq 30$

for  $i = 1 \dots K$

$S_i$  is TEST SET

$T_i = \{D - S_i\}$  is TRAINING SET

$h_i = L(T_i)$  is hypothesis obtain applying that algorithm

$\delta_i = \text{error}_{S_i}(h_i)$

$$\text{error}_{\text{avg}}(h) = \frac{1}{K} \sum_{k=1}^K \delta_i$$

$$\text{Accuracy}_{L,D} = 1 - \text{error}_{S,L}(h)$$

With  $K$ -Fold Cross validation we can estimate accuracy of particular learning algorithm. Split dataset in Test and Training set and train your model on training set that give to us hypothesis that we can use to determine the sample error, see the accuracy of that particular learning algorithm, better is the accuracy, better should be our algorithm.

2) We can also use K-Fold Cross Validation to compare two different learning algorithms  $L_A$  and  $L_B$

Pseudocode

We split dataset  $D$  in  $K$  sets when  $|S_i| \geq 30$

for  $i = 1 \dots K$

$S_i$  is TEST set

$T_i$  is TRAINING set =  $\{D - S_i\}$

$$h_A = L_A(T_i)$$

$$h_B = L_B(T_i)$$

$$\delta_i = \text{error}_{S_i}(h_A) - \text{error}_{S_i}(h_B)$$

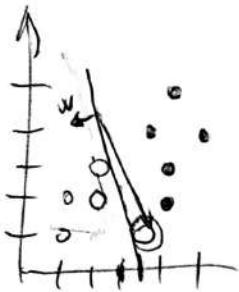
$$\bar{\delta} = \frac{1}{K} \sum_{i=1}^K \delta_i$$

If  $\bar{\delta} < 0$  this means that  $L_A$  is better than  $L_B$

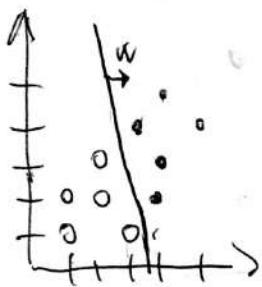
So we evaluate ~~the~~ two ~~learning~~ learning algorithms, considering some training dataset  $T_i$  and we generate two hypotheses and we consider difference between two error of two hypotheses of two learning algorithms and  $\text{error}_{S_i}(h_A) < \text{error}_{S_i}(h_B) \Rightarrow \delta_i < 0$

So Accuracy of ~~of~~  $h_A$  is better than accuracy of  $h_B$  OR vice versa.

notice that we found instance that is misclassified, we update weights and change our decision boundary



another instance  
misclassified, update weights



Data set is linearly  
separable so we converge in  
a finite N° of steps

b) SVM we want to MAXIMIZE margin, distance between decision boundary and closest points from it closest points are support vectors, and at first we have only just one, but when we maximize margin we have 2

$$x_{k+} \text{ and } x_{k-} \quad \text{where} \quad w^T x_{k+} + w_0 = 1$$

$$\text{and} \quad w^T x_{k-} + w_0 = -1$$

$$\text{and distance is } \frac{1}{\|w\|} \quad \text{where } g(x) = w^T x + w_0$$

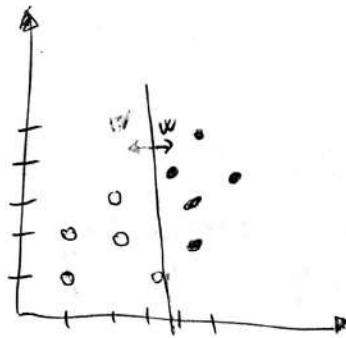
this is the reason why we obtain that solution for SVM

3) SVM is prefer respect to perceptron because if a dataset is not linearly separable, perceptron will never converge, instead SVM use some slack

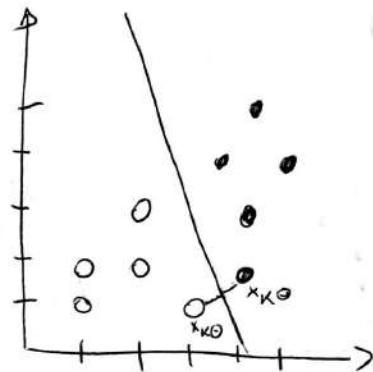
Ex 2

1)

2) PERCEPTRON



b) SUM



2) Both method wants to classify correctly inputs or data

Perceptron use weights and implement a function  
 $\delta(x) = \begin{cases} 1 & \text{if } w^T x > 0 \\ -1 & \text{otherwise} \end{cases}$

$$y = w^T x \quad o = w^T x$$

When start choosing random value of  $w$  and we update weight until there is a terminal condition

$$w^{new} \leftarrow w^{old} - \eta \nabla E$$

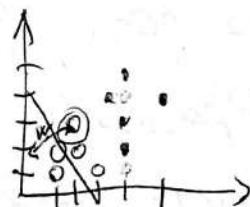
$\eta$  = learning rate

$$\nabla E = -\frac{1}{N} \sum_{n=1}^N (t_n - w^T x_n) (x_n)$$

= gradient error

We update our weight only if we found instance that is missclassified.

For example



on bottom we have instances classified in "0" and one top instances of class "1"

ble to try to find solution of  $W$  and  $w^*$ , so  
and maximum margin.

EX3

ONE STATE MDP

$\langle \{x_0\}, A, S, r \rangle$

$x_0$  is unique state

$A$  is set of actions

$S(x_0, a_i) = \forall a_i \in A$  transition function

$r(x_0, a_i, x_0)$  reward function

We need to determine optimal policy

but in K-armed bandit problem you don't know

$S$  or  $r$  so you need to find some strategy

to find this optimal policy.

In particular here we use value function iteration,  
more specific Q-learning  $Q(x_0, a)$ , after compute this  
 $\pi^*(x_0) = \text{optimal policy} = \arg \max_{a \in A} Q(x_0, a)$

2)

We have  $K$  actions  $a_1, \dots, a_K$   $v(a_i) = N(\mu_i, \sigma_i)$   
gaussian distribution (consider stochastic case)

We choose  $\epsilon$ -greedy strategy so we have uniform  
random choice ~~one~~ with probability  $\epsilon$  and best  
choice with probability  $1 - \epsilon$

the强化 rule is

$$Q_n(a_i) \leftarrow Q_{n-1}(a_i) + \alpha [r_{+/-}(a_i) - Q_{n-1}(a_i)]$$

$$\text{where } \alpha = \frac{1}{V_{n-1}(a_i)}$$

$V_{n-1}(a_i)$  = no of executions of action  $a_i$  up to time  $n-1$

$$P^* = \text{argmax } Q_n(a_i)$$



- 1) Both probabilist model want to compute (estimate)  
 $P(C_i | x)$  when  $i = 1, \dots, h$

the probability that given new istance  $x$  is sampled  
of class  $C_i$ :

generative probabilist model using bayes to estimate  
this value, in particular

$$P(C_i | x) = \frac{P(x|C_i) \cdot P(C_i)}{\sum_{i=1}^h P(x|C_i) P(C_i)}$$

So we compute  $P(x|C_i)$  and we have prior probability  
of that particular class.

Indeed with discriminative we use method  
of logistic regression to estimate  $P(C_i | x)$ , so  
we estimate that value directly from model

2D Detest of  $C_1$  and  $C_2$

$$P(C_i|x) = \frac{P(x|C_i) P(C_i)}{\overline{P(x|C_1) P(C_1) + P(x|C_2) P(C_2)}}$$

We assume  $P(x|C_1) = N(x|\mu_1, \Sigma)$  and  
 $P(x|C_2) = N(x|\mu_2, \Sigma)$

they have gaussian distribution with some covariance,  
and  $N_1$  are samples of  $C_1$  and  $N_2$  samples of  $C_2$ ,  $N = N_1 + N_2$

$$P(C_1) = \pi$$

$$P(C_2) = 1 - \pi$$

We have likelihood func

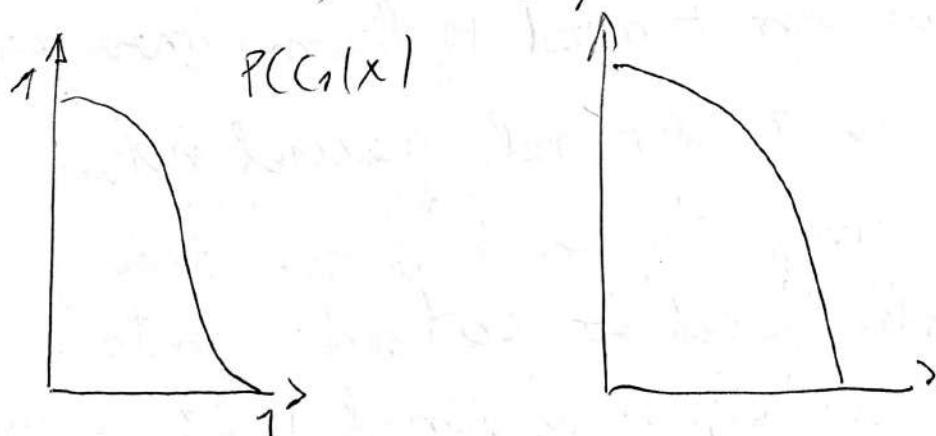
$$P(+|\pi, \mu_1, \mu_2, \Sigma) = \prod_{n=1}^N \pi N(x_n | \mu_1, \Sigma) (1-\pi) N(x_n | \mu_2, \Sigma)$$

Maximize likelihood and  
we obtain

$$\mu_1 = \frac{1}{N_1} \sum_{n=1}^N t_n x_n \quad \mu_2 = \frac{1}{N_2} \sum_{n=1}^N (1-t_n) x_n$$

$$\Sigma = \frac{N_1 S_1}{N} + \frac{N_2 S_2}{N} \quad \pi = \frac{N_1}{N}$$

$$S_i = \frac{1}{N_i} \sum_{n \in C_i} (x_n - \mu_i)^T (x_n - \mu_i) \quad i = 1, 2$$



Ex 5

1) CNN is a neural network that use a convolution function between input and kernel.

It is composed by 3 stage; the first one is convolution stage.

Convolution Stage is

$$3) \text{this is convolution} \quad (I * K)(i, j) = \sum_m \sum_n I(i+m, j+n) K(m, n)$$

For this stage we apply sparse connectivity in which output depends only of a few input, thus kernel  $K$  is much smaller than input, in this way we can reduce memory requirement and improve statistically efficient by reducing number of operations.

Another thing that is done in this stage is parameter sharing in which units share of layer share a set of parameters, in this way we increase efficiency because we don't need to learn more parameters, such as in traditional neural network.

The last thing to consider is the padding which we consider "folded" so outputs contains only valid value (the overlap of kernel) or "some" as input

er is padded with zero and output size independent of kernel size.



1) In a supervised learning we have dataset

$$D = \{(x_n, t_n)\}_{n=1}^N$$

where  $x_n$  are ~~output~~ input and  $t_n$  is the target value associated to that input.

In a unsupervised learning we have a dataset

$$D = \{x_n\}_{n=1}^N$$

We have only input of our dataset, we don't know target value of them

2)

Unsupervised learning we can model problem to ~~choose~~ ~~good~~ partition of dataset in  $k$  cluster. In particular with K-means we can ~~get~~ obtain means of this cluster may a gaussian distribution.

Unsupervised learning uses technique such as clustering, representation and density estimation, with them we hope to learn the inherent structure of our input data without using target value of them.