

Books →

1)

Machine Learning - Tom Mitchell

2)

Pattern Classification - Duda, Hart and Stork.

\* 3)

Introduction to Machine Learning - E Alpaydin

\* 4)

The elements of statistical learning -  
Hastie, Tibshirani, Friedman

19<sup>th</sup> July, 2019

classmate

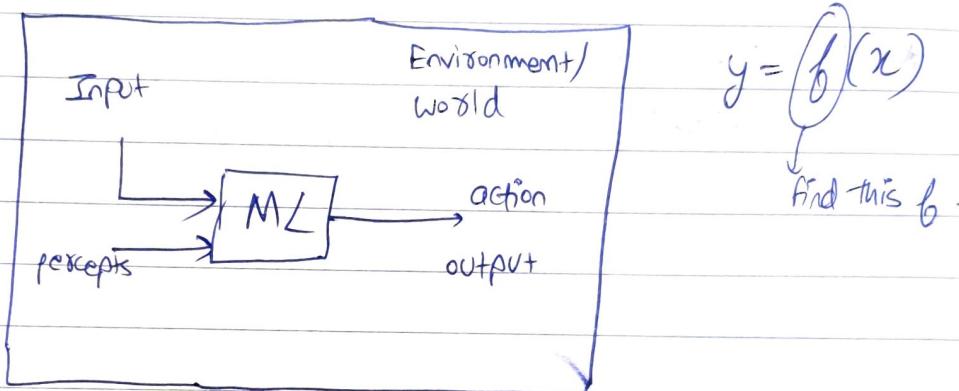
Date \_\_\_\_\_

Page \_\_\_\_\_

## Machine learning

A learning system is one which improves its performance at a task with experience:

Task:



We need a performance measurement to measure the goodness of our function  $f$ .

Learning is changing our function  $f$  when our outputs don't match the expected outputs.

Deductive Knowledge:

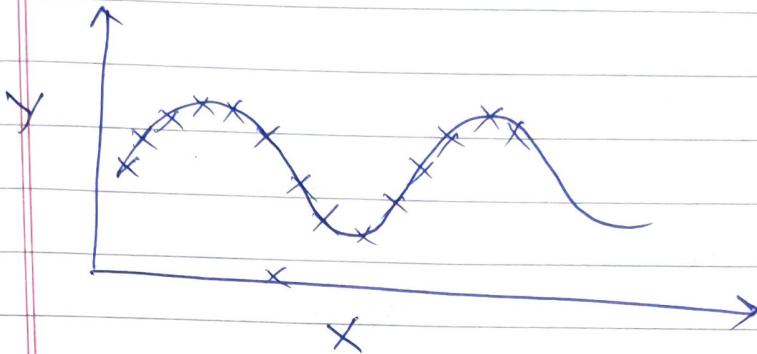
general → special  
(axioms)              (theorems)

Inductive reasoning

- special

general

→ Machine learning



Deductive - given a sinusoidal curve, we know how it will be, always a sinusoidal

Inductive - given a set of points ( $x$ ), we need to fit a wave such that most of these points lie on that.

Types of  $y$  (outputs):

1. According to nature of  $y$ :

—  $y$  is binary 0/1.

—  $y$  is a finite set of integers

1, 2, 3, 4

A, B, C, D

malaria, influ...

Classification / categorization [into one or more classes]

—  $y$  is a real number

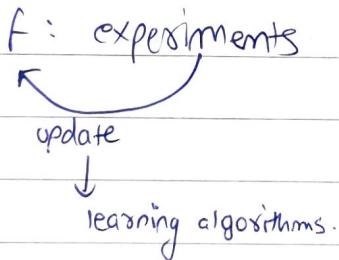
— regression problem.

2. According to nature of  $x$ :

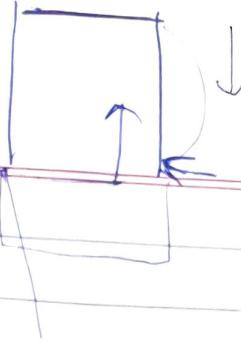
- ✓ —  $x$  is a vector (feature)
- $x$  is a sequence / signal.
- $x$  is a graph
- $x$  is a more complex object

3. According to structure of  $f$ :

Architecture of  $f$ :



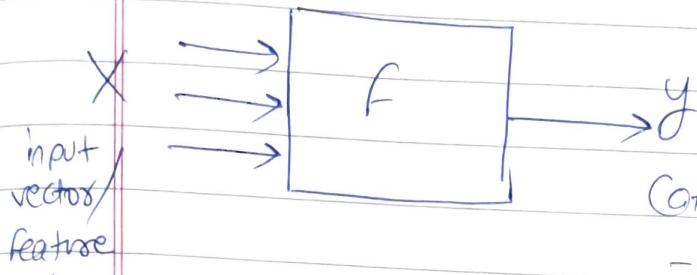
24<sup>th</sup> July



classmate \_\_\_\_\_

Date \_\_\_\_\_

Page \_\_\_\_\_



Category / integers / classes  
- classification  
binary / k-class

Real

- regression

Training -  
supervised learning

$$\left\{ \begin{array}{l} x_1, y_1 \\ x_2, y_2 \\ \vdots \\ x_n, y_n \end{array} \right.$$

Unsupervised learning / exploratory learning

$$\left\{ \begin{array}{l} x_1 \\ x_2 \\ \vdots \\ x_n \end{array} \right\}$$

Semi-supervised learning

$$\left\{ \begin{array}{l} x_1 \\ x_2 \\ x_n \\ x_{n+1} \\ x_{n+2} \end{array} \right\} \quad \left\{ \begin{array}{l} y_{n+1} \\ y_{n+2} \end{array} \right\}$$

Reinforcement learning -

$$\left\{ \begin{array}{ll} x_1 & R(f(x_1)) \\ x_2 & R(f(x_2)) \\ x_3 & \vdots \\ x_n & R(f(x_n)) \end{array} \right.$$

Learning problem

Given a training set ( supervised / unsupervised / reinforcement )

Find  $f$ , which generalizes well to a new  $X$ .

Optimization problem  $\rightarrow$

We decide a functional form for  $f$ .

Error function -  $\ell(f(x), y)$

$$\mathcal{L} = \sum_{i \in TS} \ell(f(x_i), y_i)$$

Find out  $f$  belonging to  $F$ , which has minimum loss  $\mathcal{L}$  on the training set.

The choice of the function  $f$  depends on the nature of  $x$ , nature of  $y$  and how complex the mapping between them is.

If we are constrained by the resources, we would like to find the best possible solution while still being constrained.

Eg - what is the best 10 layer neural network?

26<sup>th</sup> July

Predicting the weather on a day by

- 1) Sky Condition [ sunny, cloudy, raining ]
- 2) Air temp [ warm, cold ]
- 3) Humidity [ normal, high ]
- 4) Wind [ strong, weak ]
- 5) Water [ cool, warm ]
- 6) Forecast [ same, change ]

Output Enjoy sports yes, no binary classification

Data table →

|               | Sky   | Air Temp | Humidity | Wind   | Water | Forecast | Output |
|---------------|-------|----------|----------|--------|-------|----------|--------|
| an instance → | Sunny | Warm     | normal   | strong | cool  | same     | Yes    |

A binary classification problem can be seen as a set learning concept.

c - target concept

h - hypothesis function

$(\text{sky} = \text{sunny}) \wedge (\text{temp} = \text{low})$  ~~some~~

$h \in H \rightarrow$  hypothesis space

cardinality of  $H = 4 \times 3^5$

since any feature may or may not be included in the hypothesis.

? - don't care

$\langle \text{sunny}, ?, \text{normal}, ?, \text{cool}, \text{same} \rangle$

This should be read as -

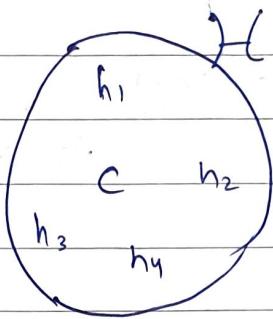
$(\text{sky} = \text{sunny}) \wedge (\text{humidity} = \text{normal}) \wedge (\text{winds} = \text{Gol}) \wedge (\text{forecast} = \text{same})$

But the form to write ~~attribute~~ hypothesis : →

$\langle \text{sunny}, ?, \text{normal}, ?, \text{cool}, \text{same} \rangle$

In a  $\text{G-D}$  space, the hypothesis becomes a set as the don't cares can change about.

- 1) We ~~can~~ make an assumption that the target concept  $c \in \mathcal{H}$  [Realisability assumption]



[Inductive hypothesis assumption]

- 2) If  $c$  matches  $h$  for all the training set, it will match for all  $x$ .

In other words, the learning set has reduced ~~with all~~ to a search problem.

A hypothesis for which  $c$  matches  $h$  over the training set, it is called a consistent  $h$ .

A consistent hypothesis is ~~THE~~ actual hypothesis, according to inductive hypothesis.

Given a set of instances, find the ~~best~~ consistent hypothesis from the family of hypothesis space.

$$T - \text{training} = \{x_1, x_2, \dots, x_n\}$$

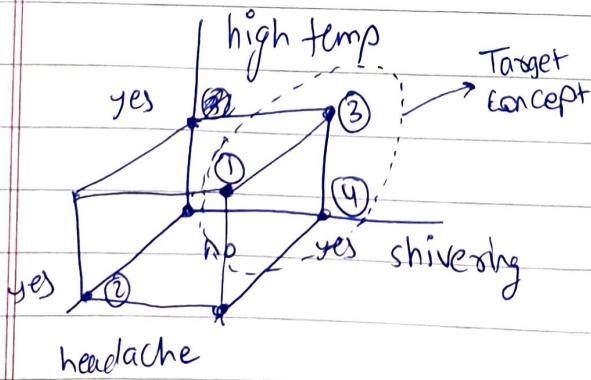
$$\text{if } h(x_i) = c(x_i) \quad \forall x_i \in T \\ \text{then } h(x) = c(x) \quad \forall x$$

31<sup>st</sup> July 2019

| High temp | Headache | shivering | outcome | malaria |
|-----------|----------|-----------|---------|---------|
| Y         | Y        | Y         | Y       |         |
| N         | Y        | N         | N       |         |
| Y         | N        | Y         | Y       |         |
| N         | N        | Y         | Y       |         |

|   |   |   |   |
|---|---|---|---|
| Y | Y | Y | Y |
| N | Y | N | N |
| Y | N | Y | Y |
| N | N | Y | Y |

Attribute space / feature space



①

$$( \text{high temp} = Y \wedge \text{headache} = Y \wedge \text{shivering} = Y )$$

one possible formula, not the correct one.

②

$$(\text{shivering} = Y)$$

$3^3$  possible hypothesis -  $h \rightarrow$  hypothesis  $\mathcal{H}$   
space

The target concept will not ~~ever~~ always be part of our hypothesis space.

If the target concept IS a part of the hypothesis space  $\mathcal{H}$ , then the target concept is realizable.

If  $c \in \mathcal{H}$ , then  $c$  is realizable using  $\mathcal{H}$ .  
*↳*

Instance set:

Set of given examples (instances) with their labels is called instance set.

High temp      Head-      Shivering      malaria  
ache

|   |   |   |   |
|---|---|---|---|
| Y | Y | Y | Y |
| N | N | N | N |
| Y | N | Y | Y |

$$\left. \begin{array}{l} C_1 \\ C_2 \\ C_3 \end{array} \right\} \left. \begin{array}{l} \langle \text{high temp} = Y \wedge \text{shivering} = Y \rangle \\ \langle \text{high temp} = Y \rangle \\ \langle \text{shivering} = Y \rangle \end{array} \right\} \checkmark$$

$C_1, C_2, C_3 \subset \mathcal{H}$

All these  $C_1, C_2, C_3$  are consistent with the

$\vee \in \mathcal{H}$  is called the version space of  $\mathcal{H}$ ,  
w.r.t instance set  $I$ .

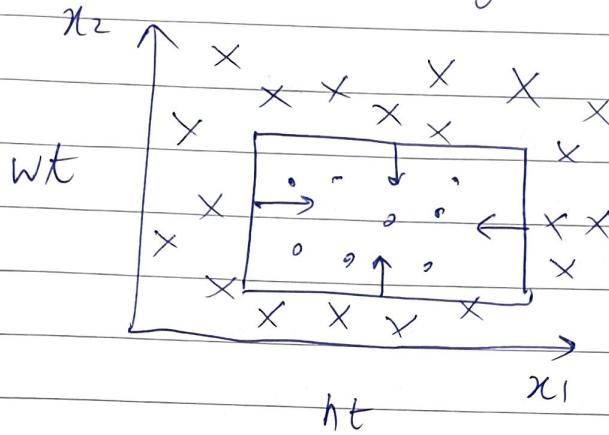
A set of all consistent hypothesis = version space

Given an instance set  $I$  and a hypothesis space  $\mathcal{H}$  find a hypothesis  $h^* \in \mathcal{H}$  which is the closest approximation of target concept  $c$ .

Learning problem

## Inductive principle:

Any consistent hypothesis is a good approximation of the target concept if the training space is large and realizability holds.



high temp    headache    shivering    malacia .

|   |   |   |   |
|---|---|---|---|
| N | Y | N | N |
| Y | N | Y | Y |
| N | N | Y | Y |

high temp = Y  $\wedge$  head ache =  $\text{?}$   $\wedge$  shivering = Y | h<sub>1</sub>

high =  $\text{?}$   $\wedge$  shivering = Y | h<sub>2</sub>

shivering = Y | h<sub>3</sub>

1<sup>st</sup> August

classmate

Date \_\_\_\_\_

Page \_\_\_\_\_

## Training examples for Enjoy Sport -

Learning means learning a subset of the features space.

- If formula  $f_1$  satisfies more examples than formula  $f_2$  then  $f_1$  is said to be more general than  $f_2$ .

Find Specific algorithm:

We start with a very specific hypothesis and then make the minimum incremental changes so that the resulting hypothesis satisfies our positive examples.

Using this algorithm, we obtain the most specific hypothesis which is consistent with the given training examples.

| Sunny | Warm | Normal | Strong | Warm | Same   | Yes |
|-------|------|--------|--------|------|--------|-----|
| Su    | W    | High   | S      | W    | S      | Y   |
| Rainy | Cold | H      | S      | W    | Change | No  |
| S     | W    | H      | S      | Cold | C      | Y   |

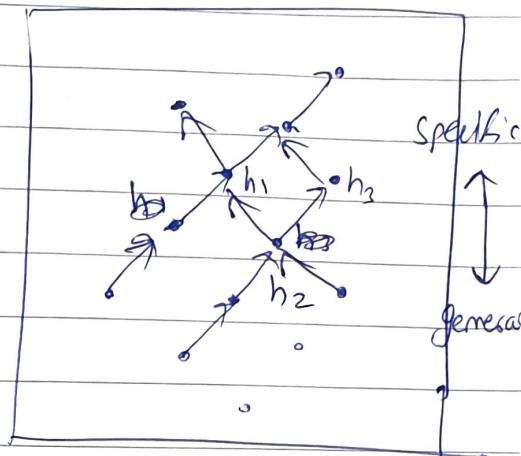
Wind = strong

classmate

Date \_\_\_\_\_

Page \_\_\_\_\_

hypothesis H



$$h_1 = \langle \text{sunny? ? cool? ?} \rangle$$

$$h_2 = \langle \text{sunny? ? ? ? ?} \rangle$$

$$h_3 = \langle \text{sunny? ? ? ? ? humid} \rangle$$

general → specific

The findSpecific algorithm will go down the lattice (depending on the order of training examples faced) and stops at the least general sol<sup>n</sup> satisfying all examples.

Version space for this example

$$\langle \text{sunny, warm? strong? ?} \rangle$$

$$\langle \text{sunny? ? strong? ?} \rangle \rightarrow \langle \text{sunny, warm? ? ? ?} \rangle \langle \text{? warm? strong? ?} \rangle$$



① We are ensuring only the positive examples amount to a positive result.

② We are ensuring this for all positive examples

① & ② implies → all negative examples will have negative output

$\langle S, W, N, S, W, S \rangle$   $\langle S, W H S W S \rangle$   $\langle S W H S C C \rangle$

$\langle S W ? S W S \rangle$ .

$\langle S W ? S ?? \rangle$

$\langle S W ??? \rangle$

$\langle S, ??? \rangle$      $\langle ? W ??? \rangle$

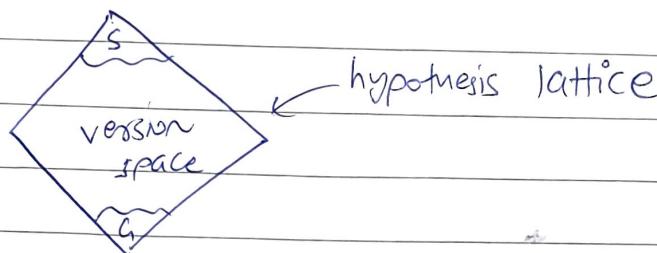
2<sup>nd</sup> August

classmate

Date \_\_\_\_\_

Page \_\_\_\_\_

To find all members of the version space, it is enough to find the most specific consistent example and the most general consistent example.



Candidate elimination algorithm →

$C \leftarrow$  maximally general hypothesis in  $H$   
 $S \leftarrow$  maximally specific hypothesis in  $H$

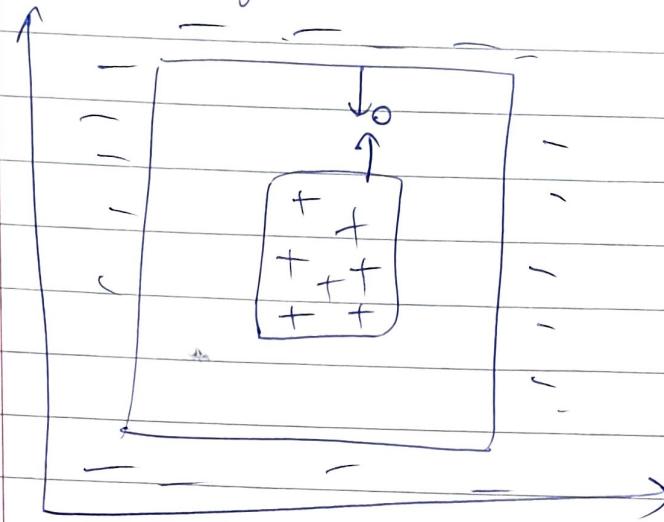
For each training example  $d$ , do

IF  $d$  is positive

— remove from  $G$  every hypothesis inconsistent with  $d$

— for each hypothesis  $s$  in  $S$  that is

To minimise the version space we will choose a new training example which is most ambiguous.



~~Since~~ since we don't know if the new training example is +ive or -ive, so we choose the new example is closer at the middle.

If there is no bias (in our example, the knowledge that the hypothesis is purely conjunctive), then we will require infinite training examples.

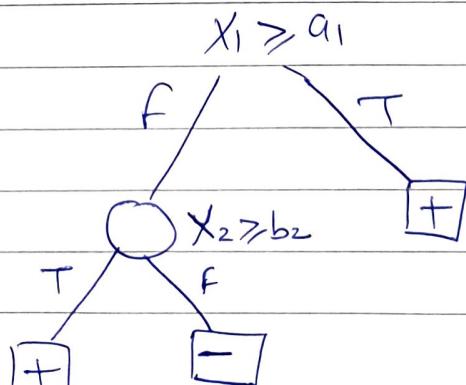
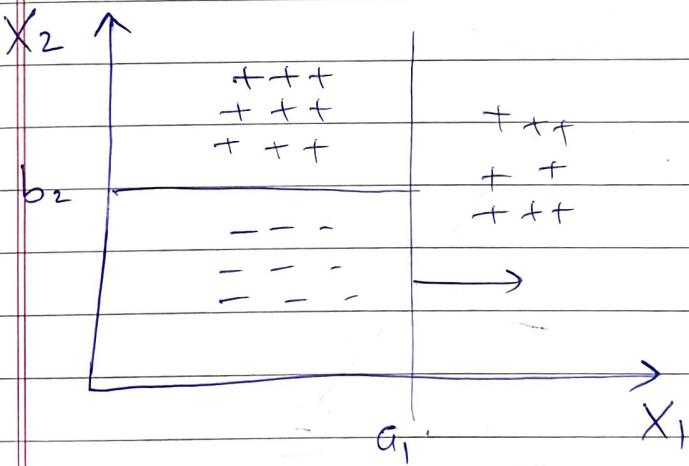
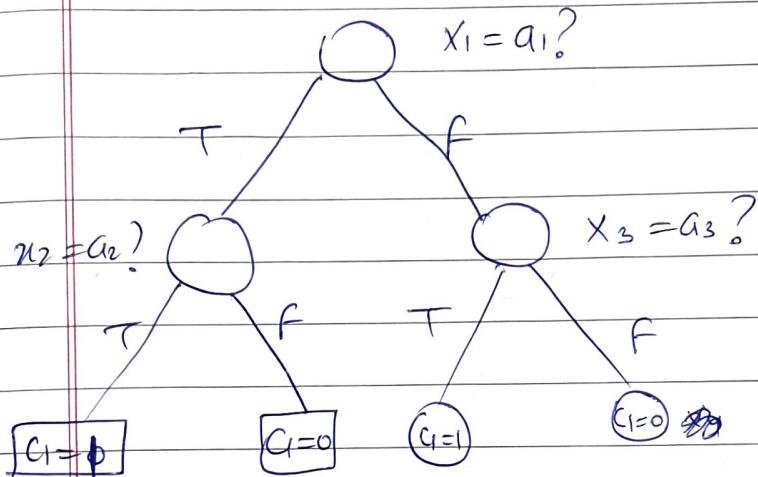
7<sup>th</sup> August '19

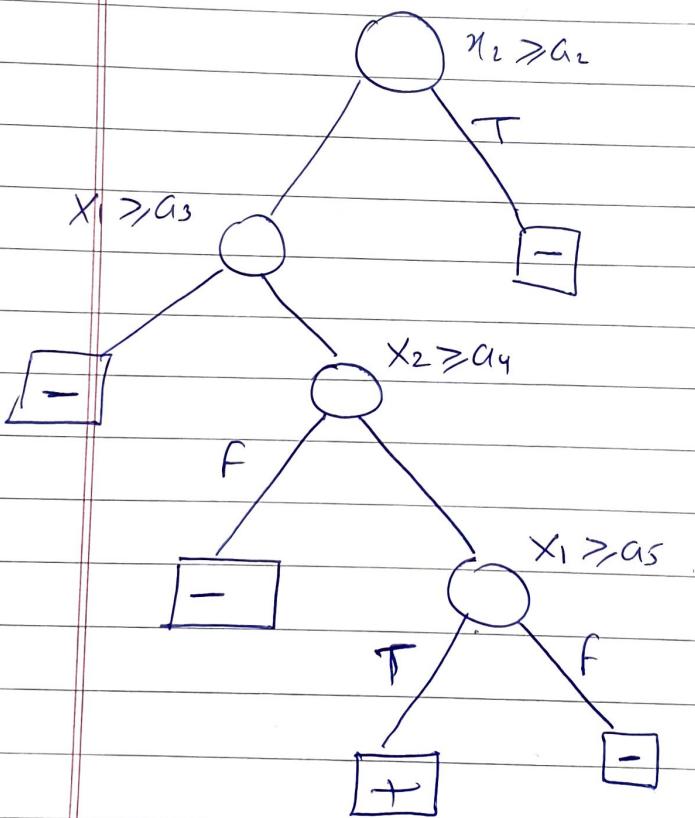
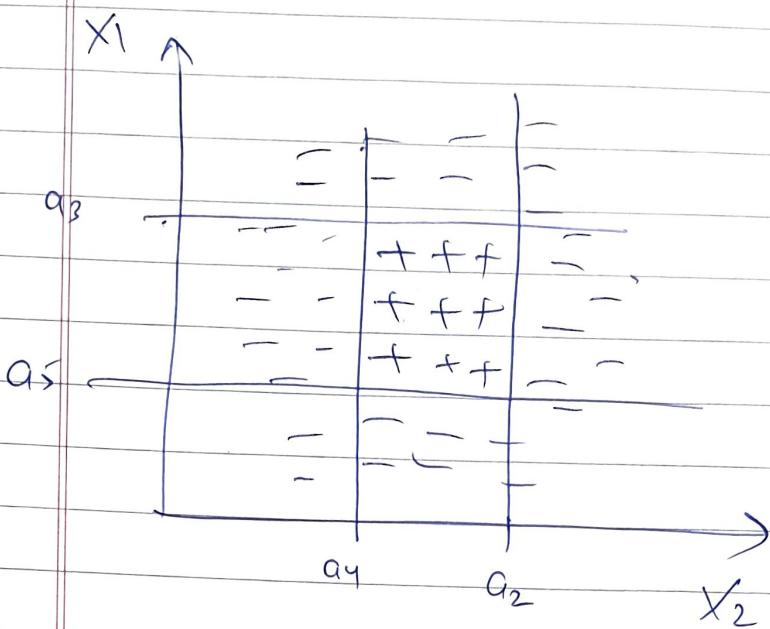
## Machine learning

$$(X_1 = a_1) \wedge (X_2 = a_2) \wedge (X_3 = a_3)$$

$$(X_1 = a_1) \wedge X_2 = a_2 \vee (X_3 = a_3)$$

Decision tree classifier →

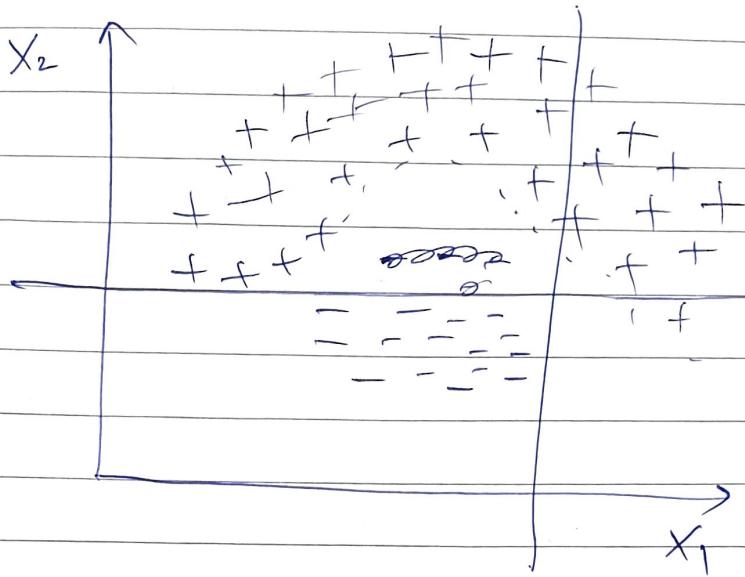




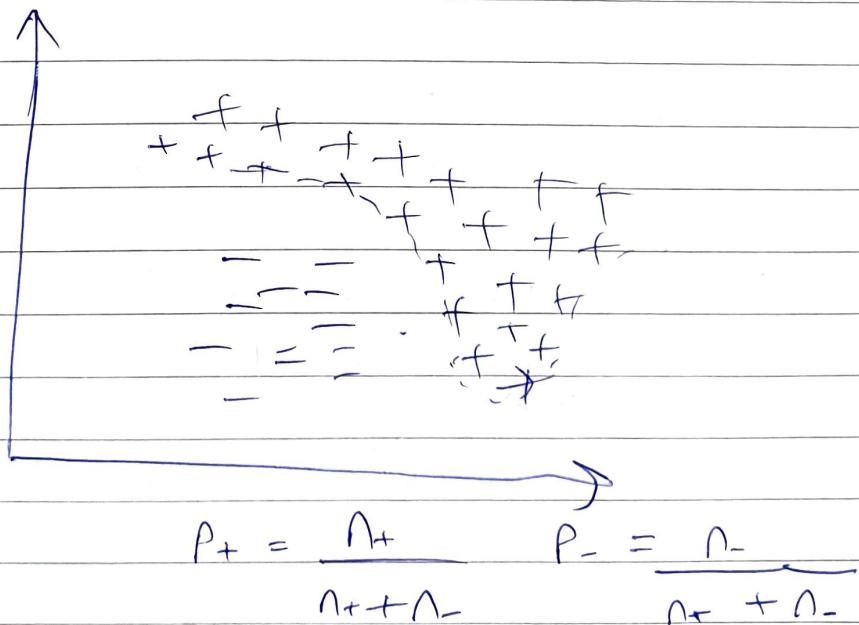
\* Each node partitions the input space into 2 classes.

\* We can obtain multiple decision trees for the same decision problem.

We would like to find the smallest decision tree.



- \* We would want to recursively divide our resultant data-set into pure and non-pure partitions.

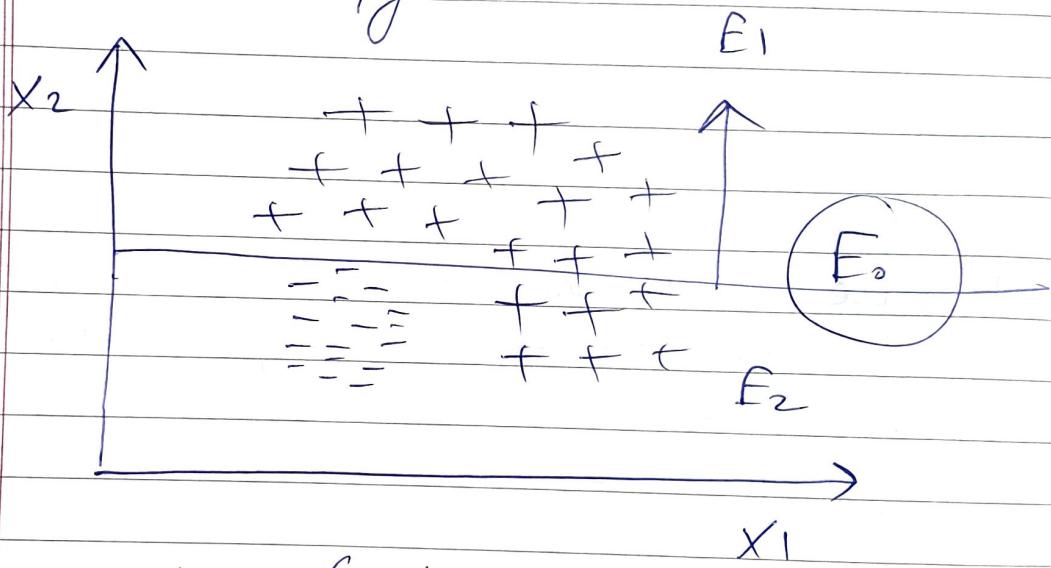


Measure of purity of classes  $\rightarrow$  entropy

$$\text{Entropy} = -p_+ \log p_+ - p_- \log p_-$$

at  $p_+ = 1$  and  $p_- = 0$   
or  $p_+ = 0$  and  $p_- = 1$

then entropy = 0



goodness of cut

= reduction in entropy

$$\begin{aligned}
 &= E_0 - (E_1 + E_2) \\
 &= -( (E_1 + E_2) - E_0 )
 \end{aligned}$$

information gain

We greedily choose the cut with maximum information gain.

8<sup>th</sup> August, '19

$\langle F, I, CS, M \rangle$

$\langle F, I, EE, M \rangle$

$\langle F, I, \cancel{EE}, ?, M \rangle$

$\langle F, ?, ?, M \rangle$

$G = \{ \langle ?, ?, ?, ?, ? \rangle \}$

$S = \{ \langle F, I, CS, M \rangle, \langle F, I, EE, M \rangle,$   
 $\langle F, \cancel{I}, CS, M \rangle, \langle F, I, ?, M \rangle,$   
 $\langle F, ?, CS, M \rangle, \langle F, I, ?, M \rangle,$   
 $\langle F, ?, ?, M \rangle, \langle F, ?, ?, M \rangle,$   
 $\langle F, ?, ?, M \rangle, \langle F, ?, ?, M \rangle \}$

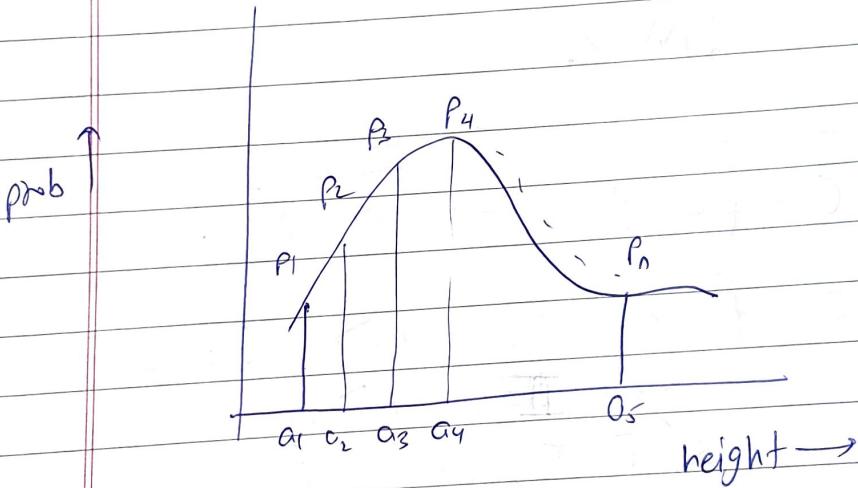
$G = \{ ???M \} \quad \langle F ??? \rangle$

9<sup>th</sup> August, Friday

## Information gain

- $\text{Gain}(S, A)$  : reduction in entropy after choosing attribute A.

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



$$\text{entropy} = - \sum_{i=1}^n p_i \log(p_i)$$

This greedy approach will not result in a decision tree of shortest length.

But since our algorithm is biased towards shorter trees.

\* Olchon's razor principle - use the minimal instrument required to solve the problem.

Methods to avoid overfitting -

\* Pre-pruning

\* Post-pruning

Methods to evaluate trees to prune -

### Gini index

Another sensible measure of impurity

$$Gini = \sum_{i \neq j} p(i)p(j)$$

B 14<sup>th</sup>

August, 2019

classmate

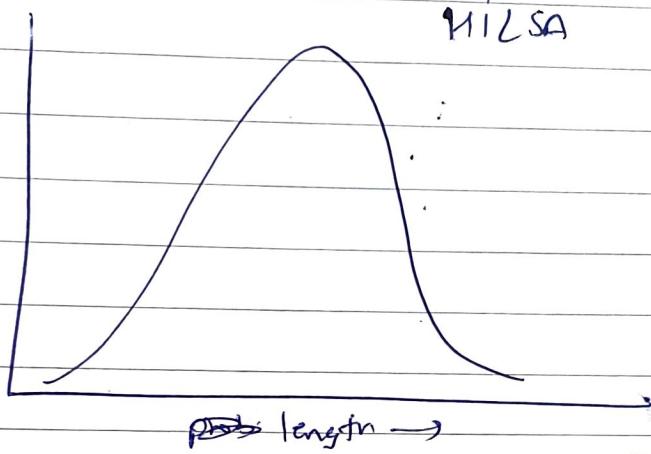
Date \_\_\_\_\_

Page \_\_\_\_\_

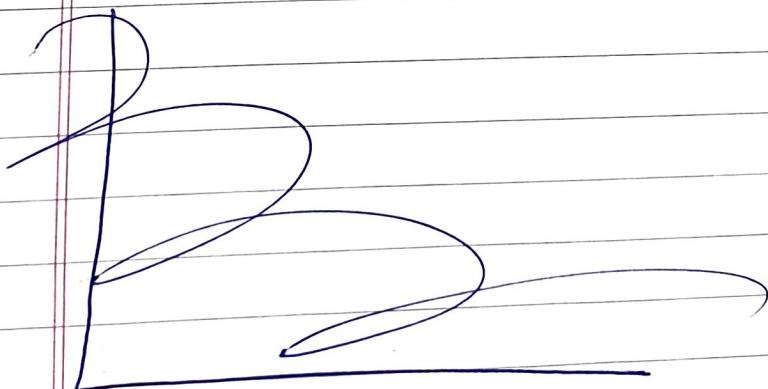
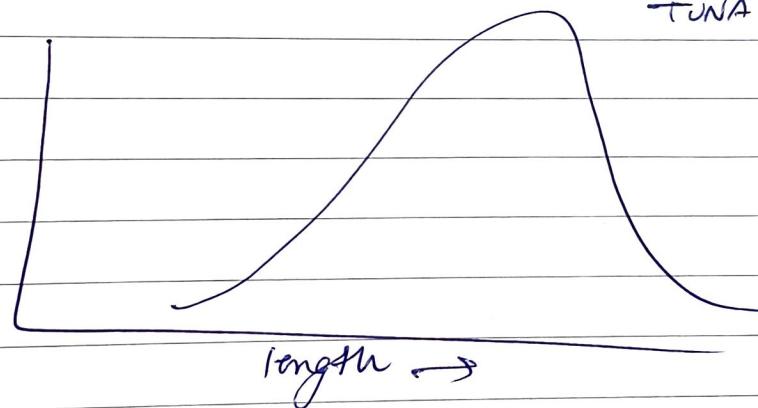
## A simple species classification problem

Measure the length of ~~this~~ a fish and decide its class.

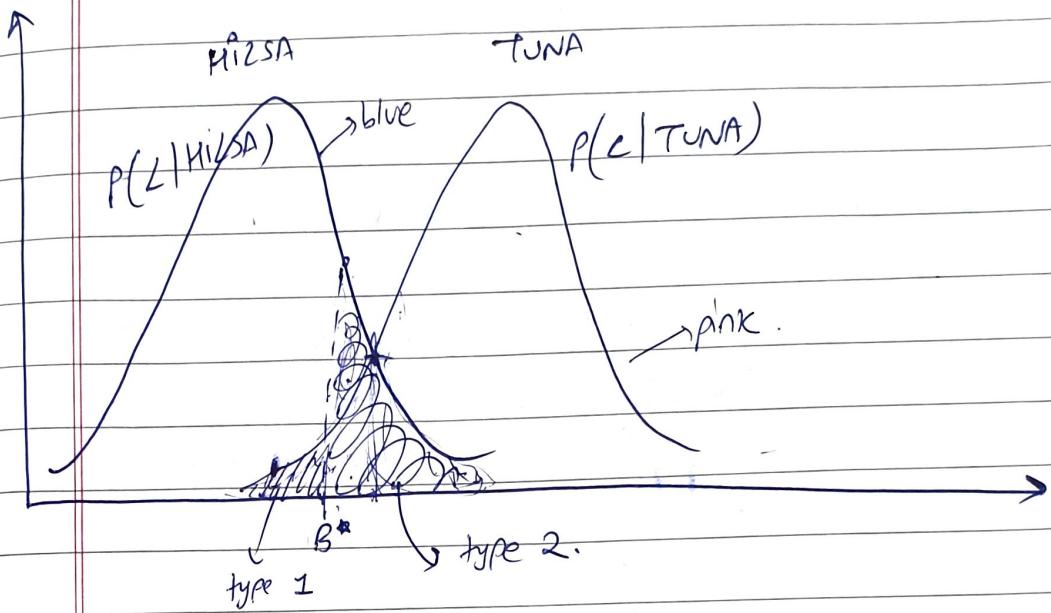
prob.  
dist.



prob.  
dist.



## Error of decision rule -



Type I  $\rightarrow$  Actually HILSA, classified as TUNA

Type II  $\rightarrow$  Actually TUNA, " " " HILSA

$B^*$   $\rightarrow$  threshold value

Minimum possible error

$$P(B^* | \text{HILSA}) = P(B^* | \text{TUNA})$$

If type I error and type II error  
are different,  $\Rightarrow$  boundary shifts.

Location / time of experiment.

more

In Colcutta  $\rightarrow$  ~~more~~ hilsas

In California  $\rightarrow$  more TUNAS.

### Apriori probability

If we were to guess the type of fish in CCU, we would guess hilsa.

Similarly, in California, we would guess Tuna.

We multiply our ~~conditional~~ class Conditional probabilities with our prior probabilities.

Posteriori  $\approx$  apriori  $\times$  class probability

Apriori

$$P(\text{Hilsa})$$

$$P(\text{Tuna})$$

$$P(A|B) = \frac{P(B)P(B|A)}{P(B)}$$

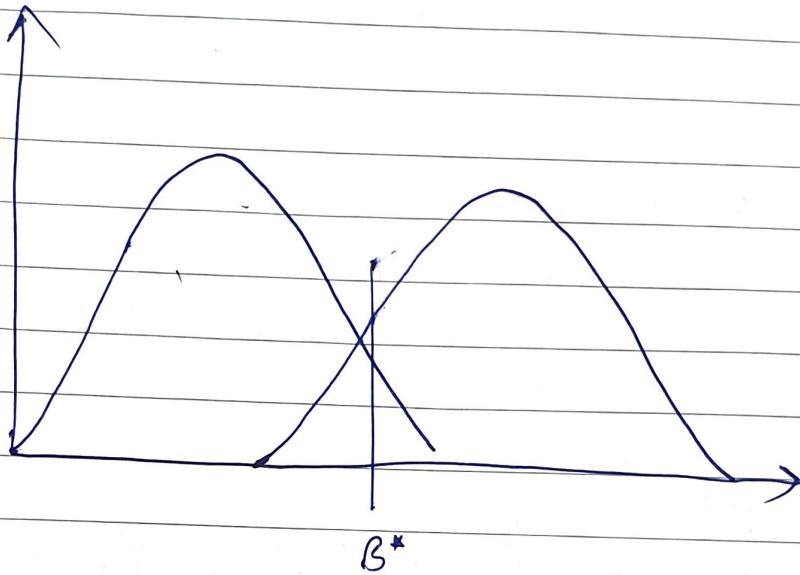
Posteriori

$$\frac{P(L=2ft | \text{Hilsa}) P(\text{Hilsa})}{P(L=2ft)}$$

$$= P(\text{Hilsa} | L=2ft)$$

$$\frac{P(L = \text{2ft} | \text{Tuna}) P(\text{Tuna})}{P(L = \text{2ft})}$$

$$= P(\text{Tuna} | L = \text{2ft})$$



Bayes optimal classifier

$$\text{Aposterior} = \frac{P(L | \text{Hilsa}) P(\text{Hilsa})}{P(L)}$$

this does not change the boundary point since it is just a constant in both Nilsa and Tuna

So it is generally dropped.

15<sup>th</sup> August, 19

TM

classmate

Date \_\_\_\_\_

Page \_\_\_\_\_

~~MAP Classified~~

Entropy

$$H(X) = - \sum_{i=1}^n p(x=i) \log_2(p(x=i))$$

Entropy  $\propto$  impurity of data.

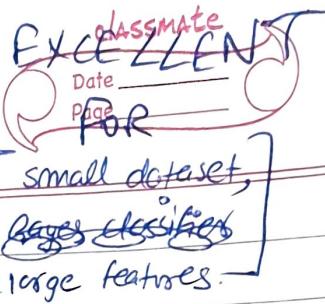
Specific conditional entropy of  $H(X|Y=v)$   
of  $X$  given  $X=v$

$$H(X|Y=v) = - \sum_{i=1}^n p(x|y=v) \log p(x|y=v)$$

Conditional entropy  $H(X|Y)$  of  $X$  given  $Y$ :

$$H(X|Y) = \sum_{v \in \text{values}(Y)} p(y=v) H(X|y=v)$$

Information gain = reduction in entropy  
 $= H(X) - H(X|Y)$



16<sup>th</sup> August, 19.

Bayes classifier

small dataset,  
large features.  
Bayes classifier

$$P(A \wedge B) = P(A) \cdot P(B) \quad \text{if } A \perp B.$$

if  $A \perp B \rightarrow$  then joint prob = prod. of marginals

$$P(A_1 \wedge A_2 \wedge \dots \wedge A_n | M) = P(A_1 | M) \cdot P(A_2 | M) \cdot \dots \cdot P(A_n | M)$$

if  $A \perp B$ .

There is no way to know if the attributes  $A$  and  $B$  are independent. We assume them suitably using our world knowledge.

$$A : [a_1, a_2, a_3, a_4] \xrightarrow{\text{given by}} [a_1, a_2, a_3, a_4 \rightarrow \text{features}]$$

$$P(A | M) = P(a_1 | M) \cdot P(a_2 | M) \cdot P(a_3 | M) \cdot P(a_4 | M)$$

$$P(A | M) P(M) > P(A | N) P(N)$$

Often, we compose the log of conditional probabilities instead of composing the actual conditional probabilities.

Smoothing -

If one of the conditional probability is 0, then the entire expression becomes 0.

Prob. estimation -

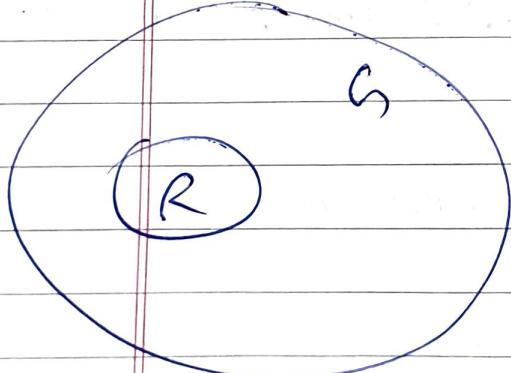
$$\text{Original: } P(A_i | c) = \frac{N_{ic}}{N_c}$$

$$\text{Laplace: } P(A_i | c) = \frac{N_{ic} + 1}{N_c + C} \quad \begin{matrix} C: \text{no. of} \\ \text{classes} \end{matrix}$$

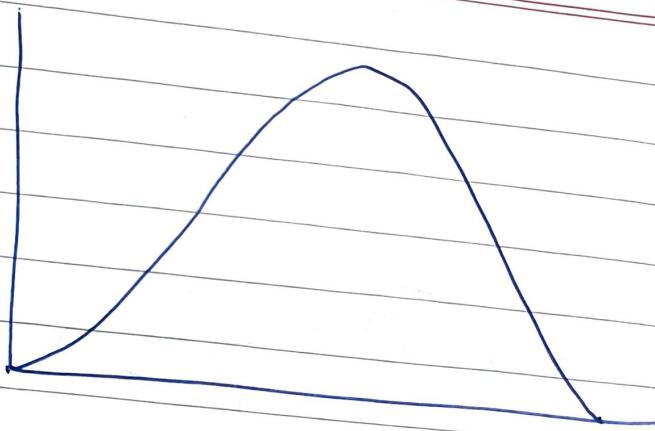
$$\text{m-estimate: } P(A_i | c) = \frac{N_{ic} + mp}{N_c + m} \quad \begin{matrix} m: \text{parameter} \\ p: \text{prior prob.} \end{matrix}$$

Shrinkage -

$$R \subseteq G$$



$$\hat{P}(R) = \lambda \hat{P}(R) + (1-\lambda) \hat{P}(G)$$



$$N(x, \mu, \sigma) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

Q. find the decision boundary when  
 $N(x_1, \mu_1, \sigma_1)$  and  $N(x_2, \mu_2, \sigma_2)$

21<sup>st</sup> August, 2019

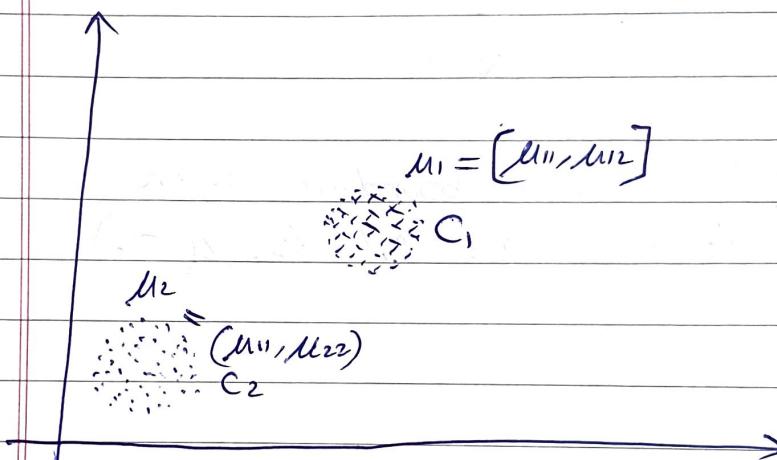
~~AAGA~~

Two class Bayes classifier

Each class is a multivariate gaussian.

$$N(x, \mu, \sigma) = \frac{1}{\sqrt{2\pi} \sigma} e^{-\frac{1}{2} \left( \frac{x-\mu}{\sigma} \right)^2}$$

$$\mu =$$



covariance matrix ( $\Sigma$ )

$$\begin{bmatrix} \sigma_{11} & \sigma_{12} \\ \sigma_{21} & \sigma_{22} \end{bmatrix}$$

for d-dimensional case.

$$X_1 = \{x_{11}, x_{12}, \dots, x_{1d}\} \quad C_1$$

$$X_2 = \{x_{21}, x_{22}, \dots, x_{2d}\} \quad C_2$$

$$\vdots$$

$$X_n = \{x_{n1}, x_{n2}, \dots, x_{nd}\} \quad C_n$$

$$\mu_1 = \frac{1}{n_1} \sum_{i=1}^{n_1} x_i \quad , \quad x_i \in C_1$$

$$= E(x_i) \quad , \quad x_i \in C_1$$

$$\sigma_{11} = E[(x_{i,1} - E(x_{i,1}))^2]$$

~~$$\sigma_{12} = E[(x_{i,1} - E(x_{i,1}))^2]$$~~

$$\Sigma = \begin{bmatrix} \sigma_{11} & \sigma_{12} & \dots & \sigma_{1d} \\ \sigma_{21} & \sigma_{22} & \dots & \sigma_{2d} \\ \vdots & & & \vdots \\ \sigma_{d1} & \sigma_{d2} & \dots & \sigma_{dd} \end{bmatrix}_{d \times d}$$

for 2 dimensional case:-

$$\sigma_{11} = E[(x_{i,1} - E(x_{i,1}))^2]$$

$$\sigma_{22} = E[(x_{i,2} - \mu_2)^2]$$

$$\sigma_{12} = E[(x_{i,1} - \mu_1)(x_{i,2} - \mu_2)]$$

$$= E[x_{i,1} x_{i,2}] - E[x_{i,1}] \cdot E[x_{i,2}]$$

If  $\sigma_{12} \rightarrow$  if 1 increases, then 2 increases  
is high

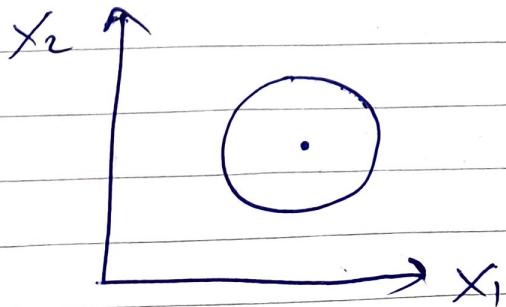
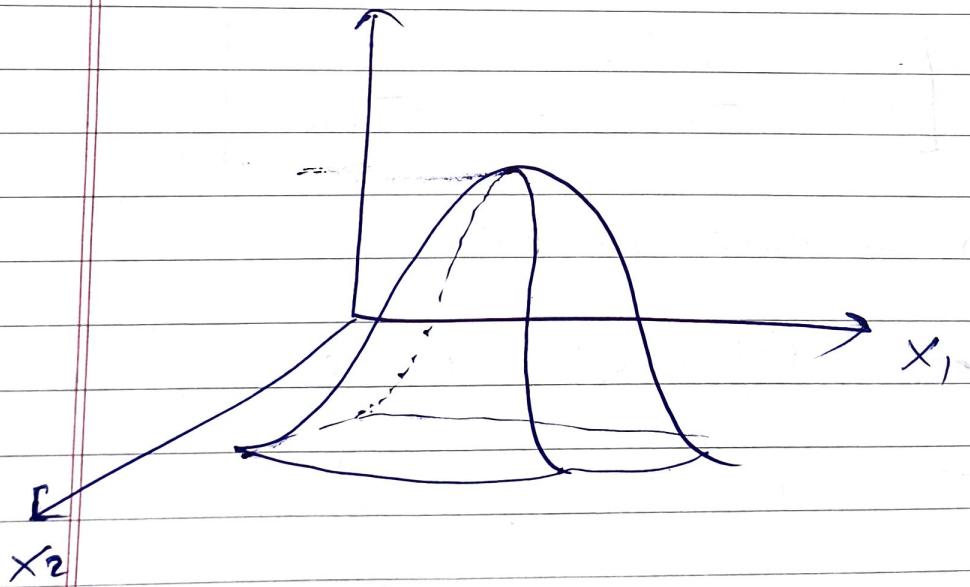
$$\frac{\text{Cov}(X_1, X_2)}{\text{Var}(X_1) \text{Var}(X_2)} = \rho$$

$$N(\mathbf{X}, \boldsymbol{\mu}, \boldsymbol{\Sigma})$$

$$\boldsymbol{\mu} = [\mu_1, \mu_2]$$

$$\boldsymbol{\Sigma} = \begin{bmatrix} \sigma_{11} & \sigma_{12} \\ \sigma_{21} & \sigma_{22} \end{bmatrix}$$

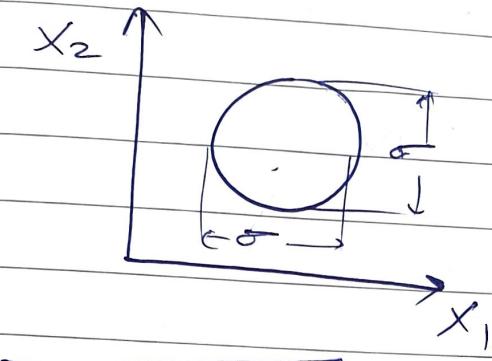
always symmetric



$$\Sigma = \begin{bmatrix} \sigma_{11} & 0 \\ 0 & \sigma_{22} \end{bmatrix}$$

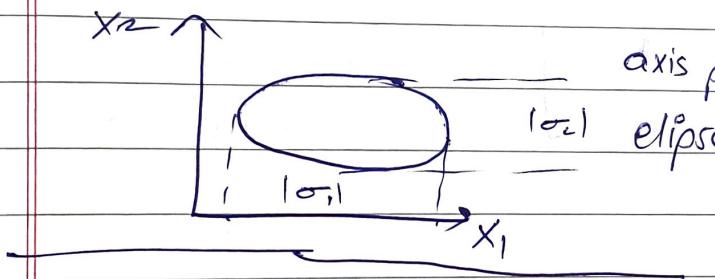
$\sigma_{11} \rightarrow$  dispersion along  $x_1$  direction

$\sigma_{22} \rightarrow$  dispersion along  $x_2$  direction.



$$\Sigma = \begin{bmatrix} \sigma_{11} & 0 \\ 0 & \sigma_{22} \end{bmatrix}$$

axis parallel ellipse.



$$\begin{bmatrix} \sigma_{11} & \sigma_{12} \\ \sigma_{12} & \sigma_{22} \end{bmatrix}$$



Bayes classifier  $\rightarrow$  (posterior prob. distribution)

$$C_1 = P_1(x) = N(x_0, \mu_1, \Sigma_1)$$

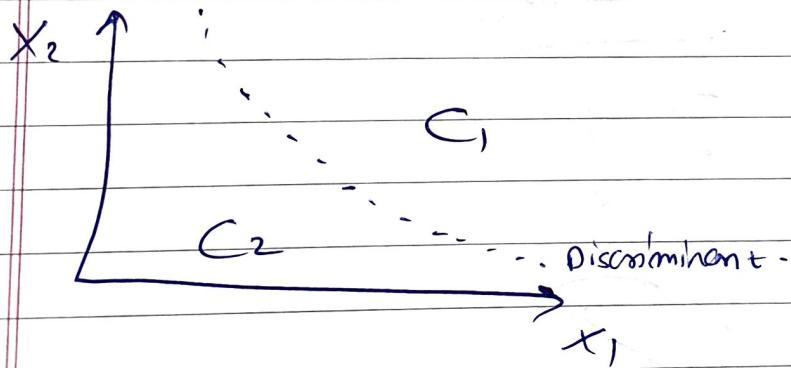
$$C_2 = P_2(x) = N(x, \mu_2, \Sigma_2)$$

log Likelihood ratio.

$$\frac{P_1(x)}{P_2(x)} = \log \left[ \frac{N(x_0, \mu_1, \Sigma_1)}{N(x_1, \mu_2, \Sigma_2)} \right]$$

$$\text{sign} \left[ \log \frac{P_1(x)}{P_2(x)} \right]$$

if +ive  $\rightarrow$  class 1  
-ive  $\rightarrow$  class 2



$$G_1 = \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}(x - \mu_1)^T \Sigma_1^{-1} (x - \mu_1)}$$

$$G_2 = \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}(x - \mu_2)^T \Sigma_2^{-1} (x - \mu_2)}$$

$$-\frac{1}{2} (x - \mu_1) (\varepsilon_1 \varepsilon_1^T)^{-1} (x - \mu_1)^T$$

$$C_1 = \frac{1}{\sqrt{2\pi}} (\varepsilon_1 \varepsilon_1^T)^{-1} e$$

$$-\frac{1}{2} (x - \mu_2) (\varepsilon_2 \varepsilon_2^T)^{-1} (x - \mu_2)^T$$

$$C_2 = \frac{1}{\sqrt{2\pi}} (\varepsilon_2 \varepsilon_2^T)^{-1} e$$

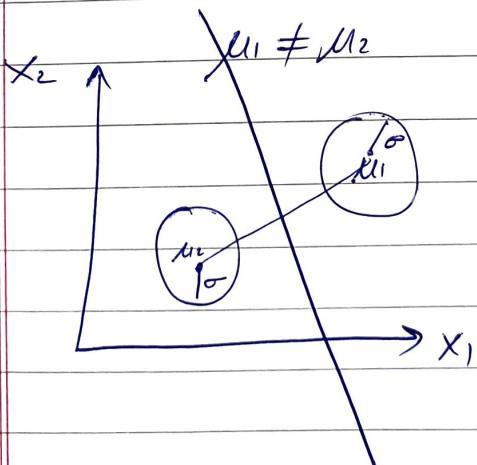
$$C_1 = C_2$$

$$-\frac{1}{2} (x - \mu_1) (\varepsilon_1 \varepsilon_1^T)^{-1} (x - \mu_1)^T + \log \frac{1}{\sqrt{2\pi}} + \log (\varepsilon_1 \varepsilon_1^T)^{-1}$$

$$= -\frac{1}{2} (x - \mu_2) (\varepsilon_2 \varepsilon_2^T)^{-1} (x - \mu_2)^T + \log \frac{1}{\sqrt{2\pi}} + \log (\varepsilon_2 \varepsilon_2^T)^{-1}$$

~~for case 1 :~~

$$\Sigma_1 = \Sigma_2 = \begin{bmatrix} \sigma & 0 \\ 0 & \sigma \end{bmatrix}$$



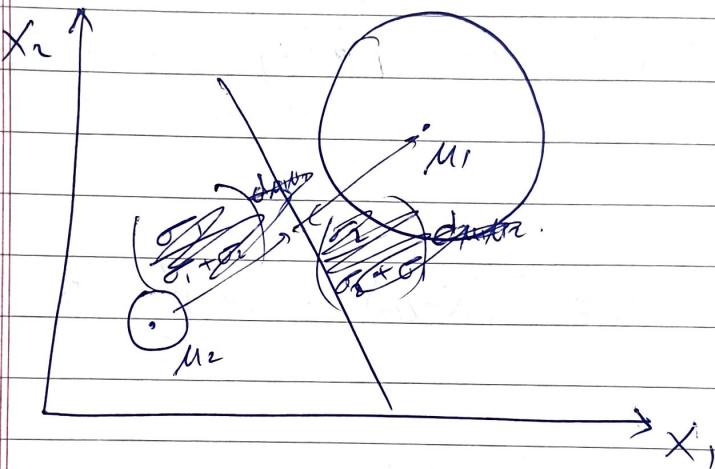
$$|x - \mu_1|^2 = |x - \mu_2|^2$$

Case 2 :

$$\Sigma_1 = \begin{bmatrix} \sigma_1 & 0 \\ 0 & \sigma_1 \end{bmatrix}$$

$$\Sigma_2 = \begin{bmatrix} \sigma_2 & 0 \\ 0 & \sigma_2 \end{bmatrix}$$

$\mu_1 \neq \mu_2$

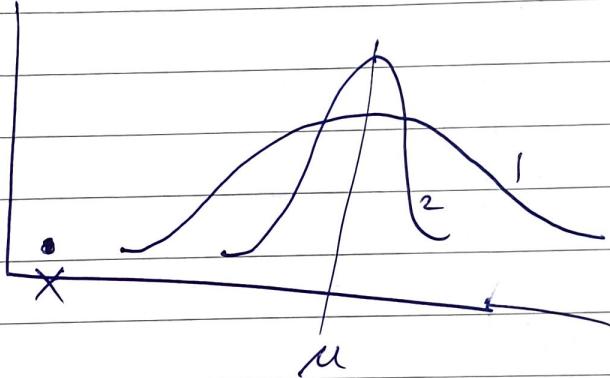
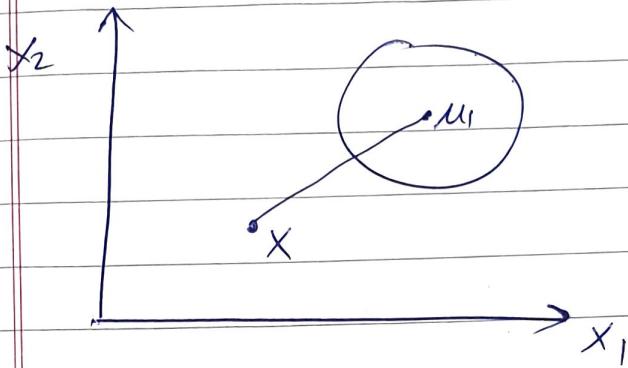


$$\left( \frac{X_1 - \mu_1}{\sigma_1} \right)^2 = \left( \frac{X_2 - \mu_2}{\sigma_2} \right)^2$$

Distance b/w points  $\rightarrow$

$$x_1 - \text{---} - x_2$$

Distance b/w a point and a distribution.



Both distributions have same  $\mu$  but we can see that 2 is closer.

so using only mean is not a good metric

so we use Mahalanobis distance -

$$\frac{x_1 - \mu_1}{\sigma_1}, \quad x_2 - \frac{x - \mu_2}{\sigma_2}$$

Case 3:

$$\Sigma_1 = \begin{bmatrix} \sigma_{11} & \sigma_{12} \\ \sigma_{21} & \sigma_{22} \end{bmatrix}$$

$$\Sigma_2 = \begin{bmatrix} \quad & \quad \end{bmatrix}$$

Quadratic discriminant

discriminant



22<sup>nd</sup> August 2019

Bayes Classifies.

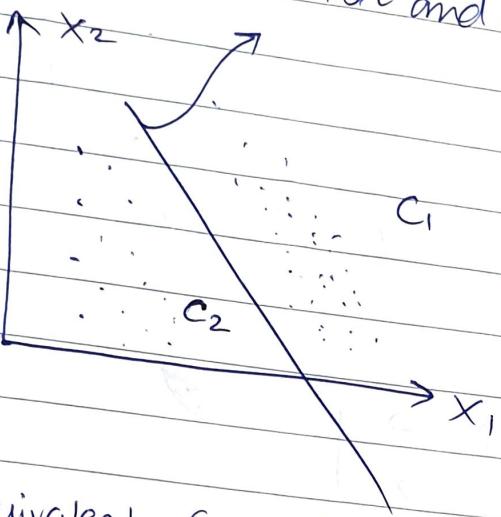
Maximum a posteriori classifier

if  $P(X|C_1) P(C_1) > P(X|C_2) P(C_2)$

→ class  $= C_1$  for  $X$

else  $= C_2$

Discrimination function and class boundary



Equivalent form of  $p(x|C_1)P(C_1) > p(x|C_2)P(C_2)$   
is

$$\frac{p(x|C_1)P(C_1)}{p(x|C_2)P(C_2)} > 1 \rightarrow C_1$$

else  $\rightarrow C_2$ .

Take log

$$\frac{p(x|C_1)}{p(x|C_2)} < \frac{P(C_2)}{P(C_1)}$$

$$\log(p(x|C_1)P(C_1)) - \log(p(x|C_2)P(C_2)) > 0$$

log likelihood ratio.

$$\mathbf{X} = [X_1, X_2]$$

$$P(C_1) = P(C_2) = 0.5$$

$$\boldsymbol{\mu}_1 = [\mu_{11}, \mu_{12}] = [3, 3]$$

$$\boldsymbol{\mu}_2 = [\mu_{21}, \mu_{22}] = [3, -3]$$

$$\boldsymbol{\Sigma}_1^* \boldsymbol{\Sigma}_1^T = \begin{bmatrix} 1/2 & 0 \\ 0 & 2 \end{bmatrix}$$

$$\boldsymbol{\Sigma}_2^* \boldsymbol{\Sigma}_2^T = \begin{bmatrix} 2 & 0 \\ 0 & 2 \end{bmatrix}$$

$$P(X_1 | C_1) \sim N(X_1, \boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1)$$

$$P(X_2 | C_2) \sim N(X_2, \boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2)$$

$$N(\mathbf{x}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{\sqrt{2\pi}} \boldsymbol{\Sigma}^{-1} e^{-\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu})}$$

$$(\boldsymbol{\Sigma}_1 \boldsymbol{\Sigma}_1^T)^{-1} = \begin{bmatrix} 2 & 0 \\ 0 & 1/2 \end{bmatrix}$$

$$= \frac{1}{\sqrt{2\pi} |\boldsymbol{\Sigma}_1|} e^{-\frac{1}{2} \mathbf{x}^T \boldsymbol{\Sigma}_1^{-1} \mathbf{x}}$$

$$= \frac{1}{\sqrt{2\pi} |\boldsymbol{\Sigma}_1|} e^{-\frac{1}{2} \begin{bmatrix} 2 & 0 \\ 0 & 1/2 \end{bmatrix} \begin{bmatrix} 1 & 2 \\ 2 & 1 \end{bmatrix}}$$

$$= \begin{bmatrix} x_1 - 3 & x_2 - 3 \end{bmatrix} \begin{bmatrix} 2 & 0 \\ 0 & 1/2 \end{bmatrix} \begin{bmatrix} x_1 - 3 \\ x_2 - 3 \end{bmatrix}$$

$$= \frac{1}{\sqrt{2\pi}} (1 \sum_i | = 1) e$$

$$= \left[ 2(x_1 - 3) \quad \frac{1}{2}(x_2 - 3) \right] \begin{bmatrix} x_1 - 3 \\ x_2 - 3 \end{bmatrix}$$

$$= \frac{1}{\sqrt{2\pi}} e^{-\left[ 2(x_1 - 3)^2 + \frac{1}{2}(x_2 - 3)^2 \right]}$$

$$= \frac{1}{\sqrt{2\pi}} e$$

$$= \frac{1}{\sqrt{2\pi}} e^{-\left[ 2(x_1 - 3)^2 + \frac{1}{2}(x_2 - 3)^2 \right]} \checkmark$$

On equating prob dist. of  $N_1$  and  $N_2$ .

$$(x - \mu_1)(\Sigma_1^{-1})(x_1 - \mu_1)^T + \log |\Sigma_1| =$$

$$(x - \mu_2)(\Sigma_2^{-1})(x_2 - \mu_2)^T + \log |\Sigma_2|$$

~~$$x^T (\Sigma_1^{-1} x_1 + \Sigma_2^{-1} x_2) + -\mu_1^T (\Sigma_1^{-1}) + -\mu_2^T (\Sigma_2^{-1})$$~~

$$w_1 x^2 + w_2 x + w_0 = 0$$

$$w_1?$$

$$w_2?$$

$$w_0?$$

$(x - \mu_1) \Sigma_i^{-1} (x - \mu_1)^T \rightarrow \text{mahalanobis distance between } x \text{ and } N(X, \mu_1, \Sigma_1)$

$$X = [x_1, x_2]$$

$$\mu_1 = [3, 4]$$

$$\mu_2 = [2, -4]$$

$$\Sigma_1 = \Sigma_2 = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

What is the discriminant function obtained by Bayes classifier?

We find the decision boundary so that we don't have to evaluate the probabilities for every test example.

Steps :-

- (1) Assume a distribution
- (2) Estimate parameters
- (3) find decision boundary.

✓  
30<sup>th</sup> August, 2019

classmate

Date \_\_\_\_\_

Page \_\_\_\_\_

## Classifier evaluation

- Parameter tuning
- Comparison
- defining a score function
- estimating a score function
- combining score function of two classifiers.

Define score function

Confusion matrix [with test set]

C-class problem

actual labels

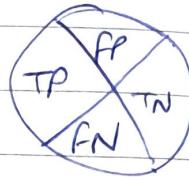
|             |  | C <sub>1</sub> | C <sub>2</sub> | .               | . | . | C <sub>c</sub> |
|-------------|--|----------------|----------------|-----------------|---|---|----------------|
|             |  | C <sub>1</sub> | C <sub>2</sub> | .               | . | . | C <sub>c</sub> |
|             |  | C <sub>1</sub> | C <sub>2</sub> | n <sub>ij</sub> |   |   |                |
| predictions |  | .              | .              |                 |   |   |                |
|             |  | .              | .              |                 |   |   |                |
|             |  | C <sub>1</sub> | C <sub>2</sub> |                 |   |   |                |
|             |  | C <sub>c</sub> |                |                 |   |   |                |

cost of ~~this~~ classifying an i-class example as a j-class example.

$$= b_{ij}$$

actual

|   | P  | N  |
|---|----|----|
| P | TP | FP |
| N | FN | TN |



$$\text{accuracy} = \frac{TP + TN}{TP + TN + FP + FN}$$

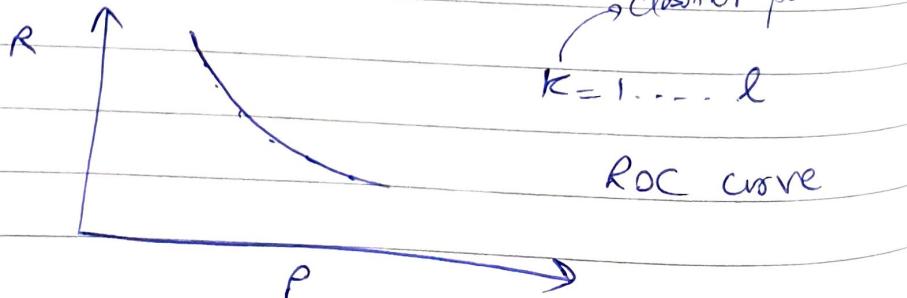
(specificity) Precision =  $\frac{TP}{TP + FP}$  (only care about correctly classifying true examples)

only getting in getting most <sup>actually true</sup> of our positive predictions to be true.

Recall (sensitivity) =  $\frac{TP}{TP + FN}$  most are classified as -ive

F-score  $\Rightarrow$  harmonic measure of precision, recall

$$\frac{2}{F} = \frac{1}{P} + \frac{1}{R}$$



Curve which has a larger area under the ROC curve is considered to be a better classifier.

$$\begin{aligned} P_{\text{correct}} &= E(\text{Accuracy}) \Big|_{|S| \rightarrow \infty} \\ &\approx \hat{\text{Accuracy}} \Big|_{|S|=n \leftarrow \infty} \end{aligned}$$

S - ~~test~~ test sample

How do we efficiently estimate  $E(\text{Accuracy})$  using a small test set?

Estimation techniques -

- bootstrap
- hold out
- $k$ -fold cross validation
- leave out one estimate.

bootstrapping  $a_i, \hat{a}_i$   
 Instead  $\hat{a}_1, \hat{a}_2, \hat{a}_3, \dots, \hat{a}_k$

$$(1 - \bar{a}_i) - \text{bias}$$

$$\sigma(a_i) - \text{variance}$$

- ① Bias is a better estimate of a than individual samples.

[ Wednesday ]

4<sup>th</sup> September 2019

CLASSMATE

Date \_\_\_\_\_

Page \_\_\_\_\_

## Classifier evaluation

- Score definition
- Estimating the score

### Score definition

|           |   | Confusion matrix |    |
|-----------|---|------------------|----|
|           |   | 1                | 0  |
| predicted | 1 | TP               | FP |
|           | 0 | FN               | TN |

$$\text{accuracy} = \frac{TP + TN}{n}$$

$n_1 + n_2 = n$

$n_1 \gg n_2 \Rightarrow \text{Recall} = \frac{TP}{n_1}$

$n_2 \gg n_1 \Rightarrow \text{Precision} = \frac{TP}{TP + FP}$

as num-training-examples denotes the prob. of a correct classification of a test data

~~✓~~ Accuracy tends to have a high-variance over separate test datasets, which is undesirable.

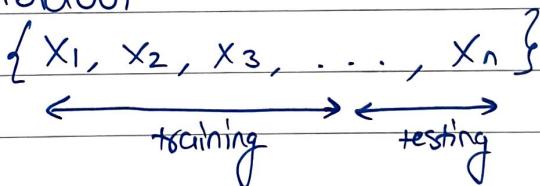
## Cross-validation

$$\begin{array}{lcl} \text{Experiment 1} & : & a_1 \\ 2 & : & a_2 \\ \vdots & : & \vdots \\ K & : & a_K \end{array} \quad \left. \begin{array}{c} \bar{a} \uparrow \\ \sigma(a) \downarrow \end{array} \right\}$$

Given a small training dataset, we would like to repeat our experiment several times to gauge a good estimate of the accuracy.

⇒ We resample our population  $K$ -times.

### - Holdout

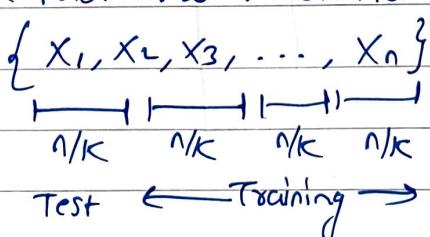


### - Bootstrapping

Repeats hold-out  $K$  times

Lots of random selections have to be made

### - $K$ -fold Cross Validation



Groups have to be made only once, no randomness involved

Every batch gets to be the test data once while the other batches act as training data.

$$\text{Error} = \text{bias} + \text{variance}$$

All classifiers fight the trade-off between accuracy and variance.

### Probabilistic graphical model

Bayesian network.

$x_1 \ x_2 \ x_3 \ x_4 \ y$

|   |   |   |   |   |
|---|---|---|---|---|
| 1 | 0 | 0 | 1 | 1 |
| 0 | 0 | 0 | 1 | 0 |

1 1 1 1 ?

$$P(y=1) = ?$$

$$P(y=0) = ?$$

Probability table

| $x_1$ | $x_2$ | $x_3$ | $x_4$ | $y$ | $P(.)$ |
|-------|-------|-------|-------|-----|--------|
| 1     | 1     | 1     | 1     | 1   | 0.2    |
| 1     | 1     | 1     | 1     | 0   | 0.3    |
| 0     | 0     | 0     | 0     | 1   | 0.1    |

If we can factor some possibilities, then we can save some calculations.

For

$$P(1, 1, 1, 1|1) = P(1|1) P(1|1) P(1|1) P(1|1)$$

instead, we could save

$$= P(1, 1|1) P(1, 1|1) \cancel{P(1|1)}$$

↳ only 2 factors have to be calculated and stored.