

11th August, 2017

Day 2

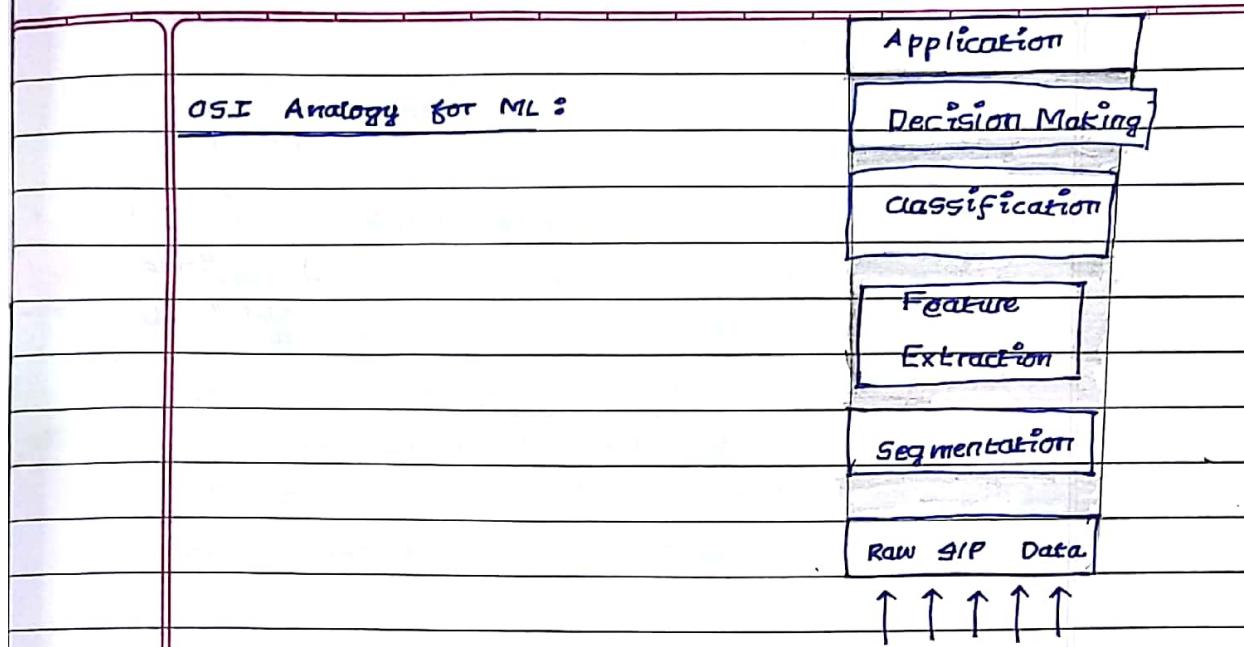
10:45 - 12:45

PAGE

DATE

DATE:

PAGE:



Critical state for Board 'B':

Row 1	*	Δ	Δ	Δ	Δ	Δ	Δ	Δ	Δ
Row 2	Δ	*	Δ	Δ	Δ	Δ	Δ	Δ	Δ
Row 3	*	*	*	*	*	*	*	*	*
Row 4	*	*	*	*	*	*	*	*	*
Row 5	*	*	*	*	*	*	*	*	*

QUESTION: Is the given board a critical state? If yes, then what is the minimum number of moves required to make it a non-critical state?

ANSWER: The given board is a critical state. It can be solved in 5 moves.

		P2			
		Δ	Δ	Δ	Δ
		Δ	Δ	Δ	Δ
					Δ
		(2, 0)	*	*	*
		*	*	*	*
		*	*	*	*

P1

We can have multiple successive state of the board after the first move, itself.

P1 could have chosen to move to anywhere else diagonally on 3rd ($i=2$) row.

Similarly, P2 could have moved to some other position

What is experience?

What we want to learn is

all such successive board state corresponding to each move such that as a player I win the game.

Teacher versus Learner's control over Experience

- Learning is done taking help of Teacher.
- Learner makes a move & gets feedback from Teacher.
- Sometimes, you play game against oneself & learn.

Direct Training Experience :

Assuming an ideal scenario,

You have data of all possible combination of successive board state per move, from we know at ALL TIMES

the move we need to make to WIN the game.

We don't have such a data in reality.

So training go towards and build strategy.

Indirect Training Experience: Training board state

chose optimal next move given a current board state, so that you MAY WIN the game. It's not for sure that you'll win the game.

Since the best move is not known, we can assign a value for every board state for move and it reduces to an optimization problem.

we can define a target function $V(b)$ for an arbitrary board state:

1. If board state b is final + won, $V(b) = +100$

If b is final + fail, $V(b) = -100$

Together provides us information about which b is a final state.

for the FIRST MOVE

itself, we could have

had another

"successive board state"

for FIRST MOVE itself.

where

b is a

final

state,

you

can

assign

such a

value.

similarly do Back Propagation

"if we know how final score is"

2. If final & draw, $V(b) = 0$

b' = base final score

looks like

a or

already

3. If b is not final, $V(b) = V(b')$

optimally. we can't

define

just

For

it

is

base

final

score

we can

recursively

not a

final

state

that

value

to plausible

we can

choose

base

states

we need to learn the fn approximately (i.e., \hat{V})

How to choose a representation for the fn approximation?

We have several possibilities:

we can choose linear,

Quadratic,

Polynomial,

ANN etc.

In case we choose Turner's π , we'll have to see how much does each of the following affect the chances of winning?

i) # of black points for current board state b. $\approx \chi_1$

ii) # " red " $\approx \chi_2$

iii) # " black " $\approx \chi_3$

that are at threat.

iv) same with red points $\approx \chi_4$

So, forming linear equation to minimize $V(b)$.

$\hat{V}(b)$ can be exp written as:

$$\hat{V}(b) = w_1 \chi_1 + w_2 \chi_2 + w_3 \chi_3 + w_4 \chi_4$$

How much does each of these board up to it

contribute depends on $\chi_1, \chi_2, \chi_3, \chi_4$ which

contribution is determined by turning

w_1, w_2, w_3, w_4 .

Intuitively, w_1 is weight of black points

and w_2 is weight of red points.

May be w_3 is

points which are at threat to black

and w_4 is points which are at threat to red.

Turner's weights are not necessarily unique.

Intuitively, w_1, w_2, w_3, w_4 are

weights assigned to each point based on its

status in board.

For example, a black point in front of a red

point has higher weight than a black point behind a red

point in front of a red point.

Similarly, a red point in front of a black

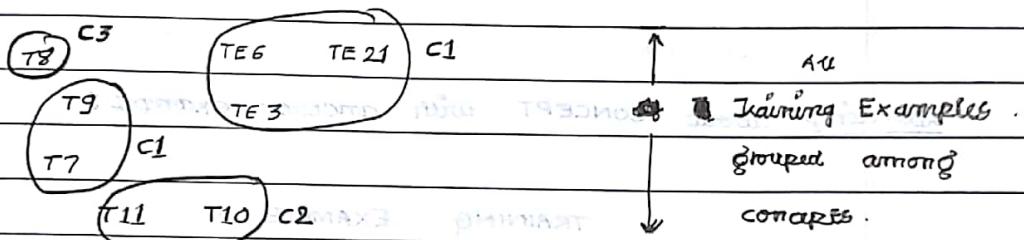
point has higher weight than a red point behind a black

point in front of a black point.

Concept learning

Ex 1.	beak	wings	legs	fly	Bird (concept)
1	Yes	Yes	Yes	Yes	Member of 1 (denoted by a boolean value)
2	Yes	Yes	Yes	No	Non-member of 1 (still a bird denoted by 1)
3	Yes	No	Yes	No	Non-member of 0 (not a " " " 0) of concept class.

$$X = \{ \text{Set of Training Examples} \}$$



Domain is set of all training example $\rightarrow f: \mathbb{R}_{m \times n} \rightarrow \text{Boolean}$
 Concept is a boolean valued function mathematically.

Otherwise, set of examples that can be grouped together,

forms a concept

Each instance has set of attributes.

Conjunction and uses of concept :

Ex 1 From above :

Mathematically,

The object has beak

x_1

AND wings conjunction

\wedge

x_2

AND legs

\wedge

x_3

AND can fly

\wedge

x_4

is a bird.

give 1

or

$(x_1, x_2, x_3, x_4, 1)$

Ex 2 From abv :

If the object has beak

Mathematically, x_1

AND has wings

 x_2

AND legs

 x_3

AND cannot fly

 $\neg x_4$

is a bird

gives 1

or

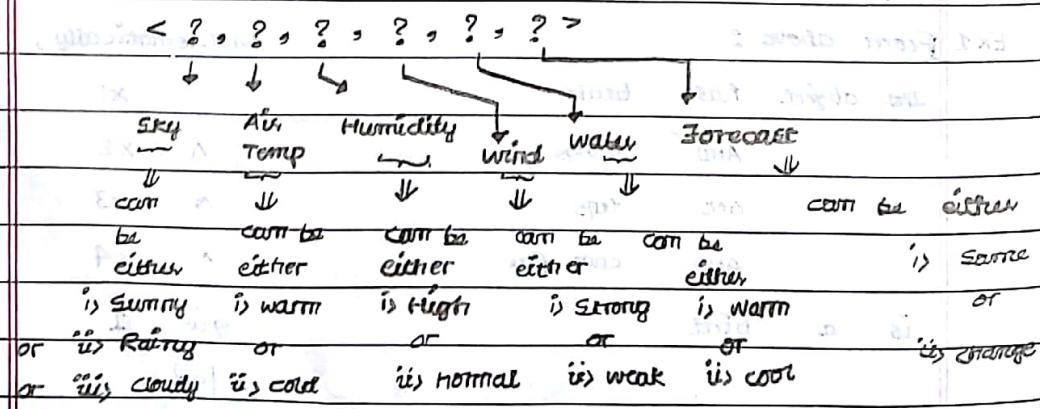
 $(x_1, x_2, x_3, \neg x_4, 1)$ Learning about CONCEPT with another example:

TRAINING EXAMPLE

Day	Sky	Air Temp	Humidity	Wind	Water	Forecast	Sport
1	Sunny	Warm	Normal	Strong	Warm	Same	Yes
2	"	"	High	"	"	"	"
3	Rainy	Cold	"	"	"	Change	No
4	Sunny	Warm	"	"	Cool	Same	Yes

whether the plays sport or not is concept.

Most general representation of any instance (or training example) is,



written for an instance, thus attribute (here, air temp)
DATE is missing PAGE:

$\langle ? = \phi, ? = \phi, ?, ?, ? \rangle$

Notations:

$X \rightarrow$ Set of all instances

$c(x) \rightarrow$ concept for a given concept.

$f_i(x) \rightarrow$ The hypothesis for a given instance. This may or
may not agree with $c(x)$

Distributions

can't be represented

Each instance is represented
as

by our chosen hypothesis.

$\langle x_1, x_2, x_3, x_4, x_5, x_6 \rangle$

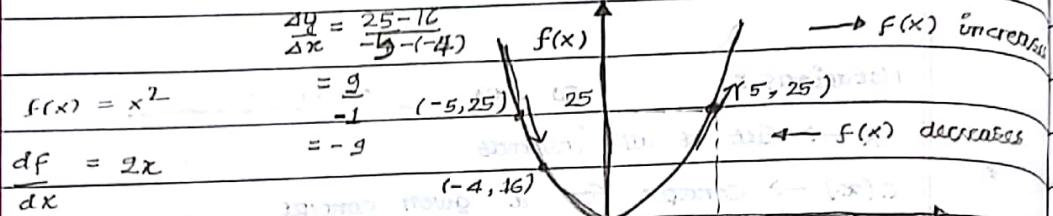
where

it means

if this instance has x_1 and x_2 and x_3 and
 x_4 and x_5 and x_6

then it may or may not belong to the
concept.

Inductive Hypothesis Learning:



i) @ $x = 5$

$$\frac{df}{dx} = 2(5) = 10$$

$$\frac{dx}{dt}$$

ii) @ $x = -5$

$$\frac{df}{dx} = 2(-5) = -10$$

$$\frac{dx}{dt}$$

For $x = 5$,

we should move in the direction of the slope
direction opposite to that of slope.

For $x = -5$,

we should move in the direction of the slope

$$(1, 1)$$

$$\text{as } x \uparrow y \uparrow$$

$$\frac{dy}{dx} = 1$$

$$(1, 1) \quad 1+1 = 2, \quad y = 2$$

Slope always
points in the
direction on

which function value
increases.

$$f(x) \quad y = -x$$

$$1, -1$$

$$2, -2$$

$$\text{as } x \uparrow y \downarrow$$

$$\frac{dy}{dx} = -1$$

$$2-1 = 1$$

$$y = -1$$

$$3-1 = 2$$

$$y = -2$$

$$\text{Earlier, } y = -3$$

$$\text{now } y = -1$$

Hypotheses,Hypothesis for a given instance,

concept " " " "

Suppose,

my hypothesis $f_1 = \langle \text{Sunny}, ?, ?, ?, \text{Strong}, ?, ?, ? \rangle$

if any instance such as

 $x_1 = \langle \text{Sunny}, \text{Warm}, \text{Normal}, \text{Strong}, \text{Warm}, \text{Same} \rangle$ $x_2 = \langle \text{Sunny}, \text{Warm}, \text{High}, \text{Strong}, \text{Cold}, \text{Same} \rangle$ which satisfies the constraints of the hypothesis f_1

irrespective of their concept value

will satisfy the hypothesis f_1 .

It means its possible that

$$\text{For concept } C, \quad c(x_1) = c(\langle \text{Sunny}, \text{Warm}, \text{Normal}, \text{Strong}, \text{Warm}, \text{Same} \rangle)$$

$$= \langle \text{Sunny}, \text{Warm}, \text{Normal}, \text{Strong}, \text{Warm}, \text{Same} \rangle, 1 \rangle$$

x₁ is considered

as a positive example

where as

$$\& \quad c(x_2) = c(\langle \text{Sunny}, \text{Warm}, \text{High}, \text{Strong}, \text{Cold}, \text{Same} \rangle)$$

$$= \langle \text{Sunny}, \text{Warm}, \text{Normal}, \text{Strong}, \text{Warm}, \text{Same}, 0 \rangle$$

x₂ is considered as a -ve example.

where

$$C \circ \text{Plays Game} : X \rightarrow \{0, 1\}$$

but both x₁ & x₂ falls under or

satisfy the constraints of

If we do $2 - (-1) = 3$ hypothesis f_1 at $x = 3$,

and

$$y = -3 \quad \therefore f_1(x_1) = 1 \quad [\text{though, } c(x_1) = 1]$$

f₁ value decreases

&

$$f_1(x_2) = 1 \quad [\text{", } c(x_2) = 0]$$

Calculating # of distinct instances?

We can classify instances based on semantic and syntax.

Syntactically,

	Sky	Air Temp	Humidity	Wind	Water	Forecast
corn	Sunny	warm	Normal	Strong	warm	Same
rainy day	Cloudy	cold	High	weak	High	change
of these	Rainy	φ	φ	φ	φ	φ
values	φ	?	?	?	?	?
	?					

Syntactically, # of distinct instances that we can get

$$= (3+2) \times (2+2) \times (2+2) \times (2+2) \times (2+2) \times (2+2)$$

$$= 5 \times 4 \times 4 \times 4 \times 4 \times 4$$

$$= 5 \times 2^2 \times 2^2 \times 2^2 \times 2^2 \times 2^2$$

$$= 5 \times 2^{10}$$

$$= 5 \times 1024$$

= 5120 syntactically distinct instances

Semantically,

we assume any instance having

at least one φ is considered

same semantically.

Semantic could be,

if any instance has even one

φ (that is even if one of the

attribute value for that

instance is missing)

we'll consider all such instance
as redundancy.

For any such instance W_i which is redundant
(as per our defined semantic),

we do not worry about
whether it satisfies any concept C
or if " " constraint of any hypothesis

Here we defined 1 semantic for all
instances having atleast one ϕ (though
all such instances are semantically different) - ①

Apart from this,

Apart from this,

	Sky	Air Temp	Humidity	Wind	Water	Forecast
can take	Sunny	warm	Nominal	Strong	Warm	Sunny
any of	cloudy	cold	High	Weak	High	change
these values	Rainy	?	?	?	?	?
&	?					

each
will
be

semantically
different

Semantically free, # of distinct instances

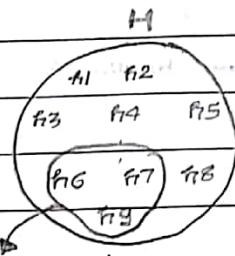
$$\begin{aligned}
 & \text{Semantically} \\
 & \text{distinct} \\
 & \text{instances} \\
 & \text{need not} \\
 & \text{satisfy} \\
 & \text{some} \\
 & \text{concept or hypothesis}
 \end{aligned}
 = (3+1) \times (2+1) \times (2+1) \times (2+1) \times (2+1) \\
 = 4 \times 3 \times 3 \times 3 \times 3 \\
 = 4 \times 81 \times 3 \\
 = 4 \times 243 \\
 = 972 - ②$$

Using ① & ②, total # of semantically distinct

$$\text{instances} = 1 + 972$$

$$= 973 \text{ semantically}$$

distinct instances.

Version SpaceVersion

space = set of hypotheses

which

satisfies

all the training

examples in D

see D,

such that

$$h(x) = c(x)$$

Candidate Elimination Algorithm

1. Japan Honda Blue 1980 Economy tree

Starting Point :

So, Most Specific $\emptyset \quad \emptyset \quad \emptyset \quad \emptyset \quad \emptyset$

6.0, Most General ? ? ? ? ?

General

For ex 1 : Since it is a +ve example, start with Generic Side

S0	∅	∅	∅	∅	∅	
S1	Japan	Honda	Blue	1980	Economy	- (2) New find more general compared to S0, specific hypothesis S1.

G1, G0 ∵ G0 satisfies ex 1
retain it as G1 - (1)

2. Japan Toyota Blue 1970 Sports +ve

Since ex2 is -ve, start with specific side.

S0	∅	∅	∅	∅	∅	
S1	Japan	Honda	Blue	1980	Economy	- (1) ∵ S1 designates ex2 as -ve
S2	"	"	"	"	"	∴ Retain it as S2

G2 gives ex2 as +ve
but in fact it is -ve
∴ we need to ~~accept~~ get a new G2. - (2)

Getting G2 :

$\langle I, ?, ?, ?, ?, ?, ? \rangle$	\times	Not satisfied by: ex1 & ex2
$\langle ?, H, ?, ?, ?, ? \rangle$	\checkmark	
$\langle ?, ?, B, ?, ?, ? \rangle$	\checkmark	
$\langle ?, ?, ?, 1980, ? \rangle$	\checkmark	
$\langle ?, ?, ?, ?, ?, Eco \rangle$	\checkmark	

$$\therefore G2 = \{ \langle ?, H, ?, ?, ?, ? \rangle, \\ \langle ?, ?, B, ?, ?, ? \rangle, \\ \langle ?, ?, ?, 1980, ? \rangle, \\ \langle ?, ?, ?, ?, ?, Eco \rangle \}$$

3. Japan, Toyota, Blue, 1990, Economy +ve.
 \because It is +ve,
start from general side.

Check if ex3 satisfies G2 - ①

$$G2 = \{ \langle ?, H, ?, ?, ?, ? \rangle \times \\ \langle ?, ?, B, ?, ?, ? \rangle \checkmark \\ \langle ?, ?, ?, 1980, ? \rangle \times \\ \langle ?, ?, ?, ?, ?, Eco \rangle \checkmark \}$$

$$\therefore G3 = \{ \langle ?, ?, B, ?, ?, ? \rangle, \\ \langle ?, ?, ?, ?, ?, Eco \rangle \}$$

② Going to specific side we get,

$$S3 = \langle Japan, ?, Blue, ?, ?, Economy \rangle$$

4. USA, Chrysler, Red, 1980, Economy -ve



ex2 $t(x) \neq c(x)$

for both

ex1 & ex2

\downarrow

$t(x)$

$= c(x)$

S3: < Japan, ?, Blue, ?, Economy >

Start

with

specific

side.

\uparrow

$t(x)$

$\neq c(x)$

Anomalous

ex4 as -ve

$\therefore S3$ is retained as $S4$

S4: < J, ?, B, ?, E > - (1)

Gathering G4:

< ?, ?, B, ?, ?, ? > satisfied by ex1, 2 & 3

G3

also

satisfies ex4

< ?, ?, ?, ?, Eco > *

satisfied by ex1, ex2, ex3

but

not by ex4.

\therefore Find a more specific hypothesis than G3 (i.e., G4)

that satisfies

ex1, ex2, ex3, ex4.

we'll revisit from ex1 to ex4

& found out that

< J, ?, ?, ?, ?, Eco >

satisfies ex1 to ex4

$\therefore G4 = \{ < ?, ?, B, ?, ?, ? >$

$< J, ?, ?, ?, ?, Eco > \} - (2)$

5.

Japan Honda White 1980 Economy +ve



Do general side

 $G_4 = \{ < ? ? B ? ? ? >$
 $< T ? ? ? Eco > \}$
Get them
specific side
 $\textcircled{2} \times 5$ does not satisfy $< ? ? B ? ? ? >$
 but satisfies $< T ? ? ? Eco >$

∴ Update

 G_4 os $G_5 = \{ < T ? ? ? Eco > \}$ Getting S_5 : $S_4 = \{ < T ? B ? E >$

becomes

 $S_5 = \{ < T ? ? ? E >$

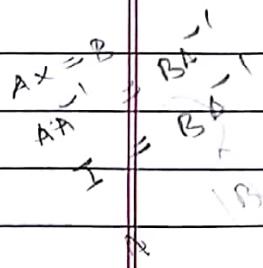
We stop when

$$S_i = G_i$$

or

when we run out of

kraining example



Imp for Exam: Rank, Singular Matrix, Eigenvalue /vector,
Structure of Matrix, Cofactor, determinant
Final Version Space we get for DATE: PAGE:

Enjoy Spore Example :

$$S = \{ \text{Sunny, Warm, ?, Strong, ?, ?} \}^{f_1}$$

f₂

f₃

f₄

$$\langle \text{Sunny, ?, ?, Strong, ?, ?} \rangle \quad \langle \text{Sunny, Warm, ?, ?, ?} \rangle \quad \langle \text{? W ? S ? ?} \rangle$$

f₅

f₆

$$G = \{ \langle S ? ? ? ? ? \rangle \quad \langle ? W ? ? ? ? \rangle \}$$

There are 6 Hypothesis in the ~~Hypotheses~~ Space final
version space.

True Data :

A	Sunny	Warm	Normal	Strong	Cool	Change
B	Rainy	Cold	"	Light	Warm	Same
C	Sunny	Warm	Normal	Light	Warm	Same
D	Sunny	Cold	Normal	Strong	Warm	Same

	f ₁	f ₂	f ₃	f ₄	f ₅	f ₆
A	✓	✓	✓	✓	✓	✓
B	✗	✗	✗	✗	✗	✗
C	✗	✗	✗	✗	✓	✓
D	✗	✓	✗	✗	-	-

$x_1 \ x_2 \ x_3 \ \bar{x}_4 \ x_5 \ \bar{x}_6$

↓

Rule is conjunction of attributes

So : $(x_1 \vee x_2 \vee x_3 \vee x_5) \rightarrow \text{play}$

↓
disjunction of conjunction

GO : $(\neg(x_4 \vee x_5)) \wedge$

Algorithm

Examples :

Enjoy Sport Example :Training Examples

Day	Sky	Air Temp	Humidity	Wind	Water	Forecast	Enjoy Sport
1	Sunny	warm	Normal	Strong	Warm	Same	Yes (+ve)
2	Sunny	warm	High	Strong	Normal	Same	Yes (+ve)
3	Rainy	cold	High	Strong	Normal	Change	No (-ve)
4	Sunny	warm	High	Strong	Cold	Change	Yes (+ve)

∴

Applying Candidate Elimination,

$$\begin{aligned} G_0(x) &= \phi \quad \phi \quad \phi \quad \phi \quad \phi \quad \phi \\ G_0(x) &= ? \quad ? \quad ? \quad ? \quad ? \quad ? \end{aligned}$$

↑ Specific
↓ General

$\leftarrow 1, C(1) \leftarrow \text{Sunny, Warm, Normal, Strong, Warm, Same}$
 $, \text{True / Yes / +ve / 1} \rightarrow$

Considering

Going from General to Specific,

 $\leftarrow 2, C(2) \leftarrow$ $\equiv \leftarrow \leftarrow \text{SUN}$ Hypothesis $G_0(1) = C(1)$

$$1 = 1$$

$$\therefore G_1 = G_0$$

Going from Specific to General.

$$S_0(1) \neq C(1)$$

∴ we need to get a more general hypothesis

S_1 such that

$$S_1(1) = C(1)$$

$$S_1 : < \text{Sunny} \quad \text{warm} \quad \text{normal} \quad \text{strong} \quad \text{warm} \quad \text{Same} >$$

Fancy Sportz

$$\text{Now, } S_1(1) = C(1)$$

Yes (+ve)

yes (+ve)

No (-ve)

Yes (+ve)

∴ For ex1,

$$S_0 : < \emptyset \emptyset \emptyset \emptyset \emptyset \emptyset >$$

$$S_1 : < \text{s} \quad \text{w} \quad \text{n} \quad \text{s} \quad \text{w} \quad \text{s} >$$

∴

$$G(1) : < ? \quad ? \quad ? \quad ? \quad ? \quad ? >$$

$$G(0) : < ? \quad ? \quad ? \quad ? \quad ? \quad ? >$$

Considering ex 2,

$$< 2, C(2) >$$

▲

$$= < \text{sunny} \quad \text{warm} \quad \text{high} \quad \text{strong} \quad \text{warm} \quad \text{same} >$$

,
True / Yes / +ve / 1

Since it's a +ve example,

Start from most general to specific,

$G_1 < ? \quad ? \quad ? \quad ? \quad ? \quad ? >$

$$G_1(?) = C(2)$$

$$\therefore G_2 = G_1$$

$G_2 < ? \quad ? \quad ? \quad ? \quad ? \quad ? >$

From specific to general,

$S_1 < \text{sunny} \quad \text{warm} \quad \text{normal} \quad \text{strong} \quad \text{warm} \quad \text{same} >$

ex2 : $< \text{sunny} \quad \text{warm} \quad \text{high} \quad \text{strong} \quad \text{warm} \quad \text{same} >$

$$S_1(2) \neq C(2)$$

\therefore Get a more general hypothesis

such that it is consistent

with both ex1 and ex2



And we know that

S_1 is already consistent

with ex1

\therefore Comparing S_1 & ex2 attribute wise

we get S_2 ,

$S_2 : < \text{sunny} \quad \text{warm} \quad ? \quad \text{strong} \quad \text{warm} \quad \text{same} >$

From ex1 & ex7 we get,

S₀ <∅ ∅ ∅ ∅ ∅ ∅ >

S₁ <S W N S W S>

S₂ <S W ? S W S>

General



G₂ <? ? ? ? ? ? >

G₁ <? ? ? ? ? ? >

G₀ <? ? ? ? ? ? >

Specific



Considering ex3,

ex3 : Rainy cold High Strong warm change -ve

<3, C(3) >

< Rainy cold High Strong warm change ,

False/No/-ve /0 >

Since it's a -ve example,

we'll start from specific to General,

S₂ and ex3 are compared w.r.t attribute wise

? Rainy ≠ Rainy

∴ we get 0

and ∵ we have conjunctions

0 ∨ anything (0 or 1) ∨ anything ∨ anything ∨ anything

∨ anything

= 0

i.e., S₂ designates ex3 as -ve example $S_2(3) = 0$

and so is

$C(3) = 0$

∴ $S_2(3) = C(3) = 0$

∴ $S_3 = S_2$

S₃ < sunny warm ? strong warm same >

From general to specific,

$G_2 \vdash ? ? ? ? ? ? ? >$

ex 3 \langle Rainy cold High strong warm change \rangle

$$c(3) = 0$$

$$G_2(3) = 1$$

$$G_2(3) \neq c(3)$$

i.e., G_2 is not consistent with ex 3

∴ we'll have to find a more
more specific (than G_2) hypothesis (G_3)

which is consistent with

ex 1, 2 & 3.

i.e.,

$$G(1) = c(1) \quad [= 1]$$

and

$$G(2) = c(2) \quad [= 1]$$

and

$$G(3) = c(3) \quad [= 0]$$

$$G_3 \vdash \langle \text{sunny } ? ? ? ? ? ? >$$

So we get :

$$G_3 \vdash \langle \text{sunny warm ? Strong warm same} >$$

$$G_3 \vdash \{ \langle \text{sunny ? ? ? ? ? ? ? ? >} , \\ \langle ? \text{ warm ? ? ? ? ? ? same} > , \\ \langle ? ? ? ? ? ? same > \}$$

Both S_3 & G_3 are consistent with ex 1, 2 & 3.

ex 4: <sunny warm high strong cold change> DATE: PAGE:

<4, c(4)> need to generate a context for sentence A

<sunny warm high strong cold change, 1>

Since it's a +ve example,

we start from most general to specific.

q5 ex4 consistent with G3 ? but B consistent and A -ve

which means there are no contradictions between them

? = {and could G3} → inconsistent

<sunny ? ? ? ? ?> <? warm ? ? ? ?> <? ? ? ? ? same>

consistent

consistent

not consistent

so I check numbers → more of box A +ve

so no contradictions in the code

∴ G4 = { <sunny ? ? ? ? ?> } (finalizing)

<? warm ? ? ? ?> }

q5 ex4 consistent with G3 ?

<Sunny Warm ? ? ? ?>

<Sunny ? ? Strong ? ? ?>

<? Warm ? Strong ? ? ?>

S3

<sunny warm ? strong warm same>

so and could be No + consistent & and +

so possibility & could be a +ve + consistent

S4 = <sunny warm ? strong ? ?>

Final Version Space, V = { } + 4 top

S4 = <sunny warm ? strong ? ?> (finalizing)

<sunny warm ? ? ? ?> <sunny ? ? strong ? ?> <? warm ? strong ? ?>

G4 = <sunny ? ? ? ? ?> <? warm ? ? ? ?>

1. A coin is tossed 4 times. $P(\text{exactly 1 head}) = ?$

$$\text{Sol. } P(HHTT) + P(THHT) + P(TTHH) + P(TTTH)$$

$$= \left(\frac{1}{2}\right)\left(\frac{1}{2}\right)\left(\frac{1}{2}\right)\left(\frac{1}{2}\right) \times 4$$

$$= \frac{1}{4}$$

2. A bag contains 6 Red & 4 Black balls. 2 balls are drawn at random one after the other without replacement. $P(2 \text{ black balls}) = ?$

$$\text{Sol. } \frac{4 \times 3}{10 \times 9} = \frac{2 \times 1}{5 \times 3} = \frac{2}{15}$$

3. A card is drawn randomly from deck. I won Rs 1000 if it is a diamond or an ace.

$$P(\text{winning}) = ?$$

$$\text{Sol. } P(\text{winning}) = \frac{13+3}{52} = \frac{16}{52} = \frac{4}{13}$$

$$\frac{13}{52}$$

4. Bag 1 contains 4 white & 6 black balls. Bag 2 contains 4 white & 3 black balls. $P(\text{picking black ball from bag 1})$

$$\text{Sol. } P = \frac{1}{2} \left(\frac{6}{10}\right) = \frac{3}{10}$$

$$P(\text{picking ball from bag 1}) =$$

$$P(\text{choosing bag 1} \cap \text{black ball}) = \frac{1}{2} \left(\frac{6}{10} \right)$$

$$P(\text{'' bag 2} \cap \text{black ball}) = \frac{1}{2} \left(\frac{3}{7} \right)$$

$$P(\text{black ball}) = P(\text{choose bag 1}) P(\text{picking black ball} | \text{choose bag 1})$$

$$= P(\text{choose bag 2}) P(\text{picking black ball} | \text{choose bag 2})$$

$$= \frac{1}{2} \left(\frac{6}{10} \right) + \frac{1}{2} \left(\frac{3}{7} \right)$$

$$= \frac{1}{2} \left[\frac{6}{10} + \frac{3}{7} \right] = \frac{1}{2} \left[\frac{3}{5} + \frac{3}{7} \right]$$

$$= \frac{3}{2} \left[\frac{12}{35} \right] = \frac{3 \cdot 6}{35} = \frac{18}{35}$$

$P(\text{picking black ball from bag 1})$

$$= \left(\frac{1}{2} \right) \left(\frac{6}{10} \right) = \left(\frac{3}{10} \right) \div \left(\frac{18}{35} \right)$$

$P(\text{black ball})$

$$= \frac{3}{10} \times \frac{35}{18} = \frac{3 \times 7}{2 \times 18} = \frac{3 \times 7}{2 \times 6} = \frac{7}{12}$$

5. Tomorrow is Ramash's wedding in Rajasthan.

Weather forecast predicted that it rains 5 days a year.

It rained 5 days per year in the last decade.

Prediction is not perfect.

Whenever it rains, it perfectly predicts 90% of the time.

" it doesn't rain for " " 10% " " , ?

What is the probability it'll rain tomorrow?

6. Scores in MATH contest follows Normal Distribution.

Arvind scores 75 / 100.

Avg. score for the test is 60.

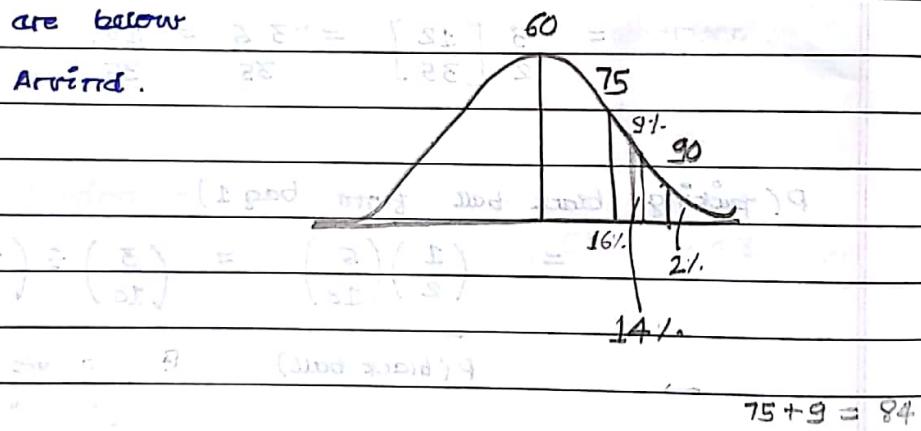
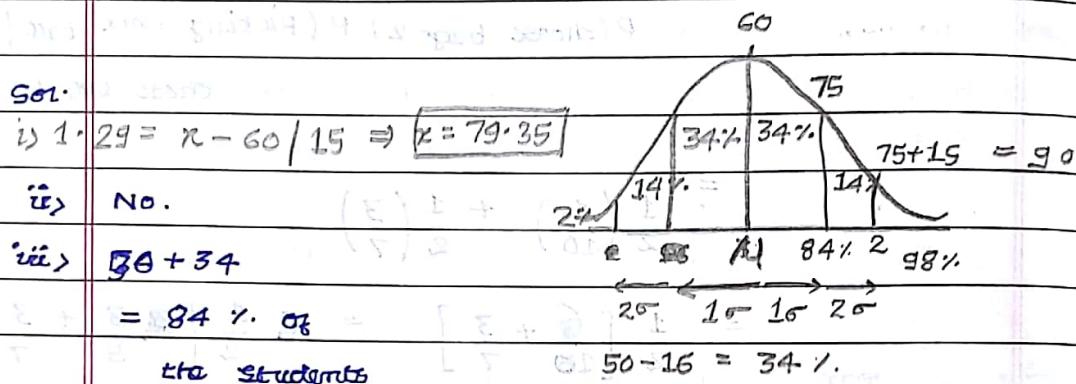
Std dev = 15.

Top 10% of the students are eligible for scholarship.

i) Least score to qualify for scholarship?

ii) Will Arvind get scholarship?

iii) How many students are below Arvind?



$$75 + 9 = 84 \quad 84 \times 2 = 168 \quad 4 \times 2 = 8$$

$$84 + 16 = 100 \quad 100 - 84 = 16 \quad 16 - 8 = 8$$

Top 10% of the students are eligible for scholarship.
Students with less than 84% are not eligible for scholarship.
Arvind's score is 75, which is less than 84%.
Therefore, Arvind is not eligible for scholarship.

Decision Trees

Day	Outlook	Temperature	Humidity	Wind	Play	Tennis
D1	Sunny	Hot	High	Weak	No	
D2	Sunny	Hot	High	Strong	No	
D3	Overcast	Hot	High	Weak	Yes	
D4	Rain	Mild	High	Weak	Yes	
D5	Rain	Cool	Normal	Weak	Yes	
D6	Overcast	Cool	Normal	Strong	No	
D7	Sunny	Cool	Normal	Strong	Yes	
D8	Sunny	Mild	High	Weak	No	
D9	Rain	Cool	Normal	Weak	Yes	
D10	Rain	Mild	Normal	Weak	Yes	
D11	Sunny	Mild	Normal	Strong	Yes	
D12	Overcast	Mild	High	Strong	Yes	
D13	Overcast	Hot	Normal	Weak	Yes	
D14	Rain	Mild	High	Strong	No	

NIR =

MR = $\frac{5}{14}$

- Decision Trees works with noisy training data as well.

$$E(s) = - \left[\frac{9}{14} \log \left(\frac{9}{14} \right) + \frac{5}{14} \log \left(\frac{5}{14} \right) \right]$$

$$\log_a$$

$$= \frac{\log a}{10}$$

$$\log_b$$

$$= - \left[\frac{9}{14} (-0.63743) \right]$$

$$+ \frac{5}{14} (-1.48543)$$

$$\log_2 a$$

$$= \frac{\log a}{\log 2}$$

$$\log_{10}$$

$$= -[-0.40977 - 0.53651]$$

$$= -[-0.9403]$$

$$= 0.94$$

Given feature space V for Enjoy Sport Example:

$\text{< ? sunny ? warm ? strong ? ? } \rightarrow \text{?}$

$\text{< ? sunny ? ? strong ? ? } \rightarrow \text{< ? sunny warm ? ? ? } \rightarrow \text{< ? warm ? strong ? ? }$

$\text{< ? sunny ? ? ? ? ? } \rightarrow \text{< ? warm ? ? ? ? }$

classify the TEST Data in

Instance	Key	AirTemp	Humidity	Wind	Water	Score	Enjoy Sport
A	Sunny	Warm	Normal	Strong	Cool	Chang	?
B	Rainy	Cold	Normal	Light	Warm	Same	?
C	Sunny	Normal	Normal	Light	Warm	Same	?
D	Sunny	Cold	Normal	Strong	Warm	Same	?

Sol - Instance A :

All 6 hypotheses give +ve. \therefore classification is +ve with 100% confidence.

Instance B :

All 6 hypotheses give -ve. \therefore " -ve " 100%

Instance C :

3 give -ve.

3 give +ve.

This is the case of ambiguity.

Our model (or learner) cannot classify this example with confidence.

NOTE : If this example is were available as

training data (i.e., if we knew the label of this instance)

then it would have been the most optimal training example

at the
rate
of $\log_2 ||vs||$

we can
reduce
the VS
by half w/

em
g

Meet
inse
com

Since
we

all
the
sen
clo

to reduce the current version space size
& to come closer to the target fit.

at the rate of $\log_2 ||VS||$ This is to be expected, because those instances whose classification is most ambiguous are

\Rightarrow precisely the instances whose true classification we can reduce would provide the most new information for the VS refining the version space.

by trial with such a training example (Here it would have been

Instance D : from 6 hypotheses in VS to
4 give -ve.
3 " " VS)

D gives +ve.

if Entity sport Here we can classify it as -ve with $\frac{4}{6} \times 100$

ge ?

?

o ?

o ?

$$= \frac{2}{3} \times 100$$

$$= 66.66\%$$

confidence

with 100% confidence

" 100% "

Method II :

Instance A :

compare it with the most specific hypothesis.

< sunny warm ? strong ? ? >

Since this classifies instance A as +ve

we need not check with other hypothesis, as

all other hypothesis are

more general than

this hypothesis. By definition of more general than (& by common

sense) we can say that all other hypothesis will also classify it as +ve.

Instance B

Comparing it with most specific hypothesis,

its classification is -ve.

We can't say about the true classification of instance B just because most specific hypothesis gave classification is -ve.

What we can do is

compare instance B with hypotheses in General set.

If ~~General set~~ ~~contains~~ all hypotheses " " " "

classifies instance B as -ve,

we can say for sure that the true classification of instance B is -ve.

Why?

Because every other hypothesis is more specific than this hypothesis. So if these say -ve, all others will surely classify it as -ve.

True classification of instance B is -ve.

with confidence 100%.

Instance C

Most specific hypothesis classifies it as -ve

Hypotheses in general set " " " " +ve & -ve

Other three hypotheses " " " " -ve, +ve & -ve

No classification can be inferred.

Instance D

Most specific hypothesis classifies it as -ve

Hypotheses in general set " " " " +ve & -ve

Other three hypotheses " " " " +ve, -ve & -ve

True classification is -ve with $4 \times 100 = 66.66\%$.

Entropy characterizes the (im) purity of an arbitrary collection of examples.

It is in the collection of samples there,

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

where P_+ is proportion of all +ve examples in S

In cases:

$$1^{\circ} \text{ All +ve examples } \Rightarrow P_{\oplus} = 1 \\ P_{\ominus} = 0$$

Trust,

$$\begin{aligned}
 \text{Entropy } (S) &= - [P_{\oplus} \log P_{\oplus} + P_{\ominus} \log P_{\ominus}] \\
 &= - [1 \log 1 + 0] \\
 &= - [1(0) + 0] \\
 &= -0 \\
 &= 0
 \end{aligned}$$

$$2. \text{ All -ve examples } \Rightarrow P_{\Theta} = 0 \\ P_{\bar{\Theta}} = 1$$

Turri,

$$\begin{aligned}
 \text{Entropy } (S) &= - [0 + 1 \log 1] \\
 &= - [0 + 0] \\
 &= -0 \\
 &= 0
 \end{aligned}$$

~~Equal~~ # of +ve & -ve examples :

$$P_{\oplus} = 0.5, \quad P_{\ominus} = 0.5$$

$$F_{\text{ENTROPY}}(S) = - [P_{\oplus} \log P_{\oplus} + P_{\ominus} \log P_{\ominus}]$$

$$= - [0.5 \log_2\left(\frac{1}{2}\right) + 0.5 \log_2\left(\frac{1}{2}\right)]$$

$$= -2 \times 0.5 \log_2 \left(\frac{1}{2} \right)$$

$$= -\frac{1}{2} [\log_2 1 - \log_2 2] = +\frac{\log_2 2}{2} = 1$$

4. Unequal # of +ve & -ve examples?
 $0 < E(S) < 1$

Information Gain (Sample S, for attribute A) :

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

$$\text{Gain}(S, \text{outlook}) = 0.94 - \left[\frac{|S_{\text{sunny}}|}{|S|} \times E(S_{\text{sunny}}) \right]$$

$$+ \left[\frac{|S_{\text{overcast}}|}{|S|} \times E(S_{\text{overcast}}) \right]$$

$$+ \left[\frac{|S_{\text{rainy}}|}{|S|} \times E(S_{\text{rainy}}) \right]$$

$$= 0.94 - \left[\frac{5}{14} \times - \left[\frac{3}{5} \log_2 \left(\frac{3}{5} \right) + \frac{2}{5} \log_2 \left(\frac{2}{5} \right) \right] \right]$$

$$+ \left[\frac{4}{14} \times - \left[\frac{4}{4} \log_2 \left(\frac{4}{4} \right) + \frac{0}{4} \log_2 \left(\frac{0}{4} \right) \right] \right]$$

$$+ \left[\frac{5}{14} \times - \left[\frac{3}{5} \log_2 \left(\frac{3}{5} \right) + \frac{2}{5} \log_2 \left(\frac{2}{5} \right) \right] \right]$$

$$= 0.94 - \left[\frac{-5}{14} (-0.44218 - 0.528771) \right]$$

$$= 0.94 - \left[-0.44218 - 0.528771 \right]$$

$$= 0.94 - \left[\frac{-5}{14} (-0.97095) \times 2 \right]$$

$$= 0.94 - [0.34677 \times 2]$$

$$= 0.94 - 0.6935365$$

$$= 0.24646$$

$$\text{Gain}(S, \text{outlook}) \approx 0.247$$

where,

$$E(S_{\text{suny}}) = -0.97$$

$$E(S_{\text{overcast}}) = 0$$

$$E(S_{\text{rain}}) = -0.97$$

)

To compute Gain ($S : \text{temp}$) we need,

$$E(S_{\text{hot}}) = E[2+, 2-] = - \left[\frac{2}{4} \log_2 \left(\frac{2}{4} \right) + \frac{2}{4} \log_2 \left(\frac{1}{2} \right) \right]$$

$$= - \left[2 \times \frac{1}{2} \times \left(\log_2 1 - \log_2 2 \right) \right]$$

$$= -[(0-1)]$$

]

$$E(S_{\text{hot}}) = 1$$

$$E(S_{\text{cool}}) = E[3+, 1-] = - \left[\frac{3}{4} \log_2 \left(\frac{3}{4} \right) + \frac{1}{4} \log_2 \left(\frac{1}{4} \right) \right]$$

$$= -[-0.31128 - 0.5]$$

$$= 0.81128 \approx 0.81$$

$\left(\frac{2}{5}\right)$

$$E(S_{\text{mild}}) = E[4+, 2-] = - \left[\frac{4}{6} \log_2 \left(\frac{4}{6} \right) \right.$$

$$\left. + \left(\frac{2}{6} \right) \log_2 \left(\frac{2}{6} \right) \right]$$

$$= -[-0.3898 - 0.5283] = 0.9181$$

$$\approx 0.92$$

$$\text{Gain}(S, \text{temp}) = E(S) - \frac{\|S_{\text{hot}}\| \times E(S_{\text{hot}})}{\|S\|}$$

$$+ \frac{\|S_{\text{cool}}\| \times E(S_{\text{cool}})}{\|S\|}$$

$$+ \frac{\|S_{\text{mild}}\| \times E(S_{\text{mild}})}{\|S\|}$$

$$= 0.94 - \left[\frac{4}{14} (1) + \frac{4}{14} (0.81) + \frac{6}{14} (0.92) \right]$$

$$= 0.94 - \left[0.2857 + \frac{2}{7} (0.81) + \frac{3}{7} (0.92) \right]$$

$$= 0.94 - \left[0.2857 + 0.2857 (0.81) + 0.42857 (0.92) \right]$$

$$= 0.94 - [0.2857 + 0.2314 + 0.39428]$$

$$= 0.94 - [0.9113844]$$

$$\text{Gain}(S, \text{temp}) \approx 0.0296$$

To compute $\text{Gain}(S, \text{Humidity})$:

$$E[S_{\text{high}}] = E[3+, 4-] = - \left[\frac{3}{7} \log_2 \left(\frac{3}{7} \right) \right]$$

+

$$\frac{4}{7} \log_2 \left(\frac{4}{7} \right)$$

$$= - \left[0.4857 \log_2 \left(\frac{3}{7} \right) + 0.5114 \log_2 \left(\frac{4}{7} \right) \right]$$

$$= - [-0.5937 - 0.46132] = +0.985$$

$$E[S_{\text{Normal}}] = E[6+, 1-]$$

$$= - \left[\frac{6}{7} \log \left(\frac{6}{7} \right) + \frac{1}{7} \log \left(\frac{1}{7} \right) \right]$$

$$= - \frac{1}{\log 2} \left[-0.0574 + 0.142857 \log \left(\frac{1}{7} \right) \right]$$

$$= - \frac{1}{\log 2} \left[-0.0574 - 0.120728 \right]$$

$$= - \frac{1}{\log 2} \left[-0.178128 \right]$$

$$= 0.5917$$

$$\text{Gain}(S, \text{Humidity}) = E(S) - \left[\frac{\|S_{\text{Normal}}\| \times E(S_{\text{Normal}})}{\|S\|} \right]$$

$$= E(S) - \left[\frac{\|S_{\text{High}}\| \times E(S_{\text{High}})}{\|S\|} \right]$$

$$= E(S) - \left[\frac{7}{14} (0.5917) + \frac{7}{14} (-0.985) \right]$$

$$= 0.94 - \left[\frac{1}{2} (0.5917) + \frac{0.985}{2} \right]$$

$$= 0.94 - [(0.5)(0.5917) + (0.5)(0.985)]$$

$$= 0.94 - [0.29585 + 0.4925]$$

$$= 0.94 - [0.78835]$$

$$= 0.15165$$

$$4 \log_2 \left(\frac{4}{7} \right)$$

$$= +0.985$$

To compute Gain (S , word) we need :

$$E[S_{\text{weak}}] = E[5+, 2-]$$

$$= - \left[\frac{5}{8} \log_2 \left(\frac{5}{8} \right) \right]$$

+

$$\left. \frac{2}{8} \log_2 \left(\frac{2}{8} \right) \right]$$

$$= - \frac{1}{\log 2} \left[-0.0937 \left(-0.1505 \right) \right]$$

+

$$= - \frac{1}{\log 2} \left[-0.244215 \right]$$

$$= \log 0.84126$$

$$= -0.244215$$

$$E[S_{\text{strong}}] = E[3+, 3-] = - \left[\frac{3}{6} \log_2 \left(\frac{3}{6} \right) \right]$$

$$\frac{3}{6} \log_2 \left(\frac{3}{6} \right)$$

$$\left(\frac{3}{6} \right)$$

Maximum

$$= - \left[\frac{1}{2} \log_2 \left(\frac{1}{2} \right) + \frac{1}{2} \log_2 \left(\frac{1}{2} \right) \right]$$

$$= - \left[\frac{2 \times 1}{2} \left[\log_2 (1) - \log_2 2 \right] \right]$$

$$= -[0 + 1]$$

$$= +1$$

$$= 1$$

Information
Gain (S)
Gain (S)
Gain (S)
Gain (S)

Information Gain:

$$\text{Gain}(S, \text{Outlook}) = 0.247$$

$$\text{Gain}(S, \text{Temp}) = 0.0286$$

$$\text{Gain}(S, \text{Humidity}) = 0.15165$$

$$\text{Gain}(S, \text{Wind}) = E(S) - \left[\frac{\|S_{\text{weak}}\| \times E(S_{\text{weak}})}{\|S\|} \right]$$

$$+ \left[\frac{\|S_{\text{strong}}\| \times E(S_{\text{strong}})}{\|S\|} \right]$$

$$= 0.94 - \left[\frac{8}{14} (0.81126) + \frac{6}{14} (1) \right]$$

$$= 0.94 - \left[\frac{4}{7} (0.81126) + \frac{6}{7} (1) \right]$$

$$= 0.94 - \left[0.5714 (0.81126) + \frac{6}{14} (1) \right]$$

$$= 0.94 - [0.46357 + 0.42857]$$

$\left(\frac{3}{6}\right)$

$$= 0.94 - [0.89214]$$

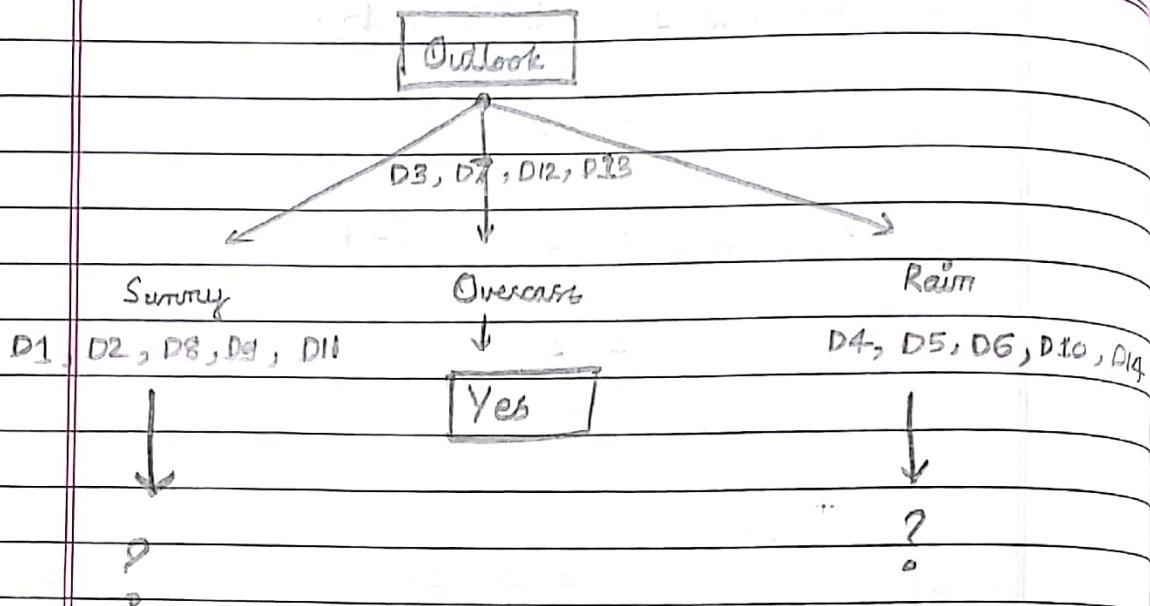
$$= 0.0478$$

$\left(\frac{3}{6}\right)$

Maximum Information Gain is with Outlook.

$$+ \frac{1}{2} \log_2 \left(\frac{1}{2} \right) \quad \therefore \text{Root Node will be Outlook.}$$

$$\left[\frac{1}{2} - \frac{1}{2} \log_2 \frac{1}{2} \right]$$



Calculating for Sunny:

	Temp	Humidity	wind	Play Tennis
D1	HOT	High	weak	No
D2	HOT	High	Strong	No
D8	Mild	High	Weak	No
D9	Cool	Normal	Weak	Yes
D11	Mild	Normal	Strong	Yes

$$\begin{aligned}
 E(\text{Sunny}) &= - \left[\frac{2}{5} \log\left(\frac{2}{5}\right) + \frac{3}{5} \log\left(\frac{3}{5}\right) \right] \\
 &= - \frac{1}{\log 2} \left[-0.15918 - 0.13311 \right] \\
 &\approx 0.971
 \end{aligned}$$

Calculating Gain(Sunny, Temp) we need:

$$E(\text{Sunny}_{\text{HOT}}) = E[0+, 2-] = 0$$

$$E[\text{Sunny Mild}] = E[1+, 1-]$$

$$= 1$$

$$E[\text{Sunny cool}] = E[1+, 0-]$$

$$= 0$$

$$\text{Gain}(\text{Sunny}, \text{Temp}) = E(\text{Sunny})$$

$$- \left[\frac{\|\text{Sunny Hot}\|}{\|\text{Sunny}\|} \times E(\text{Sunny Hot}) \right]$$

+

$$\frac{\|\text{Sunny Mild}\| \times E(\text{Sunny Mild})}{\|\text{Sunny}\|}$$

+

$$\frac{\|\text{Sunny cool}\| \times E(\text{Sunny cool})}{\|\text{Sunny}\|}$$

$$= 0.971 - \left[\frac{2}{5}(0) + \frac{2}{5}(1) + \frac{1}{5}(0) \right]$$

$$= 0.971 - (0.40)$$

$$= 0.571$$

Calculating Gain(Sunny, Humidity) we need:

$$E(\text{Sunny High}) = - \left[0 + \frac{3}{3} \log_2 \left(\frac{3}{3} \right) \right]$$

$$E[0+, 3-] = 0$$

$$E(Sunny_{Normal}) = E[2+, 0-]$$

$$= 0$$

Gain (Sunny, Humidity)

$$\begin{aligned} &= E(Sunny) - [E(Sunny_{High}) \times \frac{1}{2}] \\ &\quad + [E(Sunny_{Normal}) \times \frac{1}{2}] \\ &= 0.971 - [0 + 0] \\ &= 0.971 \end{aligned}$$

Calculating Gain (Sunny, Wind) we need :

$$E(Sunny_{Strong}) = E[1+, 1-]$$

$$= 0$$

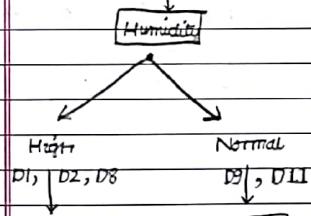
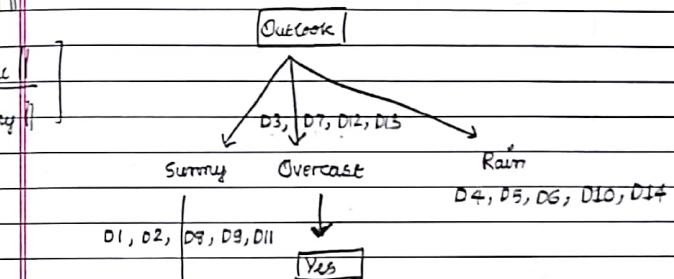
$$\begin{aligned} E(Sunny_{weak}) &= E[1+, 2-] \\ &= -\left[\frac{1}{3} \log_2\left(\frac{1}{3}\right) + \frac{2}{3} \log_2\left(\frac{2}{3}\right)\right] \\ &= -\frac{1}{\log 2} [-0.159 - 0.1174] \\ &= 0.498 \end{aligned}$$

$$\begin{aligned} \text{Gain (Sunny, wind)} &= E(Sunny) - \left[\frac{1}{5} E(Sunny_{Strong}) + \frac{3}{5} E(Sunny_{weak}) \right] \\ &= 0.971 - \left(\frac{0}{5} + \frac{3}{5} (0.498) \right) = 0.4882 \end{aligned}$$

$$\text{Gain (Sunny, Humidity)} = 0.971$$

$$\text{Gain (Sunny, Wind)} = 0.4882$$

Maximum Information Gain is with Humidity
∴ Next Node in the tree is Humidity



Calculating for Rainy :

	Temp	Humidity	Wind	Play Tennis
D4	Mild	High	Weak	Yes
D5	Cool	Normal	Weak	Yes
D6	Cool	Normal	Strong	No
D10	Mild	Normal	Weak	Yes
D14	Mild	High	Strong	No

$$E(\text{Rainy}) = E[3+, 2-] = - \left[\frac{3}{5} \log_2 \left(\frac{3}{5} \right) \right]$$

$$+ \left[\frac{2}{5} \log_2 \left(\frac{2}{5} \right) \right] \\ = - \frac{1}{\log 2} [0.6(-0.22185) \\ + 0.4(-0.3979)]$$

$$= - \frac{1}{\log 2} [-0.13311 - 0.15916] \\ = 0.97089 \\ \approx 0.979$$

Calculating Gain (Rainy, Temp) we need :

$$E(\text{Rainy}_{\text{Mild}}) = E[2+, 1-] \\ = - \left[\frac{2}{3} \log_2 \left(\frac{2}{3} \right) + \frac{1}{3} \log_2 \left(\frac{1}{3} \right) \right]$$

$$= -\frac{1}{\log 2} [-0.116 - 0.1574]$$

$$= 0.9084$$

$$\begin{aligned} E(\text{Rainy}_{\text{Cool}}) &= E[1+, 1-] \\ &= 1 \end{aligned}$$

$$\text{Gain}(\text{Rainy}, \text{Temp}) = E(\text{Rainy})$$

$$= \left[\frac{3}{5}(0.9084) + \frac{2}{5}(1) \right]$$

$$= 0.979 - [0.54504 + 0.4]$$

$$= 0.979 - 0.94504$$

$$= 0.03396$$

916]

Calculating Gain (Rainy, Humidity) we need :

$$\begin{aligned} E(\text{Rainy}_{\text{High}}) &= E[1+, 1-] \\ &= 1 \end{aligned}$$

$$E(\text{Rainy}_{\text{Normal}}) = E[2+, 1-]$$

$$= 0.9084$$

$$\begin{aligned} \text{Gain}(\text{Rainy}, \text{Humidity}) &= 0.979 - \left[\frac{2}{5}(1) + \frac{3}{5}(0.9084) \right] \\ &= 0.979 - [0.4 + 0.54504] = 0.03396 \end{aligned}$$

Calculating Gain (Rainy, Wind) we need :

$$E(\text{Rainy}_{\text{Weak}}) = E[3+, 0-]$$

$$= 0$$

$$E(\text{Rainy}_{\text{Strong}}) = E[0+, 2-]$$

$$= 0$$

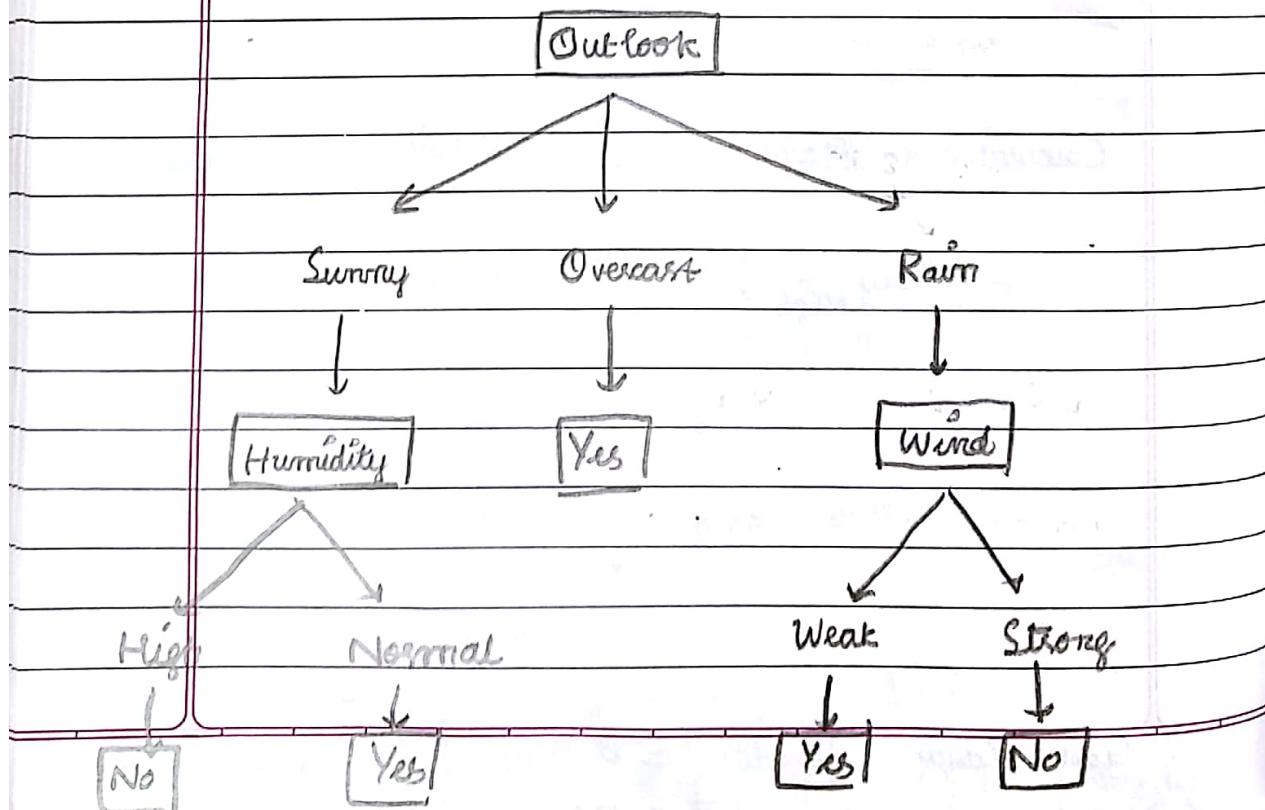
$$\text{Gain}(\text{Rain}, \text{Wind}) = E(\text{Rain}) - [0+0]$$

$$= 0.979 - 0$$

$$= 0.979$$

Information Gain is max. with wind

\therefore Next node in the tree will be wind.



Instance	Gender	Age	Country	Recommendation
D1	Male	Young	USA	Action (+)
D2	Female	Young	India	Action (+)
D3	Male	Old	USA	Drama (-)
D4	Female	Old	India	Action (+)
D5	Female	Young	USA	Action (+)
D6	Male	Old	India	Drama (-)
D7	Female	Old	China	Drama (-)
D8	Male	Young	China	Action (+)

$$\begin{aligned}
 \text{Entropy, } E(S) &= - [p_{\text{Action}} \log_2 p_{\text{Action}} + p_{\text{Drama}} \log_2 p_{\text{Drama}}] \\
 &= - \left[\left(\frac{5}{8} \right) \log_2 \left(\frac{5}{8} \right) + \left(\frac{3}{8} \right) \log_2 \left(\frac{3}{8} \right) \right] \\
 &= - \left[0.625 \log_2 \left(\frac{5}{8} \right) + (0.375) \log_2 \left(\frac{3}{8} \right) \right] \\
 &= - \frac{1}{\log_2 10} \left[-0.1275 - 0.15974 \right] \\
 &= -0.28724 = 0.9542
 \end{aligned}$$

Information Gain (Sample S, For attribute A):

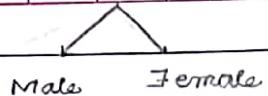
$$\text{Gain}(S, A) = E(S) - \sum_{v \in \text{values for } A} \frac{\|S_v\|}{\|S\|} \times E(S_v)$$

For attribute Gender, Sample S

DATE:

PAGE:

Given sur



[2+, 2-] [3+, 1-]

$$E(S, \text{Gender}_{\text{Male}}) = E(S, \text{Gender}_{\text{Female}})$$

$$= 1 - \left[\frac{3}{4} \log_2 \left(\frac{3}{4} \right) \right]$$

+

$$\frac{1}{4} \log_2 \left(\frac{1}{4} \right)$$

$$= -\frac{1}{\log_2 10} [0.75 (-0.1249)]$$

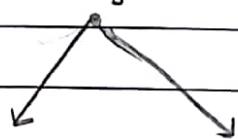
+

$$0.25 (-0.60205)$$

$$= -\frac{1}{\log_2 10} [-0.093705 - 0.150525]$$

$$= 0.81132$$

Given sample S, for attribute Age



[4+, 0-] [1+, 3-]

$$E(S, \text{Age}_{\text{Young}}) = 0$$

$$E(S, \text{Age}_{\text{Old}}) = 0.81132$$

Given sample S, for attribute Country

DATE:

PAGE:



USA(3) India(3) China(2)

[2+, 1-] [2+, 1-] [1+, 1-]

$$E(S, \text{Country}_{\text{USA}}) = - \left[\frac{2}{3} \log_2 \left(\frac{2}{3} \right) \right] = 0.9091$$

$$E(S, \text{Country}_{\text{India}}) = \frac{1}{3} \log_2 \left(\frac{1}{3} \right) = 1$$

$$= - \frac{1}{\log 2} \left[0.66 \log_{10} \left(\frac{2}{3} \right) \right]$$

$$+ 0.93 \log_{10} \left(\frac{1}{3} \right) \Big]$$

$$0.1505 \quad 25 \quad = - (3.3219) \left[-0.11622 - 0.15745 \right]$$

$$= - (3.3219) (0.27367)$$

$$= 0.9091$$

Information Gain (Sample S, for attribute A)

$$\text{Info Gain } (S, \text{Gender}) = E(S) - \left[\frac{\|S_{\text{Male}}\|}{\|S\|} \times E(S, \text{Gender}_{\text{Male}}) \right]$$

$$+ \frac{\|S_{\text{Female}}\|}{\|S\|} \times E(S, \text{Gender}_{\text{Female}})$$

$$= 0.9542 - \left[\frac{4}{8} (1) + \frac{4}{8} (0.81132) \right]$$

$$= 0.9542 - [0.5 + 0.40516]$$

$$= 0.9542 - 0.90516 = 0.04904$$

$$\text{Info Gain } (S, \text{Age}) = 0.9542 - \left[\frac{4}{8}(0) + \frac{4}{8}(0.81182) \right]$$

$$= 0.9542 - 0.40566$$

$$= 0.54854$$

$$\text{Info Gain } (S, \text{Country}) = 0.9542 - \left[\frac{3}{8}(0.9091) \right]$$

+

$$\frac{3}{8}(0.9091)$$

+

$$\frac{2}{8}(1)$$

Same

$$= 0.9542 - \left[0.375(0.9091) \right]$$

+

$$0.375(0.9091)$$

+

$$0.25$$

]

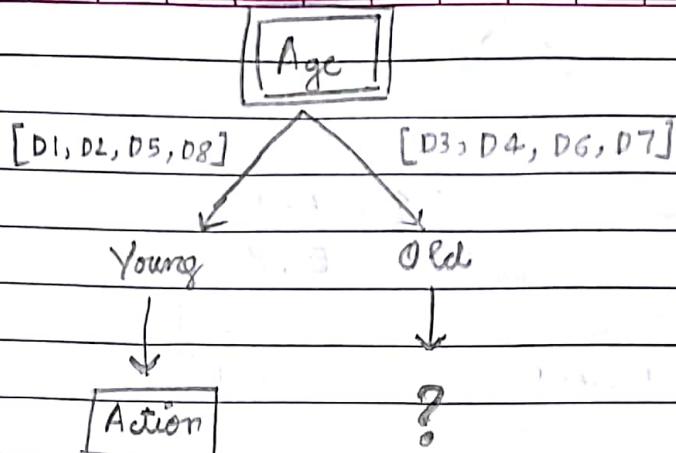
$$= 0.9542 - [0.681825 + 0.25]$$

$$= 0.9542 - [0.931825]$$

$$= 0.022375$$

Info Gain is maximum for Age

\therefore Age becomes the root node.



Given Old, the Sample S becomes :

Sample T

Insurance	Gender	Country	Recommendation
D3	Male	USA	Drama (-)
D4	Female	India	Action (+)
D6	Male	India	Drama (-)
D7	Female	China	Drama (-)

$$E(T) = - \left[p_{\text{Action}} \log_2 p_{\text{Action}} + p_{\text{Drama}} \log_2 p_{\text{Drama}} \right]$$

$$= - \left[\frac{1}{4} \log_2 \left(\frac{1}{4} \right) + \frac{3}{4} \log_2 \left(\frac{3}{4} \right) \right]$$

$$= - \frac{1}{4} \left[0.25 (\log_2 1 - 2) + 0.75 (\log_2 \frac{3}{4}) \right]$$

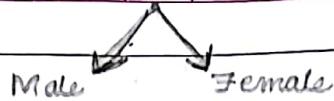
$$= - 0.3219 (-0.5 + 0.75 (\log_2 1.58496 - 2))$$

$$= 0.81128$$

For sample P, Gender (4)

DATE:

PAGE:



(♂) (♀)

[0+, 2-] [1+, 1-]

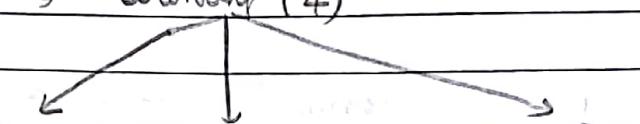
E (P, Gender_{Male})

= 0

E (P, Gender_{Female})

= 1

For sample P, Country (4)



(1)

[0+, 2-]

E (P, Country_{USA})

= 0

(2)

[1+, 1-]

E (P, Country_{India})

= 1

(1)

[0+, 1-]

E (P, Country_{China})

= 0

Information Gain (P, Gender)

$$= 0.81128 - \left[\frac{2}{4} (0) + \frac{2}{4} (1) \right]$$

$$= 0.81128 - (0 + 0.5)$$

$$= 0.81128 - 0.5$$

$$= 0.31128$$

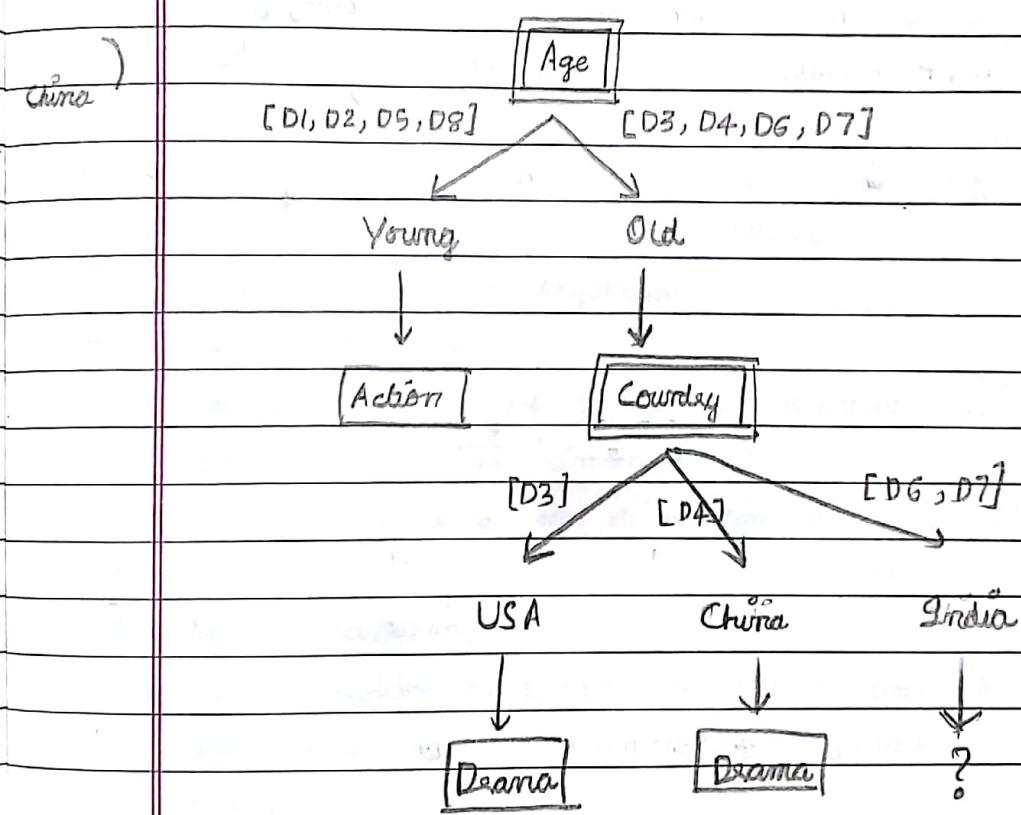
$$\text{Info Gain } (P, \text{Country}) = 0.81128 - \left[\frac{1(0)}{4} + \frac{1(0)}{4} + \frac{9(1)}{4} \right]$$

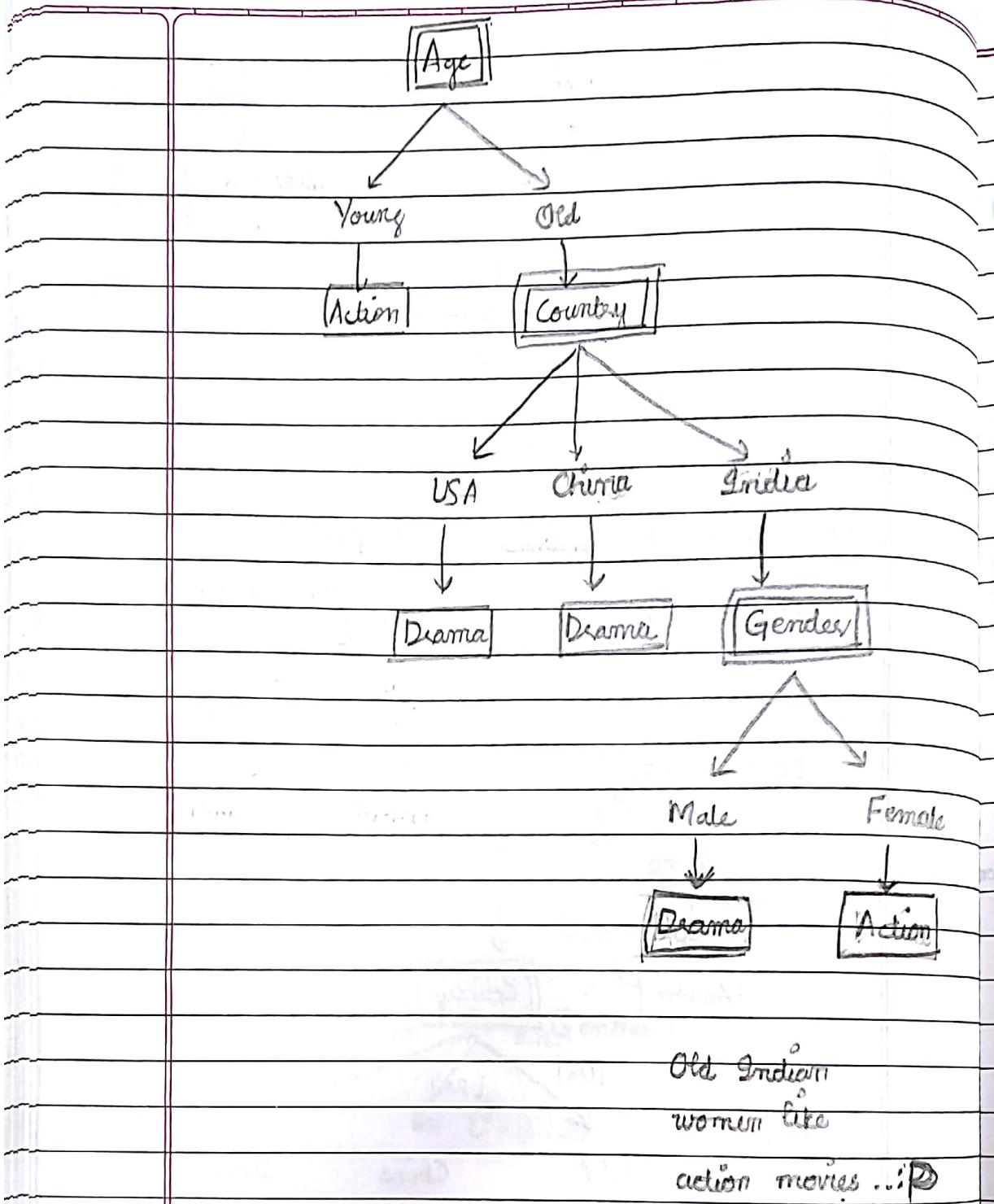
$$= 0.81128 - 0.5$$

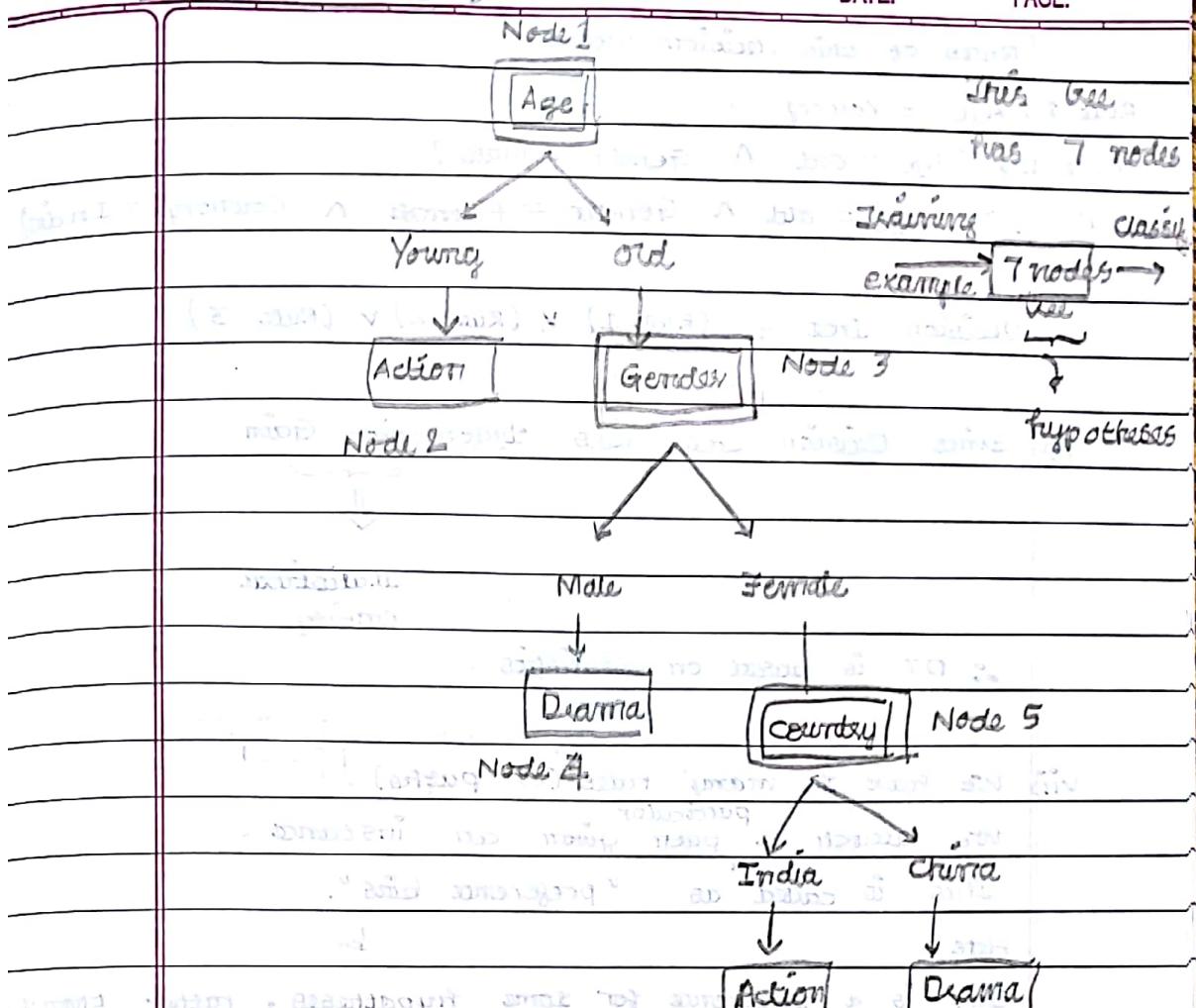
$$= 0.31128$$

$$\text{Info Gain } (P, \text{Country}) = \text{Info Gain } (P, \text{Gender})$$

Choosing Country arbitrarily, we get







Decision Tree = ID3 Algorithm

- i, Path from root to leaf : Rule or antecedent . Each rule is conjunction.
- ii, Leaf is called as consequent .
- iii, Decision Tree is a disjunction of conjunctions . (path or rule)

iv, Bias : Searching

- v, Once a decision to proceed towards a path is taken you can't backtrack (- Greedy Approach)

i.e.,
if your Test Structure is

old male India

and if old is present in the tree

once you have

gone down the tree

after looking at male & old

we can't go back . We will go to Gender now .

Rules of this decision tree :

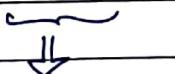
Rule 1 : Age = Young

Rule 2 : (Age = old \wedge Gender = Male)

Rule 3 : (Age = old \wedge Gender = Female \wedge Country = India)

Decision Tree = (Rule 1) \vee (Rule 2) \vee (Rule 3)

vi) Since Decision Tree uses Information Gain



Statistical entity

\therefore DT is based on statistics.

vii) We have so many rules (or paths).

We search a path given an instance.

This is called "preference bias".

Here,

Bias is a preference for some hypotheses, rather than restriction of the hypotheses space (Candidate Elimination Algorithm)

viii) Training of 40 examples

500 nodes \rightarrow classify just 40

hypotheses f'

If for same 40 training examples, if we get a 20 nodes tree is obtained

hypotheses f

then we prefer f over f' . This is the bias

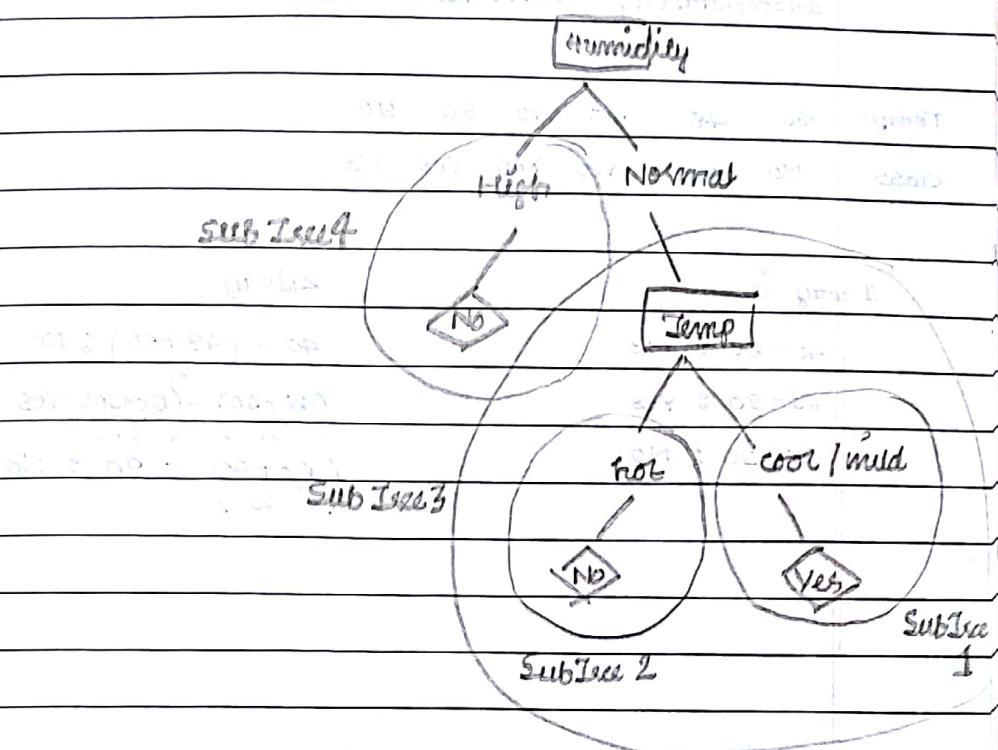
we have on preferring f over f' based

on Information Gain.

2 methods to improve accuracy:

- 1- Reduced Error Pruning : Removing subtrees to improve accuracy.
- 2- Rule based "": rule (or path or antecedent)

Reduced Error Pruning :



And so on we find multiple subtrees.

we choose subtree 1 & remove it.

check with all the training examples again.

• See if this pruned tree classifies well.

Similarly, we iteratively try to prune all subtrees one by one & try with training examples.

See if after pruning subtrees if accuracy is better.

This way we get better accuracy.

From we keep the shortest tree.

Rule based Pruning :

We pruned rules one by one (rule or path).

If accuracy is better, remove that rule (or path).

Else retain.

Incorporating Continuous Values :

Temp.	40	48	60	70	80	90
CLASS	No	No	Yes	Yes	Yes	No

1 way is

$$40 - 60 \Rightarrow \text{No}$$

$$60 - 80 \Rightarrow \text{Yes}$$

$$80 - 90 \Rightarrow \text{No}$$

2 way

$$40 - \frac{(48+60)}{2} \Rightarrow \text{No}$$

$$\frac{(48+60)}{2} - \frac{(80+90)}{2} \Rightarrow \text{Yes}$$

$$\frac{(80+90)}{2} - 90 \Rightarrow \text{No}$$

(Both are Supervised Learning Algo)

Inductive Learning

DATE:

PAGE:

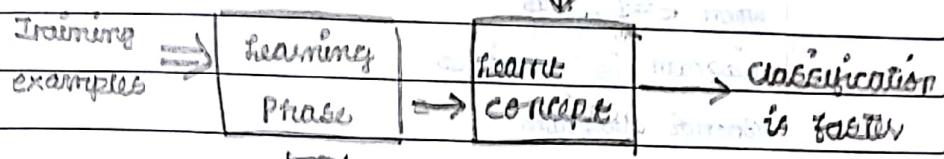
vs

Instance Based Learning

Inductive Learning

Training & Testing

Test Data



Consumes / Takes time

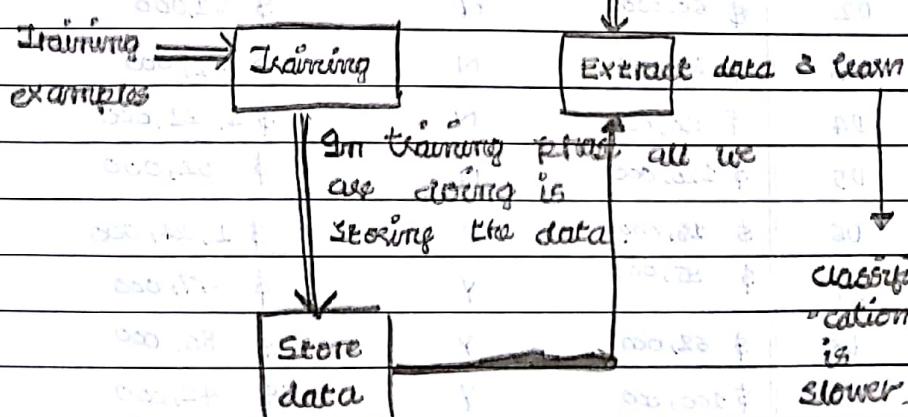
Example algorithms are :

- i) Concept Learning / Inductive Learning using Candidate Elimination Algorithm.
- ii) " " " " " " Decision Tree.

Instance Based Learning's (Lazy Learning)

Training & Testing Separated

For every Test Data



Example algorithm is

KNN or K - nearest neighbors

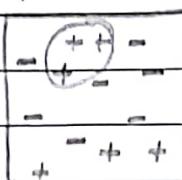
k nearest neighbors :

Requires 3 things : i) is the set of stored records -

$k = 3$

ii) Distance metric to compute

distances b/w records



iii) The value of k , the # of nearest neighbors to retrieve.

when $k=1$, this

diagram is called as

Voronoi diagram

For classification w.r.t., $k \geq 1$:

regression

i) In classification,

- take the most frequent

classification.

ii) In regression,

find the average of all

true values.

Example :

(problem part) & solution based problem

	Loan	Default	Distance from \$ 142,000
D1	\$ 40,000	N	\$ 1,40,000
D2	\$ 60,000	N	\$ 82,000
D3	\$ 80,000	N	\$ 62,000
D4	\$ 20,000	N	\$ 1,22,000
D5	\$ 120,000	N	\$ 22,000
D6	\$ 18,000	Y	\$ 1,24,000
D7	\$ 95,000	Y	\$ 47,000
D8	\$ 62,000	Y	\$ 80,000
D9	\$ 100,000	Y	\$ 42,000
D10	\$ 220,000	Y	\$ 78,000
D11	\$ 150,000	Y	\$ 8000

If $k=5$,

8K 47K 42K

D11, D7 & D9 are Y

DATE:

PAGE:

D5 & D3 are N

∴ classified as Y

rest

22K 62K

Ex 2 :

#	Action Scene	Romantic Scene	Emotional Scene	Classify (Movie)
D1	5	25	5	Romantic
D2	10	5	20	Emotional
D3	15	2	20	Emotional
D4	15	20	18	Romantic
D5	25	5	20	Action
D6	30	15	5	Action
Test Case	15	15	10	?

Distance of Test Case from Standard Instances:

$$D_1 = \sqrt{(15-5)^2 + (15-25)^2 + (10-5)^2} = \sqrt{10^2 + 10^2 + 5^2} = 15$$

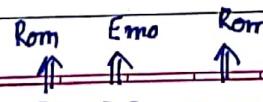
$$D_2 = \sqrt{(15-10)^2 + (15-15)^2 + (10-20)^2} = \sqrt{5^2 + 0^2 + 10^2} = 15$$

$$D_3 = \sqrt{(15-15)^2 + (15-2)^2 + (10-20)^2} = \sqrt{0 + 13^2 + 10^2} = 16.4$$

$$D_4 = \sqrt{(15-15)^2 + (15-9)^2 + (10-18)^2} = \sqrt{0 + 5^2 + 8^2} = 9.43$$

$$D_5 = \sqrt{(15-25)^2 + (15-5)^2 + (10-20)^2} = \sqrt{10^2 + 10^2 + 10^2} = 17.32$$

$$D_6 = \sqrt{(15-30)^2 + (15-15)^2 + (10-5)^2} = \sqrt{15^2 + 0 + 5^2} = 15.81$$

For $k = 3$, we get

D1, D2 & D4

∴ Test case is classified as Romantic

Algorithmically,

after we got D1, D2 & D4, (\because we decided $k = 3$)

how do we classify the test case? \Rightarrow

$\delta(A)$

For attribute Action Score (A),

D1 classifies Romantic Movie $\delta(A, RM) +$

D2 " Emotional " $\delta(A, EM) +$

D4 " Romantic " $\delta(A, RM)$

$$\text{sum} = 0 + 0 + 0 = 0$$

Similarly for attribute Romantic Score (R),

D1 classifies Romantic Movie $\delta(R, RM) +$

D2 " Emotional " " $\delta(R, EM) +$

D4 " Romantic " " $\delta(R, RM)$

$$\text{sum} = 1 + 0 + 1 = 2$$

Similarly for attribute Emotional Score (E),

$$\delta(E, RM) + \delta(E, EM) + \delta(E, RM)$$

=

$$\delta(E, RM) + \delta(E, EM) = +1 + 0 + 1 = 2$$

=

$$\delta(E, RM) + \delta(E, EM) + \delta(E, RM) = 1 + 1 + 1 = 3$$

$f(x_{\text{query}}) = \max \{0, 2, 1\} = 2$ (\therefore Romantic Movie is the classification.)

$$f(x_{\text{query}}) = f(x_{\text{test}}) = \arg \max_v \sum_{i=1}^k \delta(v, f(x_i))$$

$$f(x_{\text{query}}) = \delta(v, f(x_{\text{test}})) + \sum_{v \in \text{values}}$$

$$f(x_{\text{query}}) = \delta(v, f(x_{\text{test}})) + \delta(v, f(x_1)) + \delta(v, f(x_2)) + \dots$$

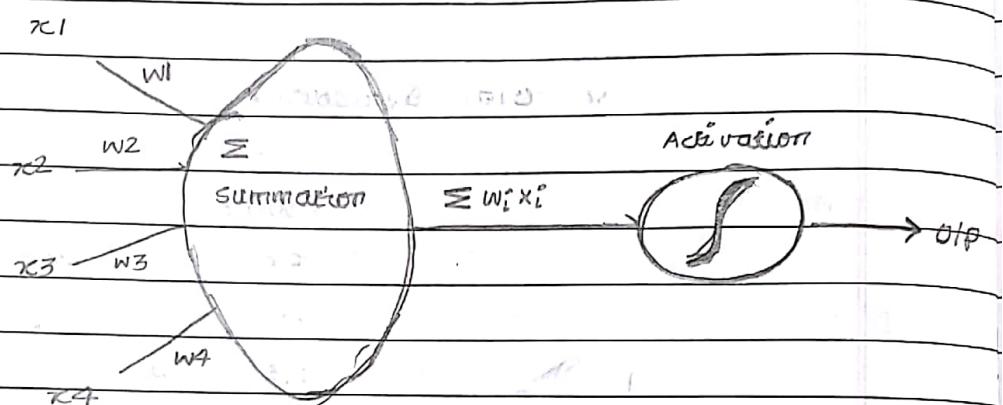
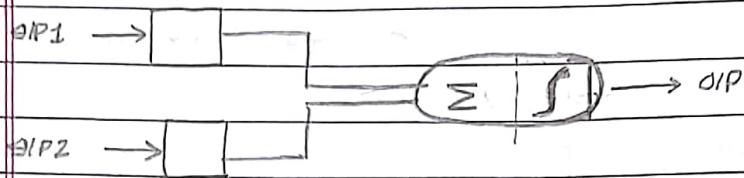
$$f(x_{\text{query}}) = \delta(v, f(x_{\text{test}})) + \delta(v, f(x_1)) + \delta(v, f(x_2)) + \dots$$

Better measures to calculate class from k neighbours?

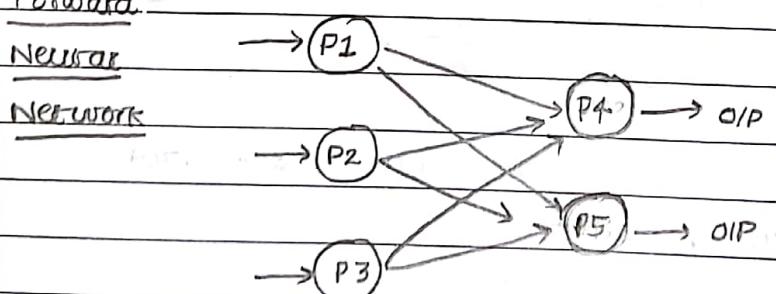
$$w_i = \frac{1}{d(x_i, x_{query})}$$

$$\hat{f}(x_{query}) = \underset{v \in \text{values}}{\operatorname{argmax}} \sum_{i=1}^k w_i \delta(v, f(x_i))$$

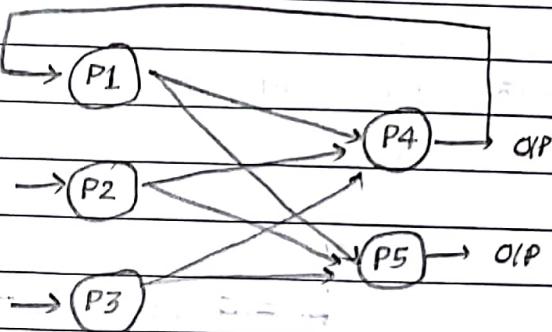
Exam question

Network (Concept Learning)Perception & Artificial Neural

$$\text{Activation } \exists i \begin{cases} = 1, \text{ if } w_i x_i > 0 \\ = -1, \text{ if } w_i x_i \leq 0 \end{cases}$$

ANN &Type 1° Forward

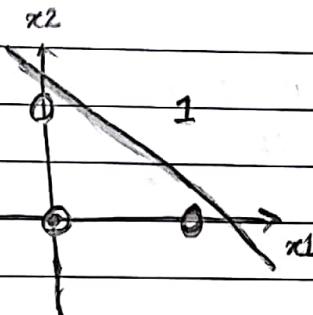
Type 2 : Feedback Forward Network



Linear Separation using Perceptron :

1. AND Gate :

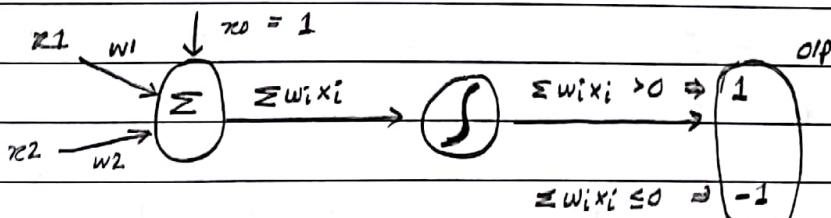
Instance	x_1	x_2	O/P
d1	0	0	0
d2	0	1	0
d3	1	0	0
d4	1	1	1



we see that this data is linearly separable.

∴ Perceptron can be used to implement AND Gate.

$$w_0 = -0.5 \text{ (Assume)}$$



initializing w_0 & w_2 with random values we get :

$$\text{take } w_1 = 1$$

$$w_2 = 0.5$$

$$d_i = w_0 x_0 + w_1 x_1 + w_2 x_2$$

For instance 1

$$\sum w_i x_i = -0.5(1) + 1(0) + (0.5)(0)$$

$$= -0.5$$

$$\text{Activation } (-0.5) \int = -1$$

when -1 , O/P $= 0$

$$\Rightarrow \hat{y}_1 = 0$$

For instance d1,

True label, $y_1 = 0$ [0 & 0 = 0]

Predicted label, $\hat{y}_1 = 0$

Error $= 0$ \therefore No need to update the weights

For instance 2 : (0, 1)

$$\begin{aligned} \sum w_i x_i &= w_0 x_0 + w_1 x_1 + w_2 x_2 \\ &= (-0.5)(1) + (1)(0) + (0.5)(1) \\ &= -0.5 + 0.5 \\ &= 0 \end{aligned}$$

$$\text{Activation } (0) = -1$$

Updating

when O/P or activation is -ve,

$$\text{O/P} = 0$$

$$\Rightarrow \hat{y}_2 = 0$$

For instance d2,

Updating

True label, $y_2 = 0$ [0 & 1 = 0]

Predicted label, $\hat{y}_2 = 0$

$$\text{Error} = 0$$

\therefore No need to update the weights.

Updating

For instance $d_3 (1, 0)$

$$\begin{aligned}\sum w_i x_i &= (-0.5)(1) + (1)(1) + (0.5)(0) \\ &= -0.5 + 1 \\ &= 0.5\end{aligned}$$

$$\text{Activation} (\sum w_i x_i) = 1$$

If Activation ($\sum w_i x_i$) > 0 , O/P = 1

$$\text{+ve} \Rightarrow \hat{y}_3 = 1$$

For instance d_3 ,

$$\text{true label, } y_3 = 0 \quad [\text{1 and } 0 = 0]$$

$$\text{predicted label, } \hat{y}_3 = 1$$

$$\text{Error} = -1$$

∴ Perceptron did not predict correctly

∴ Update the weights.

To update the weights for instances $d_3 (x_0, x_1, x_2)$

$$w_i' = w_i + (\text{true label} - \text{predicted label}) x_i \quad (1,$$

$$\text{Assuming } \eta = 0.2 \quad (1,$$

$$\begin{aligned}\text{Updating } w_0, \quad w_0 &= w_0 + \eta (y_3 - \hat{y}_3) x_0 \\ &= (-0.5) + 2(0 - 1)(1) \\ &= (-0.5) - (0.2) \\ &= -0.7\end{aligned}$$

$$\begin{aligned}\text{Updating } w_1, \quad w_1 &= w_1 + \eta (y_3 - \hat{y}_3) x_1 \\ &= 1 + 0.2(0 - 1)(1) \\ &= 1 - 0.2 \\ &= 0.8\end{aligned}$$

$$\begin{aligned}\text{Updating } w_2, \quad w_2 &= w_2 + \eta (y_3 - \hat{y}_3) x_2 \\ &= 0.5 + (0.2)(0 - 1)(0) \\ &= 0.5\end{aligned}$$

Checking for instances d1, d2 and d3,

For instance 1 (1, 0, 0) with updated weights :

$$\sum w_i x_i = (-0.7)(1) + (0.8)(0) + (0.5)(0)$$

$$= -0.7$$

$$\text{Activation} (\sum w_i x_i) = \text{Activation} (-0.7)$$

$$= -1$$

When OIP of activation is -ve ,

Predicted label $\hat{y}_1 = 0$

True label $y_1 = 0$

Error = 0 :)

For instance 2 (1, 0, 1) with updated weights :

$$\sum w_i x_i = (-0.7)(1) + (0.8)(0) + (0.5)(1)$$

$$= -0.7 + 0.5$$

$$= -0.2$$

$$\text{Activation} (\sum w_i x_i) = -1$$

When OIP of activation is -ve , predicted label = 0

True label = 0

Error = 0 :)

For instance 3 (1, 1, 0) with updated weights :

$$\sum w_i x_i = (-0.7)(1) + (0.8)(1) + (0.5)(0)$$

$$= -0.7 + 0.8 + 0$$

$$= 0.1$$

$$\text{Activation} (0.1) = -1$$

When OIP of activation is +ve, predicted label $\hat{y}_3 = 1$
 true label $y_3 = 0$
 $[0 \text{ & } 0 = 0]$
 Error = -1

\therefore we need to update the weights with instance 3
 one more time.

Updating weights with instance $d_3 (1, 1, 0)$ &

$$\begin{aligned} w_0 &= w_0 + \eta (y_3 - \hat{y}_3) x_0 \\ &= (-0.7) + 0.2(0-1)(1) \\ &= (-0.7) - 0.2 \\ &= -0.9 \end{aligned}$$

$$\begin{aligned} w_1 &= w_1 + \eta (y_3 - \hat{y}_3) x_1 \\ &= (0.8) + 0.2(0-1)(1) \\ &= (0.8) - 0.2 = 0.6 \end{aligned}$$

$$\begin{aligned} w_2 &= w_2 + \eta (y_3 - \hat{y}_3) x_2 \\ &= 0.5 + 0.2(0-1)(0) \\ &= 0.5 - 0 \\ &= 0.5 \end{aligned}$$

0

0 [0 & 1 = 0]

with $(w_0, w_1, w_2) = (-0.9, 0.6, 0.5)$

Instance (d_k)	$w_0 = -0.9$	$w_1 = 0.6$	$w_2 = 0.5$	$\sum_k A_k(z) y_k$
$d_1 (1, 0, 0)$	-0.9	0	0	-0.9 -1 0
$d_2 (1, 0, 1)$	-0.9	0	0.5	-0.4 -1 0
$d_3 (1, 1, 0)$	-0.9	0.6	0	-0.3 -1 0

For d_1, d_2, d_3

predicted label = true label

Using new $(w_0, w_1, w_2) = (-0.9, 0.6, 0.5)$
for instance $d_4 (1, 1, 1)$

Structure (d_4)	w_0	$w_1 + w_2$	Σ	Activation (Σ)	\hat{y}_4	y_4	Error
$d_4(1, 1, 1)$	-0.9	0.6 + 0.5	$-0.9 + 1.1$ $= 0.2$	1	1	1	0 :-)

\therefore Final learnt weights,

$$w_0 = -0.9$$

$$w_1 = 0.6$$

$$w_2 = 0.5$$

Approximate target f_{in} , $\hat{y} = -0.9 + 0.6x_1 + 0.5x_2$
for AND gate with initial $w_0 = 0.5$

$$w_1 = 1$$

$$w_2 = 0.5$$

$$\text{and } \eta = 0.5$$

weights were updated twice.

Ex 2 :	x_1	x_2	True label
	1	1	1 (*)
	2	-2	0 (Δ)
	-1	-1.5	0 (Δ)
	-2	-1	0 (Δ)
	-2	1	1 (*)
	1.5	-0.5	1 (*)

Assume initial values as,

$$w_0 = 0, x_0 = 1$$

$$w_1 = 1,$$

$$w_2 = 0.5$$

$$\eta = 0.2$$

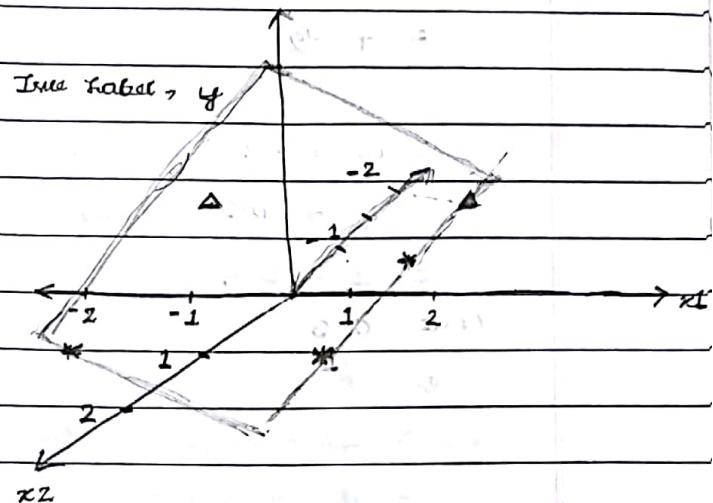
use $w_i = w_i + \eta (\text{true label} - \text{predicted label})^n$

using a single perceptron.

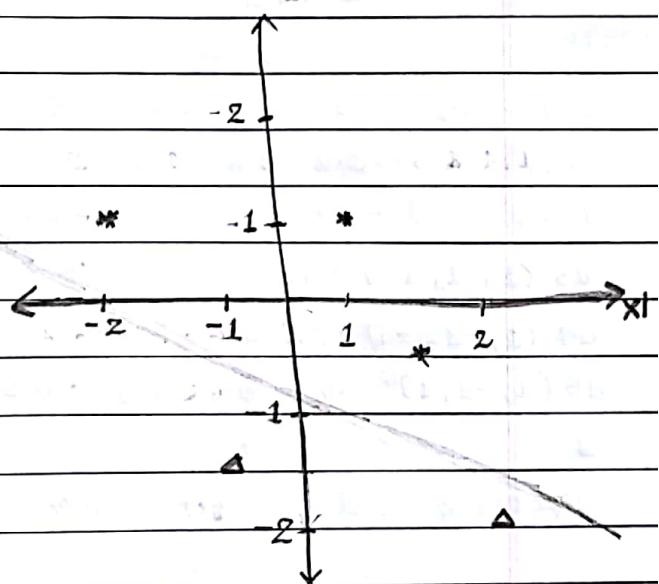
So, a single perceptron can only approximate target fn for data which is linearly separable.

Q: Checking if the given data is linearly separable,

Representation 1:



Representation 2:



Yes, it is linearly separable.

∴ A single perceptron can be used to approximate the target fn.

∴ Target fn is of the form, $\hat{y} = w_0 + w_1x_1 + w_2x_2$

DATE: PAGE:

Instance	d_k	w_0	w_1	w_2	Σ Activation (Σ)	Predicted label, \hat{y}_k	True label, y_k	Error ($y_k - \hat{y}_k$)
$d_1(1, 1, 1)$	0	1	0.5	1.5	1 ($\Sigma > 0$)	1	1	0 :
$d_2(1, 2, -2)$	0	1	0.5	1	1 ($\Sigma > 0$)	1	0	-1 :

updation of weights required with instance $d_2(1, 2, -2)$

$$\begin{aligned} w_0 &= w_0 + \eta(y_2 - \hat{y}_2)x_0 \\ &= 0 + 0.2(-1)(1) \\ &= -0.2 \end{aligned}$$

$$\begin{aligned} w_1 &= w_1 + \eta(y_2 - \hat{y}_2)x_1 \\ &= 1 + 0.2(-1)(2) \\ &= 1 - 0.4 \\ &= 0.6 \end{aligned}$$

$$\begin{aligned} w_2 &= w_2 + \eta(y_2 - \hat{y}_2)x_2 \\ &= 0.5 + 0.2(-1)(-2) \\ &= 0.5 + 0.4 \\ &= 0.9 \end{aligned}$$

Instance d_k w_0 w_1 w_2 Σ Activation (Σ) \hat{y}_k y_k Error ($y_k - \hat{y}_k$)

$d_1(1, 1, 1)$ -0.2 0.6 0.9 1.3 1 ($\Sigma > 0$) 1 1 0 :

$d_2(1, 2, -2)$ -0.2 0.6 0.9 0.8 0 ($\Sigma \leq 0$) 0 0 0 :

$d_3(1, -1, -1.5)$ -0.2 0.6 0.9 -2.15 0 ($\Sigma \leq 0$) 0 0 0 :

$d_4(1, -2, -1)$ -0.2 0.6 0.9 -2.3 0 ($\Sigma \leq 0$) 0 0 0 :

$d_5(1, -2, 1)$ -0.2 0.6 0.9 -0.5 0 ($\Sigma \leq 0$) 0 1 -1 :

Updating weights for instance $d_5(1, -2, 1)$:

$$\begin{aligned} w_0 &= w_0 + \eta(y_5 - \hat{y}_5)x_0 \\ &= -0.2 + 0.2(-1)(1) \\ &= -0.2 + 0.2 \\ &= 0 \end{aligned}$$

$$\begin{aligned} w_1 &= w_1 + \eta(y_5 - \hat{y}_5)x_1 \\ &= 0.6 + 0.2(1)(-2) \end{aligned}$$

$$= 0.6 - 0.4 = 0.2$$

Instance	d_k	w_0	w_1	w_2	Σ Activation (Σ)	\hat{y}_k	True label, y_k	Error ($y_k - \hat{y}_k$)
$d_1(1, 1, 1)$	1	1	1	1	1 ($\Sigma > 0$)	1	1	0 :
$d_2(1, 2, -2)$	0	0.2	1.1	-1.8	-1 ($\Sigma \leq 0$)	0	0	0 :

$$\begin{aligned} w_2 &= w_2 + \eta(y_5 - \hat{y}_5)x_2 \\ &= 0.9 + 0.2(1)(1) \\ &= 0.9 + 0.2 \\ &= 1.1 \end{aligned}$$

Instance	d_k	w_0	w_1	w_2	Σ Activation (Σ)	\hat{y}_k	y_k	Error ($y_k - \hat{y}_k$)
$d_1(1, 1, 1)$	-0.2	0.2	1.1	1.3	1 ($\Sigma > 0$)	1	1	0 :
$d_2(1, 2, -2)$	-0.2	0.2	1.1	-1.8	-1 ($\Sigma \leq 0$)	0	0	0 :
$d_3(1, -1, -1.5)$	0	0.2	1.1	-1.85	-1 ($\Sigma \leq 0$)	0	0	0 :
$d_4(1, -2, 1)$	0	0.2	1.1	-1.5	-1 ($\Sigma \leq 0$)	0	0	0 :
$d_5(1, -2, 1)$	0	0.2	1.1	0.7	1 ($\Sigma > 0$)	1	1	0 :
$d_6(1, 1.5, -0.5)$	0	0.2	1.1	-0.25	-1 ($\Sigma \leq 0$)	0	1	1 :

Updating weights for instance $d_6(1, 1.5, -0.5)$:

$$\begin{aligned} w_0 &= w_0 + \eta(y_6 - \hat{y}_6)x_0 = 0 + 0.2(1)(1) = 0.2 \\ w_1 &= w_1 + \eta(y_6 - \hat{y}_6)x_1 = 0.2 + 0.2(1)(1.5) = 0.2 + 0.3 = 0.5 \\ w_2 &= w_2 + \eta(y_6 - \hat{y}_6)x_2 = 1.1 + 0.2(1)(-0.5) = 1.1 - 0.1 = 1.0 \end{aligned}$$

Instance	d_k	w_0	w_1	w_2	Σ Activation (Σ)	\hat{y}_k	y_k	Error ($y_k - \hat{y}_k$)
$d_1(1, 1, 1)$	0.2	0.5	1.0	1.7	1 ($\Sigma > 0$)	-1	-1	0 :
$d_2(1, 2, -2)$	0.2	0.5	1.0	-0.8	-1 ($\Sigma \leq 0$)	0	0	0 :

Instance	d_k	w_0	w_1	w_2	Σ Activation (Σ)	\hat{y}_k	y_k	Error ($y_k - \hat{y}_k$)
$d_3(1, -1, -1.5)$	0.2	0.5	1.0	-1.8	-1 ($\Sigma \leq 0$)	0	0	0 :
$d_4(1, -2, 1)$	0.2	0.5	1.0	-1.8	-1 ($\Sigma \leq 0$)	0	0	0 :

Instance	d_k	w_0	w_1	w_2	Σ Activation (Σ)	\hat{y}_k	y_k	Error ($y_k - \hat{y}_k$)
$d_5(1, -2, 1)$	0.2	0.5	1.0	0.2	+1 ($\Sigma > 0$)	+1	+1	0 :
$d_6(1, 1.5, -0.5)$	0.2	0.5	1.0	0.45	1 ($\Sigma > 0$)	1	1	0 :

∴ Final weights learnt are: $w_0 = 0.2$, $w_1 = 0.5$ & $w_2 = 1.0$

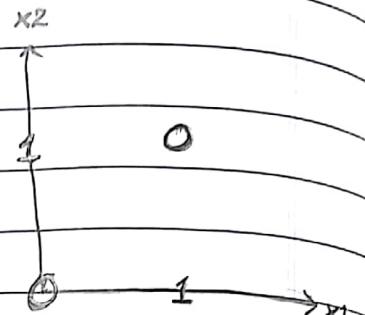
Approximate target fn: $\hat{y} = 0.2 + 0.5x_1 + 1.0x_2$

with initial $w_0 = 0$, $w_1 = 1$ & $w_2 = 0.5$ & $\eta = 0.2$

Why XOR gate can't be modeled using a single perceptron?

XOR Truth Table

Instance	x_1	x_2	True label
d1	0	0	0
d2	0	1	1
d3	1	0	1
d4	1	1	0



On plotting we find that

Algorithmically,

if we will be able

to say that this

dataset can not

be separated by a

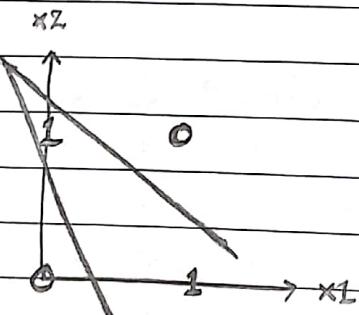
linear fit

which

you feed this data as input to a single perceptron

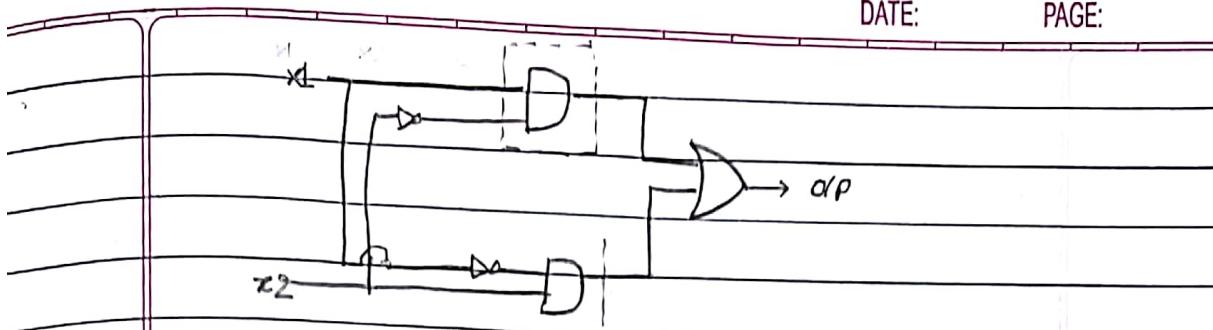
& the algorithm will never converge (infinite loop).

We would need more than one linear classifier (maybe a pair of lines)



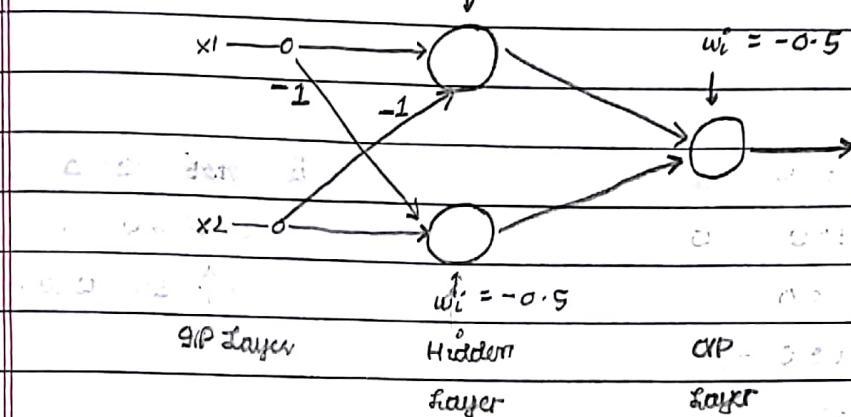
Modeling XOR Gate using more than one perceptron:

XOR Gate :



$$\text{XOR Gate} = x_1 x_2' + x_1' x_2$$

$$w_i^1 = -0.5$$



14th August, 2017

Gradient & Differentiating Linear Mean Square

DATE:

PAGE:

For a particular board state b' ,

$$\hat{v}(b') = \hat{w}_{0b'} + \hat{w}_{1b'} x_{b'}$$

where

$x_{b'}$ is the # of black pawns on
the board at board state b' .

For ~~successor~~ intermediate board state b' ,

$v_{train}(b')$ will not be given directly.

$$\therefore v_{train}(b') = \hat{v}(\text{successor board state of } b')$$

Error for this intermediate board state b' ,

$$E_{b'} = v_{train}(b') - \hat{v}(b')$$

Since b' is an intermediate board state b'

$$E_{b'} = [\hat{v}(\text{successor board state of } b') - \hat{v}(b')]^2$$

$$\text{Value of } [\hat{w}_{0(b'+1)} + \hat{w}_{1(b'+1)} x_{1(b'+1)}]$$

$$[\hat{w}_{0b'} + \hat{w}_{1b'} x_{1b'}]^2$$

We use Error calculated at board state b'

w_{ob}' to update \hat{w}_{ob}' & \hat{w}_{ib}' such that

This will do something like

w_{ob}' & w_{ib}' attain the ideal value.

$$\hat{w}_{ob}' + \eta = \hat{w}_{ob}' - \eta E_b'$$

↓ ↓

If we update this such that becomes equal to

or

close to equal to

w_{ob}'

$$= \hat{w}_{ob}' - \eta [\text{Value of } [\hat{w}_{o(b'+1)} + \hat{w}_{i(b'+1)} x_{ib'}]]$$

$$(\hat{w}_{ob}' + \hat{w}_{ib'} x_{ib'})^2$$

$w_{ob}' \approx \hat{w}_{ob}'$ when

$$\eta [\text{Value of } [\hat{w}_{o(b'+1)} + \hat{w}_{i(b'+1)} x_{ib'}]]$$

$$- [\hat{w}_{ob}' + \hat{w}_{ib'} x_{ib'}]^2 \approx 0$$

$$\Rightarrow \eta E_b' \approx 0$$

E_b' is minimum w.r.t \hat{w}_{ob}'

if when

$$\frac{dE_b'}{d\hat{w}_{ob}'} = 0$$

$$E_b' = \hat{w}_{ob'} + \hat{w}_{ib'} x_{ib'}$$

DATE: PAGE:

$$\frac{dE_b'}{dw_{ob'}} = 2E_b' (0 - 1 + 0) = -2E_b'$$

For $w_{ib'}$

$$\frac{dE_b'}{dx_{ib'}} = 2E_b' (0 - (0 + x_{ib'})) = -2E_b' x_{ib'}$$

So,

$$w_{ob'} = \hat{w}_{ob'} - \eta E_b'$$

$$w_{ib'} = \hat{w}_{ib'} - \eta E_b'$$

$\hat{w}_{ob'}$ becomes $\approx w_{ob'}$

when E_b' is min. E_b' is min wrt $w_{ob'}$

$$w_{ob'} = w_{ob'} - \eta (-2E_b')$$

$$= \hat{w}_{ob'} + 2\eta E_b'$$

$\therefore 2$ is a constant.

$$2\eta \approx \eta$$

$$w_{ob'} = \hat{w}_{ob'} + \eta E_b'$$

&

Similarly for $w_{ib'}$ we get,

$$w_{ib'} = \hat{w}_{ib'} - \eta (-2E_b' x_{ib'})$$

$$= w_{ib'} + 2\eta E_b' x_{ib'}$$

Considering $2\eta \approx \eta$

$$w_{16'} = \hat{w}_{16'} + \eta E_{16'} x_{16'}$$

DATE:

PAGE:

$$\overline{w_{16'}} = \hat{w}_{16'} + \eta [v_{train}(b') - \hat{v}(b)] x_{16'}$$

A single perceptron and OR Gate :

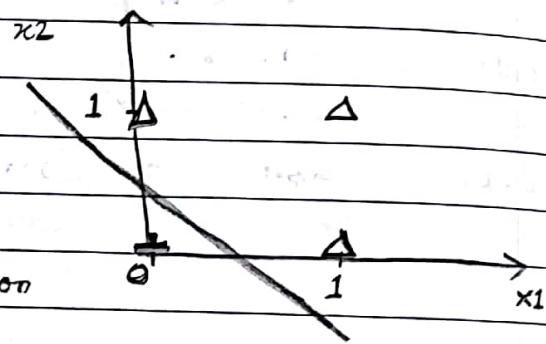
Instance	x_1	x_2	True Label
d1	0	0	0 (+)
d2	0	1	1 (-)
d3	1	0	1 (-)
d4	1	1	1 (-)

Is the data linearly separable?

Yes. The data
is linearly
separable.

So if a single perceptron
is used,
we'll be able to classify

the data (learning algorithm : gradient descent will
converge)



Set initial values for $w_0 = 1$, $\eta = 0.2$

$$w_1 = 0.5, w_2 = 0.5, x_0 = 1$$

Target for representation, $\hat{y} = w_0 + w_1 x_1 + w_2 x_2$

Instance $d_k(x_0, x_1, x_2)$	w_0	w_1	w_2	\sum	Activation	\hat{y}	y	Cost / Error (Ques of happiness)
$d1(1, 0, 0)$	-0.5	1	0.5	-0.5	-1 ($\Sigma \leq 0$)	0	0	0 :-)
$d2(1, 0, 1)$	-0.5	1	0.5	0	-1 ($\Sigma \leq 0$)	0	1	1 :-)

Updating weights with instance $d2(1, 0, 1)$

$$\begin{aligned} w_0 &= w_0 + \eta (\text{True Label} - \text{Predicted Label}) x_0 \\ &= -0.5 + 0.2(1)(1) \\ &= -0.5 + 0.2 = -0.3 \end{aligned}$$

$$w_1 = w_1 + \eta (\text{True Label} - \text{Predicted Label}) x_1$$

$$= 1 + 0.2(-1)(0)$$

$$= 1$$

$$w_2 = w_2 + \eta (\text{True Label} - \text{Predicted Label}) x_2$$

$$= 0.5 + 0.2(1)(1)$$

$$= 0.5 + 0.2$$

$$= 0.7$$

Instance	$d_k(x_0, x_1, x_2)$	w_0	w_1	w_2	\sum	Activation(Σ)	\hat{y}	y	Cost/Error
$d_1(1, 0, 0)$	-0.3	1	0.7	-0.3	-1 ($\Sigma \leq 0$)	0	0	0	0
$d_2(1, 0, 1)$	-0.3	1	0.7	0.4	+1 ($\Sigma > 0$)	1	1	0	0
$d_3(1, 1, 0)$	-0.3	1	0.7	0.7	+1 ($\Sigma > 0$)	1	1	0	0
$d_4(1, 1, 1)$	-0.3	1	0.7	1.4	+1 ($\Sigma > 0$)	1	1	0	0

The algorithm converged with

$$w_0 = -0.3$$

$$w_1 = 1$$

$$w_2 = 0.7$$

Approximated target fn, $\hat{y} = -0.3 + x_1 + 0.7x_2$

started with initial $w_0 = -0.5$

$$w_1 = 1$$

$$w_2 = 0.5$$

$$\& \eta = 0.2$$

σ (Degree of flappiness)

x

y

100

Algorithm

starting with

-ve Training DataTraining Data

Instance	Size	Color	Shape	Class / Label
d1	big	red	circle	No
d2	small	red	triangle	No
d3	small	red	circle	Yes
d4	big	blue	circle	No
d5	small	blue	circle	Yes

S0 $\{ \phi \mid \phi \subseteq \text{G} \}$ Go $\{ \phi \mid \phi \subseteq \text{G} \}$

When you are presented with a negative example, you need to remove from S_0 any hypothesis inconsistent with the current observation and replace any "hypotheses in G with its minimal specializations (all) that are consistent with the observation but still more general than some member of S_0 .

Considering d1 $\{ \phi \mid \phi \subseteq \text{G} \}$ No

-ve example;

Is consistent with all hypotheses in S_0 ?Yes \Rightarrow Do nothing

2. consistent with all " Go ?

No \Rightarrow Replace the inconsistent hypothesis in G with all its minimal specializations that are $G_1 = \{ \langle \text{small} \ ? \ ? \rangle, \langle \ ? \ , \text{blue} \ ? \rangle, \langle \ ? \ ? \ \text{triangle} \rangle \}$ $\langle \ ? \ , \ ? \ , \ ? \rangle$ $\langle \ ? \ ? \ ? \rangle$

iii) Remove from G_1 any hypothesis that is less general than another " in G_1
No such hypothesis.

∴ After d_1 ,

$$S_1 = S_0 = \{ \phi \quad \phi \quad \phi \}$$

$$G_1 = \{ \langle \text{small} \quad ? \quad ? \rangle, \\ \langle ? \quad \text{blue} \quad ? \rangle, \\ \langle ? \quad ? \quad \text{triangle} \rangle \}$$

Instance d_2 :

$$d_2 \quad \langle \text{small} \quad \text{red} \quad \text{triangle} \quad \text{No} \rangle$$

Since it's a -ve example,

1. Consistent with S_1 ? Yes \Rightarrow Do nothing $S_2 = S_1$
2. " " all hypothesis in G_1 ?
No. ϕ is not consistent with $\{f_i(d_k) \neq c(d_k)\}$
 $\langle \text{small} \quad ? \quad ? \rangle \quad \{f_1(d_2) = 1, c(d_2) = 0\}$

and

$$\langle ? \quad ? \quad \text{triangle} \rangle \quad \{f_1(d_2) = 1, c(d_2) = 0\}$$

i) Replace $\langle \text{small} \quad ? \quad ? \rangle$

with all its minimal specializations

such that they are consistent with the observation

(instance d_2) but still more general than

some member of S_0 .

Do the same for $\langle ?, ?, \text{triangle} \rangle$.

are consistent with the observation but still more general than
some member of S_0

Doing it for

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

we get $\{ \langle ? \rangle, \langle ? \rangle \}$

$\langle \text{small blue} ? \rangle \vee \langle \text{small ? circle} \rangle$

for $\langle ? \rangle ? \text{ triangle} \rangle$

we get $\{ \langle ? \rangle, \langle ? \rangle \}$

$\langle \text{big ? triangle} \rangle \vee \langle ? \rangle \text{ blue triangle}$

$$\therefore G_2 = \{ \langle ? \rangle \text{ blue} ? \rangle, \langle ? \rangle$$

$\langle \text{small blue} ? \rangle, \langle \text{small ? circle} \rangle$

$\langle \text{big ? triangle} \rangle \langle ? \rangle \text{ blue triangle} \rangle \}$

4) Remove $\langle ? \rangle$ from G_2 if any hypotheses which
is less general than any other hypothesis

in G_2 is less general than

$\langle \text{small blue} ? \rangle \wedge \langle ? \rangle \text{ blue triangle}$

is less general than

$\langle ? \rangle \text{ blue} ? \rangle$

$$\therefore G_2 = \{ \langle ? \rangle \text{ blue} ? \rangle \langle \text{small ? circle} \rangle$$

$\langle \text{big ? triangle} \rangle \}$

After instance d2,

$$S_2 = \{ \langle \phi \rangle, \langle \phi \rangle \}$$

$$G_2 = \{ \langle ? \rangle \text{ blue} ? \rangle \langle \text{small ? circle} \rangle$$

$\langle \text{big ? triangle} \rangle \}$

Instance d3, $\langle \text{small red circle} \rangle$ Yes?

For the example,

1. Consistent with all hypothesis in G_2 ? No.

Remove the inconsistent hypotheses from G_2 .

$$G_3 = \{ \text{<small ? circle>} \}$$

2. Consistent with all hypothesis in S_2 ? No.

i) Replace with all its (inconsistent hypothesis)
minimal generalizations but
still it should be more specific than
some number of G_3 .

$$S_3 = \text{<small red circle>}$$

ii) Remove from S_3 any hypothesis that is
more general than any other hypothesis in S_3 .

$$S_3 = \text{<small red circle>}$$

After instance d3,

$$G_3 = \{ \text{<small red circle>} \}$$

$$G_3 = \{ \text{<small ? circle>} \}$$

Similarly doing it for d4 <big blue circle No>

$$S_4 = \{ \text{<small red circle>} \}$$

$$G_4 = \{ \text{<small ? circle>} \}$$

Doing it for d5 <small blue circle Yes> we get:

$$S_5 = \{ \text{<small ? circle>} \}$$

$$G_5 = \{ \text{<small ? circle>} \}$$

$$\therefore S_5 = G_5 = \{ \text{<small ? circle>} \}$$

Approximation of target concept = < small ? video >

Finally we have learnt the concept of "small circles".

Derivation of the Backpropagation Rule:

$$E = \frac{1}{2} \sum_{k=1}^n (y_k - \hat{y}_k)^2$$

where n

$$\hat{y}_k = w_0 + \sum_{i=1}^m w_i x_i$$

Considering $\hat{y}_k = w_0 + w_1 x_1$

then,

$$E = \frac{1}{2} \sum_k (y_k - (w_0 + w_1 x_1))^2$$

Update rule,

w

E on the entire training dataset,

$$E = \frac{1}{2} \sum_{k=1}^{\# \text{ of examples}} (y_k - \hat{y}_k)^2$$

For a particular training instance d_k

$$E_k = \frac{1}{2} (y_k - \hat{y}_k)^2$$

$$= \frac{1}{2} (y_k - (w_0 + w_1 x_1))^2$$

c>

To update the weight, for instance k

$$w'_k = w_k - \eta \frac{\partial E_k}{\partial w_k}$$

$$w_0 = w_0 - \eta \frac{\partial}{\partial w_0} \left(\frac{1}{2} (y_k - (w_0 + w_1 x_1))^2 \right)$$

$$= w_0 - \frac{1}{2} \eta \left(2 \times (y_k - \hat{y}_k) \times (0 - 1) \right)$$

$$= w_0 + \eta (y_k - \hat{y}_k)$$

$$w_1 = w_1 - \eta \frac{\partial}{\partial w_1} \left(\frac{1}{2} (y_k - (w_0 + w_1 x_1))^2 \right)$$

$$= w_1 - \eta \frac{2 \times 1}{2} \times (y_k - \hat{y}_k) (0 - (0 + x_1))$$

$$= w_1 + \eta (y_k - \hat{y}_k) (x_1)$$

i) Step fn vs Sigmoid fn

DATE: PAGE:
Single vs Neural
Perception Network

LINEARLY VS NON
SEPARABLE LINEARLY
SEPARABLE

ii) Momentum

Gradient Descent:

$$E = \frac{1}{2} \sum_{i=0}^n (\text{True Label} - \text{Predicted Label})^2$$

DATE: PAGE:

Gradient Descent Up

date Rule:

Perception Update Rule: Discrete Value
if Step fn is used

$$w_i \leftarrow w_i + \eta (\text{True Label} - \text{Predicted Label}) x_i$$

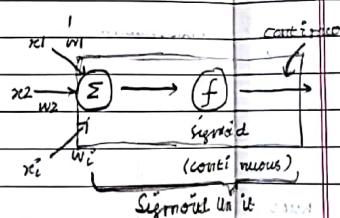
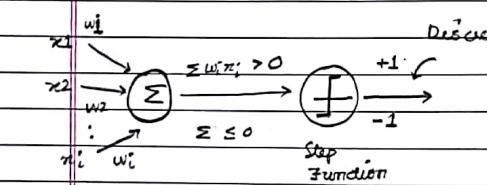
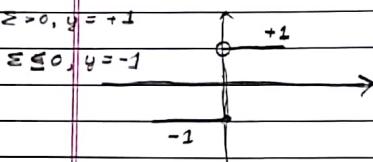
$$w_i \leftarrow w_i$$

$$+ \eta \sum_{i=0}^n (\text{True Label} - \text{Predicted Label}) x_i$$

General Gradient

Descent

Step Function is Non-Differentiable



$$\text{Differentiation of } \sigma(x) = \frac{1}{1+e^{-x}}$$

$$\frac{d}{dx} (\sigma(x)) = \frac{d}{dx} \left(\frac{1}{1+e^{-x}} \right)$$

$$= -\frac{1}{(1+e^{-x})^2} (0 - e^{-x})$$

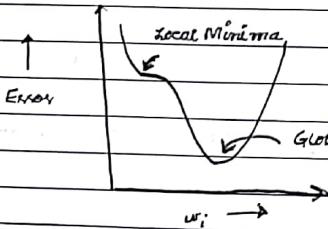
$$= \frac{e^{-x}}{(1+e^{-x})^2} = \frac{1+e^{-x}-1}{(1+e^{-x})^2} = \frac{1+e^{-x}-1}{(1+e^{-x})^2}$$

$$= \frac{1}{(1+e^{-x})} - \frac{1}{(1+e^{-x})^2}$$

$$= \sigma(x) - \sigma(x)^2$$

$$= \sigma(1-\sigma)$$

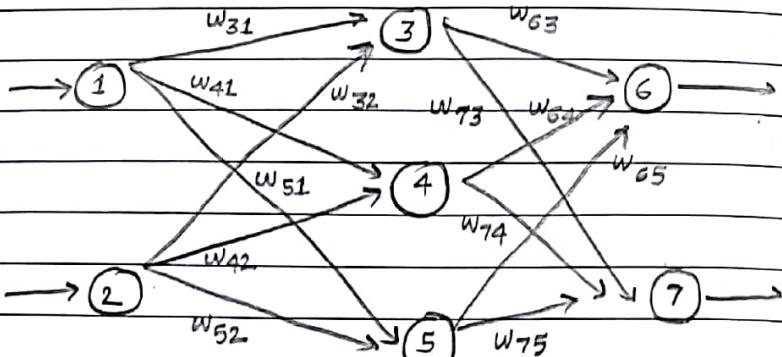
Momentum:



To prevent
from getting stuck
at local minima,
 $w_i \leftarrow \eta \sum_i x_i w_i + \alpha w_i (n-1)$,
 $0 \leq \alpha < 1$

$$\frac{d\sigma}{dx} = \sigma(1-\sigma)$$

Neural Network : # of oP from hidden layer < # of neurons



Input
Layer

Hidden
Layer

Output
Layer

(Each neuron
in the hidden
layer is a
sigmoid unit)

(Each
neuron in
the hidden
layer is a
sigmoid unit)

where each sig

we can update the weight of the neural network
either by using General Gradient Descent Rule

or

by using Stochastic Gradient Descent Rule.

Though, Stochastic Gradient Descent Rule gives
better accuracy generally.

Backpropagation

$E(\vec{w})$

General Gradient

Descent Rule

$$w_i \leftarrow w_i + \eta \frac{1}{2} \sum_{i=1}^n (\text{True Label} - \text{Predicted Label}) x_i$$

Stochastic Gradient Descent

$$w_i \leftarrow w_i + \eta (\text{True Label} - \text{Predicted Label}) x_i$$

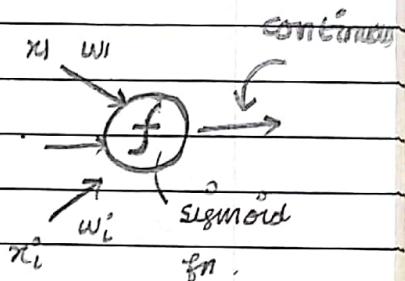
of form we do this after every example

of layer

∴ This is "Downstream"

neurons

where each sigmoid unit is

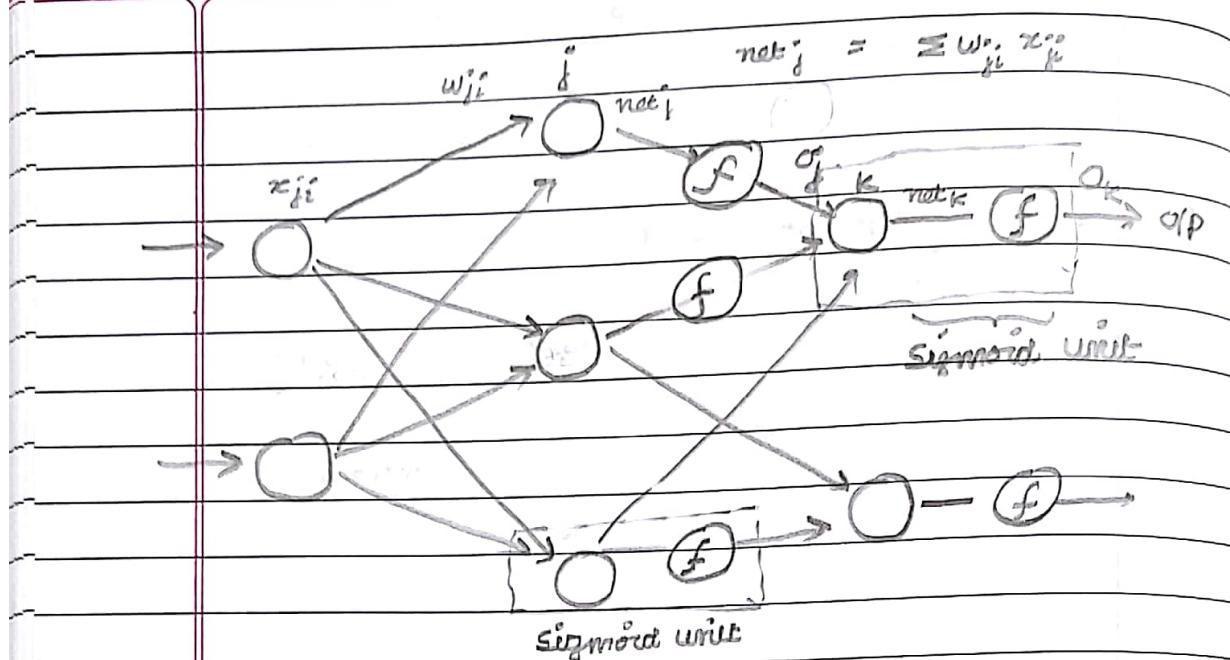


$$o(z) = \frac{1}{1 + e^{-z}}$$

Backpropagation Algorithm for Output Layer to Update weights:

$$E(\vec{w}) = \frac{1}{Z} \sum_d \sum_{k \in \text{# of labels/classes}} (t_{kd} - o_{kd})^2$$

$$\Delta w_{ij} = -\eta \frac{\partial E_d}{\partial w_{ij}} - (1B)$$



$$\frac{\partial E_d}{\partial w_{ji}} = \frac{\partial E_d}{\partial net_j} \times \frac{\partial net_j}{\partial w_{ji}}$$

$$\frac{\partial E_d}{\partial w_{ji}} = \frac{\partial E_d}{\partial net_j} \times x_{ji} - (1A)$$

Case 1 : Training rule for O/P units

$$\frac{\partial E_d}{\partial net_j} = \frac{\partial E_d}{\partial o_j} \times \frac{\partial o_j}{\partial net_j}$$

$$= \frac{1}{2} \times 2 \times \sum_{k \in \# \text{ of labels}} (t_k - o_k)^2$$

$$= \frac{1}{2} \times (t_j - o_j) \times \frac{\partial}{\partial o_j} (t_j - o_j)^2$$

$$= (t_j - o_j)(-1)$$

$$\frac{\partial E_d}{\partial o_j} = -(\epsilon_j - o_j) \quad \text{--- (2)}$$

Since $o_j = \sigma(\text{net}_j)$

$$\text{Now, 2nd term in } \frac{\partial o_j}{\partial \text{net}_j} = \frac{\partial (\sigma(\text{net}_j))}{\partial \text{net}_j}$$

$$\frac{\partial o_j}{\partial \text{net}_j} = \sigma(\text{net}_j) [1 - \sigma(\text{net}_j)]$$

$$\text{Substituting } \sigma(\text{net}_j) = o_j$$

$$\frac{\partial o_j}{\partial \text{net}_j} = o_j(1 - o_j) \quad \text{--- (3)}$$

Substituting (3) in (1A) we get

$$\frac{\partial E_d}{\partial w_{ji}} = -(\epsilon_j - o_j)o_j(1 - o_j)x_{ji} \quad \text{--- (5)}$$

Substituting eqn (5) in eqn (1B)

$$\Delta w_{ji} = \eta(\epsilon_j - o_j)o_j(1 - o_j)x_{ji}$$

$$\Delta w_{ji} = \Delta w_{ji}^e + \eta (t_i - o_i^e) o_i^e (1 - o_i^e) x_{ji}^e$$

$$\frac{\partial E}{\partial net_j} = \sum_k \frac{\partial E}{\partial net_k} \cdot \frac{\partial net_k}{\partial net_j}$$

We represent, $\frac{\partial E}{\partial net_k} = -\delta_k$

$$\frac{\partial E}{\partial net_j} = \sum_{k \in ds(j)} -\delta_k \cdot \frac{\partial net_k}{\partial net_j}$$

$$= \sum_{k \in ds(j)} -\delta_k \cdot \frac{\partial net_k}{\partial o_f} \cdot \frac{\partial o_f}{\partial net_j}$$

Rewriting, $\frac{\partial net_k}{\partial o_f} = w_{fi}$

we get :

$$\sum -\delta_k w_{fi} \frac{\partial o_f}{\partial net_j}$$

$$net_j = \sum w_f x_f = w_0 x_0 + w_1 x_1 + \dots +$$

$$\frac{\partial net_j}{\partial x_f} = w_0 + w_1 + \dots + \frac{\partial}{\partial x_f} = \sum w_f$$

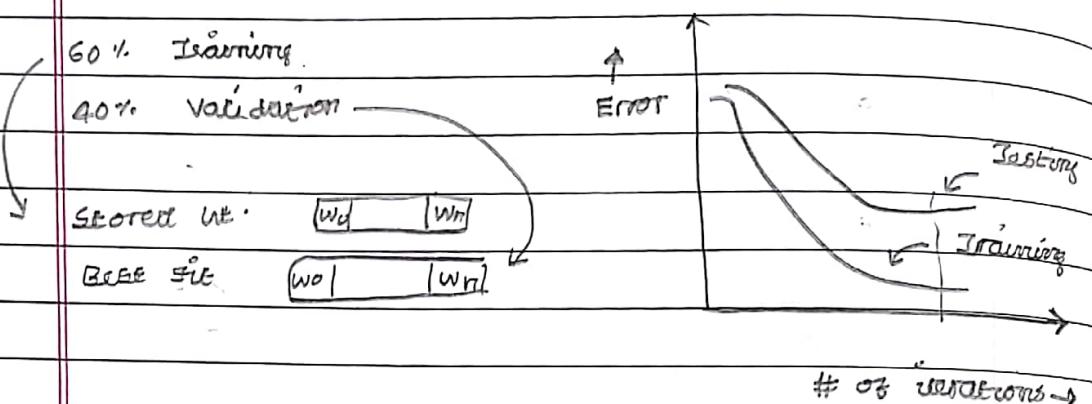
K-cross Validation : USED when # of training examples is less than # of testing examples
DATE: PAGE: is less

i) 100 examples .

ii) Divide 100 examples into k-set of 25

iii) Some mean .

calculating ideal # of iterations:



Bias in Neural Network?

We assume that there exists a smooth surface to fit the data set.

Machine

Magnitude of \vec{OA}

=

NOTE: of \vec{OA}

=

$$\sqrt{(4-0)^2 + (3-0)^2}$$

=

$$\sqrt{16+9}$$

=

$$\sqrt{25}$$

=

$$5$$

y-axis

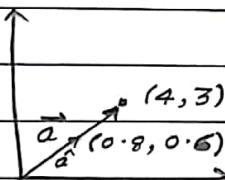
O



x-axis

Direction of \vec{OA} (or \vec{a}) = $(\cos \theta, \cos \alpha)$

$$= [\cos(\text{angle made with horizontal axis}), \cos(\text{angle made with vertical axis})]$$



$$= \left[\frac{4}{5}, \frac{3}{5} \right] = [0.8, 0.6]$$

$$= \left[\frac{a_1}{\sqrt{a_1^2 + a_2^2}}, \frac{a_2}{\sqrt{a_1^2 + a_2^2}} \right]$$

$$\hat{a} = \left[\frac{a_1}{\|a\|}, \frac{a_2}{\|a\|} \right]$$

NOTE: Direction vector of \vec{a} is a unit vector.

$$\hat{a} = [0.8, 0.6]$$

$$\|\hat{a}\| = \sqrt{(0.8-0)^2 + (0.6-0)^2}$$

$$= \sqrt{0.8^2 + 0.6^2}$$

$$= \sqrt{0.64 + 0.36}$$

$$= \sqrt{1.0}$$

$$= 1$$

Given a vector $\vec{a} (a_1, a_2)$, calculate :

i) Magnitude / norm of \vec{a} ,

$$\|\vec{a}\| = \sqrt{a_1^2 + a_2^2}$$

ii) Unit vector of \vec{a} ,

$$\hat{a} = \left(\frac{a_1}{\|\vec{a}\|}, \frac{a_2}{\|\vec{a}\|} \right)$$

$$\text{Also, } \vec{a} = \|\vec{a}\| \hat{a}$$

Dot Product :

$$\cos \alpha = \frac{a_1}{\|\vec{a}\|}, \cos \beta = \frac{b_1}{\|\vec{b}\|}$$

$$\cos(\beta - \alpha) = \cos \theta$$

$$= \cos \beta \cos \alpha \\ +$$

$$\sin \beta \sin \alpha$$

$$= \frac{b_1}{\|\vec{b}\|} \frac{a_1}{\|\vec{a}\|} + \frac{b_2}{\|\vec{b}\|} \frac{a_2}{\|\vec{a}\|}$$

$$\cos \theta = \frac{b_1 a_1 + b_2 a_2}{\|\vec{b}\| \|\vec{a}\|}$$

$$\cos \theta \|\vec{b}\| \|\vec{a}\| = b_1 a_1 + b_2 a_2 = \text{Product of 1st co-ordinates of } \vec{a} \& \vec{b}$$

= $\sum b_i a_i$

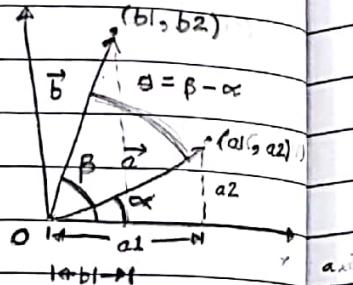
+ Product of 2nd co-ordinates of \vec{a} & \vec{b}

= Dot Product

or

Inner Product

Sum of product of elements of 1st row & 1st column.



Projection

is
in the

Direction

Example :

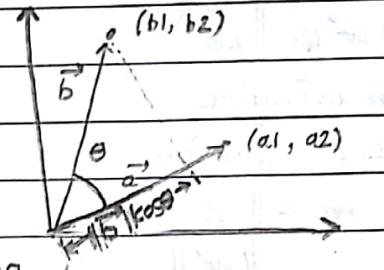
$\vec{w} =$ (10, 10)

(1) Projection

Projection of a vector on another vector :

Projection of \vec{b} on \vec{a}

is in the direction of \vec{a} .



$$\text{Direction of } \vec{a}, \hat{a} = \left[\frac{a_1}{\|\vec{a}\|}, \frac{a_2}{\|\vec{a}\|} \right]$$

$$\text{Projection of } \vec{b} \text{ on } \vec{a} = \|\vec{b}\| \cos \theta \cdot \hat{a}$$

= Magnitude Direction

$$\|\vec{a}\| \times \|\vec{b}\| \cos \theta \cdot \hat{a}$$

$$\|\vec{a}\| \quad \|\vec{a}\|$$

$$= 1 - (\vec{a} \cdot \vec{b}) \vec{a} = (\vec{a} \cdot \vec{b}) \hat{a}$$

Example : Total projection of

\vec{a} on \vec{w}

$$\hat{w} = \left(\frac{w_1}{\|\vec{w}\|}, \frac{w_2}{\|\vec{w}\|} \right) = \left(\frac{2}{\sqrt{5}}, \frac{1}{\sqrt{5}} \right)$$

product of 2nd

o-ordinate of \vec{a} & \vec{b}

$$\|\text{Projection of } \vec{a} \text{ on } \hat{w}\| = \|\vec{a}\| \cos \theta$$

$$= \frac{\|\vec{a}\| \|\hat{w}\| \cos \theta}{\|\hat{w}\|}$$

$$= \vec{a} \cdot \hat{w}$$

$$= \|\hat{w}\|$$

$$= 3 \left(\frac{2}{\sqrt{5}} \right) + 4 \left(\frac{1}{\sqrt{5}} \right)$$

$$\frac{1}{1}$$

$$= \frac{10}{\sqrt{5}}$$

Support Vector Machine DATE: PAGE:

Minimizing

$\|\vec{w}\|$

to

maximize

the margin

Critical

Point

γ -Point

$m = 2$

$\|\vec{w}\|$

using

Lagrange

Multipier:

Objective is

$$f(w, b) = \frac{\|\vec{w}\|^2}{2} \quad \text{Primal Function}$$

constraint

$$g(w, b) = y_i(\vec{w}x_i + b) - 1 \geq 0$$

Solutions: Note that we are trying to minimize the value of $\|\vec{w}\|$. $\therefore \vec{w}$ is a variable.

According to Lagrange,

$$f(w, b) = \alpha [g(w, b)]$$

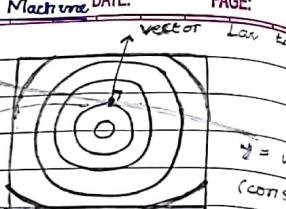
where

α : Lagrange Multiplier

Lagrange representation in the form of above eqn:

$$\Rightarrow f(w, b) - \alpha g(w, b) = 0$$

$$L(w, b, \alpha) = \nabla_{(w, b)} f(w, b) - \alpha \nabla_{(w, b)} g(w, b) = 0$$



DATE: PAGE:

$y = \vec{w}x + b$

We know that \vec{w} is also L(w) to the line $\vec{w}x + b$. \vec{w} is the vector L(w) to the point will be parallel.

$$\begin{aligned} L(w, b, \alpha) &= \frac{1}{2} \|\vec{w}\|^2 - \sum_{i=1}^m \alpha_i (y_i(\vec{w}x_i + b) - 1) = 0 \\ &= \frac{w_i^2}{2} - \sum_{i=1}^m \alpha_i (y_i(\vec{w}x_i + b) - 1) = 0 \quad (1) \\ &\approx 0^2 \end{aligned}$$

Finding Change in expression (1) w.r.t. w

$$\nabla_w L(w, b, \alpha) = \nabla_w f - \alpha \nabla_w g$$

$$\text{where } f = \frac{1}{2} \|\vec{w}\|^2 = \frac{w_i^2}{2}$$

$$g = \sum_{i=1}^m y_i (\vec{w}x_i + b) - 1$$

$$\nabla_w f = \frac{\partial f}{\partial w} = \sum_{i=1}^m w_i$$

$$\nabla_w g = \frac{\partial g}{\partial w} = \sum_{i=1}^m y_i x_i$$

$$\therefore \nabla_w L(w, b, \alpha) = w_i - \alpha_i \sum_{i=1}^m y_i x_i - 1 \quad (2)$$

Finding Change in expression (1) w.r.t. b

$$\Rightarrow w_i = \alpha_i \sum_{i=1}^m y_i x_i = 0$$

$$\nabla_b L(w, b, \alpha) = \nabla_b f - \alpha \nabla_b g$$

$$\text{where } f = \frac{1}{2} \|\vec{w}\|^2 = \frac{w_i^2}{2}$$

$$g = \alpha_i \sum_{i=1}^m [y_i (\vec{w}x_i + b) - 1]$$

$\nabla_b F$

$$= \frac{\partial F}{\partial b} = 0 \rightarrow \nabla_b g = \frac{\partial g}{\partial b}$$

$$= \frac{\partial}{\partial b} \left(\alpha_i \sum_{i=1}^m [y_i(w_i x_i + b) - 1] \right)$$

$$= \alpha_i \sum_{i=1}^m [y_i]$$

$$= \alpha_i \sum_{i=1}^m y_i$$

$$\therefore \nabla_b L = \nabla_b F - \nabla_b g = 0$$

$$= 0 - \alpha_i \sum_{i=1}^m y_i = 0$$

$$\Rightarrow \boxed{\alpha_i \sum_{i=1}^m y_i = 0} - (3)$$

From (1),

$$L(w, \alpha, b) = \frac{w_i^2}{2} - \alpha_i \sum_{i=1}^m [y_i(w_i x_i + b) - 1] = 0$$

$$\Rightarrow \text{using } = \frac{w_i^2}{2} - \alpha_i \sum_{i=1}^m y_i w_i x_i - \alpha_i \sum_{i=1}^m y_i b - \alpha_i \sum_{i=1}^m 1$$

$$\text{using (2), } w_i = \sum_{i=1}^m \alpha_i y_i x_i$$

$$\therefore (3), \alpha_i \sum_{i=1}^m y_i = 0$$

$$= \left[\sum_{i=1}^m \alpha_i y_i x_i \right]^2 - \alpha_i \sum_{i=1}^m y_i w_i x_i - b \alpha_i$$

$$= \frac{1}{2} \left[\sum_{i=1}^m \alpha_i y_i x_i \right]^2 - \alpha_i \sum_{i=1}^m y_i w_i x_i - 0 -$$

$$\sum_{i=1}^m \alpha_i$$

0 using (3)

$$\sum_{i=1}^m y_i - \sum_{i=1}^m \alpha_i$$

$$\Rightarrow \Theta(\alpha) = \frac{1}{2} \left[\sum_{i=1}^m \alpha_i y_i x_i \right]^2$$

$$\text{Now } L = \left[\sum_{i=1}^m \alpha_i y_i x_i - \sum_{j=1}^m \alpha_j y_j x_j - \sum_{i=1}^m \alpha_i \right]$$

$$\{ \text{So } w_i = \sum_{i=1}^m \alpha_i x_i \}$$

$$= \frac{1}{2} \left[\sum_{i=1}^m \alpha_i y_i x_i - \sum_{j=1}^m \alpha_j y_j x_j \right]$$

$$\sum_{i=1}^m \alpha_i y_i x_i - \sum_{j=1}^m \alpha_j y_j x_j$$

$$\sum_{i=1}^m \alpha_i = 1 \quad + (1-\alpha_i)$$

$$= \sum_{i=1}^m \alpha_i - \frac{1}{2} \left[\sum_{i=1}^m \alpha_i y_i x_i - \sum_{j=1}^m \alpha_j y_j w_j \right]$$

$$- \alpha_i \sum_{i=1}^m 1$$

$$\Theta(\alpha) = \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y_i y_j x_i x_j$$

0 using (3)

$$\sum_{i=1}^m y_i - \sum_{i=1}^m \alpha_i$$

$$\sum_{i=1}^m \alpha_i$$

Done using $\Theta(\alpha)$

Sequential to minimize $\|\vec{w}\|$

Minimization which maximizes margin m.

Optimization

Algo

(SMA Algo)

out of syllabus.

Consider the sample data :

	α_1	α_2	True Label (y_i)
x_1	1	3	+1 (y_1)
x_2	3	1	-1 (y_2)

Find optimal plane using SVM to separate the data.

Sol: Using eqn 3,

$$\sum_{i=1}^2 \alpha_i y_i = 0$$

$$\therefore \text{Here, } \sum_{i=1}^2 \alpha_i y_i = 0$$

$$\alpha_1 y_1 + \alpha_2 y_2 = 0$$

$$\alpha_1 (+1) + \alpha_2 (-1) = 0$$

$$\alpha_1 - \alpha_2 = 0$$

$$\boxed{\alpha_1 = \alpha_2} \quad \text{--- (1)}$$

$$\theta(\alpha) = \sum_{i=1}^2 \alpha_i - \frac{1}{2} \left[\sum_{i=1}^2 \sum_{j=1}^2 \alpha_i \alpha_j y_i y_j x_i x_j \right]$$

$$= [\alpha_1 + \alpha_2] - \frac{1}{2} \left[\begin{array}{l} \cancel{\alpha_1 \alpha_1 y_1 y_1 x_1 x_1} \\ \cancel{\alpha_1 \alpha_2 y_1 y_2 x_1 x_2} \\ \cancel{\alpha_2 \alpha_1 y_2 y_1 x_2 x_1} \\ \cancel{\alpha_2 \alpha_2 y_2 y_2 x_2 x_2} \end{array} \right]$$

$$\alpha_1 \alpha_1 y_1 y_1 x_1 x_1 \\ +$$

$$\alpha_1 \alpha_2 y_1 y_2 x_1 x_2 \\ +$$

$$\alpha_2 \alpha_1 y_2 y_1 x_2 x_1 \\ +$$

$$\alpha_2 \alpha_2 y_2 y_2 x_2 x_2$$

Using $\alpha_1 = \alpha_2$,

$$= (\alpha_1 + \alpha_2) - \frac{1}{2} \left[\alpha_1^2 y_1^2 x_1^2 + 2\alpha_1^2 y_1 y_2 x_1 x_2 + \alpha_1^2 y_2^2 x_2^2 \right]$$

$$= 2\alpha_1 - \frac{1}{2} \left[\alpha_1^2 (+1)^2 (1, 3)(1, 3) \right]$$

+

$$2\alpha_1^2 (-1)(+1)(1, 3)(3, 1)$$

+

$$\alpha_1^2 (-1)^2 (3, 1) \left[\begin{matrix} \\ (3, 1) \end{matrix} \right]$$

$$= 2\alpha_1 - \frac{1}{2} \left[\alpha_1^2 (1+9) + 2\alpha_1^2 (-1)(3+3) + \alpha_1^2 (9+1) \right]$$

$$= 2\alpha_1 - \frac{1}{2} [10\alpha_1^2 - 12\alpha_1^2 + 10\alpha_1^2]$$

$$= 2\alpha_1 - \frac{8\alpha_1^2}{2} = 0$$

$$\Rightarrow -8\alpha_1^2 + 2\alpha_1 = 0$$

$$2\alpha_1 [1 - 4\alpha_1^2]$$

$$\Rightarrow -4\alpha_1^2 = 2\alpha_1 \Rightarrow 2\alpha_1 - 4\alpha_1^2 = 0$$

$$\Rightarrow \alpha_1 = 2$$

$$\Rightarrow 2\alpha_1 [1 - 2\alpha_1] = 0$$

x_1^0

2 roots are

$\alpha_1 = 0$
or
$\alpha_1 = 1/2$

$$\alpha_1 \alpha_1 y_1 y_1 x_1 x_1$$

+

$$\alpha_1 \alpha_2 y_1 y_2 x_1 x_2$$

maximizing $(\alpha(\alpha))$

\Rightarrow Diff wrt α

$$-8\alpha_1 + 2 = 0$$

$$\Rightarrow -8\alpha_1 = -2$$

$$\Rightarrow \alpha_1 = 1/4$$

$$\text{At } \alpha_1 = 1/4,$$

$$\theta(\alpha) = 2\left(\frac{1}{4}\right) - \frac{1}{2} \times 8 \times \left(\frac{1}{4}\right)^2$$

$$= \frac{1}{2} - \frac{4 \times 1 \times 1}{4 \cdot 4}$$

$$= \frac{1}{2} - \frac{1}{4}$$

$$= \frac{1}{2} - \frac{1}{4} > 0$$

\therefore at $\alpha_1 = 1/4$, $\theta(\alpha)$ is max.

$$\therefore \alpha_1 = \alpha_2 = 1/4$$

Calculating w ,

$$w = \sum y_i x_i \alpha_i$$

$$= y_1 x_1 \alpha_1 + \alpha_2 y_2 x_2$$

$$= (+1)(1, 3)(1) + \left(\frac{1}{4}\right)(-1)(3, 1)$$

$$= \left(\frac{1}{4}, \frac{3}{4}\right) - \left(\frac{3}{4}, \frac{1}{4}\right)$$

$$= \left(-\frac{1}{2}, \frac{1}{2}\right)$$

$$\vec{w} = \left(-\frac{1}{2}, \frac{1}{2}\right)$$

$$\|\vec{w}\| = \sqrt{\frac{1}{4} + \frac{1}{4}} = \sqrt{\frac{2}{4}} = \frac{1}{\sqrt{2}}$$

$$m = \frac{\|\vec{w}\|}{\|\vec{u}\|} = \frac{1}{2} = \frac{1}{\sqrt{2}}$$

Eqn of 2 hyperplanes are

$$w_1x_1 + b = +1$$

$$+ w_2x_2 + b = -1$$

$$w_1x_1 + w_2x_2 + 2b = 0$$

$$2b = - \left[\left(-\frac{1}{2}, \frac{1}{2} \right) (1, 3) + \left(-\frac{1}{2}, \frac{1}{2} \right) (3, 1) \right]$$

$$2b = - \left[\left(\frac{-1+3}{2}, \frac{1}{2} \right) + \left(\frac{-3+1}{2}, \frac{1}{2} \right) \right]$$

$$2b = - \left[1 - 1 \right] = 0$$

$$b = 0$$

Hyperplane eqn: $wx + b = 0$

$$\Rightarrow (w_1, w_2) \begin{pmatrix} a_1 \\ a_2 \end{pmatrix} + d = 0$$

$$\Rightarrow \left(-\frac{1}{2}, \frac{1}{2} \right) (a_1, a_2) = 0$$

$$\Rightarrow -\frac{a_1}{2} + \frac{a_2}{2} = 0$$

$$\Rightarrow \boxed{a_2 = a_1}$$

Hyperplane 1: $wx + b = +1$

$$\left(-\frac{1}{2}, \frac{1}{2} \right) (a_1, a_2) + 0 = +1$$

$$-\frac{a_1}{2} + \frac{a_2}{2} = +1$$

$$\boxed{a_2 = 2 + a_1}$$

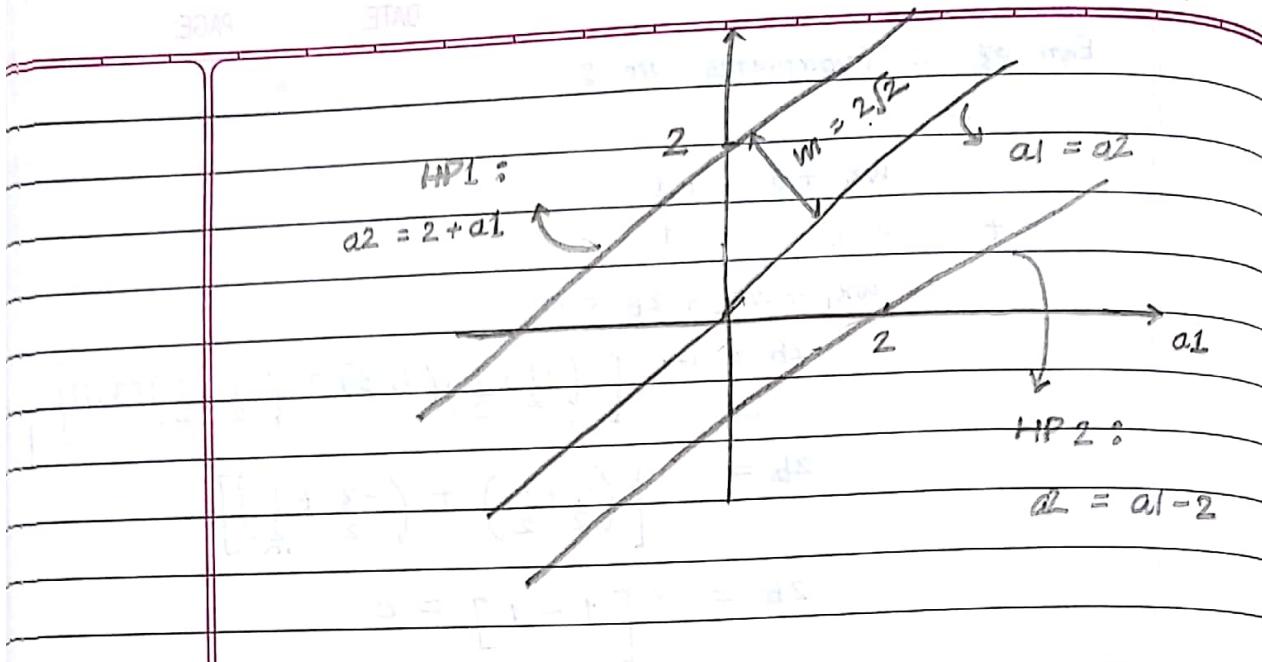
∴ 2: $wx + b = -1$

$$\left(-\frac{1}{2}, \frac{1}{2} \right) (a_1, a_2) + 0 = -1$$

$$\boxed{a_2 = a_1 - 2}$$

3049

3740



$$a_2 = a_1 - 2 \Rightarrow a_1 = a_2 + 2$$

$$a_2 = a_1 - 2 \Rightarrow a_1 = a_2 + 2$$

$$a = b_0 + (18) (\text{SW}) \quad \text{HPL}$$

$$a = b_0 + 18 \quad \text{SW}$$

$$i = 0 = (20, 10) (1, 1)$$

SW

$$a = b_0 + 18 - \frac{2}{2}$$

$$a = b_0 + 18 - [10 - 20]$$

$$18 - 10 = 8$$

$$L = 8 + 20 \Rightarrow L = 28$$

$$L = 8 + (20 - 18) (1, 1)$$

$$L = 8 + 2(1, 1) \Rightarrow L = 10 + 8 = 18$$

$$10 + 8 = 18$$

$$L = 8 + 20 \Rightarrow L = 28$$

$$L = 8 + (20 - 18) (1, 1)$$

$$10 + 8 = 18$$

Genetic Algorithm

DATE:

PAGE:

- i) Search Based Technique
- ii) Optimization is reqd., we use Genetic

Application:

1. Neural N/w

Chromosome

2. Economics is also called

3. Image Processing

as

4. Parallelization individual

5. Robotics split.

Computational Background:

1 0 1 0 1 0 ← P1

↓ Chromosome P2

P3

P4

Gene

0

Value /

Allele

1. Bit Pattern

0/1

2. Real Values

Population

{ Chromosomes }

Objective Function referred as Fitness Function.

→ is the sum of all the

Fitness Value

calculated

for every

chromosome.

Reqd. to get the

next set of (population)

chromosomes

2 methods:

1. CROSS-OVER

2. Mutation

1. Population / Set of Chromosomes / Set of Individual strings :

Initialization is done. How? Randomly.

2. Compute the fitness value for every individual string/chromosome.

3. Select the best parent chromosomes :

One of the following ways can be chosen :-

i) Roulette wheel

ii) Tournament Technique

iii) Rank Selection

iv) Random " is used in Tournament

Tournament Technique.

4. Cross over : new offsprings

5. Mutation is done :- a) Single Point mutation cross over

b) Multipoint " "

c) Uniform " "

5. Mutation : i) Bit flip ii) Swap iii) Scramble (Shuffle) iv) Reverse

6. Replacement Policy : To replace the new offsprings in the old population to generate next "

Repeat the steps from 2 to 6.

Until the objective is satisfied.

Step 1 : Initialization of Population :

can be done either using bit patterns :-

1	0	0	1	0
1	0	1	0	0
:				

or

can be done using Real Numbers :-

0.3	0.2	0.0	0.3
0.7	0.6	1.2	0.0
:			

Step 2 : For each of these chromosomes, we calculate the fitness value. Suppose, the fitness $f(x)$ is

$$ax + by + c = 0$$

And population size : P_1

	1	2	3
P_2	3	4	2
P_3	4	9	8
P_4	3	2	6
P_5	6	9	8
P_6	3	2	0

Corresponding fitness value be : $0.25 - \textcircled{1}$

$0.29 - \textcircled{2}$

$0.3 - \textcircled{3}$

$0.23 - \textcircled{4}$

$0.05 - \textcircled{5}$

$1 - \textcircled{6}$

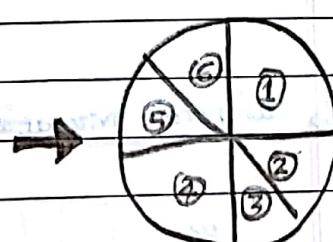
Step 3 : Choosing Parents :

i) Roulette wheel :

Role the wheel twice

& whichever comes

become the parents.



Suppose Parent 1 = P_5

we got $\rightarrow z = P_1$

ii) Tournament Technique :

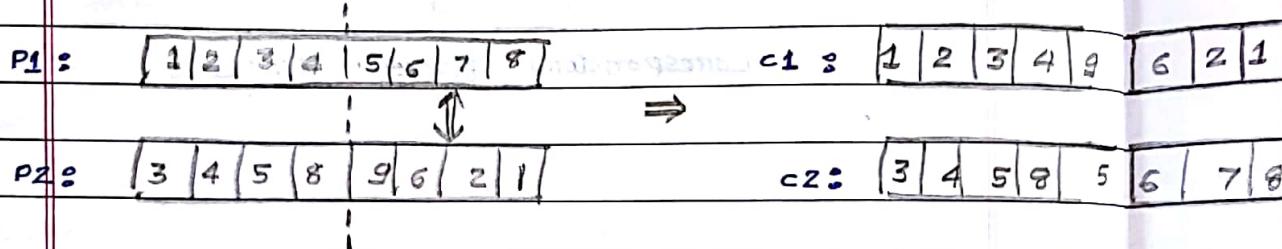
Randomly choose 2 chromosomes.

iii) Rank Selection:

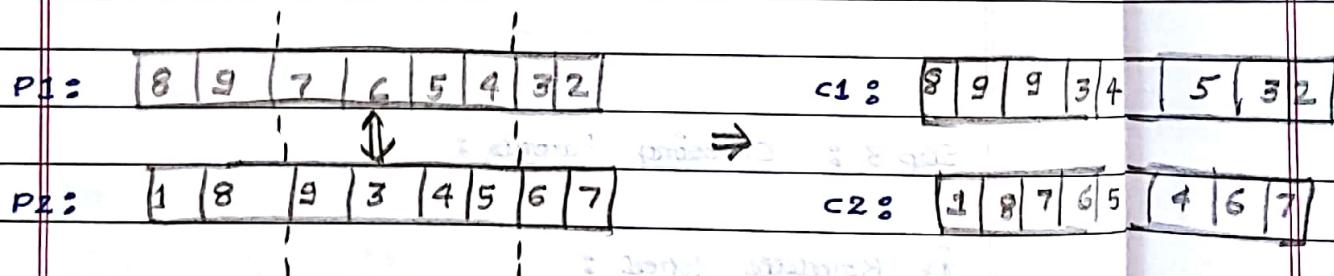
Arrange the fitness values in descending order.

Choose chromosomes corresponding to top 2 values as parents.

Step 4 : cross-over = i) Single Point Cross Over



ii) Multipoint Cross Over:



iii) Uniform Mutation: Every gene is treated independently.

Step

P1	0.2	0.2	0.5	0.5	0.3	0.3
----	-----	-----	-----	-----	-----	-----

P2	0.1	0.8	0.6	0.8	0.7	0.1
----	-----	-----	-----	-----	-----	-----

To get c_1 's 1st posn,
toss a coin. On tossing, if we get head
fill it with P_{15} 1st posn
else
 P_{25} 1st posn.

Suppose we tossed a coin 6 times & we got :

H T T H H H

Then $c_1 = [0.2 \quad 0.9 \quad 0.6 \quad 0.5 \quad 0.3 \quad 0.3]$

c_2 will follow

T H H T T T

6 2 1	0.1	0.2	0.5	0.8	0.7	0.1
-----------	-----	-----	-----	-----	-----	-----

6 | 7 | 8

iv) One of Syllabuses :

$$c_1 = ax + (1-a)y$$

$$c_2 = ax + (1-a)y$$

where, suppose $a = 0.5$

Step 5 : Mutation : applied to both the offsprings

1. BIT FLOP :

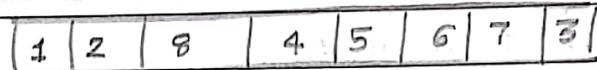
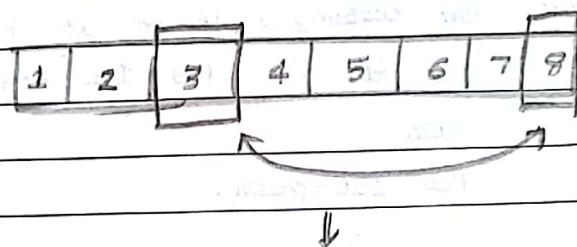
0	1	0	1	0	1	0	1
---	---	---	---	---	---	---	---



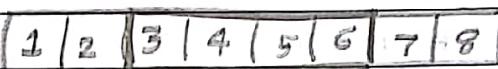
0	1	0	0	0	1	0	1
---	---	---	---	---	---	---	---

Position
is
chosen
randomly.

2. SWAP :



3. Scramble / shuffle :



shuffle



4. Reverse :



Step 5 : Replacement Policy

Termination ways : i) Terminate after fixed # of iterations.

ii) If previous & current population becomes

same i.e., there is no improvement

any more, then terminate.

iii) When the set threshold for objective

function fit is exceeded then terminate.

Numerical Problem :

Find 3 nos b/w 1 and 5 (both inclusive) such that

$$a + 2b + 3c = 15 \text{ where } a, b, c \text{ are the 3 nos.}$$

$$\text{Population size} = 4$$

Initialisation & Fitness Value (Step 1 & 2) :

	a	b	c	$ a+2b+3c - 15 = \text{Fitness Value}$
P1	5	2	3	$ 5+4+9-15 = 3$
P2	4	3	3	$ 4+6+9-15 = 4$
P3	2	5	2	$ 2+10+6-15 = 3$
P4	2	4	5	$ 2+8+15-15 = 10$

Step 3 : P4 has too high fitness value.

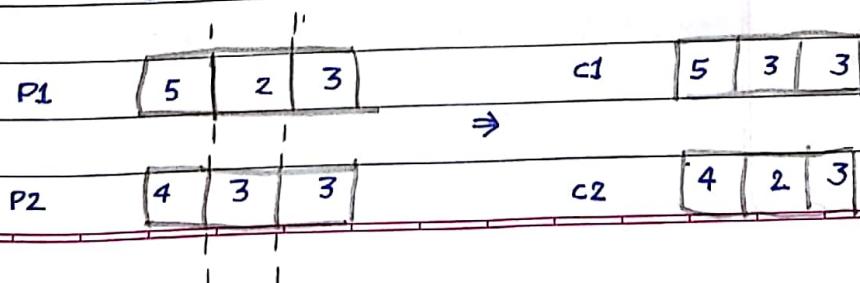
Roulette wheel b/w P1, P2 & P3

Suppose we get : P1

Now Roulette wheel b/w P2 & P3

Suppose we get : P2.

Step 4 : Using Single Point crossover b/w P1 & P2



17th October, 2017

little diff. from which we discussed

Here, we are finding the new population.

Technique is : i) keep the parents that made
dominant traits in offsprings

at least one in the offsprings

New Population :

P1: 5 1 2 2 3

P2: 2 4 1 3 3

C1 = P3: 5 3 3

C2 = P4: 4 2 3

Now make this population

	a	b	c	Fitness	Value
C1 = P3	5	1	2	'a + 2b + 3c - 15 = 0	4 + 6 + 9 - 15 = 4
C2 = P4	2	4	1		4 + 6 + 6 - 15 = 1
P1 = C1	5	3	1		5 + 6 + 3 - 15 = -11 = 1
C3 = P2	3	2	3		3 + 6 + 9 - 15 = 3

(Done by random replacement)

mutation technique

Also will repeat the process

Booster

Exp

Adaptive

Boosting

or

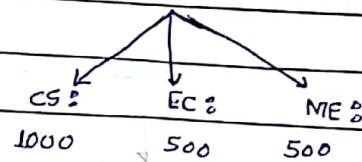
AdaBooster

17th October, 2017

Improving Performance DATE:
of (Binary) Classifier

PAGE:

Total # of structures in PESU : 2000



Suppose, 60% of students are recognized by Faculty F1

$$\downarrow = 600 \\ \text{Classifier 1, } 1000$$

$$\text{Performance on the entire dataset} = \frac{600}{2000} \\ = 30\%$$

30% of EC students are recognized by Faculty F2

$$\downarrow = 150 \\ \text{Classifier 2, } 500$$

$$\text{Performance on the entire dataset} = \frac{150}{2000} \\ = 7.5\%$$

35% of ME students are recognized by Faculty F3

$$\downarrow = 175 \\ \text{Classifier 3, } 500$$

$$\text{Performance on the entire dataset} = \frac{175}{2000} \\ = 8.75\%$$

Boosting : is Technique to improve performance on the entire dataset.

Ex: If it is a technique which improves performance by combination of hypothesis (classifiers) in a sequence

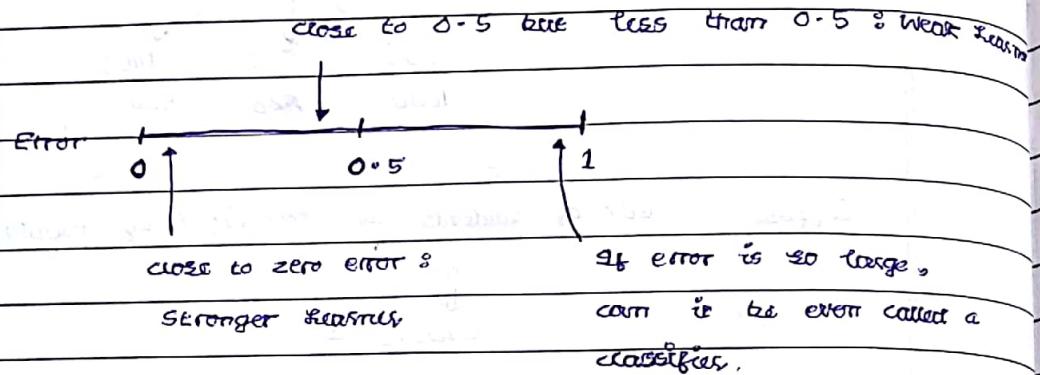
Adaptive
Boosting

or
Adaboost

Boosts the
performance

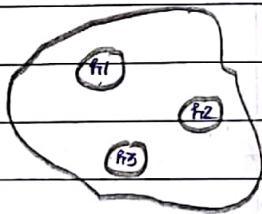
Different types of classifiers/learners based on the Error %

Error %



Combination of classifiers to improve performance :

If f_1, f_2, f_3 are three parts of dataset being wrongly classified by classifier c_1, c_2, c_3 respectively then



If we combine results of

classifier c_1, c_2 & c_3 then

Take the votes then,
we are likely to get less wrong classification.

can classify

$f_1 \& f_2$

correctly.

c_1 classify

$f_1 \& f_3$

correctly

$f_2 \& f_3$

correctly

less # of wrong classification means better performance.

Does combination of classifiers (generally odd # of classifiers)
always improve performance? No

2 out of

3 wrongly

classifiers.

So if you

go by majority

Vote on combination of
classifiers in this case

will give incorrect classification.

c1 & c2

wrongly

classify

this region.

But c3 classifier

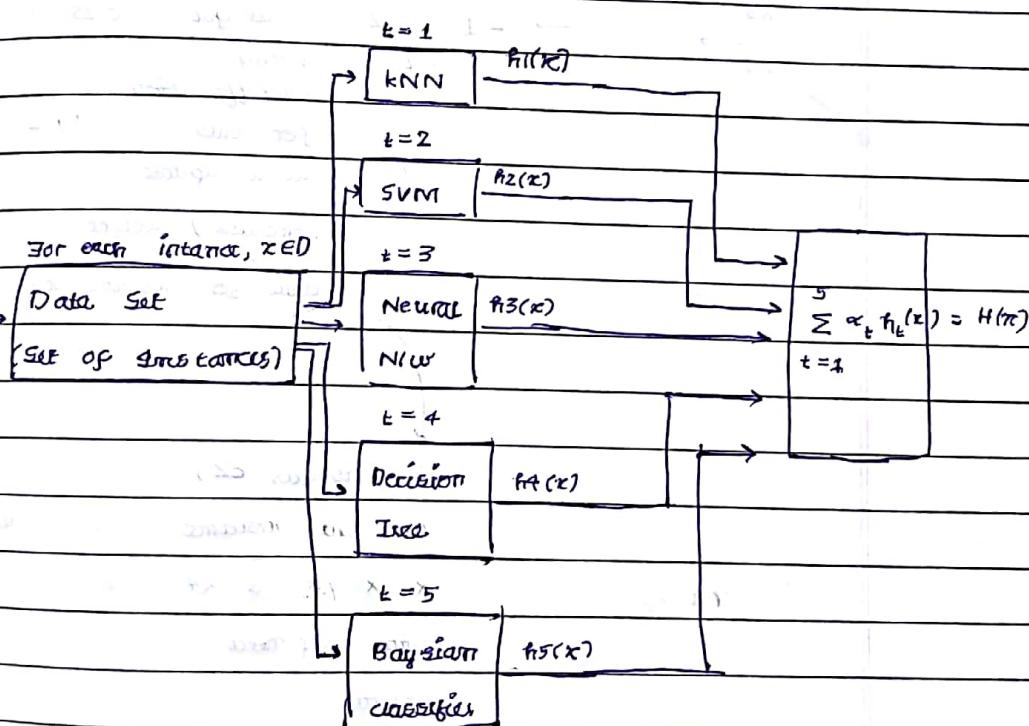
classifies this set

of examples correctly.

R1 R2

R3

How do we combine classifiers?



$$H(x) = \alpha_1 h_1(x) + \alpha_2 h_2(x) + \dots + \alpha_T h_T(x) = \sum_{t=1}^T \alpha_t h_t(x)$$

$$\text{where } \alpha_t = \frac{1}{2} \log_e \left(\frac{1 - E_t}{E_t} \right)$$

where E_t is the error

How is $(H(x))$ actually calculated?

Consider

True Label (y_i)

x_1	0.25	+1
x_2	0.25	-1
x_3	0.25	-1
x_4	0.25	+1

$$\text{Initial value of } w = \frac{1}{\# \text{ of instances}} = \frac{1}{4} = 0.25$$

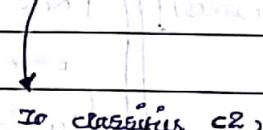
considering classifier c_1

	$f_i(x_i)$	\hat{y}_i	Error	Weight values (w)
x_1		+1	0	0.25
x_2		-1	0	0.25
x_3		-1	0	0.25
x_4		-1	2	0.25

" we got
wrong
classification
for this
we'll update

(increase) weight

value for instance x_4



To classifier c_2 ,

we feed instance

x_1, x_2, x_3 & x_4

with updated

weights

NOTE: ~~0~~ value is associated with instance values.

Ex: Out of 10 samples,
6 are correct
4 are wrong

$$\text{Then, } E = \frac{4}{10} = 0.4$$

$$\alpha = \frac{1}{2} \log_e \left(\frac{1-0.4}{0.4} \right) = \frac{1}{2} \log_e \left(\frac{0.6}{0.4} \right).$$

Plotting E_t vs α : $E_t \in [0, 1]$

$$\alpha_t = \frac{1}{2} \log_e \left(\frac{1-E_t}{E_t} \right) \quad \text{--- (2)}$$

If $E_t = 0$,

$$\alpha - \infty = \frac{1}{2} \log_e \left(\frac{1}{0} \right) = \infty$$

If $E_t = 1/2$,

$$\alpha = \frac{1}{2} \log_e \left(\frac{1-1/2}{1/2} \right) = 0$$

If $E_t > 0.5$,

$$\alpha = \frac{1}{2} \left[\log_e \left(\frac{\text{(small)}}{\text{(large)}} \right) \right] < 0$$

$$\left(\frac{\text{(small)}}{\text{(large)}} \right)^2 < 1 \quad \Rightarrow \quad \left(\frac{\text{(small)}}{\text{(large)}} \right) < 1$$

$$(1-\alpha)(\alpha) < (1-\beta)(\beta) \quad \Rightarrow \quad (1-\alpha)(\beta) < \alpha$$

$$T_{\text{min}} + S = 5$$

Derivation 2

$$w_t^i = \frac{w_t^i}{z} e^{-\alpha_t f_t(i) y_t^i}$$

Substituting α in above eqn,

$$w_{t+1}^i = \frac{w_t^i}{z} e^{-1/2 \log_e \left(\frac{1-E_t}{E_t} \right) f_t(i) y_t^i}$$

$$\Rightarrow w_{t+1}^i = \frac{w_t^i}{z} e^{-1/2}$$

NOTE:

$$f_t(i) y_t^i$$

$$\equiv \hat{y}(i) y_t^i$$

From eqn ②,

for correct classification
is always equal to 1

$$\alpha_t = \frac{1}{2} \log_e \left(\frac{1-E_t}{E_t} \right)$$

$$e^{-\alpha_t} = e^{-1/2 \log_e \left(\frac{1-E_t}{E_t} \right)} = \sqrt{\frac{1-E_t}{E_t}}$$

For

$$e^{-\alpha_t} = \left(\frac{1-E_t}{E_t} \right)^{1/2}$$

wrong classification,

$$z \Rightarrow (1-E_t) \sqrt{E_t} + E_t \sqrt{1-E_t} = \infty$$

$$= (1-E_t) \sqrt{\frac{E_t}{1-E_t}} + E_t \sqrt{\frac{1-E_t}{E_t}}$$

$$= (\sqrt{1-E_t})(\sqrt{E_t}) + (\sqrt{E_t})(\sqrt{1-E_t})$$

$$= 2(\sqrt{E_t})(\sqrt{1-E_t})$$

$$z = 2 \sqrt{E_t(1-E_t)}$$

Similarly,
 $\sum w_{t+1}^i = 1$
 wrong sample

Simplifying we get : $2 \sqrt{E_t(1-E_t)} \cdot \sum_{t=1}^T w_{t+1}^i = \sum_{t=1}^T w_t^i \sqrt{\frac{E_t}{1-E_t}}$

correct samples DATE:

PAGE: 42

$$\sum_{t=1}^T w_{t+1}^i = \sum_{t=1}^T w_t^i \sqrt{\frac{E_t}{1-E_t}} - \textcircled{a} \text{ for all correct}$$

$$\sum_{t=1}^T w_{t+1}^i = \sum_{t=1}^T w_t^i \sqrt{\frac{1-E_t}{E_t}} - \textcircled{b} \text{ for all wrong}$$

similarly
 $\sum w_{t+1}^i = 1/2$
only simple

\textcircled{a} + \textcircled{b}

$$= \left\{ \sum_{\substack{i \text{ correct} \\ \text{instances}}} w_{t+1}^i + \sum_{\substack{i \text{ wrong} \\ \text{instances}}} w_{t+1}^i \right\}$$

separation

at EO #1.

$$= \left\{ \sum_{t=1}^T w_t^i \sqrt{\frac{E_t}{(1-E_t)}} \right\}$$

correct instances

$$+ \left\{ \sum_{t=1}^T w_t^i \sqrt{\frac{(1-E_t)}{E_t}} \right\}$$

wrong instances

DATE: PAGE:

Example : Consider the classification training as follows :

x	0	1	2	3	4	5	6	7	8	9
$y:$ classification	+1	+1	+1	-1	-1	-1	+1	+1	+1	-1

check the accuracy against single hypotheses (classifiers)

use 3 classifiers

$$h_t(x_i) = \text{sign} \left\{ \sum_{t=1}^T \alpha_t f_t(x_i) \right\}$$

$$\# \text{ of instances } N = 10$$

Since we have 3 classifiers?

T = 3

50

$$H_k(x_i) = \text{sign} \left\{ \alpha_1 f_{k1}(x_i) + \alpha_2 f_{k2}(x_i) + \dots + \alpha_m f_{km}(x_i) \right\}$$

$$\alpha = \frac{1}{2} \log e \left(\frac{1 - E_t}{E_L} \right)$$

Assuming classifier 1, does the following classification:

$x \rightarrow$	0	1	2	3	4	5	6	7	8	9
$\hat{y} \rightarrow$	+1	+1	+1	-1	-1	-1	-1	-1	-1	-1

$$\text{split is } \frac{2+3}{2} = 2.5 \quad \text{i.e.,} \quad +\text{ve} \quad " \quad " \quad " \quad " \quad \text{right}$$

initial weights associated with each of these instances,

$$w = \frac{1}{N} = \frac{1}{10} = 0.1$$

$x_i \rightarrow$	0	1	2	3	4	5	6	7	8	9
$w_i \rightarrow$	0.1	0.1	0.1	0.1	0.1	0.1	0.1	0.1	0.1	0.1
$\hat{y}_i \rightarrow$	+1	+1	+1	-	-	-	-	-1	-	-
$u_i \rightarrow$	+	+	+	-	-	+	+	+	-	-

		\therefore we have 3
		$\therefore E_t = 0.1$
		$1 - E_t =$
		$So \propto$ value
		$\alpha_1 =$
		$=$
		$=$
		$\therefore 1st$ turn
		H
		Now updating
		$w_{updated} \sum w_i$
		for all
		+ve
		n_i will
		correct
	2.5	$w_{updated} [0.1]$
0.5	2.5	for +ve
		$w_{updated}$
		WS
		classification by classifier 1

∴ we have 3 wrong classification DATE: PAGE:

$$\therefore E_t = 0.1 + 0.1 + 0.1 = 0.3$$

$$1 - E_t = 1 - 0.3 = 0.7$$

So α value for classifier 1,

$$\alpha_1 = \frac{1}{2} \log_e \left(\frac{1-E_t}{E_t} \right)$$

$$= \frac{1}{2} \log_e \left(\frac{0.7}{0.3} \right)$$

$$= 0.424$$

∴ 1st term of

$$H_t(x_i) \Rightarrow \alpha_1 f_1(x_i)$$

becomes

$$0.424 f_1(x_i)$$

Now updating (decreasing w_i values in case classifier 1

C1 made correct classification &

increasing w_i values in case classifier 1

C1 made correct classification),

$$w_{\text{updated}} = \sum w_i = 1/2$$

for all

+ve x_i with

correct classification

$$w_{\text{updated}} [0.1 + 0.1 + 0.1 + 0.1 + 0.1 + 0.1 + 0.1] = 1/2$$

for 2.5
+ve

$$w_{\text{updated}} = \frac{1}{2 \times 0.7} = 0.0714$$

w for +ve samples of classifier 1, C1 = 0.1×0.714

$$= 0.0714$$

action by classifier 1

-①

$$w' \geq w_i = 1/2$$

for all

 x_i for which

classification 1

wrongly classified

$$w' [0.1 + 0.1 + 0.1] = 1/2$$

$$w' \times 0.1 \times 3 = 1/2$$

$$w' = \frac{1}{2 \times 3 \times 0.1} = \frac{1}{0.6} = \frac{10}{6} = 1.666$$

updated w values for the wrong classification

done by classifier 1,

$$0.1 \times 1.666$$

$$= 0.166 \quad - \textcircled{2}$$

Assuming classification done by 2nd classifier is

x_i	0	1	2	3	4	5	6	7	8	9
\hat{y}_i	+	+	+	+	+	+	+	+	-	-
w_i	+	+	+	<u>-</u>	<u>-</u>	<u>-</u>	+	+	+	-

of wrong classification

made by classifier 2 = 3

Using

~~$w_1 = 0.07 \ 0.07 \ 0.07 \ 0.07 \ 0.07 \ 0.07 \ 0.07 \ 0.07 \ 0.07$~~

~~$w_2 = 0.07 \ 0.07 \ 0.07 \ 0.07 \ 0.07 \ 0.07 \ 0.166 \ 0.166 \ 0.166 \ 0.07$~~

∴ we have 3 wrong classification,

$$E_t = 0.07 + 0.07 + 0.07$$

$$= 0.21$$

$$1 - E_t = 1 - 0.21 = 0.79$$

So α value for classifier 2,

$$\alpha_2 = \frac{1}{2} \log_e \left(\frac{1 - E_t}{E_t} \right) = \frac{1}{2} \log_e \left(\frac{0.79}{0.21} \right) = 0.653$$

\therefore 2 term of $H_L(x_i)$ is $\kappa_2 f_2(x_i)$

$$= 0.663 f_2(x_i)$$

On the next classification,

w initial values of w needs to be updated.

Before next classification happens

w value for $x_1, x_2, x_3, x_4, x_5, x_6$ & x_7

should be $1 \times \begin{pmatrix} 1 \\ 1 \\ 1 \\ 1 \\ 1 \\ 1 \\ 1 \end{pmatrix}$

$$2 [0.071 \times 4 + 0.166 \times 3]$$

= 1

$$2 [0.284 + 0.498]$$

$$= \frac{1}{2 \times 0.782}$$

$$= 0.639$$

$$\approx 0.64$$

$$w \text{ for } x_0, x_1, x_2, x_3, x_4, x_5, x_6, x_7 = 0.071 \times 0.64$$

$$x_8 = 0.045$$

$$w \text{ for } x_0, x_1, x_2, x_3, x_4, x_5 = 0.071 \times 2.34$$

$$= 0.16614$$

$$w \text{ for } x_6, x_7, x_8, x_9 = 0.166 \times 0.64$$

$$= 0.106$$

For classification by C3, (suppose split at 5.5)

x_i	0	1	2	3	4	5	6	7	8	9
-------	---	---	---	---	---	---	---	---	---	---

y_i	+	+	+	+	+	-	-	-	-
-------	---	---	---	---	---	---	---	---	---

y_i	+	+	+	-	-	-	+	+	+	-
-------	---	---	---	---	---	---	---	---	---	---

w_i	0.045	0.045	0.045	0.166	0.166	0.166	0.106	0.106	0.106	0.106
-------	-------	-------	-------	-------	-------	-------	-------	-------	-------	-------

Now # of correct classification = 4
 # of incorrect " " = 6

$$E_L = \frac{1}{6} (0.166 + 0.166 + 0.166 + 0.166 + 0.166 + 0.166) = 0.498 + 0.318 = 0.816$$

$$1 - E_L = 0.184$$

So α value for classifier 3,

$$\alpha_3 = \frac{1}{2} \log_e \left(\frac{1 - E_L}{E_L} \right)$$

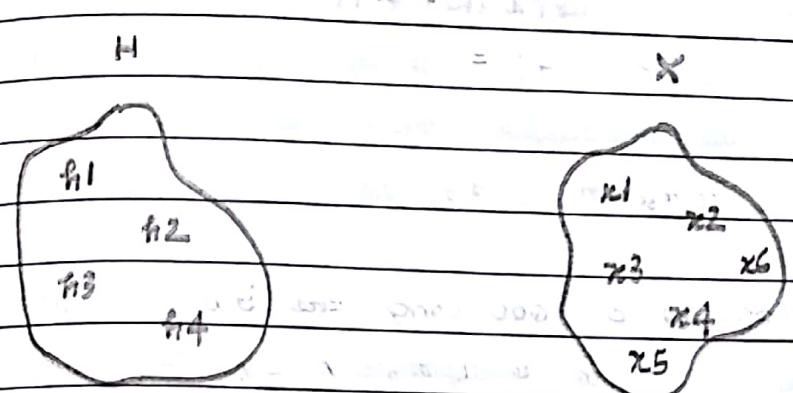
$$= \frac{1}{2} \times \log_e \left(\frac{0.184}{0.816} \right) = -0.744$$

$\therefore H_t$

$$= \begin{cases} 1 & \text{if } \alpha_1 h_1(x_i) + \alpha_2 h_2(x_i) + \alpha_3 h_3(x_i) \geq 0 \\ 0 & \text{otherwise} \end{cases}$$

$$= \text{sign} \{ 0.494 \alpha_1 h_1(x_i) + 0.663 \alpha_2 h_2(x_i) + -0.744 \alpha_3 h_3(x_i) \}$$

How is Bayesian learning different from something like concept learning w.r.t such as list them Eliminate or Candidate Elimination Algs?



In algorithms such as,

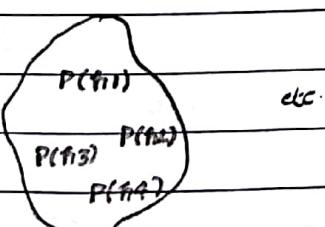
List them Eliminate or Candidate Elimination, from hypothesis set H we try to eliminate all $h \in H$ that is not consistent with all $x \in X$. The hypothesis left in H is referred to as Version Space.

In Bayesian Learning,

instead of eliminating $h \in H$ which are not consistent with all $x \in X$,

- i) we consider all $h \in H$
- ii) initialize each one them with equal probability
- iii) the req. Version Space will consist of all the hypotheses $h \in H$ with max. probabilities

H



Computation of Probabilities :

i) For all $h \in H$, initially $P(h) = \frac{1}{|H|}$
 Initially $P(h) = \frac{1}{|H|}$
 where $|H| = \# \text{ of hypotheses in } H$
 So, all hypotheses $h \in H$ are initialized with uniform probabilities.

ii) Probability of observing the data, $P(D)$
 This " is independent of the hypotheses
 iii) Prior probability, $P(D|h)$
 is the " that h has been observed
 For data D

Rule of Naive Method : This method assumes that the disease D used for training was/is COMPLETELY ERROR FREE
 To find our version space, we will calculate Probability of h given data set D for each $h \in H$ & then find the hypothesis (or set of hypotheses) having max. probability.

Computation of Posterior Probability of h given dataset D , $P(h|D)$
 = $\frac{P(D|h)P(h)}{P(D)}$
 {using Baye's Theorem }

We want to find,

$$\begin{aligned} \text{Maximum a posteriori hypothesis MAP} &= \arg\max_{h \in H} P(h|D) \\ &= \arg\max_{h \in H} \frac{P(D|h)P(h)}{P(D)} \\ h_{\text{MAP}} &= \arg\max_{h \in H} P(D|h)P(h) \quad \left[\because P(D) \text{ is constant} \right] \end{aligned}$$

we take it off.

we assumed that our training data was completely error free.

Numerical on Bayesian Learning :

0.008 of the population is suffering from cancer.
 A lab conducts test on patients to confirm this.
 When the patient has cancer, 98% of the time the lab results are +ve, & when the patient doesn't have cancer 97% of the time the lab results are -ve. What is the probability that the patient has cancer if the lab result is +ve?

patient has

$$P(\text{cancer}) = 0.008$$

$$P(\text{lab result} | \text{patient has cancer}) = 0.98$$

is +ve

$$P(\text{lab result} | \text{patient does not have cancer}) = 0.97$$

is -ve

$$P(\text{patient has cancer} | \text{lab result is +ve})$$

$$= \frac{P(\text{lab result is +ve} | \text{patient has cancer}) \times P(\text{patient has cancer})}{P(\text{lab result is +ve})}$$

$$\text{Observation: } P(h|D) \text{ for any } h \in H \propto P(D|h)$$

$$\propto P(h)$$

$$\propto \frac{1}{P(D)}$$

$$P(\text{lab result is +ve})$$

$$= P(\text{person has cancer}) P(\text{lab result is +ve} | \text{person has cancer}) + P(\text{person does not have cancer}) P(\text{lab result is +ve} | \text{person does not have cancer})$$

$$P(\text{lab result is +ve}) = P(\text{person has cancer}) P(\text{lab result is +ve} | \text{person has cancer}) + P(\text{person does not have cancer}) P(\text{lab result is +ve} | \text{person does not have cancer})$$

$$= 0.008(0.98)$$

+

$$(1 - 0.008)(1 - 0.97)$$

$$= 0.008(0.98) + (0.992)(0.03)$$

$$= 0.00784 + 0.02976$$

$$= 0.0376$$

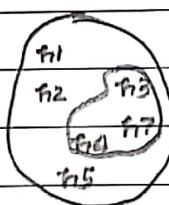
$\therefore P(\text{patient has cancer} | \text{lab result is } +ve)$

$$= \frac{0.98(0.008)}{0.0376}$$

$$= \frac{0.00784}{0.0376}$$

$$= 0.2085$$

Derivation :

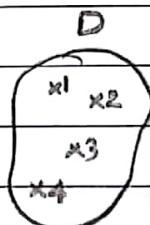


Suppose Version Space

VS contains

f_3, f_4 & f_7

then



$$P(D|f_7) = 1$$

for

all $x \in VS$

&

$$P(D|f_i) = 0$$

for

all $f_i \notin VS$

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

Initially, $P(A_i) = \frac{1}{\|H\|}$

$$\therefore P(A_i|D) = \frac{P(D|A_i) \times \frac{1}{\|H\|}}{P(D)}$$

$$= P(D|A_i) \times \frac{\frac{1}{\|H\|}}{\sum_{h \in VS} P(D|h) P(h) + \sum_{h \notin VS} P(D|h) P(h) \times \frac{1}{\|H\|}}$$

$$= P(D|A_i) \times \frac{\frac{1}{\|H\|}}{\sum_{h \in VS} 1 \cdot \frac{1}{\|H\|} + \sum_{h \notin VS} 0 \cdot \frac{1}{\|H\|}}$$

$$= P(D|A_i) \times \frac{\frac{1}{\|H\|}}{\frac{1}{\|H\|} \sum_{h \in VS} 1} \times \frac{1}{\|H\|}$$

$$= P(D|A_i) \times \frac{\frac{1}{\|VS\|}}{\frac{\|VS\|}{\|H\|}} \times \frac{1}{\|H\|}$$

$$= \frac{1}{\|VS\|} \quad \left\{ \because P(D|H) = 1 \right.$$

if constant

Bayesian Rule

(continued)

Now we calculated f_{MAP} under the assumption that the training data that we had was completely error free.

Now we would calculate f_{ML} , knowing that the training data contains error.

In case our training data is

x_i	Attr1	Attr2	...	True Label (d_i & not y_i)
x_0				d_0
x_1				d_1
x_2				d_2
:				:
x_n				d_n

Error could either be i) in the attribute value of a particular instance.

ii), true label of a particular instance.

While calculating f_{ML} we work under the assumption that few instances of the training data may have wrong true labels, i.e., for few examples in the training data,

$$d_i = f(x) + e_i$$

True label

which should

have been there

Lagrange's Theorem:

is used for minimizing / maximizing a multivalued

i) If $f(x, y)$ subject to constraint $g(x, y)$

ii) This theorem states that,

gradient vector of $f(x, y)$

is parallel to the

gradient vector of $g(x, y)$

iii) Lagrange's fn,

$$L(x, y) = f(x, y) - \alpha g(x, y) \quad \text{where } \alpha \geq 0$$

and

$$\nabla L = \nabla f - \alpha \nabla g$$

$$\nabla L = \nabla f - \alpha \nabla g$$

$L(x, y, \alpha)$ is expressed as $\Theta(\alpha)$

where x & y are expressed
in terms of α .

Also,

$$\nabla L = 0$$

$$\Rightarrow \nabla f - \alpha \nabla g = 0 \quad \text{where } \alpha \geq 0$$

Example 2 Minimize $f(x, y) = x^2 + y^2$

subject to

constraint

$$g(x, y) = y + x - 1 = 0$$

$$\text{Sot. } L(x, y) = f(x, y) - \alpha g(x, y)$$

where $\alpha \geq 0$...

&

$$\nabla L = \nabla f - \alpha \nabla g = 0$$

$$\text{So, } L(x, y) = (x^2 + y^2) - \alpha(y + x - 1)$$

$$\nabla f = \begin{bmatrix} 2x \\ 2y \end{bmatrix}$$

$$\nabla g = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

By Lagrange,

$$\nabla L = 0; \quad \alpha \geq 0$$

$$\Rightarrow \nabla f - \alpha \nabla g = 0$$

$$\Rightarrow \begin{bmatrix} 2x \\ 2y \end{bmatrix} - \alpha \begin{bmatrix} 1 \\ 1 \end{bmatrix} = 0$$

$$\Rightarrow 2x - \alpha = 0 \quad \& \quad 2y - \alpha = 0$$

$$\text{So, } x = \frac{\alpha}{2}$$

$$\text{So, } y = \frac{\alpha}{2}$$

$$\text{Expressing } L(x, y) = f(x, y) - \alpha g(x, y)$$

where $\alpha \geq 0$

as

$$\Theta(\alpha)$$

$$\text{So, } L(x, y) = (x^2 + y^2) - \alpha(y + x - 1)$$

$$\Theta(\alpha) = \left(\frac{\alpha^2}{4} + \frac{\alpha^2}{4} \right) - \alpha \left(\frac{\alpha}{2} + \frac{\alpha}{2} - 1 \right)$$

$$= \frac{\alpha^2}{4} + \frac{\alpha^2}{4} - \frac{\alpha^2}{2} - \frac{\alpha^2}{2} + \alpha$$

$$= \frac{\alpha^2}{2} - \frac{\alpha^2}{2} - \frac{\alpha^2}{2} + \alpha$$

$$= \alpha - \frac{\alpha^2}{2}$$

According to Lagrange,
maximizing $\Theta(\alpha)$

is

some α minimizing $F(x, y)$ Minimizing $\Theta(\alpha)$,

$$\Theta'(\alpha) = 0$$

$$\frac{1 - 2\alpha}{2} = 0$$

$$\boxed{\alpha = 1}$$

At $(x = \frac{\alpha}{2}, y = \frac{\alpha}{2})$, $f(x, y)$ is minimum

$\therefore \left(\frac{1}{2}, \frac{1}{2}\right)$ $f(x, y)$ is minimum.

subjected to constraint

$$g(x, y)$$

Ex 2: $f(x, y) = x^2y$ subjected to constraint $g(x, y) :$

$$x^2 + y^2 - 1 = 0$$

Sol. $L(x, y) = f(x, y) - \alpha g(x, y)$ where $\alpha \geq 0$

&

$$\nabla L = 0$$

$$\Rightarrow \nabla f - \alpha \nabla g = 0 \Rightarrow \begin{bmatrix} 2xy \\ x^2 \end{bmatrix} - \alpha \begin{bmatrix} 2x \\ 2y \end{bmatrix} = 0$$

$$\text{So, } 2x[y - \alpha] = 0 \quad \begin{array}{l} x^2 - 2\alpha y = 0 \\ x^2 = 2\alpha y \end{array}$$

$$\left(\frac{x}{\alpha}\right)^2 =$$

$$\Delta$$

$$x=0 \text{ or } [y = \alpha]$$

$$\alpha^2 = 2\alpha^2$$

$$\alpha^2 = 2\alpha^2$$

$$\alpha = \pm \sqrt{2}$$

$$\begin{aligned}
 \Theta(\alpha) &= 2\alpha^2(\alpha) - \alpha [2\alpha^2 + \alpha^2 - 1] \\
 &= 2\alpha^3 - \alpha [2\alpha^2 - 1] \\
 &= 2\alpha^3 - 2\alpha^3 + \alpha \\
 \Theta(\alpha) &= \alpha - \alpha^3
 \end{aligned}$$

Minimizing $f(x, y)$ w.r.t $g(x, y)$
is same as
maximizing $\Theta(\alpha)$

$$\begin{aligned}
 \Theta'(\alpha) &\Rightarrow 1 - 3\alpha^2 = 0 \\
 \Rightarrow 3\alpha^2 &= 1 \\
 \Rightarrow \alpha^2 &= \frac{1}{3} \\
 \Rightarrow \alpha &= \pm \frac{1}{\sqrt{3}}
 \end{aligned}$$

But $\alpha \geq 0$

$$\therefore \alpha = \frac{1}{\sqrt{3}}$$

\therefore At $(2\alpha^2, \alpha)$ $f(x, y)$ is min.
subjected to constraint
 $g(x, y)$

$\left(\frac{2}{3}, \frac{1}{\sqrt{3}}\right)$ $f(x, y)$ is min
subjected to constraint
 $g(x, y)$

(contd)

Given : $P(f_1|D) = 0.4$, going to +ve.

$P(f_2|D) = 0.3$, " " -ve.

$P(f_3|D) = 0.3$, " " -ve.

Suppose a new instance x_D is given.

If we use f_{MAP} , then classification of this sample is +ve.

classification of this sample is +ve.

$\because P(f_1|D)$ is max.

If we use "Bayes Optimal classifier" then, $V = \{+, -\}$

$$P(V_i|D) = \sum_{f_i \in H} P(V_i|f_i) \times P(f_i|D)$$

$$\begin{aligned} \text{So, } P(+|D) &= P(+|f_1) P(f_1|D) + P(+|f_2) P(f_2|D) \\ &\quad + P(+|f_3) P(f_3|D) \\ &= 1(0.4) + 0 + 0 \\ &= 0.4 \end{aligned}$$

$$P(+|D) = 0.4$$

$$\begin{aligned} P(-|D) &= P(-|f_1) P(f_1|D) + P(-|f_2) P(f_2|D) \\ &\quad + P(-|f_3) P(f_3|D) \\ &= 0(0.4) + 0.3(1) + 0.3(1) \\ &= 0.6 \end{aligned}$$

$$P(-|D) = 0.6$$

$P(+ D) = 0.4$	} max is -ve
$P(- D) = 0.6$	

\therefore Classification is -ve

Bayes Optimal classifier is better than f_{MAP} .

1. Choose one hypothesis ' h_i ' at random.
2. The classification of the new instance is w.r.t. the hypothesis ' h_i ' in Step 1.

Error $\left(\frac{\# \text{ of wrong classification}}{\text{Total } \# \text{ of }} \right)$ because of Gibbs Algorithm

cannot be more than 2 times Error in Bayes

Optimal

Classifer Hyp.

$$(0.1)^2 \times 0.1^2 + 0.1 = (0.1)^2$$

1/2

$$(0.1)^2 \times 0.1^2 + (0.1)^2 \times 0.1^2 = (0.1)^2$$