

28/8/21

Pedestrian algorithm:

Training experience: data set of pedestrians

(E) Task (T): Recognition of pedestrians crossing the road.

Performance measure: percentage of pedestrians

(P) classified accurately.
CNN: convolution, ReLU, Pooling (max, avg, global), fully connected layer

Data set:

Collection of many examples, where, an example is a collection of features.

Machine Learning Types:

1) Supervised (Guidance Eg: actual O/P already fed as data set).

2) Unsupervised (Eg: clustering of data set based on features)

3) Reinforcement

Supervised:

→ Classification, regression, brute force method → all features measured.

→ Machine learns examples along with class

labels or targets

→ Test data → O/P for a [given I/P], Training → O/I/P, I/I/P.

Classification → Which class it belongs to

Regression → Predicting a value in future.

→ O/P → real number, class label.

2

ML	DL	ML	DL required
data prep etc. giving feature ext.	directly features	→ IP required (already trained)	Classifies on its own more data required
→ Find patterns in data.	→ grouping similar feature objects together		
→ No class labels provided, learn prob. model of data.	→ CLUSTERING, decision making, future values.		
→ No supervised O/P target, no reward from environment.	→ Build representation of I/Ps.		

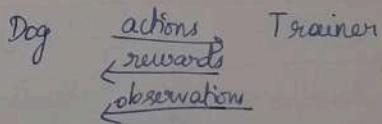
Reinforcement learning:

→ Goal oriented algorithm, knowledge to learn how to match a complex goal.

→ Automatically determine the ideal behavior within a specific context to maximize the computer's performance.

Eg: Dog, trainer

Action matches command → +ve reward
Action doesn't match " → -ve or 0 reward



→ Mapping b/w observation, action → [policy], which has to be made maximum.

→ Reward signal provided to evaluate the goodness of a trial and guide learning process.

Dog	Computer	Sample space:
→ Training is inside the dog's brain	→ Training is supervised by training algorithm	→ List of all possible outcomes
→ Agent is dog	→ Agent is computer	→ List is mutually exclusive, collectively exhaustive, and at the right granularity (noise / uncertainty)
→ Environment is surrounding the dog	→ Environment is everything outside the vehicle including car dynamics, weather, nearby vehicles etc.	
→ Observations is what the dog observes	→ Observations from car's sensors	
→ Action is how the dog reacts/responds	→ Action is generation of steering, braking, acceleration.	
→ Reward is getting a treat from the trainer of the dog	→ Reward is correct parking with right orientation, (a signal to mark the same).	
→ ML, probability go hand in hand.		

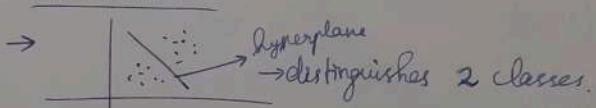
PROBABILITY AXIOMS:

Event: Subset of the sample space
- Probability is assigned to events.

AXIOMS FOR LEGITIMATE PROBABILITY:

Nonnegativity: $P(A) \geq 0$ consequences $P(A) \leq 1$
Normalisation: $P(\Omega) = 1$ sample space $P(\emptyset) = 0$.
(Finite) Additivity: If $A \cap B = \emptyset$, $P(A \cup B)$
(Disjoint). $= P(A) + P(B)$

HYPERPLANE:



AXIOMS:

- $P(A \cap A^c) = 0$ $\rightarrow P(A \cup B) = P(A) + P(B) - P(A \cap B)$
- $P(A \cup A^c) = 1$
- If $A \subset B$, $P(A) \leq P(B)$ $B = A \cup (B \setminus A')$
 $A \cap B = A'$

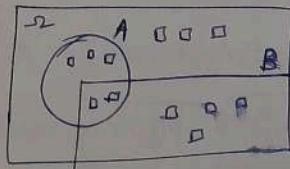
$$\begin{aligned}
 a &= P(A \cap B') \\
 b &= P(AB) \\
 c &= P(B \cap A') \\
 a+b+c &= P(A \cup B) \\
 &= P(A \cap B') + P(AB) + P(B \cap A')
 \end{aligned}$$

$$\begin{aligned}
 \text{Now, } P(A) + P(B) - P(AB) \\
 &= (a+b) + (b+c) - b \\
 &= a+b+c //
 \end{aligned}$$

$$\begin{aligned}
 \Leftrightarrow P(A \cup C \cap B) &= P(A \cap B) \cup (C \cap B) \quad \text{if } A, C \text{ are independent} \\
 \Leftrightarrow P(A \cup B \cup C) &= P(A) + P(A' \cap B \cap C)
 \end{aligned}$$

$$A \cup B \cup C = A \cup [B \cap A'] \cup [C \cap A' \cap B']$$

1/9/21 CONDITIONAL PROBABILITY

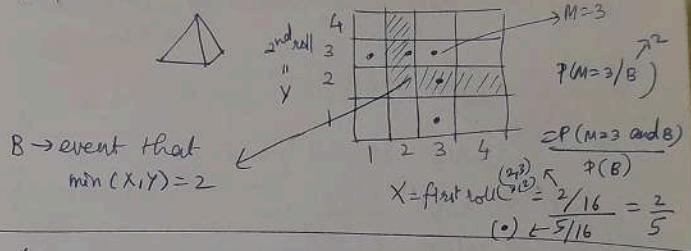


$$\begin{aligned}
 P(A) &= \frac{5}{12} \\
 P(B) &= \frac{6}{12}
 \end{aligned}$$

$$\begin{aligned}
 P(A/B) &= \frac{P(AB)}{P(B) \neq 0} \\
 P(A/B) &= \frac{2}{2+4} = \frac{1}{3} \Rightarrow P(AB) \\
 P(B/B) &= \frac{6}{6} = 1
 \end{aligned}$$



Q) Tetrahedral die rolled 2x. Let B event has occurred.

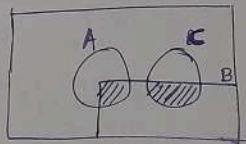


B → event that
min(X, Y) = 2

AXIOM:

$$P(A \cup C/B) = P(A/B) + P(C/B)$$

Proof:



$$\begin{aligned}
 P(A/B) + P(C/B) &= \frac{P(AB)}{P(B)} + \frac{P(CB)}{P(B)} \\
 &= \frac{P(AB) + P(CB)}{P(B)} \quad \text{--- (1)}
 \end{aligned}$$

$$P(A \cup C/B) = \frac{P(A \cup C \cap B)}{P(B)} = \frac{P(AB) + P(CB)}{P(B)}$$

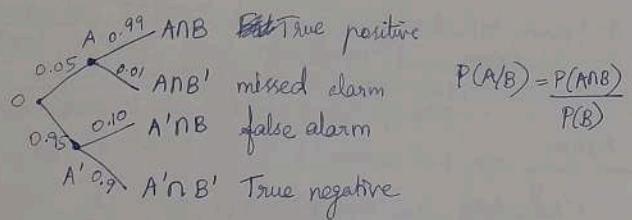
since $A \cap B$, $C \cap B$ are disjoint,

$$= \frac{P(AB) + P(CB)}{P(B)} \quad \text{--- (2)}$$

(1) = (2) proved.

Q) Event A : Aeroplane is flying above.
 Event B : Something registered on radar.

$$P(A|B) = \frac{P(AB)}{P(B)} ; P(A) = 0.05$$



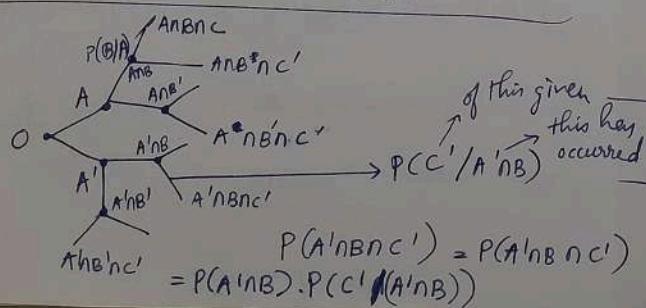
$$P(AB) = P(A)P(B/A) = 0.05 \times 0.99$$

$$P(B) = 0.05 \times 0.99 + 0.95 \times 0.1 = 0.145$$



$$P(A|B) = \frac{P(AB)}{P(B)} = \frac{0.05 \times 0.99}{0.145} = 0.34$$

2/9/21 MULTIPPLICATION RULE



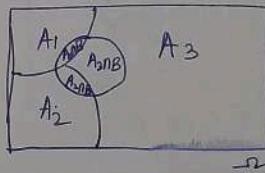
$$= P(A') P(B'|A') \cdot P(C'|A'NB)$$

$$\boxed{P(A \cap B) = P(B) P(A|B)}$$

(Def)
 $P(A) P(B|A)$

TOTAL PROBABILITY THEOREM

$$\boxed{P(B) = P(B|A_1)P(A_1) + P(B|A_2)P(A_2) + P(B|A_3)P(A_3)}$$



Sample space is partitioned into A_1, A_2, A_3
 $P(A_i)$ for every i

$$P(B) = P(A_1 \cap B) + P(A_2 \cap B) + P(A_3 \cap B)$$

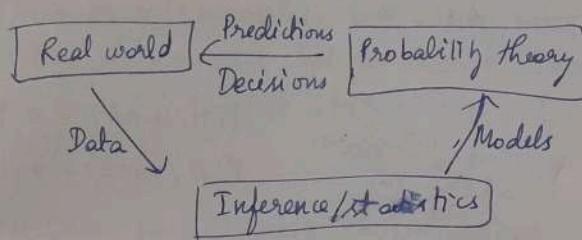
$$= P(A_1)P(B|A_1) + P(A_2)P(B|A_2) + P(A_3)P(B|A_3)$$

$$P(B) = \sum_i P(A_i)P(B|A_i) = P(B) = \sum_i P(A_i \cap B)$$

4/9/21

PROBABILITY THEOREM

- Framework to analyse phenomena with uncertain outcomes.
- Rules for consistent reasoning
- Predictions, decisions.

CONDITIONAL, BAYES THEOREM:

- 3 important rules
- Bayes' theorem → inference

$$P(A_i | B) = \frac{\{P(A_i) P(B | A_i)\}}{\sum_j P(A_j) P(B | A_j)}$$

$P(H) = \frac{3}{12}$
 $P(T) = \frac{4}{12}$

INFERENCE!

→ Initial beliefs $P(A_i)$ on possible causes of observed B .

$B \xrightarrow{\text{inference}} A_i \quad P(A_i | B) \quad (\text{revised belief})$
already B has occurred
what is $P(A_i)$?

from base $\leftarrow A_i \xrightarrow{\text{model}} B \quad P(B | A_i) \quad (\text{given } A_i, \text{ what is } P(B))$
partition $A_i, \text{ model } B$
 for every i , there will be $A_i, P(A_i) \rightarrow \text{initial belief}$

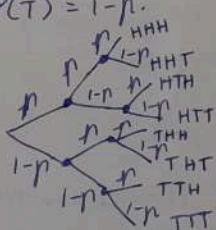
6/9/21

MODEL BASED ON CONDITIONAL PROB

Head, Tail.

$$P(H) = p$$

$$P(T) = 1-p$$



$$P(H) + P(T) = 1$$

(*independent, mutually exclusive, collectively exhaustive*)

$$\begin{aligned} P(1H) &= p(1-p)^2 + p(1-p)^2 + p(1-p)^2 \\ &= 3p(1-p)^2 \end{aligned}$$

$$\begin{aligned} P(\text{1st } H | 1H) &= \frac{P(1H)}{P(1H)} \\ &= \frac{p(1-p)^2}{3p(1-p)^2} = \boxed{\frac{1}{3}} \end{aligned}$$

=

$$\frac{p(1-p)^2}{3p(1-p)^2} = \boxed{\frac{1}{3}}$$

$$\frac{1}{3}$$

H_2 or T_2 is independent of H_1 or T_1 ,
(Getting 2nd place H/T is ind. of 1st place H/T).

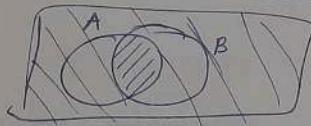
$$\therefore P(H_2 \text{ or } T_2 / H_1 \text{ or } T_1) = P(H_2 \text{ or } T_2 / T_1)$$

INDEPENDENT EVENTS:

$$P(A/B) = P(A) \quad \text{if } A, B \text{ independent.}$$

$$\Rightarrow P(A \cap B) = P(A)P(B)$$

\Leftrightarrow If A, B are independent, A, B' are independent.



$$A = (A \cap B) \cup (A \cap B')$$

$$P(A) = P(A \cap B) + P(A \cap B')$$

$$= P(A)P(B) + P(A)P(B') \quad [P(A \cap B) = P(A)P(B)]$$

$$= P(A) - P(A \cap B)$$

$$= P(A)P(B) + P(A) - P(A)P(B)$$

$$= P(A)$$

$$P(A \cap B') = P(A) - P(A \cap B)$$

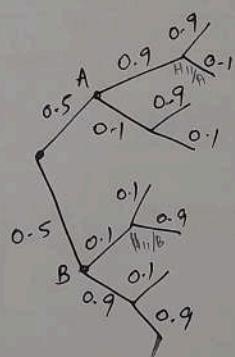
$$= P(A)[1 - P(B)]$$

$$\boxed{P(A \cap B') = P(A)P(B')}$$

$\rightarrow A, B$ are conditionally independent.
that is, A, B are not independent, but
 $(A/C), (B/C)$ are independent.

7/9/21 INDEPENDENCE VS. CONDITIONAL INDEPENDENCE

2 coins A, B. Are coin tosses independent?



$$P(H/\text{coin } A) = 0.9$$

$$P(H/\text{coin } B) = 0.1$$

$$P(\text{toss II} = H) = P(A)P(H_{II}/A) + P(B)P(H_{II}/B)$$

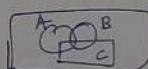
$$= 0.5 \times 0.9 + 0.5 \times 0.1 = 0.5$$

$P(\text{toss II} = H / \text{first 10 tosses are Heads})$

$$\approx P(H_{II}/A) = 0.9$$

\therefore A has more prob of H.

CONDITIONAL INDEPENDENCE:



$$\text{If } A, B \text{ independent, } P(A \cap B / C) = P(A/C) \cdot P(B/C)$$

\rightarrow Here A, B are not independent.

INDEPENDENCE OF COLLECTION OF EVENTS:

Eg: A_1, A_2, \dots are independent events

$$P(A_3 \cap A_4') = P(A_3 \cap A_4' / A_1 \cap A_2 \cap A_3) \\ P(A_1 \cap A_2 \cap A_3) = P(A_1)P(A_2)P(A_3)$$

PAIRWISE INDEPENDENCE:

If there are 3 events A_1, A_2, A_3
 Then, they are pairwise independent if $P(A_i \cap A_j) = P(A_i)P(A_j)$, $P(A_2 \cap A_3) = P(A_2)P(A_3)$, $P(A_1 \cap A_3) = P(A_1)P(A_3)$.

INDEPENDENCE VS. PAIRWISE INDEPENDENCE:

$H_1 \rightarrow$ first H
 $\vdots \vdots$

$$P(H_1) = P(H_2) = \frac{1}{2}$$

H_n, T_n are ~~pairwise~~ independent.

Event C \rightarrow 2 tosses having same result.
 (TT or HH).

$$P(H_1 \cap C) = P(H_1 \cap H_2) = \frac{1}{4} \\ P(H_1)P(C) = \frac{1}{2} \times \frac{1}{2} = \frac{1}{4} \quad \left. \begin{array}{l} \text{Same} \end{array} \right\}$$

thus H_1, C are independent.

$$P(H_1 \cap H_2 \cap C) = P(HH) = \frac{1}{4} \quad \left. \begin{array}{l} \text{different} \end{array} \right\} \\ P(H_1)P(H_2)P(C) = \frac{1}{2} \times \frac{1}{2} \times \frac{1}{2} = \frac{1}{8} \\ P(C/H_1) = P(H_2/H_1) = P(H_2) = \frac{1}{2} = P(C)$$

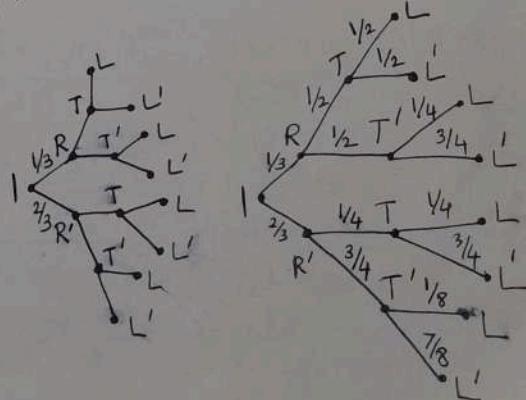
	$\frac{1}{4}$	$\frac{1}{4}$
	TT	TH
	HT	$\frac{1}{4}$

Thus H_1, H_2, C are pairwise independent but not independent.

8/9/21

- Q) It's raining $\frac{1}{3}$ of the days. Given that it is rainy, there will be heavy traffic with prob. $\frac{1}{4}$. Given it's not rainy, there will be heavy traffic with prob. $\frac{1}{4}$. If it's rainy and there is heavy traffic, I arrive late with prob. $\frac{1}{4}$. Prob. of late reduced to $\frac{1}{8}$ if it is not rainy, no heavy traffic. In other sit, prob. of being late = 0.25. You pick a random day.
- (a) What is prob that it's not rainy, there is heavy traffic, I am not late?
- (b) What is prob I am late?
- (c) Given I am late, what is prob that it rained that day?

Ans:



$$\textcircled{a} \quad P(R' \cap L') = P(R' \cap T \cap L')$$

$$= \frac{2}{3} \times \frac{1}{4} \times \frac{3}{4}$$

$$= \frac{1}{8}$$

$$\textcircled{b} \quad P(L) = ?$$

By total probability theorem,

$$P(L) = P(R \cap T \cap L) + P(R \cap T' \cap L)$$

$$+ P(R' \cap T \cap L) + P(R' \cap T' \cap L)$$

$$= \left(\frac{1}{3} \times \frac{1}{2} \times \frac{1}{2} \right) = \frac{1}{12} + \frac{1}{24} + \frac{1}{24} + \frac{1}{16}$$

$$+ \left(\frac{1}{3} \times \frac{1}{2} \times \frac{1}{4} \right)$$

$$+ \left(\frac{2}{3} \times \frac{1}{4} \times \frac{1}{4} \right) = \frac{4+2+2+3}{48} = \frac{11}{48}$$

$$+ \left(\frac{2}{3} \times \frac{3}{4} \times \frac{1}{8} \right)$$

$$\textcircled{c} \quad P(R|L) = \frac{P(R \cap L)}{P(L)} = \frac{P(R \cap L)}{P(L)}$$

$$= \frac{\left(\frac{1}{12} + \frac{1}{24} \right)}{\frac{11}{48}} = \frac{\left(\frac{4+2}{48} \right)}{\frac{11}{48}}$$

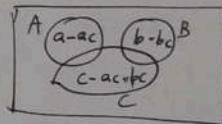
$$= \frac{6}{11} //$$

13/9/20

Q) ~~A, B~~ are disjoint events.

$$P(B \cup C) = \frac{3}{4} \quad P(A \cup C) = \frac{2}{3}$$

$$P(A \cup B \cup C) = \frac{11}{12} \quad P(C) = ?$$



$$P(A \cup C) = \frac{2}{3} = a + c - ac \quad \textcircled{1}$$

$$P(B \cup C) = \frac{3}{4} = b + c - bc \quad \textcircled{2}$$

$$P(A \cup B \cup C) = \frac{11}{12} = a + b + c - ac - bc - bc \quad \textcircled{3}$$

$$\textcircled{1} + \textcircled{2} \Rightarrow a + b + 2c - ac - bc = \frac{17}{12} \quad \textcircled{4}$$

$$\textcircled{4} - \textcircled{3} \Rightarrow c = \frac{17-11}{12} = \frac{6}{12} = \frac{1}{2}$$

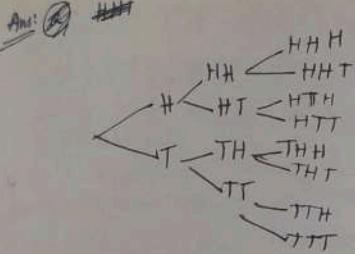
$$P(C) = \frac{1}{2} //$$

Q) Coin tossed 3 times.

(a) Probability of 3 heads = ?

(b) Prob that you get exactly one head

(c) Given a ~~head~~ at least one head, what is prob that you get at least 2 heads?



$$\textcircled{a} \quad P(HHH) = \frac{1}{8} = P(H)P(H)P(H)$$

$$\textcircled{b} \quad P(\text{exactly one head}) = P(HTT) + P(\cancel{HTH}+T) + P(TTH) \\ = \frac{1}{8} + \frac{1}{8} + \frac{1}{8} = \frac{3}{8}$$

$$\textcircled{c} \quad P(\geq 2H / \geq 1H) = \frac{P(\geq 2H \cap \geq 1H)}{P(\geq 1H)} \\ = \frac{P(HHT) + P(HTH) + P(THH) + P(HHH)}{1 - P(TTT)} \\ = \frac{\frac{1}{8} + \frac{1}{8} + \frac{1}{8} + \frac{1}{8}}{1 - \frac{1}{8}} = \frac{4}{7}$$

14/9/21

COVARIANCE, VARIANCE

COVARIANCE

→ Gives information about how X, Y variables are statistically related.

$$\begin{aligned} \text{cov}(X, Y) &= E[(X - E(X))(Y - E(Y))] \\ &= E[XY - XE(Y) - YE(X) + E(X)E(Y)] \\ &= E[XY] - E(X)E(Y) - E(X)E(Y) + E(X)E(Y) \\ &= E(XY) - E(X)E(Y) \end{aligned}$$

If X, Y are directly proportional,

$\text{cov}(X, Y)$ is +ve

If X, Y are indirectly proportional, $\text{cov}(X, Y)$ -ve.

PROPERTIES:

- 1) $\text{cov}(X, X) = \text{var}(X) = E(X^2) - E^2(X)$
 - 2) X, Y are independent, $\text{cov}(X, Y) = 0 \because E(XY) = E(X)E(Y)$
 - 3) $\text{cov}(X, Y) = \text{cov}(Y, X)$
 - 4) $\text{cov}(\alpha X, Y) = \alpha \text{cov}(X, Y) = \alpha E(XY) - \alpha(E(X)E(Y))$
 - 5) $\text{cov}(X+c, Y) = \text{cov}(X, Y)$
 - 6) $\text{cov}(X+Y, Z) = \text{cov}(X, Z) + \text{cov}(Y, Z)$
- $$\begin{aligned} &= \cancel{\text{cov}(X, E(X+Y)Z)} - E(X+Y)E(Z) \\ &= E(XZ + YZ) - E(X+Y)E(Z) \\ &= \cancel{E(XZ) + E(YZ)} - E(X)E(Y)E(Z) \\ &= \cancel{E(X)E(Z) + E(Y)E(Z)} - E(X)E(Y)E(Z) \\ &= E(X)E(Z) \cancel{[f + f(Y)]} \\ &= E(XZ) + E(YZ) - E(X)E(Z) - E(Y)E(Z) \\ &= \text{cov}(X, Z) + \text{cov}(Y, Z). \end{aligned}$$

7) covariance of a sum: $z = x+y$

$$\begin{aligned}\text{Var}(z) &= \text{cov}(z, z) = \text{cov}(x+y, x+y) \\ &= \text{cov}(x, x) + \text{cov}(x, y) + \text{cov}(y, x) + \text{cov}(y, y) \\ &= \text{var}(x) + 2\text{cov}(x, y) + \text{var}(y)\end{aligned}$$

8) $\text{Var}(ax+by) = a^2\text{var}(x) + 2ab\text{cov}(x, y) + b^2\text{var}(y)$

CORRELATION COEFFICIENT:

$$r(x, y) = \frac{\text{cov}(x, y)}{\sqrt{\text{var}(x)} \sqrt{\text{var}(y)}}$$

$$\# \left[v = \frac{x - E(x)}{\sigma_x}; u = \frac{y - E(y)}{\sigma_y} \right]$$

$$\text{cov}(x+a, y) = \text{cov}(x, y)$$

$$\Rightarrow r(x, y) = \text{cov}\left(\frac{x}{\sigma_x}, \frac{y}{\sigma_y}\right)$$

$$r(x, y) = \frac{\text{cov}(x, y)}{\sqrt{\text{var}(x)} \sqrt{\text{var}(y)}} \rightarrow \text{constants.}$$

PROPERTIES:

1) $-1 \leq r(x, y) \leq 1$

2) $r(x, y) = 1, y = ax + b, a > 0$

3) $r(x, y) = -1, y = ax + b, a < 0$

4) $r(ax+b, cy+d) = r(x, y) \text{ for } a, c > 0$

5) $r(x, y) = 0 \Rightarrow x, y \text{ not correlated.}$

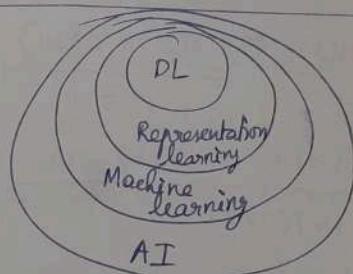
6) $r(x, y) \text{ +ve or -ve} \Rightarrow x, y \text{ +ve or -ve correlated.}$

7) If x, y are random var, x, y are independent, then x, y are not correlated.

If x, y are not correlated, then x, y may or may not be independent rand. var.

15/9/21

DEEP LEARNING



DL:

- To extract useful patterns from data.
- Neural N/w + optimisation
- Python + Tensorflow + friends
- Good questions + good data
- Data, hardware, community, tools, investment
-

DL \rightarrow ML with multilayer NN

RL \rightarrow ML extracting features of datasets

ML \rightarrow Machines improve with experience.

AI \rightarrow Computers mimic human intelligence.

Alternative Data:

\rightarrow Bank, Checking, Employment, Income, Insurance, Tenant, Utilities.

\rightarrow Cash flow underwriting.

\rightarrow App usage, Email receipt, Browsing history, Geolocation, Social Media Data.

Natural Language Processing:

\rightarrow Computer I/p, interpretation, O/p of Human Lang.

\rightarrow Audio, image, text, video including spoken, written, gestured.

\rightarrow Content generation, speech recognition, translation.

\rightarrow Chatbots, conversational interfaces, voice assistant.

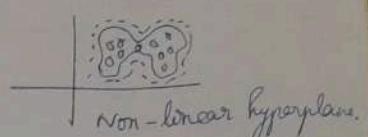
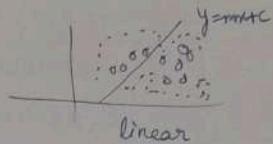
BIG DATA:

\rightarrow Describes storage, analysis of large / complicated data sets using techniques like AI.

ML:
 \rightarrow Method of designing sequence of actions to solve problem (algorithms) that optimise automatically through experience.

ACTIVATION FOR:

\rightarrow Introduce non-linearities into N/w.



ANN:

I/p, Hidden layer, O/p
bias, weights, activ. fn

11/9/21

VARIANCE

\rightarrow Defined as the measure of probability mass function spread.

\rightarrow Distance from mean (avg) $= E[X - \mu] = E[X] - \mu = 0$.

$$\text{Var}(X) = E[(X - \mu)^2] \geq 0.$$

$$E[g(X)] = \sum_x g(x) P[X(x)]$$

$$\text{Var}(X) = E[g(X)] = \sum_x (x - \mu)^2 P[X(x)]$$

$$S.D. = \sigma_X = \sqrt{\text{Var}(X)}$$

PROPERTIES OF VAR(X):

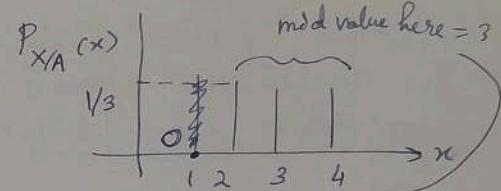
$$1) \text{Var}(\alpha X + b) = \alpha^2 \text{Var}(X)$$

$$2) \text{Var}(X + b) = E[(Y - E(Y))^2] = E[(X + b - (E(X) + b))^2]$$
$$(Y = X + b) = E[(X - E(X))^2] = E[(X - \mu)^2] = \text{Var}(X)$$

$$3) \text{ var}(ax) = E[(x-\mu)^2] = E[x^2 - 2\mu x + \mu^2]$$

$$\begin{aligned} &= E[x^2] - 2\mu E(x) + \mu^2 \\ &= E[x^2] - 2E(x)E(x) + [E(x)]^2 \\ &= E[x^2] - E[x]^2 \\ &= E[x^2] - (E[x])^2 \end{aligned}$$

Let $A = x \geq 2$



$$\begin{aligned} E[X/A] &= 3 \\ \text{var}[X/A] &= \frac{1}{3}[(4-3)^2] + \frac{1}{3}(3-3)^2 + \frac{1}{3}(2-3)^2 \\ &= \frac{1}{3} + 0 + \frac{1}{3} = \frac{2}{3} \end{aligned}$$

CONDITIONAL PROBABILITY MASS FUNCTION:

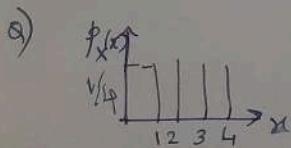
$$P_X(x) = P(X=x)$$

$$P(A) > 0.$$

$$P(x) = P(x = x/A)$$

$$\sum_x P_{x/A}(x) = 1 ; E[x/A] = \sum_x x P_{x/A}(x)$$

$$\sum_x P_{x/A}(x) = 1 ; \boxed{E[x/A] = \sum_x x P_{x/A}(x)}$$



$$E(x) = 2.5 = \frac{10}{4}$$

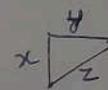
$$\text{Var}(x) = \frac{1}{12} (b-a)(b-a+2)$$

$$= \frac{1}{12} (4-1)(4-1+2) \quad \text{Euclidean distance:}$$

$$= \frac{1}{12} \times 5 \times 3 = \frac{5}{4}$$

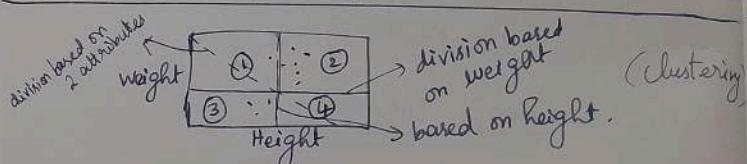
Eg:	Snake names	egg laying	scales	poison	Gold blooded	legs
	1	True				
	2	False				
	3		True	False	True	Four
	:					

$$Z = \sqrt{x^2 + y^2}$$



$$\text{Manhattan distance: } Z = x + y.$$

16/9/21 HEIGHT, WEIGHT



OBJECTIVE FUNCTION:

→ A threshold set to cluster the individuals.

FEATURES:

Eg:	Snake names	egg laying	scales	poison	Gold blooded	legs
	1	True				
	2	False				
	3		True	False	True	Four
	:					

Eg: alligator = Animal ('alligator', [1, 1, 0, 1, 1])

animals.append(alligator)

comparAnimals(animals, 3)

frog = (1, 0, 1, 1, 0, 1)

dist b/w frog, alli \Rightarrow 110101 ? features

$$\sqrt{0^2 + 1^2 + 1^2 + 1^2 + 0^2} = \sqrt{3} = 1.732 \text{ is dist}$$

difference squares b/w alligator, frog.

\rightarrow When dist ~~b/w~~ of features b/w 2 elements is less, they are closely related.

18/9/21 GRADIENT DESCENT OPTIMISATION

EPOCH \rightarrow Complete training of model.

$$J(w) = \frac{1}{2} \sum_i (\text{target}_i - \text{expected output}_i)^2$$

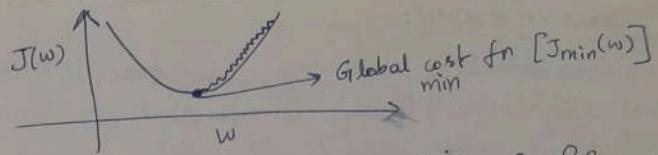
(Sum of squared error).

$$\Delta w_j = -\eta \frac{d J}{d w_j} = -\eta \frac{\partial J}{\partial w_j}$$

learning rate

$$w = w + \Delta w ; \Delta w = -\eta \sum_i (\text{target}_i - \text{output}_i) (-x_j^i)$$

$$\Delta w = \eta \sum_i (\text{target}_i - \text{output}_i) x_j^i$$



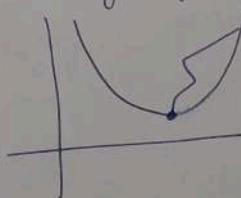
weight coefficient \rightarrow hiker; leg length of hiker \rightarrow learning rate
 $J(w) \rightarrow$ aim of climbing down.

DISADVANTAGE:

- \rightarrow Only after each epoch, weight updated.
- \rightarrow Large data set takes much time to reach global min (least cost).

STOCHASTIC GRADIENT DESCENT:

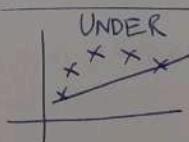
- \rightarrow Weights are updated only after completion of training dataset, in normal gradient descent.
- \rightarrow Here, weight of 1st data updated, then it goes on.



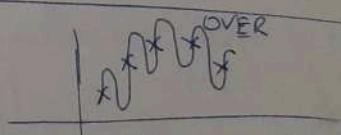
Path to global min need not be direct.
 η

21/9/21

UNDER, OVERFITTING



low complexity
high bias
high variance



high complexity
high bias
high variance

CROSS VALIDATION:

- Over, underfitting have very clear ~~data~~ difference.
- Dividing entire training set into parts for training, testing, so as to validate the model.
- This is maximum likelihood estimation.

$$y = mx + c \rightarrow \text{Line of linear regression.}$$

UNDERFITTING:

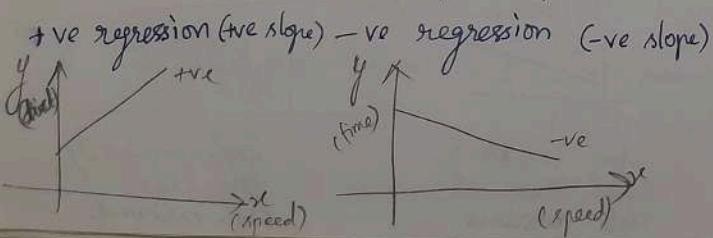
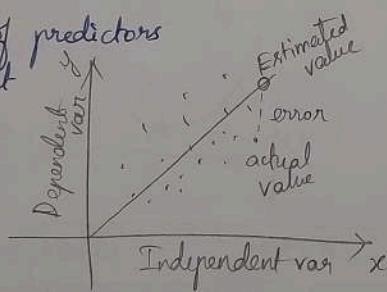
- If all the weights have the same value, then there is no learning happening in the neural network.

22/9/21 LINEAR REGRESSION

- For evaluating trends, sales estimates.
- Analysing impact of price changes.
- Risk in financial services.

3 main uses:

- Finding strength of predictions
- Forecasting an effect
- Trend forecasting.



BEST FIT:

- Best fitted line has least error.
- Aim is to have least possible error.

FINDING REGRESSION EQUATION:

x	1	2	3	4	5	Mean
y	3	4	2	4	5	3.6

$y = mx + c$ actual

$$\bar{x} = \text{mean of } x = \frac{\sum x}{5} = 3$$

$$\bar{y} = \text{mean of } y = \frac{\sum y}{5} = 3.6$$

$$\text{slope } m = \frac{\sum (x - \bar{x})(y - \bar{y})}{\sum (x - \bar{x})^2} = \frac{4}{10} = 0.4 = m$$

$$\bar{y} = m\bar{x} + c \Rightarrow 3.6 = 0.4(3) + c \\ \Rightarrow c = 2.4$$

$$y = 0.4x + 2.4$$

~~Actual~~ Regressed value:

x	y
1	2.8
2	3.2
3	3.6
4	4
5	4.4

} predicted

GOODNESS OF FIT:

(Least square estimation)
 $R^2 \rightarrow$ represents how close the data are being fixed.
 $R^2 = \frac{\sum (y_p - \bar{y})^2}{\sum (y - \bar{y})^2}$ (predicted) (mean) (actual)

$R^2 \propto$ goodness of fit.

$$\bar{x}^{(i+1)} = x^i + h \nabla f$$

$$\bar{x}^{(i+1)} = \begin{Bmatrix} 2 \\ 1 \end{Bmatrix} + h \begin{Bmatrix} 6 \\ 2 \end{Bmatrix} = \begin{Bmatrix} 2+6h \\ 1+2h \end{Bmatrix}$$

$$f(x, y) = x^2 + y^2 + 2x + 4$$

$$= (2+6h)^2 + (1+2h)^2 + 2(2+6h) + 4 = g(h)$$

gradient of h

Golden section search
Newton method.

23/9/21 GRADIENT DESCENT OPTIMISATION

Q Find minimum value of the function $f(x, y) = x^2 + y^2 + 2x + 4$

$$\bar{x}^0 = \begin{Bmatrix} x \\ y \end{Bmatrix} = \begin{Bmatrix} 2 \\ 1 \end{Bmatrix} \Rightarrow \text{Initial estimate of the optimal solution.}$$

$$g(h) = 40h^2 + 40h + 13$$

$$g_{\min} \Rightarrow \frac{dg}{dh} = 0 \Rightarrow 80h + 40 = 0$$

$$h = -\frac{1}{2} = 0.5 = l$$

Ans
Obj. function: $f(x, y) = y^2 + x^2 + 2x + 4$

To find $f_{\min}(x, y)$

Iteration 1: You calculate the gradient,

$$\frac{\partial f}{\partial x} = 2x + 2$$

$$\uparrow$$

$$x_{\min} = 2$$

$$\Rightarrow 2(2) + 2$$

$$= 6$$

$$\nabla f = 6i + 2j$$

$$\frac{\partial f}{\partial y} = 2y$$

$$\uparrow$$

$$y_{\min} = 1 \Rightarrow 2$$

$$\bar{x}^{(i+1)} = \begin{Bmatrix} 2+6(-0.5) \\ 1+2(-0.5) \end{Bmatrix} = \begin{Bmatrix} -1 \\ 0 \end{Bmatrix}$$

$$\bar{x}^{(0)} \Rightarrow \{2, 1\} \Rightarrow (4^2) + (1) + 2(2) + 4 = 13 = f$$

$$\bar{x}^{(i+1)} \Rightarrow \{-1, 0\} \Rightarrow 1 + 0 + (-2) + 4 = 3 = f$$

To find if objective function value is global minimum or not,

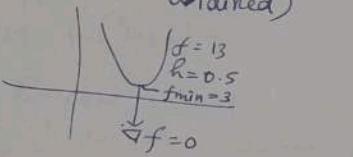
Iteration 2:

$$\text{Initial points: } (\bar{x}^{i+1}, \bar{y}^{i+1}) = (-1, 0)$$

$$\begin{aligned}\frac{\partial f}{\partial x} &= 2x+2 \\ &= 2(-1)+2 \\ &= 0\end{aligned}$$

$$\frac{\partial f}{\partial y} = 2y \Rightarrow (2(0)) = 0$$

$\nabla f = \vec{a}_i + \vec{a}_j$ (Gradient is 0, so global min obtained)
 f_{\min} at $(-1, 0) = 3$.

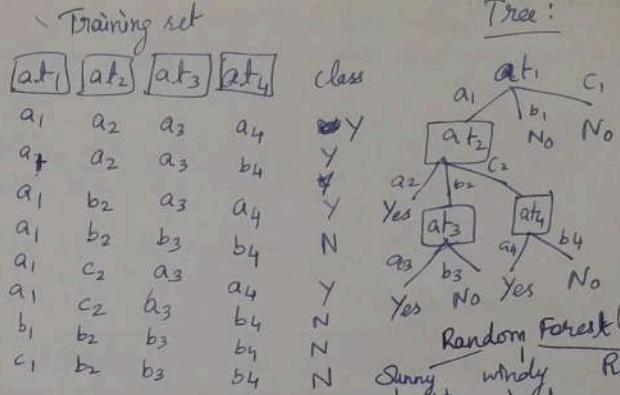


27/9/21

ML

Training set

a_1	a_2	a_3	a_4	class
a_1	a_2	a_3	b_4	Y
a_1	a_2	a_3	a_4	Y
a_1	b_2	a_3	a_4	Y
a_1	b_2	b_3	b_4	N
a_1	c_2	a_3	b_4	Y
a_1	c_2	b_3	b_4	N
b_1	b_2	b_3	b_4	N
c_1	b_2	b_3	b_4	N



Depending on the data, these decisions are made.

Vertical \rightarrow features.

Decision tree with least no. of branches is good, they are pruned (Eg: Here, $b_1, c_1 \Rightarrow$ No, ~~so~~ no need of other branches)

\rightarrow Learn from data

\rightarrow Establishes relationship b/w multiple features.

\rightarrow Extract statistical patterns

\rightarrow Reasoning under uncertain conditions.

SUPERVISED LEARNING ALGORITHMS:

\rightarrow Data set collection \rightarrow Brute force method (all relevant information is collected).

\rightarrow There might be noisy data, missing of important data

\rightarrow Algorithm selection based on prediction accuracy.

Cross validation present.

1) DECISION TREE SUPERVISED LEARNING ALGORITHM:

\rightarrow Classify instances by sorting them based on its features in training set.

2) KNN ALGORITHM (K nearest neighbours):

\rightarrow It is a lazy algorithm.

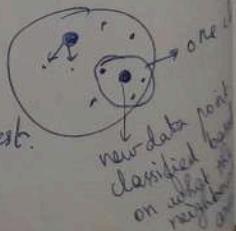
\rightarrow Takes a lot of time for computation.

\rightarrow Calculates distance b/w K neighbours, find which 2 are nearest to this point.

$$\text{Euclidean dist. } D(x, y) = \sqrt{\sum_{i=1}^n (x_i - y_i)^2}$$

It chooses which K point is nearest.

lot of computation

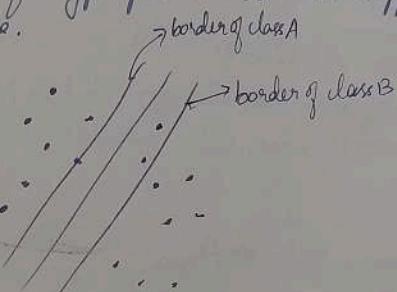


$$\text{Minkowsky dist } D(x, y) = \left(\sum_{i=1}^n (x_i - y_i)^r \right)^{1/r}$$

$$\text{Manhattan dist: } D(x, y) = \sum_{i=1}^n |x_i - y_i|$$

3) SUPPORT VECTOR MACHINE: (Used for both classification, regression).

- Uses Hyperplane (a line that separates two classes with max distance possible).
- Optimisation
- Maximises separation b/w classes, maximises distance of hyperplane to the support vector on either side.



- Quadratic polynomial hyperplane.
- Radial basis function.
- Gaussian plane kernel.

25/9/21 MAXIMUM LIKELIHOOD ESTIMATION

Random variable: x_1, x_2, \dots, x_n

Likelihood fn: $L(\theta) = P(x_1 = x_1, x_2 = x_2, \dots, x_n = x_n)$

Aim: To find the value of θ , for which max likelihood fn.

Example of MLE:

$x_1, x_2, \dots, x_n \Rightarrow$ Random samples.

$x_i = 0 \Rightarrow$ Any selected student doesn't own a car.

$x_i = 1 \Rightarrow$ Any selected student owns a car.

$x_i \rightarrow$ independent Bernoulli R.V with unknown parameter p .

(Estimator of p to be found)

Ans: Probability mass function of each x_i is

$$f(x_i, p) = p^{x_i} (1-p)^{1-x_i} \text{ such that-}$$

unknown parameter

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

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

$$1) L(p) = \prod_{i=1}^n f(x_i, p) = p^{x_1} (1-p)^{1-x_1} p^{x_2} (1-p)^{1-x_2} \dots p^{x_n} (1-p)^{1-x_n}$$

$$2) L(p) = p^{\sum x_i} (1-p)^{n-\sum x_i} \quad 0 < p < 1 \quad \boxed{\text{PMF}}$$

$$3) \ln(L(p)) = [\sum x_i \ln(p)] + [(n - \sum x_i) \ln(1-p)]$$

$$4) \text{diff} \Rightarrow 0, 1 \cdot \frac{d}{dp} \ln(L(p)) =$$

$$\frac{1}{L(p)} \frac{dL(p)}{dp} = \left[\frac{1}{p} \sum x_i \right] + \frac{(n - \sum x_i) (-1)}{1-p} = 0$$

$$\frac{\frac{p(1-p)}{L(p)}}{dp} = (1-p)\sum x_i - p(n-\sum x_i) = 0$$

$$\sum x_i - p\sum x_i - np + p\sum x_i = 0$$

$$\boxed{\sum x_i = np}$$

$\hat{p} = \frac{\sum x_i}{n}$

max: likelihood estimator

Likelihood Function:

x_1, x_2, \dots, x_n are random variables that depend on one or more unknown parameters $\theta_1, \theta_2, \dots, \theta_m$.

PMF $f(x; \theta_1, \theta_2, \dots, \theta_m)$

$\theta_1, \theta_2, \dots, \theta_m$ are restricted to a parameter space of Ω .

Joint PMF of x_1, \dots, x_n

$$L(\theta_1, \theta_2, \dots, \theta_m) = \prod_{i=1}^n f(x_i; \theta_1, \theta_2, \dots, \theta_m)$$

$(\theta_1, \theta_2, \dots, \theta_m)$ in Ω is called as likelihood fn.

Let $U_1(x_1, x_2, \dots, x_n), U_2(x_1, \dots, x_n), \dots, U_m(x_1, \dots, x_n)$ be a m -by tuple that maximises likelihood fn.

The estimator (max) $\theta_i^* = U_i(x_1, \dots, x_n)$ for $i=1 \dots m$

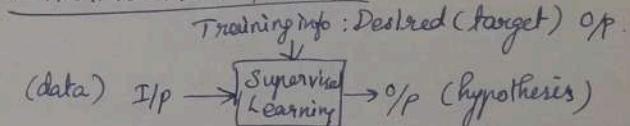
28/9/21

BAYESIAN LEARNING/INFERENCE

→ Bayes rule is transformed into naive Bayes classifier.

→ ML models are based on probability of observed data.

CLASSIFICATION PROBLEM:



$$\text{Error} = [\text{Target o/p} - \text{Actual o/p}]$$

Training data:

→ Of the form $(d, h(d))$

→ $d \rightarrow$ data objects to classify (I/p.)

→ $h(d) \rightarrow$ correct class info for d

GOAL:

→ Give d_{new} , get $h(d)_{\text{new}}$ or $h(d_{\text{new}})$

BAYESIAN FRAMEWORK:

→ Combining observed data, prior knowledge

→ Provides practical learning algorithm.

→ Generative (model based) approach, offers useful conceptual framework.

→ Any kind of objects (time series, trees etc.) can be classified based on probabilistic model specs. (Bayesian framework)

BAYES RULE:

$$P(R/d) = \frac{P(d/R) P(R)}{P(d)}$$

$d \rightarrow$ data

$R \rightarrow$ hypothesis

$P(R) \rightarrow$ Prior belief (prob. of hypothesis before seeing data) $P(+/\text{cancer}) = 0.98$

$P(d/R) \rightarrow$ likelihood (prob. of data if hypothesis is true) $P(-\text{cancer}) = 0.008$

$P(d) = \sum_h P(d/R) P(R) \rightarrow$ Data evidence

(Marginal probability of data).

Joint probability: $P(A, B) = P(AB) = P(A/B) P(B)$

$$P(X) = \sum_i P(X, Y) = \sum_i P(X/Y) P(Y) = 1$$

→ Conditional probabilities sum to 1 provided that their conditions are the same.

Q) Patient takes lab test, result is positive. Test returns correct positive result 98% of the cases, correct negative in 97% cases. 0.008 of entire population has this disease.

(i) What is the probability that this person has cancer?

(ii) What is the probability that he doesn't have cancer?

(iii) What is the diagnosis?

Ans

$$(i) P(\text{cancer}/+) = \frac{P(+/\text{cancer}) P(\text{cancer})}{P(+)} \\ = \frac{P(\text{cancer} \cap +)}{P(+)}$$

$$\begin{aligned} P(+) &= P(+/\text{cancer}) P(\text{cancer}) + P(+/-\text{cancer}) P(-\text{cancer}) \\ P(\text{cancer}/+) &= \frac{0.98 \times 0.008}{0} \end{aligned}$$

$$P(+/-\text{cancer}) = 0.03 \Rightarrow (1 - 0.97)$$

$$P(+) = P(+/\text{cancer}) (0.98 \times 0.008) + (0.03) (0.992) \\ = 0.0376$$

$$P(\text{cancer}/+) = \frac{0.98 \times 0.008}{0.0376} = 0.2085$$

$$(ii) P(-\text{cancer}/+) = \frac{P(-\text{cancer} \cap +)}{P(+)} \\ = \frac{P(+/-\text{cancer}) P(-\text{cancer})}{P(+)}$$

$$= \frac{0.03 \times 0.992}{0.0376} = 0.79148 = (1 - 0.2085)$$

(iii) The diagnosis is probably false since
(i) $\cancel{(ii)}$, the person doesn't have cancer.

CHOOSING HYPOTHESES:

Maximum likelihood hypothesis : $h_{ML} = \arg\max_{R \in H} P(d|R)$

Generally, most probable hypothesis given training data, which is maximum a posteriori hypothesis.

Useful observation: It doesn't depend on denominator $P(d)$

NAIVE BAYES CLASSIFIER:

Assumptions:

Attributes that describe data instances are conditionally independent given the classification hypothesis.

$$P(d|h) = P(a_1, \dots, a_n | h) = \prod_i P(a_i | h)$$

$a_i \rightarrow$ attributes

→ It is a simplifying assumption, it is violated in reality.

→ Bayes classifier that uses Naive Bayes assumption, computes MAP hypothesis is Naive Bayes classifier.

→ It is a practical learning method.

Applications:

- Medical diagnosis
- Text classification

$h_{MAP} = \arg\max_{R \in H} P(h|d)$ NAIVE BAYES SOLUTION:

→ classify any new datum instance $x = (a_1, \dots, a_n)$ as:

$$\begin{aligned} h_{\text{Naive Bayes}} &= \arg\max_h P(h) P(x|h) \\ &= \arg\max_h P(h) \prod_i P(a_i | h) \end{aligned}$$

→ To do this based on training examples, the parameters have to be estimated from training examples.

→ For each target value (hypothesis) h

$$\hat{P}(h) = \text{estimate } P(h)$$

$$\hat{P}(a_i | h) = \text{estimate } P(a_i | h)$$

→ For each value at each of the datum (data row)

Eg: $x = (\text{outlook} = \text{sunny}, \text{temp} = \text{cool}, \text{humidity} = \text{high})$

$$h_{NB} = \arg\max_h P(h) P(x|h) = \arg\max_h P(h) \prod_i P(a_i | h)$$

$$\arg\max_h P(h) P(\text{out} = \text{sunny}|h) P(\text{temp} = \text{cool}|h) P(\text{humidity} = \text{high}|h)$$

Here,
 $P(\text{play tennis} = \text{yes}) = \frac{9}{14} = 0.64$ \rightarrow no. of rows with yes
 $P(\text{play tennis} = \text{no}) = \frac{5}{14} = 0.36$
 $P(\text{wind} = \text{strong} / \text{Play tennis} = \text{yes}) = \frac{3}{9} = 0.33$
 $P(\text{wind} = \text{strong} / \text{play tennis} = \text{no}) = \frac{3}{5} = 0.6$

Transition		
H	C	
H	0.7	0.3
C	0.4	0.6

Emision			
I	2	3	
H	0.2	0.4	0.4
C	0.5	0.4	0.1

$T_1 = \begin{bmatrix} 0.8 & 0.2 \\ H & C \end{bmatrix} \text{ initial}$

$P(\text{yes}) P(\text{sunny}/\text{yes}) P(\text{cool}/\text{yes}) P(\text{high}/\text{yes}) = 0.0053$
 $P(\text{no}) P(\text{sunny}/\text{no}) P(\text{cool}/\text{no}) P(\text{high}/\text{no}) = 0.0206$

18/12/21 BACKWARD PROPAGATION

answer: ~~if~~ playTennis (x) = no (\because argmax w.r.t highest prob)

→ Bayes rule can be turned into a classifier.

→ Maximum A posteriori (MAP) hypothesis

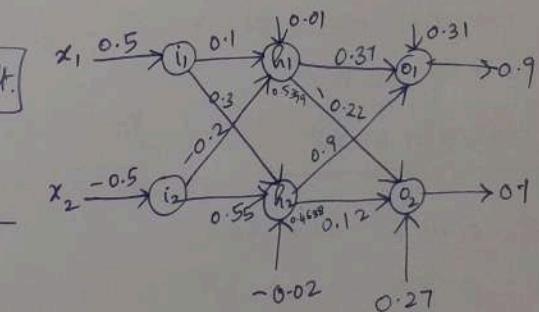
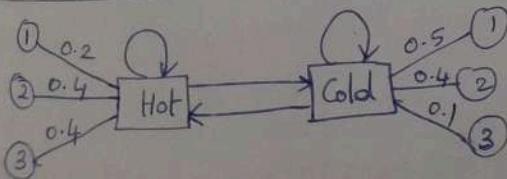
estimation uses prior knowledge, max likelihood doesn't.

→ Naive Bayes classifier assumes that attributes are independent given the class.

11/12/21

REVISION

a) Hidden Markov model:



$h_1 = (0.5)(0.1) + 0.01 + (-0.5 \times 0.2)$

$= 0.05 + 0.01 + 0.1$

$= 0.16$

$h_1' = \frac{1}{1+e^{-h_1}} = 0.5399$

$$\begin{aligned} h_2 &= -0.5 \times 0.55 + (-0.02) + (0.5 \times 0.3) \\ &= -0.275 - 0.02 + 0.15 \\ &= -0.145. \end{aligned}$$

$$h_2' = \frac{1}{1+e^{-h_2}} = \frac{1}{1+e^{0.145}} = 0.46381$$

$$\begin{aligned} o_1 &= (0.5399 \times 0.37) + 0.31 + (0.9 \times 0.46381) \\ &= 0.19976 + 0.31 + 0.4174 \\ &= 0.927163 \end{aligned}$$

$$[o_1' = 0.7164]$$

$$\begin{aligned} o_2 &= (-0.22 \times 0.5399) - (0.12 \times 0.46381) + 0.27 \\ &= -0.11877 - 0.05565 + 0.27 \\ &= 0.09557 \end{aligned}$$

$$[o_2' = 0.5238]$$

$$\begin{aligned} y_1 &= 0.9 \\ y_2 &= 0.1 \end{aligned}$$

$$e_1 = 0.9 - 0.7164 = 0.1836$$

$$e_2 = 0.1 - 0.5238 = -0.4238$$

Backward pass:

O/P to hidden layer

$$\begin{aligned} q'' &= e_1 * \delta(h_2') \delta(o_1') \\ &= 0.1836 \times 0.7164 o_1' (1-o_1') \\ &= 0.1836 \times 0.7164 (1-0.7164) = 0.0373 \end{aligned}$$

$$\begin{aligned} o_2'' &= e_2 \delta(o_2') \\ &= -0.4238 (o_2')(1-o_2') \\ &= -0.4238 (0.5238)(1-0.5238) \\ &= -0.1057 \end{aligned}$$

$$\boxed{\Delta w_{\text{new}} = \Delta w_{\text{old}} + \eta \frac{de}{dw_{ij}}}$$

$$\begin{aligned} h_1'' &= [(-0.0373 \times 0.37) + (-0.1057 \times -0.22)] \times \\ &\quad \downarrow o_1 \text{ to } h_1 \quad \cancel{o_1} \cancel{h_1} \quad \cancel{(0.5399)(1-0.5399)} \\ &= \cancel{s'(h_1')} = \cancel{0.7164} \\ &= 0.0092 \quad 0.0092 \quad h_1' (1-h_1') \end{aligned}$$

$$\begin{aligned} h_2'' &= [(-0.0373 \times 0.9) + (-0.1057 \times -0.12)] \times \\ &= [0.03357 + 0.012684] \times \cancel{0.46381} \\ &= 0.0115 \end{aligned}$$

$$\Delta w_{o_1} = 0.31 - 0.5(0.0373 \times 1) = 0.29135$$

$$\Delta w_2 = 0.37 - 0.5(0.0373 \times 0.5399) = 0.3599$$

$$\Delta w_3 = 0.9 - 0.5(-0.438 \times 0.0373) = 0.8918$$

$$\Delta w_4 = 0.27 - 0.5(-0.1057) = 0.32285$$

$$\Delta w_5 = -0.12 - 0.5(-0.1057 \times 0.4638) = 0.0954$$

$$\Delta w_6 = -0.22 - 0.05 \left(\frac{0.438}{0.5399} \times -0.1057 \right)$$

$$= -0.19146$$

$$\alpha = P(3, H) = P(3/H) P(H)$$

$$= 0.4 \times 0.8 = 0.32$$

(omission)

update all weights, then do forward pass. $b = P(3, C) = P(3/C) P(C)$

check o'_1, o'_2 , till it becomes equal to y_1, y_2 .

$$= 0.1 \times 0.2 = 0.02$$

~~$c = P(1, H) = P(1/H) P(H)$~~

~~$= 0.2 \times 0.7 = 0.14$~~

initially Hot, then cold

~~$d = P(1, C) = P(1/C) P(C/H)$~~

~~$= 0.5 \times 0.3 = 0.15$~~

~~$e = P(1, H) = P(1/H) P(H/C)$~~

~~$= 0.2 \times 0.4 = 0.08$~~

~~$f = P(1, C) = P(1/C) P(C)$~~

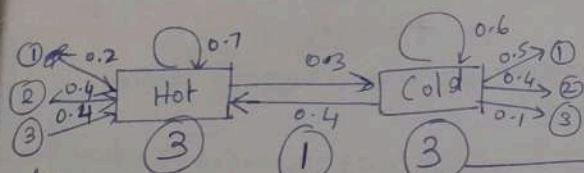
~~$= 0.5 \times 0.2 = 0.1$~~

$$c = P(1/H) P(H/H)$$

$$d = P(1/C) P(C/H)$$

Q) Hidden Marco Model.

$$T_1 = \begin{bmatrix} H & C \\ 0.8 & 0.2 \end{bmatrix}$$



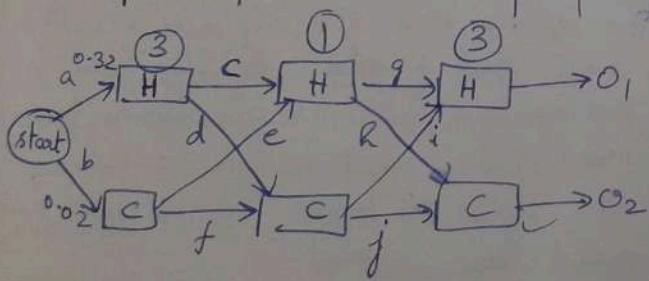
Ans

Transition matrix:

	H	C
H	0.7	0.3
C	0.4	0.6

Emission matrix:

	1	2	3
H	0.2	0.4	0.4
C	0.5	0.4	0.1



UNIT I PROB. SUMS

Q) A, B play a game. A, B toss coins. If both H, A gets 200, if both T B gets 100. One H, One T, no money given. A → biased coin with 60% prob, B fair coin. Find $E(A)$, $E(B)$, variances.

Ans:

$$P(H/A) = 0.6, P(T/A) = 0.4$$

$$P(H/B) = P(T/B) = 0.5$$

$$P(H/A)P(T/B) = 0.6 \times 0.5 = 0.3 \quad (\text{both head}) \quad A \text{ gets } 200$$

$$P(T/A)P(H/B) = 0.4 \times 0.5 = 0.2 \quad (\text{both tail}) \quad B \text{ gets } 100$$

$$\begin{aligned} E(A) &= P(H/A)P(T/B) + P(H/B)P(T/A) \\ &= 0.6(0.5) + 0.5 \times 0.2 \\ &= 0.3 + 0.2 \\ &= 0.5. \end{aligned}$$

$$\begin{aligned} E(A) &= AP(A) = (200 \times 0.3) = 60 + (0 \times 0.5) + (-100 \times 0.2) \\ &= 60 + 0 - 20 = 40. = E(A) = \mu. \end{aligned}$$

Q) $P(H/B)$

$$\begin{aligned} E(B) &= \text{Var}(A) = E[(A - E(A))^2] = E[A^2] - [E(A)]^2 \\ &= 14000 - \frac{40^2}{712400} = 200^2(0.3) + (100 \times 0.2) - 40^2 \end{aligned}$$

$$\checkmark \text{Var}(A), S.D = \sqrt{\text{Var}(A)}$$

$$\text{Var}(B) = E[(B - E(B))^2] = E(B^2) - [E(B)]^2$$

$$\begin{aligned} E(B) &= (-200 \times 0.3) + (0 \times 0.5) + (100 \times 0.2) \\ &= -60 + 20 = -40 \end{aligned}$$

$$\begin{aligned} E(B^2) &= 40000 \times 0.3 + 100^2 \times 0.2 = 12000 \\ E(B)^2 &= 40^2, E(B^2) - [E(B)]^2 = 12000 - 40^2 = 12400. \end{aligned}$$

Q) Decide b/w England, India to host next matches. No. of people attending matches, per And:

X	England	India	Probability
10000	30000	0.2	
20000	40000	0.5	
30000	60000	0.3	

Each ticket \Rightarrow £30 in England, R.8.600. Find Expected value, SD, compare.

$$\begin{aligned} E(X) &= \sum x \cdot p(x) = (10000)(0.2) + 20000(0.5) \\ &\quad + (30000)(0.3) \\ &= 2000 + 10000 + 9000 = 11000 + 10000 = 21000 \end{aligned}$$

England:

$$\begin{aligned} E(X^2) &= \sum x^2 \cdot p(x) = (10000)^2(0.2) + (20000)^2(0.5) + (30000)^2(0.3) \\ &= 490000000 \end{aligned}$$

$$\begin{aligned} \text{Var}(X) &= E(X^2) - [E(X)]^2 = 490000000 - (21000)^2 \\ &= 490000000 \end{aligned}$$

Similarly India

Q) Investor buys stock of 2 companies spending 10000 in each. Stock of each company goes up by 50% with a prob. of 0.6 or goes down by 40% with prob. of 0.4. Let random variable X present value after one month. 1) Find prob of, 2) Find Expected value. 3) Find S.D

ANS: $P(X) = \begin{cases} 0.4(10000 + 5000) + (0000 + 5000) = 30000, & P(X) = 0.6 \\ 0.4(10000 - 4000) + (10000 - 4000) = 6000 + 6000 = 12000 \\ 0.6(10000 - 4000) + (10000 + 5000) = 21000 \end{cases}$

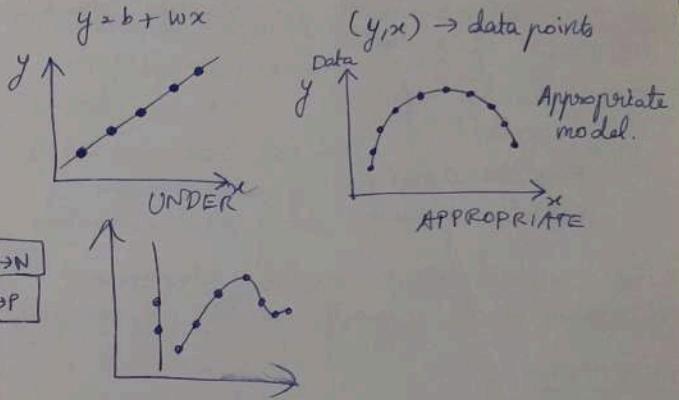
UNDER, OVER FITTING

Under Fitting:

- Not able to learn from the data very well.
- Low variance, high bias.

Over Fitting:

- High variance, low bias.



CONFUSION MATRIX:

		Positive	Negative
Predicted	+	TP	FP
	-	FN	TN

Actual

TP

TN

FP

FN

P → P, A → T

P → P, A → N

P → N, A → P

P → N, A → N

OVERFITTING: (More than the requirement)

→ Training error ↓, Testing error ↑.

→ System performs well during training, not during testing.

OVERCOMING IT:

- Reduce the number of features (feature selection)
- Regularisation (Keep all features but reduce magnitude of weights).

UNDERFITTING: (Less than requirement)

→ High training error, low testing error.

OVERCOMING IT:

→ (capacity has to be ~~increased~~ appropriate, it needs to solve the present task.)

OVERCOMING IT:

- Early stopping (Stop in those iterations where the accuracy is nearly 100%, and stop with appropriate ↑ in accuracy).

REGULARISATION: (L_1, L_2)

L_1 : LASSO: Minimise the absolute value of weight
(Some values of weights made 0.)

L_2 : RIDGE: Minimise the squared magnitude of weight

Ignore none. Sync all values.

CROSS VALIDATION:

→ Split the data into 5 equal parts.
 Allows to tune the hyperparameters, with the training set.
 → Same as testing.

HYPERPARAMETERS:

→ High level parameters that describe structural information of a model.

Eg: Learning rate (α)
 Weights
 Regularisation (λ)
 Activation Functions

→ These control ML model's behavior.

→ Maximum model capacity resulting in overfitting.

VALIDATION:

→ Testing the model with random data

→ Testing the model with datasets that are already not used while training.

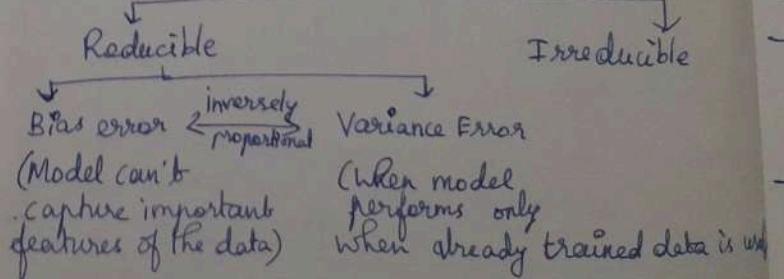
TESTING:

→ Part of the dataset that were already used while training the model, are used but in random manner.

MACHINE LEARNING PROCESS:

- Data collection (small amount of data)
- Data Preprocessing (To enhance the usable quality of data).
(Data augmentation)
- Split the data (Training, Cross validation)
- Train the model
- Evaluate the model (Unseen data)
- Deploying the model (Implement as an app or interface).

ERRORS IN MACHINE LEARNING:



SUPERVISED LEARNING ALGORITHMS

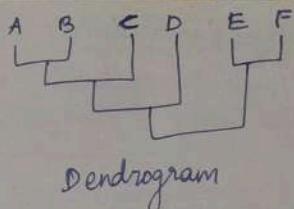
→ Labelled data are used.

1) CLASSIFICATION: (SVM, Decision Trees, Random Forest, K-nearest Neighbors)

→ To classify the data into different groups.
→ Just discriminating the features.

2) REGRESSION: (Linear reg, logistic, polynomial)

→ Based on past data, predicting the future data.
→ Determines the strength and character of relationship between one dependent variable and other independent variables.



DIMENSIONALITY REDUCTION:

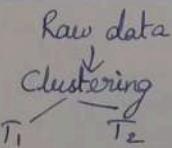
→ Principle component analysis (Based on Eigen values, Eigen vectors)

→ Similar value decomposition (Independent component analysis)

→ Auto encoders.

With compressed data, F/O/P → O/P
new representation of dataset, reducing feature set

UNSUPERVISED LEARNING ALGORITHM:



1) HIERARCHICAL CLUSTERING: (HCA)

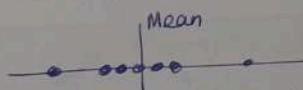
→ Agglomerative or divisive.

→ Bottleneck approach

→ Different clusters combined together based on similarities.

Ward's linkage:

→ Distance b/w 2 clusters defined by increased in sum of squared after the clusters are merged.



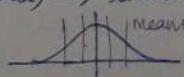
Eg: Here the predicted values are closer to mean, so normal distribution.

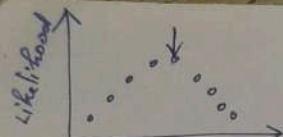
Prob. vs. likelihood:

$P(\text{data}/\text{distribution}) \rightarrow \text{probability}$

mean value

$P(\text{distribution}/\text{data}) \rightarrow \text{likelihood}$ (fitting a dist by finding mean for given data)





Bayesian statistics:
Frequentist statistics
E.g. Tossing a fair coin,

BAYESIAN STATISTICS:

$$\text{Difference} \rightarrow 0.5 \times (\text{No. of tosses}) - \text{No. of Heads}$$

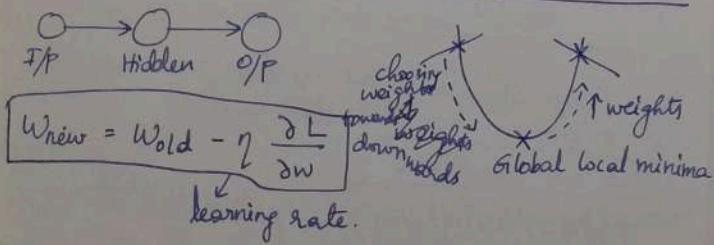
Actual prob. of getting Head in each toss. Actual no. of heads for all the tosses.

Dependence of the result \rightarrow No. of times exp. is repeated. (\because for each no. of trials, prob. changes).

p-Value:

- Sample size fixed
- Stopping intension - different.
- CI (Confidence)

STOCHASTIC GRADIENT DESCENT



- 1) Gradient descent \rightarrow When n data points considered (total population)
- 2) Stochastic gradient descent \rightarrow When only 1 point is considered.

③ Minibatch stochastic gradient descent \rightarrow When only a sample (K of no. of points) are taken for consideration, then it is Minibatch stochastic G.D.

(It takes time to jump to the centre).

LINEAR REGRESSION:

$$y = mx + c \quad (\text{One dep., ind var.})$$

$$\hat{y} = b_0 + b_1 x$$

$$b_1 = \frac{\sum (x - \bar{x})(y - \bar{y})}{\sum (x - \bar{x})^2}$$

$$b_0 = \bar{y} - b_1 \bar{x}$$

y, x values \rightarrow given as table

$$\text{Standard error} = \sqrt{\frac{\sum (\hat{y} - y)^2}{n-2}}$$

(Difference b/w actual, predicted) no. of samples

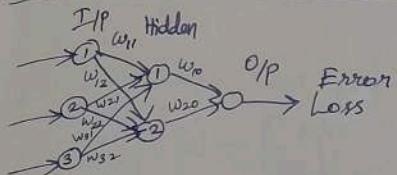
Multiple linear regression $\rightarrow y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots$

Polynomial linear regression $\rightarrow y = \beta_0 + \beta_1 x_1 + \beta_2 x_1^2 + \dots$

UNIT 2

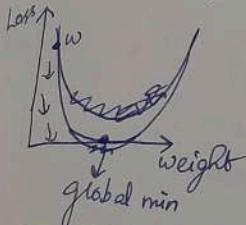
OPTIMIZERS :

1) GRADIENT DESCENT OPTIMIZER:



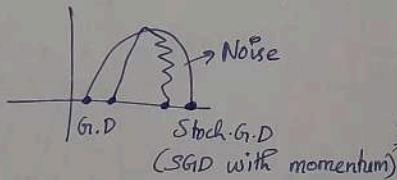
$$\hat{y} = \frac{1}{2} \sum_{i=1}^3 (y_i - \hat{y}_i)^2$$

Epoch → No. of times of visiting the algorithm.
Iteration → No. of times of processing it (no. of dataset passed once).



$$w_{\text{new}} = w_{\text{old}} - \eta \times \frac{\partial L}{\partial w_{\text{old}}} \rightarrow \text{G.D}$$

$$t \Rightarrow \text{batch size} \\ w_t = w_{t-1} - \eta \times v_{dw} \rightarrow \begin{array}{l} \text{Exponential} \\ \text{weighted avg.} \end{array}$$



$$\eta'_t = \frac{\eta}{\sqrt{\alpha_t + \epsilon}} \quad \text{initial learning rate } \alpha_t = \frac{1}{t} \sum_{i=1}^t \left(\frac{\partial L}{\partial w_{t-i}} \right)^2$$

3) RMS Prop (Root Mean Squared Propagation) (Ada delta)

$$\eta' = \frac{\eta}{\sqrt{s_{dw} + \epsilon}} \quad s_{dw} = \beta \times s_{dw} + (1-\beta) \left(\frac{\partial L}{\partial w_{t-1}} \right)^2$$

↓ weighted avg.

$$\begin{cases} w_t = w_{t-1} - \eta'_t \frac{\partial L}{\partial w_{\text{old}}} \\ b_t = b_{t-1} - \eta'_t \left(\frac{\partial L}{\partial b} \right) \end{cases} \rightarrow \begin{array}{l} v_{dw} \\ v_{db} \end{array}$$

Both weight, Learning rate change dynamically

4) ADAM OPTIMISER: (Better than others)

$$w_0 = w_{t-1} - [\eta] \times v_{dw} \rightarrow \text{Exponential weighted avg.}$$

$$(\text{Dynamic LR}) \quad \eta_t \leftarrow \sqrt{s_{dw} + \epsilon} \rightarrow \text{Small +ve no.}$$

$$w_t = w_{t-1} - \eta_t \times v_{dw}$$

→ For both output function and learning rate, if the weighted avg. is used, then the adam optimiser is the best.

2) ADAPTIVE GRADIENT DESCENT OPTIMISER: (ADA)

$\eta \rightarrow$ dynamic learning rate.



LINEAR DISCRIMINANT ANALYSIS (Dimensionality reduction)

- Dimensionality reduction technique, for supervised learning.
- Preprocessing step for pattern classification and ML applications.
- Finds axes that ~~maximize~~ separation between multiple classes.

GOAL: To project a feature space (n -dimensional data) onto smaller subspace K ($K \leq n-1$) while maintaining class discrimination.

Eg: 2D data set.

$$C_1 \Rightarrow X_1 = (x_1, y_1) = \{(4, 1), (2, 4), (2, 3), (3, 6), (4, 4)\}$$

$$C_2 \Rightarrow X_2 = (x_2, y_2) = \{(9, 10), (6, 8), (9, 5), (8, 7), (10, 8)\}$$

Step 1: Compute within class scatter matrix.

total covariance matrix $\leftarrow S_W = S_1 + S_2$

$S_1 \rightarrow$ Covariance matrix for C_1
 $S_2 \rightarrow$ Covariance matrix for C_2

Step 2: Compute covariance matrix of each class

$$S_1 = \left[\sum_{x \in C_1} (x - \mu_1)(x - \mu_1)^T \right] / n$$

\downarrow
mean of class C_1

$$\mu_1 = \left\{ \frac{4+2+2+3+4}{5}, \frac{1+4+3+6+4}{5} \right\} = [3.0, 3.6]$$

$$\mu_2 = \left\{ \frac{9+6+9+8+10}{5}, \frac{10+8+5+7+8}{5} \right\} = [8.4, 7.6]$$

$$(x - \mu_1) = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 4-3 & 2-3 & 2-3 & 3-3 & 4-3 \\ 1-3.6 & 4-3.6 & 3-3.6 & 6-3.6 & 4-3.6 \end{bmatrix}$$

$$= \begin{bmatrix} 1 & -1 & -1 & 0 & 1 \\ -2.6 & 0.4 & -0.6 & 2.4 & 0.4 \end{bmatrix}$$

$$\left[\sum (x - \mu_1)(x - \mu_1)^T \right] / 5 = S_1$$

$$\begin{bmatrix} 1 \\ -2.6 \end{bmatrix} \begin{bmatrix} 1 & -2.6 \\ -2.6 & 6.76 \end{bmatrix} = \begin{bmatrix} 1 & -2.6 \\ -2.6 & 6.76 \end{bmatrix} \quad (1)$$

$$\begin{bmatrix} -1 \\ 0.4 \end{bmatrix} \begin{bmatrix} -1 & 0.4 \\ 0.4 & 0.16 \end{bmatrix} = \begin{bmatrix} 1 & -0.4 \\ -0.4 & 0.16 \end{bmatrix} \quad (2)$$

$$\begin{bmatrix} -1 \\ -0.6 \end{bmatrix} \begin{bmatrix} -1 & -0.6 \\ -0.6 & 0.36 \end{bmatrix} = \begin{bmatrix} 1 & -0.6 \\ -0.6 & 0.36 \end{bmatrix} \quad (3)$$

$$\begin{bmatrix} 0 \\ 2.4 \end{bmatrix} \begin{bmatrix} 0 & 2.4 \\ 2.4 & 5.76 \end{bmatrix} = \begin{bmatrix} 0 & 0 \\ 0 & 5.76 \end{bmatrix} \quad (4)$$

$$\begin{bmatrix} 1 \\ 0.4 \end{bmatrix} \begin{bmatrix} 1 & 0.4 \\ 0.4 & 0.16 \end{bmatrix} = \begin{bmatrix} 1 & 0.4 \\ 0.4 & 0.16 \end{bmatrix} \quad (5)$$

$$\frac{(1) + (2) + (3) + (4) + (5)}{5} \Rightarrow S_1 = \begin{bmatrix} \frac{4}{5} & -\frac{2}{5} \\ -\frac{2}{5} & 2.6 \end{bmatrix}$$

$$S_1 = \begin{bmatrix} 0.8 & -0.4 \\ -0.4 & 2.6 \end{bmatrix}$$

$$S_2 = \begin{bmatrix} 1.84 & -0.04 \\ -0.04 & 2.64 \end{bmatrix}$$

within class scatter matrix

$$S_W = S_1 + S_2 = \begin{bmatrix} 2.64 & -0.44 \\ -0.44 & 5.28 \end{bmatrix}$$

Step 3: Between class scatter matrix:

$$S_B = (\mu_1 - \mu_2)(\mu_1 - \mu_2)^T$$

$$= \begin{bmatrix} -5.4 \\ -4 \end{bmatrix} \begin{bmatrix} -5.4 & -4 \end{bmatrix} = \begin{bmatrix} 29.16 & 21.6 \\ 21.6 & 16.00 \end{bmatrix}$$

Step 4: Find best LDA projection vector.

$$\mathbf{V}^T \mathbf{V} = S_W^{-1} S_B \mathbf{V}$$

projection vector

Find Eigen vector having largest Eigen value.

$$\begin{bmatrix} v_1 & v_2 \end{bmatrix} \begin{bmatrix} v_1 \\ v_2 \end{bmatrix} = S_W^{-1} (\mu_1 - \mu_2)$$

$$\begin{bmatrix} v_1 & v_2 \end{bmatrix} = \begin{bmatrix} 0.1921 & -0.032 \\ -0.032 & 0.38 \end{bmatrix} \begin{bmatrix} -5.4 \\ -4 \end{bmatrix} = \begin{bmatrix} -0.91 & -0.39 \end{bmatrix}$$

Step 5: Dimension reduction

$$y = \mathbf{V}^T x \rightarrow I/P \text{ data samples.}$$

↓
projection vector

$$= \begin{bmatrix} -0.91 \\ -0.39 \end{bmatrix}$$

PCA (PRINCIPLE COMPONENT ANALYSIS):

→ High dimension data is complex to process due to inconsistencies in features that increase computation time, make data processing complex.

→ This is a dimensionality reduction technique that identifies correlations, patterns in data set, so that it can be transformed into dataset of significantly lower dimension without loss of important info.

STEP 1: Scaling of data

$$z = \frac{\text{variable value} - \text{Mean}}{\text{S.D}}$$

STEP 2: Find covariance matrix (expresses correlation b/w different variables in dataset)

"Identify heavily dependent variables, as they contain biases", "Redundant info reduces overall performance of model".

NXN matrix

N dimensions

COVARIANCE VALUE:

$$\begin{bmatrix} \text{cov}(x, x) & \text{cov}(x, y) \\ \text{cov}(y, x) & \text{cov}(y, y) \end{bmatrix}$$

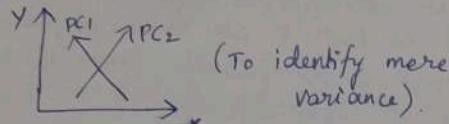
STEP 3: Find Eigen values, Eigen vectors.

PCA → new set of data derived from original data.

→ Compress, preserves most of the useful info that was scattered among initial variables. $\text{cov}(x,y) = \frac{1}{4-1} [(4-8)(11-8.5) + (8-8)(4-8.5) + (13-8)(5-8.5) + (7-8)(14-8.5)]$

→ Eigen vectors don't change direction (linear transformation).

→ Eigen values just denote scalars of the respective Eigen vectors.



Eg: Movie review dataset.

Q) Given following data, use PCA to reduce dimension from 2 to 1.

Feature	E-1	E-2	E-3	E-4
x	4	8	13	7
y	11	4	5	14

No. of features $n=2$

No. of samples $N=4$

Mean: $\bar{x} = 8$, $\bar{y} = 8.5$

Covariance matrix:

$$\text{cov}(x,y) = \left(\frac{1}{N-1} \right) \sum_{k=1}^N (x_{ik} - \bar{x}_i)(y_{ik} - \bar{y}_j)$$

$$= \left(\frac{1}{4-1} \right) [(4-8)^2 + (8-8)^2 + (13-8)^2 + (7-8)^2]$$

$$= \frac{1}{3} [16 + 0 + 25 + 1] = 14$$

$$= -11$$

$$\text{cov}(x,y) = \frac{1}{N-1} \sum_i (x_i - \bar{x})(y_i - \bar{y}) = \text{cov}(y,x)$$

$$\text{cov}(y,y) = \frac{1}{N-1} \sum_i (y_i - \bar{y})^2 = 23$$

$$S = \begin{bmatrix} \text{cov}(x,x) & \text{cov}(x,y) \\ \text{cov}(y,x) & \text{cov}(y,y) \end{bmatrix} = \begin{bmatrix} 14 & -11 \\ -11 & 23 \end{bmatrix}$$

Eigen value:

$$\det [S - \lambda I] = 0$$

$$\begin{vmatrix} 14-\lambda & -11 \\ -11 & 23-\lambda \end{vmatrix} = 0 \Rightarrow \lambda^2 - 37\lambda + 201 = 0$$

$$\lambda = 30.3849, 6.615$$

$$\text{Hence, } \lambda_1 = 30.3849 \rightarrow \text{PC}_1$$

Eigen vector of λ_1 :

$$(S - \lambda_1 I) v_1 = 0$$

$$\begin{bmatrix} 14-\lambda_1 & -11 \\ -11 & 23-\lambda_1 \end{bmatrix} \begin{bmatrix} v_1 \\ v_2 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$$(14-\lambda_1)v_1 - 11v_2 = 0$$

$$-11v_1 + (23-\lambda_1)v_2 = 0$$

$$\frac{U_1}{11} = \frac{U_2}{14-11} = t$$

When $t=1$,
 $U_1 = 11; U_2 = 14 - 11$

Eigen vector U_1 of γ = $\begin{bmatrix} 11 \\ 14-11 \end{bmatrix} = \begin{bmatrix} U_1 \\ U_2 \end{bmatrix}$

$$\begin{bmatrix} 11 \\ 14-30+3849 \end{bmatrix} = \begin{bmatrix} 11 \\ -16 \cdot 3849 \end{bmatrix}$$

Normalised Eigen Vector:

$$e_1 = \begin{bmatrix} 11 \\ -16 \cdot 3849 \end{bmatrix} \Rightarrow \begin{bmatrix} \overset{U_1}{\uparrow} & \overset{\sqrt{U_1^2 + U_2^2}}{\rightarrow} \\ \downarrow \sqrt{11^2 + 16 \cdot 3849^2} & \\ -16 \cdot 3849 & \downarrow \sqrt{11^2 + 16 \cdot 3849^2} \end{bmatrix}$$

$$e_1 = \begin{bmatrix} 0.5574 \\ -0.8303 \end{bmatrix}$$

$$\text{For } e_2 = \begin{bmatrix} 0.8303 \\ 0.5574 \end{bmatrix}$$

NEW dataset:

E_1	E_2	E_3	E_4
P_{11}	P_{12}	P_{13}	P_{14}

$$P_{11} = e_1^T \begin{bmatrix} x_i - \bar{x}_i \\ y_i - \bar{y}_i \end{bmatrix} = [0.5574 \quad -0.8303] \begin{bmatrix} 4-8 \\ 11-8.5 \end{bmatrix} = -4.3052$$

$$P_{12} = 3.7361$$

$$P_{13} = 5.6928$$

$$P_{14} = -5.1238$$

FISHER'S LINEAR DISCRIMINANT FUNCTION: (FLD)

→ First introduced to discriminate 2 plant species ~~iris~~ ^{iris} setosa, iris versicolor

d → dimensional data points $x_1, \dots, x_2, \dots, x_n$

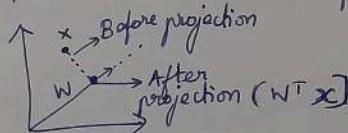
c_1, c_2
 \downarrow
 class 1 class 2

$N_1, N_2 \rightarrow$ No. of samples in each class.
 $w \rightarrow$ unit vector to project data points.

$$N = N_1 + N_2$$

If $x[n]$ are samples of feature space,
 then $[w^T x[n]] \rightarrow$ Data points after projection

$m_i \rightarrow$ mean of classes before projection
 $M_i \rightarrow$ " " " after " ". $M_i = w^T m_i$



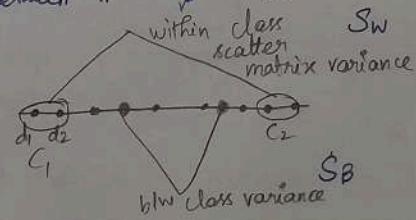
Scatter matrix ($m \times m$)

↳ +ve semi definite matrix

Sample variance \propto No. of samples.

→ Within class scatter matrix.

→ Between " "



$$\arg \max j(w) = \frac{(M_1 - M_2)^2}{S_1^2 + S_2^2} \rightarrow \begin{matrix} \text{Between class} \\ \text{scatter} \end{matrix} \quad \rightarrow \begin{matrix} \text{Within class scatter} \end{matrix}$$

Numerator:

$$\begin{aligned} (M_1 - M_2)^2 &= (w^T m_1 - w^T m_2)(w^T m_1 - w^T m_2)^T \\ &= w^T (m_1 - m_2)(m_1 - m_2)^T w \\ &= w^T S_b w \end{aligned}$$

$S_b \rightarrow b/w$ class scatter.

Denominator:

$$S_1^2 + S_2^2 \neq$$

Before projection $S_i = \sum_{x[n]} (x[n] - m_i)(m[n] - m_i)^T$

After projection $S_i^2 = \sum_{x[n]} (w^T x[n] - M_i)(w^T x[n] - M_i)^T$

$$M_i = w^T x[n]$$

$$\begin{aligned} S_i^2 &= w^T \left(\sum (x[n] - m_i)(x[n] - m_i)^T \right) w \\ &= w^T S_i w \end{aligned}$$

$$S_1^2 + S_2^2 = w^T (S_1 + S_2) w = w^T S_w w$$

$S_w \rightarrow$ within class scatter

$$J(w) = \frac{(M_1 - M_2)^2}{S_1^2 + S_2^2} \Rightarrow \frac{w^T S_b w}{w^T S_w w} \rightarrow \begin{matrix} \text{differentiate} \\ \text{equate to 0,} \\ \text{V} \rightarrow \text{Eigen value} \\ \text{W} \rightarrow \text{Eigen vector} \end{matrix}$$

$$S_b w = V \times S_w w$$

$$S_w \rightarrow \text{full rank matrix} \Rightarrow S_w^{-1} S_b w = V w$$

LDA for multiple classes:

$$S_K = \sum_{k=1}^C (x[n] - m_i)(x[n] - m_i)^T$$

$S_w \rightarrow$ within class scatter

$C \rightarrow$ No. of distinct classes

$$S_w = \sum_{k=1}^C S_k \Rightarrow m_i = \frac{1}{N_i} \sum x[n]$$

$$S_b = \sum_{i=1}^C N_i (m_i - m)(m_i - m)^T$$

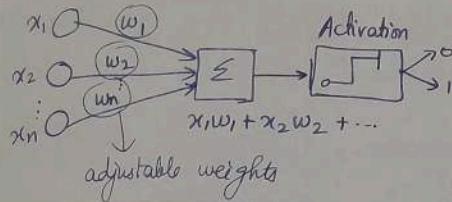
$$m = \frac{1}{N} \sum_{i=1}^n x_i \quad ; \quad w = \text{eig}(S_w^{-1} S_b)$$

$$S_b \leq C-1$$

LDA \rightarrow reducing dimension, computing time
 ↳ Multiple class classification

LDA	FLD
-----	-----

PERCEPTRON LEARNING ALGORITHM: (Supervised algorithm)



$$y = \sum x_i w_i + b$$

$w_{\text{new}} = w_{\text{old}} + \alpha \Delta w$ learning rate
error

(actual - predicted) \rightarrow error

$$b_{\text{new}} = b_{\text{old}} + \alpha t$$

target

updating weights and bias:

$$w_{\text{new}} = w_{\text{old}} + \alpha t x_i$$

Q) Implementing AND (bipolar I/P, targets) using perceptron.

x_1	x_2	y
-1	-1	-1
-1	1	-1
1	-1	-1
1	1	1

Ans Initialisation $\Rightarrow w_1 = w_2 = 0, b = 0$
 $\alpha = 1$ (default)

$$y = x_1 w_1 + x_2 w_2 + b$$

$$(i) x_1 = x_2 = 1; y_{\text{in}} = 0 + 0 + 0 = 0$$

$$\text{sigmoid } f(y_{\text{in}}) = \begin{cases} 0 & y_{\text{in}} \leq 0 \\ 1 & y_{\text{in}} > 0 \end{cases}$$

$$\text{Sigmoid} = \frac{1}{1+e^{-x}}$$

\rightarrow outlier drawback

Generative approach → Eg: To differentiate or find a language from all sets, learn everything then discriminate. → Focus on how data was generated	Discriminative approach → To find the language just know the difference, and how it sounds, without learning the lang. → Draw boundaries, then separate one class from another.
--	---

$P(Y/X) = P(Y) \times P(X/Y) \rightarrow$ Likelihood

$P(X) \times P(Y/X)$ generative model

$P(Y) \times P(X/Y)$ discriminative model

$P(Y/X) \rightarrow$ posterior

$P(X/Y) \rightarrow$ prior

$P(X) \rightarrow$ evidence

$P(Y) \rightarrow$ prior over labels

- Discriminative models separate classes instead of modeling conditional prob.
- separate one class from another.

Outliers → Something irrelevant in the I/P set

* Discriminative work better with outliers

* Generative ~~not~~ have disadvantage with outliers.

Maximum Likelihood Estimation (MLE)

- Probability distribution function
- Likelihood function.
- MLE
- Parameter estimation of Bernoulli density function
 $f(x|\theta)$
 ↓
 data $x \rightarrow$ Data (Source)
 $y \rightarrow$ Label (Target)

→ To estimate parameters of prob. distribution, to test whether given dataset follows some particular distribution.

→ Method to estimate parameters of prob. distn, given the observations.

→ MLE attempts to find parameter values that maximise likelihood function.

→ Resulting estimate is maximum likelihood estimate.

→ Random sample $X = \{x_1, x_2, \dots, x_n\}$ taken from prob. distribution having prob. density fn $p(x|\theta)$

x → Value of random variable

θ → Set of parameters that appear in function

→ Likelihood of sample X is a fn. of parameters, defined as the chance that the parameter θ would generate the observed data.

$$l(\theta) = P(X|\theta) = P(x_1|\theta)P(x_2|\theta)\dots P(x_n|\theta)$$

Generative	Discriminative
→ More biased (have stronger assumptions).	→ Less biased, and have lesser assumptions.
→ Can tolerate with missing data, and still produce good results.	→ Can't produce efficient results with missing data.
→ Generative less accurate	→ Assign values to data
→ For fewer data sets.	→ More accurate
→ Assign values to data, detect outliers	→ Classifying data with some conditions, so can be used for classifying, not for regression.
→ Learns $p(x)$	→ Learns $p(y x)$

PROBABILISTIC GENERATIVE MODELS:

(CONDITIONAL GENERATIVE MODEL) → Assign labels (unsupervised learning while rejecting outliers)

Eg: Rain prediction (~~if~~ less data)
 Stock market - prediction (~~if~~ for regression).

$$\log(L(\theta)) = \log [P(x_1/\theta) P(x_2/\theta) \dots P(x_n/\theta)]$$

$$L(\theta) = \log(P(x_1/\theta)) + \log P(x_2/\theta) + \dots + \log(P(x_n/\theta))$$

BERNOULLI DENSITY FUNCTION:

$$\begin{array}{ll} \rightarrow \text{Success} & x=1 \quad (p) \\ \rightarrow \text{Failure} & x=0 \quad (1-p) \\ f(x/p) = p^x (1-p)^{1-x} & x=1, 0. \end{array}$$

Estimation of parameter p of Bernoulli density function using MLE

$$f(x/p) = p^x (1-p)^{1-x}$$

$$L(p) = \log(f(x_1/p)) + \log(f(x_2/p)) + \dots + \log(f(x_n/p))$$

$$\begin{aligned} p &= \log[p^{x_1} (1-p)^{1-x_1}] + \log[p^{x_2} (1-p)^{1-x_2}] \\ &\quad + \dots + \log[p^{x_n} (1-p)^{1-x_n}] \end{aligned}$$

$$= x_1 \log p + (1-x_1) \log(1-p) + \dots + x_n \log p + (1-x_n) \log(1-p)$$

$$= (x_1 + x_2 + \dots + x_n) \log p + [(1-x_1) + (1-x_2) + \dots + (1-x_n)] \log(1-p)$$

$$\text{Let } x_1 + x_2 + \dots + x_n = Y$$

$$L(p) = Y \log p + (n-Y) \log(1-p)$$

To find ~~$\frac{d}{dp}$~~ p for max value of $L(p)$, diff $= \log e^{-\lambda} + \log \lambda^{x_1} - \log(x_1!) + \dots + \log e^{-\lambda} + \log \lambda^{x_n} - \log(x_n!)$

$$\frac{dL}{dp} = 0, \quad \frac{d(Y \log p + (n-Y) \log(1-p))}{dp} = 0$$

$$Y \left(\frac{1}{p}\right) + (n-Y) \left(\frac{1}{1-p}\right) (-1) = 0$$

$$\boxed{\frac{Y}{p} = \frac{n-Y}{1-p}} \quad \boxed{\frac{1-p}{p} = \frac{n-y}{y}}$$

$$\frac{1}{p} \times 1 = \frac{n}{Y} \times 1$$

$$\frac{1}{p} = \frac{n}{Y} \Rightarrow \boxed{p = \frac{Y}{n}}$$

~~$\frac{1}{p} = \frac{n}{Y}$~~

$$p = \frac{1}{n} (x_1 + x_2 + \dots + x_n)$$

$f(x/\theta) \rightarrow$ Prob. dist. f_n .

~~$f(x/\theta)$~~

$$\theta = \frac{1}{n} (x_1 + x_2 + \dots + x_n)$$

PARAMETER ESTIMATION USING MLE METHOD:

$$\text{Poisson dist: } f(x, \lambda) = \frac{e^{-\lambda} \lambda^x}{x!}$$

$$\text{Geometric dist: } f(x/p) = p(1-p)^{x-1}$$

~~$f(x, \lambda) = \frac{e^{-\lambda} \lambda^x}{x!}$~~

$$\text{POISSON: } f(x, \lambda) = \frac{e^{-\lambda} \lambda^x}{x!}$$

$$L(\lambda) = \log f(x_1/\lambda) + \log f(x_2/\lambda) + \dots + \log f(x_n/\lambda)$$

$$= \log \left[\frac{e^{-\lambda} \lambda^{x_1}}{x_1!} \right] + \log \left[\frac{e^{-\lambda} \lambda^{x_2}}{x_2!} \right] + \dots + \log \left[\frac{e^{-\lambda} \lambda^{x_n}}{x_n!} \right]$$

$$= -n\lambda + \sum_{i=1}^n x_i \log \lambda - \sum_{i=1}^n \log x_i!$$

$$\frac{dL(\lambda)}{d\lambda} = 0 \Rightarrow -n + \sum_{i=1}^n x_i \left(\frac{1}{\lambda}\right) - 0 = 0$$

$$\sum_{i=1}^n x_i = n ; \quad \boxed{\lambda = \frac{\sum_{i=1}^n x_i}{n}}$$

~~GEOMETRIC~~ GEOMETRIC DISTRIBUTION:

$$f(x, n) = p(1-p)^{x-1}$$

$$L(p) = \log f(x_1, p) + \log f(x_2, p) + \dots + \log f(x_n, p)$$

$$= \log [p(1-p)^{x_1-1}] + \dots + \log [p(1-p)^{x_n-1}]$$

$$L(p) = n \log p + \log(1-p) \left[\sum_{i=1}^n x_i - n \right]$$

$$\frac{dL(p)}{dp} = 0 \Rightarrow n \left(\frac{1}{p} \right) + \left[\sum_{i=1}^n x_i - n \right] \left(\frac{1}{1-p} \right) (-1) = 0$$

$$\frac{n}{p} = \frac{\sum_{i=1}^n x_i - n}{1-p} ; \frac{1-p}{p} = \frac{\sum_{i=1}^n x_i - n}{n}$$

$$\frac{1}{p} - 1 = \frac{\sum_{i=1}^n x_i}{n} - 1 \Rightarrow p = \frac{n}{\sum_{i=1}^n x_i}$$

To handle the outliers, and in cases where the values are unbounded (predicted values), logistic regression used.

Why LOGISTIC?

Mathematics Intuition:

Geometric Intuition:

(i) Mathematics Intuition

(ii) Geometric Intuition

(iii) Binary classification

(iv) Logistic Function

(v) Decision Boundary

(vi) Loss Function

(vii) Gradient Descent

(viii) Regularization

(ix) Overfitting

(x) Underfitting

(xi) Model Selection

(xii) Model Comparison

(xiii) Model Evaluation

(xiv) Model Selection

(xv) Model Comparison

(xvi) Model Evaluation

(xvii) Model Selection

(xviii) Model Comparison

(xix) Model Evaluation

(xx) Model Selection

(xxi) Model Comparison

(xxii) Model Evaluation

(xxiii) Model Selection

(xxiv) Model Comparison

(xxv) Model Evaluation

(xxvi) Model Selection

(xxvii) Model Comparison

(xxviii) Model Evaluation

(xxix) Model Selection

(xxx) Model Comparison

(xxxi) Model Evaluation

(xxxii) Model Selection

(xxxiii) Model Comparison

(xxxiv) Model Evaluation

(xxxv) Model Selection

(xxxvi) Model Comparison

(xxxvii) Model Evaluation

(xxxviii) Model Selection

(xxxix) Model Comparison

(xxxx) Model Evaluation

(xxxxi) Model Selection

(xxxxii) Model Comparison

(xxxxiii) Model Evaluation

(xxxxiv) Model Selection

(xxxxv) Model Comparison

(xxxxvi) Model Evaluation

(xxxxvii) Model Selection

(xxxxviii) Model Comparison

(xxxxix) Model Evaluation

(xxxxx) Model Selection

(xxxxxi) Model Comparison

(xxxxxii) Model Evaluation

(xxxxxiii) Model Selection

(xxxxxiv) Model Comparison

(xxxxxv) Model Evaluation

(xxxxxvi) Model Selection

(xxxxxvii) Model Comparison

(xxxxxviii) Model Evaluation

(xxxxxix) Model Selection

(xxxxxx) Model Comparison

(xxxxxi) Model Evaluation

(xxxxxii) Model Selection

(xxxxxiii) Model Comparison

(xxxxxiv) Model Evaluation

(xxxxxv) Model Selection

(xxxxxvi) Model Comparison

(xxxxxvii) Model Evaluation

(xxxxxviii) Model Selection

(xxxxxix) Model Comparison

(xxxxxx) Model Evaluation

(xxxxxi) Model Selection

(xxxxxii) Model Comparison

(xxxxxiii) Model Evaluation

(xxxxxiv) Model Selection

(xxxxxv) Model Comparison

(xxxxxvi) Model Evaluation

(xxxxxvii) Model Selection

(xxxxxviii) Model Comparison

(xxxxxix) Model Evaluation

(xxxxxx) Model Selection

(xxxxxi) Model Comparison

(xxxxxii) Model Evaluation

(xxxxxiii) Model Selection

(xxxxxiv) Model Comparison

(xxxxxv) Model Evaluation

(xxxxxvi) Model Selection

(xxxxxvii) Model Comparison

(xxxxxviii) Model Evaluation

(xxxxxix) Model Selection

(xxxxxx) Model Comparison

(xxxxxi) Model Evaluation

(xxxxxii) Model Selection

(xxxxxiii) Model Comparison

(xxxxxiv) Model Evaluation

(xxxxxv) Model Selection

(xxxxxvi) Model Comparison

(xxxxxvii) Model Evaluation

(xxxxxviii) Model Selection

(xxxxxix) Model Comparison

(xxxxxx) Model Evaluation

(xxxxxi) Model Selection

(xxxxxii) Model Comparison

(xxxxxiii) Model Evaluation

(xxxxxiv) Model Selection

(xxxxxv) Model Comparison

(xxxxxvi) Model Evaluation

(xxxxxvii) Model Selection

(xxxxxviii) Model Comparison

(xxxxxix) Model Evaluation

(xxxxxx) Model Selection

(xxxxxi) Model Comparison

(xxxxxii) Model Evaluation

(xxxxxiii) Model Selection

(xxxxxiv) Model Comparison

(xxxxxv) Model Evaluation

(xxxxxvi) Model Selection

(xxxxxvii) Model Comparison

(xxxxxviii) Model Evaluation

(xxxxxix) Model Selection

(xxxxxx) Model Comparison

(xxxxxi) Model Evaluation

(xxxxxii) Model Selection

(xxxxxiii) Model Comparison

(xxxxxiv) Model Evaluation

(xxxxxv) Model Selection

(xxxxxvi) Model Comparison

(xxxxxvii) Model Evaluation

(xxxxxviii) Model Selection

(xxxxxix) Model Comparison

(xxxxxx) Model Evaluation

(xxxxxi) Model Selection

(xxxxxii) Model Comparison

(xxxxxiii) Model Evaluation

(xxxxxiv) Model Selection

(xxxxxv) Model Comparison

(xxxxxvi) Model Evaluation

(xxxxxvii) Model Selection

(xxxxxviii) Model Comparison

(xxxxxix) Model Evaluation

(xxxxxx) Model Selection

(xxxxxi) Model Comparison

(xxxxxii) Model Evaluation

(xxxxxiii) Model Selection

(xxxxxiv) Model Comparison

(xxxxxv) Model Evaluation

(xxxxxvi) Model Selection

(xxxxxvii) Model Comparison

(xxxxxviii) Model Evaluation

(xxxxxix) Model Selection

(xxxxxx) Model Comparison

(xxxxxi) Model Evaluation

(xxxxxii) Model Selection

(xxxxxiii) Model Comparison

(xxxxxiv) Model Evaluation

(xxxxxv) Model Selection

(xxxxxvi) Model Comparison

(xxxxxvii) Model Evaluation

(xxxxxviii) Model Selection

(xxxxxix) Model Comparison

(xxxxxx) Model Evaluation

(xxxxxi) Model Selection

(xxxxxii) Model Comparison

(xxxxxiii) Model Evaluation

(xxxxxiv) Model Selection

(xxxxxv) Model Comparison

(xxxxxvi) Model Evaluation

(xxxxxvii) Model Selection

(xxxxxviii) Model Comparison

(xxxxxix) Model Evaluation

(xxxxxx) Model Selection

(xxxxxi) Model Comparison

(xxxxxii) Model Evaluation

(xxxxxiii) Model Selection

(xxxxxiv) Model Comparison

(xxxxxv) Model Evaluation

(xxxxxvi) Model Selection

(xxxxxvii) Model Comparison

(xxxxxviii) Model Evaluation

(xxxxxix) Model Selection

(xxxxxx) Model Comparison

(xxxxxi) Model Evaluation

(xxxxxii) Model Selection

(xxxxxiii) Model Comparison

(xxxxxiv) Model Evaluation

(xxxxxv) Model Selection

(xxxxxvi) Model Comparison

(xxxxxvii) Model Evaluation

(xxxxxviii) Model Selection

(xxxxxix) Model Comparison

(xxxxxx) Model Evaluation

(xxxxxi) Model Selection

(xxxxxii) Model Comparison

(xxxxxiii) Model Evaluation

(xxxxxiv) Model Selection

(xxxxxv) Model Comparison

(xxxxxvi) Model Evaluation

(xxxxxvii) Model Selection

(xxxxxviii) Model Comparison

(xxxxxix) Model Evaluation

(xxxxxx) Model Selection

(xxxxxi) Model Comparison

(xxxxxii) Model Evaluation

(xxxxxiii) Model Selection

(xxxxxiv) Model Comparison

(xxxxxv) Model Evaluation

(xxxxxvi) Model Selection

(xxxxxvii) Model Comparison

(xxxxxviii) Model Evaluation

(xxxxxix) Model Selection

(xxxxxx) Model Comparison

(xxxxxi) Model Evaluation

(xxxxxii) Model Selection

(xxxxxiii) Model Comparison

(xxxxxiv) Model Evaluation

(xxxxxv) Model Selection

(xxxxxvi) Model Comparison

(xxxxxvii) Model Evaluation

(xxxxxviii) Model Selection

(xxxxxix) Model Comparison

(xxxxxx) Model Evaluation

(xxxxxi) Model Selection

(xxxxxii) Model Comparison

(xxxxxiii) Model Evaluation

(xxxxxiv) Model Selection

(xxxxxv) Model Comparison

(xxxxxvi) Model Evaluation

(xxxxxvii) Model Selection

(xxxxxviii) Model Comparison

(xxxxxix) Model Evaluation

(xxxxxx) Model Selection

(xxxxxi) Model Comparison

(xxxxxii) Model Evaluation

(xxxxxiii) Model Selection

(xxxxxiv) Model Comparison

(xxxxxv) Model Evaluation

(xxxxxvi) Model Selection

(xxxxxvii) Model Comparison

(xxxxxviii) Model Evaluation

(xxxxxix) Model Selection

(xxxxxx) Model Comparison

(xxxxxi) Model Evaluation

(xxxxxii) Model Selection

(xxxxxiii) Model Comparison

(xxxxxiv) Model Evaluation

(xxxxxv) Model Selection

(xxxxxvi) Model Comparison

(xxxxxvii) Model Evaluation

(xxxxxviii) Model Selection

(xxxxxix) Model Comparison

(xxxxxx) Model Evaluation

(xxxxxi) Model Selection

(xxxxxii) Model Comparison

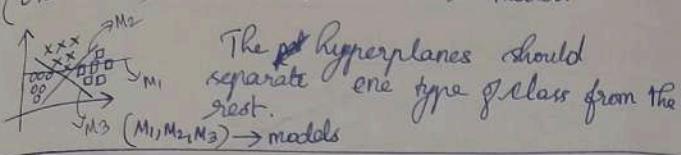
(xxxxxiii) Model Evaluation

(xxxxxiv) Model Selection

(xxxxxv) Model Comparison

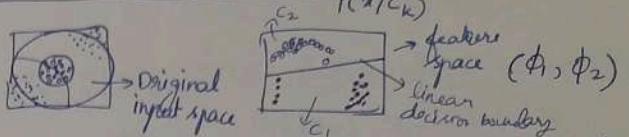
LOGISTIC REGRESSION FOR MULTICLASS CLASSIFICATION:

(One vs. rest) or (one vs. all) model.



The ~~per~~ hyperplanes should separate one type of class from the rest.

FIXED BASIS FUNCTIONS:



- Transform inputs using a vector of basis functions
- Resulting decision boundaries are linear in feature space.

→ To convert non-linearly separable function to ~~for~~ linearly separable function.

3 APPROACHES FOR CLASSIFICATION

DISCRIMINANT FUNCTIONS:

Prob. Generative Model:

→ Fit class conditional densities, class priors separately.

$$\frac{P(Y|x)}{\text{posterior probability}} = \frac{\text{Likelihood} \times \text{Prior}}{\text{Evidence}} = \frac{P(x|y) P(y)}{P(x)}$$

→ Apply Bayes' theorem to find posterior class probabilities.

→ Posterior prob. of a class written as:

- Logistic sigmoid acting on a linear function of x (2 classes). (Binary)

- Softmax transformation of a linear function of x (Multi-class). (Multi-class)

→ Parameters of densities and class priors found using Maximum Likelihood.

Prob. Discriminative Model:

→ Use functional form of generalized linear model explicitly.

→ Find the parameters directly using Maximum likelihood.

Iterative reweighted least squares:

→ For linear regression models, ML solution on the assumption of Gaussian noise leads to closed-form solution, as a result of quadratic dependence of log likelihood on parameter w .

→ Not a closed form solution, but error function is concave, has unique minimum, efficient iterative technique can be used, Newton-Raphson update to minimise a function $E(w)$

$$w_{\text{new}} = w_{\text{old}} - H^{-1} \nabla E(w)$$

Hessian matrix, second derivatives of $E(w)$.

$$\text{Sum of squares Error function: } E(w) = \frac{1}{2} \sum_{n=1}^N \{ t_n - w^T \phi(x_n) \}^2$$

$$\nabla E(w) = \phi^T \phi w - \phi^T t$$

$$H = \nabla \nabla E(w) = \phi^T \phi$$

$$\text{Newton Raphson update: } w_{\text{new}} = (\phi^T \phi)^{-1} \phi^T t$$

$$\text{Cross entropy error function: } E(w) = - \sum_{n=1}^N \{ t_n \ln y_n + (1-t_n) \ln (1-y_n) \}$$

$$\nabla E(w) = \sum_{n=1}^N (y_n - t_n) \phi_n = \phi^T (y - t)$$

$$H = \sum_{n=1}^N y_n(1-y_n) \phi_n \phi_n^T = \phi^T R \phi$$

Newton-Raphson update: $w^{(\text{new})} = (\phi^T R \phi)^{-1} \phi^T R z$

$$z = \phi w^{(\text{old})} - R^{-1}(y - b)$$

$$R_{nn} = y_n(1-y_n)$$

3) MULTICLASS LOGISTIC REGRESSION:

$$p(C_k|\phi) = y_k(\phi) = \frac{e^{a_k}}{\sum e^{a_j}} ; a_k = w_k^T \phi$$

Likelihood function using 1-of-K coding scheme.

$$p(T|w_1, \dots, w_K) = \prod_{n=1}^N \prod_{k=1}^K p(C_k|\phi_n)^{t_{nk}} = \prod_{n=1}^N \prod_{k=1}^K t_{nk} y_n^{t_{nk}}$$

Cross entropy error function for multiclass classification

$$E(w_1, \dots, w_K) = -\ln p(T|w_1, \dots, w_K) = -\sum_{n=1}^N \sum_{k=1}^K t_{nk} \ln y_n$$

$$\nabla_{w_k} E(w_1, \dots, w_K) = \sum_{n=1}^N (y_n - t_{nj}) \phi_n$$

Same form, product of error times the basis function.

Hessian matrix: $(\phi^T \phi)^{-1}$

$$\nabla_{w_k} \nabla_{w_j} E(w_1, \dots, w_K) = -\sum_{n=1}^N y_n (1 - y_n) \phi_n \phi_n^T$$

IRLS algorithm is also used for batch processing. → The cost is 0 if both the predicted, actual value are same.

→ If not same, calculate loss function.

$$\min_w \lambda \|w\|^2 + \sum_{i=1}^n (1 - Y_i(x_i; w))^2$$

Loss fn for SVM

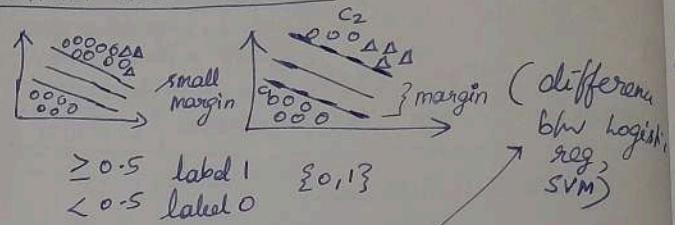
4) PROBIT REGRESSION:



→ for broad range of class conditional distributions, resulting posterior class prob are given by logistic or softmax transformation acting on a linear function of the feature variables.

→ Not the case for all choices of class-conditional density
→ Not worth exploring other types of discriminative probabilistic model.

SUPPORT VECTOR MACHINE:



→ Drawback of logistic Reg:

→ If o/p > 1 (beyond range) it won't classify

Why SVM?

→ Even if o/p > 1 it will classify. (outliers will also be classified).

HYBRID LOSS:

Helps maximise margin

$$C(x, y, f(x)) = \begin{cases} 0 & \text{if } y \neq f(x) \\ 1 - y f(x) & \text{else} \end{cases}$$

→ The cost is 0 if both the predicted, actual value are same.

→ If not same, calculate loss function.

$$\min_w \lambda \|w\|^2 + \sum_{i=1}^n (1 - Y_i(x_i; w))^2$$

Loss fn for SVM

$\Delta E(w)$

$$\frac{\partial}{\partial w_k} \Delta E(w) = 2 \gamma w_k$$

$$\frac{\partial}{\partial w_k} (1 - y_i(x_i, w))_+ = \begin{cases} 0 & \text{if } y_i(x_i, w) \geq 1 \\ -y_i x_{ik} & \text{else} \end{cases}$$

gradient

$w = w - \alpha (2 \gamma w)$ → Gradient update

When there is misclassification

$$w = w + \alpha (y_i x_i - 2 \gamma w)$$

gradient update when there is misclass.

SUPPORT VECTOR REGRESSION:

- SVM-Kernel → Combined with SVM to make the segregation of non-separable data.
- Takes low dimensional D/P, transforms into higher dimensional
- Non-separable to separable problem.

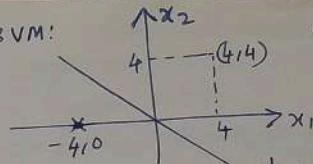
PROS OF SVM:

- Clear margin of separation
- High dimensional spaces.
- No. of dimensions > No. of samples.
- Memory efficient as it uses a subset of training points.

CONS OF SVM:

- Large dataset
- Hard noise

Q) SVM:



SVM

$$y = mx + c$$

$$y = w^T x + b$$

$$y = w^T x + \delta$$

intercept

$$y = \begin{bmatrix} -1 \\ 0 \end{bmatrix} [4 \ 0]$$

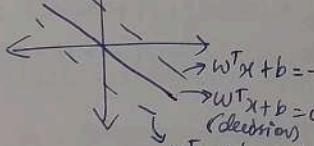
slope
intercept

$$y = \begin{bmatrix} -1 \\ 0 \end{bmatrix} [4 \ 4]$$

= 4 (+ve)
so any point below the plane,

= -4 (-ve)
so any point above plane -ve.

Q) SVM



$$w^T x_1 + b = -1 \quad \text{--- (1)}$$

$$w^T x_2 + b = +1 \quad \text{--- (2)}$$

$$w^T (x_3 - x_2) = -2$$

$$w^T (x_2 - x_1) = 2$$

$$\frac{w^T (x_2 - x_1)}{\|w\|} = \frac{2}{\|w\|} \rightarrow \text{optimisation}$$

norms

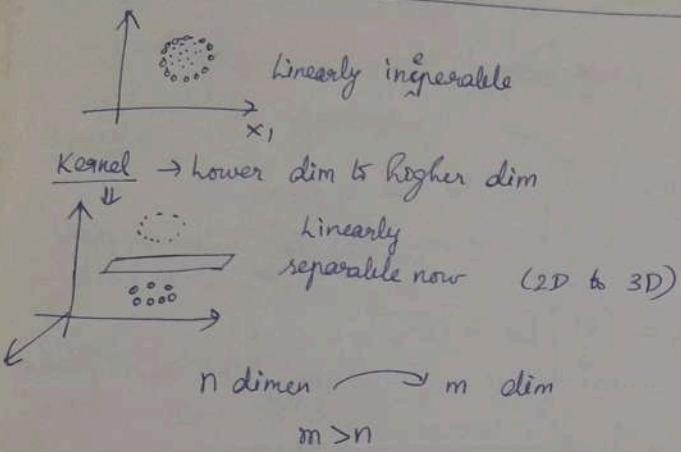
$$(w^*, b^*) \text{ max } \frac{(c)}{\|w\|} \text{ such that } y_i = \begin{cases} +1 & w^T x_i + b \geq 1 \\ -1 & w^T x_i + b \leq -1 \end{cases}$$

↓ Maximum margin ↓ $y_i x_i w^T x_i + b_i \geq 1$

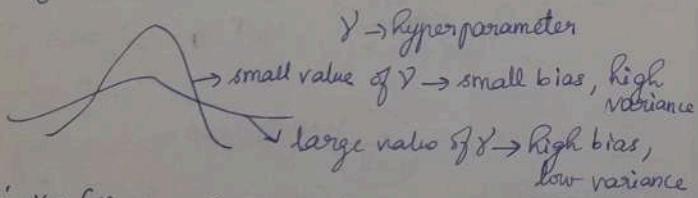
$$(w^*, b^*) = \min \frac{\|w\|^2}{2} + C \sum_{i=1}^n \xi_i \rightarrow \text{value of the error}$$

(Regulariser)

SVM KERNEL:



→ Polynomial kernel $k(x, x') = (1 + x \cdot x')^k$
 → RBF Kernel $k(x, x') = e^{-\gamma \|x - x'\|^2}$
 → Sigmoid Kernel



Eg: $x = (x_1, x_2, x_3) \quad y = (y_1, y_2, y_3)$

 $f(x) = (x/x_1, x/x_2, x/x_3); \quad k(x, y) = \langle f(x), f(y) \rangle$
 $(x_1 x_1, x_1 x_2, x_1 x_3, x_2 x_1, x_2 x_2, x_2 x_3, x_3 x_1, x_3 x_2, x_3 x_3) \quad \langle x, y \rangle \rightarrow \text{dot product.}$
 $k(x, y) = (\langle x, y \rangle)^2$
 $x = (1, 2, 3) \quad y = (4, 5, 6)$
 $f(x) = (1, 2, 3, 2, 4, 6, 3, 6, 9) \rightarrow \text{multiplying } x \text{ with } x \text{ itself.}$
 $f(y) = (4, 20, 24, 20, 25, 30, 24, 30, 36) \rightarrow \text{multiplying } y \text{ with itself.}$

$$\langle f(x), f(y) \rangle = \underbrace{16+40+72+40+100+180+72+180}_{\text{dot product}} + 324 = 1024$$

$$K(x, y) = (4+10+18)^2 = 32^2 = 1024$$

$$K(x, x') = (1+x+x')^2$$

$$f(x) = (1, 2, 3) \text{ multiplying}$$

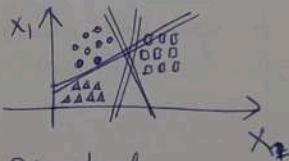
$$f(y) = (4, 5, 6)$$

SVM for multi-class classification:

Techniques:

- ① One v.s. One (OVO)
- ② One v.s. all (OVA)
- ③ Directed Acyclic Graph (DAG)

1) ONE vs. ONE:

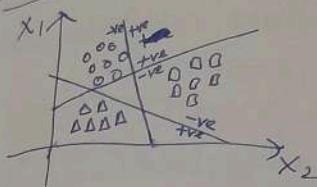


only 2 groups considered at a time and those two are separated by the hyperplane

Drawback:

- More no. of SVMs need to be trained.
- Too much computation required.
- Unbalancing (difference in the no. of points) occurs in the SVM.

2) ONE VS. ALL:



One class separated from the rest.

→ No unbalancing of data point arises as no data is left out.

Drawback:

→ Computation problem ∵ n no of classifiers required.

3) DIRECTED ACYCLIC GRAPH:

→ After logical grouping ^{only} this should be applied.
→ Requires only few no. of SVMs, compared to OVA.

SUPPORT VECTOR REGRESSION (SVR):

→ Use same principle as SVM.

→ Instead of classification, regression is used.

→ This is the only difference.

OVERRAP CLASS DISTRIBUTION:

→ Training data may not be linearly separable, which leads to misclassification of data points.

→ Resulting SVM gives exact separation in input space \mathbb{X} , although decision boundary will be non-linear.

→ Exact separation leads to poor generalization, so SVM should be allowed to misclassify some training points.

Modifying SVM:

→ Implicit error function of SVM:

$$\sum_{n=1}^N \xi_n (y(x_n) \cdot b_n - 1) + \lambda \|w\|^2$$

→ Gives infinite error if point is misclassified.

→ Gives zero error if it was classified correctly.

→ Parameters are optimised to maximise margin.

→ Points are allowed to lie on the wrong side but ensuring that penalty is proportional to distance from boundary.

Slack variables: ξ_n

$$\xi_n \geq 0, n=1, \dots, N$$

→ One slack variable for each training point.
→ Defined by $\xi_n = 0$ for data points that are on or inside correct margin boundary.

→ $\xi_n = |b_n - y(x_n)|$ for other points.

→ Data point that is on decision boundary, $y(x_n) = 0$ has $\xi_n = 1$.

→ Points with $\xi_n \geq 1$ misclassified.

Soft Margin Classification:

- Allows some data points to be misclassified.
- While slack variables allow overlapping class distributions, framework is sensitive to outliers.

Optimisation with slack variables:

- To maximise the margin while softly penalising points that lie on the wrong side of margin boundary.
- Minimise $C \sum_{n=1}^N \xi_n + \frac{1}{2} \|w\|^2$
- Parameter $C > 0$ controls trade-off b/w slack variables penalty and margin.

Subject to constraints:

$$\ln y(x_n) \geq 1 - \xi_n, n = 1, \dots, N$$

Lagrangian with slack:

$$L(w, b, \alpha) = \frac{1}{2} \|w\|^2 + C \sum_{n=1}^N \xi_n - \sum_{n=1}^N \alpha_n [\ln y(x_n) - 1 + \xi_n] - \sum_{n=1}^N \mu_n \xi_n$$

$\{\alpha_n \geq 0\}, \{\mu_n \geq 0\} \rightarrow$ Lagrange multipliers.

KKT conditions:

$$\alpha_n \geq 0$$

$$\ln y(x_n) - 1 + \xi_n \geq 0$$

$$\alpha_n \{\ln y(x_n) - 1 + \xi_n\} = 0$$

$$\mu_n \geq 0$$

$$\xi_n \geq 0$$

$$\mu_n \xi_n = 0$$

BAYESIAN LOGISTIC REGRESSION:

- Discriminative probabilistic linear classifier.

$$p(c_i | \phi) = \sigma(w^T \phi)$$

- Exact Bayesian inference for Logistic Regression $p(c_i | \phi) = \int \sigma(w^T \phi) p(w) dw$ is not easy to control, because:

- Evaluation of posterior distribution $p(w | t)$
Needs normalisation of prior $p(w) = N(w | m_0, S_0)$
times likelihood (product of sigmoids)
Solution: Use Laplace approx.

- Evaluation of predictive distribution
Convolution of sigmoid and Gaussian.

LAPLACE APPROXIMATION:

- Need mode w_0 of posterior distribution $p(w | t)$
(Done by numerical optimisation algorithm).

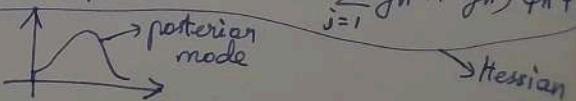
- Fit a Gaussian centered at mode

$$g(w) = \frac{1}{N} f(w) = \frac{|A|^{1/2}}{(2\pi)^{M/2}} e^{-\frac{1}{2}(w-w_0)^T A (w-w_0)}$$

$$= N(w | w_0, A^{-1})$$

- Needs second derivatives of log posterior.
Equivalent to finding Hessian matrix.

$$S_N = -\nabla \nabla \ln p(w | t) = S_0^{-1} + \sum_{j=1}^N y_n (1 - y_n) \phi_n \phi_n^T$$



EVALUATION OF POSTERIOR DISTRIBUTION:

→ Gaussian prior

$$p(w) = N(w/m_0, S_0)$$

where m_0, S_0 are hyper-parameters.

→ Posterior distribution

$$p(w|t) \propto p(w)p(t|w)$$

$$t = (t_1, \dots, t_N)^T$$

$$\text{Substituting } p(t|w) = \prod_{n=1}^N t_n^{y_n} (1-t_n)^{1-y_n}$$

$$\ln p(w|t) = -\frac{1}{2} (w - m_0)^T S_0^{-1} (w - m_0)$$

$$+ \sum_{i=1}^N (t_n \ln y_n + (1-t_n) \ln (1-y_n) + \text{const})$$

Where,

$$y_n = \sigma(w^T \phi_n)$$

GAUSSIAN APPROXIMATION OF POSTERIOR:

→ Maximise posterior $p(w|t)$ to give

- MAP solution w_{map}

(Done by numerical optimisation)

- Defines mean of Gaussian.

→ Covariance given by

- Inverse of matrix of 2nd derivatives of negative log-likelihood $S_N = -\nabla \nabla \ln p(w|t) = S_0^{-1} + \sum_{i=1}^N t_i (1-t_i) \phi_i \phi_i^T$

→ Gaussian approximation to posterior

$$q(w) = N(w/w_{\text{map}}, S_N)$$

→ Need to marginalise w.r.t this distribution to make predictions.

Predictive Distribution:

→ Predictive dist for class C_1 , given new feature vector $\phi(x)$.

→ Obtained by marginalising w.r.t posterior $p(w|t)$.

$$\begin{aligned} p(C_1 | \phi, t) &= \int p(C_1, w | \phi, t) dw \text{ (sum rule)} \\ &= \int p(C_1 | \phi, t, w) p(w | t) dw \text{ (Product rule)} \\ &= \int p(C_1 | \phi, w) p(w | t) dw \text{ (Given } \phi, w, C_1 \text{ is independent of } t) \\ &\approx \int \sigma(w^T \phi) q(w) dw \text{ Approximate } p(w | t) \\ &\quad \text{by Gaussian } q(w) \end{aligned}$$

Corresponding probability for class C_2

$$p(C_2 | \phi, t) = 1 - p(C_1 | \phi, t)$$

Predictive Dist. is a convolution:

$$p(C_1 | \phi, t) = \int \sigma(w^T \phi) q(w) dw$$

→ Function $\sigma(w^T \phi)$ depends on w only through its projection onto ϕ .

$$\rightarrow a = w^T \phi, \sigma(w^T \phi) \approx \int \delta(a - w^T \phi) \sigma(a) da \quad (\text{where } \delta \text{ is Dirac delta fn})$$

$$\rightarrow \text{Thus } \int \sigma(w^T \phi) q(w) dw = \int \sigma(a) p(a) da$$

$$p(a) = \int \delta(a - w^T \phi) q(w) dw$$

→ Can evaluate $p(a)$ because:

delta function imposes linear constraint on w , since $q(w)$ is Gaussian, its marginal is also Gaussian

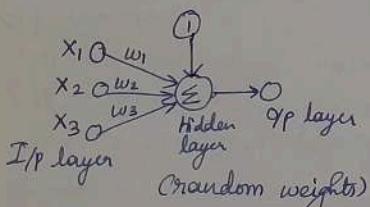
Evaluate its mean, covariance:

$$\mu_a = E(a) = \int p(a) da = \int q(w) w^T \phi dw = w_{\text{avg}}^T \phi$$

$$\sigma_a^2 = \text{var}[a] = \int p(a) \{a^2 - E[a]^2\} da$$

$$= \int q(w) \{(w^T \phi)^2 - (w_{\text{avg}}^T \phi)^2\} dw = \phi^T S_N \phi$$

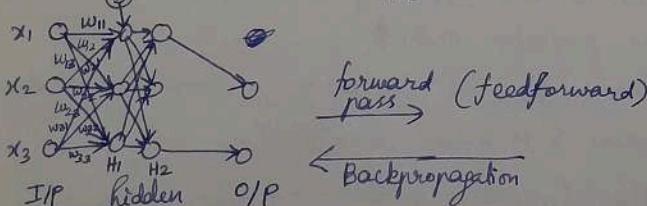
BACKPROPAGATION ALGORITHM:



$$H_1 = w_1 x_1 + w_2 x_2 + w_3 x_3 + b$$

$$Z = f(H_1) = \text{activation}(H_1)$$

$$\text{Sigmoid } \hat{z} = \frac{1}{1+e^{-x}} = \frac{1}{1+e^{-h_1}}$$



Multilayer perception

Input (x_i) Desired O/P

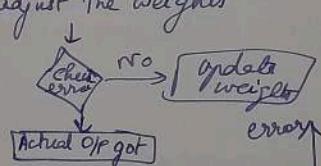
0 0
1 2
2 4

$w=3:$	$(w \times x)$	Predicted O/P (model)	$(y - wx)$	Error
		0	0	0
		3	1	1
		6	2	2

$w=4:$	I/P (x)	desired O/P (y)	$(wx \times x)$	Error
	0	0	0	0
	1	2	4	2
	2	4	8	4

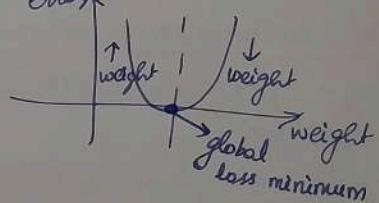
$w=2:$	I/P	y	wx	$(y - wx)$
	0	0	0	0
	1	2	2	0
	2	4	4	0

To calculate the error,
adjust the weights

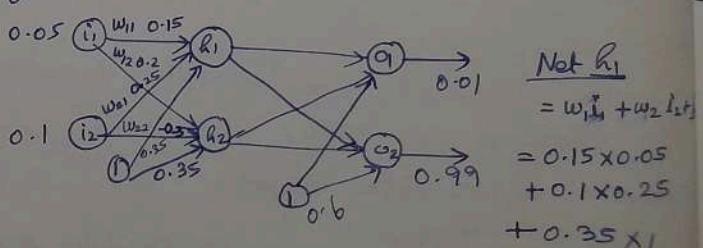


Updating weights,

Back prop.



Eg:



$$O/P \text{ of } H_1 : \frac{1}{1+e^{-H_1}} = \frac{1}{1+e^{-(\text{---})}} = 0.5932$$

$$H_2 = w_{21}i_1 + w_{22}i_2 + b_2x_1 = 0.5968$$

$$O/P \text{ of } H_2 = 0.5968$$

$$\text{Net } O_1 = w_{31}h_1 + w_{32}h_2 + b_1x_1 = 1.1059$$

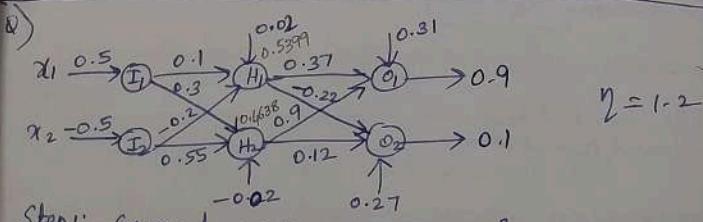
$$O/P \text{ of } O_1 = \text{Sigmoid of } O_1 = \frac{1}{1+e^{-9}} = 0.7513$$

$$O_2 = w_{41}h_1 + w_{42}h_2 + b_2x_1 = 0.7729$$

$$\begin{aligned} \text{Error } O_1 : EO_1 &= \sum \frac{1}{2} (\text{target} - O/P)^2 \\ &= \frac{1}{2} (0.99 - 0.7513)^2 = 0.274 \quad E_1 \end{aligned}$$

$$EO_2 = \frac{1}{2} (0.99 - 0.7729)^2 = 0.0235 \quad E_2$$

$$\text{Total error} = EO_1 + EO_2 = 0.274 + 0.0235 = 0.2983$$



Step 1: forward pass, $Z = \sum x_i w_{ik}$

$$\begin{aligned} Z_1' &= x_1' w_{11} + x_2' w_{21} + b_1 x_0 \\ &= 0.5 \times 0.1 + (-0.5 \times -0.2) + 0.31 \\ &= 0.16 \end{aligned}$$

$$\begin{aligned} Z_2' &= (0.5 \times 0.55) + (0.3 \times 0.5) + (-0.02 \times 1) \\ &= -0.145 \end{aligned}$$

STEP 2: O/P from hidden layers Activation fn: sigmoid

$$f(z_1') = \frac{1}{1+e^{-z_1'}} = \frac{1}{1+e^{-0.16}} = 0.5399$$

$$f(z_2') = \frac{1}{1+e^{-z_2'}} = \frac{1}{1+e^{-0.145}} = 0.4638$$

STEP 3: Hidden to O/p layer

$$\begin{aligned} y_1' &= 0.5399 \times 0.31 + 0.4638 \times 0.9 + 0.31 \times 1 \\ &= 0.9272 \end{aligned}$$

$$\begin{aligned} y_2' &= 0.4638 \times -0.12 + 0.5399 \times -0.22 + 0.27 \times 1 \\ &= -0.0955 \end{aligned}$$

STEP 4:

$$f(y_1') = \frac{1}{1+e^{-y_1'}} = \frac{1}{1+e^{-0.9272}} = 0.7164$$

$$f(y_2') = \frac{1}{1+e^{-0.0955}} = 0.5238$$

Backward pass:

$$d_1' = 0.9, d_2' = 0.1 \quad \text{Error:}$$

$$[e_1'] = d_1' - f(y_1') = 0.9 - 0.7164 = 0.1836$$

$$[e_2'] = d_2' - f(y_2') = 0.1 - 0.5238 = -0.4238$$

O/P to hidden layer:

$$\begin{aligned} \delta_1' &= e_1' \times f'(y_1) ; \delta'(y_1) = f(y_1)(1-f(y_1)) \\ &= 0.1836 \times (0.7164)(1-0.7164) = 0.0373 = \delta_1' \end{aligned}$$

$$\begin{aligned} \delta_2' &= e_2' \times f'(y_2) (1-f(y_2)) \\ &= -0.4238 \times (0.5238)(1-0.5238) \end{aligned}$$

$$\delta_2' = -0.1057$$

Hidden to I/O layer:

$$\begin{aligned}\delta'_1(z_1) &= (\delta'_1 w + \delta'_2 w) \delta' z_1 \\ &= [(0.0373 \times 0.37) + (-0.1057 \times -0.22)] (\delta(z_1)(1-\delta(z_1))) \\ &= [(0.0373 \times 0.37) + (-0.1057 \times -0.22)] \times 0.5399 (1-0.5399) \\ &= 0.0092\end{aligned}$$

$$\begin{aligned}\delta'_2(z_2) &= \delta'(z_2) (\delta_1 w + \delta_2 w) \\ &= [(0.0373 \times 0.7) + (0.1057 \times 0.12)] \times (0.4638)(1-0.4638)\end{aligned}$$

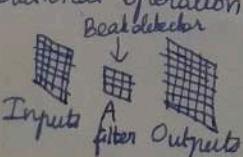
$$\delta'_2(z_2) = 0.0115$$

Weight change for pass 1:

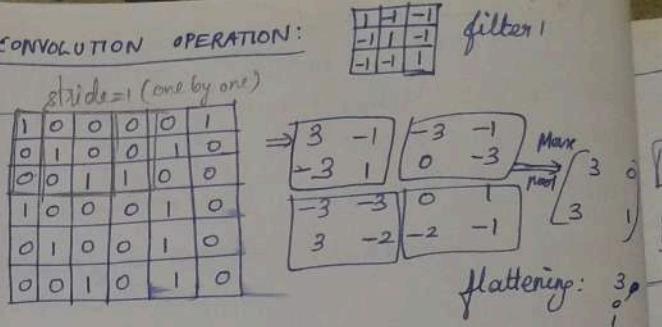
$$\begin{aligned}w_{new} &= w_{old} - \eta \frac{\delta E}{\delta w_{ij}} \\ \Delta'_{n1} &= 0.31 + 1.2 \times 0.0373 \times 1 = 0.0447 + 0.31 = 0.3547 \\ \Delta'_{n2} &= 0.31 + 1.2 \times 0.0373 \times 0.5399 = 0.0241 + 0.37 = 0.3941 \\ \Delta'_{h21} &= [0.9] + (1.2 \times 0.0373 \times 0.4638) = 0.9207 \\ \Delta n_{21} &= \end{aligned}$$

CONVOLUTIONAL LAYER

- CNN is a neural network with some convolutional layers, and some other layers.
- Convolutional layer has no. of filters that does convolutional operation.



CONVOLUTION OPERATION:



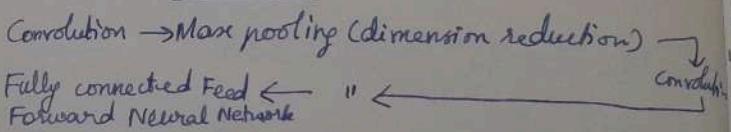
RGB → (3 channels).

CONVOLUTION vs. FULLY CONNECTED N/L:

→ In convolution, the ~~values~~ are finally taken as a single row and the dimension reduces.

→ Dropout is present (those values of less importance are not connected to the O/P).

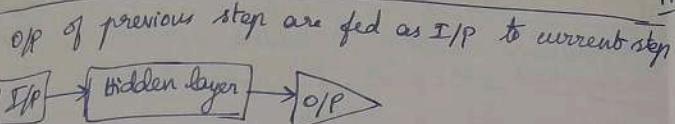
Convolution PROCESS:



COMPRESSION IN CNN:

- Reducing no. of connections
- Shared weights on the edges
- Max pooling reduces complexity.

RNN:



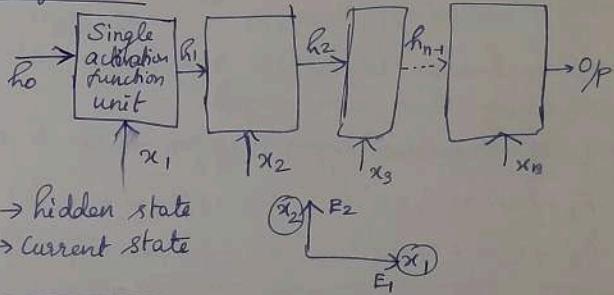
PROS: RNN → Time series problems (remembering is demanded)
Eg: Weather prediction

RNN → LSTM (Long Short Term Memory).

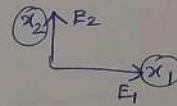
CONS:

- Gradient vanishing and exploding.
- RNN training → difficult
- Can't process very long sequences.

Flow of RNN:



h0 → hidden state
x1, x2, ..., xn → current state

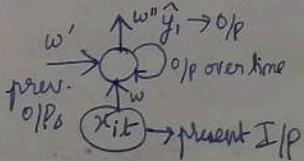


x1, x2, x3, ..., xn → given at each time instant

→ that's why RNN is time series problem.
→ O/P generated at each time instant.

NATURAL LANGUAGE PROCESSING:

- I/P (text data)
- BOW, TF-IDF, WORD2VEC → algorithms



RNN WORKING:

```

graph LR
    I[I/P] --> H1((H1))
    H1 --> H2((H2))
    H2 --> H3((H3))
    H3 --> O[O/P]
  
```

Hidden layers

$h_1 = f(h_{t-1}, x_t)$ → I/P of current state
 ↓
 (Current state) (prev. hidden layer state)

$h_2 = f(h_1, x_t)$ → I/P of current state
 ↓
 (prev. hidden layer state)

$h_3 = f(h_2, x_t)$ → I/P of current state
 ↓
 (prev. hidden layer state)

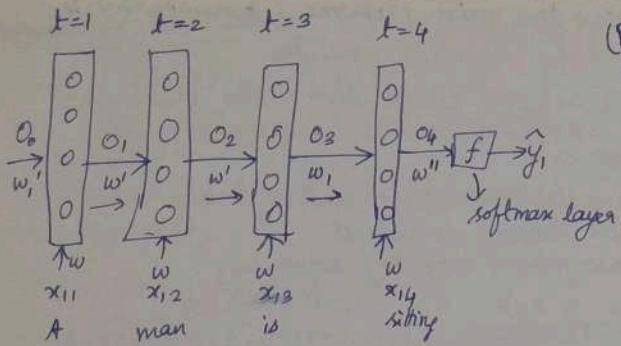
→ Each layer is independent of each other
 → 1st layer O/P → I/P of 2nd.

Formula for applying activation function (on h)

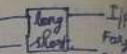
$$h_t = \text{tanh}(W_{ah} h_{t-1} + W_{xh} x_t)$$

(recurrent neuron) (I/P neuron)

O/P: $y_t = w_y \times h_t$ Way → weight at O/P layer.
 $y_t \rightarrow O/P$.

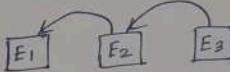


(Better accuracy) **LSTM (Long short-term memory)**

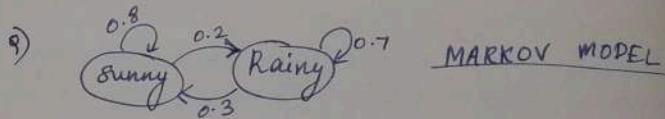


→ Vanishing/exploding gradient overcome.
→ Since RNN can retain info only for short time.

HIDDEN MARKOV MODEL!



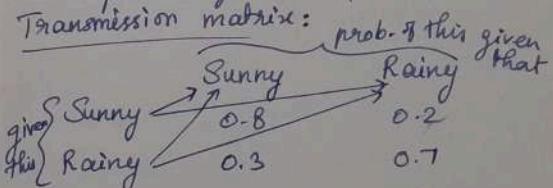
E_n (Event n) depends on E_{n-1} (previous event)



States: Sunny, Rainy

Initial probability: {0.4, 0.6}

Transmission matrix:



$Q = \{ \text{Sunny, Rainy, Rainy} \}$
Today Tom day aftertom

$$P(Q) = P(S) * P(R/S) * P(R/R)$$

$$= 0.4 * (0.2) * (0.7)$$

$$= 0.08 * 0.7 = 0.056$$

→ Same hidden layers are repeated, so only all are carrying their own weights.

→ O/P of 1st hidden layer is given as I/P to 2nd.

$$O_1 = f(x_1 w + O_0 w')$$

$$O_3 = f(x_3 w + O_2 w)$$

$$O_2 = f(x_2 w + O_1 w')$$

$$O_4 = f(x_4 w + O_3 w')$$

$\delta P = \text{prev} + \text{present state}$

$$\text{Loss} = (\hat{y} - y)$$

reduce the loss to get desired o/p.

$$\frac{\partial L}{\partial \hat{y}_1} \times \frac{\partial L}{\partial w} = \frac{\partial L}{\partial y_1} \times \frac{\partial \hat{y}_1}{\partial O_4} \times \frac{\partial O_4}{\partial w} \Rightarrow \frac{\partial L}{\partial w}$$

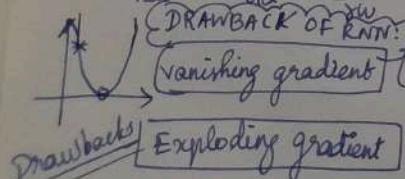
$$w = w + \frac{\partial L}{\partial w}$$

Vanishing, exploding gradients

$$w_{\text{new}} = w_{\text{old}} + \frac{\partial L}{\partial w}$$

negligible weight updation,
so waste of processing

(the weight keeps changing
randomly, doesn't reach
global optimum)



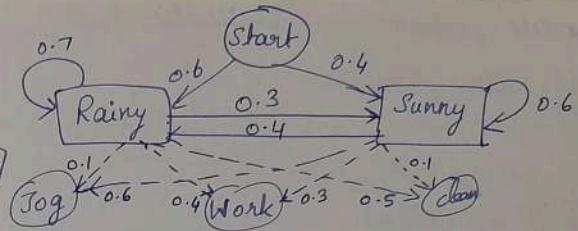
Exploding gradient

MARKOV MODEL

- Observables are visible (explicit)
- Eg: If it will be sunny or rainy is visible (explicit)

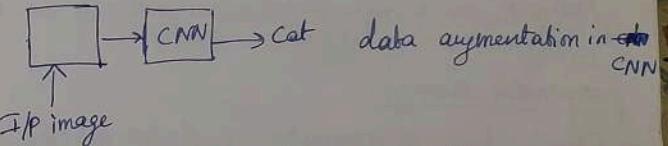
HIDDEN MARKOV MODEL

- Observables are not visible (implicit)
- Eg: If it is sunny or rainy is not visible (implicit)
- Whether the weather is sunny or rainy can be implicitly drawn from {cotton, woollen clothes}



Hidden Markov model

DATA AUGMENTATION: PARAMETER TRYING, SHARING IN CNN



→ Due to data augmentation, sample images are increased, so better model, better results, more useful features.

Parameter tying, sharing!

- To reduce complexity.
- To reduce computation time

→ Types:

- Other methods for prior knowledge of parameters
 - Parameter tying
 - Parameter sharing
- (w) " " in CNNs.

HIDDEN MARKOV MODEL:

- One of the statistical models used in ML algorithms to go for classification problems.

Eg: {cotton, woollen} \Rightarrow {sunny, rainy}

Q) States : {Rainy, Sunny}

Observations : {Jogging, office, cleaning}

Starting probability : {Rainy = 0.6, Sunny = 0.4}

Transition probability $P_{(S|S)}$: {Actual Probabilities of 2 states:
Rainy = 0.7, sunny = 0.3
Rainy = 0.4, sunny = 0.6}

Emission probability = {

Rainy : {Jog = 0.1, work = 0.4, clean = 0.5}

Sunny : {Jog = 0.6, work = 0.3, clean = 0.1}

2 models perform some classification task:

model 1 with parameters of w^A
 model 2 " " " w^B

$$\hat{y}^{(A)} = f(w^{(A)}, x)$$

$$y^{(B)} = f(w^{(B)}, x)$$

$$w_i^{(A)} \approx w_i^{(B)}$$

parameter norm penalty:

$$\rightarrow (w^{(A)}, w^{(B)}) = \|w^{(A)} - w^{(B)}\|_2^2 \quad (\text{similar features tied})$$

\rightarrow Used for regularising parameters of one model trained as a supervised classifier to be closer to parameters of another model.

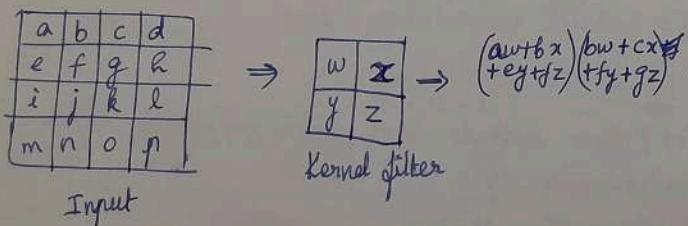
parameter sharing:

\rightarrow Single type of parameter used for all the tasks.

Eg: Parameters c_1, c_2

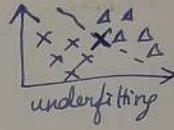
Tasks T_1, T_2

Either c_1 used for both T_1, T_2 or c_2 used for both T_1, T_2 .



\rightarrow Done because transferring the ~~image~~ pixels or scaling them or doing some little changes doesn't make a big effect on the overall image.

REGULARISATION:
EARLY STOPPING:



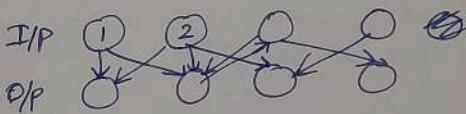
Training, testing phase

~~90% error~~ 90%
 (27. error)

\rightarrow Here the error has nearly saturated, further epoch don't change any of these parameters, hence, early stopping is necessary to save time.

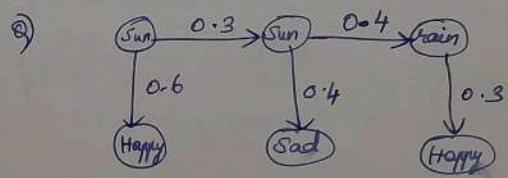
DROPOUT:

\rightarrow Helps reducing interdependent learning among the neurons.



\rightarrow When more than one neuron convey the same info, redundant (non-contributing) neurons are removed to save time, reduce complexity.

FORWARD ALGORITHM FOR HMM:



Current time step to derive the probability of the next time step.

Given a sequence of visible state V_t , what will be the probability that the hidden Markov model will be in a particular hidden state S at a particular time step t ?

DERIVATION:

$$\alpha_j(t) = P(v(1) \dots v(t), s(t)=j)$$

Solution using probabilities.

When $t=1$:

$$\begin{aligned} \alpha_j(1) &= P(v_k(1), s(1)=j) \\ &= P(v_k(1) | s(1)=j) P(s(1)=j) \\ &= \pi_j \underbrace{P(v_k(1) | s(1)=j)}_{b_{jk}} \end{aligned}$$

$$\boxed{\alpha_j(1) = \pi_j \times b_{jk}}$$

When $t \geq 2$:

$$\begin{aligned} \alpha_j(2) &= P(v_k(1), v_k(2), s(2)=j) \\ &= \sum_{i=1}^M P(v_k(1), v_k(2), s(1)=i, s(2)=j) \\ &= \sum_{i=1}^M P(v_k(2), s(2)=j, v_k(1) \cancel{s(1)=i}) \\ &\quad \times P(v_k(1), s(2)=i) \\ &= \sum_{i=1}^M P(v_k(2) | s(2)=j) P(v_k(1) | s(1)=i) \\ &\quad \times P(s(2) | v_k(1), s(1)=i) \\ &\quad \times P(v_k(1), s(1)=i) \end{aligned}$$

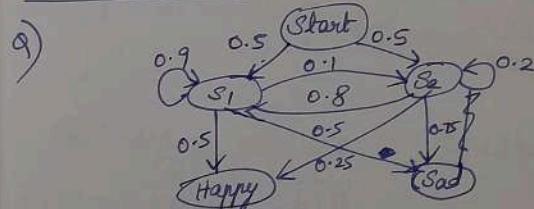
$$\begin{aligned} &= \sum_{i=1}^M P(v_k(2) | s(2)=j) P(s(2) | s(1)=i) \\ &\quad \times P(v_k(1), s(1)=i) \\ &= P(v_k(2) | s(2)=j) \underbrace{\sum_{i=1}^M P(s(2) | s(1)=i)}_{a_{i2}} \underbrace{P(v_k(1), s(1)=i)}_{\alpha_i(1)} \\ &\quad \times b_{jk} \end{aligned}$$

$$\boxed{= b_{jk} \sum_{i=1}^M a_{i2} \alpha_i(1)}$$

Any time step $t+1$:

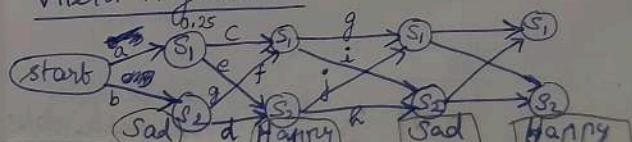
$$\begin{aligned} \alpha_j(t+1) &= P(v_k(1) \dots v_k(t+1), s(t+1)=j) \\ &= \sum_{i=1}^M P(v_k(1) \dots v_k(t+1), s(t)=i, s(t+1)=j) \\ &= b_{jk} v_k(t+1) \underbrace{\sum_{i=1}^M a_{ij} \alpha_i(t)}_{\alpha_j(t)} \end{aligned}$$

TRELLIS DIAGRAM:



Pattern \rightarrow Sad - Happy - Sad - Happy

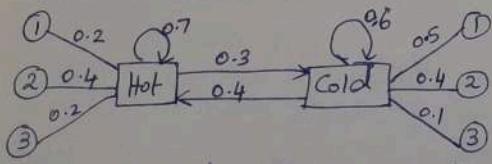
Viterbi Algorithm:



	Sad	Happy	Sad	Happy
Sad				
a \Rightarrow	$0.5 \times 0.5 = 0.25$ (Start $\rightarrow S_1$) $S_1 \rightarrow S_{\text{Sad}}$)			
c \Rightarrow	0.25×0.9 (state) ($S_1 \rightarrow S_2$) value $\times 0.5 = 0.1125$ ($S_2 \rightarrow \text{Happy}$)			
b \Rightarrow	$0.5 \times 0.75 = 0.375$			
f \Rightarrow	$0.375 \times 0.8 \times 0.75 = 0.225$			
d \Rightarrow	$0.25 \times 0.2 =$			

DECODING WITH VITERBI ALGORITHM : (HMM)

$$\pi = [0.8, 0.2] \quad H \quad C$$

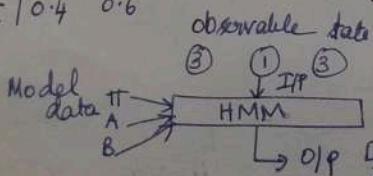
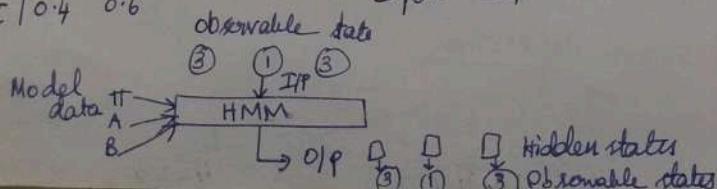


Transition matrix A

$$\begin{array}{|c|c|} \hline & H & C \\ \hline H & 0.7 & 0.3 \\ \hline C & 0.4 & 0.6 \\ \hline \end{array}$$

$$\begin{array}{|c|c|c|} \hline & 1 & 2 & 3 \\ \hline H & 0.2 & 0.4 & 0.4 \\ \hline C & 0.5 & 0.4 & 0.1 \\ \hline \end{array}$$

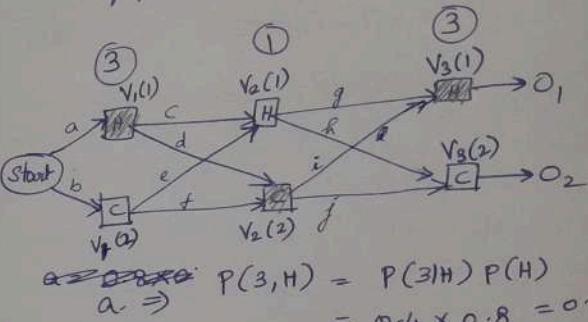
1, 2, 3 \rightarrow Prob of taking that many icecreams.



$N=2, T=3 \downarrow$
no. of states \rightarrow No. of observables (1, 2, 3)
No. of possible outcomes $= N^T = 2^3 = 8$

$$P(3, H) = P(3/H) P(H)$$

$$a \Rightarrow 0.4 \times 0.8 = 0.32 \quad \text{Initial state}$$



$$b \Rightarrow P(3, C) = P(3/C) P(C) = 0.1 \times 0.2 = 0.02$$

$$c \Rightarrow P(1, H) = [P(1/H) \times P(H)] = P(1, H) = 0.2 \times 0.8 = 0.16$$

$$d \Rightarrow P(1/C) \times P(C/H) = 0.5 \times 0.3 = 0.15$$

$$e \Rightarrow P(1/H) \times P(H/C) = 0.2 \times 0.4 = 0.08$$

$$f \Rightarrow P(1/C) \times P(C/C) = 0.5 \times 0.6 = 0.3$$

$$v_1(1) = 0.32, v_1(2) = 0.02$$

$$v_2(1) = 0.32 \times c \\ (0.8) \\ 0.02 \times e$$

$$v_2(2) = 0.32 \times d \\ (0.8) \\ 0.02 \times f$$

Max value taken for Viterbi algorithm among the two

$$P(3, H) = P(3/H) P(H/H) \times V_2(0) = [0.4 \times 0.7 \times 0.32c] \Rightarrow \sum_{j=0}^M p_j(t+1) b_{j,k} v(t+1) q_{ij}$$

(or) $P(3/H) P(H/C) \times V_2(2) = [0.4 \times 0.4 \times 0.32d]$

$\nabla_2(2) \rightarrow P(3,C) = P(3/C) P(C/H) \times V_2(1) = [0.1 \times 0.3 \times V_2(1)]$

(or) $P(3/C) P(C/C) \times V_2(2) = [0.1 \times 0.6 \times V_2(2)]$

Now $V_1(1)$ is maximum among $V_1(1), V_1(2)$

~~Note~~ $V_2(2)$ is max

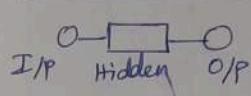
$V_3(1)$ is max

NOTE: For forward algorithm, sum of 2 values taken for a node

For Viterbi algorithm, max of 2 values taken

TYPES OF RNN:

1) ONE-ONE RNN:



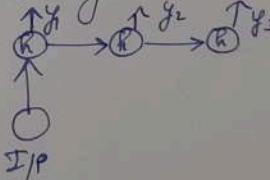
Single

→ (Vanilla's NN)

Applications:
→ Weather prediction
→ Stock prediction

2) ONE - MANY RNN:

→ Only one I/P fed, many O/Ps obtained at each stage.



Eg: Image captioning

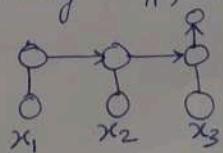
BACKWARD ALGORITHM FOR HMM:

DERIVATION?

$$\begin{aligned} B_i(t) &= P(V_k(t+1) \dots V_k(t) | s(t)=i) \\ &= \sum_{j=0}^M P(V_k(t+1) \dots V_k(t), s(t+1)=j | s(t)=i) \\ &= \sum_{j=0}^M P(V_k(t+2) \dots V_k(t), s(t+1)=j | s(t)=i) \\ &= \sum_{j=0}^M P(V_k(t+1), s(t+1)=j, s(t)=i) \\ &\quad P(s(t+1)=j | s(t)=i) \end{aligned}$$

3) MANY - ONE RNN:

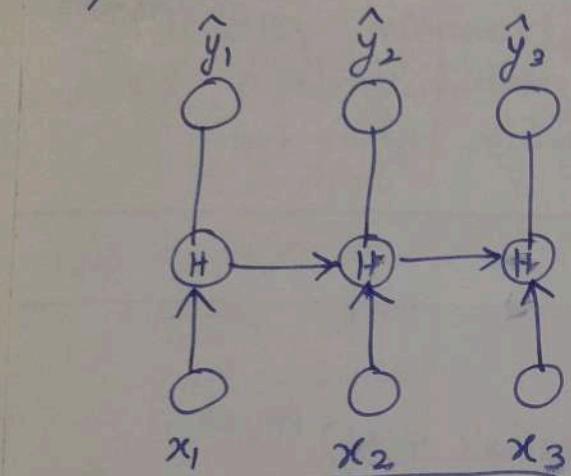
→ Many I/P, one O/P.



Eg: Sequential analysis

→ Movie review rating based on comments.

4) MANY - MANY RNN

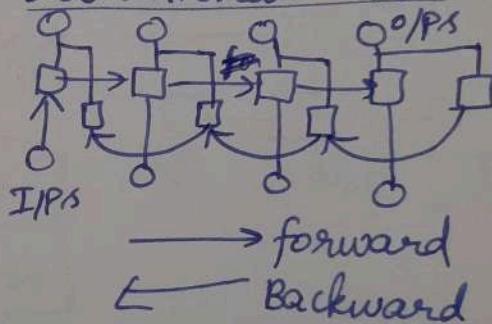


Eg: Google translate

DEEP RNN's:

^{than 1}
→ More hidden layers

Bidirectional RNN:



GRU (Gated Recurrent Units):

- → 2 gates only (update, reset)
- adds new info \downarrow removes old info
- → Computation is less, efficient processing, but less accuracy.