

^{7 simple} Naive Bayes classifier

Discriminative → all you need to know is to which side of line this datapoint belongs to. (Discriminates)

Generative → You need to know what is the underlying process which leads us to the data we are seeing. (eg → Bayesian)
 your task is to learn distribution that gave rise to this data.

Instead of knowing if this point lies here or there, we tend to know which part of feature space belongs to which class.

(In discriminative model you know predict label from features)

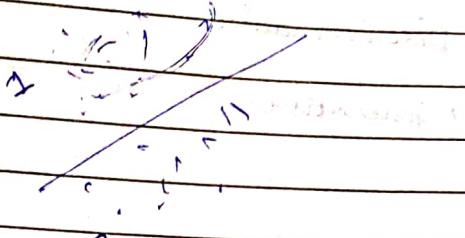
But in Generative → you see what is the probability of getting that feature vector for a class and then we deduce the probability of getting label given the feature vector.

So, basically it is the probability of feature vector given that the label is something
 $P(x|y)$ and then you find probability of label given feature vector $P(y|x)$.

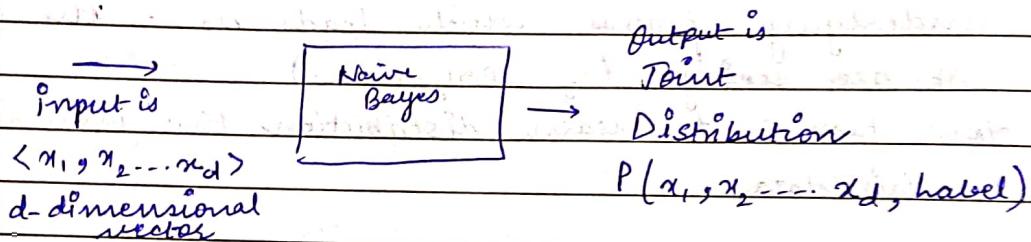
feature vector label

Bayes → Dis

Generative



+ learning distribution
Two and their parameters



Joint Distribution

Joint Probability $\rightarrow P(A \cap B)$ where A and B are two events & joint prob. is the probability of occurrence of both events simultaneously. ($P(x, y)$)

Joint Probability Distribution \rightarrow when x and y are discrete random variables and joint distribution is $P(x, y)$ for every pair of values (x, y) .

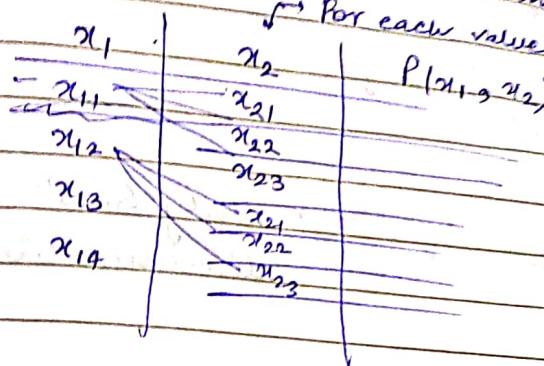
Marginal Distribution →

Marginalization refers to the process of removing the influence of one or more events from a probability.

$$g(x) = \sum_y f(x, y)$$

Marginal of x → Joint Pro. $\Rightarrow n$ removed!

So, lets say we have 2 random variables x_1 & x_2 .



For each value of x_1 , x_2 will have some value

Why Joint Distribution?

→ We can get Conditionals

→ " Marginals

→ " every damn thing!

So, we know everything about data. Anything cannot be better than this.

But not possible to do it everytime. (Hard to implement)

How can you get Marginal from Joint?

$$\text{Marginal} \Rightarrow \sum_{x_2} P(x_1, x_2) = P(x_1) \left[\begin{array}{l} \text{Summation of all} \\ \text{over all } x_2 \\ \text{from table} \end{array} \right]$$

$$\text{Conditional} \Rightarrow P(x_1 / x_2 = b)$$

→ We freeze value of x_2 & see how the distribution of x_1 looks like.
i.e., subset of Joint Distribution table!

we can ans quest' like \rightarrow (From ~~good~~ Joint Distribut' table)

$$P(x_1 < 5 \mid \text{Label} = 1) = ?$$

$$P(x_1/x_2, \text{label}) = ?$$

$P(\text{Label} \mid x_1, x_2) = ? \rightarrow$ Only Ques' that a discriminative model can ans.

Probability Axioms \rightarrow

3 imp things \rightarrow (that define Probability space)

\rightarrow Sample Space

\rightarrow Event Space

\rightarrow Probability function

$$\textcircled{1} \quad 0 \leq P(x) \leq 1$$

$$\textcircled{2} \quad \text{sum of (all probabilities)} \text{ Sample Space } P(S) = 1.$$

$$\textcircled{3} \quad P(A \cup B) = P(A) + P(B) - P(A \cap B)$$

for disjoint events $\rightarrow P(A \cup B) = P(A) + P(B)$

$$\# P(A) = P(A \cap B) + P(A \cap \bar{B})$$

Conditional Probability \rightarrow

$$P(A/B) = \frac{P(A \cap B)}{P(B)}$$

$[P(B) \neq 0]$

Chain Rule →

$$P(A_1 \wedge A_2) = P(A_1) \cdot P(A_2 | A_1)$$

$$P(A_1, A_2, A_3) = P(A_1) \cdot P(A_2 | A_1) \cdot P(A_3 | A_1, A_2)$$

Bayes Rule →

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

Given a the outcome ~~compute~~ determine the feature.

Application in health care ^{company} for clinical trials →

They have → People with Disease = True Group A
 → People without Disease = False Group B

So you will apply test on A & know what condⁿ are (features) giving outcome true & you will select them.

So, ~~P(A True | B False)~~ & ~~P(A True | B True)~~

~~conditional~~ ~~P(Pred true | have disease)~~ ✓ ~~P(Pred true | dont have disease)~~ ✓ You will know first kind of

then, for new people you will know.

Bayes → $P(\text{do this person have disease} | \text{Pred true})$

#

Bayes Rule - Interpretation

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

↳ Likelihood → Prior
↓ Normalization factor

Posterior
Probability

Biggest thing here → P(A) prior
probability ~~than before you saw a test~~

So, Prior * Evidence from test = Posterior Prob.

Eg → Suppose 1% people have ~~can~~ funny disease
so, $P(f) = 0.01$ i.e. prior prob

$$\textcircled{1} P(f / \text{test is } +ve) = P(+ve \text{ test}) / f$$

Knowing if there ~~is~~ is funny disease or not given test results
are +ve

Prior that

1% have it

$$\textcircled{2} P(+ve / \text{test is true}) = P(+ve) * P(+ve \text{ test} / +ve)$$

Is it easy to estimate joint Distribution?

e.g. → if we have d -dimensional features.

$X = \langle x_1, x_2, \dots, x_d \rangle$

each x_i takes k distinct values (of feature) from data points
and we have c distinct classes.

So how many entries in joint D Table?

$\frac{k^d}{=}$ and c classes \Rightarrow for each class k^d values are replicated. \Rightarrow $c \cdot k^d$

if $d = 100$, $k = 2$ & $c = 2 \Rightarrow 2^{100}$ i.e. Huge

So here comes Naive Bayesian Assumption →

The joint distribution of features conditioned on label becomes independent

Once they become independent, we can write joint as the product of individual ones.

$$\therefore P(x_1, x_2, x_3, \dots, x_d | y) = \prod_{i=1}^d P(x_i | y)$$

↓
Product of Prob
of each feature
conditioned on label.

i.e., Product of individual conditional probability

i.e., Now table size $\Rightarrow ?$

So, now to make table we only need $P(x_1) * P(x_2)$
 i.e. no. of values in x_1 & no. of values in x_2 .
 $\therefore \underline{kdc} = \underline{400}$. just that!

So,

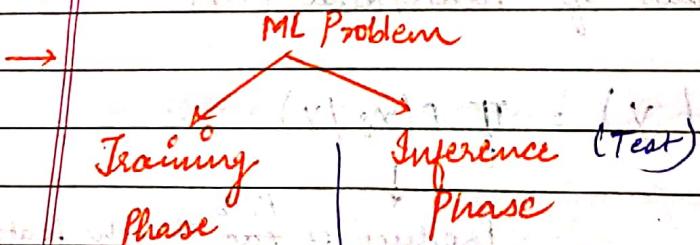
$$P(y=y_k | x_1 \dots x_n) = \frac{P(y=y_k)}{\sum_j P(y=y_j)} P(x_1 \dots x_n | y=y_k)$$

$\prod_i P(x_i^o | y=y_k)$

because → $\prod_j P(x_j^o | y=y_j)$

So you will pick up that y for which numerator is highest
 out of all possible combinations of x_i^o .

$$y^{new} = \arg \max_{y_k} P(y=y_k) \prod_i P(x_i^{new} | y=y_k)$$



In learn prob of
 occurrence of each
 the classes. & prob that
 for a given class a feature
 vector will take on a
 certain value.

You will use these values of
 features & ~~to~~ compute
 all y 's and take on y which
~~gives~~ is maximum

lot/karib?

Algo →

① First compare $P(x|y)$ and $P(y)$.

② Given new x :

$$P(y=1|x) = \frac{P(x|y=1) P(y=1)}{P(x)}$$

$$\sum_x P(x) = \sum_y P(x,y) \rightarrow \text{Marginal of } x$$

$$\text{and } P(y=0|x) = \frac{P(x|y=0) P(y=0)}{P(x)}$$

$$\text{where } P(x) = \sum_y P(x,y) = P(x|y=1) \cdot P(y=1) + P(x|y=0) \cdot P(y=0)$$

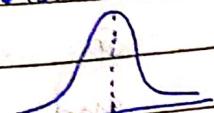
for Continuous Features → Gaussian Naive Bayes!
(this dataset)

We assume that they come from some distribution (Gaussian).

We will estimate mean & variance from that distribution.

(let say $n=100$)

Eg → $\langle x_1, \dots, x_n \rangle$ data points coming from Gaussian
Suppose → IID → Independent & Identically distributed.



Now estimate μ & σ^2 .

How will you do it?

Use fact IID to simplify joint.

Identical \rightarrow

i.e. these points are coming from a Gaussian with same μ & σ^2

Independent \rightarrow

Joint is the product of individual probabilities.

$$\text{Joint} = P(x_1, \dots, x_n | \text{class}) = P(x_1) * P(x_2) * P(x_3) * \dots * P(x_n)$$

$$P(x_i) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x_i - \mu)^2}{2\sigma^2}}$$

[pdf of Gaussian]

$$P(x_i) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x_i - \mu)^2}{2\sigma^2}}$$

- $\infty < x < \infty$ i.e. continuous.

and so on

$$\therefore P(x_1, \dots, x_n | y) = \left(\frac{1}{\sqrt{2\pi}\sigma} \right)^n e^{-\sum_{i=1}^n \frac{(x_i - \mu)^2}{2\sigma^2}}$$

Likelihood function

So this function we need to maximize why?

Because we want maximum likelihood.

So we will differentiate it w.r.t. μ .

and equate it to zero.

This is how you will get best parameters for Gaussian.

\Rightarrow we will get value of μ

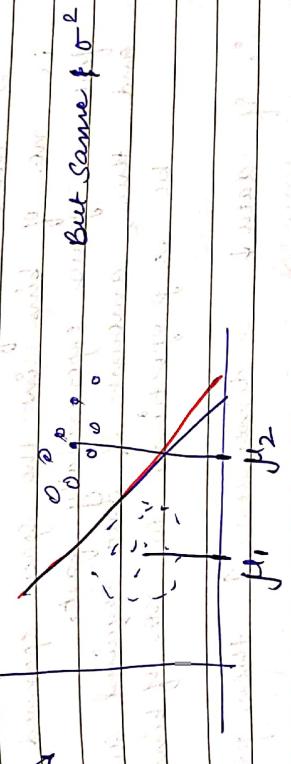
\Rightarrow same mean is the maximum likelihood.

same thing for Variance.

Continued \rightarrow
Open Questions

- ↳ Features - Continuous \rightarrow Gaussian NB
 ↳ " - Discrete \rightarrow If Multinomial NB
 ↳ " - Binary $(0|1) \rightarrow$ Bernoulli NB.

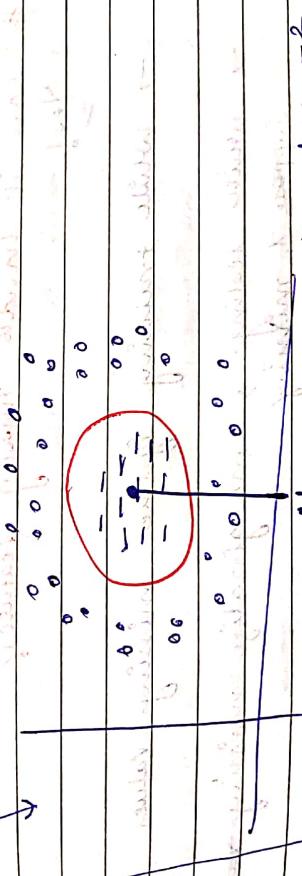
- ↳ What does a decision boundary look like?
 ↳ (All classes have \rightarrow)
 ↳ Different means & same variance \rightarrow straight line / plane
 ↳ same mean & different variance \rightarrow circle / ellipse
 ↳ General case \rightarrow parabolic curve



But same μ , σ^2

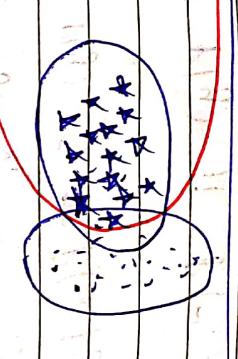
$$\mu_1$$

$$\mu_2$$



μ same, But diff σ^2

General case!



PA

What is the impact of Dimensions on the classifier?

$$d = 100 \quad \& \quad d = 1000 ?$$

for larger $n \rightarrow P(x_1, x_2, \dots, x_n | \text{class})$ takes exponential amount of time.

as dimensions \uparrow , data becomes more sparse
 \Rightarrow more no. of bins needs to be created.

i.e. no. of bins in the space spanned by features grows exponentially fast.
and thus, amount of data required should be proportional to no. of bins

\therefore if data \downarrow & dimensions \uparrow then problem!

Can it handle missing features?
Yes!

\rightarrow while training ignore missing values.

\rightarrow while testing, marginalize missing feature.
Suppose x_i missing

$$P(y | x_1, x_2, \dots, x_d) \propto P(y) P(x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_d | y)$$

$$= P(y) \underset{x_1}{\underset{\dots}{\underset{x_i}{\underset{\dots}{\underset{x_d}{\prod}}}}} P(x_1, \dots, x_d | y)$$

$$= P(y) \prod_{i=2, 3, \dots, d} P(x_i | y)$$

$$= P(y) \prod_{i=2, 3, \dots, d} P(x_i^0 | y)$$

Can it perform model updating?

Suppose you trained your model & now I give you new data. How will you update your model? Will you train it once again all over again?

$$\text{eg} \rightarrow f_{w,i} = \frac{\text{frequency of word}}{\text{document length}}$$

Let say your model categorize it & make 3 bins \rightarrow

$$0 \leq f_{w,i} < 0.001 \quad \text{very rare}$$

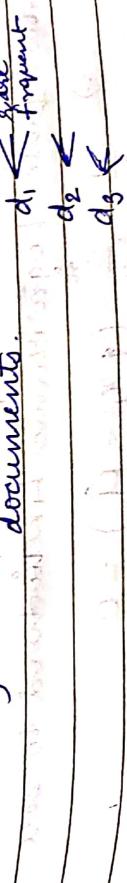
$$0.001 \leq f_{w,i} < 0.01 \quad \text{rare}$$

$$0.01 \leq f_{w,i} <= 1 \quad \text{frequent}$$

This feature value for word w must be within these 3 intervals.

In new example comes you can update your parameters, by \rightarrow

maintaining \rightarrow ① counts over three bins for all documents.



$$\text{Prob of } w = \frac{\text{this count}}{\sum_{\text{bin}} \text{Sum of counts}} \quad \text{in subsequent document}$$

Continued calc. of. $\hat{\mu}_1$ & $\hat{\sigma}^2$

Finding $\hat{\mu}$

$$\text{likelihood } L = \left(\frac{1}{\sqrt{2\pi}\sigma} \right)^n e^{-\frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2}$$

where we have exponents. Then we take log
 \therefore log likelihood \rightarrow

$$\log L = n \log \left(\frac{1}{\sqrt{2\pi}\sigma} \right) + \underbrace{-\frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2}_{\text{const.}}$$

$$\frac{d}{d\mu} \left(\log L \right) = 0 + \frac{1}{2\sigma^2} \frac{d}{d\mu} \sum_{i=1}^n (x_i - \mu)^2$$

$$0 = \frac{1}{2\sigma^2} \frac{d}{d\mu} \sum_{i=1}^n (x_i^2 - 2x_i\mu + \mu^2)$$

$$\Rightarrow -\frac{1}{2\sigma^2} \frac{d}{d\mu} \sum_{i=1}^n (x_i^2 - 2x_i\mu + \mu^2)$$

$$\Rightarrow -\frac{1}{2\sigma^2} \cdot 2 \cdot \sum_{i=1}^n (x_i^2 - 2x_i\mu + \mu^2)$$

Now equate logarithmic likelihood to zero.

$$\Rightarrow \frac{1}{\sigma^2} \sum_{i=1}^n (x_i^2 - 2x_i\mu + \mu^2) = 0$$

$$\Rightarrow \sum_{i=1}^n x_i^2 - n\mu^2 = \frac{\sum_{i=1}^n x_i^2}{n}$$

$$\frac{\partial^2}{\partial \mu^2} \ln L = -\frac{1}{\sigma^2}$$

$$\frac{d}{dx} \log f(x) =$$

classmate
Date _____
Page _____

take second order derivative \rightarrow

$$\Rightarrow -\frac{n}{\sigma^2}$$

which is -ve

which shows it is Maximum Likelihood estimation.

Finding σ

Differentiate wrt σ^2

$$\frac{d(\log L)}{d\sigma^2} = \frac{d}{d\sigma} \left(\left(\frac{n}{\sqrt{2\pi}\sigma} \right)^n e^{-\frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2} \right)$$

$$\frac{d(\log L)}{d\sigma^2} = \frac{-n}{\sigma^2} \frac{d}{d\sigma} \log \left(\frac{1}{\sqrt{2\pi}\sigma} \right) + \frac{-1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2$$

$$= -\frac{n}{\sigma^2} \frac{1}{\sqrt{2\pi}\sigma} + \frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2$$

$$= +\frac{n}{2\sigma^2} \frac{d \log}{d\sigma^2}$$

$$= -\frac{n}{2\sigma^2} \frac{d}{d\sigma^2} \log \left(\frac{1}{\sqrt{2\pi}\sigma} \right) + \frac{-1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2$$

$$= -\frac{n}{2} \frac{1}{2\pi\sigma^2} \frac{1}{\sigma^2} + \frac{1}{2\sigma^4}$$

$$= -\frac{n}{2\sigma^2}$$

$$\begin{aligned}
 &= \frac{d}{d\sigma^2} \left(-\frac{n}{2} \log(\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2 \right) \\
 &= \frac{\partial}{\partial \sigma} \left(-\frac{n}{2} \log(2\pi\sigma) - \frac{1}{2\sigma} \sum_{i=1}^n (x_i - \mu)^2 \right) \\
 &= -\frac{n}{2} \times \frac{-1}{2\sigma} = \frac{n}{2\sigma} \left(\sum_{i=1}^n (x_i - \mu)^2 \right) \\
 &= -\frac{n}{2\sigma} + \frac{1}{2} \left(\frac{\partial^2}{\partial \sigma^2} \left(\sum_{i=1}^n (x_i - \mu)^2 \right) \right) \\
 &= -\frac{n}{2\sigma} + \frac{1}{2} \left(\frac{\partial^2}{\partial \sigma^2} \left(\sum_{i=1}^n (x_i - \mu)^2 \right) \right) \\
 &\Rightarrow \sigma^2 = \frac{n}{2\sigma} + \frac{1}{2} \sum_{i=1}^n (x_i - \mu)^2 = 0 \\
 &\sigma^2 = \frac{2}{2\sigma^4} \sum_{i=1}^n (x_i - \mu)^2 = \sigma^2
 \end{aligned}$$

So best estimate is variance!

- # Eg of Time Series data →
 - ① → Audio, speech
 - ② → Weather forecasting or 10 days history of weather.
- We make set of states.
- Data point is a time-ordered set of states
 one day can take one state out of 3
- $\{K \rightarrow PC, CR\}$
- Clear partly cloudy Rain

classmate
Date _____
Page _____

Sequence of States \rightarrow Ordering [for each day]

KK, PC, CR, PC, PC, PC, KK, KK [we cannot change its order]

① Radio Broadcast \rightarrow

What you want from you

what is state? (freq? No) which lang?

= $\Delta\alpha$

State Space = can be very large

eg \rightarrow teen talk, sports talk, movie talk etc.

② Text Data \rightarrow

State Space = All possible sentences of a language

③ Word data \rightarrow

State Space = Sylabols. $\rightarrow \{ \text{aa}, \text{aa}, \text{aa}, \text{aa}, \text{aa}, \text{aa} \}$

note

Random Variable → function from sample space to \mathbb{R} [Real values]

Random Vector → funct' from Sample Space to \mathbb{R}^n [n-dim vec]

Random process → A random process is a collection of indexed random variables defined over a probability space.

{Triplet of} → Sample Space → P.d.f.
Event Space

Where index is time in most cases (time series)

Markov chain → It is a special type of random process where probability of the present state depends only on the previous state.

It must be "memory-less".

$$P(\text{Current state} / \text{All history}) = P(\text{Current s} / \text{previous})$$

↓
(same only)

Hidden Markov Model → It is a markov model which models a Markov chain with unobservable (hidden) states.

classmate
Date _____
Page _____

ple space

ce to R^N
dim rect)

a
defined

in series)

random state

reviews)

e model
observable

classmate
Date _____
Page _____

H. Markov Model \rightarrow

Set of states must be finite and discrete

$$S = \{s_1, s_2, \dots, s_N\}$$

when you go from one state to other, it will generate sequence of states.

$$s_1, s_2, s_3, \dots, s_k$$

Markovian Property \rightarrow

$$P(s_k | s_1, s_2, \dots, s_{k-1}) = P(s_k | s_{k-1})$$

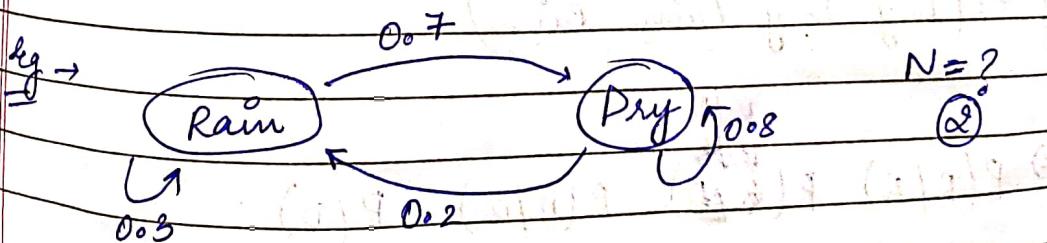
① Transition Probability \rightarrow $a_{ij} = P(s_j | s_i)$

It's a matrix of size $N \times N$.

If entry = 0, we cannot directly reach there.

② Initial State / Vector / Probabilities \rightarrow

$$\pi_i^0 = P(s_i^0)$$



		Rain	Dry	Transition Matrix
Rain	0.3	0.7		
Dry	0.2	0.8		

But we should know initial Probabilities →

$$P(\text{Rain}) = 0.4 \quad P(\text{Dry}) = 0.6$$

→ Joint Probability using Markovian Property →

$$P(S_{i_1}, S_{i_2}, \dots, S_{i_k}) = P(S_{i_k} | S_{i_1}, S_{i_2}, S_{i_3}, \dots, S_{i_{k-1}}) P(S_{i_1}, \dots, S_{i_{k-1}})$$

$$P(A \cap B) = P(A|B) * P(B) \quad [\text{conditional probability}]$$

Now by using joint Markovian Property

$$P(S_{i_1}, S_{i_2}, \dots, S_{i_k}) = P(S_{i_k} | S_{i_1}, S_{i_2}, \dots, S_{i_{k-1}}) \cdot P(S_{i_1}, \dots, S_{i_{k-1}})$$

$$= P(S_{i_k} | S_{i_{k-1}}) \cdot P(S_{i_{k-1}} | S_{i_1}, \dots, S_{i_{k-2}}) \cdot P(S_{i_1}, \dots, S_{i_{k-2}})$$

Probability of Occurrence of a Sequence = $P(S_{i_k} | S_{i_{k-1}}) \cdot P(S_{i_{k-1}} | S_{i_{k-2}}) \cdot P(S_{i_{k-2}} | S_{i_{k-3}}) \cdots P(S_{i_2} | S_{i_1}) P(S_{i_1})$

e.g. → $P(\text{Dry Dry Rain Rain})$

$$\Rightarrow P(R|R) \cdot P(D|R) \cdot P(D|D) \cdot P(D)$$

$$\Rightarrow 0.3 \times 0.7 \times 0.8 \times 0.6$$

Nothing hidden so far.

Hidden Markov Model

States are not visible, but each state randomly generates one of M observations (or visible state). True states are not known. $\{v_1, v_2, \dots, v_N\}$ (I DK about temp but I can see you wearing sweater or not)

Eg → What gene is driving disease but we can't see symptoms
Tent analysis S ~~ass~~ → 5 or S

Connectⁿ b/w visible states & sequence of states?

Visible & hidden states both are entirely different!

If I'm in state s_1 (hidden) what is the probability of going to $v_{21}, v_{22} \dots$ etc.

Eg → If true thing was 'A' what is prob of writing 'S'?
Very low

Very high. \uparrow S_{ij}^{oo} \uparrow V_{ij}^{oo} .

① Matrix of Transition Probabilities →

$$A = (a_{ij}^{oo}), \quad a_{ij}^{oo} = P(S_j^o | S_i^o)$$

Going from i to j

② Matrix of Observation Probabilities →

$$B = (b_i(v_m)), \quad b_i(v_m) = P(v_m | s_i)$$

Prob of getting observation v_m when you are in state s_i .

Size $\rightarrow (N \times M)$ [$N \rightarrow$ true states (Hidden)
 $M \rightarrow$ visible observations]

③ Vector of Initial Probabilities →

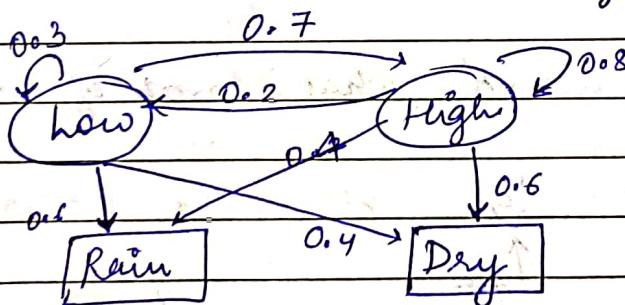
$$\pi = (\pi_i), \quad \pi_i = P(s_i) \quad [\text{Prob of starting from } s_i \text{ state}]$$

Model → $M = (A, B, \pi)$ [Representation]

leg →

Two Observations : Rain & Dry

Hidden : Low & High Atmospheric pressure



Transitions are not b/w visible states!

$A \rightarrow$		L	H
	L	0.3	0.7
	H	0.2	0.8

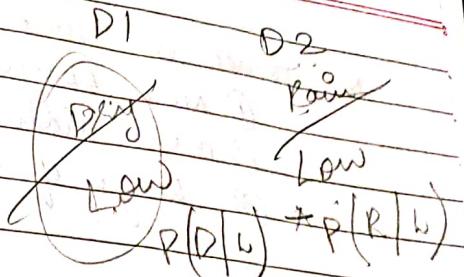
$B \rightarrow$		Rain	Dry
	L	0.6	0.4
	H	0.4	0.6

$$\pi \Rightarrow P(\text{'Low'}) = 0.4 \\ P(\text{'High'}) = 0.6$$

Q → Sequence → Dry, Rain

$$P(D, R) = P(D) P(R)$$

Prob.



$$P(D, R) = P(D|R) P(\{D, R\}, \{L, H\}) +$$

$$P(\{D, R\}, \{L, H\}) +$$

$$P(\{D, R\}, \{H, L\}) +$$

$$P(\{D, R\}, \{H, L\})$$

$$\begin{aligned} P(\{D, R\}, \{L, H\}) &= P(L|H) * P(\{D, R\} | \{L, H\}) \\ &= P(how) P(how | how) * P(D|L) * P(R|L) \\ &= 0.4 * 0.3 * 0.4 * 0.6 \end{aligned}$$

Do same for other 3 terms too.

Q → Computational complexity of $P(\text{Observed Sequence})$ depends on?

If Sequence of states are T in length
and ' c ' are the no. of hidden states
then Complexity = $O(c^2 T)$

~~These will be the no. of~~

Q-1 What is more tolerable?

- ① More Observable States
- ② more hidden states
- ③ length of Observed State

Initial state found in family = (α, α)

$\rightarrow (\text{child}, \text{child}) \alpha$

$\rightarrow (\text{child}, \text{adult}) \alpha$

$\rightarrow (\text{adult}, \text{adult}) \alpha$

$\rightarrow (\text{child}, \text{adult}) \alpha * (\text{child}, \alpha) = (\text{child}, \text{adult}) \alpha$

$\rightarrow (\text{child}, \text{adult}) \alpha * (\text{adult}, \text{adult}) \alpha = (\text{adult}, \text{adult}) \alpha$

Learning the HMM ?

We will need to learn A , B and π from data

Modeling limitations \rightarrow first we need to fix no. of hidden states, Observed states etc. then find A , B , π

[Eg \rightarrow in linear regression we fix model $y = \theta_0 + \theta_1 x + \theta_2 x^2$ then we find θ_0 , θ_1 & θ_2 .]

Maximize $P(O|M)$

\rightarrow you'll take initial model
and then ~~your~~ model will give size
to observation & you maximize it.

No closed form soln available

classmate

Date _____

Page _____

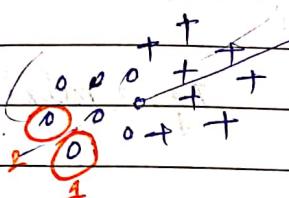
Baum-Welch Algo \rightarrow Iterative method which gives you local maximum (& not global)
It is based on EM Algo.
 \downarrow
Expectation Maximization.

EM Algo \rightarrow Given set of points knowing they are coming from 2 or more Gaussians.
Now you need to figure out parameters of those Gaussians.

EM Algo \rightarrow k-means

Given GMM \rightarrow Gaussian Mixture Model

(Two Gaussians)

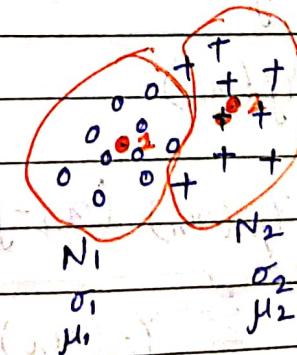


1) Select 2 representatives of 2 classes.

2) K-Means \rightarrow Representatives keeps on moving & getting refined.

At convergence,

- one centroid.



Each of the representatives will be assigned to one

3) We are trying to estimate parameters through EM.
Now you got N_1 points of class 1 & N_2 points of class 2.

Now find μ & σ^2 . that's it!

Through logarithmic likelihood.

E step

- 4) Now, evaluate responsibilities using the current parameters values.
Assuming that these are parameters you'll compute value is the prob that the point was generated by the Gaussian.
- But generated by π_k Σ_k μ_k
Probability: If point generated by π_k then Σ_k gaussian then high
 π_k else low weightage.

$$\gamma_j^o(x) = \pi_k N(x | (\mu_k, \Sigma_k))$$

$$\sum_{j=1}^K \pi_j^o N(x | \mu_j, \Sigma_j)$$

initial → mixed model
 $\mu_j^o \rightarrow$ Mean $\Sigma_j^o \rightarrow$ Covariance $\pi_j^o \rightarrow$ Mixing coefficient

- 5.) ~~M step~~ Re-estimate parameters using current responsibilities.

$$\rightarrow \mu_j^o = \frac{\sum_{n=1}^N \gamma_n^o(x_n) x_n}{\sum_{n=1}^N \gamma_n^o(x_n)}$$

$$\rightarrow \Sigma_j^o = \sum_{n=1}^N \gamma_n^o(x_n) (x_n - \mu_j^o)(x_n - \mu_j^o)^T$$

$$\sum_{n=1}^N \gamma_n^o(x_n)$$

$$\rightarrow \pi_j^o = \frac{1}{N} \sum_{n=1}^N \gamma_n^o(x_n)$$

6.) evaluate log likelihood

$$\ln p(x | \mu, \Sigma, \pi) = \sum_{n=1}^N \ln \left\{ \sum_{k=1}^K \pi_k N(x_n | \mu_k, \Sigma_k) \right\}$$

If no convergence, return to step 2.