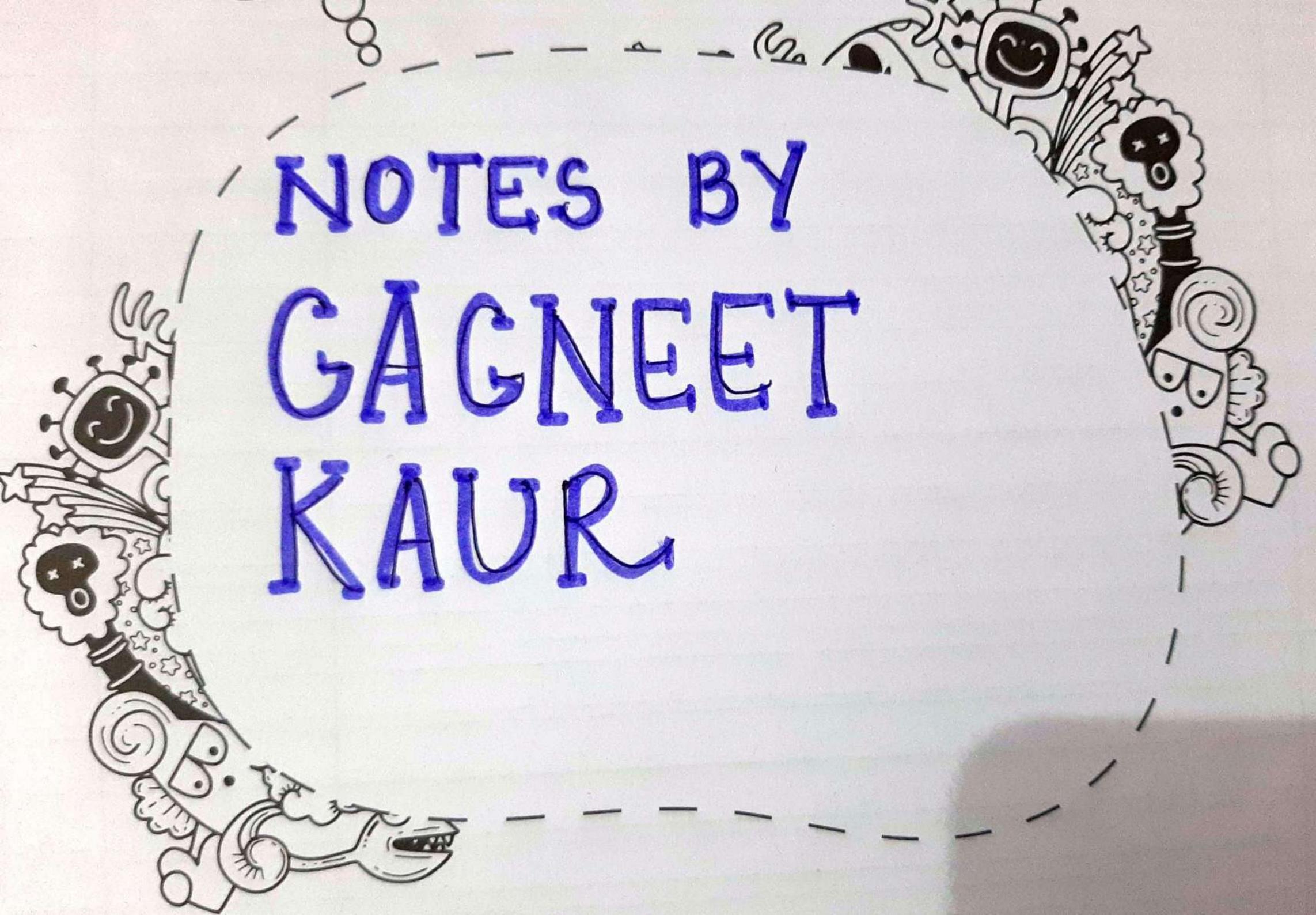


NOTES BY GAGNEET KAUR



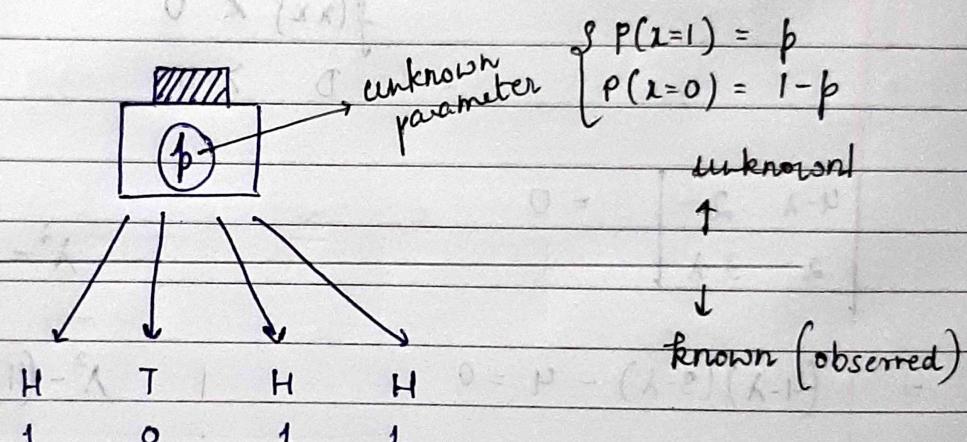
Week 4 MLT

Unsupervised Learning - ESTIMATION

Lecture 1 → Introduction to Estimation

"There is some probabilistic mechanism that generates data"
- about which we don't know something!!

Goal is : Given data find/estimate what we don't know.



Steps of Estimation Problem : (General Idea Of Estimation)

- Observe data
- 'Assume' a probabilistic model that generates data
- Estimate unknown parameters using data.

ASSUMPTIONS

Data $\rightarrow \{x_1, x_2, \dots, x_n\}$

Observations are :

- independent
- identically distributed

INDEPENDENCE : $P(x_i/x_j) = p(x_i) \quad \forall i \neq j$

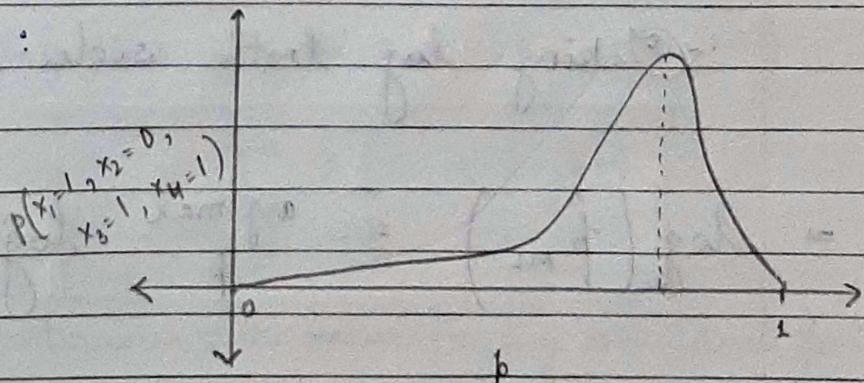
IDENTICAL DIST. : $p(x_i=1) = p(x_j=1) = p \quad \forall i, j$

lecture 2 → Maximum Likelihood Estimation

The principled way to get estimators from data - LIKELIHOOD FUNCTION

Likelihood Function :

Eg: $\{1, 0, 1, 1\}$



For every value of P , we ask the question :

I have seen this data. If the true p was this, then what is the chance that I would have seen the state?

- CURVE PEAK \rightarrow the value of p for which the chance of seeing the data is highest.
- \rightarrow the data is most likely for which parameter of P , which choice of p .

Fischer's Principle of Maximum Likelihood

$$\begin{aligned}
 L(p, \{x_1, x_2, \dots, x_n\}) &= P(x_1, x_2, \dots, x_n; p) \\
 &= P(x_1; p) \cdot P(x_2; p) \cdots \cdot P(x_n; p) \\
 &= \prod_{i=1}^n p^{x_i} (1-p)^{(1-x_i)}
 \end{aligned}$$

↴ underlying parameter

↴ parameter

↴ if $\begin{cases} x_i = 1 \Rightarrow p^1(1-p)^0 = p \\ x_i = 0 \Rightarrow p^0(1-p)^1 = 1-p \end{cases}$

$$\text{ESTIMATOR } \hat{p}_{ML} = \arg \max_p \prod_{i=1}^n p^{x_i} (1-p)^{(1-x_i)}$$

Taking log both sides,

$$\log(\hat{p}_{ML}) = \arg \max_p \log \left(\prod_{i=1}^n p^{x_i} (1-p)^{(1-x_i)} \right)$$

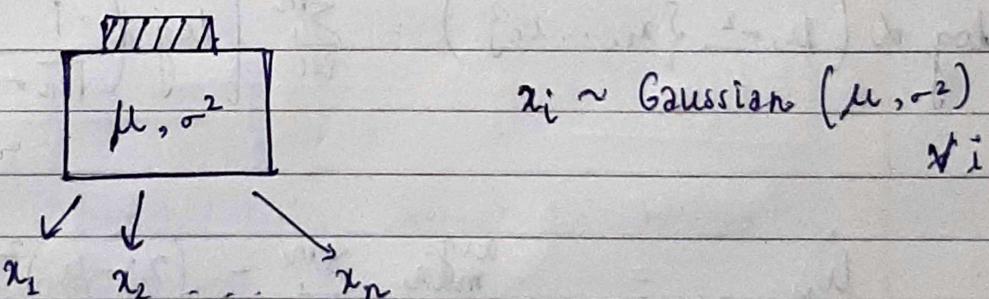
$$= \arg \max_p \sum_{i=1}^n [x_i \log p + (1-x_i) \log (1-p)]$$

Take derivative of $\log L(p)$, set it to 0 to get \hat{p}_{ML} .

$$\hat{p}_{ML} = \frac{1}{n} \sum_{i=1}^n x_i$$

- Guess for the true 'p' by looking at the data is that value of 'p' which maximizes the chance that we observe this data.

~~Different form of Data~~ Data = $\{x_1, x_2, \dots, x_n\}$ $x_i \in \mathbb{R}$ $\forall i$



$\mu \rightarrow$ unknown ; $\sigma^2 \rightarrow$ known

$$\begin{aligned} L(\mu, \sigma^2, \{x_1, \dots, x_n\}) &= P(x_1, x_2, \dots, x_n ; \mu, \sigma^2) \\ &= \prod_{i=1}^n P(x_i ; \mu, \sigma^2) \end{aligned}$$

→ because Gaussian is continuous, the value for any individual point is going to be ZERO.

→ Only intervals will have non-zero values for continuous distributions.

Thus the above fn will give us the value ZERO for any mean.

So, Fisher's proposal was to not use the probabilities

Instead, replace the probabilities and define the likelihood of the parameter as :

$$\mathcal{L}(\mu, \sigma^2, \{x_1, \dots, x_n\}) = f_{x_1, \dots, x_n}(x_1, \dots, x_n; \mu, \sigma^2)$$

$$= \prod_{i=1}^n f_{x_i}(x_i; \mu, \sigma^2)$$

$$= \prod_{i=1}^n \left[\frac{1}{\sqrt{2\pi}\sigma} e^{-(x_i - \mu)^2 / 2\sigma^2} \right]$$

$$\log \mathcal{L}(\mu, \sigma^2, \{x_1, \dots, x_n\}) = \sum_{i=1}^n \left[\log \left(\frac{1}{\sqrt{2\pi}\sigma} \right) - \frac{(x_i - \mu)^2}{2\sigma^2} \right]$$

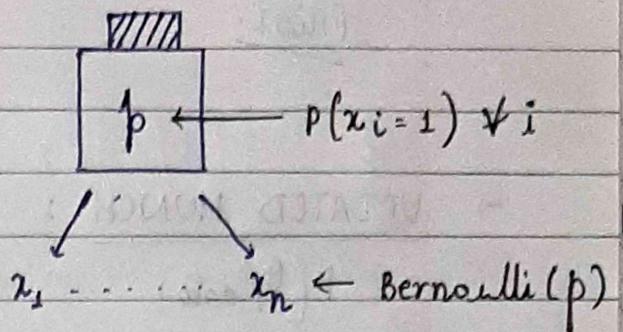
$$\hat{\mu}_{ML} = \arg \max_{\mu} \sum_{i=1}^n - (x_i - \mu)^2$$

Take derivative & set it to zero !!

$$\hat{\mu}_{ML} = \frac{1}{n} \sum_{i=1}^n x_i$$

Lecture 3 : Bayesian Estimation

Consider the coin example :



Let's say someone says,

'I believe the bias p is somewhere close to 1'

We may have 'HUNCH' about parameters.

- Is there a way we can somehow codify our hunch into mathematically more precise mechanisms that can be incorporated into our estimation procedure?

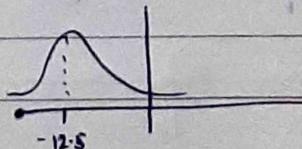
→ BAYESIAN MODELLING APPROACH

GOAL : Incorporate "hunch/belief" about parameters of interest into the estimation procedure.

APPROACH : Think of the parameter to estimate as a 'random' variable.

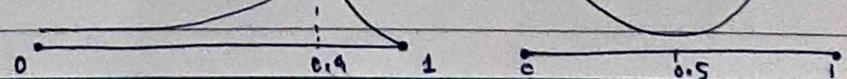
Earlier

μ



p

Now



- HUNCH : Codified using a probability distribution over Θ
 $P(\Theta)$
Prior
 \Downarrow DATA
- UPDATED HUNCH : Codified using prob. distribution
 $P(\Theta/\text{data})$
Posterior

BAYES LAW : $P(A/B) = \frac{P(B/A) \cdot P(A)}{P(B)}$

A → Parameters Θ

B → Data $\{x_1, \dots, x_n\}$

$$P(\Theta / \{x_1, \dots, x_n\}) = \frac{\left(P(\{x_1, \dots, x_n\} / \Theta) \right) \cdot P(\Theta)}{P(\{x_1, \dots, x_n\})}$$

Posterior

likelihood
 Prior
 does not depend on Θ
EVIDENCE

What is Bayes Theorem : You have some initial belief about your prior distribution. To go to a posterior distribution, you have to review your prior dist. You have to make a multiplicative update to this prior.

DATA — Bernoulli (p)

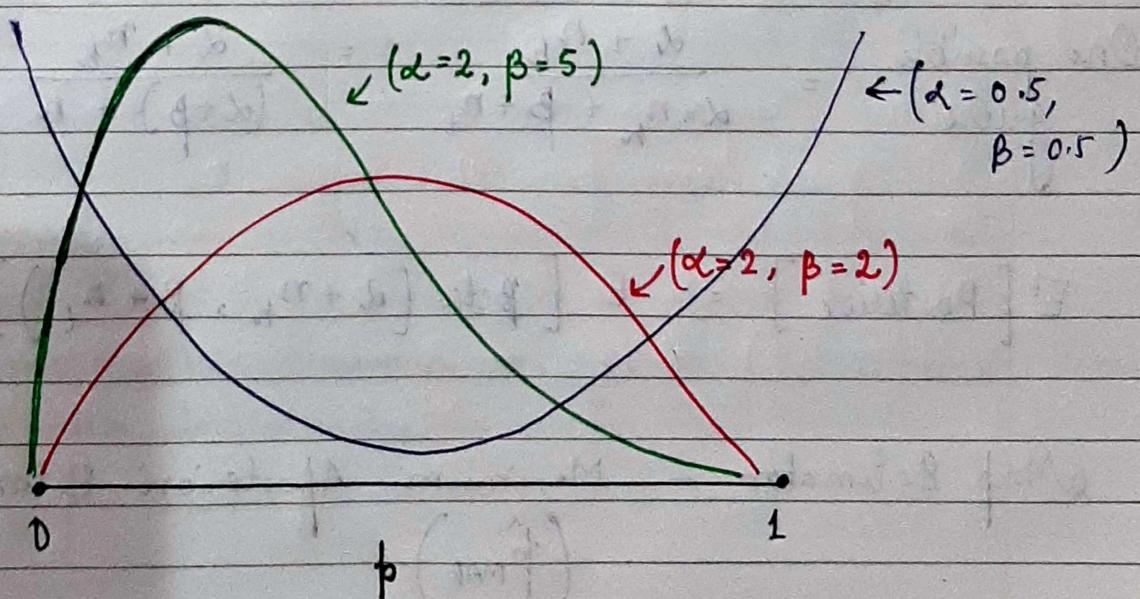
Prior? $P(\theta)$

Beta Prior $\rightarrow f(p; \alpha, \beta)$

- We want some continuous distribution whose values are between 0 & 1.
- We can't use Gaussian here because it can give me any value between $-\infty$ and ∞ , but then the parameter I am trying to estimate (P) takes values b/w 0 & 1.
- So any distribution that we use as a prior to encode our prior knowledge should be supported only in 0 & 1.
- A good choice — BETA PRIOR.

BETA PRIOR

$$f(p; \alpha, \beta) = \begin{cases} \frac{p^{\alpha-1} (1-p)^{\beta-1}}{Z} & \forall p \in [0, 1] \\ 0 & \text{otherwise} \end{cases}$$



It kind of says that when you have small α & big β , then true p is somewhere closer to 0 than 1.

$$P(\theta/\text{data}) \propto P(\text{data}/\theta) \cdot P(\theta)$$

$$f_{\theta/\text{data}}(p) \propto \left(\prod_{i=1}^n \left[p^{x_i} (1-p)^{1-x_i} \right] \right) \cdot \left[p^{\alpha-1} (1-p)^{\beta-1} \right]$$

Likelihood Prior

$$f_{\theta/\text{data}}(p) \propto p^{\sum x_i + \alpha - 1} (1-p)^{\sum (1-x_i) + \beta - 1}$$

Same functional form as PRIOR !!

$$\text{Beta Prior } (\alpha, \beta) \xrightarrow[\text{Bernoulli}]{\text{DATA}} \text{Beta posterior } \left(\alpha + n_h, \beta + n_t \right)$$

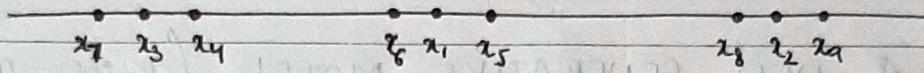
$$\text{One possible guess} = \frac{\alpha + n_h}{\alpha + n_h + \beta + n_t} = \frac{\alpha + n_h}{(\alpha + \beta) + n}$$

$$E[\text{Posterior}] = E[\text{Beta}(\alpha + n_h, \beta + n_t)] = \frac{\alpha + n_h}{(\alpha + \beta) + n}$$

MAP Estimator - Maximum A posteriori Estimator
 (\hat{p}_{MAP})

Lecture 4 : Gaussian Mixture Models

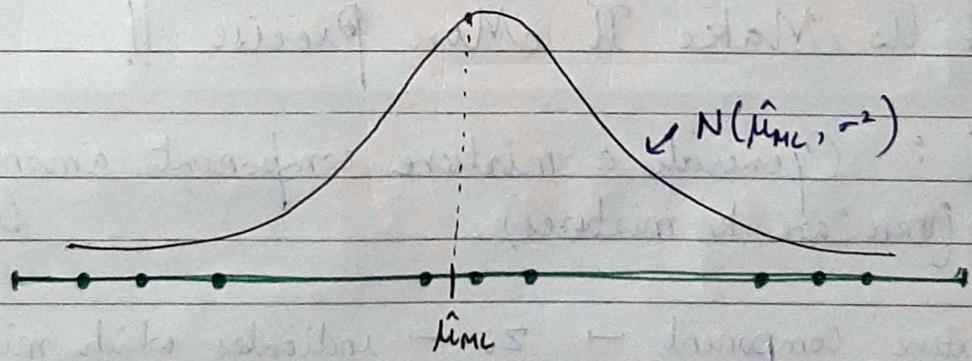
GOAL : Estimation for Slightly complicated data



→ Goal is to come up with a model that explains the data.

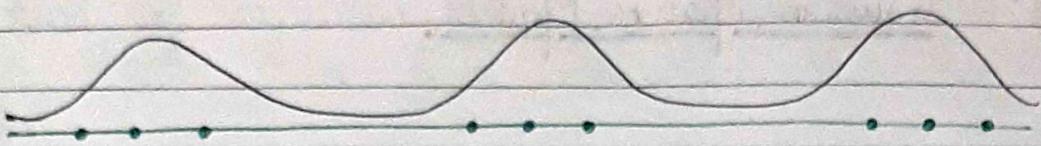
* What could be a good generative story?

- Try measuring or modelling this using a Gaussian distribution.
- Let us say we did a principle of maximum likelihood to get the best Gaussian



- The sample mean is the maximum likelihood estimator.

* Our explanation via Gaussian is not very satisfactory because there are two low density regions for the Gaussian.



- We want something, a density like above to explain this data.
- Basically, we want three different models or three places where the density should peak.

A NEW GENERATIVE MODEL (mechanism that generates the data)

The name of this new model that we are going to come up is called as mixtures of Gaussians.

STEP 1: Pick which mixture a data point comes from.

STEP 2: Generate data points from that mixture.

① Let us Make It More Precise !!

STEP I : Generate a mixture component among $\{1, \dots, k\}$
(There are k mixtures)

Mixture Component $\rightarrow z_i \rightarrow$ indicates which mixture the i^{th} data point comes from that value z_i

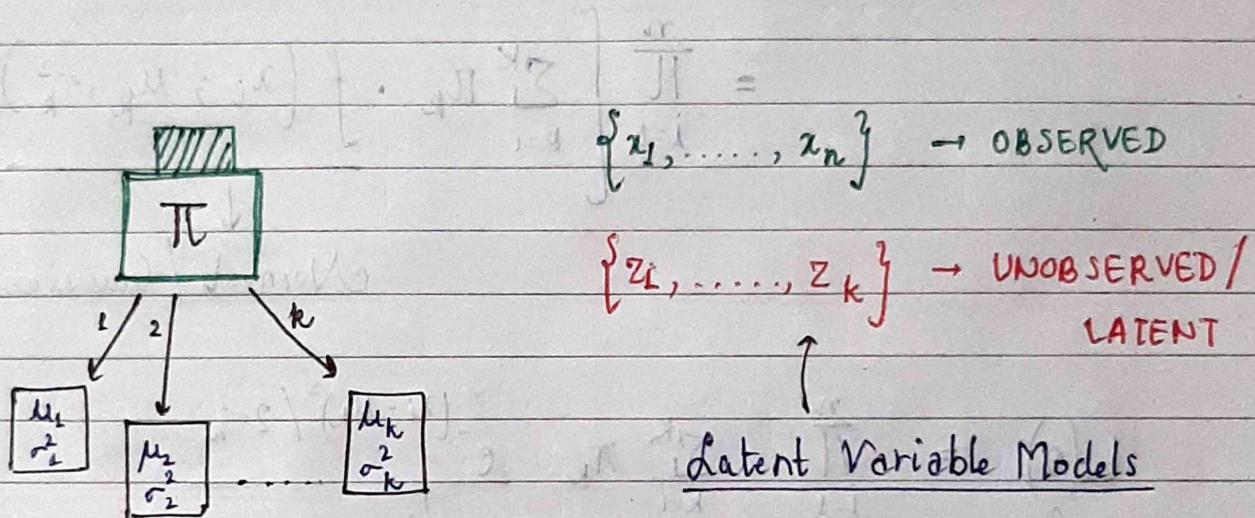
Let us formulate it as Probability Distribution :

$P(z_i = l) = \pi_l$ \leftarrow It means that the probability that the i^{th} data point comes from l^{th} mixture is given by π_l . $\left(\sum_{i=1}^k \pi_i = 1 \right)$

So now we have decided which cluster or mixture the data has to come from !!

STEP II : Generate $x_i \sim N(\mu_{z_i}, \sigma^2_{z_i})$

for each mixture we have a Gaussian with its own mean & variance.



Our Gaussian Mixture Model is a latent variable model, because the final output that you are observing depends not only on some parameters that you want to estimate, but also on some unobserved latent variables.

HOW MANY PARAMETERS WILL DETERMINE THIS MODEL COMPLETELY ?

$$\pi = [\pi_1 \ \pi_2 \ \dots \ \pi_k]$$

for each K , there are 2 parameters $\underbrace{\pi_k}_{(\pi)} \rightarrow \mu_k$ & σ^2_k

$$\text{In Total : } 2K + k - 1 = 3K$$

We need to estimate $(3k - 1)$ parameters from data.

Lecture 5 : Likelihood of GMM

$$L \left(\begin{matrix} \mu_1, \dots, \mu_k \\ \sigma_1^2, \dots, \sigma_k^2 \\ \pi_1, \dots, \pi_k \end{matrix} ; \begin{matrix} x_1, \dots, x_n \end{matrix} \right) = \prod_{i=1}^n f_{\text{mix}} \left(x_i ; \begin{matrix} \mu_1, \dots, \mu_k \\ \sigma_1^2, \dots, \sigma_k^2 \\ \pi_1, \dots, \pi_k \end{matrix} \right)$$

$$= \prod_{i=1}^n \left[\sum_{k=1}^K \pi_k \cdot f(x_i; \mu_k, \sigma_k^2) \right]$$

↓
Normal / Gaussian Density.

$$L(\theta) = \prod_{i=1}^n \sum_{k=1}^K \pi_k \frac{e^{-(x_i - \mu_k)^2 / 2\sigma_k^2}}{\sqrt{2\pi} \sigma_k}$$

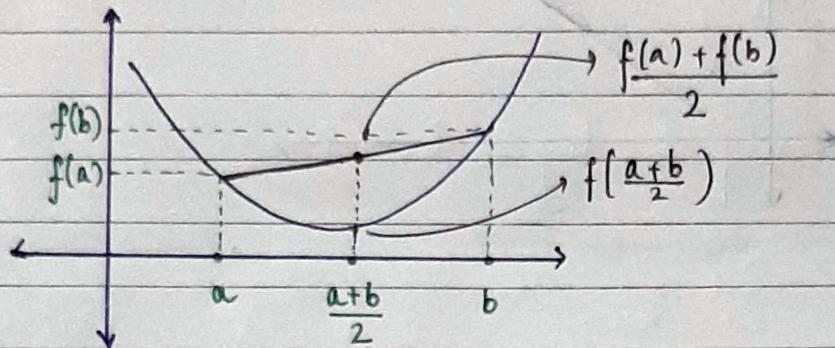
↑
all parameters

$$\log L(\theta) = \sum_{i=1}^n \log \left(\sum_{k=1}^K \pi_k \frac{e^{-(x_i - \mu_k)^2 / 2\sigma_k^2}}{\sqrt{2\pi} \sigma_k} \right)$$

- Not possible to solve the above eqn Analytically.
- Need an alternate way to solve this efficiently!

Lecture 6 : Convex Functions & Jensen's Equality

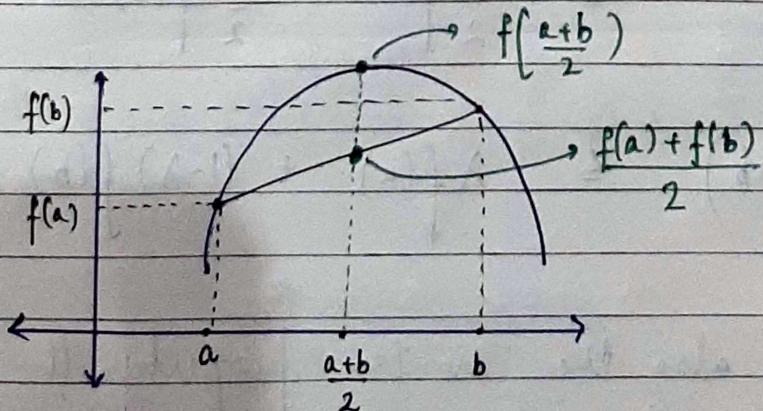
- What is a Convex Function ?



If the linear interpolation at two points has a strictly higher value than the function itself \rightarrow such a function is called CONVEX FUNCTION.

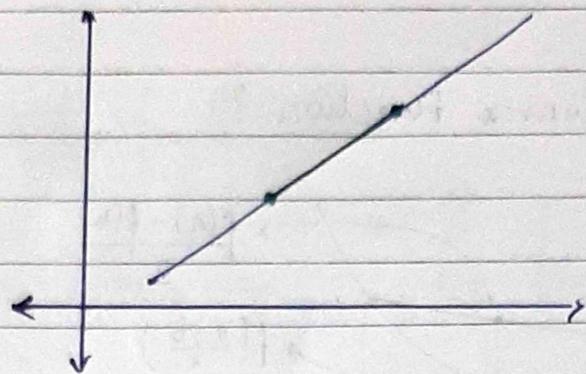
$$f\left(\frac{a+b}{2}\right) \leq \frac{f(a) + f(b)}{2} \quad \forall a, b$$

- What is CONCAVE function ?



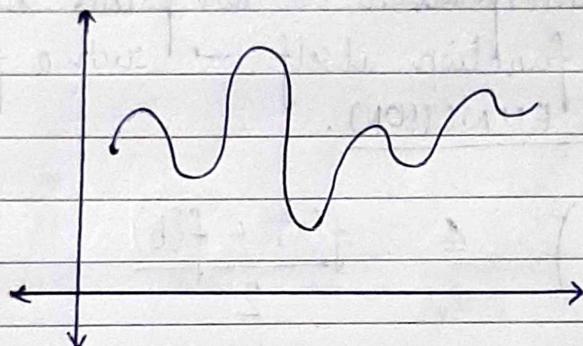
$$\forall a, b \quad f\left(\frac{a+b}{2}\right) \geq \frac{f(a) + f(b)}{2}$$

- If we have a linear fn \rightarrow Both Concave & Convex



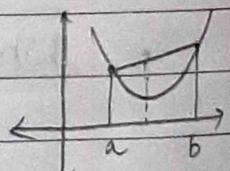
$$f\left(\frac{a+b}{2}\right) = \frac{f(a) + f(b)}{2}$$

- Neither concave nor convex :



Convex :

$$f\left(\frac{1}{2}a + \frac{1}{2}b\right) \leq \frac{1}{2}f(a) + \frac{1}{2}f(b)$$



$$\Rightarrow f(\lambda a + (1-\lambda)b) \leq \lambda f(a) + (1-\lambda)f(b) \quad \lambda \in [0,1]$$

Concavity has also the similar property !!

- The line segment joining a and b . You can generally use any lambda λ as you vary lambda from a to b .

~~Generalization~~

General: Now you can extend it to multiple points :

$$f(\lambda_1 a_1 + \lambda_2 a_2 + \dots + \lambda_k a_k) \geq \lambda_1 f(a_1) + \dots + \lambda_k f(a_k)$$

$$f\left(\sum_{k=1}^K \lambda_k a_k\right) \geq \sum_{k=1}^K \lambda_k f(a_k)$$

JENSEN'S
INEQUALITY

$$\left(\sum_{k=1}^K \lambda_k = 1\right) \quad (0 \leq \lambda_k \leq 1, 1 \leq k \leq K)$$

- Log is a concave function !
- How can we exploit Jensen's for performing maximum likelihood ?

Lecture 7 : Estimating The Parameters

Recall,

$$\log L(\theta) = \sum_{i=1}^n \log \left(\sum_{k=1}^K \left(\pi_k \frac{e^{-(x_i - \mu_k)^2 / 2\sigma_k^2}}{\sqrt{2\pi} \sigma_k} \right) \right)$$

- Introduce - for every data point i , the parameters

$$\{\lambda_1^i, \dots, \lambda_K^i\} \text{ such that } \forall i \sum_{k=1}^K \lambda_k^i = 1, 0 \leq \lambda_k^i \leq 1$$

$\forall i, k$

$$\log L(\theta) = \sum_{i=1}^n \log \left(\frac{\sum_{k=1}^K \lambda_k^i \left(\frac{\pi_k e^{-(x_i - \mu_k)^2 / 2\sigma_k^2}}{\sqrt{2\pi} \sigma_k} \right)}{\lambda_k^i} \right)$$

By Jensen's Inequality :

$$\log L(\theta) \geq \text{modified-log } L(\theta, \lambda)$$

$$= \sum_{i=1}^n \sum_{k=1}^K \lambda_k^i \log \left(\frac{\pi_k e^{-(x_i - \mu_k)^2 / 2\sigma_k^2}}{\lambda_k^i \cdot \sqrt{2\pi} \sigma_k} \right)$$

- Note that the above modified log likelihood gives a lower bound for the true log likelihood at θ for any choice of λ

$$\left\{ \lambda_1^i, \dots, \lambda_K^i \right\} \quad \left\{ \mu_1, \dots, \mu_K \right. \\ \vdots \quad \left. \sigma_1^2, \dots, \sigma_K^2 \right. \\ \left\{ \pi_1, \dots, \pi_K \right\}$$

KEY INSIGHT :

- If we fix λ , it is easy to maximize wrt θ
- If we fix θ , it is easy to maximize wrt λ .

(I) FIX λ and MAXIMIZE OVER θ

$$\max_{\theta} \sum_{i=1}^n \sum_{k=1}^K \lambda_k^i \left[\log \left(\pi_k \frac{e^{-(x_i - \mu_k)^2 / 2\sigma_k^2}}{\sqrt{2\pi} \sigma_k} \right) \cdot \frac{1}{\lambda_k^i} \right]$$

$\lambda_k^i \rightarrow$ Probability that the i^{th} point goes into the k^{th} cluster !!

Shrikanth
PAGE NO.
DATE:

$$= \max_{\theta} \sum_{i=1}^n \sum_{k=1}^K \left[\lambda_k^i \log \pi_k - \lambda_k^i \frac{(x_i - \mu_k)^2}{2\sigma_k^2} - \lambda_k^i \log \sqrt{2\pi} \sigma_k \right]$$

Take derivative wrt μ, σ to get :

$$\begin{aligned} \hat{\mu}_k^{\text{MML}} &= \frac{\sum_{i=1}^n \lambda_k^i x_i}{\sum_{i=1}^n \lambda_k^i} & \hat{\sigma}_k^2 \text{ MML} &= \frac{\sum_{i=1}^n \lambda_k^i (x_i - \hat{\mu}_k^{\text{MML}})^2}{\sum_{i=1}^n \lambda_k^i} \\ \text{mean of the} \\ \text{data points assigned} \\ \text{to a particular cluster} & & \text{sample variance} \\ & & \text{of data points} \\ & & \text{assigned to a particular cluster} \end{aligned}$$

Maximizing π :

$$\max_{\pi_1, \dots, \pi_K} \sum_{i=1}^n \left(\sum_{k=1}^K \lambda_k^i \log \pi_k \right) \quad \text{such that} \\ \sum_k \pi_k = 1 ; \pi_k \geq 0$$

can solve using method of Lagrange multipliers,

$$\hat{\pi}_k^{\text{MML}} = \frac{\sum_{i=1}^n \lambda_k^i}{n}$$

From fixing λ , we get :

$$\hat{\mu}_k^{\text{MML}} = \frac{\sum_{i=1}^n \lambda_k^i x_i}{\sum_{i=1}^n \lambda_k^i} ; \hat{\sigma}_k^2 \text{ MML} = \frac{\sum_{i=1}^n \lambda_k^i (x_i - \hat{\mu}_k^{\text{MML}})^2}{\sum_{i=1}^n \lambda_k^i}$$

$$\hat{\pi}_k^{\text{MML}} = \frac{\sum_{i=1}^n \lambda_k^i}{n}$$

II) Fix θ & Maximize λ
wrt

$$= \sum_{i=1}^n \left[\sum_{k=1}^K \lambda_k^i \log(a_{ik}) - \lambda_k^i \log \lambda_k^i \right]$$

$$\text{where } a_{ik} = \left(\pi_k e^{-(x_i - \mu_k)^2 / 2\sigma_k^2} \cdot \frac{1}{\sqrt{2\pi}\sigma_k} \right)$$

Fix any i .

$$\max_{\lambda_1^i, \dots, \lambda_K^i} \sum_{k=1}^K \left[\lambda_k^i \log(a_{ik}) - \lambda_k^i \log \lambda_k^i \right]$$

$$\text{such that } \sum_{k=1}^K \lambda_k^i = 1, \quad 0 \leq \lambda_k^i \leq 1$$

$$\text{Can be solved analytically, } P(z_i=k | x_i) = \frac{P(x_i | z_i=k) \cdot P(z_i=k)}{P(x_i)}$$

$$\lambda_k^i \text{ MML} = \frac{\left(\frac{1}{\sqrt{2\pi}\sigma_k} e^{-(x_i - \mu_k)^2 / 2\sigma_k^2} \right) \cdot \pi_k}{\sum_{k=1}^K \left(\frac{1}{\sqrt{2\pi}\sigma_k} e^{-(x_i - \mu_k)^2 / 2\sigma_k^2} \cdot \pi_k \right)}$$

$$P(z_i=k | x_i) = \frac{\sum_{k=1}^K \left(\frac{1}{\sqrt{2\pi}\sigma_k} e^{-(x_i - \mu_k)^2 / 2\sigma_k^2} \cdot \pi_k \right)}{\sum_{k=1}^K \left(\frac{1}{\sqrt{2\pi}\sigma_k} e^{-(x_i - \mu_k)^2 / 2\sigma_k^2} \cdot \pi_k \right)}$$

Lecture 8 : EM Algorithm

- Initialize $\theta^0 = \{\mu_1^0, \dots, \mu_K^0, \sigma_1^0, \dots, \sigma_K^0, \pi_1^0, \dots, \pi_K^0\}$ Tolerance Parameter
- until convergence ($\|\theta^{t+1} - \theta^t\| \leq \epsilon$)

Expectation step $\lambda^{t+1} = \arg \max_{\lambda} \text{modified log L}(\theta^t, \lambda)$

Maximization step $\theta^{t+1} = \arg \max_{\theta} \text{modified log L}(\theta, \lambda^{t+1})$

→ end

- EM produces "soft clustering"
- EM takes variances into account.
- EM need not be in linear regions!