

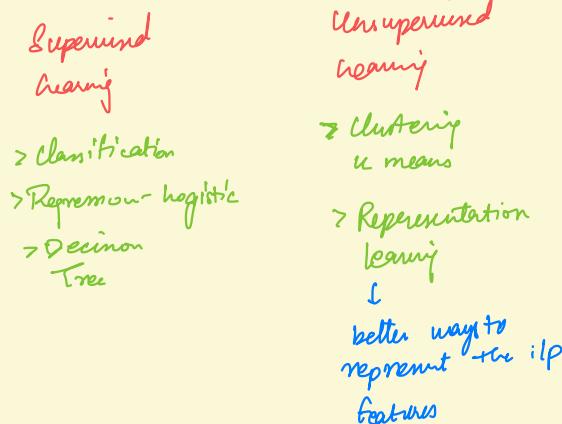
- (1) Generative vs Discriminative
- (2) Maximum Likelihood Estimation
- (3) Covariance Matrix - Closed form  
some Examples for GMM - motivation

- (4) GMM for 1D  
EM for 1D
- (5) GMM for  $d \geq 2$   
EM for  $d \geq 2$

#### References :

- (i) NPTEL - IIT Madras BS Deg Prog.  
MLT Dr Anirudh Rajkumar  
(YOUTUBE)
- (ii) Rituik math - GMM
- (iii) Machine Learning TV - GMM

# Broad Paradigms of ML



Sequential learning → make a decision, observe feedback based on it make the next decision

> Reinforcement learning

## - generative Model

- $P(x, y)$

given  $(x^i, y^i)$  we must be able to associate it with a density/probability func when features are continuous/discrete

If we make assumptions on  $p(x, y)$  how should the density look like then we're doing generative modelling

here we are modelling how the features themselves are generated & also how the labels are classified

[How sentences in English are generated]

Cares about how features are also generated  
Models feature generation also

If you know how to generate well you know how to discriminate well.

∴ If you know how data could have been generated, then gen. model is great.

## - Discriminative Model

- $P(y|x)$

Doesn't model the joint distribution of  $(x, y)$  but it only models  $P(y|x)$

( $P(y \text{ given } x)$ )

to be able to acquire this ability to discriminate b/w English & non-English sentences, you don't necessarily need to know how sentences themselves are generated

How is  $y$  given  $x$  is generated.  
, not how  $x$  itself is generated

Examples: K-means Decision Tree

## \* Estimation:

"There is some probabilistic mechanism that generates the data"

Goal: Given data find/estimate the distribution

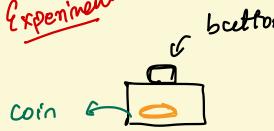
Assumptions: All observations are independent  
identically distributed

$$P(x_i | n_j) = x_i^p \quad \forall i, j$$

$$P(x_i = 1) = P(x_j = 1) = p \quad (\text{same coin is tossed})$$

i.i.d.: how determining the outcome is same  
everytime.

Experiment



every time you press  
the button the coin  
gets flipped. The

coin is an unbiased coin with  
prob. of head occurring to be 'p'.  
What is the choice of 'p' which  
generates data like

Σ 1 0 1 1 3

1 - head  
0 - tail

in 4 tosses?

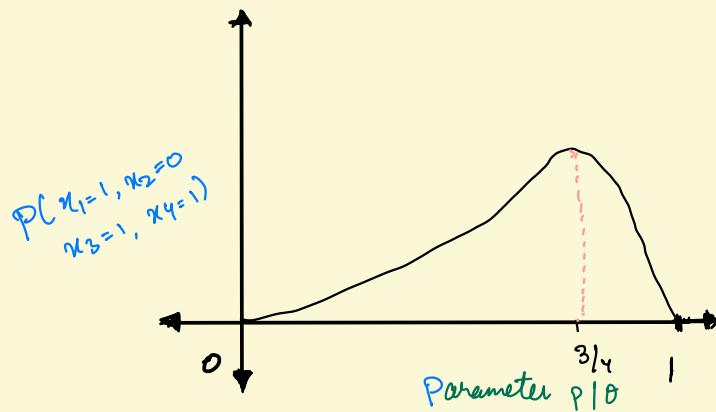
Among all choices of p,  
which choice has max. prob.

of generating the above observation.

How can we resolve it?

For each p between 0 & 1 calculate the chance of getting Σ 1 0 1 1 3.

i.e plot



$p = 3/4$  is most likely to generate the dp

Fisher's Principle of Maximum Likelihood.

Likelihood  
function

$$L(p, \{x_1, x_2, \dots, x_n\}) = P(x_1, x_2, \dots, x_n; p)$$

$$= P(x_1; p) P(x_2; p) \dots P(x_n; p) \quad (\text{By independence})$$

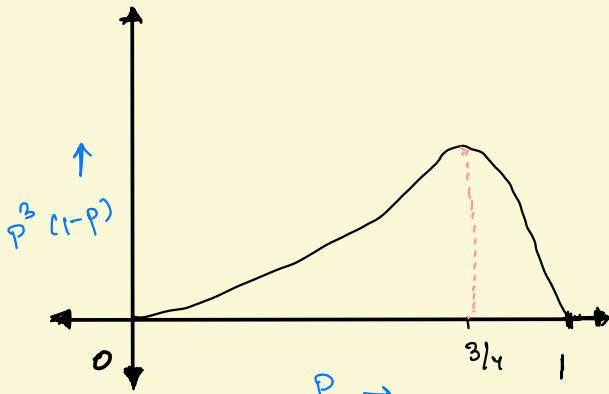
$$= \prod_{i=1}^n p^{x_i} (1-p)^{1-x_i}$$

(for this prob)

Prob. of  $\{x_1, x_2, \dots, x_n\}$  given the  
parameter (here prob. of head) is  $P$ .

$$P^{x_i^0}(-P)^{x_i^0} = \begin{cases} P, & x_i^0 = 1 \\ 1-P, & x_i^0 = 0 \end{cases}$$

$$\therefore P(x_1=1, x_2=0, x_3=1, x_4=1) = (P^1)(1-P)^0 P \cdot P = P^3(1-P)$$



$$\hat{P}_{ML} = \arg \max_P h(p, \{x_1, x_2, \dots, x_n\})$$

$$\hat{P}_{ML} = \arg \max_P \prod_{i=1}^n P^{x_i^0} (1-P)^{x_i^0}$$

↓ easier to look at

$$\hat{P}_{ML} = \arg \max_P \log \left( \prod_{i=1}^n P^{x_i^0} (1-P)^{x_i^0} \right)$$

$\log$  is monotone ↑

log likelihood fun.

$$\log [L(\{x_1, x_2, \dots, x_n\}, p)]$$

$$\hat{P}_{ML} = \arg \max_P \sum_{i=1}^n x_i^0 \log P + (1-x_i^0) \log (1-P)$$

$\underbrace{\log [L(p)]}$

$$\frac{\partial}{\partial P} (\log [L(p)]) = 0 \Rightarrow \sum_{i=1}^n \frac{x_i^0}{P} - \frac{(1-x_i^0)}{1-P} = 0$$

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

$$\Rightarrow \sum_{i=1}^n x_i^0 - np = 0 \Rightarrow \hat{P}_{ML} = \frac{\sum_{i=1}^n x_i^0}{n}$$

\* If Data =  $\{x_1, x_2, \dots, x_n\}$ ,  $x_i^0 \in \mathbb{R}$

Find  $\theta$  s.t.  $x_i \sim \text{Gaussian } (\mu, \sigma^2)$  &

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

Assume  $\sigma^2$  is known.

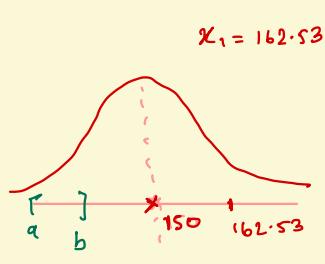
Eg: If  $\mu = 150$  what is prob. that  $x_1 = 162.53$ ?

for continuous dist.  $P(162.53, \mu=150)=0$  ∵ it's a point

Only intervals will have non zero values for probability.

$$\therefore P(x_i, \mu, \sigma^2) = 0 \quad \text{for any } x_i^0$$

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



$$\begin{aligned}
 h(\mu, \sigma^2, \{x_1, x_2, \dots, x_n\}) &= f_{x_1, x_2, \dots, x_n}^{(n_1, n_2, \dots, n_n, \mu, \sigma^2)} \quad \text{where } f \text{ is the} \\
 &= \prod_{i=1}^n f_{x_i}^{(x_i, \mu, \sigma^2)} \\
 &= \prod_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x_i - \mu)^2}{2\sigma^2}}
 \end{aligned}$$

$$\therefore \log(L(\mu, \sigma^2, \{x_1, x_2, \dots, x_n\})) = \sum_{i=1}^n \log\left(\frac{1}{\sqrt{2\pi}\sigma}\right) - \sum_{i=1}^n \frac{(x_i - \mu)^2}{2\sigma^2}$$

$$\frac{\partial \log L}{\partial \mu} = -2 \sum_{i=1}^n \frac{(x_i - \mu)}{2\sigma^2} (-1) = 0.$$

$$\Rightarrow \sum_{i=1}^n x_i - n\mu = 0 \quad \Rightarrow \quad \hat{\mu} = \frac{\sum_{i=1}^n x_i}{n}$$

If we assume the data  $\{x_1, x_2, \dots, x_n\}$  is generated acc. to a Gaussian with unknown mean  $\mu$ , the best guess acc. to ML principle is the sample mean

If I change the assumption that underlying dist. is Gaussian instead something else like Laplace then I get a different mean

## \* Covariance Matrix

variables →

|      | A  | B  |
|------|----|----|
| Sub1 | 1  | 1  |
| Sub2 | 3  | 0  |
| Sub3 | -1 | -1 |

link b/w happiness one derives from eating A  
happiness they derive by eating B       $A, B = 2$  Features

$$\begin{matrix} & \begin{matrix} A & B \end{matrix} \\ \begin{matrix} A \\ B \end{matrix} & \left[ \begin{array}{cc|cc} \text{cov}(A, A) & \text{cov}(A, B) & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot \\ \text{cov}(B, A) & \text{cov}(B, B) & \cdot & \cdot \end{array} \right] \end{matrix}$$

$$\text{cov}(A, A) = \text{Variance}(A)$$

$$\text{cov}(B, B) = \text{Variance}(B)$$

Correlation: Bounded b/w -1 & 1

Covariance: Not necessarily bdd b/w -1 & 1

Cov. matrix: Square matrix

Symmetric

high degree B covariance:

If Both are highly true or  
highly -ive

low degree of cov:

One is very true, other is very -ive

$$\text{Cov}(A, B) = E(AB) - E(A)E(B)$$

→ expected value of the prod. of the variables →  
the prod. of exp. value of each of the variables

$$E(B) = 0$$

$$E(A) = \frac{-1 + 3 + 1}{3} = 1$$

$$\text{Cov}(A, B) = \frac{2}{3} = \text{Cov}(B, A)$$

AB

$$E(AB) = \frac{1 + 0 + 1}{3} = \frac{2}{3}$$

| S <sub>1</sub> | 1 |
|----------------|---|
| S <sub>2</sub> | 0 |
| S <sub>3</sub> | 1 |

$$\text{Cov}(A, A) = E(A^2) - (E(A))^2$$

$$A^2 = \frac{11}{3} - 1 = 8/3$$

1  
9  
1

$$\text{Cov}(B, B) = E(B^2) - (E(B))^2$$

$$\frac{8}{3}$$

$$\text{Cov} = \begin{bmatrix} 8/3 & 2/3 \\ 2/3 & 8/3 \end{bmatrix}$$

\* Closed Form

Samples:  $x_1, x_2, \dots, x_N$  in  $\mathbb{R}^D$

$$X = N \times D$$

$\text{Cov. Blw } i^{\text{th}} \text{ & } j^{\text{th}}$  variables  $\equiv$  Cov. b/w  $i^{\text{th}}$  fruit &  $j^{\text{th}}$  fruit, the happiness index of those fruits

Cov Blw  $i^{\text{th}}$  &  $j^{\text{th}}$  variable:

$$\text{Cov}(x_i, x_j) = \frac{1}{N} \sum_{k=1}^N (x_{ki} - \bar{x}_i)(x_{kj} - \bar{x}_j) \quad x_{ki} = k^{\text{th}} \text{ sample } i^{\text{th}} \text{ feature}$$

$$X =$$

$$\begin{matrix} x_1 & x_2 & \dots & x_i & x_j & \dots & x_D \\ X^{(1)} & & & & & & \\ x^{(2)} & & & & & & \\ \vdots & & & & & & \\ x^{(n)} & & & & & & \\ \bar{x} & \bar{x}_1 & \bar{x}_2 & \dots & \bar{x}_i & \bar{x}_j & \dots & \bar{x}_D \end{matrix}$$

D Dimensional

Cov Matrix:

$$S = \left[ \begin{array}{cccc} \frac{1}{N} \sum_{k=1}^N (x_{k1} - \bar{x}_1)(x_{k1} - \bar{x}_1) & \frac{1}{N} \sum_{k=1}^N (x_{k1} - \bar{x}_1)(x_{k2} - \bar{x}_2) & \dots & \frac{1}{N} \sum_{k=1}^N (x_{k1} - \bar{x}_1)(x_{kD} - \bar{x}_D) \\ \frac{1}{N} \sum_{k=1}^N (x_{k2} - \bar{x}_2)(x_{k1} - \bar{x}_1) & \frac{1}{N} \sum_{k=1}^N (x_{k2} - \bar{x}_2)(x_{k2} - \bar{x}_2) & \dots & \frac{1}{N} \sum_{k=1}^N (x_{k2} - \bar{x}_2)(x_{kD} - \bar{x}_D) \\ \vdots & \vdots & \ddots & \vdots \\ \frac{1}{N} \sum_{k=1}^N (x_{kD} - \bar{x}_D)(x_{k1} - \bar{x}_1) & \frac{1}{N} \sum_{k=1}^N (x_{kD} - \bar{x}_D)(x_{k2} - \bar{x}_2) & \dots & \frac{1}{N} \sum_{k=1}^N (x_{kD} - \bar{x}_D)(x_{kD} - \bar{x}_D) \end{array} \right]$$

$$S = \frac{1}{N} \begin{bmatrix} x_{11} - \bar{x}_1 & x_{21} - \bar{x}_1 & \dots & x_{N1} - \bar{x}_1 \\ x_{12} - \bar{x}_2 & x_{22} - \bar{x}_2 & \dots & x_{N2} - \bar{x}_2 \\ \vdots & \vdots & \ddots & \vdots \\ x_{1D} - \bar{x}_D & x_{2D} - \bar{x}_D & \dots & x_{ND} - \bar{x}_D \end{bmatrix} \begin{bmatrix} x_{11} - \bar{x}_1 & x_{12} - \bar{x}_2 & \dots & x_{1D} - \bar{x}_D \\ x_{21} - \bar{x}_1 & x_{22} - \bar{x}_2 & \dots & x_{2D} - \bar{x}_D \\ \vdots & \vdots & \ddots & \vdots \\ x_{N1} - \bar{x}_1 & x_{N2} - \bar{x}_2 & \dots & x_{ND} - \bar{x}_D \end{bmatrix}^T \quad N \times D$$

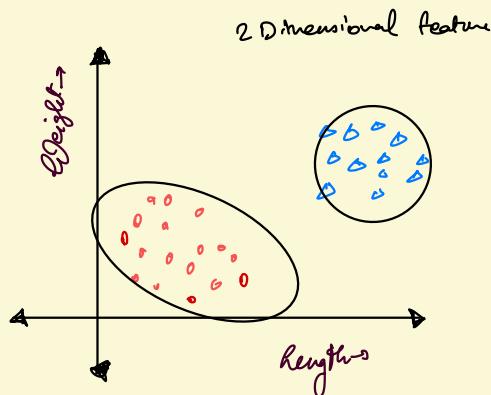
Python

$$S = \frac{1}{N} (X - \bar{X})^T (X - \bar{X}) = \frac{1}{N} \text{np.dot}((X - \bar{X})^T, (X - \bar{X}))$$

An Example.

### \* Gaussian Mixture

○ = Tuna  
○ = Salmon

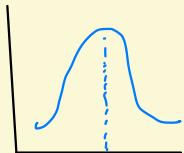


Based on weight & length classify the fish as Tuna or Salmon

GMM  $\Rightarrow$  all classes are distributed in Gaussian Distribution.

One dimension

Normal



$\pi_i$

For higher dimension

for each of class (Tuna / Salmon) the mean is  $\mu$  & instead of variance we have covariance  $\Sigma$  (this describes the interaction b/w weight & length & hence decides the shape - either circular / elliptical, tilted or straight)

$\pi_i$  = prob. for each of the class to occur.

Goal: GMM tries to figure out optimal  $[\mu, \Sigma, \pi]$  for each of the class

Tuna Distribution

$$N(x | \mu_T, \Sigma_T)$$

Salmon Distribution.

$$N(x | \mu_S, \Sigma_S)$$

If  $x = \text{tuna}$  then it's distributed according to  $\mu_T$  &  $\Sigma_T$

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

$$P(\text{Salmon}) = \pi_S$$

$$x = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix}$$

$$P(x) = \pi_T N(x | \mu_T, \Sigma_T) + \pi_S N(x | \mu_S, \Sigma_S)$$

$x$  = training example

2 classes

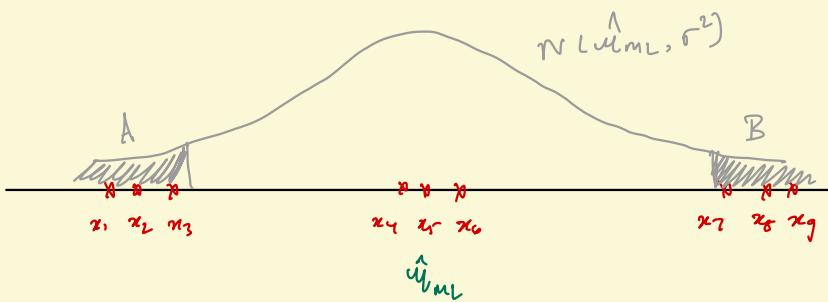
Bayes

Q How do we pick the best  $\pi, \mu, \Sigma$ ?

We want to maximize  $P(x | \pi, \Sigma, \mu)$

seeing the data  $X$  given the parameters  $\mu, \pi, \Sigma$

## Gaussian Mixture Model



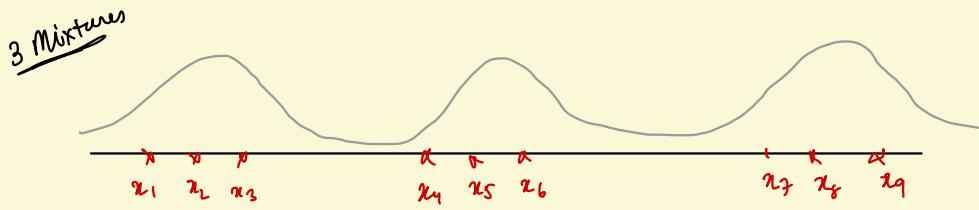
What could be a generative story?

Is Gaussian a good model for the data with mean  $\hat{\mu}_{ML} = \frac{3x_i}{9}$

The problem is, this Gaussian with  $\hat{\mu}_{ML}$  has very less prob. in the regions A & B. But there are data points in the region.

∴ Problem is we have assumed model is Gaussian.

We want a density like  to explain the data points.



\* Univariate

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

\* Multivariate

$$N(x | \mu, \Sigma) = \frac{1}{(2\pi)^{K/2} |\Sigma|^{1/2}} e^{-\frac{1}{2}(x-\mu)^T \Sigma^{-1} (x-\mu)}$$

why Gaussian?

- ≥ most of obs. follow Gaussian
- ≥ easy to do math. manipulation eg- diff.

Step 1 Generate a mixture component among  $\xi_1, \xi_2, \dots, \xi_K$  call  $\xi_i^*$ ,  $1 \leq i \leq n$

For each  $\xi_i^*$  assign a mixture  $\pi_i$  according to probabilities

$$P(\xi_i^* = k) = \pi_k \quad \left( \sum_{i=1}^K \pi_i = 1 \quad 0 \leq \pi_i \leq 1 \quad \forall i \right) \text{ Pick the one with highest prob.}$$

Step 2 Generate  $x_i^* \sim N(\mu_{\xi_i^*}, \sigma_{\xi_i^*}^2)$

$x_1, x_2, \dots, x_n$  = observed  
 $\xi_1, \xi_2, \dots, \xi_n$  = Unobserved / latent variable

$$\Pi = [\pi_1, \pi_2, \dots, \pi_K] \quad i.e. [\mu_K, \sigma_K^2]$$

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

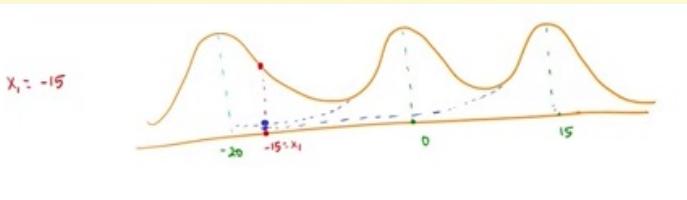
Determine/  
Estimate  $3K-1$  parameters

## \* Likelihood function for HMM

$$L \left( \begin{matrix} \mu_1, \dots, \mu_K \\ \sigma_1^2, \dots, \sigma_K^2 \\ \pi_1, \dots, \pi_K \end{matrix}; x_1, x_2, \dots, x_n \right) = \prod_{i=1}^n f(x_i; \begin{matrix} \mu_1, \mu_2, \dots, \mu_K \\ \sigma_1^2, \sigma_2^2, \dots, \sigma_K^2 \\ \pi_1, \pi_2, \dots, \pi_K \end{matrix})$$

→ Product since  
 $x_1, x_2, \dots, x_n$  are  
indep & iid.

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



For one dimension

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

*neg likelihood*

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

Not possible to solve this  
analytically  
hence we cannot arrive at a  
closed form or explicit expression  
for  $\pi_k$  or  $\mu_k$  or  $\sigma_k$

And also constraints are involved.

It's a maximization problem.

Need an alternate way to solve the above

## \* Convex functions and Jensen's Inequality

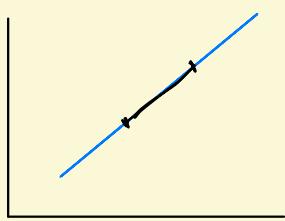
Convex functions

$$x, y \in D$$

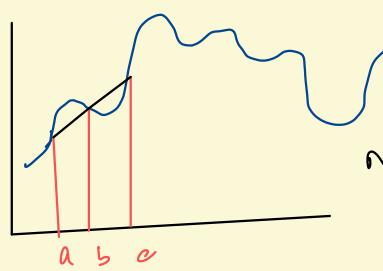
$$f(\theta x + (1-\theta)y) \leq \theta f(x) + (1-\theta)f(y), \quad \theta \in [0,1]$$

Concave func.

$$f(\theta x + (1-\theta)y) \geq \theta f(x) + (1-\theta)f(y), \quad \theta \in [0,1]$$



Both convex & concave



Neither concave nor convex

\* Jensen's Inequality:  $f \left( \sum_{k=1}^K \lambda_k a_k \right) \geq \sum_{k=1}^K \lambda_k f(a_k)$  with  $\sum_{k=1}^K \lambda_k = 1$

- $\log$  is a concave function.

$\pi_k$  = prob. of each cluster  $k$

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

For every data point  $i$  introduce the parameters  $\{\lambda_1^i, \lambda_2^i, \dots, \lambda_K^i\}$  s.t.  $\forall i \sum_{k=1}^K \lambda_k^i = 1$

$$0 \leq \lambda_k^i \leq 1 \quad \forall i, k$$

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

By Jensen's Inequality,

$$\log(h(\theta)) \geq \underbrace{\sum_{i=1}^n \sum_{k=1}^K \lambda_k^i \log \left( \frac{\pi_k e^{-\frac{(x_i - \mu_k)^2}{2\sigma_k^2}}}{\sqrt{2\pi} \sigma_k \lambda_k^i} \right)}_{\text{modified log likelihood obtained by Jensen's}}$$

By Jensen's for each:

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

gives a lower bound for the true log likelihood @  $\theta$ .

for any choice of  $\lambda = \{\lambda_1^1, \dots, \lambda_K^1, \lambda_1^2, \dots, \lambda_K^2, \dots, \lambda_1^n, \dots, \lambda_K^n\}$  we are just getting a lower bound for the likelihood.

But what are we gaining?

Key insight: If we fix some  $\lambda$ , it's easy to maximize w.r.t.  $\theta$ .  
If we fix some  $\theta$ , it's easy to maximize w.r.t.  $\lambda$ .

\* Fix  $\lambda$  and maximize over  $\theta$

$$l = \underset{\theta}{\text{maximize}} \sum_{i=1}^n \sum_{k=1}^K \lambda_k^i \log \pi_k - \lambda_k^i \frac{(x_i - \mu_k)^2}{2\sigma_k^2} - \lambda_k^i \log(\sqrt{2\pi} \sigma_k) - \lambda_k^i \log(\lambda_k^i)$$

$$\textcircled{1} \frac{\partial l}{\partial \mu_k} = 0 \Rightarrow \sum_{i=1}^n \frac{\lambda_k^i (x_i - \mu_k)}{\sigma_k^2} = 0 \Rightarrow$$

$$\hat{\mu}_k^{\text{ML}} = \frac{\sum_{i=1}^n \lambda_k^i x_i}{\sum_{i=1}^n \lambda_k^i}$$

$$\textcircled{2} \quad \frac{\partial J}{\partial \pi_k} = 0 \Rightarrow \sum_{i=1}^n \left( \frac{\lambda_k^i (\mathbf{x}_i - \boldsymbol{\mu}_k)^2}{\sigma_k^2} - \frac{\lambda_k^i}{\sqrt{2\pi} \sigma_k} \right) = 0$$

$$\Rightarrow \sum_{i=1}^n \lambda_k^i ((\mathbf{x}_i - \boldsymbol{\mu}_k)^2 - \sigma_k^2) = 0 \Rightarrow$$

$$\hat{\pi}_k^{\text{ML}} = \frac{\sum_{i=1}^n \lambda_k^i (\mathbf{x}_i - \boldsymbol{\mu}_k)^2}{\sum_{i=1}^n \lambda_k^i}$$

$\lambda_k^i$  = prob. that  $i^{\text{th}}$  point goes to the  $k^{\text{th}}$  cluster  
what you think is the prob.

$$\textcircled{3} \quad \max_{\theta} \sum_{i=1}^n \sum_{k=1}^K \lambda_k^i \log \pi_k \quad \text{s.t.} \quad \sum_{k=1}^K \pi_k = 1, \quad \pi_k \geq 0$$

can solve this using the method of Lagrange multipliers  $\gamma$  = multiplier

$$L^* = \sum_{i=1}^n \sum_{k=1}^K \lambda_k^i \log (\pi_k) - \gamma \left( \sum_{k=1}^K \pi_k - 1 \right) \quad \text{And } \pi_k^*, \pi_k^*$$

$$\textcircled{1} \quad \frac{\partial L^*}{\partial \pi_k^*} = 0 \Rightarrow \sum_{i=1}^n \frac{\lambda_k^i}{\pi_k^*} - \gamma^* = 0 \Rightarrow \gamma^* = \frac{\sum_{i=1}^n \lambda_k^i}{(\pi_k^*)}$$

$$\therefore \pi_k^* = \left( \frac{\sum_{i=1}^n \lambda_k^i}{\gamma^*} \right)$$

$$\textcircled{2} \quad \text{Substituting in the constraint,} \quad \sum_{k=1}^K \pi_k^* = 1 \Rightarrow \sum_{k=1}^K \frac{\sum_{i=1}^n \lambda_k^i}{(\gamma^*)} = 1$$

$$\therefore (\gamma^*)^{-1} = \sum_{k=1}^K \sum_{i=1}^n (\lambda_k^i) \quad \text{And, } \sum_{k=1}^K (\lambda_k^i) = 1$$

$$\therefore \gamma^* = n$$

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

\* Now fix  $\theta$  & maximize  $\lambda$

modified

likelihood function:

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

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

for any  $i$ ,

$$\max_{\lambda_1^i, \lambda_2^i, \dots, \lambda_K^i} \sum_{k=1}^K \lambda_k^i \log \left( \frac{\pi_k e^{-\frac{(x_i - \mu_k)^2}{2\sigma_k^2}}}{\sqrt{2\pi} \sigma_k} \right) - \sum_{k=1}^K \lambda_k^i \log (\lambda_k^i)$$

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

logrange multiplier,

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

$$\textcircled{1} \quad \frac{\partial L^*}{\partial \lambda_k^i} = \log \left( \frac{\pi_k e^{-\frac{(x_i - \mu_k)^2}{2\sigma_k^2}}}{\sqrt{2\pi} \sigma_k} \right) - \log (\lambda_k^i) - 1 + \alpha^* = 0$$

$$\Rightarrow \frac{\pi_k e^{-\frac{(x_i - \mu_k)^2}{2\sigma_k^2}}}{\sqrt{2\pi} \sigma_k (e^{\alpha^*} - 1)} = \lambda_k^i$$

$$\textcircled{2} \quad \sum_{k=1}^K \lambda_k^i = 1 \Rightarrow \sum_{k=1}^K \left( \frac{\pi_k e^{-\frac{(x_i - \mu_k)^2}{2\sigma_k^2}}}{\sqrt{2\pi} \sigma_k} \right) = e^{\alpha^*} - 1$$

hence,

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

$$\hat{\lambda}_k^i = P(z_i = k | x_i)$$

\* Algorithm:

- ① Initialize  $\theta^0 = \{ \mu_1^0, \dots, \mu_K^0, \sigma_1^2, \dots, \sigma_K^2, \pi_1^0, \dots, \pi_K^0 \}$
- ② until convergence  $\| \theta^{t+1} - \theta^t \| \leq \epsilon$  Tolerance

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

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

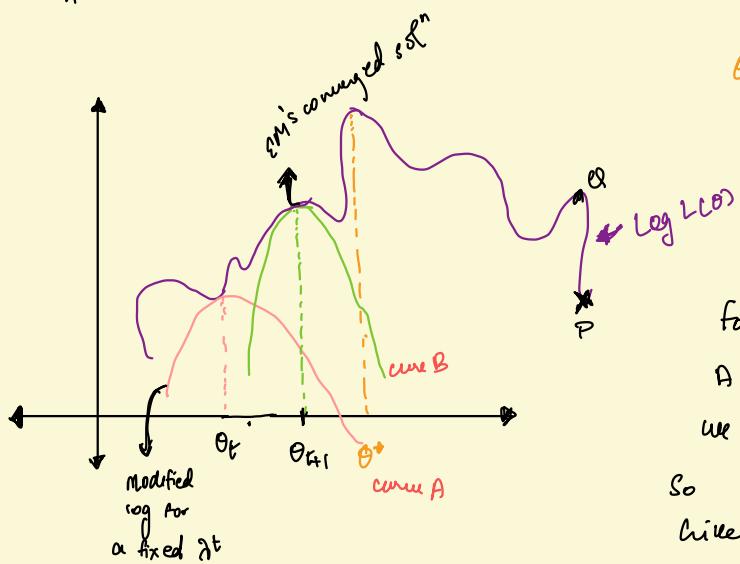
① K-Means Algorithm.:
 

- Step 1 → Compute Means
- Step 2 → Reassignment

② EM produces soft clustering.

$\pi_k^t$  = probability that  $i$ th datapoint goes to cluster  $K$   
 EM takes variance into account

But are we solving for the original problem? We are maximizing only the modified log function



$\theta^*$  is the actual value that we wish to get.

But the  $\theta_t$  that we get from any choice of  $\pi_t$ ?  
 Solving the modified log likelihood is always a lower bound for  $\theta^*$ .

for a fixed  $\pi_t$ , we find the max  $\theta_t$  for curve A with that  $\theta_t$  we find  $\pi^{t+1}$ . For  $\pi^{t+1}$  we find  $\theta_{t+1}$  of curve B  
 So EM converges to a local maximum of log likelihood.

So if I start with an initialization @ P I can only reach the local maxima @ ④  
 which is not close to the actual maximum  $\theta^*$

so given the data run K means clustering, get the hard clusters. Once we have the clusters compute sample mean and sample variance for each cluster which will give the  $\mu_k$  &  $\sigma_k^2$  for each cluster  $k$  & to get the  $\pi_k$ 's look at the fraction of data points in each of the clusters.

$\therefore \theta^0$  is got from K-means.

### \*EM Algorithm.

Run K means on the data & find  $\theta^0$  from K-means algorithm

① Initialize  $\theta^0 = \{ \mu_1^0, \dots, \mu_K^0, \sigma_1^2, \dots, \sigma_K^2, \pi_1^0, \dots, \pi_K^0 \}$

② until convergence  $\|\theta^{t+1} - \theta^t\| \leq \epsilon$  <sup>Tolerance</sup>

$$\lambda^{t+1} = \underset{\lambda}{\arg \max} \text{ modified } \log L(\theta^t, \lambda) \rightarrow \text{Expectation Step}$$

$$\theta^{t+1} = \underset{\theta}{\arg \max} \text{ modified } \log L(\theta, \lambda^{t+1}) \rightarrow \text{Maximization Step}$$

end.

As the no. of clusters increase ( $K \uparrow$ ) GMM can approximate almost any distribution in  $\mathbb{R}^n$

\* Derivation for case where  $\dim = 2$

$$\begin{pmatrix} \text{Six Modes} \\ \pi_k = \omega_k \\ \lambda_k^i = \frac{\omega_k}{\sum_k} \end{pmatrix}$$

Instead of variance we have the covariance matrix.

For each cluster  $k \in \{1, 2, \dots, K\}$  we have  $\{\mu_k, \pi_k, \Sigma_k\}$   
 $\downarrow$   
 $\text{dimensional vector}$        $\downarrow$   
 $\frac{\partial}{\partial x}$

$$L(\frac{u_1, \dots, u_n}{\Sigma_1, \dots, \Sigma_K}, \pi_1, \dots, \pi_K, x_1, x_2, \dots, x_n) = \prod_{i=1}^n f(x_i; \frac{u_1, u_2, \dots, u_n}{\Sigma_1, \Sigma_2, \dots, \Sigma_K}, \frac{\pi_1, \pi_2, \dots, \pi_K}{\pi_1 + \pi_2 + \dots + \pi_K}) \rightarrow \text{Product since } u_1, u_2, \dots, u_n \text{ are independent & iid.}$$

$$= \prod_{i=1}^n \left( \sum_{k=1}^K \pi_k f(x_i; \mu_k, \Sigma_k) \right)$$

$$= \prod_{i=1}^n \left( \sum_{k=1}^K \frac{\pi_k}{(2\pi)^{D/2} |\Sigma_k|^{1/2}} e^{-\frac{1}{2}(x_i - \mu_k)^T \Sigma_k^{-1} (x_i - \mu_k)} \right)$$

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

Introduce for each  $i$ :  $\sum_{k=1}^K \lambda_k^i = 1$

latent variable  
 $\lambda_k^i = \text{prob. of sample } i \text{ occurring in cluster } k$

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

$\log$  is a concave func. Jensen's Ineq  $\Rightarrow$

$$\log(L(\theta, \pi)) \geq \sum_{i=1}^n \sum_{k=1}^K \lambda_k^i \log \left( \frac{\pi_k e^{-\frac{1}{2}(x_i - \mu_k)^T \Sigma_k^{-1} (x_i - \mu_k)}}{(2\pi)^{D/2} |\Sigma_k|^{1/2}} \right)$$

$$\log(L(\theta, \pi)) \geq \sum_{i=1}^n \sum_{k=1}^K \lambda_k^i (\log \pi_k - \frac{1}{2}(x_i - \mu_k)^T \Sigma_k^{-1} (x_i - \mu_k) - \frac{1}{2} \ln(2\pi) - \frac{1}{2} \ln(|\Sigma_k|) - \ln(\lambda_k^i))$$

modified log function

(A) Maximize modified log function  $\mathbb{Q}$  given  $\pi_k^*$

$$\textcircled{1} \quad \frac{\partial \mathbb{Q}}{\partial \mu_k} = 0 \Rightarrow \sum_{i=1}^n (\pi_i - \mu_k) \sum_k \lambda_i^k = 0 \quad (\sum_k \neq 0)$$

$$\boxed{\frac{\sum_{i=1}^n (\pi_i) \lambda_i^k}{\sum_{i=1}^n \lambda_i^k} = \mu_k}$$

$$\begin{aligned} \frac{\partial}{\partial A} (x^T A^{-1} y) &= \\ -(A^T)^T x y^T (A^T)^{-1} &= \\ \frac{\partial}{\partial A} (\ln |A|) &= \\ (A^{-1})^T = (A^T)^{-1} & \end{aligned}$$

$$\textcircled{2} \quad \frac{\partial \mathbb{Q}}{\partial \Sigma_k} = 0 \Rightarrow \sum_{i=1}^n \frac{1}{2} (\Sigma_k^{-1})^{-1} (\pi_i - \mu_k) (\pi_i - \mu_k)^T (\Sigma_k^{-1})^{-1} \lambda_i^k - \frac{1}{2} \sum_{i=1}^n (\Sigma_k^{-1})^{-1} \lambda_i^k = 0$$

Postmul. by  $(\Sigma_k^T)$

$$\Rightarrow \frac{(\Sigma_k^T)^{-1}}{2} \left( \sum_{i=1}^n (\pi_i - \mu_k) (\pi_i - \mu_k)^T \lambda_i^k - \sum_{i=1}^n \lambda_i^k (\Sigma_k^T) \right) = 0$$

$$\Rightarrow \boxed{\Sigma_k^T = \frac{\sum_{i=1}^n (\pi_i - \mu_k) (\pi_i - \mu_k)^T \lambda_i^k}{\sum_{i=1}^n \lambda_i^k}}$$

$$\textcircled{3} \quad \text{with constraint } \sum_{k=1}^K \pi_k = 1 \quad d \in \mathbb{R} \text{ is a log multip.}$$

$$\mathbb{Q}^* = \mathbb{Q} + \alpha \left( \sum_{k=1}^K \pi_k - 1 \right)$$

$$\textcircled{i} \quad \frac{\partial \mathbb{Q}^*}{\partial \pi_k} = 0 \Rightarrow \sum_{i=1}^n \frac{\lambda_i^k}{\pi_k} - d^* = 0$$

$$\Rightarrow \pi_k^* = \frac{\left( \sum_{i=1}^n \lambda_i^k \right)}{d^*}$$

$$\textcircled{ii} \quad \sum_{k=1}^K \pi_k^* = 1$$

$$\Rightarrow \sum_{k=1}^K \sum_{i=1}^n \frac{\lambda_k^i}{\alpha^*} = 1 \quad \text{Now for each } i, \quad \sum_{k=1}^K \lambda_k^i = 1$$

$$\sum_{i=1}^n \frac{1}{\alpha^*} = 1 \Rightarrow \alpha^* = \sum_{i=1}^n 1 = n$$

$\therefore$  we get,

$$\Pi_k^* = \frac{\sum_{i=1}^n \lambda_k^i}{n}$$

(B) Maximize  $\lambda_k^i$  given parameters  $(\mu_k, \Sigma_k, \Pi_k)$

$$Q \text{ with } \sum_{k=1}^K \lambda_k^i = 1 \quad \text{for each } i \quad \alpha^* \Pi_k \text{ is a lag multiplier}$$

for each  $i$ ,

$$Q^* = Q + \alpha^* \left( \sum_{k=1}^K \lambda_k^i - 1 \right)$$

$$(i) \quad \frac{\partial Q^*}{\partial \lambda_k^i} \Bigg|_{\lambda_k^{i*}, \alpha^*} = 0 \Rightarrow \frac{\partial}{\partial \lambda_k^i} \left( \lambda_k^i (\log \Pi_k) - \left( \frac{1}{2} (x_i - \mu_k)^T \Sigma_k^{-1} (x_i - \mu_k) - \frac{1}{2} \ln(2\pi) - \frac{1}{2} \ln(\det \Sigma_k) \right) + \frac{\partial}{\partial \lambda_k^i} [\alpha^* \left( \sum_{k=1}^K \lambda_k^i - 1 \right)] \right) = 0$$

$$\Rightarrow \log \left( \frac{\Pi_k^* e^{-\frac{1}{2} (x_i - \mu_k)^T \Sigma_k^{-1} (x_i - \mu_k)}}{\sqrt{2\pi} |\Sigma_k|} \right) + \frac{\partial}{\partial \lambda_k^i} [\lambda_k^i \ln(\lambda_k^i)] + \alpha^* = 0$$

$$\Rightarrow \log \left( \frac{\Pi_k^* e^{-\frac{1}{2} (x_i - \mu_k)^T \Sigma_k^{-1} (x_i - \mu_k)}}{\sqrt{2\pi} |\Sigma_k|} \right) - \log \lambda_k^i - 1 + \alpha^* = 0 \Rightarrow$$

$$\lambda_k^{i*} = \frac{\Pi_k^* e^{-\frac{1}{2} (x_i - \mu_k)^T \Sigma_k^{-1} (x_i - \mu_k)}}{\sqrt{2\pi} |\Sigma_k| (e^{\alpha^*} - 1)}$$

$$(ii) \quad \sum_{k=1}^K \lambda_k^i = 1 \Rightarrow \sum_{k=1}^K \frac{\Pi_k^* e^{-\frac{1}{2} (x_i - \mu_k)^T \Sigma_k^{-1} (x_i - \mu_k)}}{\sqrt{2\pi} |\Sigma_k| (e^{\alpha^*} - 1)} = 1$$

$$e^{\alpha^*} = \frac{1}{\sqrt{2\pi}} \sum_{k=1}^K \pi_k e^{-\frac{1}{2} (\mathbf{x}_i - \boldsymbol{\mu}_k)^T \Sigma_k^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_k)}$$

Hence

$$\pi_k^* = \frac{\frac{1}{\sqrt{2\pi}} \pi_k e^{-\frac{1}{2} (\mathbf{x}_i - \boldsymbol{\mu}_k)^T \Sigma_k^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_k)}}{\sum_{k=1}^K \frac{1}{\sqrt{2\pi}} \pi_k e^{-\frac{1}{2} (\mathbf{x}_i - \boldsymbol{\mu}_k)^T \Sigma_k^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_k)}}$$

## EM Algorithm for GMMs

given a training set as  $\mathcal{D} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\}$

### EM algorithm for GMMs

initialize  $\{w_m^{(0)}, \mu_m^{(0)}, \Sigma_m^{(0)}\}$ , set  $n = 0$

**while** not converged **do**

**E-step:** for all  $m = 1, \dots, M$  and  $i = 1, \dots, N$ :

$$\text{Given } \{w_m^{(n)}, \mu_m^{(n)}, \Sigma_m^{(n)}\} \cup \{\mathbf{x}_i\} \longrightarrow \{\xi_m^{(n)}(\mathbf{x}_i)\}$$

Bind Expectation

**M-step:** for all  $m = 1, \dots, M$ :

$$\text{Given } \{\xi_m^{(n)}(\mathbf{x}_i)\} \cup \{\mathbf{x}_i\} \longrightarrow \{w_m^{(n+1)}, \mu_m^{(n+1)}, \Sigma_m^{(n+1)}\}$$

Bind Maximization

$$n = n + 1$$

**end while**

$$w_m = \pi_m$$

$$\xi_m^n = \Delta_k^p$$

- k-means only considers the mean to update the centroid while GMMs takes into account the mean as well as the variance of the data!