

NPTEL Week-9 Live Session

on Machine Learning and Deep Learning - Fundamentals and Applications (noc24_ee146)

A course offered by: Prof. Manas Kamal Bhuyan, IIT Guwahati

NPTEL Quiz Solution: Week-8

→ GMM
→ K-means clustering.



$D = \{x, y\}$ → supervised
↑
Feature label
 $D = \{x\}$ → unsupervised case.

By

Arka Roy

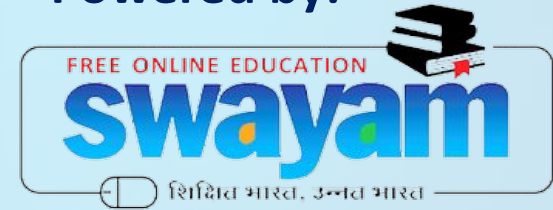
NPTEL PMRF TA

Prime Minister's Research Fellow

Department of Electrical Engineering, IIT Patna

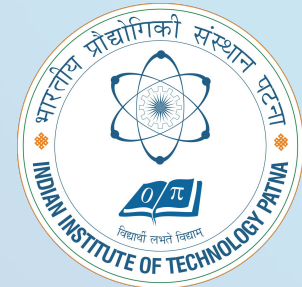
Web: <https://sites.google.com/view/arka-roy/home>

Powered by:



PMRF

Prime Minister's Research Fellows
Ministry of Education
Government of India



1) The EM algorithm is based on the principle of

- ✓ ☒ Maximum likelihood estimation
- ☐ Bayesian theory
- ☐ Feature selection
- ☐ Regularization

$$GMM = \sum_k \pi_k N(x | \mu_k, \Sigma_k)$$

- * If we have to evaluate these parameters using MLE then, we have to do partial differentiation for all such k w.r.t μ_k and Σ_k .
- * Therefore solving such k w.r.t differential solutions is difficult as for analytical solution is (impossible)
- * Therefore we go by computational method of solving through EM algorithm

GMM

= Linear Combination of different gaussian distribution

$$= \sum_{i=1}^K \pi_k N(x | \mu_k, \Sigma_k)$$

Multi Variate Gaussian Distribution

MLE:-

$$f(x|\theta)$$
$$\ln f(x|\theta)$$
$$\frac{\partial}{\partial \theta} \ln f(x|\theta) = 0$$
$$\theta = \{\mu, \Sigma\}$$

2) Each component in a GMM is characterized by

○ Mean, ~~median, and mode~~

✓ ○ Mean and covariance matrix

○ Variance, ~~skewness, and kurtosis~~

○ ~~Standard deviation and correlation coefficient~~

Mean = 1st order statistical moment
Variance = 2nd order statistical moment
In case of Gaussian or normal dist = $N(\mu, \Sigma)$

$x \in \mathbb{R}^D$

$\rightarrow GMM = \sum_{k=1}^K \pi_k N(x | \underline{\mu}_k, \underline{\Sigma}_k)$

weighting / mixing coefficients

Normal distribution \rightarrow multi variate normal distribution.

Parameterised by

\rightarrow Mean, & Covariance matrix

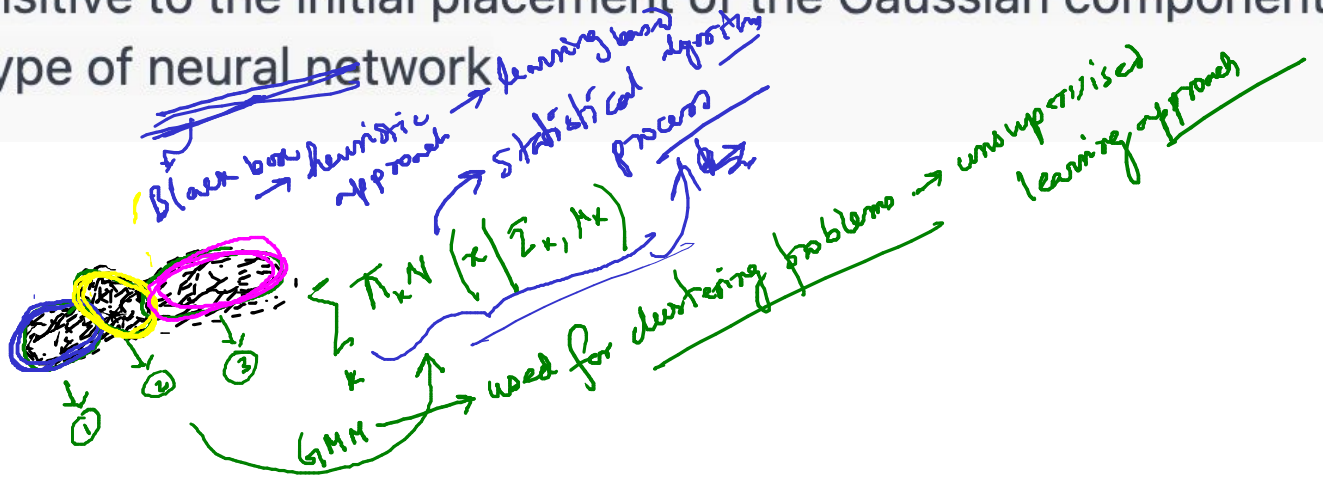
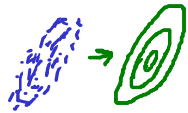
μ_k, Σ_k of multivariate Gaussian distribution

$GMM \rightarrow$ weighting factor, μ_k, Σ_k

$CC = \frac{Covariance(x_d)}{Var(x_d) \cdot Var(x_d)}$

3) Which of the following statements is true about GMMs?

- ~~GMMs are only used for regression tasks~~ *clustering Task.*
- ~~GMMs can only model data with a single cluster~~
- ~~GMMs are sensitive to the initial placement of the Gaussian components~~
- ~~GMMs are a type of neural network~~



4) What does each component in a GMM represent?

- ☐ ~~A data point in the dataset~~
- ☐ ~~A principal component~~
- ☐ ~~A cluster center~~ (μ_k)
- ☒ A Gaussian distribution with its mean and covariance

$$GMM = \sum_k \pi_k \mathcal{N}(x | \mu_k, \Sigma_k)$$

Handwritten red annotations: Arrows point from μ_k and Σ_k in the formula to the corresponding terms in the list above.

5) Let there be 10 data points: $\{-1, 27, 31, 2, 59, 3, 61, 34, 0, 12\}$. Use K-means clustering for 3 iterations to cluster them into 3 clusters. The initial labels are $\{1, 1, 1, 2, 2, 2, 3, 3, 3, 1\}$. List out the labels after 3 iterations. $\rightarrow \{1, 2, 2, 1, 3, \dots\}$

☐ $\{1, 2, 2, 1, 3, 1, 3, 2, 1, 1\}$

☐ $\{1, 1, 2, 1, 3, 1, 3, 2, 1, 2\}$

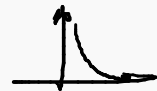
☐ $\{1, 2, 3, 1, 3, 1, 3, 2, 3, 1\}$

☐ $\{1, 1, 2, 1, 3, 1, 3, 3, 1, 1\}$

$D = \{-1, 27, 31, 2, 59, 3, 61, 34, 0, 12\}$

$\mathcal{L} = \{1, 1, 1, 2, 2, 2, 3, 3, 3, 1\}$

(C_1)



$D_1 = \{-1, 27, 31, 12\} \in \text{Class 1} \Rightarrow \text{Centroid of class 1 data cloud} = \frac{-1 + 27 + 31 + 12}{4}$

$= 17.25$

(ϕ)
Distance = $\sqrt{(x_1 - x_2)^2}$ (P^2 dist)
 $= |x_1 - x_2|$ (P^1 dist)

Centroid-2 = $(2 + 59 + 3)/2 = 21.33$ (C_2)
Centroid-3 = $(0 + 61 + 34)/3 = 31.66$ (C_3)

$\phi(P_3, C_1) = 13.75$
 $\phi(P_3, C_2) = 9.67$
 $\phi(P_3, C_3) = 0.67$

$\Rightarrow P_3 \in W_3$

It-1:- $\phi(P_1, C_1) = |-1 - 17.25| = 18.25 \rightarrow P_1 \in W_1$

$\phi(P_1, C_2) = 22.33$
 $\phi(P_1, C_3) = 32.66$

$\phi(P_2, C_1) = 9.75$
 $\phi(P_2, C_2) = 5.67$
 $\phi(P_2, C_3) = 4.66$
 $\Rightarrow P_2 \in W_3$

$\phi(P_4, C_1) = 15.75, \phi(P_4, C_2) = 19.33, \phi(P_4, C_3) = 29.67$
 $P_4 \in W_1$

$\phi(P_5, C_1) = 41.75$
 $\phi(P_5, C_2) = 37.65$
 $\phi(P_5, C_3) = 27.23$
 $\Rightarrow P_5 \in W_3$

$\phi(P_6, C_1) = 14.25$
 $\phi(P_6, C_2) = 18.33$
 $\phi(P_6, C_3) = 28.67$
 $\Rightarrow P_6 \in W_1$

- 5) Let there be 10 data points: $\{-1, 27, 31, 2, 59, 3, 61, 34, 0, 12\}$. Use K-means clustering for 3 iterations to cluster them into 3 clusters. The initial labels are $\{1, 1, 1, 2, 2, 2, 3, 3, 3, 1\}$. List out the labels after 3 iterations.

$D = \{-1, 27, 31, 2, 59, 3, 61, 34, 0, 12\}$

$\{1, 2, 2, 1, 3, 1, 3, 2, 1, 1\}$

$\{1, 1, 2, 1, 3, 1, 3, 2, 1, 2\}$

$\{1, 2, 3, 1, 3, 1, 3, 2, 3, 1\}$

$\{1, 1, 2, 1, 3, 1, 3, 3, 1, 1\}$

$d_0 = \{1, 1, 1, 2, 2, 2, 3, 3, 3, 1\}$

$\phi(p_7, c_1) = 43.75$
 $\phi(p_7, c_2) = 39.67$
 $\phi(p_7, c_3) = 29.33$

$\phi(p_8, c_1) = 16.75$
 $\phi(p_8, c_2) = 12.67$
 $\phi(p_8, c_3) = 2.33$

$\phi(p_9, c_1) = 17.25$
 $\phi(p_9, c_2) = 21.33$
 $\phi(p_9, c_3) = 31.67$

$\phi(p_{10}, c_1) = 5.25$
 $\phi(p_{10}, c_2) = 9.33$
 $\phi(p_{10}, c_3) = 19.67$

Therefore after iteration ①
 $d_1 = \{1, 3, 3, 1, 3, 1, 3, 3, 1, 1\}$

After 1st:-

Centroid-1 = $-1 + 2 + 3 + 0 + 12 / 5 = 3.2$

Centroid-2 = 21.33

Centroid-3 = $27 + 31 + 59 + 61 + 34 / 5 = 42.4$

→ start iteration 2

- 5) Let there be 10 data points: $\{-1, 27, 31, 2, 59, 3, 61, 34, 0, 12\}$. Use K-means clustering for 3 iterations to cluster them into 3 clusters. The initial labels are $\{1, 1, 1, 2, 2, 2, 3, 3, 3, 1\}$. List out the labels after 3 iterations.

☐ $\{1, 2, 2, 1, 3, 1, 3, 2, 1, 1\}$

☐ $\{1, 1, 2, 1, 3, 1, 3, 2, 1, 2\}$

☐ $\{1, 2, 3, 1, 3, 1, 3, 2, 3, 1\}$

☐ $\{1, 1, 2, 1, 3, 1, 3, 3, 1, 1\}$

$$D = \{-1, 27, 31, 2, 59, 3, 61, 34, 0, 12\}$$

$$x_1 = \{1, 3, 3, 1, 3, 1, 3, 3, 1, 1\}$$

$$c_1 = 3.2, c_2 = 21.33, c_3 = 42.4$$

$$\begin{aligned} \phi(p_1, c_1) &= 4.2 \\ \phi(p_1, c_2) &= 22.33 \\ \phi(p_1, c_3) &= 43.4 \end{aligned} \rightarrow p_1 \in W_1$$

$$\begin{aligned} \phi(p_2, c_1) &= 24.2 \\ \phi(p_2, c_2) &= 6.67 \\ \phi(p_2, c_3) &= 15.4 \end{aligned} \rightarrow p_2 \in W_2$$

$$\begin{aligned} \phi(p_3, c_1) &= 27.8 \\ \phi(p_3, c_2) &= 9.67 \\ \phi(p_3, c_3) &= 11.4 \end{aligned} \rightarrow p_3 \in W_2$$

$$R_2 = \{1, 2, 2, 1, 3, 1, 3, 3, 1, 1\}$$

$$\begin{aligned} \phi(p_8, c_1) &= 30.8 \\ \phi(p_8, c_2) &= 12.67 \\ \phi(p_8, c_3) &= 8.4 \end{aligned}$$

$$\begin{aligned} \text{Centroid-1} &= \frac{-1+2+3+0+12}{5} = 3.2 \\ \text{Centroid-2} &= \frac{27+31}{2} = 29 \\ \text{Centroid-3} &= \frac{59+61+34}{3} = 51.33 \end{aligned}$$

$$\begin{aligned} p_4 &\in W_1 & p_7 &\in W_3 \\ p_5 &\in W_3 & p_8 &\in W_3 \\ p_6 &\in W_1 & p_9 &\in W_1 \\ & & p_{10} &\in W_1 \end{aligned}$$

$$\begin{aligned} |12 - 3.2| &= 8.8 \\ |12 - 21.33| &= 9.33 \\ |12 - 42.4| &= 30.4 \end{aligned} \rightarrow p_{10} \in W_1$$

- 5) Let there be 10 data points: $\{-1, 27, 31, 2, 59, 3, 61, 34, 0, 12\}$. Use K-means clustering for 3 iterations to cluster them into 3 clusters. The initial labels are $\{1, 1, 1, 2, 2, 2, 3, 3, 3, 1\}$. List out the labels after 3 iterations.

☒ $\{1, 2, 2, 1, 3, 1, 3, 2, 1, 1\}$

☐ $\{1, 1, 2, 1, 3, 1, 3, 2, 1, 2\}$

☐ $\{1, 2, 3, 1, 3, 1, 3, 2, 3, 1\}$

☐ $\{1, 1, 2, 1, 3, 1, 3, 3, 1, 1\}$

$$D = \{-1, 27, 31, 2, 59, 3, 61, 34, 0, 12\}$$

$$L_2 = \{1, 2, 2, 1, 3, 1, 3, 3, 1, 1\}$$

$$C_1 = 3.2$$

$$C_2 = 29$$

$$C_3 = 51.33$$

Third iteration:-

$$P_1 \in W_1$$

$$P_9 \in W_1$$

$$P_2 \in W_2$$

$$P_{10} \in W_1$$

$$P_3 \in W_2$$

$$P_4 \in W_1$$

$$P_5 \in W_3$$

$$P_6 \in W_1$$

$$P_7 \in W_3$$

$$P_8 \in W_2$$

$$L_3 = \{1, 2, 2, 1, 3, 1, 3, 2, 1, 1\}$$

$$\text{Centroid-1} = \frac{-1 + 2 + 3 + 0 + 12}{5} = 3.2$$

$$\text{Centroid-2} = \frac{27 + 31 + 34}{3} = 30.66$$

$$\text{Centroid-3} = \frac{61 + 59}{2} = 60$$

$$\text{Centroid} = \{3.2, 30.66, 60\}$$

6) For the above question, find the centroids after 3 iterations.

$$D = \{-1, 27, 31, 2, 59, 3, 61, 34, 0, 12\}$$

$$\{5, 16.5, 56\}$$

$$\{3, 21.33, 42.4\}$$

$$\underline{3.2}, \underline{30.67}, \underline{60}$$

$$3.2, 29.51, 33$$

$$L_3 = \{1, 2, 2, 1, 3, 1, 3, 2, 1, 1\}$$

$$\text{Centroid-1} = \frac{-1 + 2 + 3 + 0 + 12}{5} = 3.2$$

$$\text{Centroid-2} = \frac{27 + 31 + 34}{3} = 30.66$$

$$\text{Centroid-3} = \frac{61 + 59}{2} = 60$$

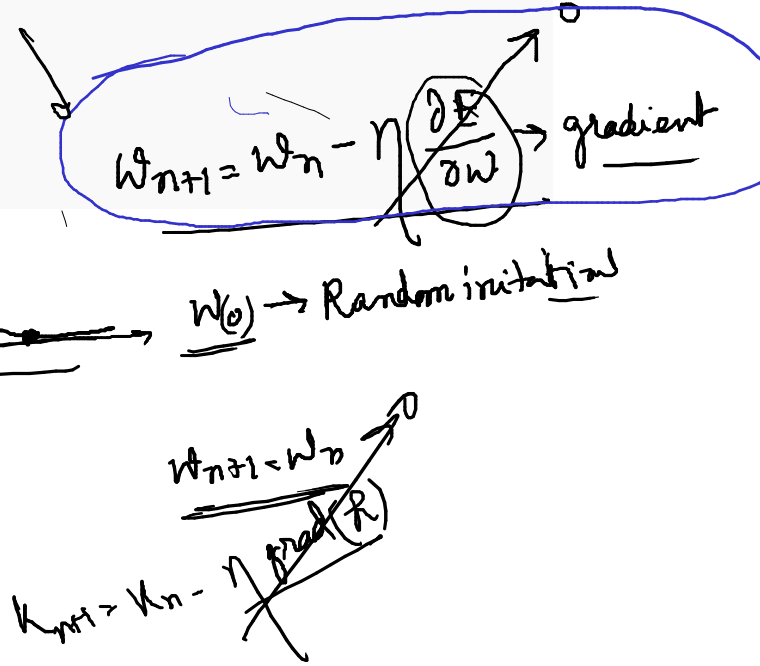
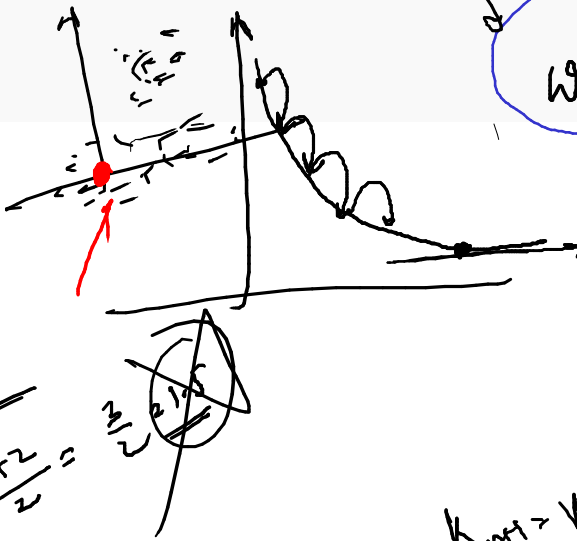
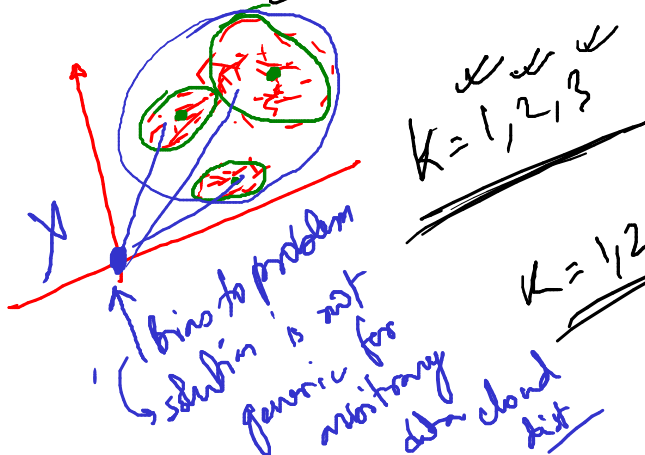
$$\text{Centroid} = \{3.2, 30.66, 60\}$$

7) How does K-Means initialize cluster centroids?

- At the maximum data point values
- Based on the mean of the data points
- At the origin (0,0)

L → Random

Randomly

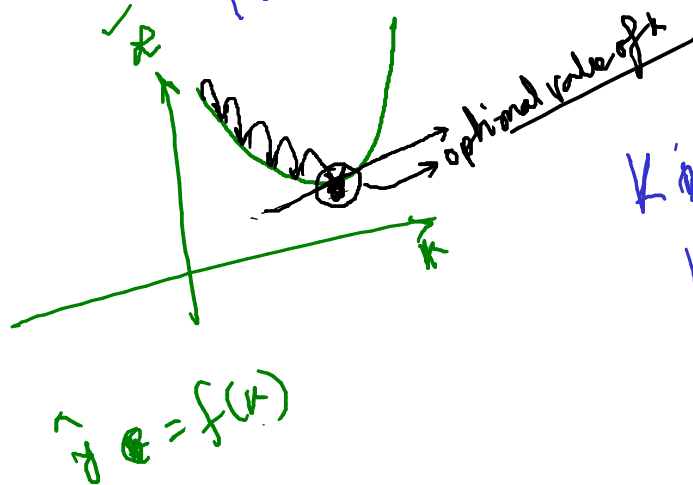


8) What is the best way to determine the optimal number of clusters (K) in K-Means?

- ☒ Trial and error
- ☒ Using the Elbow Method
- ☒ Setting K equal to the number of data points
- ☐ Using the Silhouette score

Sum of square error less.
$$d = \sqrt{(y - \hat{y})^2}$$

True output \rightarrow Target output from K-Means clustering.
Quadratic equation (Parabola)
Convex



100 groups of Pen

100 data points

Pen \rightarrow different atom different length

K may vary
 \rightarrow Problem specific
 \rightarrow Heuristically chosen variable



