

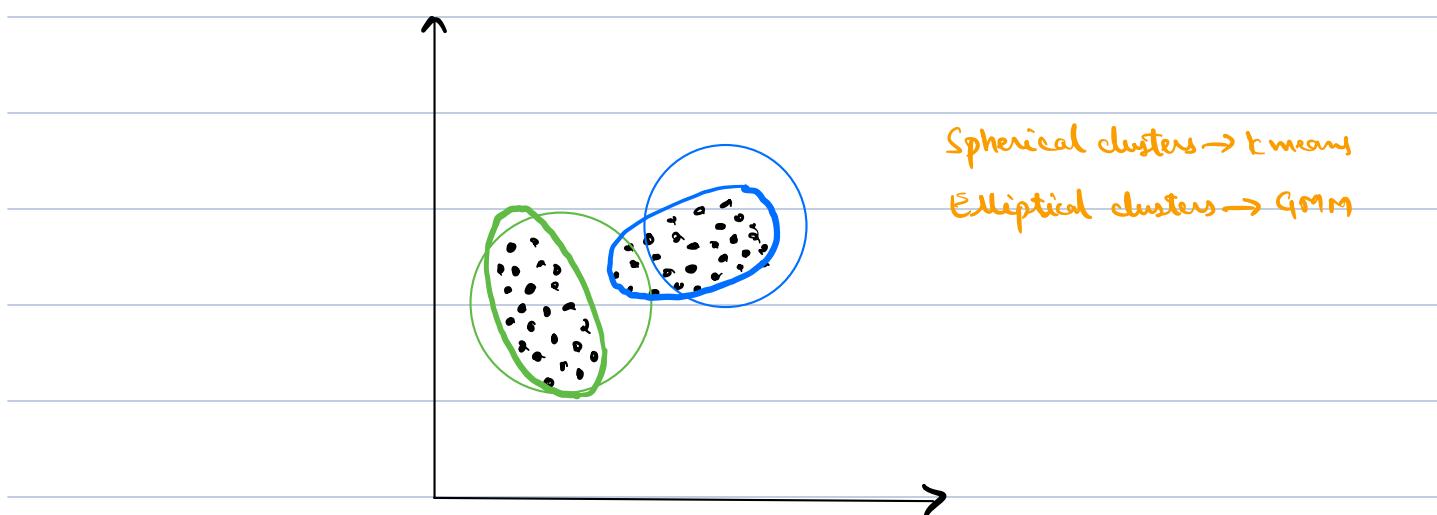
"Clustering using GMM"

'k-means'

- ① Initialization of centroids
 - ② Data points Assignment
 - ③ Re-calculation of centroids
 - ④ Till Convergence
- } Repeated

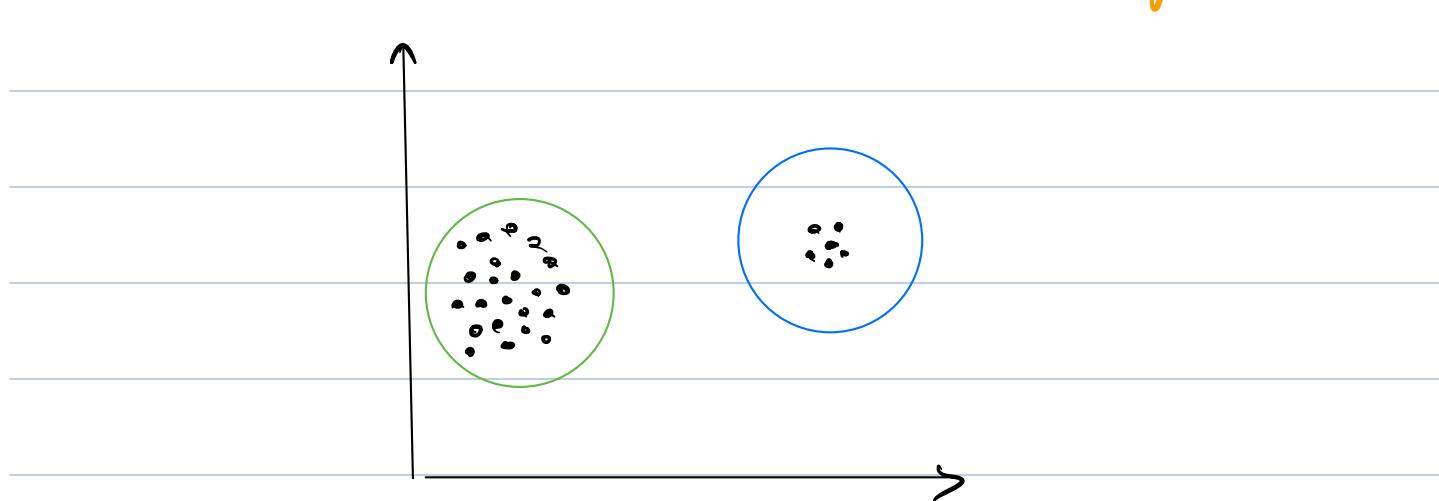
How to decide 'k' ?

- ① Domain knowledge
- ② Silhouette Score
- ③ Elbow Curve

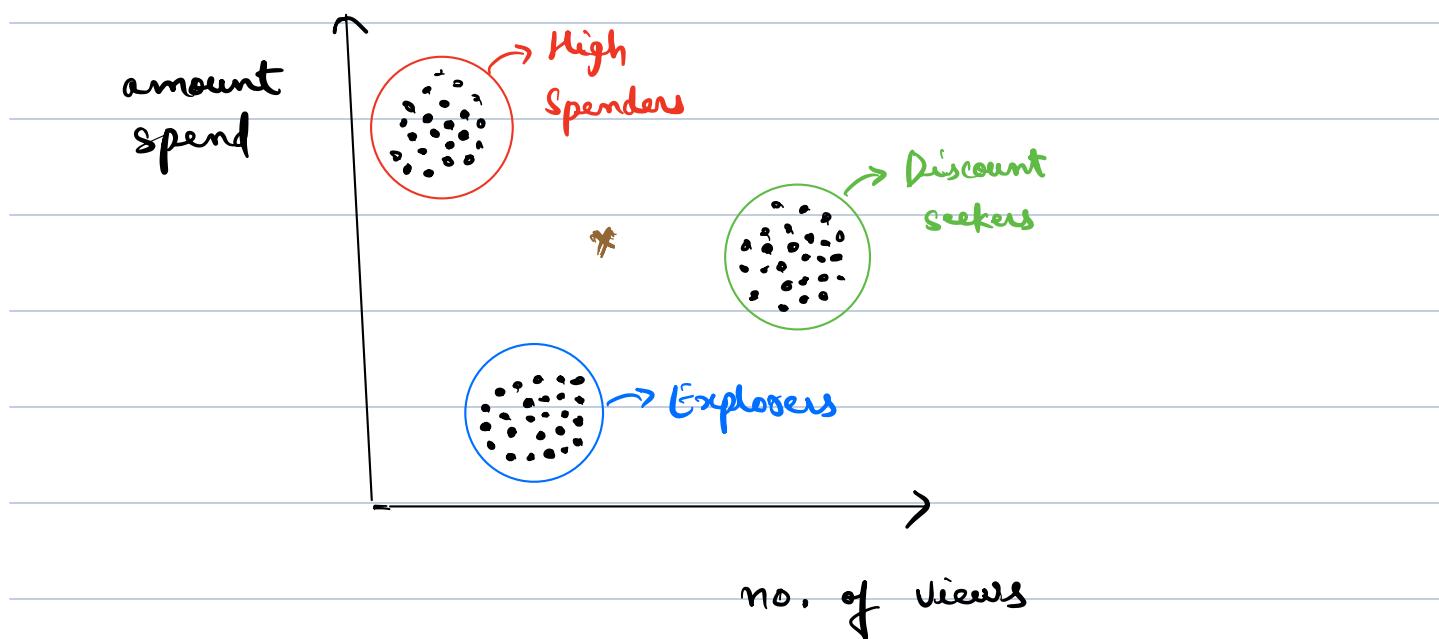


Limitations of k-means:

- ① Always gives spherical / globular clusters.
- ② Always gives almost equal sized clusters.
- ③ Hard Clustering : For each data point only 1 cluster is assigned.



"ecommerce data"



$x \rightarrow$ New customer

(Wealthy but also price-conscious)

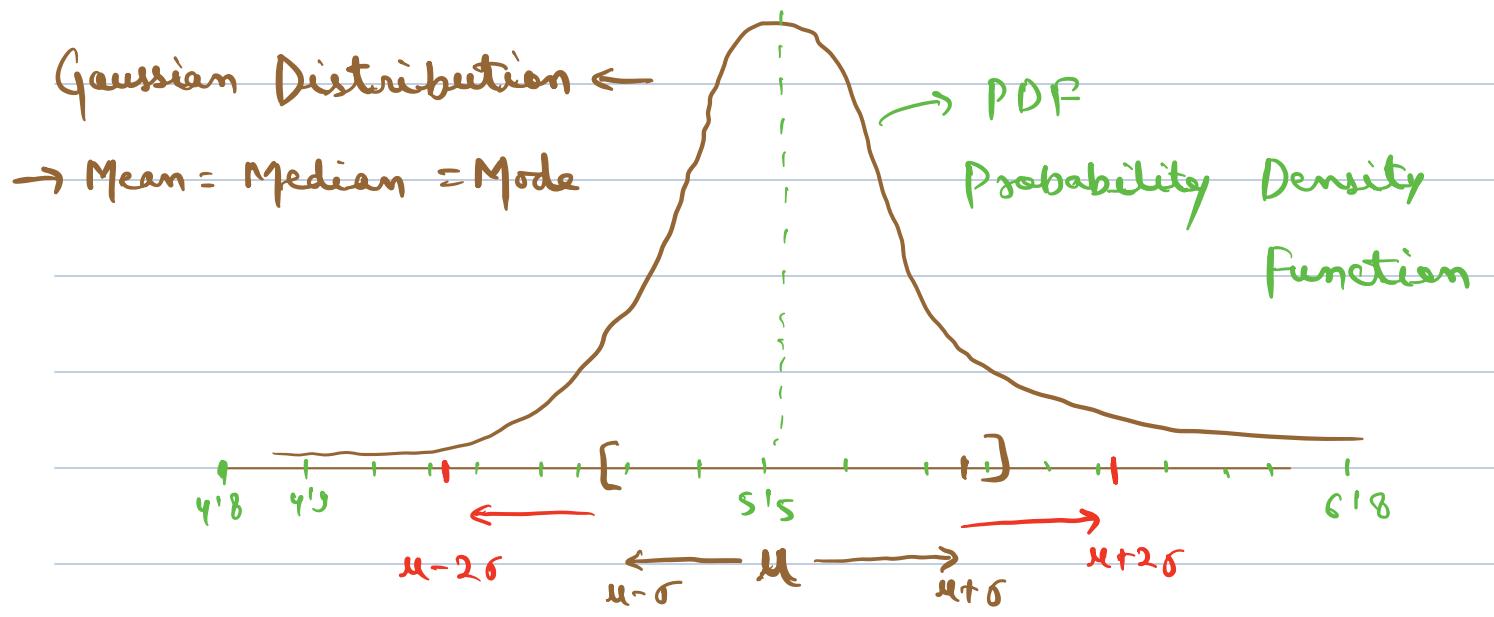
$\rightarrow 60\% \rightarrow$ High spenders

$\rightarrow 40\% \rightarrow$ Discount seekers

\hookrightarrow Calculate probabilities to see where a data point belongs

[Soft clustering]

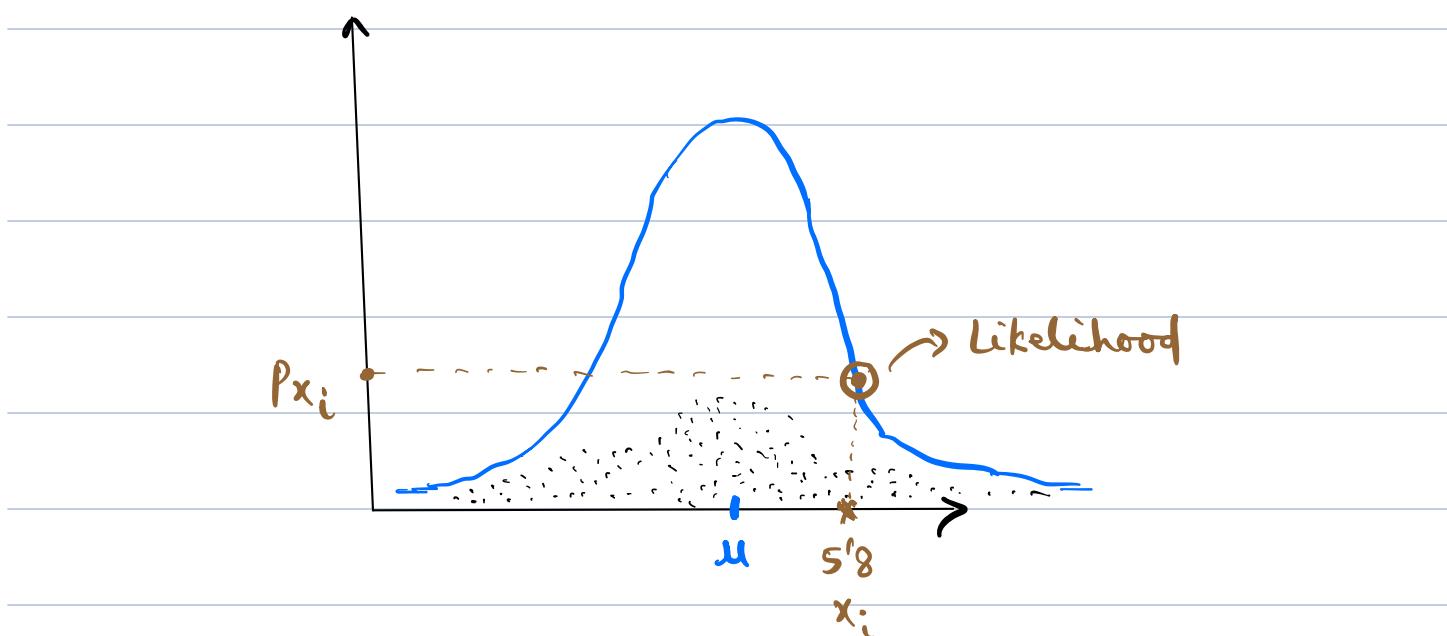
Example: We take a survey of heights of Indian citizens.



On average, height of Indian citizen is 5'5

Empirical formula :

- $\mu \pm \sigma \rightarrow 68\% \text{ of data } \checkmark$
- $\mu \pm 2\sigma \rightarrow 95\% \text{ of data } \checkmark$
- $\mu \pm 3\sigma \rightarrow 99\% \text{ of data}$



$$\text{PDF} = p(x; \mu, \sigma) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$



Data $\rightarrow \mu$ and σ

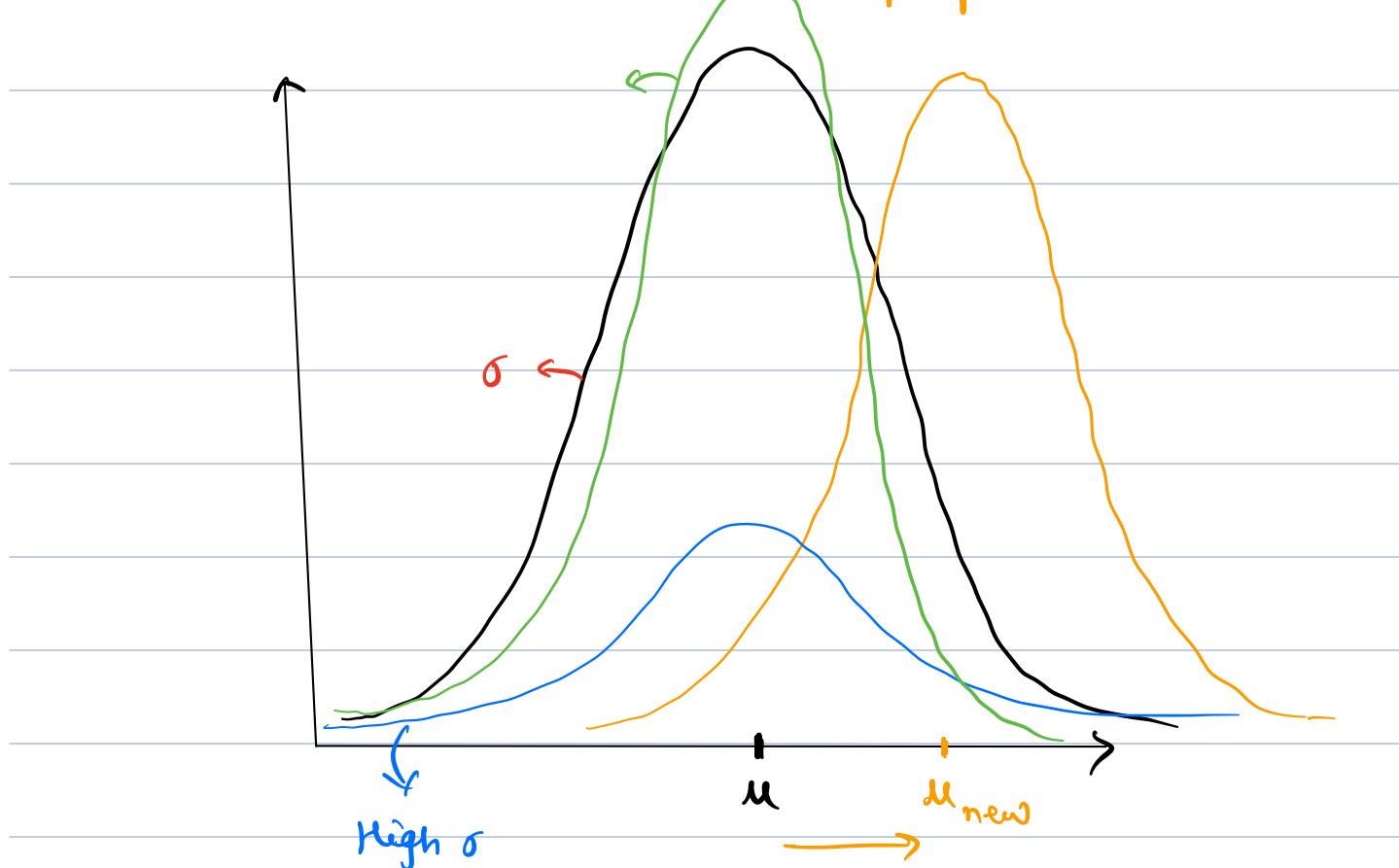
Data Point $\rightarrow x$

"1-D Gaussian"

• PARAMETERS

→ Mean : Position of the Distribution

→ σ → Standard Deviation : Shape of the Distribution



$$\text{PDF} = p(x; \mu, \sigma) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

2-D Gaussian !

↳ 2 Distributions

$\rightarrow \mu_1$

Position \rightarrow Means $\xrightarrow{\mu_1, \mu_2}$

Shape \rightarrow S.D of Distributions $\xrightarrow{\sigma_1, \sigma_2}$

* Covariance \rightarrow Orientation / Angle of the Ellipse

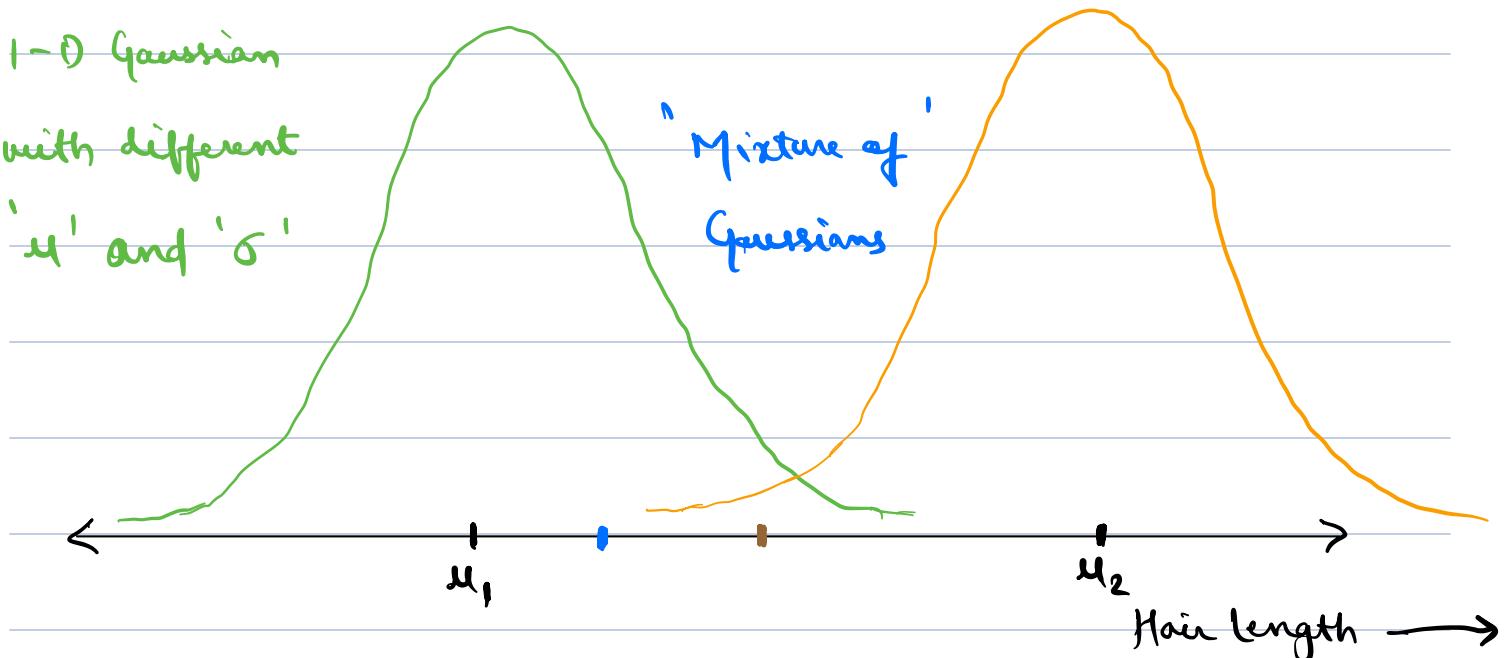
Ex:-

Survey was done \rightarrow to find hair length of all Indian citizens

1-D Gaussian

with different
' μ ' and ' σ '

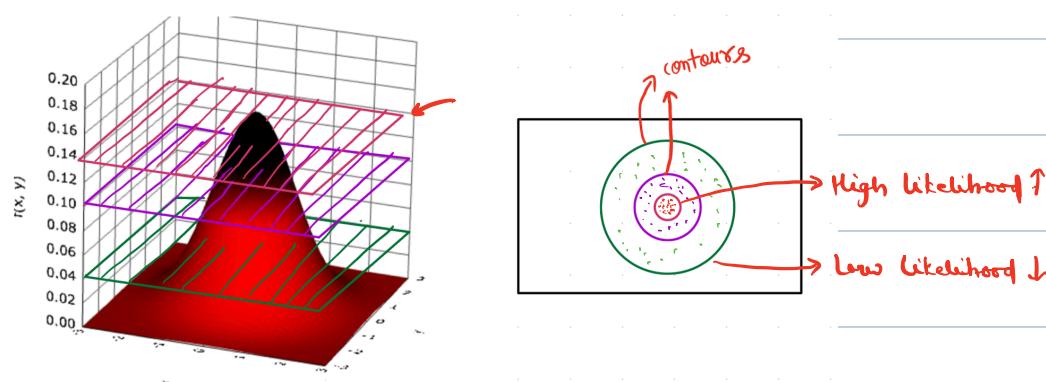
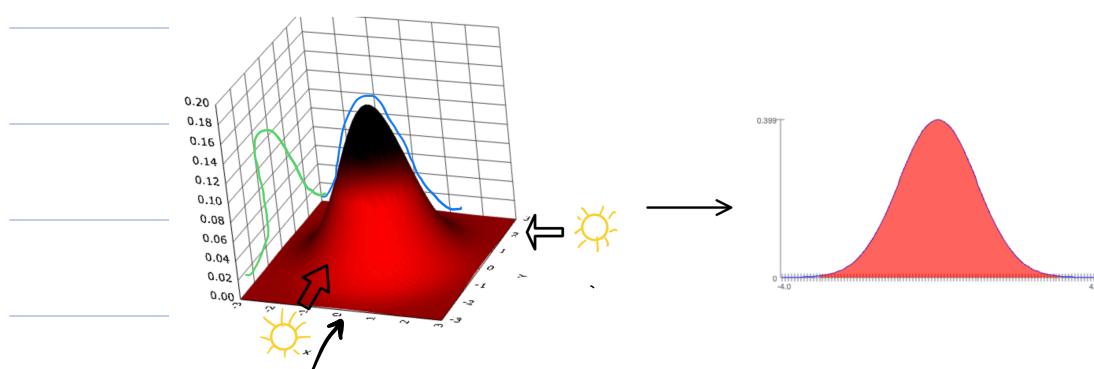
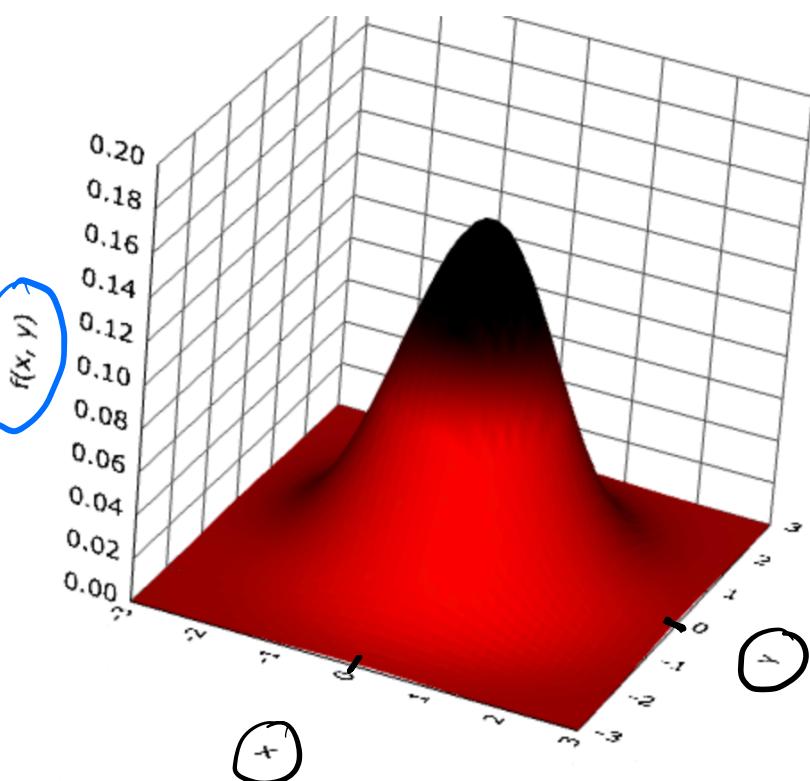
"Mixture of
Gaussians"



Uni-Modal \rightarrow 1 gaussian

Bi- Modal \rightarrow 2 gaussians

Multi- Modal \rightarrow Multi- gaussians



Steps of GMM:

- ① Initialize 'k' gaussians
- ② Expectation Step

↳ Calculate likelihood / probability

③ Maximization Step

↳ Update mean and standard deviation

↳ Weighted avg. approach

④ Tell the convergence.

Sample Calculation:

$$D = \{x_1, x_2, x_3, \dots, x_n\}$$

Step 1 : Choose 'k' random points as Mean.

Initialize
Means $\rightarrow \mu_1 = x_1$

$$\rightarrow \mu_2 = x_3$$

Initialize the variance

$$\sigma_1 = \frac{1}{n} \sum (x_i - \mu_1)^2$$

$$\sigma_2 = \frac{1}{n} \sum (x_i - \mu_2)^2$$

Step-2 : Expectation Step

$$p(x, \mu_1)$$

$$p(x, \mu_2)$$

$$p_1(x, \mu_1)$$

$$p_n(x, \mu_2)$$

x_1	$0.1 / 0.4 \Rightarrow 0.25$	$0.3 / 0.4 \Rightarrow 0.75$	0.25	0.25
x_2	0.1	0.1	0.5	0.5
x_3	0.2	0.4	0.33	0.67
x_4	\vdots	\vdots	\vdots	\vdots
x_5	\vdots	\vdots	\vdots	\vdots
\vdots	\vdots	\vdots	\vdots	\vdots

Normalization of Probabilities

↳ so that sum of all prob. is equal to 1

③ Maximization

likelihood of data point in u_i ,

$$\sum_{i=1}^n p(x, u_i) \cdot x_i \rightarrow \text{data point}$$

$$\sum_{i=1}^n p(x, u_i) \rightarrow$$

Weighted average

$$x_1 = 10$$

$$x_2 = 20$$

$$x_3 = 30$$

$$u_1 \text{ new} = \frac{10 \times 0.25 + 20 \times 0.5 + 30 \times 0.33}{0.25 + 0.5 + 0.33}$$

$$u_2 \text{ new} =$$

$$\text{Variance} \rightarrow \sigma^2$$

Variance

$$\sigma^2 = \frac{\sum p(\mu_{x_i}) \overbrace{(x_i - \mu)}^2}{\sum p(\mu_{x_i})}$$

} → Weighted average of likelihood

GMM

Repeats Step 2 and Step 3 till Convergence.

k-means

- ① Hard Clustering
→ Distance

GMM

- ① Soft Clustering
→ Probability

- ② Spherical Clusters
- ③ Equal sized Clusters
- ④ Use Centroids

- ② Elliptical Clusters
- ③ Size can vary
- ④ Use Mean and Variance