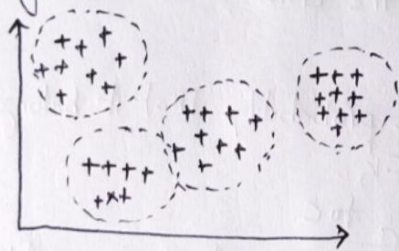


GAUSSIAN MIXTURE MODEL (GMM)

- Gaussian mixture model are a probabilistic model for representing normal distributed subpopulation within an overall population.
- Mixture models in general don't require which subpopulation a data point belongs to, allowing the model to learn the subpopulation automatically.
- For example, in modeling human being data, height is typically modelled as normal distribution for each gender with a mean approximately 5'10" for males and 5'5" for females. Given only height data and not gender assignment for each data point, distribution of all height would follow the sum of two scaled (different variance) normal distribution.
- A model making this assumption, the parameters follow normal distribution.
- GMM assume there are a certain number of Gaussian distributions and each of these distribution represent a cluster. Hence a Gaussian Mixture model tends to group the data points belong to a single distribution together.

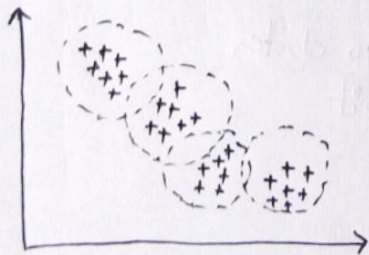
Motivating GMM: Weakness of K Means.



- In this case, data is well separated data.
- K means algorithm can quickly label those clusters.
- There might be slightly overlap between two clusters' assignment of point between them.

- Unfortunately, K means has no measure of probability.

→ K means model is that it places a circle at the center of each cluster, the radius acts as a hard cutoff for cluster assignment, any point outside the circle is not considered a member of cluster.



→ By eye, we can recognized that transformed clusters are non-circular thus circular clusters would be poor fit.

→ K means try to force-fit data into four circular clusters, resulting circles overlap.

→ GMM address this issue, since it is probabilistic model, it find possible probabilistic cluster assignment.

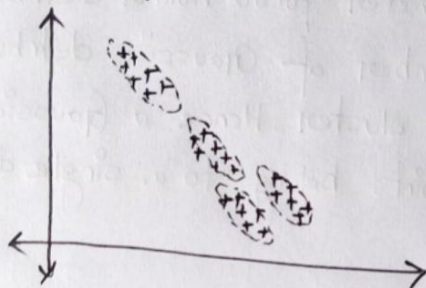
What is Expectation - Maximization? (EM)

→ EM is a statistical algorithm for finding the right model parameters. We typically use EM when the data has missing values or in other words when the data is incomplete.

Broadly, EM algorithm has two steps:

- i) E-step → In this step, available data is used to estimate (guess) the value of the missing variables.
- ii) M-step → Based on the estimated values generated in E-step the complete data is used to update the parameters.

After applying GMM Clustering



- Clusters are organised diagonally.
- Gaussian mixture model worked better in this case (non-spherical data)

→ k means only considers mean to update the centroid while GMM takes into account the mean as well as variance of the data.

Normal EM follows the steps -

E-step → For each point x_i , calculate the probability that it belongs to cluster divide by distribution c_1, c_2, \dots, c_k .

$$E\text{-step} = \frac{\text{Probability } x_i \text{ belong to } c}{\text{Sum of probability } x_i \text{ belong to } c_1, c_2, \dots, c_k}$$

Value will be high when point is assigned to right cluster & lower otherwise.

M step →

New density is defined by ratio of the number of points in the cluster & total number of points.

$$= \frac{\text{No. of points assigned to cluster}}{\text{Total number of point}}$$