

Gaussian Mixture Models (GMMs)

Setting: unsupervised learning.

$$\{x^{(1)}, x^{(2)}, \dots, x^{(n)}\}$$

No ground-truth labels available.

Recap of k-means.

clustering approach \Rightarrow goal \rightarrow group similar data points together (in a cluster)

$$\min \sum_{i=1}^n \|x_i - \mu_i\|^2$$

\downarrow datapoint \rightarrow centroid (mean).

outcome was x_i belongs to cluster k
 \Downarrow resulted in hard clustering.

And also, no model assumption was made.
all we did was to try and minimize the distance between x_i and centroid (μ_i) of the cluster to which x_i belongs to.

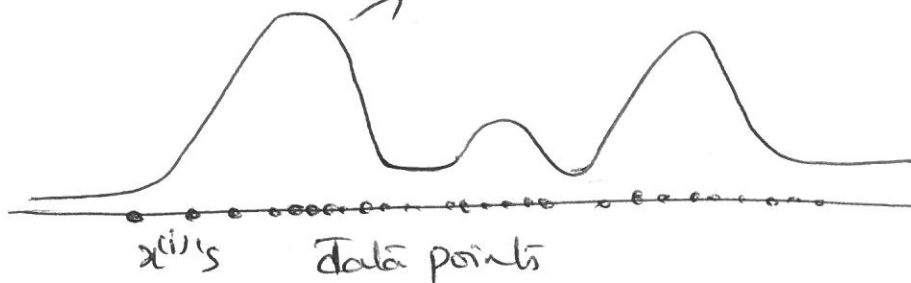
What if we want soft clustering??

Can we come up with a probability based approach?
and also can we assume a model from which data points $x^{(i)}$ can be drawn.

Graphical representation

Goal \Rightarrow find this # model distribution.

Fig. 1



Lets make some assumptions.

Based on ~~the~~ Fig 1 we view $x^{(1)}, x^{(2)} \dots x^{(n)}$ as random variables drawn from an unknown distribution with density $p(x)$

Such that
$$p(x) = \sum_{j=1}^K \pi_j \underbrace{w_j(x; \mu_j, \Sigma_j)}_{\text{Component}}$$

some distribution

~~However~~ coming up

For GMM, this distribution $p(x)$ has a mixture of K components

Each one is a multivariate Gaussian distribution.

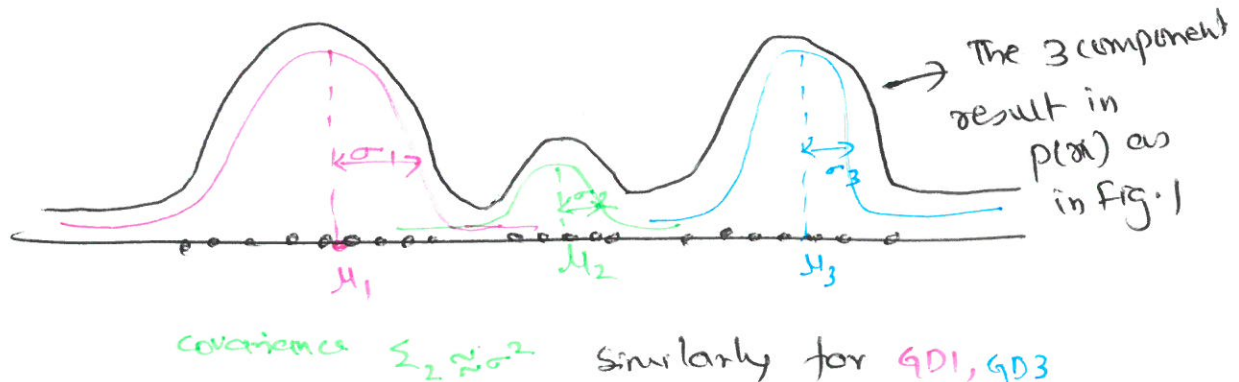
$\rightarrow \mu(\text{normal})$
 $w(x; \mu_j, \Sigma_j) \rightarrow$ Gaussian distribution with unknown parameters (μ_j, Σ_j)

$\pi_j \rightarrow$ unknown probability of selecting component j

Such that $\sum_{j=1}^K \pi_j = 1$ (Total probability should add up to 1).

Graphical rep.

Fig 2



Let us now consider a latent component indicator $z^{(i)}$
 (hidden/unobserved) (latent class/
 cluster for datapoint x_i)

See z_i 's as labels in supervised learning setting.

But in our case these z_i 's are not available or we don't
 (an supervised setting) get to see them.

$z^{(i)} \sim \text{multinomial}(w_k) \rightarrow \{w_1, w_2, \dots, w_k\}$ k components
 \downarrow in binary classification $z^{(i)}$ takes 1 or 0.

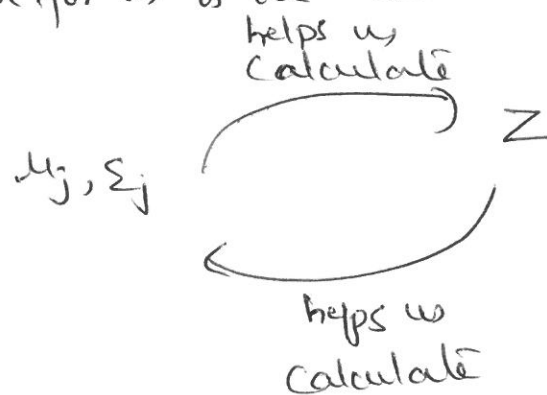
Here it takes probability of datapoint $z^{(i)}$
 coming from cluster a specific

~~$x^{(i)}$~~ $x^{(i)} | z^{(i)} \approx \mathcal{N}(\mu_{z^{(i)}}, \Sigma_{z^{(i)}})$
 \downarrow
 $z^{(i)}$ conditioned to $z^{(i)}$ \downarrow Gaussian/normal distribution
 \downarrow prob. of $x^{(i)}$ being observed in cluster $z^{(i)}$

Now if we knew $z^{(i)}$, it becomes supervised problem
 we can go ahead and calculate joint distribution

$$p(x^{(i)}, z^{(i)}) = p(x^{(i)} | z^{(i)}) p(z^{(i)})$$

The problem for us is we do not know $z^{(i)}, \mu_j, \Sigma_j$



So how do we solve this??

Estimation - Maximization algorithm (EM)

Similar to what we did in k-means

- ① estimate centroids (random select)
- ② compute new centroids based on data point assignments.

In EM we have 2 steps

Initialize w, μ, Σ randomly.

E Step :- compute values of $z^{(i)}$

for each i, j set (training & component set)

$$w_j^{(i)} = p(z^{(i)} = j | x^{(i)}; w, \mu, \Sigma)$$

gives the class memberships (soft allocations)

gives the prob. of $x^{(i)}$ observed by Gaussian (μ_j, Σ_j)

$$= \frac{p(x^{(i)} | z^{(i)} = j; \mu_j, \Sigma_j) p(z^{(i)} = j; w)}{\sum_{l=1}^K p(x^{(i)} | z^{(i)} = l; \mu_l, \Sigma_l) p(z^{(i)} = l; w)}$$

probability of Gaussian (w_j) (prior)

Total probability over all components.

M-step :- update parameters based on Gaussian.

$$w_j := \frac{1}{n} \sum_{i=1}^n w_j^{(i)} ; \mu_j := \frac{\sum_{i=1}^n w_j^{(i)} x^{(i)}}{\sum_{i=1}^n w_j^{(i)}}$$

updated weights

no. of samples

update means

$$\Sigma_j := \frac{\sum_{i=1}^n w_j^{(i)} (x^{(i)} - \mu_j)(x^{(i)} - \mu_j)^T}{\sum_{i=1}^n w_j^{(i)}}$$

update covariance

Repeat.

until no new assignments of data points to k components.