

Gaussian Mixture Model

clustering

(soft)
hard

mixtures are a PROBABILISTICALLY GROUNDED way of soft clustering.

each cluster is a gaussian model.

A gaussian mixture is a LINEAR SUPERPOSITION OF MULTIPLE GAUSSIANS
(we can assume a weight for each gaussian)

caz unimodal dist. not always practical. ← Say we're k gaussians w/ a weight π_k .

$$P(x) = \int \sum \pi_k N(x|\mu_k, \Sigma_k) \leftarrow p(x) = \sum_{k=1}^k \pi_k \cdot N(x|\mu_k, \Sigma_k)$$

= $\sum \pi_k \underbrace{N(x|\mu_k, \Sigma_k)}_{\text{PDF}}$ when $0 \leq \pi_k \leq 1$
= $(1) \cdot (1)$
= 1 $\sum \pi_k = 1$

if we take log likelihood:

EM works instead!

$$\ln p(x|\mu, \Sigma, \pi) = \sum \ln p(x_n) = \sum_{n=1}^N \ln \left[\sum_{k=1}^k \pi_k \cdot N(x_n|\mu_k, \Sigma_k) \right] \quad \times$$

EM for GMM, where each G is 1D.

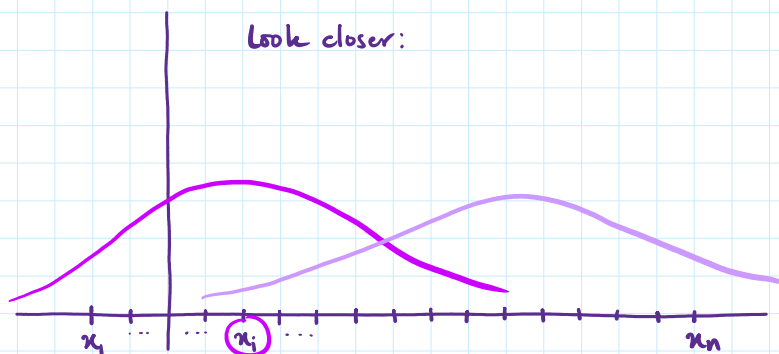
① Init step: π_j, μ_j, σ_j^2 for each gaussian j .
weight → defining characteristics of gaussian.

No closed form solution.
⇒ MLE will not work.

② Σ step: Calculate posterior distribution of latent variable
(probability that this gaussian generated this data)

repeat till convergence.

③ M step: Update parameter for each gaussian — π_j, μ_j, σ_j^2



x_1, x_2, \dots, x_n (observations)
 $k=2$ (2 gaussians)

If the source of x_i is known, estimating parameters is easy.
What if the source isn't known?
Then it's a latent variable!

Let the two gaussians be (μ_a, σ_a) and (μ_b, σ_b)
we'll look@ b but it's the same for a .

$$P(x_i|b) = \frac{1}{\sqrt{2\pi}\sigma_b} \exp\left(-\frac{(x_i - \mu_b)^2}{2\sigma_b^2}\right)$$

We don't know the latent var. ← $P(b|x_i) = \frac{P(x_i|b) \cdot P(b)}{P(x_i|b) \cdot P(b) + P(x_i|a) \cdot P(a)}$
⇒ we need (μ_b, σ_b^2) to get

So we take an EM based approach.

① Init step: Start w/ 2 random placings of gaussians — (μ_a, σ_a^2) and (μ_b, σ_b^2)

so we take an EM based approach.

- ① Init step: Start w/ 2 random placings of gaussians — (μ_a, σ_a^2) and (μ_b, σ_b^2)
- ② E-step: $p(b|x_i)$ — from b^0
 $p(a|x_i)$ — from a^0 } generate new dataset using soft weight.
- ③ M-step: adjust $(\mu_a, \sigma_a^2) + (\mu_b, \sigma_b^2)$ till it fits

do until convergence

Zoom in.
 Consider the joint probability $p(X = \bar{x}, Z = k)$
 $p(\bar{x}, k)$ — probability of input vector \bar{x} occurring together w/ class k .

Using bayes, $P(\bar{x}, k) = p(\bar{x}|k)p(k)$

PDF of k^{th} gaussian component.

$$p(\bar{x}|k) = N(\bar{x}|\mu_k, \Sigma_k) \quad k \in [1, K]$$

$$\text{Let } p(k) = \pi_k \quad k \in [1, K]$$

$\pi_k \approx \frac{N_k}{N}$ → no. of instances belonging to k
 total no. of instances

$$\text{So } p(\bar{x}, k) = p(k)p(\bar{x}|k) = \pi_k \cdot N(\bar{x}; \mu_k, \Sigma_k)$$

So, from this, we get marginal probability of \bar{x} :

$$p(\bar{x}) = \sum p(\bar{x}, k) = \sum \pi_k N(\bar{x}; \mu_k, \Sigma_k)$$

linear superposition of multiple gaussians.

What do we take from this?

① GMM = weighted sum of K gaussian components.

② π_k is the selection probability of the k^{th} gaussian

③ Each of the K gaussians models a single class.

↓
 k^{th} gaussian $N(\bar{x}; \mu_k, \Sigma_k)$ can be modeled as condⁿ prob $p(\bar{x}|k)$
 probability of \bar{x} being in k^{th} sub class.

④ The sum of all joint subclass probabilities = marginal probability $p(\bar{x})$ of data value (\bar{x})

Classification via GMM

Set of unlabeled input X is given

fit a GMM for all K classes.

given an arbitrary \bar{x} , compare $p(k|\bar{x})$

probability that it's class k given \bar{x} .

Value of k yielding highest $p(k|\bar{x})$ is the class corresponding to \bar{x} .

$$p(k|\bar{x}) = \frac{p(\bar{x}, k)}{p(\bar{x})}$$

$$= \frac{\pi_k N(\bar{x}; \mu_k, \Sigma_k)}{\sum_{i=1}^K \pi_i N(\bar{x}; \mu_i, \Sigma_i)}$$

Whole process is super easy if we know GMM params.

given \bar{x} , just classify by] — that yields max value of $p(Z=k|X=\bar{x})$

given x , just directly by assigning it to cluster k — that gives the value of $p(z=k|x=x)$

MLE of GMM Parameters.

A GMM is fully described by parameter set:

$$\theta = \{\pi_k, \mu_k, \Sigma_k \mid k \in \{1 \dots K\}\}$$

π_k : weight of k^{th} gaussian
 μ_k : mean of k^{th} gaussian
 Σ_k : variance of k^{th} gaussian

* we're not estimating k for the MLE.

Now, when we had a simple gaussian:

* a closed form soln has unknowns on LHS and knowns on RHS.

- ① we took joint log likelihood of all training data given by gaussian PDF.
- ② took gradient on both sides & set it to 0.

But that doesn't work w/ the GMM — we get an equation that has NO closed form soln

So we go for an ITERATIVE APPROXIMATION:

Let $x^{(1)}, x^{(2)}, x^{(3)} \dots x^{(n)}$ be the set of observed data points.

$$p(x^{(i)}|\theta) = \sum_{k=1}^K \pi_k \cdot N(x^{(i)}; \mu_k, \Sigma_k)$$

likelihood
→ we'll refer to it as $p(x^{(i)})$

All datapoints $(x^{(i)})$ are independent.

$$\Rightarrow \text{Joint probability} = p(x^{(1)}) \cdot p(x^{(2)}) \cdot p(x^{(3)}) \dots p(x^{(n)})$$

$$\begin{aligned} \Rightarrow \text{Joint loglikelihood} &= \log(p(x^{(1)}) \cdot p(x^{(2)}) \cdot p(x^{(3)}) \dots p(x^{(n)})) \\ &= \sum_{i=1}^n \log(p(x^{(i)})) \\ &= \sum_{i=1}^n \log\left[\sum_{k=1}^K \pi_k \cdot N(x^{(i)}; \mu_k, \Sigma_k)\right] \end{aligned}$$

logarithm of a sum!

⇒ We've to identify $\mu_1, \Sigma_1, \mu_2, \Sigma_2 \dots$ that will maximize the log joint likelihood. — take gradient wrt param & solve.

Say we take μ_1 for gradient:

$$\nabla_{\mu_1} \log(p(x^{(1)}) \cdot p(x^{(2)}) \dots p(x^{(n)})) = 0$$

$$\nabla_{\mu_1} \sum_{i=1}^n \log(p(x^{(i)})) = 0$$

$$\nabla_{\mu_1} \sum_{i=1}^n \log\left[\sum_{k=1}^K \pi_k \cdot N(x^{(i)}; \mu_k, \Sigma_k)\right] = 0$$

grad is a linear operator so move it in.

$$\sum_{i=1}^n \nabla_{\mu_1} \log\left[\sum_{k=1}^K \pi_k \cdot N(x^{(i)}; \mu_k, \Sigma_k)\right]$$

$$\frac{d}{dx} \log f(x) = \frac{1}{f(x)} \frac{\partial f}{\partial x}$$

$$\sum_{i=1}^n \frac{\nabla_{\mu_1} \sum_{k=1}^K \pi_k \cdot N(x^{(i)}; \mu_k, \Sigma_k)}{\sum_{k=1}^K \pi_k \cdot N(x^{(i)}; \mu_k, \Sigma_k)} = 0$$

$$\frac{d}{dx} \log f(x) = \frac{1}{f(x)} \frac{df}{dx} \rightarrow \sum_{i=1}^K \frac{V_{\mu_i} \sum_{k=1}^K \pi_k \cdot N(x^{(i)}; \mu_k, \Sigma_k)}{\sum_{k=1}^K \pi_k \cdot N(x^{(i)}; \mu_k, \Sigma_k)} = 0$$

If x_1, x_2 are independent, $\partial x_2 / \partial x_1 = 0$.
 So the gradient wrt μ_i will be 0 for all $k \neq i$.