

Mixture densities & EM algorithm.

In general, densities can be multi-modal.

A single parametric density model might not be a good choice for class-conditional densities.

Solution: Use convex-combination of multiple parametric densities.

One example: Gaussian mixture models.

$$f(x_i) = \sum_{j=1}^K \alpha_j f_j(x_i)$$

$$\text{Each } f_j \sim \mathcal{N}(\mu_j, \sigma_j).$$

$$l(\alpha, \mu, \sigma) = \sum_{i=1}^n \log \sum_{j=1}^K \alpha_j f_j(x_i).$$

Likelihood cannot be maximized directly
(we have log of sum).

Wish to model these kinds of densities using another latent variable & specifying the joint distribution of the data & the latent variable.

specifically,

Define a new 'un-observed' RV z

s.t., $z_i \sim \text{Multinomial}(\alpha) \in \{1, 2, \dots, k\}$.

$$\alpha_j \geq 0 \quad \sum_{j=1}^k \alpha_j = 1.$$

$$\alpha_j = P(z_i = j).$$

Model: Each x_i was generated by 'selecting' one of k values for z_i & ~~z_i~~ drawing x_i from that Gaussian.

Now, we have

$$l(\alpha, \mu, \sigma) = \sum_{i=1}^n \log p(x_i; \alpha, \mu, \sigma).$$

$$= \sum_{i=1}^n \log \sum_{z_i=1}^K p(x_i | z_i; \mu, \sigma) p(z_i; \alpha)$$

Now suppose z_i 's were known, then optimizing the above would have been easy.

$\therefore z_i$ was fixed for every i .

$$\text{Thus, } l = \sum_{i=1}^n \log \left(p(x_i | z_i; \sigma, \mu) \right) + \log p(z_i; \alpha)$$

$$\text{Now, } \frac{\partial l}{\partial \alpha_j} = \frac{1}{n} \sum_{i=1}^n I_{[z_i=j]}$$

$$\frac{\partial l}{\partial \mu_j} = \frac{\sum_{i=1}^n I_{[z_i=j]} x_i}{\sum_{i=1}^n I_{[z_i=j]}}$$

$I \sim$
Indicator
of $z_i=j$

(57).

$$\sigma_j = \frac{\sum_{i=1}^m I[z_i=j] (x_i - \mu_j)^2}{\sum_{i=1}^n I[z_i=j]}$$

Basically, the estimates are similar to MLE for a ^{single} gaussian but they are weighted by the proportion of occurrence of each component.

But the problem: We do not have any information about z_i .

Question: How to estimate MLE for GMMs in that case?

Answer: use an iterative algo:

- i) Guess a value for z_i
- ii) update the parameters with that guess
- iii) Repeat (i) & (ii) until convergence.

(58)

Procedure: For every x_i , calculate

$$i) \quad w_{ij}^t = P[z_i = j \mid x_i; \alpha^t, \mu^t, \sigma^t]$$

(Guess for z_i)

$$ii) \quad \alpha_j^{t+1} = \frac{1}{n} \sum_{i=1}^n w_{ij}^t$$

$$\mu_j^{t+1} = \frac{\sum_{i=1}^n w_{ij}^t x_i}{\sum_{i=1}^n w_{ij}^t}$$

$$\sigma_j^{t+1} = \frac{\sum_{i=1}^n w_{ij}^t (x_i - \mu_j^t)^2}{\sum_{i=1}^n w_{ij}^t}$$

[update the parameters].

How to find w_{ij}^t ?

Answer: use Bayes rule.

$$w_{ij} = P(z_i = j | x_i; \alpha^t, \mu^t, \sigma^t)$$

$$= \frac{P(x_i | z_i = j; \mu^t, \sigma^t) P(z_i = j; \alpha^t)}{\sum_{l=1}^K P(x_i | z_i = l) P(z_i = l; \alpha^t)}$$

We can compute all the above terms.

Thus, EM algorithm looks like giving a 'weighted' averages with 'weights' being 'soft' in the sense that they are probabilistic.

[This accounts for the fact that z_i 's are only guesses]

Question: Why should this procedure work?

The Generalized EM algorithm.

The general setting.

$D = \{x_1, x_2, \dots, x_n\}$ are sampled IID

from a GMM. [seen data]

z_i indicates the component from which x_i was drawn. [latent/unseen/hidden variable]

Goal: Find a procedure to obtain the MLE for such a model.

$$l(\theta) = \sum_{i=1}^n \log f(x_i; \theta).$$

$$= \sum_{i=1}^n \log \sum_z f(x_i, z; \theta).$$

General strategy of EM:
repeatedly.

Construct a lower bound on $l(\theta)$ &
repeatedly optimize the lower bound.

$\forall i$, Let ~~Φ_i~~ ϕ_i be a distribution over Z_i .

$$\therefore \sum_i \phi_i(z) = 1, \phi_i(z) \geq 0.$$

$$l = \sum_i \log f(x_i; \theta)$$

$$= \sum_i \log \sum_{z_i} f(x_i, z_i; \theta)$$

$$= \sum_i \log \sum_{z_i} \phi_i(z_i) \frac{f(x_i, z_i; \theta)}{\phi_i(z_i)} \quad - (2).$$

(62).

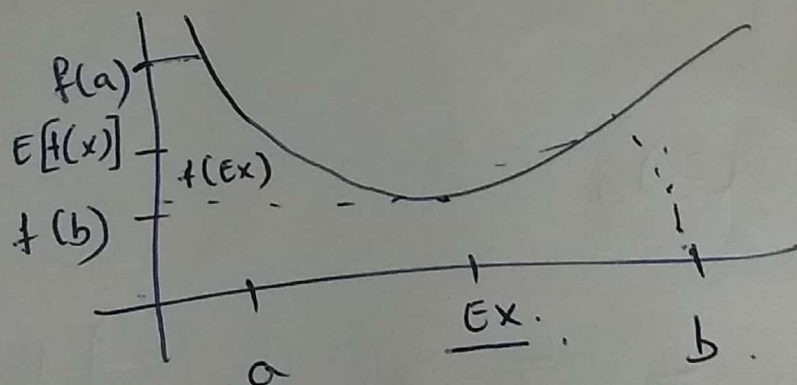
Jensen's inequality :

convex .

f be a real-valued function
i.e, $f''(x) \geq 0 \quad \forall x \in \mathbb{R}$.

If x is a RV, then

$$E[f(x)] \geq f(Ex).$$



$$x \in \{a, b\}$$

$$P(x=a) = \frac{1}{2}$$

For concave functions, $E[f(x)] \leq f(Ex)$.

In (2), since \log is a concave fn,

$$\sum_{z_i} \phi_i(z_i) \left[\frac{p(x_i, z_i; \theta)}{\phi_i(z_i)} \right] =$$

(A)

$$E_{z_i \sim \phi_i}(A)$$

[Hence the name EM].

$$\text{Thus, } p(E_{z_i}(A)) \geq E_{z_i \sim \phi_i} [A]$$

Thus in (2),

$$\ell \geq \sum_i \sum_{z_i} \phi_i(z_i) \log \frac{p(x_i, z_i; \theta)}{\phi_i(z_i)}.$$

Any ϕ_i would give a lower bound on \mathcal{L} . But we want a tight one

Jensen's inequality would be tight if $E[x]$ is a constant.

in our case,

$$\frac{f(x_i, z_i)}{\phi_i(z_i)} = c. \quad [\text{not depending upon } z_i].$$

$$\Rightarrow \phi_i(z_i) \propto f(x_i, z_i; \theta).$$

$$\text{Also, } \sum_z \phi_i(z_i) = 1 \quad \because \text{ it is a distribution}$$

$$\text{Thus, } \phi_i(z_i) = \frac{f(x_i, z_i; \theta)}{\sum_z f(x_i, z; \theta)}$$

(65).

$$= \frac{f(x_i, z_i; \theta)}{f(x_i; \theta)}$$

$$= \underline{f(z_i | x_i; \theta)}.$$

$\phi_i(z_i)$ is posterior of z_i given x_i

[~~is~~ same as our previous guess].

Now, EM:

for each i , set

$$\phi_i(z_i) = P[z_i | x_i; \theta^t]. \quad [\text{E-step}]$$

$$\text{M-step: } \theta^{t+1} = \underset{\theta}{\operatorname{argmax}} \sum_i \sum_{z_i} \phi_i(z_i) \frac{\log f(x_i, z_i; \theta^{t+1})}{\phi_i(z_i)}.$$

(6.6).

Question: Does this procedure converge?

Consider: E-step. Ensures that

$$l(\theta^t) = \sum_i \sum_{z_i} \underbrace{\phi_i^t(z_i) \log f(x_i, z_i; \theta^t)}_{\phi_i^t(z_i)}.$$

θ^{t+1} is obtained by maximizing the above.

$$\therefore l(\theta^{t+1}) \geq \sum_i \sum_{z_i} \underbrace{\phi_i^t(z_i) \log f(x_i, z_i; \theta^{t+1})}_{\phi_i^t(z_i)}.$$

[Jensen's inequality].

$$\geq \sum_i \sum_{z_i} \underbrace{\phi_i^{(t)}(z_i) \log f(x_i, z_i; \theta^t)}_{\phi_i^t(z_i)},$$

$$= l(\theta^{(t)})$$

[by choice]
(67).

[$\because \theta^{t+1}$ is explicitly obtained to maximize $l(\theta^t)$.]

$$\therefore \underline{l(\theta^{t+1}) \geq l(\theta^t)}.$$

Note that global convergence is not guaranteed.

Thus, EM guarantees the increase in likelihood for each iteration.

EM for the GMM case:

As before, let's define E-step at t^{th} iteration

$$\phi_i^t(z_i=j) = P[z_i=j | x_i; \mu^t, \sigma^t].$$

Let's call it w_{ij}^t

M-Step:

$$\ell = \sum_{i=1}^n \sum_j \phi_i^t(z_i) \frac{\log f(x_i, z_i; \mu^t, \sigma^t)}{\phi_i^t(z_i)}$$

$$= \sum_{i=1}^n \sum_{j=1}^K \phi_i^t(z_i=j) \frac{\log \left(\frac{1}{\sigma_i^t} \exp\left(-\frac{1}{2} \frac{(x_i - \mu_i^t)^2}{\sigma_i^t}\right) \right)}{\phi_i^t(z_i=j)}$$

$$= \sum_{i=1}^n \sum_{j=1}^K \omega_{ij}^t \log \frac{1}{\sqrt{2\pi\sigma_i^t}} \exp\left(-\frac{1}{2} \frac{(x_i - \mu_i^t)^2}{\sigma_i^t}\right) \cdot \frac{1}{\omega_{ij}^t}$$

$$\frac{\partial \ell}{\partial \mu_j^t} = 0 \Rightarrow \mu_j^{t+1} = \frac{\sum_{i=1}^n \omega_{ij}^t x_i}{\sum_{i=1}^n \omega_{ij}^t}$$

$$\text{Similarly for } \sigma_j^{t+1} = \frac{\sum_{i=1}^n \omega_{ij}^t (x_i - \mu_j^t)^2}{\sum_{i=1}^n \omega_{ij}^t}$$

Eq (69).

We have to estimate α_i^s too.

in ℓ , the terms that depend upon α

are is
$$\ell' = \sum_{i=1}^n \sum_{j=1}^k w_{ij}^t \log \alpha_j^t$$

This again is a constrained optimization problem $\because \sum_{j=1}^k \alpha_j = 1 \therefore \alpha_j = P(Z_i = j)$.

$$\therefore L(\alpha) = \sum_{i=1}^n \sum_{j=1}^k w_{ij}^t \log \alpha_j^t + d \left(\sum_{j=1}^k \alpha_j^t - 1 \right)$$

$$\Rightarrow \alpha_j^t = \frac{\sum_{i=1}^n w_{ij}^t}{d}$$

since $\sum_{j=1}^k \alpha_j^t = 1,$

$$d = \sum_{i=1}^n \sum_{j=1}^k w_{ij}^t$$

Thus,
$$\alpha_j^t = \frac{1}{d} \sum_{i=1}^n w_{ij}^t$$

$$\sum_{j=1}^k \sum_{i=1}^n \frac{w_{ij}^t}{d} = 1$$

$$\frac{n}{d} = 1 \Rightarrow d = n$$