

Expectation-Maximization

notes from notes from Stanford Machine Learning, Ng and Ma 2023

https://cs229.stanford.edu/main_notes.pdf

Mixture of Gaussians model is an example of EM for missing data

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

Observed, incomplete-data

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

the model is the joint probability that each $x^{(i)}$ was generated by randomly choosing $z^{(i)}$ from $\{1 \dots k\}$.

$$z^{(i)} \sim \text{Multinomial}(\phi)$$

Latent, unobserved cluster assignments, needed for complete-data

$$\phi_j \geq 0, \sum_{j=1}^k \phi_j = 1,$$

$$\phi_j \text{ gives } p(z^{(i)} = j),$$

Cluster (gaussian) weights

$$p(x^{(i)} | z^{(i)} = j) \sim \mathcal{N}(\mu_j, \Sigma_j)$$

we adopt the gaussian distribution to represent x

ϕ , μ and Σ are the model parameters.

k is a hyper-parameter

Number of clusters (number of gaussians)

we need the likelihood of \mathbf{x} , $\mathbf{p}(\mathbf{x}; \boldsymbol{\phi}, \boldsymbol{\mu}, \boldsymbol{\Sigma})$ to estimate the parameters:

$$\begin{aligned}\ell(\boldsymbol{\phi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) &= \sum_{i=1}^n \log p(x^{(i)}; \boldsymbol{\phi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) \\ &= \sum_{i=1}^n \log \sum_{z^{(i)}=1}^k p(x^{(i)}|z^{(i)}; \boldsymbol{\mu}, \boldsymbol{\Sigma})p(z^{(i)}; \boldsymbol{\phi}) = \sum_{i=1}^n \log \sum_{z^{(i)}=1}^k p(x^{(i)}, z^{(i)})\end{aligned}$$

If we knew $z^{(i)}$'s, the problem would be solvable using:

$$\ell(\boldsymbol{\phi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \sum_{i=1}^n \log p(x^{(i)}|z^{(i)}; \boldsymbol{\mu}, \boldsymbol{\Sigma}) + \log p(z^{(i)}; \boldsymbol{\phi}).$$

then from setting derivatives to 0 to find the maxima, we would have:

$$\begin{aligned}\phi_j &= \frac{1}{n} \sum_{i=1}^n 1\{z^{(i)} = j\}, \\ \mu_j &= \frac{\sum_{i=1}^n 1\{z^{(i)} = j\}x^{(i)}}{\sum_{i=1}^n 1\{z^{(i)} = j\}}, \\ \Sigma_j &= \frac{\sum_{i=1}^n 1\{z^{(i)} = j\}(x^{(i)} - \mu_j)(x^{(i)} - \mu_j)^T}{\sum_{i=1}^n 1\{z^{(i)} = j\}}.\end{aligned}$$

since we do NOT know the $z^{(i)}$'s, we use expectation-maximization:

in Expectation step: guess the $z^{(i)}$'s

in Maximization step: update the model parameters, based upon the $z^{(i)}$ estimates

step 0: init parameters k , and ϕ_j, μ_j, Σ_j for $k=1$ to k

prev_log_l = log(1.E-323); tol=1E-4

while (true):

#E-step:

for each $j=1,n$, for each $i=1,k$: **#estimate the posterior**

#set $Q_i(z^{(i)} = j) := p(z^{(i)} | x^{(i)}; \theta)$ so that $ELBO(x; Q, \theta) = \log p(x; \theta)$ for x and the current θ

$Q_i(z^{(i)} = j) = \frac{w^{(i)}_j}{\sum_{\ell=1}^n w^{(i)}_{\ell}} = \frac{p(z^{(i)} = j | x^{(i)}; \phi, \mu, \Sigma)}{\sum_{\ell=1}^n p(z^{(i)} = \ell | x^{(i)}; \phi, \mu, \Sigma)}$ from Bayes Rule = $p(x|z) * p(z) / p(x)$
= $N(x^{(i)} | \mu_j, \Sigma_j) * \phi_j / \sum_{\ell=1}^n (\phi_{\ell} * N(x^{(i)} | \mu_{\ell}, \Sigma_{\ell}))$
#note that these should sum to 1 for a single data point $x^{(i)}$

#calculate log likelihood too as it is used to define convergence.

$\log_l = l(\theta) = \sum_{j=1,n} (\log p(x; \theta)) = \sum_{j=1,n} (\log (\sum_{i=1,k} (w^{(i)}_j)))$

if $(\log_l - \text{prev_log_l} \leq \text{tol})$:

converged = true

break

#M-step:

estimate ϕ, μ, Σ by maximizing ELBO:

EQN 11.12 $\theta = \operatorname{argmax} \operatorname{ELBO}(x^{(i)}; Q_i, \theta)$ w.r.t θ while fixing the choice of Q_i

$$\begin{aligned} & \sum_{i=1}^n \sum_{z^{(i)}} Q_i(z^{(i)}) \log \frac{p(x^{(i)}, z^{(i)}; \phi, \mu, \Sigma)}{Q_i(z^{(i)})} \\ &= \sum_{i=1}^n \sum_{j=1}^k Q_i(z^{(i)} = j) \log \frac{p(x^{(i)} | z^{(i)} = j; \mu, \Sigma) p(z^{(i)} = j; \phi)}{Q_i(z^{(i)} = j)} \\ &= \sum_{i=1}^n \sum_{j=1}^k w_j^{(i)} \log \frac{\frac{1}{(2\pi)^{d/2} |\Sigma_j|^{1/2}} \exp\left(-\frac{1}{2}(x^{(i)} - \mu_j)^T \Sigma_j^{-1} (x^{(i)} - \mu_j)\right) \cdot \phi_j}{w_j^{(i)}} \end{aligned}$$

$$\phi_j := \frac{1}{n} \sum_{i=1}^n w_j^{(i)},$$

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

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

hold $w^{(i)}_j$ fixed while taking the gradient of the rest of the ELBO. set gradient to zero and solve for parameters.

note that : $w^{(i)}_j = Q_i(j)$

Generalization of details of $p(\mathbf{z}|\mathbf{x})$, **choosing** Q as estimate and the ELBO:

choose $Q(z)$ proportional to the joint distribution as the model. the posterior follows.

$$Q(z) \propto p(x, z; \theta) = p(z|x; \theta) \quad \text{posterior of z's} \quad (11.8)$$

$$Q_i(z^{(i)}) = p(z^{(i)}|x^{(i)}; \theta)$$

$$\log p(x^{(i)}; \theta) \geq \text{ELBO}(x^{(i)}; Q_i, \theta) = \sum_{z^{(i)}} Q_i(z^{(i)}) \log \frac{p(x^{(i)}, z^{(i)}; \theta)}{Q_i(z^{(i)})}$$

$$\ell(\theta) \geq \sum_i \text{ELBO}(x^{(i)}; Q_i, \theta) \quad (11.11)$$

$$= \sum_i \sum_{z^{(i)}} Q_i(z^{(i)}) \log \frac{p(x^{(i)}, z^{(i)}; \theta)}{Q_i(z^{(i)})}$$