

Lecture 9: 4th Feb 2026

E M Algorithm:

Recall

$$\text{ELBO} \quad F_\theta(q) = \mathbb{E}_{q(z)} \log \frac{p_\theta(x, z)}{q(z)}$$

the optimal value

$$q(z) = p_\theta(z|x)$$

Note :

$q(z)$ is a func of θ

suppose we evaluate $F_\theta(q)$ at some θ^t

$$F_\theta(q) \Big|_{\theta=\theta^t} = \mathbb{E}_{q(z)} \log \frac{p_\theta(x, z)}{q(z)}$$

$$\text{with } q(z) = p_{\theta^t}(z|x)$$

$q(z)$ and $F_\theta(q)$ can be evaluated at different values of parameter θ

E M :

Initialize θ'

for $t=1$ to Convergence:

$$q^{t+1}(\gamma) = p_{\theta^t}(\gamma | x)$$

compute

E-step $F_{\theta}(q^{t+1}) = \mathbb{E}_{q^{t+1}(\gamma)} \frac{\log p_{\theta}(x, \gamma)}{q^{t+1}(\gamma)}$

M-step $\theta^{t+1} = \underset{\theta}{\operatorname{argmax}} F_{\theta}(q^{t+1})$

End For

Convergence of EM Algorithm

Claim with EM, $\lambda(\theta^{t+1}) \geq \lambda(\theta^t)$

proof: Consider θ^t

$$q^{t+1}(\gamma) = p_{\theta^t}(\gamma | x)$$

$$\lambda(\theta^t) = F_{\theta}(q^{t+1}) \Big|_{\theta=\theta^t}$$

In M step

$$\theta^{t+1} = \underset{\theta}{\operatorname{argmax}} F_{\theta}(q^{t+1}) \Big|_{\theta=\theta^t}$$

$$\therefore F_{\theta}(q^{t+1}) \Big|_{\theta=\theta^{t+1}} \geq F_{\theta}(q^{t+1}) \Big|_{\theta=\theta^t}$$

also $F_{\theta}(q^{t+1}) \Big|_{\theta=\theta^{t+1}} \leq L(\theta^{t+1})$

By Jensen's inequality

$$\Rightarrow L(\theta^{t+1}) \geq L(\theta^t) \approx$$

→ likelihood func never decreases when we run

EM algo

→ EM can stuck in local minima

→ stopping criteria based on change in the

likelihood func

EM for GMMs:

$$p_{\theta}(x) = \sum_z p_{\theta}(x|z)$$

$$= \sum_z \alpha_j p_{\theta_j}(x)$$

$$p_{\theta_j}(x) \sim N(x; \mu_j, \Sigma_j)$$

$$= \sum_z p_{\theta}(z) p_{\theta}(x|z)$$

$$z \sim \text{Discrete} \{ \alpha_1, \alpha_2, \dots, \alpha_m \}$$

$$p_{\theta}(x|z=j) \sim N(x; \mu_j, \Sigma_j)$$

$$\Theta = \{ \alpha_j, \mu_j, \Sigma_j \}_{j=1}^m$$

$$\alpha_j \in \mathbb{R}$$

$$\mu_j \in \mathbb{R}^d$$

$$\Sigma_j \in \mathbb{R}^{d \times d}$$

Step 1:

Initialize θ^t

Step 2: Compute

$$\begin{aligned} q^{t+1}(z) &= \frac{p_{\theta^t}(z|x)}{\sum_j p_{\theta^t}(x|z) p_{\theta^t}(z)} \\ &= \frac{N(x; \mu_j^t, \Sigma_j^t) \alpha_j^t}{\sum_{j=1}^m N(x; \mu_j^t, \Sigma_j^t) \alpha_j^t} \end{aligned}$$

We can have a closed form sol'n for this.

Step 3: E step: Compute $F_\theta(q^{t+1})$

$$= \mathbb{E}_{q^{t+1}(z)} \log \frac{p_\theta(x, z)}{q^{t+1}(z)}$$

$$= \mathbb{E}_{\theta^t} \left[\log \frac{p_{\theta}(x, z)}{p_{\theta^t}(z|x)} \right]$$

$$= \sum_{j=1}^m \alpha_j^t \log \frac{p_{\theta}(x|z)}{p_{\theta^t}(z)}$$

=

Solve this with
proper notations.

Step 4 : M step

$$\mu_j^{t+1} = \underset{\mu_j}{\operatorname{argmax}} F_{\theta}(q_j^{t+1})$$

Solve it .

Lecture 10: 9th Feb 2026:

Bayesian Methods for density Estimation

- * Recall in MLE, the parameter set " Θ " is considered to be a deterministic variable, that is fully dependent on data.
- * When n is small, MLE might not be optimal.

e.g.: Consider MLE for a Bernoulli RV.

$$\hat{\theta}_{MLE} = \frac{\# \text{ trial has yielded "1"}}{n}$$

Suppose $D = \{1, 1, 1, 1\}$

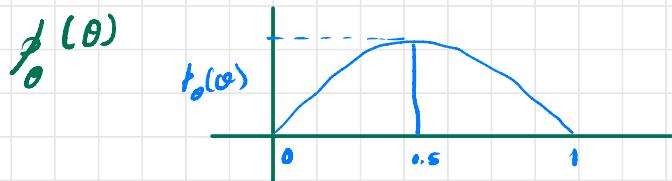
$$\hat{\theta}_{MLE} = \frac{4}{4} = 1$$

is it a good estimate?

Can we rely fully on the data?

The above estimate is not "optimal", if the true value is around 0.5.

- * How do we incorporate the prior belief/information into the estimation technique.
- * Assume the parameter " θ " to be a RV & model the prior belief in its distribution.
- * Suppose in the previous example of Bernoulli RV.



- * prior belief is completely independent of data
- * use this in addition to the Data.

* Suppose $X \sim p_x$
 $\theta \sim p_\theta$: prior on parameter.

Compute

$$p_{\theta|x} \propto \frac{p_\theta(\theta)}{p_\theta(\theta)} \frac{p_x(x|\theta)}{p_x(x)}$$

↓ ↓
prior likelihood.

* How to use $p_{\theta|x}$?

MAP Estimate: Take the Mode

of $p_{\theta|x}$ as the estimate for θ .

Conjugate Priors:

- * choice of p_θ is user choice.
- * The prior on θ , that ensures the posterior to have the same distributional form as that of the prior.

$p_\theta(\theta)$ is bounded b/w 0 & 1

$p_{\theta|x}$ can't be Gaussian with infinite Support.

\therefore choose p_θ s.t $p_{\theta|x}$ is also of a same form.

Example of MAP for Bernoulli.

$X \sim \text{Bernoulli}(\theta)$

$$p_{x|\theta} = \theta^x (1-\theta)^{1-x}$$

$$p_\theta(\theta) = \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha) \Gamma(\beta)} \theta^{\alpha-1} (1-\theta)^{\beta-1}$$

where $P(\cdot)$ is std gamma func & α, β are user chosen params.

compute : (assuming one data point)

$$P_{\theta|x} \propto p_{x|\theta} p_\theta$$

$$\begin{aligned} P_{\theta|x} &\propto [\theta^\alpha (1-\theta)^{1-\alpha}] \theta^{\alpha-1} (1-\theta)^{\beta-1} \\ &\propto \theta^{\alpha+\alpha-1} (1-\theta)^{\beta-\alpha} \end{aligned}$$

In MAP estimate, θ_{MAP} is taken to be the mode

$$\text{of } P_{\theta|x}$$

* In the above case,

$$\theta_{MAP} = \frac{x + \alpha - 1}{\alpha + \beta - 2} \quad / \quad \frac{\sum_{i=1}^n x_i + \alpha - 1}{n + \alpha + \beta - 2}$$

$$\theta_{MLE} = \frac{1}{n} \sum_{i=1}^n x_i$$

as $n \rightarrow \infty$, the belief on the data vanishes.

\Rightarrow MAP converges to MLE

Non - Parametric Density Estimate

- * These methods can evaluate the density func at a point without making any parametric func form assumption:
- * Suppose f_x is the density func that is being evaluated.

$$x \in \mathbb{R}^d$$

"R" denote a "region" on the support of f_x

- * Let P be the prob of a point "x" falling within the region. R

$$P = \int_R f_x(x) dx$$

- * Suppose we have "n" datapoints drawn from f_x ^(iid)
- * if the event of a datapoint falling within the region R , be modeled as a Bernoulli Rv.

* the MLE for that parameter is given by

$$\hat{P} = \frac{k}{n}$$

where "k" is the number of points within the region.

* From earlier what

$$\hat{P} = \int_R p_x(x) dx$$

$$\therefore \frac{k}{n} = \int_R p_x(x) dx.$$

if R is small enough

$$\frac{k}{n} = p_x(x) v \quad \text{where}$$

"v" is the volume of the region.

$$\Rightarrow p_x(x) = \frac{k}{n v}$$

Density is const over the volume

Similar to histogram.

- * 2 ways to construct Non parametric density estimate
 - 1) Fix N & Count K
 - 2) Fix K & grow N to encompass K points.

Parzen Window Estimates

- * Fix V & Comt K to evaluate density
- * Suppose we want to evaluate the density func at a point x.
- * Let R be a d-dimensional hypercube around x with length h

$$V = h^d$$

- * to count "K" define a window func
- $$\phi(u) = \begin{cases} 1 & \text{if } |u_j| \leq \frac{1}{2} \text{ } \forall j=1, \dots, d \\ 0 & \text{otherwise} \end{cases}$$

$$\phi\left(\frac{x - x_i}{h}\right) \Rightarrow \text{"x;" is within hypercube of Volume } h^d \text{ centered around } x;$$

$$p_x(x) = \frac{k}{n}$$

$$f_x(x) = \frac{\sum_{i=1}^n \phi\left(\frac{x-x_i}{h}\right)}{n \cdot h^d}$$

* $\phi()$ is very hard, we need to make it more smoother by choosing

$$\phi(u) = \exp(-\|u - u_0\|_2^2)$$

* $f(x)$: kernel, depending on the choice of kernel the estimate changes

Nearest Neighbour estimate:

+ non parametric estimate where fix n & grow N .

$$p_x(x) = \frac{k}{n}$$

Example Use Case: "m" class classification.

* Bayes classifier with NN density estimate for

$$p_{y|x}$$

- * Suppose we want to estimate the class of sample "x"
- * we place a volume V around "x" & capture k-samples.

* Let k_i be the number of points within volume V corresponding to the i^{th} class.

$$\sum_{i=1}^m k_i = k$$

* The NN estimate for $f(x, y_i) = \frac{k_i}{nV}$

* now to construct

$$\begin{aligned}
 f_{y|x}(y=y_i|x) &= \frac{f(x, y_i)}{\sum_{i=1}^m f(x, y_i)} \\
 &= \frac{k_i/nV}{\sum_{j=1}^m k_j/nV} \\
 &= \frac{k_i}{k}
 \end{aligned}$$

$$h_B(x) = i \quad \text{if} \quad p(y_i|x) > p(y_j|x) \quad \forall j \neq i$$

$$= i \quad \text{if} \quad \frac{k_i}{k} > \frac{k_j}{k} \quad \forall j \neq i$$

$$= i \quad \text{if} \quad k_i > k_j \quad \forall j \neq i$$

This is the K-Nearest neighbour Classifier.

* True risk with NN is at most twice that of Risk associated with Bayes Classifier. : Proof ??