

Lecture 1b : EM Theory.

Quiz 2 \rightarrow Up to & incl. boosting (proof)

Recall:

GMM \rightarrow easy ex. of EM alg.



\rightarrow there are many application

$$p(x; \theta) = \sum_{j=1}^n \underbrace{P(j)}_{\Sigma P(j)=1} N(x; \mu_j, \underbrace{\Sigma_j}_{P\text{-semi D}})$$

Optimization Method \rightarrow EM Alg.

① E \rightarrow posterior prob. $\Rightarrow p^{(e)}(j|t) = P(j|x_t, \theta^{(e)})$

Note: $p^{(e)}(j|t) = \mathbb{E}[\delta(j|t) | X, \theta^{(e)}]$

illustrate why we can substitute $\delta(j|t) \rightarrow p^{(e)}(j|t)$

② M : $\hat{n}(j) = \sum_{t=1}^n p^{(e)}(j|t) \Rightarrow$ effective # of samples in j -th component.

$\sum_{j=1}^m \hat{n}(j) = n$

Mixture proportions: $p^{(e+1)}(j) = \frac{\hat{n}(j)}{n}$

Component means: $\mu_j^{(e+1)} = \frac{1}{\hat{n}(j)} \sum_{t=1}^n p^{(e)}(j|t) x_t$

Component covariance matrix:

$$\Sigma_j^{(e+1)} = \frac{1}{\hat{n}(j)} \sum_{t=1}^n p^{(e)}(j|t) (x_t - \mu_j^{(e)}) (x_t - \mu_j^{(e)})^T$$

Posterior prob. of x_t belonging to component j :

$p^{(e)}(j|t) = \mathbb{E}[\delta(j|t) | X, \theta^{(e)}]$ Question

\rightarrow since it is Bernoulli r.v.

$$= \Pr(\text{sample } x_t \text{ belongs to component } j \mid \underline{\theta}^{(l)})$$

$$= \frac{P(j) P(x_t \mid j, \underline{\theta}^{(l)})}{\sum_{j'=1}^M P(j') P(x_t \mid j', \underline{\theta}^{(l)})} \quad \leadsto \text{Bayes Rule}$$

Theorem to guarantee.

Let $f(\underline{\theta}) = \sum_t \log P(x_t; \underline{\theta})$ be the log-likelihood on $X = \{x_1, \dots, x_n\}$.

Specially, $f(\underline{\theta}^{(l)})$: log-likelihood of EM alg. ($l \in \mathbb{N}$)

Theorem: [EM Thm]

Unless $\underline{\theta}^{(l)}$ is a stationary point of $\max_{\underline{\theta}} f(\underline{\theta})$,
we have $\underline{f(\underline{\theta}^{(l+1)}) > f(\underline{\theta}^{(l)})} \quad \forall l \in \mathbb{N}$ in any case!

very beautiful

Idea: Majorization - Minimization (MM) Alg.

Consider the opt. (minimum) prob. $\min_{x \in \mathbb{R}^d} f(x) \Rightarrow$ Bad function

We assume that this is difficult. (directly)

↓ Method
Successively minimize an auxiliary function $u(x, x^{(l)})$

We hope $\{x^{(l)}\}_{l=1}^{\infty}$ converges to a stationary point!

→ to strong

↓

we want to guarantee $f(x^{(l+1)}) \leq f(x^{(l)})$

Defn: An auxiliary function of $f: \mathbb{R}^d \rightarrow \mathbb{R}$ is a bivariate func.

$$u: \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R} \quad \text{s.t.}$$

① $u(x, x) = f(x)$

② $u(x, x^{(l)}) \geq f(x)$ $\forall (x, x^{(l)}) \in \mathbb{R}^d \times \mathbb{R}^d$



says that $u(\cdot, x^{(l)})$ is a Majorizer of f !

Say we have an auxiliary f $u(x, x^{(l)})$ of $f(x)$

MM procedure is

$$x^{(l+1)} = \boxed{\operatorname{argmin}_x u(x, x^{(l)})}$$



assume this is easier than prime prob.
(substantially)

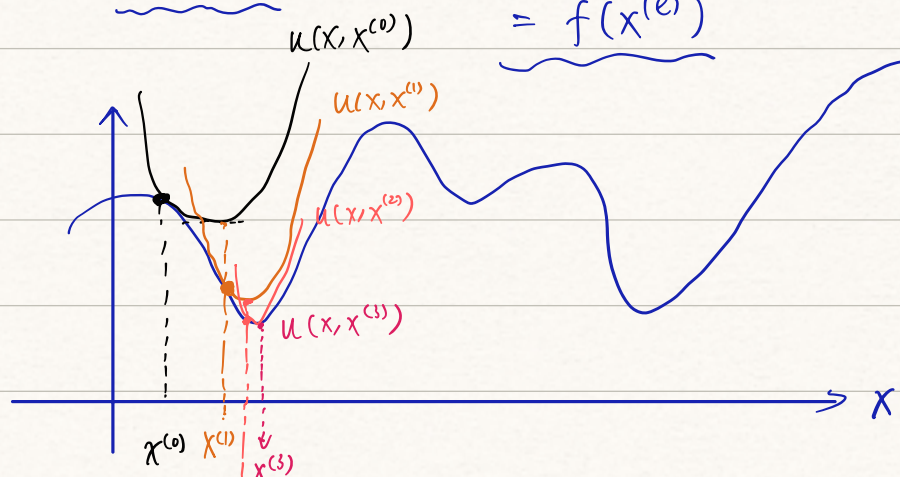
[Props]: $f(x^{(l+1)}) \leq f(x^{(l)}) \Rightarrow$ monotonically decrease!

Pf: $f(x^{(l+1)}) \leq u(x^{(l+1)}, x^{(l)})$ $= \min u(x, x^{(l)})$

↑
property 2

$$\leq u(x^{(l)}, x^{(l)})$$

$$= \underline{f(x^{(l)})}$$



Rmk: Equality holds iff $f(x^{(l+1)}) = u(x^{(l+1)}, x^{(e)})$ &
 $u(x^{(l+1)}, x^{(e)}) = u(x^{(e)}, x^{(e)})$

Rmk: would like $x^{(e)} \rightarrow$ stationary point of f (*)

Under additional condition \rightarrow continuity & differentiability
of u , we can show (*)

Rmk: How to get $u(\cdot, \cdot)$? \rightarrow

① first-order Taylor

② Jensen inequality

To be more precise, we should
find a GOOD u !

Application of MM framework to EM alg.

\nearrow observation

$X = \{x_t\}_{t=1}^n$ Latent / Hidden variable $Z = \{z_t\}_{t=1}^n$

Complete Dataset (X, Z)

\uparrow
incomplete

Goal: Estimate θ from incomplete X

$$\text{MLE} \rightarrow \hat{\theta}_{\text{MLE}} = \underset{\theta}{\operatorname{argmin}} -\log p(X|\theta)$$

Intractable since $p(X|\theta) = \sum_z p(X, Z|\theta)$

E-step: Evaluate $p(Z|X, \underline{\theta}^{(e)}) = \{p^{(e)}(j|t)\}$



Soft guess of latent variables Z given dataset X and current $\underline{\theta}^{(e)}$

M-step: Update $\underline{\theta}$ as:

$$(*) \leftarrow \underline{\theta} = \operatorname{argmin}_{\underline{\theta}} - \mathbb{E}[\log P(X, Z|\underline{\theta}) | X, \underline{\theta}^{(e)}]$$

Complete dataset with soft Z (guess)

↗ $Z \sim Z|X, \underline{\theta}^{(e)}$

We like to show that $-\log P(X|\underline{\theta}^{(t+1)}) \leq -\log P(X|\underline{\theta}^{(e)})$

Pf: Notice $(*)$ is the same as,

$$\underline{\theta}^{(t+1)} = \operatorname{argmin}_{\underline{\theta}} u(\underline{\theta}, \underline{\theta}^{(e)})$$

(not related to $\underline{\theta}$)

Const

$$\text{and } u(\underline{\theta}, \underline{\theta}^{(e)}) = -\mathbb{E}[\log P(X, Z|\underline{\theta}) | X, \underline{\theta}^{(e)}] - \boxed{H(P(\cdot | X, \underline{\theta}^{(e)}))}$$

$$H(p(\cdot)) = - \sum_z p(z) \log(p(z))$$

If we can check $u(\underline{\theta}, \underline{\theta}^{(e)})$ is an auxiliary f^n of $f(\underline{\theta}) = -\log P(X|\underline{\theta})$, then our conclusion comes from MM props.

Then, check:

$$① f(\underline{\theta}) = u(\underline{\theta}, \underline{\theta})$$

$$u(\underline{\theta}, \underline{\theta}) = -\mathbb{E}[\log P(X, Z|\underline{\theta}) | X, \underline{\theta}] - H(P(\cdot | X, \underline{\theta}))$$

$$\begin{aligned}
&= - \sum_z \left[\log P(X, Z | \theta) \right] P(Z | X, \theta) + \sum_z P(Z | X, \theta) \log P(Z | X, \theta) \\
&= - \sum_z P(Z | X, \theta) \log P(X | \theta) \\
&= - \log P(X | \theta) \\
&= f(\theta)
\end{aligned}$$

② To check $u(\theta, \theta^{(n)}) \geq f(\theta)$

[Jensen Inequality] Let $g \rightarrow \text{convex} \Rightarrow \underline{g(\mathbb{E}Y) \leq \mathbb{E}[g(Y)]}$

↓
Verify Part ②

$$P(X|\theta) \leftarrow \frac{\frac{P(X, Z|\theta)}{P(Z|X, \theta)}}{\sum P(X|Z, \theta) P(Z|\theta)}$$

$$\text{RHS} = f(\theta) = -\log P(X|\theta)$$

$$= -\log \sum_z P(X|Z, \theta) P(Z|\theta)$$

$$\begin{aligned}
&-\log \sum P(X, Z|\theta) \\
&= -\log \sum p_i \frac{P(X, Z|\theta)}{p_i} \\
&\leq -\sum p_i \log \frac{P(X, Z|\theta)}{p_i}
\end{aligned}$$

$$= -\log \sum_z P(Z|\theta) \left(\frac{P(Z|X, \theta^{(n)}) P(X|Z, \theta)}{P(Z|X, \theta^{(n)})} \right)$$

$$= -\log \sum_z P(Z|X, \theta^{(n)}) \frac{P(Z|\theta) P(X|Z, \theta)}{P(Z|X, \theta^{(n)})}$$

Jensen Inequality

$$\leq -\sum_z P(Z|X, \theta^{(n)}) \log \frac{P(Z|\theta) P(X|Z, \theta)}{P(Z|X, \theta^{(n)})}$$

$$= -\sum_z P(Z|X, \theta^{(n)}) \log P(X, Z|\theta) - H(p_{\bullet} | X, \theta^{(n)})$$

$$= - \mathbb{E}_{z \sim p(z|x, \theta^{(l)})} [\log P(x, z | \theta)] - H(P(\cdot | x, \theta^{(l)}))$$

$$= u(\theta, \theta^{(l)})$$

Hence, by the properties of Auxiliary function $u(\theta, \theta^{(l)})$ of function f .

$$\text{we have } -\log P(x | \theta^{(l+1)}) \leq -\log P(x | \theta^{(l)})$$

$$\Leftrightarrow \underline{f(\theta^{(l+1)}) \leq f(\theta^{(l)})}$$

Application of Jensen's Ineq. explained.

$$= -\log \sum_z P(z | x, \theta^{(l)}) \frac{P(x | z, \theta) P(z | \theta)}{P(z | x, \theta^{(l)})}$$

$$= -\log \mathbb{E}_{z \sim P(z | x, \theta^{(l)})} \left[\frac{P(x | z, \theta) P(z | \theta)}{P(z | x, \theta^{(l)})} \right]$$

$$= g \mathbb{E}(Y)$$

$$\leq \mathbb{E}[g(Y)]$$

$$= \sum_z P(z | x, \theta^{(l)}) - \log \frac{P(x | z, \theta) P(z | \theta)}{P(z | x, \theta^{(l)})}$$

$$= - \sum_z \left[P(z | x, \theta^{(l)}) \log \frac{P(x | z, \theta) P(z | \theta)}{P(z | x, \theta^{(l)})} \right] \quad \textcircled{2}$$

$$= \underline{f(\theta)}$$

Other Application of MM : Ranking

$m = 20$ teams

Each team has a skill level $\theta_i \in [0, 1]$

BTL Model

$$\Pr(i \text{ beats } j) = \frac{\theta_i}{\theta_i + \theta_j}$$

$$\mathcal{D} = \{b_{ij} : i, j \in [m]\} \quad \underline{\theta} = \{\theta_1, \dots, \theta_m\}$$

$b_{ij} = \#$ of times i beat j in one season

$$\mathcal{L}(\mathcal{D}; \underline{\theta}) = \prod_{i,j} \left(\frac{\theta_i}{\theta_i + \theta_j} \right)^{b_{ij}} \Rightarrow \text{Here, } i \neq j!$$

$$\text{MLE} \rightarrow \ell(\mathcal{D}; \underline{\theta}) = \sum_{i,j} b_{ij} [\log \theta_i - \log(\theta_i + \theta_j)]$$

$$\text{Minimize } -\ell(\mathcal{D}; \underline{\theta}) = f(\underline{\theta}) = - \sum_{\substack{i,j \\ i \neq j}} b_{ij} [\log \theta_i - \log(\theta_i + \theta_j)]$$

① Gradient $\Rightarrow \nabla f \Rightarrow \theta_i \leftrightarrow \theta_j$ related in second term $\log(\theta_i + \theta_j)$

$$\frac{\partial f}{\partial \theta_k} = - \sum_j b_{kj} \left(\frac{1}{\theta_k} - \frac{1}{\theta_k + \theta_j} \right)$$

$$+ \sum_{\substack{i=k \\ i \neq k, j=k}} b_{ik} \cdot \frac{1}{\theta_i + \theta_k}$$

Couple

② MM framework \rightarrow find Auxiliary $u(\cdot, \cdot)$ of f°

For a concave f° $h(\cdot)$,

$$h(y) \leq h(x) + h'(x)(y - x)$$

$u(x, y)$
 \Downarrow
 upper Bound \rightarrow property [2]

we want to apply this to $h(y) = \log(y)$

\Downarrow

$$h(y) \leq \log x + \frac{y-x}{x}$$

A majorizer of f is:

$$u(\underline{\theta}, \underline{\theta}^{(l)}) = - \sum_{i,j} b_{ij} \left[\log \theta_i - \log (\theta_i^{(l)} + \theta_j^{(l)}) - \frac{\theta_i + \theta_j}{\theta_i^{(l)} + \theta_j^{(l)}} + 1 \right]$$

$$= - \sum_{i,j} b_{ij} \log \theta_i + \sum_{i,j} b_{ij} \frac{\theta_i + \theta_j}{\theta_i^{(l)} + \theta_j^{(l)}}$$

$$\Downarrow \quad \frac{\partial}{\partial \theta_i} = - \sum_j b_{ij} \frac{1}{\theta_i} + \sum_j \frac{(b_{ij} + b_{ji})}{\theta_i^{(l)} + \theta_j^{(l)}}$$

Update Way: $\underline{\theta}_i^{(l+1)} = \underset{\theta}{\operatorname{argmin}} u(\underline{\theta}, \underline{\theta}^{(l)})$

= closed form

$$= \frac{\sum_{j \neq i} b_{ij}}{\sum_{j \neq i} (b_{ij} + b_{ji}) / (\theta_i^{(l)} + \theta_j^{(l)})}$$