

Lec 6 : Missing Data & EM Algorithm

Yanjun Han
Oct. 17, 2023



Convex conjugate

Defn (convex conjugate). Let f be convex on \mathbb{R}^d . The convex conjugate f^* of f is defined as

$$f^*(t) = \max_{x \in \mathbb{R}^d} \langle t, x \rangle - f(x).$$

Examples. 1. if $f(x) = \frac{1}{2}\|x\|_2^2$, then

$$\begin{aligned} f^*(t) &= \max_{x \in \mathbb{R}^d} \langle t, x \rangle - \frac{1}{2}\|x\|_2^2 \\ &= \max_{x \in \mathbb{R}^d} -\frac{1}{2}\|x-t\|_2^2 + \frac{1}{2}\|t\|_2^2 \\ &= \frac{1}{2}\|t\|_2^2 \end{aligned}$$

2. if $f(x) = \sum_{i=1}^d x_i \log x_i$, then

$$f^*(t) = \max_x (t_1 x_1 + t_2 x_2 + \dots + t_d x_d - \sum_{i=1}^d x_i \log x_i).$$

F.O.C. for maximization:

$$t_i - (1 + \log x_i^*) = 0 \Rightarrow x_i^* = e^{t_i - 1}$$

$$\Rightarrow f^*(t) = \sum_i (t_i x_i^* - x_i^* \log x_i^*) = \sum_i x_i^* = \frac{1}{e} (e^{t_1} + \dots + e^{t_d}).$$

Properties. Assume that f is convex & continuously differentiable.

1. $t \mapsto f^*(t)$ is convex.

Pf. $t \mapsto \langle t, x \rangle - f(x)$ is affine.

$$\Rightarrow t \mapsto \max_x \langle t, x \rangle - f(x) = f^*(t) \text{ is convex.}$$

2. The maximizer x^* satisfies $t = \nabla f(x^*)$.

Pf. F.O.C. $0 = \nabla_x (\langle t, x \rangle - f(x))|_{x=x^*} = t - \nabla f(x^*)$.

3 (Fenchel-Young inequality) $f(x) + f^*(t) \geq \langle x, t \rangle$.

Equality holds iff $t = \nabla f(x)$

Pf. $f^*(t) = \max_{x_0} \langle x_0, t \rangle - f(x_0) \geq \langle x_0, t \rangle - f(x)$,

equality holds iff $x_0 = x$ is the maximizer, i.e. $t = \nabla f(x)$.

4 If f is convex, then $f^{**} = f$.

(In other words, $f(x) = \max_{t \in \mathbb{R}^d} \langle t, x \rangle - f^*(t)$.)

Pf. Fenchel-Young $\Rightarrow f(x) \geq \langle t, x \rangle - f^*(t)$ $\forall t$
 $\Rightarrow f(x) \geq \max_{t \in \mathbb{R}^k} \langle t, x \rangle - f^*(t)$.

Equality condition in Fenchel-Young:

$$f(x) = \langle \nabla f(x), x \rangle - f^*(\nabla f(x))$$

$$\leq \max_{t \in \mathbb{R}^k} \langle t, x \rangle - f^*(t).$$

Therefore,

$$f(x) = \max_{t \in \mathbb{R}^k} \langle t, x \rangle - f^*(t) = f^{**}(x). \quad \square$$

Application in exponential families. $p_\theta(y) = \exp(\langle \theta, T(y) \rangle - A(\theta)) h(y)$.

Since $\nabla^2 A(\theta) = \text{Cov}_\theta(T(y)) \succeq 0 \Rightarrow \theta \mapsto A(\theta)$ convex.

Therefore, $A(\theta) = \max_{\mu \in \mathbb{R}^k} \langle \theta, \mu \rangle - A^*(\mu)$. (Property 4)

Maximizer $\mu^*(\theta)$: $\mu^*(\theta) = \underset{\substack{\uparrow \\ \text{property 2}}}{\nabla_\theta A(\theta)} = \mathbb{E}_\theta[T(y)]$

Missing data in exponential families

$$(x_1, y_1), \dots, (x_n, y_n) \stackrel{\text{iid}}{\sim} p_\theta(x, y) = \exp(\langle \theta, T(x, y) \rangle - A(\theta)) h(x, y).$$

(x_1, \dots, x_n) : unobserved variables (y_1, \dots, y_n) : observed variables

Goal: find the MLE for θ .

EM algorithm: derivation.

incomplete log-likelihood for y :

$$\begin{aligned} l_\theta(y_1, \dots, y_n) &= \frac{1}{n} \sum_{i=1}^n \log p_\theta(y_i) \\ &= \frac{1}{n} \sum_{i=1}^n \log \int_X p_\theta(x, y_i) dx \\ &= \frac{1}{n} \sum_{i=1}^n \log \int_X \exp(\langle \theta, T(x, y_i) \rangle - A(\theta)) h(x, y_i) dx \\ &= \frac{1}{n} \sum_{i=1}^n (A_{y_i}(\theta) - A(\theta)), \end{aligned}$$

where

$$A_{y_i}(\theta) = \log \int_X \exp(\langle \theta, T(x, y_i) \rangle) h(x, y_i) dx$$

is convex in θ .

$$\begin{aligned}
 \text{MLE: } \max_{\theta} l_{\theta}(y_1, \dots, y_n) &= \max_{\theta} \frac{1}{n} \sum_{i=1}^n (A_{y_i}(\theta) - A(\theta)) \\
 &= \max_{\theta} \frac{1}{n} \sum_{i=1}^n (\max_{\mu_i} (\mu_i, \theta) - A_{y_i}^*(\mu_i) - A(\theta)) \\
 &= \max_{\theta} \max_{\mu_1, \dots, \mu_n} \frac{1}{n} \sum_{i=1}^n (\langle \mu_i, \theta \rangle - A_{y_i}^*(\mu_i) - A(\theta)) \\
 &=: f(\theta, \mu_1, \dots, \mu_n)
 \end{aligned}$$

Key intuition: 1. for fixed θ , $(\mu_1, \dots, \mu_n) \mapsto f(\theta, \mu_1, \dots, \mu_n)$ is concave
 2. for fixed (μ_1, \dots, μ_n) , $\theta \mapsto f(\theta, \mu_1, \dots, \mu_n)$ is concave.

(Warning: $(\theta, \mu_1, \dots, \mu_n) \mapsto f(\theta, \mu_1, \dots, \mu_n)$ is NOT "jointly" concave!)

EM algorithm. Initialize any $\theta^{(0)} \in \mathbb{R}^d$. For $t = 0, 1, 2, \dots$:

$$\bullet \text{ E-step: } (\mu_1^{(t+1)}, \mu_2^{(t+1)}, \dots, \mu_n^{(t+1)}) = \arg \max_{(\mu_1, \dots, \mu_n)} f(\theta^{(t)}, \mu_1, \dots, \mu_n)$$

$$\Rightarrow \mu_i^{(t+1)} = \nabla_{\theta} A_{y_i}(\theta) = \mathbb{E}_{X \sim p_{f(\theta)}(x|y_i)} [T(X, y_i)]$$

("expectation" step)

$$\begin{aligned}
 \bullet \text{ M-step: } \theta^{(t+1)} &= \arg \max_{\theta} f(\theta, \mu_1^{(t+1)}, \dots, \mu_n^{(t+1)}) \\
 &= \arg \max_{\theta} \left\langle \frac{1}{n} \sum_{i=1}^n \mu_i^{(t+1)}, \theta \right\rangle - A(\theta)
 \end{aligned}$$

$$\Rightarrow \nabla_{\theta} A(\theta^{(t+1)}) = \frac{1}{n} \sum_{i=1}^n \mu_i^{(t+1)}$$

("maximization" step)

EM Intuition:

1. E-step: for each sample i with missing data x_i , think of a fake $\tilde{x}_i \sim p_{\theta}(x_i|y_i)$ and compute sufficient statistic $\mu_i = \mathbb{E}[T(\tilde{x}_i, y_i)]$.
2. M-step: aggregate the sufficient statistics "as if" there were no missing data problem: $\frac{1}{n} \sum_{i=1}^n \mu_i = \nabla A(\theta)$
3. Iterate the above process.

Example: Gaussian mixture model.

$$\text{For } i=1, \dots, n : P(x_i=j) = \pi_j, \quad j=1, 2, \dots, m,$$

$$y_i | x_i=j \sim N(\gamma_j, 1).$$

Unknown parameters : $(\pi_1, \dots, \pi_m, \gamma_1, \dots, \gamma_m)$

Unobserved variables : (x_1, \dots, x_n) observed variables : (y_1, \dots, y_n)

$$\begin{aligned} \text{Joint likelihood: } p_\theta(x, y) &= \prod_{j=1}^m \left(\frac{\pi_j}{\sqrt{2\pi}} \exp\left(-\frac{(y-\gamma_j)^2}{2}\right) \right)^{1(x=j)} \\ &\propto \exp\left(\sum_{j=1}^m 1(x=j) \left(\log \pi_j + \gamma_j - \frac{\gamma_j^2}{2} \right) \right. \end{aligned}$$

Reparametrization $\pi_j = \frac{e^{\alpha_j}}{\sum_{k=1}^m e^{\alpha_k}}, \quad \tilde{\alpha}_j = \alpha_j - \frac{\gamma_j^2}{2}.$

$$\begin{aligned} \text{Then } p_\theta(x, y) &\propto \exp\left(\sum_{j=1}^m 1(x=j) \left(\log \frac{e^{\alpha_j}}{\sum_{k=1}^m e^{\alpha_k}} + \gamma_j - \frac{\gamma_j^2}{2} \right) \right) e^{-y/2} \\ &= \exp\left(\sum_{j=1}^m 1(x=j) \left(\tilde{\alpha}_j + \gamma_j \right) - \log\left(\sum_{k=1}^m e^{\alpha_k}\right) \right) e^{-y/2} \\ &= \exp(\langle \theta, T(x, y) \rangle - A(\theta)) e^{-y/2}, \end{aligned}$$

where

$$\theta = (\tilde{\alpha}_1, \dots, \tilde{\alpha}_m, \gamma_1, \dots, \gamma_m), \quad \tilde{\alpha}_k = \alpha_k + \frac{\gamma_k^2}{2},$$

$$A(\theta) = \log \sum_{k=1}^m e^{\alpha_k} = \log \sum_{k=1}^m e^{\tilde{\alpha}_k + \frac{\gamma_k^2}{2}},$$

$$T(x, y) = (1(x=1), \dots, 1(x=m), y, 1(x=1), \dots, y, 1(x=m))$$

Auxiliary variable $\mu_i = (\mu_{i,1}, \dots, \mu_{i,m}, \gamma_{i,1}, \dots, \gamma_{i,m})$

$$\begin{aligned} \text{E-step. } P_\theta(x=j | y) &= \frac{P_\theta(x=j, y)}{\sum_{k=1}^m P_\theta(x=k, y)} = \frac{\pi_j \exp\left(-\frac{(y-\gamma_j)^2}{2}\right)}{\sum_{k=1}^m \pi_k \exp\left(-\frac{(y-\gamma_k)^2}{2}\right)}, \\ &\Rightarrow \begin{cases} \mu_{i,j}^{(t+1)} = \mathbb{E}_{\theta^{(t)}}[1(x=j) | y_i] = \frac{\pi_j^{(t)} \exp\left(-\frac{(y_i-\gamma_j^{(t)})^2}{2}\right)}{\sum_{k=1}^m \pi_k^{(t)} \exp\left(-\frac{(y_i-\gamma_k^{(t)})^2}{2}\right)} \\ \mu_{i,j}^{(t+1)} = \mathbb{E}_{\theta^{(t)}}[y_i 1(x=j) | y_i] = y_i \cdot \pi_{i,j}^{(t+1)} \end{cases} \end{aligned}$$

$$\left(\text{In vector form: } \mu_i^{(t+1)} = \mathbb{E}_{x \sim P_{\theta^{(t)}}(x|y_i)} [T(x, y_i)] \right)$$

$$\begin{aligned}
 \text{M-step.} \quad & \frac{1}{n} \sum_{i=1}^n \mu_i^{(t+1)} = \nabla A(\theta^{(t+1)}) \\
 & = \frac{1}{\sum_{k=1}^n e^{\tilde{x}_k + \frac{v_k^2}{2}}} \left(e^{\tilde{x}_1 + \frac{v_1^2}{2}}, \dots, e^{\tilde{x}_m + \frac{v_m^2}{2}}, v_1 e^{\tilde{x}_1 + \frac{v_1^2}{2}}, \dots, v_m e^{\tilde{x}_m + \frac{v_m^2}{2}} \right) \\
 & \quad (\text{all quantities here have superscript } (t+1)) \\
 & = (\pi_1^{(t+1)}, \dots, \pi_m^{(t+1)}, \pi_1^{(t+1)} v_1^{(t+1)}, \dots, \pi_m^{(t+1)} v_m^{(t+1)}) \\
 \Rightarrow & \begin{cases} \pi_j^{(t+1)} = \frac{1}{n} \sum_{i=1}^n p_{i,j}^{(t+1)} \\ v_j^{(t+1)} = \frac{1}{n \pi_j^{(t+1)}} \sum_{i=1}^n y_i^{(t+1)} p_{i,j}^{(t+1)} \end{cases}, \quad j = 1, 2, \dots, m.
 \end{aligned}$$

Generalization: Evidence Lower Bound (Optional)

Def. For probability distributions P, Q over X, the Kullback-Leibler (KL) divergence is

$$D_{KL}(P \parallel Q) = \sum_{x \in X} p(x) \log \frac{p(x)}{q(x)} \quad (\text{could be } +\infty)$$

Thm. $D_{KL}(P \parallel Q) \geq 0$.

Pf. Since $\log t \geq 1 - \frac{1}{t}$ for all $t \geq 0$,

$$\begin{aligned}
 D_{KL}(P \parallel Q) &= \sum_x p(x) \log \frac{p(x)}{q(x)} \geq \sum_x p(x) \left(1 - \frac{q(x)}{p(x)} \right) \\
 &= \sum_x p(x) - \sum_x q(x) = 1 - 1 = 0 \quad \square
 \end{aligned}$$

Evidence lower bound (ELBO).

$$\log p_\theta(y^n) \geq \mathbb{E}_{x^n \sim q(x)} \left[\log \frac{p_\theta(x^n, y^n)}{q(x^n)} \right], \quad \text{ELBO}$$

Pf. LHS - RHS = $D_{KL}(q(x^n) \parallel p_\theta(x^n | y^n)) \geq 0$.

Choice of q .

1. $q(x^*) = p_{\theta'}(x^* | y^*)$ for some $\theta' \in \Theta$.

$$\max_{\theta} \log p_{\theta}(y^*) = \max_{\theta} \max_{\theta'} \mathbb{E}_{x^* \sim p_{\theta'}(\cdot | y^*)} \left[\log \frac{p_{\theta}(x^*, y^*)}{p_{\theta'}(x^* | y^*)} \right]$$

$=: L(\theta, \theta')$

- for fixed θ : $\operatorname{argmax}_{\theta'} L(\theta, \theta') = \theta$.
- for fixed θ' :

$$\begin{aligned} \operatorname{argmax}_{\theta} L(\theta, \theta') &= \operatorname{argmax}_{\theta} \mathbb{E}_{x^* \sim p_{\theta}(\cdot | y^*)} [\log p_{\theta}(x^*, y^*)] \\ &= \operatorname{argmax}_{\theta} \left\langle \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{x_i \sim p_{\theta}(\cdot | y_i)} [T(x_i, y_i)], \theta \right\rangle \\ &\quad - A(\theta) \end{aligned}$$

(in exponential families)

- recovers EM iteration $\theta^{(t+1)} = \operatorname{argmax}_{\theta} L(\theta, \theta^{(t)})$
 - E-step: evaluation of $L(\theta, \theta^{(t)})$
(or equivalently, the posterior mean)
 - M-step: maximization over θ (the MLE on "fake data")

2. $q(x^*) = q_{\phi}(x^* | y^*)$ indexed by ϕ .

$$\begin{aligned} \max_{\theta} \log p_{\theta}(y^*) &\geq \max_{\theta} \max_{\phi} \mathbb{E}_{x^* \sim q_{\phi}(\cdot | y^*)} \left[\log \frac{p_{\theta}(x^*, y^*)}{q_{\phi}(x^* | y^*)} \right] \\ &= \max_{\theta} \max_{\phi} \underbrace{-D_{KL}(q_{\phi}(\cdot | y^*) \| p_{\theta}(x^*))}_{\text{typically admits closed-form expression}} + \mathbb{E}_{x^* \sim q_{\phi}(\cdot | y^*)} [\log p_{\theta}(y^* | x^*)] \\ &\quad \text{by taking } p_{\theta} = N(\mu, I) \\ &\quad q_{\phi}(\cdot | y^*) = N(\mu_{\phi}(y^*), \sigma_{\phi}^2(y^*)) \end{aligned}$$

Variational autoencoder: jointly maximize over (θ, ϕ) using SGD

- Gradient of $\mathbb{E}_{z \sim q_\phi}[f(z)]$ w.r.t. ϕ computed via

$$\nabla_\phi \mathbb{E}_{q_\phi}[f(z)] = \mathbb{E}_{q_\phi}[f(z) \nabla_\phi \log q_\phi(z)] \approx \frac{1}{L} \sum_{i=1}^L f(z^{(i)}) \nabla_\phi \log q_\phi(z^{(i)})$$

or $\nabla_\phi \mathbb{E}_{q_\phi}[f(z)] \approx \nabla_\phi \left(\frac{1}{L} \sum_{i=1}^L f(g_\phi(\varepsilon^{(i)})) \right)$ for $g_\phi(\varepsilon) \sim q_\phi$.