

Lecture 16 - scribbles

Tuesday, October 25, 2016
14:35

today: • duality (cont.)
• exp. family

dual problem for Max. Ent:

$$(P) \left[\begin{array}{l} \min_q \sum_x q(x) \log \frac{q(x)}{u(x)} \\ q(x) \geq 0 \\ \sum_x q(x) = 1 \\ \sum_x q(x) T_j(x) = \alpha_j \end{array} \right] \in \mathcal{M}$$

$$J(q, v, c) = \sum_x q(x) \log \frac{q(x)}{u(x)} + \sum_j v_j (\alpha_j - \sum_x q(x) T_j(x)) + c (1 - \sum_x q(x))$$

$$\Rightarrow \hat{q}(x) = u(x) \exp(v^T T(x) + c - 1)$$

exponential family

dual function:

$$g(v, c) = \inf_q J(q, v, c) = \mathbb{E}_{q^*} [v^T T(x) + c - 1] + v^T \alpha - \mathbb{E}_{q^*} [v^T T(x)] - c + \mathbb{E}_{q^*} [c]$$

$$= v^T \alpha + c - \sum_x u(x) \exp(v^T T(x)) e^{c-1}$$

maximize

with respect to c:

$$\nabla_c \rightarrow 1 - \underbrace{\sum_x u(x) \exp(v^T T(x))}_{\hat{=} Z(v)} e^{c-1} = 0$$

$$e^{c-1} = \frac{1}{Z(v)}$$

plug back, we get $\max_c g(v, c) = v^T \alpha - \ln Z(v) \hat{=} \tilde{g}(v)$

if $\alpha = \frac{1}{n} \sum_{i=1}^n T(x_i) = \mathbb{E}_{p_n} [T(x)]$

then $\tilde{g}(v) = \frac{1}{n} \sum_{i=1}^n [v^T T(x_i) - \ln Z(v)]$

$$\ln p(x; \nu) \text{ where } p(x; \nu) \triangleq \frac{\exp(\nu^T T(x))}{Z(\nu)}$$

dual problem was $\max_{\nu} \tilde{g}(\nu) = \max_{\nu} \log p(x_1, \dots, x_n | \nu)$ i.e. MLE

so ML in the exponential family with $T(x)$ as suff. statistic

is equivalent dual of Max Entropy with moment constraint $\mathbb{E}_{p_n}[T(x)] = \alpha$

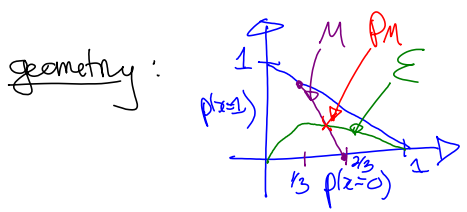
note: $\nabla_{\nu} \ln Z(\nu) = \frac{1}{Z(\nu)} \nabla_{\nu} \sum_x u(x) \exp(\nu^T T(x))$
 $= \frac{1}{Z(\nu)} \sum_x u(x) \exp(\nu^T T(x)) T(x)$
 $\underbrace{\sum_x p(x; \nu)}_{p(x; \nu)} T(x)$

$$\nabla_{\nu} \ln Z(\nu) = \mathbb{E}_{p(x; \nu)} [T(x)] \triangleq \mu(\nu) \quad \begin{matrix} \text{model} \\ \text{moment} \end{matrix}$$

so $\nabla_{\nu} \tilde{g}(\nu) = \mathbb{E}_{p_n} [T(x)] - \mu(\nu)$ $T(x) \in \mathbb{R}^1$

$$\nabla_{\nu} \tilde{g}(\nu) = 0 \Rightarrow \mathbb{E}_{p(x; \nu)} [T(x)] = \mathbb{E}_{p_n} [T(x)] \quad \text{moment matching?}$$

MCE in exponential family \Leftrightarrow moment matching in exponential family



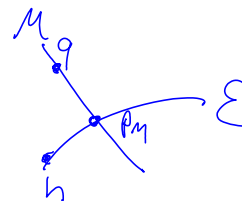
$$T(x) = \begin{cases} 1 & \text{if } x=0 \\ 1/2 & \text{if } x=1 \\ 0 & \text{if } x=2 \end{cases} \quad \Lambda = \{0, 1, 2\}$$

$$\mathbb{E}_{p_n} [T(x)] = \frac{2}{3} = \alpha$$

$$\mathcal{E} = \{ p : p(x) = u(x) \exp(\eta^T T(x) - A(\eta)) \}$$

for any $q \in \mathcal{M}$ and $h \in \mathcal{E}$

$$KL(q \parallel h) = KL(q \parallel p_M) + KL(p_M \parallel h)$$



$$\mathbb{E}_q[\log \frac{q}{h}]$$

"Information Pythagorean thm."

restate our duality result: $p_M = \arg \min_{q \in \mathcal{M}} KL(q \parallel u)$ [Max Ent]

"I-projection" I \rightarrow information

$p_M = \arg \min_{q \in \mathcal{E}} KL(\hat{p}_n \parallel q)$ MLE in exp family

"M-projection" M \rightarrow moment

Exponential Family

a (flat/canonical) exponential family on X

is a parametric family of distributions define by two quantities

I) $\underbrace{h(x)}_{\text{reference density}} d\mu(x) \rightarrow$ reference measure on X
 $\underbrace{d\mu(x)}_{\text{base measure}} \begin{cases} \text{counting (discrete)} \\ \text{Lebesgue (cts. space)} \end{cases}$

II) $T: X \rightarrow \mathbb{R}^F$ called "sufficient statistics" vector (aka. feature vector)

members of family have dist.:

$$p(x; \eta) d\mu(x) = \exp(\underbrace{\eta^T T(x)}_{\text{"canonical parameter"}} - \underbrace{A(\eta)}_{\text{log-normalizer or cumulant gen. function}}) \underbrace{h(x) d\mu(x)}_{\text{defining pieces } (\Omega_X)}$$

if Ω_X is discrete, $p(x; \eta)$ is pmf $\exp(\eta^T T(x) - A(\eta)) h(x)$

if Ω_X is cts. $p(x; \eta)$ is pdf " "

$$\neq \text{want } 1 = \int_{\Omega_X} p(x; \eta) d\mu(x) = \int_{\Omega_X} \exp(\eta^T T(x)) e^{-A(\eta)} h(x) d\mu(x)$$

$$\Rightarrow A(\eta) \triangleq \log \left(\int_X \exp(\eta^T T(x)) h(x) d\mu(x) \right)$$

domain $\Omega \triangleq \{ \eta \in \mathbb{R}^p \mid A(\eta) < \infty \}$

* more generally, consider reparameterization as a subset of the family by defining mapping

$$\eta: \Theta \rightarrow \Omega \text{ parameter}$$

consider $p(x; \theta) \triangleq p(x; \eta(\theta))$ for $\theta \in \Theta$

("curved exponential family" if $\eta(\Theta)$ is a curved manifold in Ω)

two examples not an exponential family: • Unif $[0, \theta]$
• Mixture model

but most standard distributions are (Poisson, Gamma, Gaussian, etc...)

Example: (multinomial)

$$X \sim \text{Mult}(1, \pi)$$

$$X = \{0, 1\}^k$$

$$\Omega_X = \Delta_k \cap X$$

parameter form $\pi \in \Delta_k$, suppose $\pi_j > 0 \forall j$

$$p(x; \pi) = \prod_{j=1}^k \pi_j^{x_j} = \exp \left(\sum_{j=1}^k x_j (\log \pi_j) \right) = \exp(\eta(\pi)^T x - 0)$$

we have $\eta_j(\pi) = \log \pi_j$

$$T(x) = x$$

$d\mu(x)$ = counting measure on X

$$h(x) = \mathbb{1}_{\{x \in \Omega_X\}} = \mathbb{1}_{\{x \text{ has exactly one entry equal to 1}\}}$$

here $A(\eta(\pi)) = 0$ here $\Theta = \text{int}(\Delta_k)$

$$A(\eta(\theta)) = 0 \forall \theta \in \Theta$$

remark: $\Theta \rightarrow$ dimension $k-1$

$$\perp \rightarrow \parallel \quad k$$

here, for any x s.t. $h(x) \neq 0$, $\sum_{j=1}^k T_j(x) = 1$
 i.e. $\sum_{j=1}^k T_j(x) - 1 = 0$

affine linear dep. between components of T

multiple η 's map to same distribution
 "overparametrization"

↳ not a minimal exp. family

for multinomial, minimal exp family $T(x) = \begin{pmatrix} x_1 \\ \vdots \\ x_{k-1} \end{pmatrix}$

$$Z(\eta) = \sum_{x \in \Omega_x} \exp(\eta^T T(x)) = \sum_{j=1}^{k-1} e^{\eta_j} + 1$$

$$p(x; \eta) = \exp\left(\sum_{j=1}^{k-1} \eta_j x_j - \underbrace{\log\left(\sum_{j=1}^{k-1} \exp(\eta_j) + 1\right)}_{A(\eta)}\right)$$

recall: $\nabla_{\eta} A(\eta) = \mathbb{E}_{p(x; \eta)} [T(x)]$ (valid $\eta \in \text{int}(\Omega)$)

for multinomial: $\frac{\partial A(\eta)}{\partial \eta_j} = \frac{1}{Z(\eta)} \exp(\eta_j) = p(x_j = 1 | \eta) = \mathbb{E}_{p(x; \eta)} [T_j(x)]$ as required

⊛ $T(x)$ is called "sufficient" as in statistics

$T: x \mapsto T(x)$ is 'sufficient' for a parametric model P_{θ}

$$\text{iff } \forall \theta \in \Theta \quad P_{\theta}(x) = \underbrace{h(x)}_{\text{same for family}} g(T(x); \theta)$$

i.e. dependence on θ only happens through $T(x)$

$$\text{iid. model, } \log p(x_1, \dots, x_n | \eta) = \sum_{i=1}^n [\log h(x_i) + \eta^T T(x_i) - A(\eta)]$$

→ exp. family with reference measure $\tilde{h}(x) = \prod_i h(x_i)$

$$\text{sufficient } \tilde{T}(x) = \sum_{i=1}^n T(x_i)$$

$$\text{log-partition } \tilde{A}(\eta) = n A(\eta)$$

Example 2: Gaussian

$$X \sim N(\mu, \sigma^2) \quad X = \mathbb{R}$$

$$p(x; (\mu, \sigma^2)) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$$

$$= \exp\left[-\frac{x^2}{2} \left[\frac{1}{\sigma^2}\right] + x \left(\frac{\mu}{\sigma^2}\right) - \left[\frac{\mu^2}{2\sigma^2} + \frac{1}{2} \log(2\pi\sigma^2)\right]\right]$$

$$T(x) = \begin{bmatrix} x \\ -\frac{x^2}{2} \end{bmatrix} \quad \eta(\theta) = \begin{bmatrix} \mu/\sigma^2 \\ 1/\sigma^2 \end{bmatrix} = \begin{bmatrix} \eta_1 \\ \eta_2 \end{bmatrix}$$

$$\eta_2 = \frac{1}{\sigma^2} = \text{precision} > 0$$

$$\eta_1 = \eta_2 \cdot \mu \quad \Omega = \{(\eta_1, \eta_2) : \eta_2 > 0\}$$