

Lecture 17 - scribbles

Tuesday, November 6, 2018 14:19

today: • finish MaxENT duality
• exponential family

dual problem for max. entropy

MaxENT primal form (P)

$$\min_q \sum_x q(x) \log \frac{q(x)}{u(x)}$$

absorb this constraint in domain definition of $KL(q||u)$ i.e.

$$KL(q||u) = \begin{cases} +\infty & \text{if } q(x) < 0 \text{ for some } x \\ KL(q||u) & \text{o.w.} \end{cases}$$

$$\left. \begin{aligned} & q(x) \geq 0 \quad \forall x \\ & \sum_x q(x) = 1 \end{aligned} \right\} \Delta(X) \quad \mathcal{M}$$

$$\sum_x q(x) T_j(x) = \alpha_j \quad \forall_j$$

feature fct.

$$\mathcal{L}(q, v, c) = \sum_x q(x) \log \frac{q(x)}{u(x)} + \sum_j v_j (\alpha_j - \mathbb{E}_q[T_j(x)]) + c (1 - \sum_x q(x))$$

$$\frac{\partial}{\partial q(x)} = 1 + \log \frac{q(x)}{u(x)} - \sum_j v_j T_j(x) - c = 0$$

$$\Rightarrow q^*(x) = u(x) \exp\left(\underbrace{v^T T(x)}_{\sum_j v_j T_j(x)} + c - 1\right)$$

exponential family?

dual function.

exponential family.

dual function:

plug back in $\beta(\dots)$

$$g(\nu, c) = \beta(q^*, \nu, c)$$

$$= \mathbb{E}_{q^*}[\nu^T T(x) + c - 1] + \nu^T \alpha - \mathbb{E}_{q^*}[\nu^T T(x)] + c - \mathbb{E}_{q^*}[c]$$

$$= \nu^T \alpha + c - \underbrace{\sum_x u(x) \exp(\nu^T T(x))}_{\triangleq Z(\nu)} \exp(c-1)$$

max $g(\nu, c)$

with respect to c

$$\nabla_c = 0 \Rightarrow 1 - Z(\nu) \exp(c-1) = 0$$

$$\Rightarrow \exp(c^* - 1) = \frac{1}{Z(\nu)}$$

plug back c^* :

$$\max_c g(\nu, c) = \nu^T \alpha + \underbrace{c^* - Z(\nu) \frac{1}{Z(\nu)}}_{\log Z(\nu)}$$

$$\tilde{g}(\nu) \triangleq \nu^T \alpha - \log Z(\nu)$$

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

$$\text{then } \tilde{g}(\nu) = \frac{1}{n} \sum_{i=1}^n \underbrace{[\nu^T T(x_i) - \ln Z(\nu)]}_{\log p(x_i; \nu)}$$

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

dual problem is $\max_{\nu} \tilde{g}(\nu) = \max_{\nu} \frac{1}{n} \log p(x_{1:n} | \nu)$ i.e. MLE //

to summarize: ML in the exp. family with $T(x)$ as sufficient statistics
 is equivalent to Max entropy with moment constraints on $T(x)$ where $\alpha = \mathbb{E}_{\hat{p}_n} [T(x)]$

they are Lagrangian dual of each other

MLE in exponential family \Leftrightarrow moment matching in exp. family

note: $\nabla_{\nu} \ln Z(\nu) = \frac{1}{Z(\nu)} \nabla_{\nu} \sum_x u(x) \exp(\nu^T T(x)) = \sum_x \left(\frac{1}{Z(\nu)} u(x) \exp(\nu^T T(x)) \right) T(x)$

$\nabla_{\nu} \ln Z(\nu) = \mathbb{E}_{p(x|\nu)} [T(x)] \stackrel{\text{def}}{=} \mu(\nu)$ "model moment"

$\nabla_{\nu} \tilde{g}(\nu) = \underbrace{\mathbb{E}_{\hat{p}_n} [T(x)]}_{\hat{\mu}_n \text{ "empirical moment"}} - \mu(\nu)$

$\nabla_{\nu} \tilde{g}(\nu) = 0 \Rightarrow \boxed{\mu(\nu^*) = \hat{\mu}_n}$ i.e. moment matching !

15h32

(see end of old lecture 16 for "KL Pythagorean theorem" and I-projection vs. M-projection for KL + geometry)

Exponential family

a (flat/canonical) exponential family on X

a (flat/canonical) exponential family on X

is a parametric family of distributions defined by two quantities

I) $h(x) d\mu(x)$ \rightarrow reference measure on X

reference density

base measure

Counting (discrete R.V.)
Lebesgue (cts. R.V.)

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

members of family will have dist.

$$p(x; \eta) d\mu(x) = \exp(\underbrace{\eta^T T(x)}_{\text{canonical parameter}} - \underbrace{A(\eta)}_{\text{log normalizer or cumulant generating fct.}}) h(x) d\mu(x)$$

defining pieces (\mathcal{I}_X)

log normalizer or cumulant generating fct.
log partition fct.

if \mathcal{I}_X is discrete, then $p(x; \eta)$ is a pmf

" " cts, " " " a pdf

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

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

$$\Rightarrow \int_{\mathcal{X}} \exp(\eta^T \phi(x)) \pi(x) d\mu(x)$$

domain $\Omega \triangleq \{\eta \in \mathbb{R}^p \mid A(\eta) < \infty\}$
 (set of valid canonical parameters)

note: $A(\eta)$ is convex in η
 $\Rightarrow \Omega$ is convex

⊛ more generally, consider a reparameterization of a subset of the family
 by defining the mapping $\eta: \Theta \rightarrow \Omega$

new set of parameters

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

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

↳ o.g. could consider Gaussians where $N(\mu, \mu^2)$

* note: any single distribution can be put in exponential by using $h(x) \triangleq p(x)$

- * too examples of family not an exponential family:
- mixture of Gaussians
 - $\text{Unif}(0, \theta)$

Example: (Multinomial)

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

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

$$\Omega_X = \Delta_k \cap X \quad (\text{one-hot encodings})$$

parameter $\pi \in \Delta_k$; suppose $\pi_i > 0 \forall i$

$$p(x; \pi) = \prod_{j=1}^k \pi_j^{x_j} = \exp\left(\sum_j x_j \log \pi_j\right)$$

think as "0"

$$= \exp(\eta(\pi)^T x - 0)$$

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

$$T(x) = x$$

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

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

$$\Theta = \text{int}(\Delta_k)$$

here, $A(\eta(\pi)) = 0 \forall \pi \in \Theta$

but here $\Omega = \mathbb{R}^k$

domain
(not to be confused)
with Ω_X
sample space

$$\Theta \rightarrow \text{dimension } k-1$$

$$\eta(\Theta) \rightarrow \text{ " } k-1$$

$$\Omega \rightarrow \text{dimension } k$$

here, we do not have "minimal exponential family"

Here, we do not have "minimal exponential family"

note:

for any x st. $h(x) \neq 0$

$$\text{hence, } \sum_{j=1}^k T_j(x) = 1$$

$$T_j(x) = x_j \quad \sum_j x_j = 1$$

affine linear dep. between components of T

\Rightarrow multiple η 's give rise to same distribution; "overparameterization"

\rightarrow not a minimal exp. family

* for multinomial, minimal exp family

$$T(x) = \begin{pmatrix} x_1 \\ \vdots \\ x_{k-1} \end{pmatrix}$$

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

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

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

for multinomial, $\partial_{\eta_j} A(\eta) = \frac{1}{\sum_{j=1}^{k-1} e^{\eta_j} + 1} \exp(\eta_j) = p(x=j | \eta) = \mathbb{E}_{p(x; \eta)} [x_j | x]$ as required

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

$T: x \mapsto T(x)$ is "sufficient" for a parametric family P_θ

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

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

i.i.d.s model in exp. family

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

↳ exp. family on $x_{1:n}$, with reference density $\tilde{h}(x) = \prod_{i=1}^n h(x_i)$

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

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

Example 2: (1D) Gaussian

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

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

$\Theta = (\mu, \sigma^2)$ "moment parameterization"

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

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

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

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