

Lecture 18 - MaxEnt duality

November 9, 2021 4:34 PM

today : - Max. Ent. duality
- exponential families

dual problem for max. entropy

Max. ENT. in primal form

$$\min_{q \in \mathbb{R}^K} \left\{ \begin{array}{l} \sum_x q(x) \log \frac{q(x)}{u(x)} \\ q(x) > 0 \forall x \\ \sum_x q(x) = 1 \\ \sum_x q(x) T_j(x) = \alpha_j \forall j \end{array} \right. \Delta(x) \text{ eq. constraints}$$

$$u(x) \stackrel{!}{=} \frac{1}{|X|} = \frac{1}{K}$$

absorb this constraint in domain of def. of KL(q||u)

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

$$\mathcal{L}(q, \nu, c) = \sum_x q(x) \log \frac{q(x)}{u(x)} + \sum_j \nu_j (\alpha_j - \sum_x q(x) T_j(x)) + c (1 - \sum_x q(x))$$

$\mathbb{E}[q T_j] = \alpha_j$
 $\sum_x q(x) = 1$
 k-dim d-dim scalar

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

$$\nu^T T(x) = \sum_j \nu_j T_j(x)$$

$$\Rightarrow q_{\nu, c}^*(x) = u(x) \exp(\nu^T T(x) + c - 1)$$

generalized Max-ENT KL(q||h₀)

exponential family?

dual feat.: plug $q_{\nu, c}^*$ in $\mathcal{L}(\dots)$

$$g(\nu, c) = \mathcal{L}(q_{\nu, c}^*, \nu, c)$$

$$\mathbb{E}_q[S(x)] \triangleq \sum_x q(x) S(x)$$

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

$$\sum_x u(x) \exp(\nu^T T(x)) \exp(c-1)$$

$$\triangleq Z(\nu)$$

max $g(\nu, c)$ with respect to c

$$\nabla_c = 0 \Rightarrow 1 - Z(\nu) \exp(c-1) \stackrel{\text{want}}{=} 0$$

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

$$\text{plug back } c^*: \max_c g(\nu, c) = \nu^T \alpha + c^* - \frac{1}{Z(\nu)}$$

$$c^* - 1 = -\log Z(\nu)$$

$$\text{dual problem } \max_{\nu} \tilde{g}(\nu) \quad \left[\tilde{g}(\nu) \triangleq \nu^T \alpha - \log Z(\nu) \right]$$

link with M.L.F.:

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

then $\tilde{g}(\eta) = \frac{1}{n} \sum_{i=1}^n [\underbrace{\eta^T T(x^{(i)}) - \log z(\eta)}_{\log p(x^{(i)}|\eta) + \text{const.}}]$

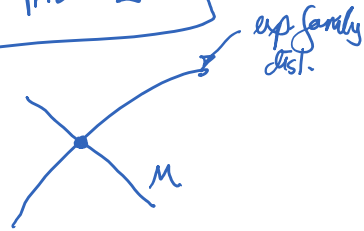
where $p(x|\eta) \triangleq \frac{1}{z(\eta)} u(x) \exp(\eta^T T(x))$

ie. dual problem is $\max_{\eta} \tilde{g}(\eta) = \max_{\eta} \frac{1}{n} \log p(x_{1:n}|\eta)$ ie. **MLE**

to summarize: ML in exp. family with $T(x)$ as suff. statistics is equivalent to Max ENT with moment constraints on $T(x)$ where $\alpha = \mathbb{E}_{p_n} [T(X)]$

they are Lagrangian dual of each other?

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



note [direct derivation]:

$$\begin{aligned} \nabla_{\eta} \log z(\eta) &= \frac{1}{z(\eta)} \nabla_{\eta} \left(\sum_x u(x) \exp(\eta^T T(x)) \right) \\ &= \sum_x T(x) \frac{u(x) \exp(\eta^T T(x))}{z(\eta)} \\ &= \sum_x T(x) p(x|\eta) \end{aligned}$$

$\nabla_{\eta} \log z(\eta) = \mathbb{E}_{p(x|\eta)} [T(X)] \triangleq \mu(\eta)$ "model moment"

$\nabla_{\eta} \tilde{g}(\eta) = \underbrace{\alpha}_{\mathbb{E}_{p_n} [T(X)]} - \mu(\eta)$

μ_n "empirical moment"

want $\nabla_{\eta} \tilde{g}(\eta) = 0 \Rightarrow \boxed{\mu(\eta^*) = \mu_n}$ i.e. moment matching?

→ see lecture 16 in 2017 for "KL Pythagorean theorem"

17h41

Exponential family

a (flat/canonical) exponential family on X

is a parametric family of dist. on X defined by two quantities:

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

reference density base measure

counting measure (discrete R.V.)
Lebesgue measure (cts. R.V.)

$\rightarrow \sum_x$ pmf
 $\rightarrow \int_x$ pdf

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

members of the family will have pmf/pdf

$$p(x; \eta) d\mu(x) = \exp(\underbrace{\eta^T T(x)}_{\text{"canonical parameter"}} - A(\eta)) \underbrace{h(x) d\mu(x)}_{\text{defining piece } (+ \Omega_X)}$$

"canonical parameter"

log normalization

or log partition fun.

or cumulant generating function

If Ω_X is discrete, then $p(x; \eta)$ is a pmf

" is cts. " " " " pdf

(for binary example $d\mu(x)$ is counting)

$$\delta_{x_0} + \delta_{x_1}$$

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

divide $\left[\sum_x p(x; \eta) \right]$

$$\Rightarrow A(\eta) \triangleq \log \left(\underbrace{\int_X \exp(\eta^T T(x)) h(x) d\mu(x)}_{Z(\eta)} \right)$$

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

Set of valid canonical parameters

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

⊛ more generally, consider a reparameterization of a subset of the flat family

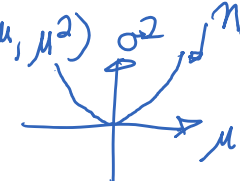
by defining $\eta: \Theta \rightarrow \Omega$

↑ new set of parameters

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

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

↳ e.g. could consider Gaussians where $\mathcal{N}(\mu, \mu^2)$



$\mathcal{N}(\Theta)$

* note: any single dist. $p(x)$ can be put in an exp. family by using $h(x) = p(x)$

* two examples of family not an exp. family: • mixture of Gaussians (latent variable model)
• unif $(0, \Theta)$

Example: (Multinoulli)

$$X \sim \text{Mult}(\pi) \quad X = \{0, 1\}^k$$

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

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's' = $\exp(\sum_j x_j \log \pi_j - 0)$

we have $\eta_j(\pi) = \log \pi_j$ $\Omega_X \subseteq \mathbb{R}^k$

$T(x) = x$ $d\mu(x) = \text{counting measure on } X$

$$h(x) = \mathbb{1}\{x \in \Omega_X\} = \mathbb{1}\{x \in \Delta_K \cap X\}$$

$$\mathcal{H} = \text{int}(\Delta_K) \quad A(\eta(\pi)) = 0 \quad \forall \pi \in \mathcal{H}$$

$$\begin{array}{l} \mathcal{H} \rightarrow \text{dimension } k-1 \\ \eta(\mathcal{H}) \rightarrow \text{" " " " } \\ \Omega_X \rightarrow \text{" " } k \end{array}$$

we do not have a "minimal exponential family"

note: ^{here} for any x st. $h(x) \neq 0$

$$\text{here } \underbrace{\sum_{j=1}^k T_j(x)} = \sum_j x_j = 1$$

affine linear dep. between components of T
 \Rightarrow multiple η 's give rise to same dist.
 (overparameterization)

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

for multinomial, $\frac{\partial A}{\partial \eta_j} = \frac{1}{z(\eta)} e^{\eta_j} = p(x=j | \eta)$

$= \mathbb{E}_{p(x; \eta)} [T_j(x)]$ as required //

verified