

Lecture 17 - scribbles

Tuesday, November 12, 2019 14:30

- max Ent duality
- exponential family

quick presentation of Lagrangian duality

convex optimization problem :

$$\min_x f(x)$$

• f, f_j are convex fct. s.t. $f_j(x) \leq 0 \quad \forall j=1, \dots, m$
 • g_k affine fct. $g_k(x) = 0 \quad \forall k=1, \dots, n$

} "primal problem"

Lagrangian fct. $\mathcal{L}(x, \lambda, \nu) \triangleq f(x) + \sum_{j=1}^m \lambda_j f_j(x) + \sum_{k=1}^n \nu_k g_k(x)$

"Lagrange multipliers"

magic trick (saddle pt. interpretation) of Lagrangian duality

$$h(x) = \sup_{\lambda \geq 0, \nu} \mathcal{L}(x, \lambda, \nu) = \begin{cases} f(x) & \text{if } x \text{ is feasible} \\ +\infty & \text{if " " not feasible} \end{cases}$$

an equivalent problem to primal problem

$$\inf_x \left(\sup_{\lambda \geq 0, \nu} \mathcal{L}(x, \lambda, \nu) \right)$$

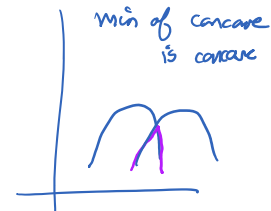
} $h(x)$ fancy non-smooth fct.

duality trick is to swap inf & sup

$$\sup_{\lambda \geq 0, \nu} \inf_x \mathcal{L}(x, \lambda, \nu) \triangleq g(\lambda, \nu)$$

} Lagrange dual fct.

→ this fct. is always concave



Lagrangian dual problem

$$\sup_{\lambda, \nu} g(\lambda, \nu) \quad \lambda \geq 0$$

} "dual variables"

weak duality (always true)

in-general $\sup \inf \mathcal{L}(x, \lambda, \nu) \leq \inf \sup \mathcal{L}(x, \lambda, \nu)$

strong duality when $\sup \inf \mathcal{L} = \inf \sup \mathcal{L}$

(or when primal is convex + constraint qualification condition (e.g. Slater's condition))

(can get optimal primal variables $z^*(\lambda^*, v^*)$ using KKT conditions)

see chapter 5 in Boyd's book for more info on duality: <http://stanford.edu/~boyd/cvxbook/>

dual problem for max. entropy

MaxENT in 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.

$$\sum_x q(x) = 1$$

feature set.

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

KL(q||u) = $\begin{cases} \text{too big } q(x) < 0 \\ \text{sum } \sum_x q(x) < 1 \\ \text{KL}(q||u) < 0 \end{cases}$ a.u.

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

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

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

exponential family?

dual f.f.:

plug back q^* in $g(\dots)$

$$g(v, c) = g(q^*, v, c)$$

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

$$= v^T \alpha + c - \underbrace{\sum_x u(x) \exp(v^T T(x)) \exp(c-1)}_{\triangleq Z(v)} + c - \mathbb{E}_{q^*} [c]$$

$$\max g(v, c)$$

with respect to c

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

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

plug back c^* :

$$\max g(v, c) = v^T \alpha + c^* - Z(v) \frac{1}{Z(v)}$$

$$C^{-1} = -\log(Z(\nu))$$

dual problem $\max_{\nu} \tilde{g}(\nu)$

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

link with MLE:

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

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

$\log p(x_i|\nu) + \text{const. where } p(x|\nu) \propto \exp(\nu^T T(x)) \cdot \log Z(\nu)$

ie. dual problem is $\max_{\nu} \tilde{g}(\nu) = \max_{\nu} \frac{1}{n} \log p(x_{1:n}|\nu)$ ie. 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 exp. family \Leftrightarrow moment matching in exp. family

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

$p(x|\nu)$

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

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

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

\rightarrow see lecture 16 in 2017 for "KL Pythagorean thm."

15h40

(see end of old lecture 16 2017 for "KL Pythagorean theorem" and I-projection vs. M-projection for KL + geometry)

Exponential family

a (flat/canonical) exponential family on \mathcal{X}

is a parametric family of distributions defined by two quantities

is a parametric family of distributions defined by two quantities

I) $h(x) d\mu(x)$ → reference measure on X
reference density measure ← counting (discrete R.V.)
Lebesgue (cts. R.V.)

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

members of family will have dist:

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

log normalizer or cumulant generating func. log partition set.

if Ω_X is discrete, then $p(x; \eta)$ is a pmf
 " " cts. " " " pdf

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

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

Ω_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

⇒ Ω is convex

⊗ more generally, consider a reparametrization of a subset of the family

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

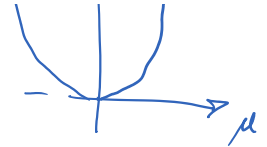
new set of parameters

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

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

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





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

* two examples of family not an exp. family:

- mixture of Gaussians (latent variable models)
- $\text{unif}(0, \sigma)$

Example: (Multinoulli)

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

$$\Omega_X = \Delta_K \cap X \quad (\text{one encoding})$$

parameter $\pi \in \Delta_K$; suppose $\pi_j > 0 \forall j$

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

think as "6"

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

$$\text{we have } \eta_j(\pi) = \log \pi_j$$

$$T(x) = x$$

$d_{\mu}(x)$ = counting measure on X

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

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

$$\text{here, } A(\eta(\pi)) = 0 \quad \forall \pi \in \Theta$$

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

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

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

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

we do not have a "minimal exponential family"

note: for any x s.t. $h(x) \neq 0$

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

$$\text{here, } \sum_{j=1}^k T_j(z) \stackrel{\leftarrow}{=} x_j = 1$$

affine linear dep. between components of T

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

\hookrightarrow 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 $\eta \in \text{int}(\Omega)$)

$$\begin{aligned} \text{for multinomial, } \frac{\partial A(\eta)}{\partial \eta_j} &= \frac{\partial \log Z(\eta)}{\partial \eta_j} = p(x=j | \eta) \\ &= \mathbb{E}_{p(x; \eta)} [T_j(x)] \text{ as required} \end{aligned}$$