

## Lecture 16 - scribbles

Friday, November 3, 2017

13:33

today: max. entropy  
equivalence with ML  
duality  
exponential family

recall: last time  $ML \Leftrightarrow \min_{\theta \in \Theta} KL(\hat{p}_n \| p_\theta)$

today: maximum entropy  $\min_{q \in M} KL(q \| \text{uniform})$

Maximum entropy principle:

idea: consider some subset of dist. on  $X$   
according to some data-driven information

get a subset  $M \subseteq \Delta_{|X|}$

principle pick  $\hat{p} \in M$  by maximizing entropy

ie,  $\hat{p} = \arg \max_{q \in M} H(q)$

$= \arg \min_{q \in M} KL(q \| \text{uniform})$

→ "generalized max entropy"  
which uses  $h_0$  instead of uniform  
& favorite a priori dist.

$$\sum_{z \in X} q(z) \log_{\text{const.}} q(z) = -H(q) + \text{const.}$$

\* example from Wainwright

$P_L = \frac{3}{4}$  kangaroos are left-handed

$P_B = \frac{2}{3}$  " drink Foster beer

question: how many " are both left-handed & drink F. beer

[ here: max entropy solution is that  $p(B, L) = P_B \cdot P_L$  (independence) ]

\* how do we get M?

typically: through empirical observation

feature functions  $T_1(x), \dots, T_d(x)$

$$\text{define } M = \left\{ q : \underbrace{\mathbb{E}_q [T_j(x)]}_{\text{model expected feature count}} = \underbrace{\mathbb{E}_{p_n} [T_j(x)]}_{\text{empirical feature count}} \quad j=1, \dots, d \right\}$$

"moment constraints"

then Max ENT:  $\min_{q \in \mathbb{R}^{1 \times X}} KL(q || \text{unif})$  some scalar

st.  $q \in M$  }  $\sum_x q(x) T_j(x) = \alpha_j \quad \forall j$

$q \in \Delta_{|X|}$  } i.e.  $\langle \vec{q}, \vec{T}_j \rangle = \alpha_j$

→ convex opt. problem over  $q \in \Delta_{|X|} \subseteq \mathbb{R}^{1 \times X}$

Lagrange duality segway:

convex optimization problem

$$\min_x f(x)$$

st  $f_i(x) \leq c_i \quad i=1, \dots, m$

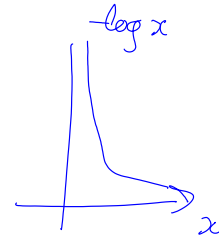
} "primal problem"

$$\text{s.t. } f_j(x) \leq 0 \quad j=1, \dots, m \\ g_k(x) = 0 \quad k=1, \dots, n$$

- $f_j, f_i$  are convex functions
- $g_k$  to affine

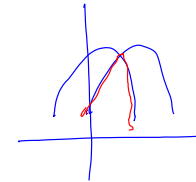
here,  $f: \mathbb{R}^d \rightarrow \mathbb{R} \cup \{\pm\infty\}$  "extended real valued function"

$$\text{dom}(f) \triangleq \{x: f(x) < \infty\} \quad \text{e.g. } f(x) = \begin{cases} -\log x & x > 0 \\ +\infty & \text{o.w.} \end{cases}$$



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 point interpretation)

$$\sup_{\lambda \geq 0, \nu} \mathcal{L}(x, \lambda, \nu) = \begin{cases} f(x) & \text{if } x \text{ is feasible} \\ +\infty & \text{o.w.} \end{cases}$$

an equivalent problem to primal problem is

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

duality trick is swap  $\inf$  &  $\sup$

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

$$\triangleq g(\lambda, \nu) \quad \text{Lagrange dual fct.}$$

dual problem

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

$$\sup_{\lambda, \nu} g(\lambda, \nu)$$

vector inequality

always true:

Weak duality

$$\sup_{\lambda \geq 0, \nu} \inf_x f(x, \lambda, \nu) \leq \inf_x \sup_{\lambda \geq 0, \nu} f(x, \lambda, \nu)$$

always true

if  $p^* = \inf_{x \text{ feasible}} f(x)$   
 global of primal

$$g(\lambda, \nu) \leq p^* \quad \forall \lambda \geq 0, \nu \text{ feasible dual variables}$$

always concave in  $\lambda, \nu$

Strong duality  $\rightarrow$  when have equality i.e.  $d^* = \sup_{\lambda \geq 0, \nu} g(\lambda, \nu) = p^*$

a sufficient condition for strong duality

is Slater's condition:  $\exists x \in \text{int}(\text{dom}(f))$

s.t.  $f_j(x) < 0 \quad \forall j$  where  $f_j$  is nonlinear  
 and  $x$  is feasible

KKT conditions:

when  $f$  etc. are differentiable

necessary conditions for strong duality

and  $x^* \{ \lambda^*, \nu^* \}$  are respectively primal and dual optimal

$$g(\lambda^*, \nu^*) = f(x^*) = \inf_x f(x, \lambda^*, \nu^*)$$

$$\nabla f(x^*) + \sum \lambda_j^* \nabla f_j(x^*) + \sum \nu_k^* \nabla g_k(x^*) = 0$$

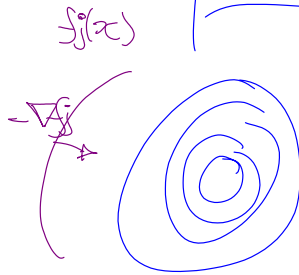
necessary conditions are

KKT conditions

see chapter 5 in Boyd's book: <http://stanford.edu/~boyd/cvxbook/>

if primal is convex  
+ Slater's condition  
they are also sufficient

complementary slackness:  $\lambda_j^* f_j(x^*) = 0$   
 $x^*$  is primal feasible  
 $(\lambda^*, v^*)$  is dual feasible (i.e.  $A^* \geq 0$ )



dual problem for Max ENT

Max ENT primal form (P)

$$\min_q \begin{cases} \sum_x q(x) \log \frac{q(x)}{u(x)} \\ q(x) \geq 0 \\ \sum_x q(x) = 1 \\ \sum_x q(x) T_j(x) = a_j \end{cases} \Delta x \quad \left. \vphantom{\sum_x} \right\} M$$

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

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

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

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

## Exponential family?

(we have strong duality by Slater  
if  $\exists q \in M$  st.  $q(x) > 0 \forall x$ )

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

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

maximize  
with respect to  
 $c$

$$\nabla_c \rightarrow 1 - z e^{c-1} \stackrel{\Delta}{=} z(v) = 0$$

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

plug back  $\max_c g(v, c) = v^T \alpha - \underbrace{\log z(v)}_{(c^*-1)} \stackrel{\Delta}{=} \tilde{g}(v)$

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

then  $\tilde{g}(v) = \frac{1}{n} \sum_{i=1}^n \underbrace{[v^T T(x_i) - \ln z(v)]}_{\log p(x_i|v)}$

where  $p(x|v) \stackrel{\Delta}{=} \exp(v^T T(x) - \ln z(v))$

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

to summarize: ML in the exponential family with  $T(x)$  as sufficient statistics

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

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

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

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

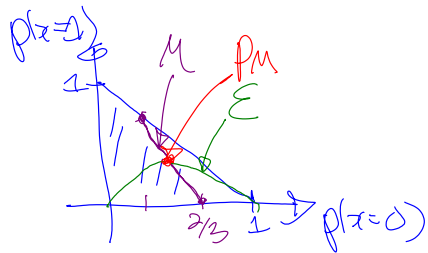
$\nabla_{\nu} \tilde{g}(\nu) = \mathbb{E}_{p_n} [T(x)] - \mu(\nu)$

$\nabla_{\nu} \tilde{g}(\nu) = 0 \Rightarrow \mu(\nu) = \mathbb{E}_{p_n} [T(x)]$  i.e. moment matching?  
 $\mathbb{E}_{p(x|\nu)} [T(x)]$

geometry:

$X = \{0, 1, 2\}$

$T(x) = \begin{cases} 1 & \text{if } x=0 \\ \nu_2 & \text{if } x=1 \\ 0 & \text{if } x=2 \end{cases}$

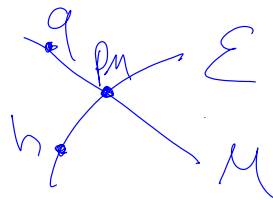


$\mathbb{E}_{p_n} [T(x)] = 2 = \alpha$

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

"information Pythagorean thm.":

for any  $q \in \mathcal{M}$  and  $h \in \mathcal{E}$   
 $KL(q||h) = KL(q||p_n) + KL(p_n||h)$



restate our duality result :

$$p_M = \operatorname{argmin}_{q \in \mathcal{M}} \text{KL}(q \| u) \quad [\text{Max ENT}]$$

"I-projection"  $I \rightarrow$  information

$p_M = \operatorname{argmin}_{q \in \mathcal{E}} \text{KL}(\hat{p}_M \  q)$	MLE in exp-family	
	"M-projection" $M \rightarrow$ moment	