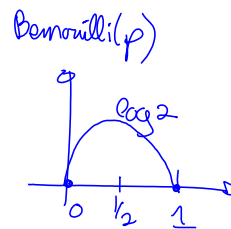


Lecture 15 - scribbles

Friday, October 21, 2016
13:35

today: KL divergence
exp family and duality

recall: Entropy $H(p) \triangleq - \sum_{x \in \mathcal{X}} p(x) \log p(x)$



in coding theory, use length of code $\propto -\log p(x)$

$\log_2 \rightarrow$ "bits"
 $\log_e \rightarrow$ "nats"

KL divergence: excess cost (in terms length of code)
to use distribution q for coding instead of p

$$KL(p, q) \triangleq \sum_x p(x) \log \frac{p(x)}{q(x)}$$

entropy of uniform on k states $\rightarrow - \sum_x \frac{1}{k} \log \frac{1}{k}$

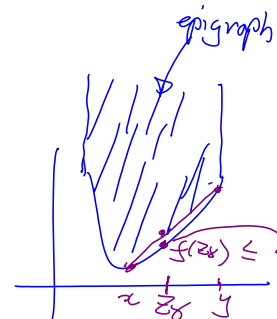
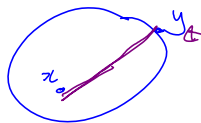
$$\hookrightarrow H(p) = \log k$$

properties of KL:

- $KL(p, q) \geq 0$
- is strictly convex in each of its argument
i.e. $KL(p, \cdot)$ is strictly convex
as fct. of 2nd argument

$KL(\cdot, q)$ " " "

convex set $A \Leftrightarrow \forall x, y \in A, (1-\delta)x + \delta y \in A$
for $\delta \in [0, 1]$



convex function $f \Leftrightarrow$ its epigraph is convex

$$\hookrightarrow \{(x, t) : x \in \text{dom}(f), t \geq f(x)\}$$

$$\Leftrightarrow f(\underbrace{(1-\delta)x + \delta y}_{z_\delta}) \leq (1-\delta)f(x) + \delta f(y)$$

convex
strictly convex



if f is diff. $\rightarrow f(u) \geq f(x) + \langle \nabla f(x), u-x \rangle$

If f is diff. $\Rightarrow f(y) \approx f(x) + \langle \nabla f(x), y-x \rangle$

ML and KL:

define empirical distribution given $\{x_i\}_{i=1}^n$ observations
 $\hat{p}_n \triangleq \frac{1}{n} \sum_{i=1}^n \delta(x, x_i)$ Kronecker-delta fct.

prop.: $\{P_\theta\}_{\theta \in \Theta}$ parametric family on X

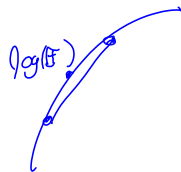
then ML for $P_\theta \Leftrightarrow \min_{\theta} KL(\hat{p}_n \parallel P_\theta)$

$$\begin{aligned} KL(\hat{p}_n \parallel P_\theta) &= \sum_x \hat{p}_n(x) \log \frac{\hat{p}_n(x)}{P_\theta(x)} \\ &= -H(\hat{p}_n) - \sum_x \hat{p}_n(x) \log P_\theta(x) \\ &= -H(\hat{p}_n) - \sum_{i=1}^n \frac{1}{n} \log P_\theta(x_i) \\ &= -H(\hat{p}_n) - \frac{1}{n} \sum_{i=1}^n \log P_\theta(x_i) \\ &= -H(\hat{p}_n) - \log \left(\prod_{i=1}^n P_\theta(x_i) \right) \end{aligned}$$

also: can see ML as ERM for $L(P, P_\theta) = \mathbb{E}_{x \sim P} [-\log P_\theta(x)]$

$$-KL(p, q) = \sum_x p(x) \log \frac{q(x)}{p(x)}$$

$$\begin{aligned} &\mathbb{E}_p \left[\log \left(\frac{q(x)}{p(x)} \right) \right] \\ &\stackrel{\text{Jensen's}}{\leq} \log \mathbb{E}_p \left[\frac{q(x)}{p(x)} \right] \\ &= \log \sum_x p(x) \frac{q(x)}{p(x)} = \log \sum_x q(x) = \log 1 = 0 \end{aligned}$$



$$-KL(p, q) \leq 0$$

Maximum entropy principle

idea: consider some subset of distribution on X according to some data-driven information

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

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

$$\hat{p} = \operatorname{argmax}_{q \in M} H(q)$$

$$= \operatorname{argmin}_{q \in M} KL(q \parallel \text{uniform})$$

could use any h_0 is fixed
 ("prior" data is general
 ("generalized")
 ("maximum-entropy")

$$q \in M \quad \sum_x q(x) \log \frac{q(x)}{cst.} \quad \text{max. entropy}$$

$$= -H(q) + cst.$$

example by Martin 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.B. . . .

(here Max-ent. solution is that $P(B, L) = P_B \cdot P_L$ (independence)]

* where do we get M ?

typically: through empirical observations

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

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

"moment constraints"

so Max-entropy principle: $\min_{\substack{q \in \Delta_{|X|} \\ q \in M}} KL(q \parallel \text{unif.})$

$$q \in M \Rightarrow \sum_x q(x) T_j(x) = \alpha_j$$

$$\langle \vec{q}, \vec{T}_j \rangle = \alpha_j$$

→ convex optimization problem over $q \in \Delta_{|X|} \subseteq \mathbb{R}^{|X|}$

Lagrangian duality:

convex optimization problem $\min_x f(x)$
 s.t. $f_j(x) \leq 0 \quad j=1, \dots, m$
 $g_k(x) = 0 \quad k=1, \dots, n$ } primal problem

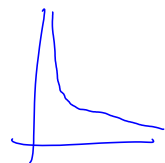
• f, f_j are convex fct.

• g_k is affine

here, $f: \mathbb{R}^d \rightarrow \mathbb{R} \cup \{+\infty\}$

$\text{dom}(f) \triangleq \{x: f(x) < +\infty\}$

$$f(x) = \begin{cases} +\infty & x \leq 0 \\ -\log x & x > 0 \end{cases}$$



(extended real valued fct.)

$$\dots \sum_{j=1}^m \lambda_j f_j(x) + \sum_{k=1}^n \mu_k g_k(x) \dots$$

maximization problem (primal)

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

\uparrow one "Lagrange" multipliers
 \uparrow

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

an equivalent problem to primal is

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

duality trick is swap inf & sup

$$\sup_{\lambda \geq 0, \nu} \inf_x L(x, \lambda, \nu) \triangleq g(\lambda, \nu) \quad \text{Lagrange dual fun.}$$

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

Weak duality:

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

always true

moreover, $g(\lambda, \nu) \leq p^* \quad \forall \lambda \geq 0, \nu$ where p^* is global min of primal

strong duality \rightarrow have equality

a sufficient condition for strong duality is

Slater's condition: $\exists x \in \text{int}(\text{dom}(f))$ s.t. $f_j(x) < 0 \quad \forall j$
and x feasible

KKT conditions:

- necessary conditions for strong duality: $g(\lambda^*, \nu^*) = f(x^*) = \inf_x L(x, \lambda^*, \nu^*)$
- and $x^* \in \{ (\lambda^*, \nu^*) \}$ are respectively primal optimal and dual optimal
- complementary slackness: $\lambda_j^* f_j(x^*) = 0$



$$\inf_x \rightarrow \nabla f(x) + \sum_j \lambda_j \nabla f_j(x) + \sum_k \nu_k \nabla g_k(x) = 0$$

in other words, for differentiable fun. necessary conditions for x^* be optimal and (λ^*, ν^*) be dual optimal is:

- x^* primal feasible
 - (λ^*, ν^*) dual " ($\lambda^* \geq 0$)
 - $\nabla f(x^*) + \sum_j \lambda_j^* \nabla f_j(x^*) + \sum_k \nu_k^* \nabla g_k(x^*) = 0$
 - $\lambda_i^* f_i(x^*) = 0$
- and if primal is convex, also sufficient

condition

↳ KKT conditions

dual problem for Max Ent:

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

$$f(x, \lambda, \nu) = \sum_x q(x) \log \frac{q(x)}{u(x)} + \sum_j \nu_j (\alpha_j - \sum_x q(x) T_j(x)) + C (1 - \sum_x q(x))$$

$$\frac{\partial}{\partial q(x)} = \begin{array}{l} 1 + \log q(x) \\ - \log u(x) \end{array} - \sum_j \nu_j T_j(x) - C = 0$$

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

(we have strong duality by Slater exponential family)
 if $\exists q \in \mathcal{M}$ s.t. $q(x) > 0 \forall x$)