

today: information theory
Max Ent.
duality

Information theory

KL divergence: for discrete dist. $p \neq q$

$$KL(p \parallel q) \triangleq \mathbb{E}_p \left[\log \frac{p(x)}{q(x)} \right] = \sum_{x \in \Omega} p(x) \log \frac{p(x)}{q(x)}$$

convention: $0 \cdot \log 0 = 0$
($\lim_{x \rightarrow 0^+} x \log x = 0$)

[if $\exists x$ s.t. $q(x) = 0$
but $p(x) \neq 0$]

$\rightarrow p(x) \log \frac{p(x)}{q(x)} = +\infty$

if support of $p \not\subseteq$ support of $q \Rightarrow KL(p \parallel q) = +\infty$

motivation from density estimation

recall statistical decision theory

(statistical) loss $L(p_\theta, a)$ ^{world} here, density estimation, say \hat{q}

(cross-entropy loss)

Standard (MLE) loss is log-loss $L(p_\theta, \hat{q}) = \mathbb{E}_{x \sim p_\theta} [-\log \hat{q}(x)]$

if use $\hat{q} = p_\theta$, then get

$$\sum_{x \in \Omega} -p_\theta(x) \log p_\theta(x) \triangleq H(p_\theta)$$

entropy of p_θ

(statistical) excess loss for action $a = \hat{q}$

$$L(p, \hat{q}) - \min_{q \in \mathcal{Q}} L(p, q) = \sum_x p(x) \log \frac{q(x)}{p(x)} = KL(p, \hat{q})$$

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

Coding theory:

use length of code $\propto -\log p(x)$

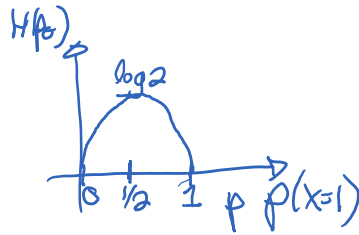
expected length of code: $\sum_x p(x) (-\log p(x)) \rightarrow$ entropy of p measure in bits

KL divergence \rightarrow interpreted as excess cost (in terms of length of code) to use dist. q to design code vs

the optimal dist. (true p)

Example of H :

entropy for a Bernoulli
 $-p \log p - (1-p) \log(1-p)$



entropy for a uniform dist. on k states

$$-\sum_{\alpha=1}^k \frac{1}{k} \log \frac{1}{k} = \log k$$

(max entropy dist. over k states)

properties of KL divergence:

• $KL(p||q) \geq 0$ ← to show this, use Jensen's ineq. $f(\mathbb{E}X) \leq \mathbb{E}f(X)$ when f convex

• KL is strictly convex in each of its arguments

ie. $KL(p||\cdot) = \Delta_K \subseteq \mathbb{R}^k \rightarrow \mathbb{R}$

• $KL(p||p) = 0 \quad \forall p \in \Delta_K$

$KL(\cdot||q)$

} strictly convex set.

• not symmetric $KL(p||q) \neq KL(q||p)$ in general

Symmetrized version $\frac{1}{2}(KL(p||q) + KL(q||p))$

MLE & KL minimization

$\{p_\theta\}_{\theta \in \Theta}$ parametric family of dist. for a discrete obs space

then $MLE \text{ for } \theta \Leftrightarrow \min_{\theta \in \Theta} KL(\hat{p}_n || p_\theta)$

empirical dist. $\hat{p}_n(x) \triangleq \frac{1}{n} \sum_{i=1}^n \delta(x, x^{(i)})$
 ↑ Kronecker-delta

proof:

$$KL(\hat{p}_n || p_\theta) = \sum_x \hat{p}_n(x) \log \frac{\hat{p}_n(x)}{p_\theta(x)}$$

$$= -H(\hat{p}_n) - \sum_x \hat{p}_n(x) \log p_\theta(x)$$

$$\frac{1}{n} \sum_{i=1}^n \sum_x \delta(x, x^{(i)}) \log p_\theta(x)$$

$$= \underbrace{-H(\hat{p}_n)}_{\text{const. w.r. to } \theta} - \frac{1}{n} \sum_{i=1}^n \log p_\theta(x^{(i)}) = \log \prod_{i=1}^n p_\theta(x^{(i)}) //$$

const. w.r. to σ $\log \prod_i p(x^{(i)})$ //

Maximum entropy principle:

idea: consider some subset of dist. over X according to some data-driven constraints

get a subset $M \subseteq \Delta(X)$ & proba samples over $|X|=k$ elements

MAXENT principle $\hat{p} \in M$ which maximizes the entropy

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

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

$KL(q \parallel u) = \sum_x q(x) \log \frac{q(x)}{u(x)} \cdot 1/k = \text{constant} = -H(q) + \text{const.}$

"generalized max. entropy" $KL(q \parallel p_0)$

preferred dist. to be biased towards

16/07

* example from Weinwright

$\hat{p}_L = \frac{3}{4}$ kangaroos are left-handed

$\hat{p}_B = \frac{2}{3}$ kangaroos drink Sabatt beer

question: how many kang. are both L.H. & drink Sab. beer? (what proportion)

[here: max-entropy solution is that $p(B=1, L=1) = \hat{p}_B \cdot \hat{p}_L$ (indep.)]

* how do we get set M ?

typically: through empirical "moments"

kangaroo
 $T_1(x) = \mathbb{1}_{\{x \text{ drinks Sabatt}\}}$
 $T_2(x) = \mathbb{1}_{\{x \text{ is L.H.}\}}$

feature functions: $T_1(x), \dots, T_d(x)$ of features

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

then $\boxed{\operatorname{Max ENT} \min_{q \in M} KL(q \parallel \text{unif.}) \rightarrow \sum_x q(x) T_j(x) = \frac{1}{n} \sum_{i=1}^n T_j(x^{(i)}) = \alpha_j}$

$$\min_{q \in \mathbb{R}^{|\mathcal{X}|}} \text{KL}(q \parallel \text{unif}) \quad \text{st. } q \in \mathcal{M} \subseteq \Delta_{|\mathcal{X}|}$$

$$\sum_x q(x) J_j(x) = \frac{1}{n} \sum_{i=1}^n J_j(x^{(i)}) = q_j$$

ie. $\langle \vec{q}, \vec{1}_j \rangle = q_j$

constraint.
 \hookrightarrow convex opt. problem over $q \in \Delta_{|\mathcal{X}|} \subseteq \mathbb{R}^{|\mathcal{X}|}$

Quick presentation of Lagrangian duality

convex min. problem

$$\min_{x \in \mathbb{R}^d} f(x)$$

convex problem $\left\{ \begin{array}{l} f, f_j \text{ are convex fct.} \\ g_k \text{ affine fct.} \end{array} \right.$ st. $\begin{array}{l} f_j(x) \leq 0 \quad j=1, \dots, m \\ g_k(x) = 0 \quad k=1, \dots, n \end{array}$ } "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)

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

an equivalent form of the primal problem: $h(x)$ is fancy non-smooth fct.

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

duality trick is to swap inf & sup

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

→ this is always concave (λ, ν)
 Lagrange dual fct.



Lagrange dual problem

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

"dual variables"

$$\sup \inf \mathcal{L} \leq \inf \sup \mathcal{L}$$

"weak duality"

in general, $\sup \inf \mathcal{L} \leq \inf \sup \mathcal{L}$

dual problem

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

→ sufficient conditions: • when primal problem is convex
 both constraint qualif. cond. (e.g. Slater's condition)

(can get optimal primal variables $x^*(x^*, \nu^*)$)



using KKT conditions)

(see ch. 5 of Boyd's book)

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