

today: info theory & KL
 max entropy
 equivalence with ML in exp. family & duality

Information theory

KL divergence:

for discrete dist. p & q

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

$$0 \cdot \log 0 = 0$$

$$\left(\lim_{x \rightarrow 0^+} x \log x = 0 \right)$$

motivation from density estimation

recall statistical decision theory
 world

(statistical) loss $L(p_0, a)$

↳ here, estimation of distribution, say \hat{q}

standard (ML) loss is log-loss $L(p_0, \hat{q}) = \mathbb{E}_{x \sim p_0} [-\log \hat{q}(x)]$

aside: "cross-entropy"

if use $\hat{q} = p_0$, then get

$$\sum_{x \in \Omega} -p_0(x) \log p_0(x) \triangleq H(p_0)$$

entropy of p_0

If use $q = p_0$, then get $\sum_{x \in \Omega_X} -p_0(x) \log p_0(x) = H(p_0)$ | entropy of p_0

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

$$L(p, \hat{q}) - \underbrace{\min_q L(p, q)}_{L(p, p)} = - \sum_x p(x) \log \frac{\hat{q}(x)}{p(x)} = KL(p \| \hat{q})$$

Coding theory:

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

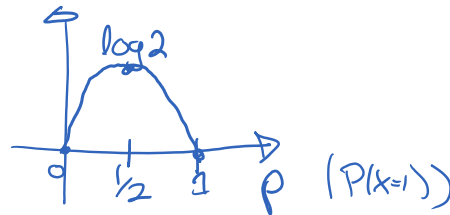
$\log \hat{=} \log_2 \rightarrow$ "bits" $\log_e \rightarrow$ "nats"

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

KL divergence \rightarrow interpreted as the excess cost (in terms of length of code) to use dist q for coding vs. the optimal dist. (true p)

example:

entropy for a Bernoulli:



entropy of uniform dist. on K states: $-\sum_{x=1}^K \frac{1}{K} \log \frac{1}{K} = \log K$

(max entropy dist. over K states)



(max entropy dist. over K states)



properties of KL:

- $KL(p||q) \geq 0$ ~~to~~ to show this, use Jensen's inequality $f(\mathbb{E}x) \leq \mathbb{E}f(x)$ when f is convex
 - is strictly convex in each argument
 - not symmetric: $KL(p||q) \neq KL(q||p)$
- i.e. $KL(p||\cdot) : \mathbb{R}^k \rightarrow \mathbb{R}$
 $KL(\cdot||q)$ } strictly convex functions
- $KL(\vec{p}||\vec{p}) = 0 \quad \forall \vec{p} \in \Delta_K$

ML and KL minimization:

$\{p_\theta\}_{\theta \in \Theta}$ parametric family

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

then ML for $\Theta \iff \min_{\theta \in \Theta} KL(\hat{p}_n || p_\theta)$

Kronecker delta

proof:

$$\begin{aligned}
 KL(\hat{p}_n || p_\theta) &= \sum_x \hat{p}_n(x) \log \frac{\hat{p}_n(x)}{p_\theta(x)} \\
 &= -H(\hat{p}_n) - \sum_x \hat{p}_n(x) \log p_\theta(x) \\
 &= \underbrace{-H(\hat{p}_n)}_{\text{constant}} - \frac{1}{n} \sum_{i=1}^n \log p_\theta(x^{(i)})
 \end{aligned}$$

14/2/21

$$\log \prod_{i=1}^n p(x^{(i)})$$



Maximum entropy principle:

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

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

MAX ENT principle: pick $\hat{p} \in M$ which maximizes the entropy

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

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

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

"generalized max. entropy" $KL(q \parallel h_0)$

↑
preferred dist. to bias towards

* example from Wainwright:

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

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

question: how many kang. are both left-handed & drink S. beer?

[here: max entropy sol'n is that $p(B, L) = P_B \cdot P_L$ (independent)]

* how do we get M ?

typically: through empirical "moments"

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

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

"moment constraints"

then MaxENT:

$$\begin{aligned} \min_{q \in \mathbb{R}^{|\mathcal{X}|}} & \text{KL}(q \parallel \text{unif}) \\ \text{s.t. } & q \in M \\ & q \in \Delta_{|\mathcal{X}|} \end{aligned} \quad \rightarrow \quad \sum_x q(x) T_j(x) = \frac{1}{n} \sum_i T_j(x^{(i)}) = \alpha_j$$

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

\leadsto convex opt. problem over $q \in \Delta_{|\mathcal{X}|} \subset \mathbb{R}^{|\mathcal{X}|}$

Lagrangian duality segway:

convex optimization problem:

$$\min_x f(x)$$

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

$$g_k(x) = 0 \quad \forall k=1, \dots, n$$

• f, f_j are convex functions

• g_k affine fets

"primal problem"

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

"extended real-valued fets"

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

$$\text{eg } f(x) = \begin{cases} -\log x & x > 0 \\ +\infty & \text{o.w.} \end{cases}$$



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

"Lagrange multipliers"

magic
trick
(saddle pt. interpretation
of Lagrangian duality)

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

an equivalent problem to primal problem

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

sup of convex fct. \Rightarrow convex
inf "concave" \Rightarrow concave

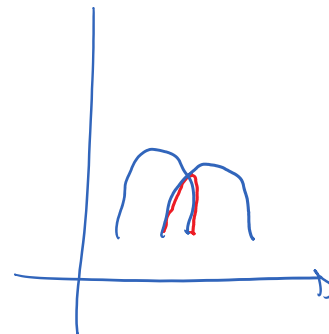
duality trick is swap inf & sup

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

duality trick is swap \inf & \sup $h(x)$ fancy non-smooth fct.

$$\sup_{\lambda \geq 0, v} \underbrace{\inf_x f(x, \lambda, v)}_{\triangleq g(\lambda, v)} \rightarrow \text{this fct. is always concave}$$

Lagrange dual fct.



Lagrangian dual problem

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

always true: weak duality: $\sup_{\lambda, v} \inf_x f(x, \lambda, v) \leq \inf_x \sup_{\lambda, v} f(x, \lambda, v)$

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

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

always concave in λ, v

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

a sufficient condition for strong duality

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

Slater condition
+ "

s.t. $f_j(x) < 0$ ($\forall j$ where f_j is non-linear)

Slater condition
 +
 convex primal
 \Rightarrow strong duality

s.t. $f_j(x) \leq 0$ ($\forall j$ where f_j is non-linear)
 and x is feasible

KKT conditions:

when f is differentiable

necessary conditions for x^* & v^*
 to be primal & dual optimal

and strong duality

[when x^* is int(dom(f))]

$$g(v^*) = f(x^*) = \inf_x \mathcal{L}(x, v^*)$$

\Downarrow

$$\nabla f(x^*) + \sum_k \nu_k^* \nabla g_k(x^*) = 0$$

$\leftarrow x^*$ is primal feasible

KKT conditions

if primal is convex problem + Slater's condition \Rightarrow sufficient conditions

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