

today:

- max Ent & duality
- exponential family

MLE & KL minimization

if $p_{\theta} \in \mathbb{G}$ parametric family for a discrete observation space

$$\begin{aligned} \text{Then ML for } \mathbb{G} \text{ for iid data} &\Leftrightarrow \min_{\theta \in \mathbb{G}} \text{KL}(\hat{p}_n \| p_{\theta}) \\ &\hat{p}_n(x) \stackrel{\Delta}{=} \frac{1}{n} \sum_{i=1}^n \delta(x, x^{(i)}) \quad \text{"empirical distribution"} \quad \text{kronencker-delta} \end{aligned}$$

$$\begin{aligned} \text{proof: } \text{KL}(\hat{p}_n \| p_{\theta}) &= \sum_x \hat{p}_n(x) \log \frac{\hat{p}_n(x)}{p_{\theta}(x)} \\ &= -H(\hat{p}_n) - \sum_x \underbrace{\hat{p}_n(x) \log p_{\theta}(x)}_{\stackrel{\Delta}{=} \sum_i \delta(x, x^{(i)})} \\ &\quad \underbrace{\sum_{i=1}^n \sum_x \delta(x, x^{(i)}) \log p_{\theta}(x)}_{\log p_{\theta}(x^{(i)})} \\ &= -H(\hat{p}_n) - \underbrace{\sum_{i=1}^n \log p_{\theta}(x^{(i)})}_{\text{constant w.r.t. } \theta} \quad \log \left(\prod_{i=1}^n p_{\theta}(x^{(i)}) \right) // \end{aligned}$$

Maximum entropy principle

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

get a subset $M \subseteq \Delta_{|X|}$ ← prob. simplex over $|X|=k$ elements

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

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

$$= \underset{q \in M}{\operatorname{argmin}} \text{KL}(q \| \text{uniform})$$

constant

$$KL(q \parallel h) = \sum_x q(x) \log \frac{q(x)}{h(x)} = -H(q) + \overline{\log(h)}$$

"generalized max. entropy" $KL(q \parallel h)$
q preferred dist to bias towards

* example from Wainwright

$$\hat{p}_L = \frac{3}{4} \quad \text{Kangaroos are left-handed}$$

$$\hat{p}_R = \frac{2}{3} \quad \text{" drunk labatt bear"}$$

question: how many kangaroos are both L.H. & drunk lab. bear

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

* how do we get set M ?

typically: through empirical moments

$$\begin{aligned} \text{kangaroo: } T_1(x) &= \#\{x \text{ drunk lab}\} \\ T_2(x) &= \#\{x \text{ is L.H.}\} \end{aligned}$$

feature functions: $T_1(x), T_2(x), \dots, T_d(x)$ d features

$$\text{define } M = \{q : \underbrace{Eq[T_j | x]}_{\substack{\text{model expected} \\ \text{feature "count"}}, \text{ "moment constraints}}} = \underbrace{E\hat{p}_n[T_j(x)]}_{\substack{\text{empirical feature "count"} \\ \text{"moment constraints" }}} \}, j=1, \dots, d$$

model expected feature "count" empirical feature "count" "moment constraints"

then

$$\begin{array}{l} \text{MaxENT} \\ \min_{q \in \mathbb{R}^{|X|}} KL(q \parallel \text{unif}) \\ q \in M \\ q \in \Delta_{|X|} \end{array}$$

$$\begin{aligned} \sum_x q(x) T_j(x) &= \frac{1}{n} \sum_{i=1}^n T_j(x^{(i)}) \\ &= \alpha_j \\ \therefore \langle \vec{q}, \vec{T_j} \rangle &= \alpha_j \end{aligned}$$

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

quick presentation of Lagrangian duality

convex min- problem

convex pt. problems $\leftarrow \begin{cases} \cdot f_1, f_2 \text{ are convex} \\ \cdot g_k \text{ affine fct.} \end{cases}$

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

$$\begin{cases} f_i(x) \leq 0 \quad i \in \{1, \dots, m\} \\ g_k(x) = 0 \quad k \in \{1, \dots, n\} \end{cases}$$

} "primal problem"

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

λ

ν

"Lagrange multipliers"

magic trick
(saddle pt. interpretation)

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

$$h(x) = f(x) + S_M(x)$$

Indicator set. $S_M(x) = \begin{cases} +\infty & \text{if } x \in M \\ 0 & \text{else} \end{cases}$

an equivalent problem to primal problem $\inf_x h(x)$

fancy non-smooth fn.

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

duality trick is to swap inf & sup

$$\sup_{\substack{\lambda \geq 0 \\ \nu}} \left(\inf_x g(x, \lambda, \nu) \right) \rightarrow \text{this fn. is always concave} \\ \equiv g(\lambda, \nu) \quad \text{Lagrange dual fn.}$$



Lagrange dual problem

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

"dual variables"

"weak duality"

$$\text{in general, } \sup_{\lambda, \nu} \inf_x g(x, \lambda, \nu) \leq \inf_x \sup_{\lambda, \nu} g(x, \lambda, \nu)$$

Strong duality when $\sup \inf g = \inf \sup g$

↳ sufficient conditions:

- primal problem is convex
- constraint qualification condition (e.g. Slater's condition)

(can get optimal primal variables $x^* = (x^*, \nu^*)$
using KKT conditions)

(see ch.5 of Boyd's book)

13h26

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

dual problem for max. entropy

Max Ent in
primal form
(P)

$$\min_{q \in \mathbb{R}^K} KL(q || w) \\ \sum_x q(x) \log \frac{q(x)}{w(x)}$$

$$u(x) = \frac{1}{K}$$

absorb this constraint
in domain of def.

primal form (P)

eq. const.

$$q \in \mathbb{R}^K$$

$$c \rightarrow \sum_x q(x) \log \frac{q(x)}{u(x)}$$

$$\nabla \rightarrow \sum_x q(x) = 1 \quad \Delta \mathbf{1}_x$$

$$\sum_x q(x) T_j(x) = v_j \quad \forall j$$

absorb this constraint in domain of ref.
 $kL(q||u)$ i.e.
 $kL(q||u) = \begin{cases} +\infty & \text{if } q_k < 0 \\ kL(q||u) & \text{otherwise} \end{cases}$

$$\mathcal{L}(q, v, c) = \sum_x q(x) \log \frac{q(x)}{u(x)} + \sum_{j=1}^d v_j (q_j - \sum_x q(x) T_j(x)) + c (1 - \sum_x q(x))$$

$$q(x) = q_x \quad \frac{\partial \mathcal{L}}{\partial q_x} = 1 + \log \frac{q(x)}{u(x)} + \sum_{j=1}^d v_j T_j(x) - c \stackrel{\text{want}}{=} 0$$

$$\Rightarrow q_{v,c}^*(x) = u(x) \exp(v^T T(x) + c - 1)$$

exponential family?

dual fit.:

plug back $q_{v,c}^*$ in $\mathcal{L}(\dots)$

$$g(v, c) = \mathcal{L}(q_{v,c}^*, v, c) = \tilde{E}_{q^*} [v^T T(x) + c - 1] + v^T \alpha - \tilde{E}_{q^*} [v^T T(x)]$$

~~shorthand~~
 $\sum_x q(x)$

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

~~max g(v, c)
with respect to c~~

$$\nabla_c = 0 \Rightarrow 1 - \exp(c-1) Z(v) \stackrel{\text{want}}{=} 0$$

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

plug back c^*

$$\max_c g(v, c) = v^T \alpha + c^* - \underbrace{Z(v)}_{c^* - 1 = -\log Z(v)}$$

dual problem $\max_v \tilde{g}(v)$

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

link with MLE:

$$\mathcal{L} \propto \frac{1}{n} \sum_{i=1}^n T(x^{(i)}) = \tilde{E}_{p_n} [T(x)]$$

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

$\log p(x^{(i)}|v) + \text{const.}$

where $p(x|v) \triangleq u(v) \exp(v^T T(x) - \log Z(v))$

$$\log p(x|v) + \text{cst.}$$

$$\text{where } p(x|v) \triangleq u(x) \exp(v^T T(x) - \log Z(v))$$

i.e. dual problem is $\max_v \tilde{g}(v) = \max_v \int_n \log p(x_i|v) + \text{cst.}$

i.e. MLE //

to summarize: ML in exp. family with $T(x)$ as sufficient statistics

is equivalent to max ENT with empirical moment constraints on $T(x)$
where $q = \mathbb{E}_{p_n} [T(x)]$

they are Lagrangian dual of each other?

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

$$\text{note: } \nabla_v \log Z(v) = \frac{1}{Z(v)} \nabla_v \sum_x u(x) \exp(v^T T(x))$$

$$= \sum_x T(x) \left(\frac{1}{Z(v)} u(x) \exp(v^T T(x)) \right)$$

KL Pythagorean theorem

→ see lecture 16 in 2017

for any $q \in M$ and $h \in E$

$$KL(q||h) = KL(q||p_n) + KL(p_n||h)$$

$$\nabla_v \log Z(v) = \mathbb{E}_{p(x|v)} [T(x)] \stackrel{\text{def}}{=} \mu(v) \quad \text{"model moment"}$$

$$\nabla_v \tilde{g}(v) = \underbrace{\mathbb{E}_{p_n} [T(x)]}_{\hat{\mu}_n} - \mu(v) \quad \text{"empirical moment"}$$

$$\nabla_v \tilde{g}(v) = 0 \Rightarrow \boxed{\mu(v^*) = \hat{\mu}_n} \quad \text{i.e. moment matching!}$$

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