

today : • max Ent & duality
• exponential family

MLE & KL minimization

$\{P_\theta\}_{\theta \in \Theta}$ parametric family for a discrete observation space

then $ML \text{ for } \theta \Leftrightarrow \min_{\theta \in \Theta} KL(\hat{P}_n \parallel P_\theta)$

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

proof : $KL(\hat{P}_n \parallel P_\theta) = \sum_x \hat{P}_n(x) \log \frac{\hat{P}_n(x)}{P_\theta(x)}$
 $= -H(\hat{P}_n) - \sum_x \hat{P}_n(x) \log P_\theta(x)$
 $= -H(\hat{P}_n) - \frac{1}{n} \sum_{i=1}^n \log P_\theta(x^{(i)})$
 $= \underbrace{-H(\hat{P}_n)}_{\text{constant w.r. to } \theta} - \frac{1}{n} \sum_{i=1}^n \log P_\theta(x^{(i)})$
 $= \text{constant} - \log \prod_{i=1}^n P_\theta(x^{(i)})$

Maximum entropy principle :

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

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

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

i.e. $\hat{p} = \underset{q \in M}{\text{argmax}} H(q)$
 $= \underset{q \in M}{\text{argmin}} KL(q \parallel \underset{u}{\text{uniform}})$

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

"generalized max. entropy" $KL(q \parallel p_0)$
 preferred dist. to be bias towards

* example from Wainwright

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

$\hat{p}_B = \frac{2}{3}$ " drink salted beer

question: how many kang are both L.H & drink salt beer?

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

* how do we get set M?

typically = through empirical "moments"

kangaroos
 $T_L(x) = \mathbb{1}_{\{x \text{ drinks salted beer}\}}$
 $T_B(x) = \mathbb{1}_{\{x \text{ is left handed}\}}$

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

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

then

Max ENT
 $\min_{q \in \mathcal{R}} \text{KL}(q \| \pi)$
 st. $q \in M$
 $q \in \Delta^{|\mathcal{X}|}$

$\mathbb{E}_q[T_j(x)]$ want
 $\sum_x q(x) T_j(x) = \frac{1}{n} \sum_{i=1}^n T_j(x^{(i)}) = \alpha_j$
 $i.e. \langle \vec{q}, \vec{T}_j \rangle = \alpha_j$

↳ convex opt. problem over $q \in \Delta^{|\mathcal{X}|} \in \mathbb{R}^{|\mathcal{X}|}$

quick presentation of Lagrangian duality

convex problem \Leftrightarrow convex min. problem

$\cdot f, f_j$ are convex fct.

$\cdot g_k$ affine fct.

$\min_{x \in \mathcal{X}} f(x)$

st. $f_j(x) \leq 0 \quad \forall j \in \{1, \dots, m\}$

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

}

"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 problem to the primal problem

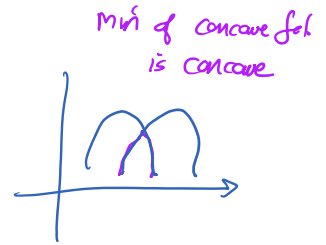
$\inf_x \left(\sup_{\lambda \geq 0, \nu} f(x, \lambda, \nu) \right)$

duality trick is to swap inf & sup

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

$\triangleq g(\lambda, \nu)$ Lagrange dual fct.

this fct. is always concave



Lagrange dual problem

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

"dual variables"

in general, $\sup \inf f(x, \lambda, \nu) \leq \inf \sup f(x, \lambda, \nu)$

weak duality

strong duality when $\sup \inf f = \inf \sup f$

sufficient conditions:

- when primal problem is convex
- + constraint qualification condition (eg. Slater's condition)

(can get optimal primal variables x^* (λ^*, ν^*) 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/>

15h40

dual problem for max. entropy

Max ENT in primal form (P)

$\min_{q \in \mathbb{R}^K} KL(q||u)$

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

$q(x) \geq 0 \forall x$

$\sum_x q(x) = 1$

$\sum_x q(x) T_j(x) = \alpha_j \forall j$

$\Delta_{|X|}$

M

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

absorbs this constraint in domain of def. of $KL(q||u)$ i.e.

$KL(q||u) = \begin{cases} +\infty & \text{if } q(x) < 0 \text{ for some } x \\ KL(q||u) & \text{o.w.} \end{cases}$

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

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

want $\sum_j \nu_j T_j(x)$

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

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

exponential family!

dual set.:

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

$$\begin{aligned} g(\nu, c) &= \mathcal{J}(q_{\nu, c}^*, \nu, c) \\ &= \mathbb{E}_{q^*}[\nu^T T(x) + c - 1] + \nu^T \alpha - \mathbb{E}_{q^*}[\nu^T T(x)] \\ &\quad + c - \mathbb{E}_{q^*}[1] \\ &= \nu^T \alpha + c - \underbrace{\sum_x u(x) \exp(\nu^T T(x))}_{\hat{=} Z(\nu)} \exp(c-1) \end{aligned}$$

$$\begin{aligned} \max_{\text{with respect to } c} g(\nu, c) \quad \nabla_c = 0 &\Rightarrow 1 - Z(\nu) \exp(c-1) \stackrel{\text{want}}{=} 0 \\ &\Rightarrow \exp(c^*-1) = \frac{1}{Z(\nu)} \end{aligned}$$

$$\begin{aligned} \text{plug back } c^* : \max_c g(\nu, c) &= \nu^T \alpha + c^* - \frac{1}{Z(\nu)} \\ & \quad c^*-1 = -\log Z(\nu) \end{aligned}$$

dual problem

$$\max_{\nu} \tilde{g}(\nu) \quad \boxed{\tilde{g}(\nu) \hat{=} \nu^T \alpha - \log Z(\nu)}$$

link with MLE:

$$\text{if } \alpha = \frac{1}{n} \sum_{i=1}^n T(x^{(i)}) = \mathbb{E}_{p_n}[T(x)]$$

$$\text{then } g(\nu) = \frac{1}{n} \sum_{i=1}^n [\underbrace{\nu^T T(x^{(i)})}_{\log p(x^{(i)}|\nu)} - \log Z(\nu)] + \text{const.}$$

$$\text{where } p(x|\nu) \hat{=} u(x) \exp(\nu^T T(x) - \log Z(\nu))$$

$$\text{i.e. dual problem is } \max_{\nu} g(\nu) = \max_{\nu} \frac{1}{n} \log p(x_{1:n}|\nu) \quad \text{i.e. MLE}$$

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

is equivalent to Max-ENT with moment constraints on $T(x)$

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

to summarize: ML in exp. family with $T(x)$ as sufficient statistics
 is equivalent to Max-ELT with moment constraints on $T(x)$
 where $\alpha = \mathbb{E}_{p_n}[T(x)]$

they are Lagrangian dual of each other!

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

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

$$= \sum_x T(x) \frac{u(x) \exp(\nu^T T(x))}{Z(\nu)}$$

$$= \sum_x T(x) p(x|\nu)$$

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

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

$$\nabla_{\nu} \tilde{g}(\nu) = 0 \Rightarrow \mu(\nu^*) = \hat{\mu}_n$$
 i.e. moment matching!

↑ "empirical moment"

→ see lecture 16 in 2017 for "KL Pythagorean theorem"

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