

today: • Max Ent & duality
• exponential family

MLE & KL minimization

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

then MLE for Θ for iid. data $\Leftrightarrow \min_{\theta \in \Theta} KL(\hat{p}_n \parallel p_\theta)$

$\hat{p}_n(x) \stackrel{\Delta}{=} \frac{1}{n} \sum_{i=1}^n \delta(x, x^{(i)})$
 "empirical distribution" Kronecker-delta

proof:

$$\begin{aligned}
 KL(\hat{p}_n \parallel p_\theta) &= \sum_x \hat{p}_n(x) \log \frac{\hat{p}_n(x)}{p_\theta(x)} \\
 &= -H(\hat{p}_n) - \sum_x \hat{p}_n(x) \log p_\theta(x) \\
 &\quad \frac{1}{n} \sum_{i=1}^n \delta(x, x^{(i)}) \\
 &= -H(\hat{p}_n) - \frac{1}{n} \sum_{i=1}^n \sum_x \delta(x, x^{(i)}) \log p_\theta(x) \\
 &\quad \log p_\theta(x^{(i)}) \\
 &= \underbrace{-H(\hat{p}_n)}_{\text{constant w.r. to } \Theta} - \frac{1}{n} \sum_{i=1}^n \log p_\theta(x^{(i)}) \\
 &\quad \log \left(\prod_{i=1}^n p_\theta(x^{(i)}) \right)
 \end{aligned}$$

Maximum entropy principle:

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

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

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

i.e. $\hat{p} = \underset{q \in M}{\text{argmax}} H(q)$

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

$KL(q \parallel u) = \sum_x q(x) \log \frac{q(x)}{1/|X|} = -H(q) + \text{constant}$

$$KL(q \| u) = \sum_x q(x) \log \frac{q(x)}{u(x)} = -H(q) + \text{constant}$$

"generalized max. entropy" $KL(q \| h_\theta)$
 preferred det. to bias towards

* example from Wainwright

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

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

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

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

* how do we get set M

typically: through empirical "moments"

kangaroo:
 $T_1(x) = \mathbb{1}_{\{x \text{ drinks Sabatt}\}}$
 $T_2(x) = \mathbb{1}_{\{x \text{ is left handed}\}}$

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

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

then

Max ENT

min $q \in \mathbb{R}^{|\mathcal{X}|}$ $KL(q \| \text{unif})$

$q \in M$

$q \in \Delta_{|\mathcal{X}|}$

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

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

quick presentation of Lagrangian duality

convex min. problem

convex opt. problem \leftarrow

- f, f_j are convex fct.
- g_k affine fct.

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

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

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

}

"primal problem"

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

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

"Lagrange multipliers"

magic trick
(saddle pt. interpretation)

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

$$h(x) = f(x) + \delta_{\mathcal{M}}(x) \text{ fct.}$$

$$\delta_{\mathcal{M}}(x) = \begin{cases} +\infty & \text{if } x \notin \mathcal{M} \\ 0 & \text{o.w.} \end{cases}$$

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

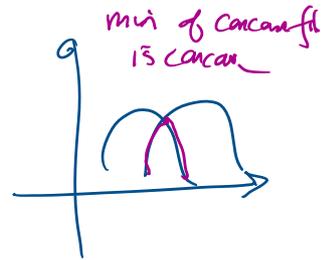
$$\inf_x \left(\sup_{\lambda \geq 0, \nu} L(x, \lambda, \nu) \right)$$

↳ fancy non-smooth fct.

duality trick is to swap \inf & \sup

$$\sup_{\lambda \geq 0, \nu} \left(\inf_x L(x, \lambda, \nu) \right) \rightarrow \text{this fct. is always concave}$$

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



Lagrangian dual problem

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

"dual variables"

"weak duality"

in general, $\sup_{\lambda, \nu} \left(\inf_x L(x, \lambda, \nu) \right) \leq \inf_x \left(\sup_{\lambda, \nu} L(x, \lambda, \nu) \right)$

$g(\lambda, \nu)$

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

- ↳ sufficient conditions:
- when primal problem is convex
 - + constraint qualification condition (e.g. Slater's condition)

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

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/>

15h33

| | | |
|-------------------------------|---------------------------------------|------------------------|
| dual problem for max. entropy | $\min_{\lambda, \nu} \text{kl}(q w)$ | $\inf_x \frac{1}{ x }$ |
| MaxEUT in | min | |

convex program for max. entropy

Max EUT in primal form (P)

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

$$q(x) \geq 0 \quad \forall x \quad \Delta_{|X|}$$

$$c \rightarrow \sum_x q(x) = 1$$

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

absorb this constraint in domain of def. KL(q||u) i.e. KL(q||u) = { +∞ if q(x) < 0 for some x, KL(q||u) o.w. }

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

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

$$\Rightarrow \frac{q(x)}{u(x)} = \exp(v^T T(x) + c - 1)$$

exponential family!

dual fct: plug back $q_{v,c}^*$ in $J(\dots)$

$$g(v, c) = J(q_{v,c}^*, v, c)$$

$$= \mathbb{E}_{q^*} [v^T T(x)] + d(-1) + v^T \alpha - \mathbb{E}_{q^*} [v^T T(x)] + c - \mathbb{E}_{q^*} [d]$$

shorthand for $\sum_x q(x)$

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

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

$$c^*-1 = -\log Z(v)$$

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

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

link with MLE!

$$i.e. \quad v = \left(\sum_{i=1}^n T(i) \right) - \mathbb{E}_{\pi} [T(i)]$$

link with MLE:

$$\text{if } \alpha = \int \sum_{i=1}^n T(x^{(i)}) = \mathbb{E}_{\hat{p}_n} [T(x)]$$

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

log p(x⁽ⁱ⁾|\nu) + cst.

where $p(x|\nu) \triangleq u(x) \exp(\nu^T T(x) - \log z(\nu))$

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

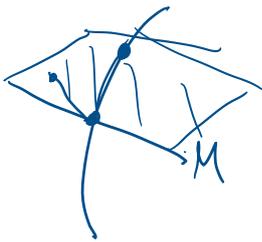
i.e. **MLE** //

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 $\alpha = \mathbb{E}_{\hat{p}_n}(T(x))$

they are Lagrangian dual of each other?

MLE in exp family \Leftrightarrow moment matching in exp family



KL Pythagorean theorem

→ see lecture 16 in 2017

$$\text{note: } \nabla_{\nu} \log z(\nu) = \int \frac{\nabla_{\nu} \sum_x u(x) \exp(\nu^T T(x))}{z(\nu)}$$

$$= \sum_x T(x) \frac{u(x) \exp(\nu^T T(x))}{z(\nu)} = \sum_x T(x) p(x|\nu)$$

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

$$\nabla_{\nu} \tilde{g}(\nu) = \underbrace{\mathbb{E}_{\hat{p}_n} [T(x)]}_{\hat{\mu}_n \text{ "empirical moment"}} - \mu(\nu)$$

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

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