Lecture 7. Monday. 4/11/2022

Regularized (discrete) OT
$\begin{cases} - \text{Entropic regularization} \\ - \text{Kullback-Leibler divergence} \\ - \text{Generalization.} \end{cases}$

Given $a \in \mathcal{P}_m$ and $b \in \mathcal{P}_n$, $C \in \mathbb{R}^{m \times n}$ with $C \geq 0$.

Def. $\mathcal{A}(a,b) = \left\{ P = [P_{ij}] \in \mathbb{R}^{m \times n} : P \geq 0, \sum_{j=1}^{m} P_{ij} = a_i, \forall i; \sum_{i=1}^{m} P_{ij} = b_j, \forall j \right\}$

OT: $W_C(a,b) := \min_{P \in \mathcal{A}(a,b)} \langle P, C \rangle.$ $\qquad E[P] \overset{\Delta}{=} \langle P, C \rangle.$

<u>Regularization</u>: Consider $E_\varepsilon[P] = \langle P, C \rangle + \varepsilon h(P)$
for some $h$. Hope that $E_\varepsilon[P]$ is easier to minimize
and as $\varepsilon \to 0$ minimizer/minimum value of $E_\varepsilon$
converge to those of $E (= E_0)$.

Call $h$ a <u>regularizer</u>. How to choose $h$?
$\odot$ (Component wise) convex $\Rightarrow$ uniqueness for each $\varepsilon > 0$.
$\odot$ (component wise) minimum of $h$ is inside $(0,1)$ so that the constraints will be satisfied.

<u>Entropic regularization</u>: For $\varepsilon > 0$, define

$E_\varepsilon[P] = \langle C, P \rangle + \varepsilon \langle P(\log P - 1), \mathbb{1}_{m \times n} \rangle$

$\qquad = \sum_{i=1}^{m} \sum_{j=1}^{n} P_{ij} C_{ij} + \varepsilon \underbrace{\sum_{i=1}^{m} \sum_{j=1}^{n} P_{ij}(\log P_{ij} - 1)}_{h(P)}.$

Notation:
$\odot$ $P(\log P - 1) \in \mathbb{R}^{m \times n}$: $[P(\log P - 1)]_{ij} = P_{ij}(\log P_{ij} - 1) \ \forall ij$
$\odot$ $\mathbb{1}_{m \times n} \in \mathbb{R}^{m \times n}$: all entries $= 1$. (Not $I_{m \times n}$)

<u>Remarks</u> $\odot$ If $f : \mathbb{R}^d \to \mathbb{R}$ satisfies $f \geq 0$ in $\mathbb{R}^d$
and $\int_{\mathbb{R}^d} f \, dx = 1$ then $H(f) = -\int_{\mathbb{R}^d} f \log f \, dx$ is
called the entropy of $f$.

If $p = (p_i) \in \mathcal{P}_d$ : all $p_i \geq 0$, $\sum_{i=1}^{d} p_i = 1$, then

$S = -\sum_{i=1}^{d} p_i \log p_i$ is called the entropy of $P$.

In thermodynics, the free energy is $F = U - TS$ with $U, T, S$ being internal energy, temperature, and entropy. $S = -\frac{\partial F}{\partial T}$.

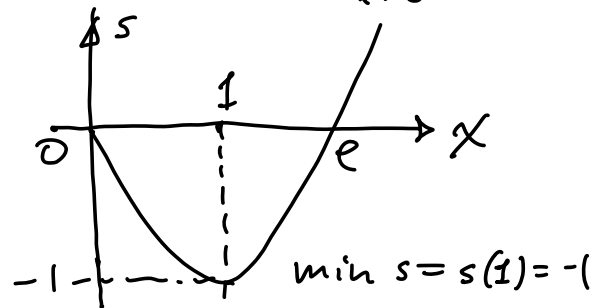In statistical mechanics, entropy $S = -k_B \log \mathcal{r}$, where $\mathcal{r}$ is the partition function.

② Define $s(x) = x \log x - x$ $(x > 0)$. $s(0) = 0 = \lim_{x \to 0^+} s(x)$.

So, $s$ is cont. on $[0, \infty)$. $s(e) = 0$

$s'(x) = \log x = 0 \implies x = 1$

$0 < x << 1 \implies s'(x) << -1$

$s''(x) = \frac{1}{x} > 0$. So, $s$ is strictly convex on $[0, 1]$.

min $s = s(1) = -1$

Fix $i, j$. denote $c = C_{ij}$, and consider

$g(x) = cx + \varepsilon x (\log x - 1)$ $(0 \leq x \leq 1)$. $g(x) = 0, x > 0 \implies x = e^{1 - c/\varepsilon}$

$g'(x) = c + \varepsilon \log x$. $g'(x) = 0 \implies \boxed{x_c = e^{-c/\varepsilon} \in (0,1)}$

$g(x_c) = c e^{-c/\varepsilon} + \varepsilon (-c/\varepsilon) = c(e^{-c/\varepsilon} - 1)$.

$0 < \varepsilon << 1, c > 0 \implies x_c \approx 0$.

$g''(x) = \frac{\varepsilon}{x}$. $g$ is convex.

Recall for $\varepsilon = 0$, the original OT has minimizers but may not be unique. Denote

$\mathcal{M}_c(a, b) = \{\hat{p} \in \mathcal{A}(a,b) : <C, \hat{p}> = W_c(a,b)\}$.

This is the set of all optimal plans. We showed before that $\mathcal{M}_c(a, b) \neq \phi$ is convex and compact.

<u>Proposition</u> There exists a unique $P_0 \in \mathcal{M}_C(a,b)$ such that
$$h(P_0) = \min_{P \in \mathcal{M}_C(a,b)} h(P).$$

<u>Proof</u> Since $h(P) = \langle P(\log P - 1), \mathbf{1}_{m \times n} \rangle$ is continuous and $\mathcal{M}_C(a,b)$ is compact, there exists $\arg\min_{\mathcal{M}_C(a,b)} h(\cdot)$. Suppose $P_0$ and $Q_0 \in \mathcal{M}_C(a,b)$ satisfy $h(P_0) = h(Q_0) \leq h(P)$ $\forall P \in \mathcal{M}_C(a,b)$. Let $R_\lambda = (1-\lambda)P_0 + \lambda Q_0$ $(0 \leq \lambda \leq 1)$. Then $R_\lambda \in \mathcal{M}_C(a,b)$ since $\mathcal{M}_C(a,b)$ is convex. Moreover, the convexity of $h$ implies that $h(R_\lambda) \leq (1-\lambda)h(P_0) + \lambda h(Q_0) \leq \min_{\mathcal{M}_C(a,b)} h(\cdot)$. Hence, $R_\lambda = \arg\min_{\mathcal{M}_C(a,b)} h(\cdot)$ and $f(\lambda) := h(R_\lambda) = $ const.

Suppose $P_0 \neq Q_0$. We examine $f'(\lambda)$, $f''(\lambda)$ for $\lambda$ close to 0. If $P_{0,ij} = Q_{0,ij} = 0$ then $s(R_{\lambda,ij}) = 0$. $(s(x) = x\log x - x)$. If $Q_{0,ij} > 0$ then $\frac{d}{d\lambda}\big|_{\lambda=0} s(R_{\lambda,ij})$ exists and is finite. For any $(i,j)$, if $Q_{0,ij} = 0$ and $P_{0,ij} > 0$, then $\frac{d}{d\lambda} s(R_\lambda) = \frac{d}{d\lambda} s(\lambda P_{0,ij}) \to -\infty$ as $\lambda \to 0$. But $f' \equiv 0$ as $f = $ const. Hence $P_{0,ij} = Q_{0,ij} = 0$ or $P_{0,ij} > 0$ and $Q_{0,ij} > 0$ for any $i,j$. Thus, $h(R_\lambda) = \sum_{ij} s(R_{\lambda,ij})$ is a $C^2((0,1))$-function of $\lambda$. Now, direct calculations lead to
$$f''(0) = \frac{d^2}{d\lambda^2}\Big|_{\lambda=0} s(R_\lambda) = {\sum_{ij}}' \frac{1}{Q_{0,ij}}(P_{0,ij} - Q_{0,ij})^2 > 0,$$
where ${\sum_{ij}}'$ sums over all $(i,j)$ with $P_{0,ij} > 0$ and $Q_{0,ij} > 0$. But $f(\lambda) = $ const. on $[0,1]$. So $f''(0) = 0$, a contradiction. <u>QED</u>

<u>Theorem</u> (1) For any $\varepsilon > 0$, there exists a unique $P_\varepsilon \in \mathcal{A}(a,b)$ such that $E_\varepsilon[P_\varepsilon] = \min_{P \in \mathcal{A}(a,b)} E_\varepsilon[P]$.

(2) $\lim_{\varepsilon \to 0^+} P_\varepsilon = P_0 = \arg\min_{P \in \mathcal{M}_C^0(a,b)} h(P)$.

(3) $\lim_{\varepsilon \to 0^+} E_\varepsilon[P_\varepsilon] = \lim_{\varepsilon \to 0^+} E[P_\varepsilon] = E[P_0] = W_C(a,b)$.

**Proof** (1) $E_\varepsilon : \mathcal{A}(a,b) \to \mathbb{R}$ is continuous and $\mathcal{A}(a,b)$ is compact. Hence, $\exists \, P_\varepsilon = \arg\min\limits_{\mathcal{A}(a,b)} E_\varepsilon$. The uniqueness can be shown by the same argument in the proof of Proposition above.

(2) It suffices to show that any convergent sequence $\{P_{\varepsilon_k}\}$ ($\varepsilon_k \downarrow 0$) has the same limit $P_0$. Assume $P_{\varepsilon_k} \to \overline{P}$. Clearly, $\overline{P} \in \mathcal{A}(a,b)$. Now, $E_{\varepsilon_k}[P_{\varepsilon_k}] \le E_{\varepsilon_k}[P_0]$, i.e.,

$$E[P_{\varepsilon_k}] + \varepsilon_k \, h(P_{\varepsilon_k}) \le E[P_0] + \varepsilon_k \, h(P_0).$$

Since $h$ is bounded and continuous, we have for $k \to \infty$ that $E[\overline{P}] \le E[P_0] = W_c(a,b)$, i.e., $\overline{P} \in \mathcal{M}_c(a,b)$. Now,

$$E_{\varepsilon_k}[P_{\varepsilon_k}] = E[P_{\varepsilon_k}] + \varepsilon_k \, h(P_{\varepsilon_k}) \le E_{\varepsilon_k}[P_0] = E[P_0] + \varepsilon_k \, h(P_0)$$
$$\le E[P_{\varepsilon_k}] + \varepsilon_k \, h(h_0).$$

Hence, $h(P_{\varepsilon_k}) \le h(h_0)$. Taking $k \to \infty$, we get $h(\overline{P}) \le h(P_0)$. By the above Proposition, $\overline{P} = P_0$.

(3) Since $P_\varepsilon \to P_0$, $E_\varepsilon[P_\varepsilon] = E[P_\varepsilon] + \varepsilon \, h(P_\varepsilon) \to E[P_0]$ $= W_c(a,b)$. $\underline{QED}$

**Theorem** (Cominetti-San Martin 1994) The convergence $P_\varepsilon \to P_0$ in the above Thm is exponential, i.e., $P_\varepsilon = P_0 + G(\varepsilon)$ for $0 < \varepsilon \le \varepsilon_0$ for some $\varepsilon_0$ and $G(\varepsilon)$ satisfies $\lim\limits_{\varepsilon \to 0+} \dfrac{G_{ij}(\varepsilon)}{e^{-\mu/\varepsilon}} = 0$ $\quad \forall i, j$.

for some $\mu > 0$. $\underline{QED}$

A different prob, $\min\limits_{P \in \mathcal{A}(a,b)} \langle P(\log P - 1), \mathbb{1}_{m \times n} \rangle$.

Answer: the unique minimizer is $a \otimes b$, and the minimum is

$$\sum_{i,j} a_i b_j [\log(a_i b_j) - 1] = \sum_{i,j} \left[ a_i b_j (\log a_i + \log b_j - 1) \right]$$

$$= \sum_i a_i \log a_i + \sum_j b_j \log b_j - 1.$$

"$\underline{Proof}$" Denote $h(P) = \langle P(\log P - 1), \mathbb{1}_{m \times n} \rangle = \sum_{i,j} s(P_{ij})$, $P \in \mathscr{A}(a, b)$

Since $h$ is continuous and $\mathscr{A}(a, b)$ is compact, there exists a minimizer, which is a critical point of the Lagrangian

$$\mathscr{L}(P, \lambda) = h(P) + \sum_{i=1}^m f_i \left( \sum_{j=1}^n P_{ij} - a_i \right) + \sum_{j=1}^n g_j \left( \sum_{i=1}^m P_{ij} - b_j \right), \quad \lambda = (f, g).$$

$$\partial_{f_i} \mathscr{L} = 0 \implies \sum_j P_{ij} = a_i \;\; \forall i, \;\; \sum_i P_{ij} = b_j \;\; \forall j.$$

$$\partial_{P_{k\ell}} \mathscr{L} = 0 \implies \log P_{k\ell} - f_k - g_\ell = 0 \;\; \text{or} \;\; P_{k\ell} = 0 \;\; \forall k, \ell.$$

So, if $P_{k\ell} > 0$ then $P_{k\ell} = e^{f_k + g_\ell} = \alpha_k \beta_\ell, \; \alpha_k = e^{f_k}, \; \beta_\ell = e^{g_\ell}$.

If $P_{k\ell} = 0$ then $P_{k\ell} = \alpha_k \beta_\ell$ with $\alpha_k = 0$ or $\beta_\ell = 0$. So, $P_{k\ell} = \alpha_k \beta_\ell$ with $\alpha_k \geq 0, \; \beta_\ell \geq 0$. Let $\alpha = \sum_k \alpha_k, \; \beta = \sum_\ell \beta_\ell$. Since $P \in \mathscr{A}(a, b)$

$$a_k = \sum_\ell P_{k\ell} = \sum_\ell \alpha_k \beta_\ell = \alpha_k \beta. \qquad b_\ell = \alpha \beta_\ell. \qquad \sum_k a_k = 1 \implies \alpha \beta = 1.$$

$$\sum_j b_j = 1 \implies \alpha \beta = 1. \quad So. \quad P_{k\ell} = \alpha_k \beta_\ell = \alpha_k \beta \cdot \beta_\ell \alpha = a_k b_\ell. \quad \underline{QED}$$

$\underline{Remark}$ The above proof is not rigorous as the inequality constraints $P_{ij} \geq 0 \;\; \forall i, j$. are not included in the Lagrangian. A similar but correct proof should use the KKT conditions for minimizing a convex function with equalities and inequalities constraints.

$\quad \lambda = (f, g)$ for ~~the~~ equality constraints

$\quad \mu = (\mu_{ij})$ for ~~the~~ inequality constraints.

Then, necessary conditions $\implies P_{k\ell} = e^{f_k + g_\ell - \mu_{k\ell}}$

$\quad \sum_\ell P_{k\ell} = a_k, \; \sum_k P_{k\ell} = b_\ell \;\; \forall k, \ell. \;\; \mu_{k\ell} \geq 0.$ ALSO,

the complementarity condition $\implies \mu_{ij} P_{ij} = 0 \;\; \forall i, j$

Hence, all $\mu_{ij} = 0 \implies P_{k\ell} = e^{f_k + g_\ell}$. As above.

$p_{k\ell} = a_k b_\ell$. Finally, check that the sufficient conditions. <u>QED</u>

<u>Corollary</u> If $P = [p_{ij}] \in \mathcal{A}(a,b)$ then

$$\sum_{i,j} p_{ij} \log p_{ij} \geq \sum_i a_i \log a_i + \sum_j b_j \log b_j.$$

<u>Pf.</u> $\sum_{i,j} p_{ij}(\log p_{ij} - 1) \geq \sum_{i,j} a_i b_j [\log(a_i b_j) - 1]$

$\sum_{i,j} p_{ij} = 1 \Rightarrow \sum_{i,j} p_{ij} \log p_{ij} \geq \sum a_i b_j \log a_i + \sum a_i b_j \log b_j$
$\qquad\qquad\qquad\qquad\qquad = \sum_i a_i \log a_i + \sum_j b_j \log b_j.$ <u>QED</u>

<u>Lagrange multiplier and Kullback-Leibler
divergence (or: relative entropy).</u>

$\min_{\mathcal{A}(a,b)} E_\varepsilon$. Define the lagrange multiplier
$\lambda = (f, g) \in \mathbb{R}^m \times \mathbb{R}^n$ and the lagrangian

$$\mathcal{L}_\varepsilon(P, \lambda) = \sum_{i,j} p_{ij} c_{ij} + \varepsilon \sum_{i,j} p_{ij}(\log p_{ij} - 1)$$
$$- \sum_{i=1}^m f_i \left( \sum_{j=1}^n p_{ij} - a_i \right) - \sum_{j=1}^n g_j \left( \sum_{i=1}^m p_{ij} - b_j \right).$$

Minimizers of $E_\varepsilon$ are extreme points of $\mathcal{L}_\varepsilon$.

$\partial_{f_k} \mathcal{L}_\varepsilon = 0 \Rightarrow \sum_{j=1}^n p_{kj} = a_k \quad \forall k.$

$\partial_{g_\ell} \mathcal{L}_\varepsilon = 0 \Rightarrow \sum_{i=1}^m p_{i\ell} = b_\ell \quad \forall \ell.$

$\partial_{p_{k\ell}} \mathcal{L}_\varepsilon = 0 \Rightarrow c_{k\ell} + \varepsilon \log p_{k\ell} - f_k - g_\ell = 0$

$\Rightarrow p_{k\ell} = e^{-\frac{1}{\varepsilon}(c_{k\ell} - f_k - g_\ell)} \quad \forall k,\ell.$

Denote $k_{\varepsilon, ij} = e^{-c_{ij}/\varepsilon} \quad \forall i,j.$ $\quad K_\varepsilon = [k_{\varepsilon, ij}] \in \mathbb{R}^{m \times n}$

Define $KL(P | K_\varepsilon) = \sum_{i,j} \left[ p_{ij} \left( \log \frac{p_{ij}}{k_{\varepsilon, ij}} - 1 \right) + k_{\varepsilon, ij} \right].$

Call it the _Kullback-Leibler divergence_ of $p$ relative to $K_\varepsilon$, or the _relative entropy_ of $p$ relative to $K_\varepsilon$.

What is this quantity?

$$KL(p|K_\varepsilon) = \sum_{i,j} P_{ij}(\log P_{ij} - 1) - \sum_{i,j} P_{ij} \log K_{\varepsilon,ij} + \sum_{i,j} K_{\varepsilon,ij}$$

$$= \sum_{i,j} P_{ij}(\log P_{ij} - 1) + \frac{1}{\varepsilon}\sum_{i,j} P_{ij} C_{ij} + \sum_{i,j} K_{\varepsilon,ij}$$

$$= \frac{1}{\varepsilon}\left[\langle P, C\rangle + \varepsilon \langle P(\log P - 1), \mathbb{1}_{m\times n}\rangle\right] + \langle K_\varepsilon, \mathbb{1}_{m\times n}\rangle$$

$$= \frac{1}{\varepsilon} F_\varepsilon[P] + \langle K_\varepsilon, \mathbb{1}_{m\times n}\rangle.$$

<u>Proposition</u>  $F_\varepsilon[P] = \varepsilon KL(P|K_\varepsilon) - \varepsilon\langle K_\varepsilon, \mathbb{1}_{m\times n}\rangle. \quad \forall p \in \mathscr{A}(a,b)$

Hence, the unique minimizer of $F_\varepsilon$ over $\mathscr{A}(a,b)$ is the unique minimizer of $KL(\cdot | K_\varepsilon)$ over $\mathscr{A}(a,b)$. <u>QED</u>

<u>General regularized OT problem</u>

Let $h \in C([0,1]) \wedge C^2((0,1))$ with $h'' > 0$ on $(0,1)$.

consider $E_\varepsilon[P] = \langle C, P\rangle + \varepsilon\langle h(P), \mathbb{1}_{m\times n}\rangle$

$$= \sum_{i,j} C_{ij} P_{ij} + \varepsilon \sum_{i,j} h(P_{ij}), \quad P \in \mathscr{A}(a,b)$$

$$E[P] = \langle C, P\rangle, \quad P \in \mathscr{A}(a,b)$$

Recall $a \in \mathcal{P}_m$, $b \in \mathcal{P}_n$, $C \in \mathbb{R}^{m\times n}$. $C \geq 0$.

$$\mathscr{A}(a,b) = \{P \in \mathbb{R}^{m\times n} : P \geq 0, \sum_i P_{ij} = b_j \,\forall j, \sum_j P_{ij} = a_i \,\forall i\}$$

$$W_C(a,b) = \min_{\mathscr{A}(a,b)} E[\cdot].$$

$$\mathcal{M}_C(a,b) = \{\hat{P} \in \mathscr{A}(a,b) : \langle C, \hat{P}\rangle = W_C(a,b)\}.$$

<u>Thm</u>. (1) $\exists! P_h \in \mathcal{M}_C(a,b)$ such that $h(P_h) \leq h(\hat{P})$ for any $\hat{P} \in \mathcal{M}_C(a,b)$.

(2) For each $\varepsilon > 0$, there exists a unique
$P_\varepsilon = \arg\min\limits_{P \in \mathcal{A}(a,b)} E_\varepsilon[P]$.

(3) $\lim\limits_{\varepsilon \to 0^+} P_\varepsilon = P_h$ and $\lim\limits_{\varepsilon \to 0^+} E_\varepsilon[P_\varepsilon] = W_C(a,b)$. QED

Examples (1) Entropic regularization.

$h(P) = \sum\limits_{i,j} P_{ij}(\log P_{ij} - 1)$.

(2) Quadratic regularization. $h(P) = \sum\limits_{i,j} \frac{1}{2}P_{ij}^2$.

(3) Binary entropic regularization

$h(P) = \sum\limits_{i,j} [P_{ij} \log P_{ij} + (1 - P_{ij})\log(1 - P_{ij})]$.