

Lecture 8, Wednesday, 4/13/2022

• Sinkhorn algorithm

Let $a \in \mathbb{R}_+^m$, $b \in \mathbb{R}_+^n$, $C \in \mathbb{R}^{m \times n}$, $C \geq 0$.

$\mathcal{A}(a, b) = \{P \in \mathbb{R}^{m \times n}: P \geq 0, \text{ row sum of } P = a, \text{ col. sum of } P = b\}$.

K-discrete OT: $\min_{P \in \mathcal{A}(a, b)} E[P]$, $E[P] = \langle P, C \rangle$.

Entropy regularized OT: $\min_{P \in \mathcal{A}(a, b)} E_\varepsilon[P]$.

$$E_\varepsilon[P] = \langle P, C \rangle + \varepsilon \langle P \log P - I \rangle = \sum_{i,j} p_{ij} c_{ij} + \varepsilon \sum_{i,j} p_{ij} (\log p_{ij} - 1).$$

$\exists! P_\varepsilon = \arg \min_{P \in \mathcal{A}(a, b)} E_\varepsilon$. $P_\varepsilon \rightarrow P_0 \in \mathcal{A}(a, b)$ with max. entropy.

Method of Lagrange multipliers

$$\mathcal{L}(P, \lambda) = \mathcal{L}(P, (f, g)) = E_\varepsilon[P] - \sum_{i=1}^m f_i \left(\sum_{j=1}^n p_{ij} - a_i \right) - \sum_{j=1}^n g_j \left(\sum_{i=1}^m p_{ij} - b_j \right).$$

$$\frac{\partial \mathcal{L}}{\partial f_k} = 0 \Rightarrow \sum_j p_{kj} = a_k \quad \forall k, \quad \frac{\partial \mathcal{L}}{\partial g_l} = 0 \Rightarrow \sum_i p_{il} = b_l \quad \forall l.$$

$$\frac{\partial \mathcal{L}}{\partial p_{kl}} = 0 \quad c_{kl} - \varepsilon \log p_{kl} - f_k - g_l = 0 \quad \forall k, l.$$

$$p_{kl} = e^{-c_{kl} - f_k - g_l}/\varepsilon = e^{-c_{kl}/\varepsilon} e^{-f_k/\varepsilon} e^{-g_l/\varepsilon}$$

$$\text{Set } K_{kl} = K_{\varepsilon, kl} = e^{-c_{kl}/\varepsilon}, \quad p_{kl} = K_{\varepsilon, kl} u_k v_l \geq 0$$

$$\left. \begin{array}{l} u_k = u_{\varepsilon, k} = e^{-f_k/\varepsilon}, \\ v_l = v_{\varepsilon, l} = e^{-g_l/\varepsilon}. \end{array} \right\} P = \text{diag}(u) K \text{diag}(v)$$

$$\begin{bmatrix} u_1 & \dots & 0 \\ 0 & \dots & u_m \end{bmatrix}$$

The problem to solve now becomes the following:

Given $K \in \mathbb{R}^{m \times n}$, $K \geq 0$. Find diagonal matrices $\text{diag}(u)$ and $\text{diag}(v)$ with $u \in \mathbb{R}^m$, $u > 0$ and $v \in \mathbb{R}^n$, $v > 0$, such that $P = \text{diag}(u) K \text{diag}(v) \in \mathcal{A}(a, b)$, i.e., row sum of $P = a$, col. sum of $P = b$.

(Note: We assume $a > 0$ and $b > 0$.)

$$\sum_{\ell=1}^n p_{k\ell} = a_k \Rightarrow \sum_{\ell=1}^n K_{k\ell} u_k v_\ell = u_k (Kv)_k = a_k$$

$$\Rightarrow u \odot (Kv) = a, \text{ similarly } v \odot (K^T u) = b.$$

Remark Strictly speaking, we should consider the inequality constraints $p_{ij} \geq 0 \forall i,j$ in the Lagrangian function, adding $\sum_{i,j} \mu_{ij} p_{ij}$. And use the KKT conditions. But, the complementary slackness: $u_{ij} p_{ij} = 0$ imply that $u_{ij} = 0$ since $p_{ij} > 0$ by the conditions $\frac{\partial L}{\partial p_{ij}} = 0 \forall i,j$.

Theorem Assume $a > 0$ and $b > 0$. Then the unique minimizer $P_\varepsilon = \arg \min_{\mathcal{A}(a,b)} F_\varepsilon[\cdot]$ is given by $P_\varepsilon = \text{diag}(u) K \text{diag}(v)$ for some $u \in \mathbb{R}_+^m$ and $v \in \mathbb{R}_+^n$ such that $P_\varepsilon \in \mathcal{A}(a,b)$. Moreover, (u,v) is unique up to a scaling factor.

Proof The characterization of P_ε follows from the calculations above using the method of Lagrange multipliers. See a theorem below for the uniqueness. QED

Now, we know $K = e^{-C/\varepsilon} = K_\varepsilon$. We need to find $u = u_\varepsilon \in \mathbb{R}_+^m$ and $v = v_\varepsilon \in \mathbb{R}_+^n$ such that $P = \text{diag}(u) K \text{diag}(v) \in \mathcal{A}(a,b)$. This means $u_i = \frac{a_i}{Kv_i}$ and $v_j = \frac{b_j}{K^T u_j}$ (i.e., $u_i = a_i / (Kv)_i$, and $v_j = b_j / (K^T u)_j$).

Sinkhorn algorithm Initialize $v^{(0)} (= I_n, \text{e.g.})$

$$u^{(k)} = \frac{a}{K v^{(k-1)}}, \quad v^{(k)} = \frac{b}{K^T u^{(k)}} \quad (k=1, 2, \dots).$$

row sum of $u^{(k)} \odot K v^{(k-1)} = a$, col.sum of $v^{(k)} \odot K^T u^{(k)} = b$.

Or: Initialize $u^{(0)} > 0$ and set

$$v^{(k)} = \frac{b}{K^T u^{(k-1)}} \quad \text{and} \quad u^{(k)} = \frac{a}{K v^{(k)}} \quad (k=1, 2, \dots).$$

Some remarks

① The alternate iteration is equivalent to row-column alternate normalization (process) as initially studied in Richard Sinkhorn (1964) (only for $m=n$, $a_i=b_i=1/n$).

Notation. $G \in \mathbb{R}^{m \times n}$:

$r_i(G)$ = row- i sum of G , $c_j(G)$ = col- j sum of G , $\forall i, j$.

Definition Let $a \in \mathbb{P}_m$, $b \in \mathbb{P}_n$ and $G \in \mathbb{R}^{m \times n}$ with $a > 0$, $b > 0$, and $G > 0$.

(1) Let $d_i = \frac{1}{a_i} r_i(G)$ ($1 \leq i \leq m$) and define $G^a \in \mathbb{R}^{m \times n}$ by

$$G_{i,j}^a = G_{i,j} / d_i, \quad 1 \leq i \leq m, \quad 1 \leq j \leq n.$$

Call G^a the row-normalization of G w.r.t. a .

(2) Let $e_j = \frac{1}{b_j} c_j(G)$ ($1 \leq j \leq n$) and define $G^b \in \mathbb{R}^{m \times n}$ by

$$G_{i,j}^b = G_{i,j} / e_j, \quad 1 \leq i \leq m, \quad 1 \leq j \leq n.$$

Call G^b the column-normalization of G w.r.t. b .

Note: $d_i = (\text{row-}i \text{ sum of } G) / a_i$.

$$\begin{bmatrix} & G^a \\ \cdot & \ddots & \ddots & \ddots \end{bmatrix} = \begin{bmatrix} & G \\ \cdot & \ddots & \ddots & \ddots \end{bmatrix} \begin{matrix} 1/\lambda_1 \\ \vdots \\ 1/\lambda_i \\ \vdots \\ 1/\lambda_m \end{matrix}$$

$$\text{row-}i \text{ sum of } G^a = \sum_j G_{ij}^a = \sum_j G_{ij} / \lambda_i \\ = a_i \cdot \sum_i G_{ij} / (\text{row-}i \text{ sum of } G) = a_i, \forall i.$$

So, Sinkhorn's alg. is just the alternate row-column normalization process. In fact, this is the original Sinkhorn process/iteration.

○ Sinkhorn (1964) mentioned that Lloyd Welch used the normalization algorithm in an unpublished report of Inst. of Defense Analysis, no time was given. Sinkhorn-Knopp (1967) extended the result. So, the method is sometimes called the Sinkhorn-Knopp algorithm.

○ Sinkhorn's alg. (with $v^{(0)}$ given) produces $A^{(k)} = \text{diag}(u^{(k)}) K \text{diag}(v^{(k-1)})$ and $B^{(k)} = \text{diag}(u^{(k)}) K \text{diag}(v^{(k)})$ ($k=1, 2, \dots$) that converge to $P \in \mathcal{A}(a, b)$; see below. But for any finite k , $A^{(k)}$ and $B^{(k)}$ may not be in $\mathcal{A}(a, b)$. One idea is to "project" some $A^{(k)}$ or $B^{(k)}$ to $\mathcal{A}(a, b)$. Altschuler et al. (2017)

suggested the following alg. to compute some $P \in \mathcal{A}(a, b)$ that is close to a positive matrix F (like $A^{(k)}$ or $B^{(k)}$):

$$\alpha \wedge \beta = \min(\alpha, \beta)$$

$$\cdot \quad x_i = \frac{\alpha_i}{r_i(F)} \wedge 1 \quad (i=1, \dots, m)$$

$$\cdot \quad F' = \text{diag}(x)F$$

$$\cdot \quad y_j = \frac{b_j}{c_j(F')} \wedge 1 \quad (j=1, \dots, n)$$

$$\cdot \quad F'' = F' \text{diag}(y)$$

$$\cdot \quad P = F'' + [a - r(F')] [b - c(F'')]^T / \|a - r(F'')\|_1.$$

Here $r_i(F) = \text{row-}i \text{ sum of } F$,

$c_j(F') = \text{col-}j \text{ sum of } F'$

$r(F') = \text{the vector of } r_i(F')$

$c(F'') = \text{the vector of } c_j(F'')$

This is not the closed distance projection. But the convergence (together with Sinkhorn) is proved.

Theorem Let $a \in \mathbb{R}_m^m, b \in \mathbb{R}_n^m, a > 0, b > 0$. Let $K \in \mathbb{R}^{m \times n}$ with $K > 0$. Let $u, \bar{u} \in \mathbb{R}^m, v, \bar{v} \in \mathbb{R}^n$, $u > 0, \bar{u} > 0, v > 0, \bar{v} > 0$. Suppose $P = \text{diag}(u) K \text{diag}(v) \in \mathcal{A}(a, b)$ and $\bar{P} = \text{diag}(\bar{u}) K \text{diag}(\bar{v}) \in \mathcal{A}(a, b)$. Then $\bar{P} = P$. Moreover, $\exists p > 0$ such that $u = p\bar{u}$ and $v = \frac{1}{p}\bar{v}$.

[This means that the pair $u > 0, v > 0$ such that $(u; K_{ij}; v_j) \in \mathcal{A}(a, b)$ is unique up to a multiplicative factor.]

Proof $P = \text{diag}(u) K \text{diag}(v) \Leftrightarrow P_{ij} = u_i K_{ij} v_j (> 0) \quad \forall i, j$

Similarly, $\bar{P}_{ij} = \bar{u}_i K_{ij} \bar{v}_j \quad \forall i, j$. $K_{ij} = (\bar{u}_i)^{-1} \bar{P}_{ij} (\bar{v}_j)^{-1}$. So,

$P_{ij} = (u_i/\bar{u}_i) \bar{P}_{ij} (v_j/\bar{v}_j)$. Let $\alpha_i = u_i/\bar{u}_i$ and $\beta_j = v_j/\bar{v}_j$.

Then $P_{ij} = \alpha_i \bar{P}_{ij} \beta_j \quad \forall i, j$. It suffices to show that $\exists p > 0$ s.t. $\alpha_i = p$ and $\beta_j = \frac{1}{p}$ $\forall i, j$.

Let $\alpha_{im} = \min_i \alpha_i$ and $\beta_{jm} = \max_j \beta_j$. Since $\bar{P} \in \mathcal{A}(a, b)$ and $P \in \mathcal{A}(a, b)$, we have

$$\begin{aligned}\alpha_{im} \beta_{jm} &= \frac{\alpha_{im} \beta_{jm}}{b_{jm}} \sum_i \bar{P}_{ijm} \stackrel{(A)}{\leq} \frac{1}{b_{jm}} \sum_i \alpha_i \bar{P}_{ijm} \beta_{jm} \\ &= \frac{1}{b_{jm}} \sum_i P_{ijm} = 1.\end{aligned}$$

$$\begin{aligned}\alpha_{im} \beta_{jm} &= \frac{\alpha_{im} \beta_{jm}}{\alpha_{im}} \sum_j \bar{P}_{imj} \geq \frac{1}{\alpha_{im}} \sum_j \alpha_{im} \bar{P}_{imj} \beta_j \\ &= \frac{1}{\alpha_{im}} \sum_j P_{imj} = 1.\end{aligned}$$

Hence, $\alpha_{im} \beta_{jm} = 1$, and all above are equal. (A) is an equality. So,

$$\alpha_{im} \sum_i \bar{P}_{ijm} = \sum_i \alpha_i \bar{P}_{ijm}, \quad \sum_i (\alpha_i - \alpha_{im}) \bar{P}_{ijm} = 0 \Rightarrow \alpha_i = \alpha_{im} \quad \forall i$$

Similarly, $\beta_j = \beta_{im} = \frac{1}{\alpha_{im}} \quad \forall j$.

Set $p = \alpha_i \quad \forall i$. Then $u = p \bar{u}$ and $v = \frac{1}{p} \bar{v}$. Clearly, then $P = \bar{P}$

QED