

Lecture 10, Monday, 4/18/2022

Today: More about Sinkhorn's algorithm.

- Finish the proof of convergence of Sinkhorn's alg.
- Convergence rate via Hilbert's metric.

Hilbert's projective metric. (Or: Cayley-Hilbert metric)

Denote $\mathbb{R}_+^n = \{x \in \mathbb{R}^n : x > 0\}$. Let $x, y \in \mathbb{R}_+^n$. Define

$$\begin{aligned} d_H(x, y) &= \log \max_{i,j} \frac{x_i/y_i}{x_j/y_j} = \log \max_{i,j} \frac{y_i/x_i}{y_j/x_j} \\ &= \log \max_{1 \leq i, j \leq n} \frac{x_i y_j}{x_j y_i} \left(= \max_{1 \leq i, j \leq n} \log \frac{x_i y_j}{x_j y_i} \right). \end{aligned}$$

One can verify:

$$(1) \quad d_H(x, y) \geq 0 \quad \forall x, y \in \mathbb{R}_+^n. \quad d_H(x, y) = 0 \iff \frac{x}{\|x\|} = \frac{y}{\|y\|}$$

($\iff \exists \lambda > 0$ such that $x = \lambda y$)

$$(2) \quad d_H(x, y) = d_H(y, x) \quad \forall x, y \in \mathbb{R}_+^n$$

$$(3) \quad d_H(x, y) \leq d_H(x, z) + d_H(z, y) \quad \forall x, y, z \in \mathbb{R}_+^n$$

Prove (3). $\max_{i,j} \frac{x_i y_j}{x_j y_i} = \frac{x_{i_0} y_{j_0}}{x_{j_0} y_{i_0}} = \frac{x_{i_0} z_{j_0}}{x_{j_0} z_{i_0}} \cdot \frac{y_{j_0} z_{i_0}}{y_{i_0} z_{j_0}}$
 $\leq \max_{i,j} \frac{x_i z_j}{x_j z_i} \cdot \max_{i,j} \frac{y_i z_j}{y_j z_i}$ (for some i_0, j_0). QED

Clearly, $\forall x, y \in \mathbb{R}_+^n \quad \forall \lambda, \mu > 0,$

$$d_H(x, y) = d_H\left(\frac{x}{\|x\|}, \frac{y}{\|y\|}\right) = d_H(\lambda x, \mu y).$$

$$d_H(x, y) = d_H\left(\frac{x}{\|x\|}, \mathbb{I}_n\right).$$

Therefore, d_H is a metric on $\mathbb{R}_+^n/\mathbb{I}$ where $x \sim y \iff \frac{x}{\|x\|} = \frac{y}{\|y\|}$.

It is called Hilbert's projective metric on \mathbb{R}_+^n .

Proposition $(\mathbb{R}_+^n/\mathbb{I}, d_H)$ is a complete metric space.

Proof Only the completeness. Assume $x^{(k)} \in \mathbb{R}^n$ ($k=1, 2, \dots$) and $d_H(x^{(k)}, x^{(\ell)}) \rightarrow 0$ as $k, \ell \rightarrow \infty$. Without loss of generality, we may assume $\|x^{(k)}\|=1$ (ℓ^2 -norm) for all k . For any $\varepsilon > 0$, $d_H(x^{(k)}, x^{(\ell)}) \rightarrow 0$ means $\exists N$ s.t. $k, \ell \geq N$ $\Rightarrow d_H(x^{(k)}, x^{(\ell)}) \leq \varepsilon$. i.e., $1 \leq \max_{i,j} \frac{x_i^{(k)} x_j^{(\ell)}}{\sqrt{x_i^{(k)} x_i^{(\ell)}}} \leq 1 + \varepsilon$. (*)

Since $\|x^{(k)}\|=1$, \exists a subseq $\{x^{(k')}\}$ of $\{x^{(k)}\}$ s.t. $\|x^{(k')}-x\| \rightarrow 0$ for some $x \in \mathbb{R}^n$ with $\|x\|=1$. We show that $x > 0$. From (*), with k' replacing k , we get $x_i^{(k')} x_j^{(\ell)} \leq x_j^{(k')} x_i^{(\ell)} (1+\varepsilon)$ $\forall i, j, k', \ell \geq N$. Setting $\ell=N$ and sending $k' \rightarrow \infty$, we get $x_i^{(N)} x_j^{(N)} \leq x_j^{(N)} x_i^{(N)(1+\varepsilon)}$ $\forall i, j$. If $x_{j_0} = 0$ for some j_0 , then all $x_i = 0$. This is in contradiction with $\|x\|=1$. Hence $x > 0$. From (*), replacing k by k' , and letting $k' \rightarrow \infty$, we get $1 \leq \max_{i,j} \frac{x_i x_j^{(\ell)}}{\sqrt{x_i x_i^{(\ell)}}} \leq 1 + \varepsilon \quad \forall \ell \geq N$. i.e., $d_H(x, x^{(\ell)}) \leq \log(1+\varepsilon) \leq \varepsilon$ if $\ell \geq N$. Hence, $d_H(x, x^{(\ell)}) \rightarrow 0$ as $\ell \rightarrow \infty$. QED.

Denote for $x \in \mathbb{R}^n$ $\|x\|_{\text{var}} = \max_i x_i - \min_i x_i$. Then

$$d_H(x, y) = \|\log x - \log y\|_{\text{var}}$$

where $\log x = (\log x_i)$. This can be shown by definition:

$$\begin{aligned} d_H(x, y) &= \max_{i,j} \log \frac{x_i y_j}{x_j y_i} \\ &= \max_{i,j} [\log x_i - \log y_i] - [\log x_j - \log y_j] \\ &= \max_i (\log x_i - \log y_i) - \min_j (\log x_j - \log y_j) \\ &= \|\log x - \log y\|_{\text{var}}. \quad \underline{\text{QED}} \end{aligned}$$

Projective diameter and contraction ratio of a positive matrix

Now, denote $\mathbb{R}_+^{m \times n} = \{A = [A_{ij}] \in \mathbb{R}^{m \times n} : \text{all } A_{ij} > 0\}$,
Definition Let $A \in \mathbb{R}_+^{m \times n}$.

(1) Denote

$$\theta(A) = \sup \{d_H(Ax, Ay) : x, y \in \mathbb{R}_+^n\},$$

and call it the projective diameter of A .

(2) Denote

$$K(A) = \sup \left\{ \frac{d_H(Ax, Ay)}{d_H(x, y)} : x, y \in \mathbb{R}_+^n, \frac{x}{\|x\|_1} \neq \frac{y}{\|y\|_1} \right\},$$

call it the contraction ratio of A with respect to Hilbert's metric.

Proposition:

$$\theta(A) = \max_{i, j, k, l} \log \frac{a_{ik} a_{jl}}{a_{il} a_{jk}} = \log \max_{i, j, k, l} \frac{a_{ik} a_{jl}}{a_{il} a_{jk}}.$$

Proof Let $M = \max_{i, j, k, l} \frac{a_{ik} a_{jl}}{a_{il} a_{jk}}$. Then $a_{ik} a_{jl} \leq M a_{il} a_{jk}$ for all i, j, k, l . Now $\forall x, y \in \mathbb{R}_+^n$. We have for any i, j'

$$\begin{aligned} \frac{(Ax)_i (Ay)_{j'}}{(Ax)_{j'} (Ay)_i} &= \frac{\left(\sum_{\alpha} a_{i\alpha} x_{\alpha} \right) \left(\sum_{\beta} a_{j'\beta} y_{\beta} \right)}{\left(\sum_{\alpha} a_{j'\alpha} x_{\alpha} \right) \left(\sum_{\beta} a_{i\beta} y_{\beta} \right)} = \frac{\sum_{\alpha, \beta} a_{i\alpha} a_{j'\beta} x_{\alpha} y_{\beta}}{\sum_{\alpha, \beta} a_{j'\alpha} a_{i\beta} x_{\alpha} y_{\beta}} \\ &\leq \frac{M \sum_{\alpha, \beta} a_{i\alpha} a_{j'\beta} x_{\alpha} y_{\beta}}{\sum_{\alpha, \beta} a_{j'\alpha} a_{i\beta} x_{\alpha} y_{\beta}} = M. \end{aligned}$$

Hence, $d_H(Ax, Ay) = \log \max_{i, j} \frac{(Ax)_i (Ay)_{j'}}{(Ax)_{j'} (Ax)_i} \leq \log M = \theta(A).$

$$\sup \{d_H(Ax, Ay) : x \in \mathbb{R}_+^n, y \in \mathbb{R}_+^n\} \leq \theta(A).$$

Conversely, assume $M = \frac{a_{ik} a_{jl}}{a_{il} a_{jk}}$ for some i, j, k, l .

Let $\varepsilon \in (0,1)$ and define $x, y \in \mathbb{R}_+^n$ by $x_k = 1$ and $x_j = \varepsilon$ if $j \neq k$, $y_\ell = 1$ and $y_j = \varepsilon$ if $j = \ell$. Then

$$\frac{(Ax)_i \cdot (Ay)_j}{(Ax)_j \cdot (Ay)_i} = \frac{\left(\sum_{\alpha} a_{i\alpha} x_{\alpha}\right) \left(\sum_{\beta} a_{j\beta} y_{\beta}\right)}{\left(\sum_{\alpha} a_{j\alpha} x_{\alpha}\right) \left(\sum_{\beta} a_{i\beta} y_{\beta}\right)} = \frac{(a_{ik} + O(\varepsilon))(a_{j\ell} + O(\varepsilon))}{(a_{jk} + O(\varepsilon))(a_{i\ell} + O(\varepsilon))} = M_{\varepsilon}$$

So, $e^{O(A)} \geq e^{d_H(Ax, Ay)} \geq \frac{(Ax)_i \cdot (Ay)_j}{(Ax)_j \cdot (Ay)_i} = M_{\varepsilon} \rightarrow M$. GED

Proposition (1) $O(A) = O(A^T)$ and $K(A) = K(A^T)$.

(2) If $A \in \mathbb{R}_+^{m \times n}$ and $B \in \mathbb{R}_+^{n \times \ell}$, then $K(AB) \leq K(A)K(B)$.

(3) If $A, B \in \mathbb{R}_+^{m \times n}$ are diagonally equivalence, defined by $A = \text{diag}(u) B \text{ diag}(v)$ for some $u \in \mathbb{R}_+^m$ and $v \in \mathbb{R}_+^n$, then $K(A) = K(B)$.

Proof (1) The fact $O(A) = O(A^T)$ follows from the formula for $O(A)$ (cf. proposition above). This together with Birkhoff-Hopf Thm implies $K(A^T) = K(A)$.

(2) $\forall x, y \in \mathbb{R}_+^{\ell}$: $\frac{d_H(ABx, ABy)}{d_H(x, y)} = \frac{d_H(ABx, ABy)}{d_H(By, By)} \cdot \frac{d_H(By, y)}{d_H(x, y)} \leq K(A)K(B)$. Hence, $K(AB) \leq K(A)K(B)$.

(3) $\forall x, y \in \mathbb{R}_+^n$. Define $\hat{x} = v \odot x$ (i.e., $\hat{x}_i = v_i x_i$), $\hat{y} = v \odot y$. Clearly, $\frac{\hat{x}_i \hat{y}_j}{\hat{x}_j \hat{y}_i} = \frac{x_i y_j}{x_j y_i}$. Hence $d_H(x, y) = d_H(\hat{x}, \hat{y})$.

Since $A_{ij} = u_i B_{ij} v_j$ $\forall i, j$, we have

$$\begin{aligned} \frac{(Ax)_i \cdot (Ay)_j}{(Ax)_j \cdot (Ay)_i} &= \frac{\left(\sum_{\alpha} A_{i\alpha} x_{\alpha}\right) \left(\sum_{\beta} A_{j\beta} y_{\beta}\right)}{\left(\sum_{\alpha} A_{j\alpha} x_{\alpha}\right) \left(\sum_{\beta} A_{i\beta} y_{\beta}\right)} \\ &= \frac{u_i \left(\sum_{\alpha} B_{i\alpha} v_{\alpha} x_{\alpha}\right) u_j \left(\sum_{\beta} B_{j\beta} v_{\beta} y_{\beta}\right)}{u_j \left(\sum_{\alpha} B_{j\alpha} v_{\alpha} x_{\alpha}\right) u_i \left(\sum_{\beta} B_{i\beta} v_{\beta} y_{\beta}\right)} \end{aligned}$$

$$= \frac{(\sum_{\alpha} B_{i\alpha} \hat{x}_{\alpha})(\sum_{\beta} B_{j\beta} \hat{y}_{\beta})}{(\sum_{\alpha} B_{j\alpha} \hat{x}_{\alpha})(\sum_{\beta} B_{i\beta} \hat{y}_{\beta})} = \frac{(B\hat{x})_i \cdot (B\hat{y})_j}{(B\hat{x})_j \cdot (B\hat{y})_i}.$$

Hence, $d_H(Ax, Ay) = d_H(B\hat{x}, B\hat{y})$, and thus

$$\frac{d_H(Ax, Ay)}{d_H(x, y)} = \frac{d_H(B\hat{x}, B\hat{y})}{d_H(\hat{x}, \hat{y})}.$$

Since $(x, y) \rightarrow (\hat{x}, \hat{y})$ is a bijection, $\kappa(A) = \kappa(B)$ QED

Birkhoff (1957) - Hopf 1963 Then let $A \in \mathbb{R}_+^{m \times n}$ and denote $\eta(A) = e^{\Theta(A)} = \max_{i,j} \frac{a_{ik} a_{je}}{a_{il} a_{jk}} (> 1)$. Then

$$\kappa(A) = \frac{\sqrt{\eta(A)} - 1}{\sqrt{\eta(A)} + 1} \in (0, 1).$$

Given $a \in \mathbb{P}_m \cap \mathbb{R}_+^m$, $b \in \mathbb{P}_n \cap \mathbb{R}_+^n$, and $K \in \mathbb{R}_+^{m \times n}$

Recall Sinkhorn's algorithm: Select $v^{(0)} \in \mathbb{R}_+^n$.

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

Set for $k=1, 2, \dots$:

$$A^{(k)} = \text{diag}(u^{(k)}) K \text{diag}(v^{(k-1)}), B^{(k)} = \text{diag}(u^{(k)}) K \text{diag}(v^{(k)}).$$

The original Sinkhorn's construction (or process) uses the row sum and column sum to alternatively normalizing the matrix to have the row sum = a and col. sum = b :

$A^{(k)} \rightarrow B^{(k)}$: calculate $\lambda_j^{(k)} = \frac{1}{b_j} \text{col-}j \text{ sum of } A^{(k)}$, set $B_{ij}^{(k)} = A_{ij}^{(k)} / \lambda_j^{(k)}$.

$B^{(k)} \rightarrow A^{(k+1)}$: calculate $\lambda_i^{(k)} = \frac{1}{a_i} \text{row-}i \text{ sum of } B^{(k)}$, set $A_{ij}^{(k+1)} = B_{ij}^{(k)} / \lambda_i^{(k)}$.

We have verified that (see Lecture 9)

$$\lambda^{(k)} = \frac{u^{(k)}}{u^{(k+1)}} \text{ and } u^{(k)} = \frac{v^{(k-1)}}{v^{(k)}} , \quad k=1,2,\dots.$$

Note that

$$\text{row-sum of } A^{(k)} = a, \quad \text{col. sum of } B^{(k)} = b.$$

$$\lambda^{(k)} \odot a = \text{row sum of } B^{(k)}, \quad u^{(k)} \odot b = \text{col. sum of } A^{(k)}.$$

By Sinkhorn's theorem (or rather its proof),

$$\lambda^{(k)} \rightarrow \mathbb{1}_m \text{ and } u^{(k)} \rightarrow \mathbb{1}_n.$$

$$u_i^{(k)} v_j^{(k)} \rightarrow u_i v_j \quad \forall i,j \text{ for some } u \in \mathbb{R}_+^m, v \in \mathbb{R}_+^n.$$

$$P := \text{diag}(u) K \text{ diag}(v) \in \mathcal{A}(a, b)$$

$$A^{(k)} \rightarrow P \text{ and } B^{(k)} \rightarrow P.$$

The following theorem provides the convergence rate of Sinkhorn's algorithm.

Theorem (Franklin-Lorenz 1989) We have

$$(1) \quad d_H(u^{(k)}, u) \leq [\kappa(K)]^{2k-1} d_H(v^{(0)}, v),$$

$$d_H(v^{(k)}, v) \leq [\kappa(K)]^{2k} d_H(v^{(0)}, v).$$

$$(2) \quad d_H(\lambda^{(k)}, \mathbb{1}_m) \leq [\kappa(K)]^{2k-2} d_H(\lambda^{(1)}, \mathbb{1}_m),$$

$$d_H(u^{(k)}, \mathbb{1}_n) \leq [\kappa(K)]^{2k} d_H(v^{(0)}, v).$$

$$(3) \quad d_H(u^{(k)}, u) \leq \frac{[\kappa(K)]^{2k-2}}{1 - [\kappa(K)]^2} d_H(\lambda^{(1)}, \mathbb{1}_m),$$

$$d_H(v^{(k)}, v) \leq \frac{[\kappa(K)]^{2k-2}}{1 - [\kappa(K)]^2} d_H(u^{(1)}, \mathbb{1}_n).$$

$$(4) \|\log B^{(k)} - \log P\|_\infty \leq d_H(u^{(k)}, u) \leq \frac{[K(K)]^{2k-2}}{1 - [K(K)]^2} d_H(\lambda^{(1)}, \mathbf{1}_m).$$

$$\|\log A^{(k+1)} - \log P\|_\infty \leq d_H(v^{(k)}, v) \leq \frac{[K(K)]^{2k-2}}{1 - [K(K)]^2} d_H(u^{(1)}, \mathbf{1}_n).$$

Proof (1) We have for any $k \geq 1$,

$$\begin{aligned} d_H(u^{(k)}, u) &= d_H\left(\frac{a}{K^T v^{(k-1)}}, \frac{a}{K v}\right) = d_H(K^T v^{(k-1)}, K^T v) \\ &\leq K(K) d_H(v^{(k-1)}, v) = K(K) d_H\left(\frac{b}{K u^{(k-1)}}, \frac{b}{K u}\right) \\ &= K(K) d_H(K u^{(k-1)}, K u) \leq [K(K)]^2 d_H(u^{(k-1)}, u) \quad (\times) \\ &\leq \dots \leq [K(K)]^{2(k-1)} d_H(u^{(1)}, u) = [K(K)]^{2k-2} d_H\left(\frac{a}{K^T v^{(0)}}, \frac{a}{K^T v}\right) \\ &= [K(K)]^{2k-2} d_H(K^T v^{(0)}, K^T v) \leq [K(K)]^{2k-2} K(K^T) d_H(v^{(0)}, v) \\ &= [K(K)]^{2k-1} d_H(v^{(0)}, v). \end{aligned}$$

Similarly,

$$\begin{aligned} d_H(v^{(k)}, v) &= d_H\left(\frac{b}{K u^{(k)}}, \frac{b}{K u}\right) = d_H(K u^{(k)}, K u) \\ &\leq K(K) d_H(u^{(k)}, u) \leq [K(K)]^{2k} d_H(v^{(0)}, v). \end{aligned}$$

$$\begin{aligned} (2) \quad d_H(\lambda^{(k)}, \mathbf{1}_m) &= d_H\left(\frac{u^{(k)}}{u^{(k+1)}}, \mathbf{1}_m\right) = d_H(u^{(k)}, u^{(k+1)}) \\ &\leq [K(K)]^2 d_H(u^{(k-1)}, u^{(k)}) \leq \dots \leq [K(K)]^{2k-2} d_H(u^{(1)}, u^{(2)}) \\ &= [K(K)]^{2k-2} d_H\left(\frac{u^{(1)}}{u^{(2)}}, \mathbf{1}_m\right) = [K(K)]^{2k-2} d_H(\lambda^{(1)}, \mathbf{1}_m). \end{aligned}$$

Similarly,

$$d_H(u^{(k)}, \mathbf{1}_n) \leq [K(K)]^{2k-2} d_H(u^{(1)}, \mathbf{1}_n).$$

(3) By the triangle inequality, we have

$$\begin{aligned}
d_H(u^{(k)}, u) &\leq d_H(u^{(k)}, u^{(k+1)}) + d_H(u^{(k+1)}, u) \\
&\stackrel{(*)}{=} d_H\left(\frac{u^{(k)}}{u^{(k+1)}}, \mathbf{1}_m\right) + [\kappa(K)]^2 d_H(u^{(k)}, u) \\
&= d_H(\lambda^{(k)}, \mathbf{1}_m) + [\kappa(K)]^2 d_H(u^{(k)}, u).
\end{aligned}$$

Hence,

$$\begin{aligned}
d_H(u^{(k)}, u) &\leq \frac{1}{1 - [\kappa(K)]^2} d_H(\lambda^{(k)}, \mathbf{1}_m) \\
&\leq \frac{[\kappa(K)]^{2k-2}}{1 - [\kappa(K)]^2} d_H(\lambda^{(1)}, \mathbf{1}_m).
\end{aligned}$$

Similarly,

$$d_H(v^{(k)}, v) \leq \frac{[\kappa(K)]^{2k-2}}{1 - [\kappa(K)]^2} d_H(u^{(k)}, \mathbf{1}_n).$$

(F) Now, fix $k \geq 1$. Let $x_i = \frac{u_i^{(k)}}{u_i^{(k+1)}}, x(m) = m \text{ in } x_i$, and $y_j = \frac{v_j^{(k)}}{v_j^{(k+1)}}$. Then, $\frac{\beta_{ij}^{(k)}}{p_{ij}} = \frac{u_i^{(k)} K_{ij} v_j^{(k)}}{u_i^{(k+1)} K_{ij} v_j^{(k+1)}} = x_i \cdot y_j$.

Let $\varepsilon = d_H(u^{(k)}, u)$. Then,

$$\varepsilon = d_H(u^{(k)}, u) = d_H\left(\frac{u^{(k)}}{u}, \mathbf{1}_m\right) = d_H(x, \mathbf{1}_m) = \log \max_{i,j} \frac{x_i}{x_j}.$$

Hence $x(m) \leq x_i \leq x(j) e^\varepsilon$. Let $x_j = x(m)$. We get

$$x(m) \leq x_i \leq x(m) e^\varepsilon. \text{ Hence } \frac{1}{x(m)} \geq \frac{1}{x_i} \geq \frac{1}{x(m)} e^{-\varepsilon}. \quad (*)$$

Since $p_{ij} y_j = x_i^{-1} \beta_{ij}^{(k)}$, col. sum of β = b , and

col. sum of $\beta^{(k)}$ = b , we have

$$y_j b_j = y_j \sum_i p_{ij} = \sum_i x_i^{-1} \beta_{ij}^{(k)},$$

$$\frac{e^{-\varepsilon} b_j}{x(m)} = \frac{1}{x(m)} e^{-\varepsilon} \sum_i \beta_{ij}^{(k)} \leq \sum_i x_i^{-1} \beta_{ij}^{(k)} \leq \frac{1}{x(m)} \sum_i \beta_{ij}^{(k)} = \frac{1}{x(m)} b_j.$$

Hence $\frac{e^{-\varepsilon}}{x(m)} b_j \leq y_j - b_j \leq \frac{1}{x(m)} b_j$, so, $\frac{e^{-\varepsilon}}{x(m)} \leq y_j \leq \frac{1}{x(m)}$.

Thus, $e^{-\varepsilon} \leq x_i y_j \leq e^{\varepsilon}$. But $x_i y_j = \beta_{ij}^{(u)} / p_{ij}$. $\forall i, j$.

Hence $-\varepsilon \leq \log \beta_{ij}^{(u)} - \log p_{ij} \leq \varepsilon, \forall i, j$.

Thus, $\| \log \beta^{(u)} - \log p \| \leq d_H(u^{(u)}, u)$.

This and (3) imply the first inequality in (4). The second one is similar. QED