

Homework 6

HANMO CHEN

December 15, 2020

-
- **Acknowledgments:** For L_2 -norm I refer to the wikipedia https://en.wikipedia.org/wiki/Matrix_norm. For 6.4, I refer to <https://ieeexplore.ieee.org/document/8849720>¹
 - **Collaborators:** I finish this homework all by myself.
 - *I certify that all solutions are entirely in my words and that I have not looked at another student's solutions. I have credited all external sources in this write up.*

Hanmo Chen

6.1. (a)

$$p_{y_1 y_2}(y_1, y_2; x) = 1, x \leq y_1, y_2 \leq x + 1 \quad (1)$$

And $p(s = \max(y_1, y_2) \leq s) = (s - x)^2, x \leq s \leq x + 1$, so

$$p_s(s; x) = 2(s - x), x \leq s \leq x + 1 \quad (2)$$

$$\frac{p_{y_1 y_2}(y_1, y_2; x)}{p_s(s; x)} = \frac{1}{2(s - x)}, x \leq y_1, y_2, s \leq x + 1 \quad (3)$$

So s is not a sufficient statistic for $p_{y_1 y_2}(y_1, y_2; x)$.

(b) $p_r(r; x) = p_r(-r; x)$ due to symmetry. Suppose $r \geq 0$,

$$p_r(r; x) = \int_x^{x+1} p_{y_1}(y + r; x) p_{y_2}(y; x) dy = 1 - r \quad (4)$$

So $p_r(r; x) = 1 - |r|$ does not depend on x , r is an ancillary statistic for $p_{y_1 y_2}(y_1, y_2; x)$.

(c) Yes. Because we can recover y_1, y_2 , with s, r .

$$(y_1, y_2) = \begin{cases} (s, s - r), & r \geq 0 \\ (s + r, s), & r < 0 \end{cases} \quad (5)$$

So $\mathbf{u} = [s, r]^T$ is a sufficient statistic for $p_{y_1 y_2}(y_1, y_2; x)$ and

¹S. Huang, X. Xu, L. Zheng and G. W. Wornell, "An Information Theoretic Interpretation to Deep Neural Networks," 2019 IEEE International Symposium on Information Theory (ISIT), Paris, France, 2019, pp. 1984-1988, doi: 10.1109/ISIT.2019.8849720.

6.2. Let $x = y$ be randomly chosen from $\{0, 1, 2\}$ with equal probability $\frac{1}{3}$ and $\xi(x) = x, \eta(y) = y^2$. $\rho(x, y) = \rho(x, x) = 1$.

$$\rho(\xi(x), \eta(y)) = \rho(x, x^2) = \frac{\mathbb{E}[x^3] - \mathbb{E}[x^2] \mathbb{E}[x]}{\sqrt{\text{var}(x) \text{var}(x^2)}} = \sqrt{\frac{12}{13}} \neq \rho(x, y) \quad (6)$$

6.3. (a) i.

$$\begin{aligned} \psi^T \mathbf{B} \phi &= \sum_{x \in \mathcal{X}, y \in \mathcal{Y}} g(y) \sqrt{P_y(y)} \frac{P_{xy}(x, y)}{\sqrt{P_x(x) P_y(y)}} f(x) \sqrt{P_x(x)} \\ &= \sum_{x \in \mathcal{X}, y \in \mathcal{Y}} f(x) g(y) P_{xy}(x, y) \\ &= \mathbb{E}[f(x) g(y)] \end{aligned} \quad (7)$$

ii.

$$\begin{aligned} \mathbf{B} \phi(y) &= \sum_{x \in \mathcal{X}} \frac{P_{xy}(x, y)}{\sqrt{P_y(y)}} f(x) \\ &= \sqrt{P_y(y)} \sum_{x \in \mathcal{X}} P_{x|y}(x|y) f(x) \\ &= \sqrt{P_y(y)} \mathbb{E}[f(x)|y] \end{aligned} \quad (8)$$

So $\mathbf{B} \phi \leftrightarrow \mathbb{E}[f(x)|y]$

iii.

$$\begin{aligned} \mathbf{B}^T \psi(x) &= \sum_{y \in \mathcal{Y}} \frac{P_{xy}(x, y)}{\sqrt{P_x(x)}} g(y) \\ &= \sqrt{P_x(x)} \sum_{y \in \mathcal{Y}} P_{y|x}(y|x) g(y) \\ &= \sqrt{P_x(x)} \mathbb{E}[g(y)|x] \end{aligned} \quad (9)$$

So $\mathbf{B}^T \psi \leftrightarrow \mathbb{E}[g(y)|x]$

(b)

$$\mathbf{B} \phi_1(y) = \sum_{x \in \mathcal{X}} \frac{P_{xy}(x, y)}{\sqrt{P_y(y)}} = \sqrt{P_y(y)} \quad (10)$$

$$\mathbf{B}^T \psi_1(x) = \sum_{y \in \mathcal{Y}} \frac{P_{xy}(x, y)}{\sqrt{P_x(x)}} = \sqrt{P_x(x)} \quad (11)$$

So $\mathbf{B} \phi_1 = \psi_1, \mathbf{B}^T \psi_1 = \phi_1$. Its means that if we let $f(x) = g(y) = 1$ in (a), we have $\mathbb{E}[1|y] = \mathbb{E}[1|x] = 1$

(c) The L_2 -norm of \mathbf{B} is equivalent to

$$\|\mathbf{B}\|_2 = \sup \|\mathbf{B} \mathbf{x}\|_2, \text{ subject to } \|\mathbf{x}\|_2 = 1 \quad (12)$$

To maximize $\|\mathbf{B} \mathbf{x}\|_2$, every component in \mathbf{x} must ≥ 0 . Thus for every $\mathbf{X} \in \mathbb{R}^{|\mathcal{X}|}$, there is a correspond distribution of $X \in \mathcal{X}$. So \mathbf{X} can be expressed as

$$\mathbf{x} = \left[\sqrt{P_{\mathbf{x}'}(1)}, \sqrt{P_{\mathbf{x}'}(2)}, \dots, \sqrt{P_{\mathbf{x}'}(|\mathcal{X}|)} \right]^T \quad (13)$$

Thus

$$\begin{aligned} \|\mathbf{B}\mathbf{x}\|_2^2 &= \sum_{y \in \mathcal{Y}} \left(\sum_{x \in \mathcal{X}} \frac{P_{xy}(x, y)}{\sqrt{P_{\mathbf{x}}(x)P_{\mathbf{y}}(y)}} \sqrt{P_{\mathbf{x}'}(x)} \right)^2 \\ &= \sum_{y \in \mathcal{Y}} P_{\mathbf{y}}(y) \left(\sum_{x \in \mathcal{X}} \frac{P_{x|y}(x|y)}{\sqrt{P_{\mathbf{x}}(x)}} \sqrt{P_{\mathbf{x}'}(x)} \right)^2 \\ &= \sum_{y \in \mathcal{Y}} P_{\mathbf{y}}(y) \mathbb{E}^2 \left[\frac{\sqrt{P_{\mathbf{x}'}(x)}}{\sqrt{P_{\mathbf{x}}(x)}} | y \right] \\ &\leq \sum_{y \in \mathcal{Y}} P_{\mathbf{y}}(y) \mathbb{E} \left[\frac{P_{\mathbf{x}'}(x)}{P_{\mathbf{x}}(x)} | y \right] = \mathbb{E} \left[\mathbb{E} \left[\frac{P_{\mathbf{x}'}(x)}{P_{\mathbf{x}}(x)} | y \right] \right] \\ &= \mathbb{E} \left[\frac{P_{\mathbf{x}'}(x)}{P_{\mathbf{x}}(x)} \right] = 1 \end{aligned} \quad (14)$$

And the "=" holds when $\mathbf{x} = \phi_1$. So $\|\mathbf{B}\|_2 = 1$

- 6.4. (a) The empirical mean $\frac{1}{n} \sum_{i=1}^n \log Q_{y|x}(y_i|x_i)$ is the realization of expectation $\mathbb{E}_{P_{xy}}[\log Q_{y|x}(x, y)]$.

$$\begin{aligned} (\mathbf{g}^*, \mathbf{b}^*) &= \arg \max_{\mathbf{g}, \mathbf{b}} D(P_{xy} \| P_{\mathbf{x}} Q_{y|x}) \\ &= \arg \max_{\mathbf{g}, \mathbf{b}} \mathbb{E}_{P_{xy}} \left[\log \frac{P_{xy}(x, y)}{P_{\mathbf{x}}(x) Q_{y|x}(y|x)} \right] \\ &= \arg \min_{\mathbf{g}, \mathbf{b}} \mathbb{E}_{P_{xy}} [\log Q_{y|x}(y|x)] \\ &= \arg \min_{\mathbf{g}, \mathbf{b}} \frac{1}{n} \sum_{i=1}^n \log Q_{y|x}(y_i|x_i) \end{aligned} \quad (15)$$

- (b) First, we will prove the first order approximation of K-L divergence, i.e. if $P_1(x) - P_2(x) = O(\varepsilon)$,

$$\begin{aligned}
D(P_1(x) \| P_2(x)) &= - \sum_{x \in \mathcal{X}} P_1(x) \log \left(\frac{P_1(x) + P_2(x) - P_1(x)}{P_1(x)} \right) \\
&= - \sum_{x \in \mathcal{X}} P_1(x) \left(\frac{P_2(x) - P_1(x)}{P_1(x)} - \frac{1}{2} \frac{(P_2(x) - P_1(x))^2}{(P_1(x))^2} + O(\epsilon^3) \right) \\
&\approx \frac{1}{2} \sum_{x \in \mathcal{X}} \frac{(P_2(x) - P_1(x))^2}{P_1(x)} = \frac{1}{2} \sum_{x \in \mathcal{X}} \frac{(P_2(x) - P_1(x))^2}{P_2(x)} \frac{P_2(x)}{P_1(x)} \\
&= \frac{1}{2} \sum_{x \in \mathcal{X}} \frac{(P_2(x) - P_1(x))^2}{P_2(x)} \left(1 - \frac{O(\epsilon)}{P_1(x)} \right) \\
&\approx \frac{1}{2} \sum_{x \in \mathcal{X}} \frac{(P_2(x) - P_1(x))^2}{P_2(x)}
\end{aligned} \tag{16}$$

Then we just need to prove that $P_{xy}(x, y) - P_x(x)Q_{y|x}(y) = O(\epsilon)$.

Because $\mathbf{f}^T(x)\mathbf{g}^*(y) + d^*(y) = O(\epsilon)$,

$$P_y(y)e^{\mathbf{f}^T(x)\mathbf{g}^*(y) + d^*(y)} = P_y(y)(1 + \mathbf{f}^T(x)\mathbf{g}^*(y) + d^*(y) + O(\epsilon^2)) \tag{17}$$

$$\sum_{y \in \mathcal{Y}} P_y(y)(1 + \mathbf{f}^T(x)\mathbf{g}^*(y) + d^*(y) + O(\epsilon^2)) = 1 + \mathbb{E}_{P_y}[\mathbf{f}^T(x)\mathbf{g}^*(y)] + \mathbb{E}_{P_y}[d^*(y)] + O(\epsilon^2) \tag{18}$$

Without loss of generality, we can assume

$$\mathbb{E}_{P_y}[\mathbf{g}^*(y)] = \mathbb{E}_{P_y}[d^*(y)] = 0.$$

$$\begin{aligned}
Q_{y|x}(x, y) &= \frac{P_y(y)(1 + \mathbf{f}^T(x)\mathbf{g}^*(y) + d^*(y) + O(\epsilon^2))}{\sum_{y' \in \mathcal{Y}} P_y(y')(1 + \mathbf{f}^T(x)\mathbf{g}^*(y') + d^*(y') + O(\epsilon^2))} \\
&= \frac{P_y(y)(1 + \mathbf{f}^T(x)\mathbf{g}^*(y) + d^*(y) + O(\epsilon^2))}{1 + O(\epsilon^2)} \\
&= P_y(y)(1 + \mathbf{f}^T(x)\mathbf{g}^*(y) + d^*(y) + O(\epsilon^2))(1 - O(\epsilon^2)) \\
&= P_y(y)(1 + \mathbf{f}^T(x)\mathbf{g}^*(y) + d^*(y)) + O(\epsilon^2)
\end{aligned} \tag{19}$$

Thus,

$$P_{xy}(x, y) - P_x(x)Q_{y|x}(y) = P_x(x)P_y(y) + O(\epsilon) - P_x(x)P_y(y)(1 + O(\epsilon)) = O(\epsilon) \tag{20}$$

So we have,

$$D(P_{xy}(x, y) \| P_x(x)Q_{y|x}(y)) \approx \frac{1}{2} \sum_{x \in \mathcal{X}, y \in \mathcal{Y}} \frac{(P_{xy}(x, y) - P_x(x)Q_{y|x}(y))^2}{P_x(x)Q_{y|x}(y)} \tag{21}$$