

Problem Set 6

Issued: Monday 7th December, 2020

Due: Sunday 20th December, 2020

Notations: We use $\mathbf{x}, \mathbf{y}, \mathbf{w}$ and $\underline{\mathbf{x}}, \underline{\mathbf{y}}, \underline{\mathbf{w}}$ to denote random variables and random vectors.

6.1. Suppose that y_1 and y_2 are independent random variables each uniformly distributed between x and $x + 1$. Let $\mathbf{s} = \max(y_1, y_2)$ and $\mathbf{r} = y_1 - y_2$.

- (a) Show that \mathbf{s} is not a sufficient statistic for $p_{y_1 y_2}(y_1, y_2; x)$.
- (b) An ancillary statistic is one whose distribution does not depend on the parameters of the model. Show that \mathbf{r} is an ancillary statistic for $p_{y_1 y_2}(y_1, y_2; x)$.

- (c) Is $\mathbf{u} = \begin{bmatrix} \mathbf{s} \\ \mathbf{r} \end{bmatrix}$ a sufficient statistic for $p_{y_1 y_2}(y_1, y_2; x)$?

6.2. Please verify that Pearson correlation coefficient does not satisfy property (5), i.e., for any one-to-one mapping $\xi(\cdot)$ and $\eta(\cdot)$, $\rho(\xi(\mathbf{x}), \eta(\mathbf{y})) = \rho(\mathbf{x}, \mathbf{y})$.

6.3. Given two random variables $\mathbf{x} \in \mathcal{X}, \mathbf{y} \in \mathcal{Y}$ with the joint distribution $P_{\mathbf{xy}}(x, y)$, the corresponding matrix $\mathbf{B} \in \mathbb{R}^{|\mathcal{Y}| \times |\mathcal{X}|}$ is defined as

$$\mathbf{B}(y, x) \triangleq \frac{P_{\mathbf{xy}}(x, y)}{\sqrt{P_{\mathbf{x}}(x)} \sqrt{P_{\mathbf{y}}(y)}}.$$

In the derivation of HGR maximal correlation analysis, given feature functions $f: \mathcal{X} \rightarrow \mathbb{R}, g: \mathcal{Y} \rightarrow \mathbb{R}$, we defined the corresponding *information vectors* as the vectors $\boldsymbol{\phi} \in \mathbb{R}^{|\mathcal{X}|}, \boldsymbol{\psi} \in \mathbb{R}^{|\mathcal{Y}|}$ with elements $\boldsymbol{\phi}(x) = f(x)\sqrt{P_{\mathbf{x}}(x)}, \boldsymbol{\psi}(y) = g(y)\sqrt{P_{\mathbf{y}}(y)}$. We denote them as $\boldsymbol{\phi} \leftrightarrow f(\mathbf{x})$ and $\boldsymbol{\psi} \leftrightarrow g(\mathbf{y})$.

- (a) Show that
 - i. $\mathbb{E}[f(\mathbf{x})g(\mathbf{y})] = \boldsymbol{\psi}^T \mathbf{B} \boldsymbol{\phi}$.
 - ii. $\mathbf{B} \boldsymbol{\phi} \leftrightarrow \mathbb{E}[f(\mathbf{x})|\mathbf{y}]$.
 - iii. $\mathbf{B}^T \boldsymbol{\psi} \leftrightarrow \mathbb{E}[g(\mathbf{y})|\mathbf{x}]$.
- (b) Suppose $\boldsymbol{\phi}_1 = [\sqrt{P_{\mathbf{x}}(1)}, \dots, \sqrt{P_{\mathbf{x}}(|\mathcal{X}|)}]^T, \boldsymbol{\psi}_1 = [\sqrt{P_{\mathbf{y}}(1)}, \dots, \sqrt{P_{\mathbf{y}}(|\mathcal{Y}|)}]^T$. Show that $\mathbf{B} \boldsymbol{\phi}_1 = \boldsymbol{\psi}_1, \mathbf{B}^T \boldsymbol{\psi}_1 = \boldsymbol{\phi}_1$, and interpret their meanings from the perspective of conditional expectation.
- (c) Prove that $\|\mathbf{B}\|_2 = 1$. The definition of L- p matrix norm is $\|\mathbf{A}\|_p \triangleq \sup_{\mathbf{x} \neq \mathbf{0}} \frac{\|\mathbf{Ax}\|_p}{\|\mathbf{x}\|_p}$.

6.4. We are now using softmax regression under a discriminative model of the form

$$Q_{\mathbf{y}|\mathbf{x}}(y|x) = \frac{\exp(\mathbf{f}(x)^T \mathbf{g}(y) + b(y))}{\sum_{y' \in \mathcal{Y}} \exp(\mathbf{f}(x)^T \mathbf{g}(y') + b(y'))},$$

to address the classification problems. Note that here \mathbf{x} and \mathbf{y} is the sample and label with alphabet \mathcal{X} and \mathcal{Y} respectively. Now we assume that feature $\mathbf{f}(\cdot)$ is decided. We hope to derive $\mathbf{g}(\cdot)$ and $b(\cdot)$ by

$$(\mathbf{g}^*, b^*) = \arg \min_{(\mathbf{g}, b)} D(P_{\mathbf{xy}} \| P_{\mathbf{x}} Q_{\mathbf{y}|\mathbf{x}}),$$

where $P_{\mathbf{xy}}$ is the distribution corresponding to the dataset.

- (a) Explain when we have data $\{(x_i, y_i)\}_{i=1}^n$, the softmax regression optimal solution (\mathbf{g}^*, b^*) is minimizing the empirical mean

$$(\mathbf{g}^*, b^*) = \arg \max_{(\mathbf{g}, b)} \frac{1}{n} \sum_{i=1}^n \log Q_{y|x}(y_i|x_i)$$

- (b) We assume a weak dependency here, which means the true distribution P_{xy} satisfies $P_{xy}(x, y) - P_x(x)P_y(y) = O(\epsilon)$. We define $d^*(y) \triangleq b^*(y) - \log P_y(y)$. We can derive an equivalent expression of the weak dependency, which also needs some mathematical tricks to prove, $\mathbf{f}^T(x)\mathbf{g}^*(y) + d^*(y) = O(\epsilon)$. Please prove that the first order approximation of $D(P_{xy} \| P_x Q_{y|x})$ is $\sum_{x \in \mathcal{X}, y \in \mathcal{Y}} \frac{(P_{xy}(x, y) - P_x(x)Q_{y|x}(y|x))^2}{P_x(x)Q_{y|x}(y|x)}$.