

**Written Assignment 4**

**Issued:** Saturday 30<sup>th</sup> November, 2019

**Due:** Saturday 14<sup>th</sup> December, 2019

---

Questions 4.1 & 4.2 will help you understand some basic properties about maximal correlation and the alternating conditional expectation (ACE) algorithm.

4.1. (2 points) Conditional Expectation A key step in the ACE algorithm is computing  $g(y)$  as the conditional expectation:

$$g(y) \triangleq \mathbb{E}[X|Y = y] = \sum_{x \in \mathcal{X}} P_{X|Y}(x|y)x, \forall y \in \mathcal{Y}.$$

Then  $\mathbb{E}[X|Y] = g(Y)$  is also a random variable. Geometrically,  $g(Y)$  is the projection of  $X$  onto the function space of  $Y$ . It is common sense that the length of projection is less than or equal to the length of the raw vector. Prove that

$$\mathbb{E}[g^2(Y)] \leq \mathbb{E}[X^2].$$

4.2. (3 points) Information Vectors In the derivation of HGR maximal correlation analysis, given a feature function  $f: \mathcal{X} \rightarrow \mathbb{R}$ , we defined the corresponding *information vector* as the vector  $\phi \in \mathbb{R}^{|\mathcal{X}|}$  with elements  $\phi(x) = f(x)\sqrt{P_X(x)}$ . This correspondence between function  $f$  and information vector  $\phi$  is denoted by  $\phi \leftrightarrow f(X)$ . Show that

- (a)  $\phi_1 \leftrightarrow 1(X)$ , where  $\phi_1 = \left(\sqrt{P_X(1)}, \dots, \sqrt{P_X(|\mathcal{X}|)}\right)^T$ , and  $1(x)$  is a constant function, i.e.  $1(x) = 1$  for all  $x \in \mathcal{X}$ .
- (b) The variance of a feature is the length of its corresponding information vector:  $\mathbb{E}[f^2(X)] = \|\phi\|^2$ , where  $\phi \leftrightarrow f(X)$ .
- (c) The covariance of two features is the inner product of their information vectors:  $\langle \phi_1, \phi_2 \rangle = \mathbb{E}[f_1(X)f_2(X)]$ , where  $\phi_1 \leftrightarrow f_1(X)$ ,  $\phi_2 \leftrightarrow f_2(X)$ .

4.3. (2 points) ICA It's well known that Gaussian variables are forbidden in ICA. To understand this limitation, let's assume that the joint distribution of two independent components, say,  $s_1, s_2$ , are Gaussian.

$$P(\mathbf{s}_i) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{s_i^2}{2}\right)$$

- (a) Please find the joint pdf  $P(s_1, s_2)$ .
- (b) Suppose that the mixing matrix  $\mathbf{A}$  is orthogonal. For example, we could assume that this is so because the data has been whitened, which means  $\mathbf{A}^{-1} = \mathbf{A}^T$  holds. Please find the joint pdf  $P(x_1, x_2)$  of the mixtures  $x_1$  and  $x_2$  and then explain why Gaussian variables are forbidden.

- 4.4. (3 points) *EM for Mixture of Gaussians (Soft k-Means)* We talked about EM for Mixture of Gaussians in class. Please repeat what have been done in this problem. Consider the case of a mixture of  $k$  Gaussians  $\theta$  is a triplet  $(\phi, \{\mu_1, \dots, \mu_k\}, \{\Sigma_1, \dots, \Sigma_k\})$ . For simplicity, we assume that  $\Sigma_1 = \dots = \Sigma_k = \mathbf{I}$ , which don't need calculations in your EM steps. We have that

$$P_{\theta^{(t)}}(Z = z) = \phi_z^{(t)}$$

$$P_{\theta^{(t)}}(\mathbf{X} = \mathbf{x} | Z = z) = \frac{1}{(2\pi)^{\frac{d}{2}} |\Sigma_z|^{\frac{1}{2}}} \exp \left( -\frac{1}{2} (\mathbf{x} - \mu_z^{(t)})^T \Sigma_z^{-1} (\mathbf{x} - \mu_z^{(t)}) \right)$$

- (a) Please write down the E-step and M-step. *Hint: The E-step need to write down  $P_{\theta^{(t)}}(Z = z | \mathbf{X} = \mathbf{x}_i)$*
- (b) Write down the updated parameter  $\theta^{(t+1)}$  and compare your procedures with K-means.

- 4.5. (3 points) (Extra credit) *Weyl's Theorem* This problem introduces you to the perturbation theory in PCA. Perturbation theory is useful in many real world problems, for instance, suppose we have computed the largest eigenvalue of the covariance matrix of some original samples. Then suddenly a bunch of new data come in and the covariance matrix should be like

$$\Sigma = \frac{n_{origin} \Sigma_{origin} + n_{new} \Sigma_{new}}{n_{origin} + n_{new}}$$

Let's note it as

$$\Sigma = \mathbf{A} + \mathbf{B}$$

Define  $\lambda(M)$  as the eigenvalue operator of matrix  $M$ . Our target is to bound the eigenvalues  $\lambda(\Sigma)$  given some knowledge about  $\lambda(\mathbf{A})$ .

Let  $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{n \times n}$  and their eigenvalues denoted by  $\{\lambda_i(\mathbf{A})\}_{i=1}^n, \{\lambda_i(\mathbf{B})\}_{i=1}^n$  with  $\lambda_1 > \dots > \lambda_n$ . Please prove that for any  $1 \leq k \leq n$

$$\lambda_k(\mathbf{A}) + \lambda_n(\mathbf{B}) \leq \lambda_k(\mathbf{A} + \mathbf{B}) \leq \lambda_k(\mathbf{A}) + \lambda_1(\mathbf{B})$$

*Hint: you should first prove that for any  $\mathbf{v} \in \mathbb{R}^n$*

$$\lambda_n(\mathbf{B}) \leq \frac{\mathbf{v}^T \mathbf{B} \mathbf{v}}{\|\mathbf{v}\|^2} \leq \lambda_1(\mathbf{B})$$