

**Problem Set 2**

**Issued:** Monday 22<sup>nd</sup> October, 2018

**Due:** Monday 29<sup>th</sup> October, 2018

---

- 2.1. A data set consists of  $m$  data pairs  $(\mathbf{x}^{(1)}, y^{(1)}), \dots, (\mathbf{x}^{(m)}, y^{(m)})$ , where  $\mathbf{x} \in \mathbb{R}^n$  is the independent variable, and  $y \in \mathbb{R}$  is the dependent variable. Denote the design matrix by  $\mathbf{X} \triangleq [\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}]^\top$ , and let  $\mathbf{y} \triangleq [y^{(1)}, \dots, y^{(m)}]^\top$ . The least-squares method then minimizes the square loss  $J(\boldsymbol{\theta})$  defined as

$$J(\boldsymbol{\theta}) = \frac{1}{2} \|\mathbf{y} - \mathbf{X}\boldsymbol{\theta}\|_2^2,$$

where  $\boldsymbol{\theta} \in \mathbb{R}^n$  is the parameter to be estimated. To find the optimal  $\boldsymbol{\theta}$ , let  $\nabla J(\boldsymbol{\theta}) = 0$ , and we can get the normal equation:

$$\mathbf{X}^\top \mathbf{X} \boldsymbol{\theta} = \mathbf{X}^\top \mathbf{y}. \quad (1)$$

When  $\mathbf{X}^\top \mathbf{X}$  is invertible, we have  $\boldsymbol{\theta} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$ .

Now, suppose  $\mathbf{X}^\top \mathbf{X}$  is singular. Does the solution of (1) still exist? Prove your result, and explain its meaning in plain words.

- 2.2. A data set consists of  $m$  data pairs  $(\mathbf{x}^{(1)}, y^{(1)}), \dots, (\mathbf{x}^{(m)}, y^{(m)})$ , where  $\mathbf{x} \in \mathbb{R}^n$  is the independent variable, and  $y \in \{1, \dots, k\}$  is the dependent variable. The conditional probability  $P_{y|\mathbf{x}}(y|\mathbf{x})$ <sup>1</sup> estimated by the softmax regression is

$$P_{y|\mathbf{x}}(l|\mathbf{x}) = \frac{\exp(\boldsymbol{\theta}_l^\top \mathbf{x} + b_l)}{\sum_{j=1}^k \exp(\boldsymbol{\theta}_j^\top \mathbf{x} + b_j)}, \quad l = 1, \dots, k,$$

where  $\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_k \in \mathbb{R}^n$  and  $b_1, \dots, b_k \in \mathbb{R}$  are the parameters of softmax regression. The term  $b_i$  is called a bias term.

The log-likelihood of the softmax regression model is

$$\ell = \sum_{i=1}^m \log P_{y|\mathbf{x}}(y^{(i)}|\mathbf{x}^{(i)}).$$

- (a) Evaluate  $\nabla_{b_l} \ell$ .

The data set can be described by its empirical distribution  $\hat{P}_{\mathbf{x},y}(\mathbf{x}, y)$  defined as

$$\hat{P}_{\mathbf{x},y}(\mathbf{x}, y) = \frac{1}{m} \sum_{i=1}^m \mathbb{1}\{\mathbf{x} = \mathbf{x}^{(i)}, y = y^{(i)}\},$$

where  $\mathbb{1}\{\cdot\}$  is the indicator function. Similarly, the empirical marginal distributions of this data set are

$$\hat{P}_{\mathbf{x}}(\mathbf{x}) = \frac{1}{m} \sum_{i=1}^m \mathbb{1}\{\mathbf{x} = \mathbf{x}^{(i)}\}, \quad \hat{P}_y(y) = \frac{1}{m} \sum_{i=1}^m \mathbb{1}\{y = y^{(i)}\}.$$

---

<sup>1</sup>The notation  $P_{y|\mathbf{x}}(y|\mathbf{x})$  stands for  $\mathbb{P}(y = y|\mathbf{x} = \mathbf{x})$ , i.e., the conditional probability of  $y = y$  given  $\mathbf{x} = \mathbf{x}$ .

(b) Suppose we have set the biases  $(b_1, \dots, b_k)$  to their optimal values, prove that

$$\hat{P}_{\mathbf{y}}(l) = \sum_{\mathbf{x} \in \mathcal{X}} P_{\mathbf{y}|\mathbf{x}}(l|\mathbf{x}) \hat{P}_{\mathbf{x}}(\mathbf{x}),$$

where  $\mathcal{X} = \{\mathbf{x}^{(i)} : i = 1, \dots, m\}$  is the set of all samples of  $\mathbf{x}$ .

*Hint: The optimality implies  $\nabla_{b_1} \ell = \nabla_{b_2} \ell = \dots = \nabla_{b_k} \ell = 0$ .*

2.3. The multivariate normal distribution can be written as

$$p_{\mathbf{y}}(\mathbf{y}; \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{\frac{n}{2}} |\boldsymbol{\Sigma}|^{\frac{1}{2}}} \exp \left( -\frac{1}{2} (\mathbf{y} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{y} - \boldsymbol{\mu}) \right),$$

where  $\boldsymbol{\mu}$  and  $\boldsymbol{\Sigma}$  are the parameters. Show that the family of multivariate normal distributions is an exponential family, and find the corresponding  $\eta$ ,  $b(\boldsymbol{\eta})$ ,  $T(\mathbf{y})$ , and  $a(\boldsymbol{\eta})$ .

*Hints: The parameters  $\eta$  and  $T(\mathbf{y})$  are not limited to be vectors, but can also be matrices. In this case, the [Frobenius inner product](#) can be used to define the inner product between two matrices, which is represented as the trace of their products. The properties of [matrix trace](#) might be useful.*