

Homework 1

TIAN Chenyu

October 8, 2019

-
- **Acknowledgments:** This template takes some materials from course CSE 547/Stat 548 of Washington University:
<https://courses.cs.washington.edu/courses/cse547/17sp/index.html>.
 - **Collaborators:** I finish this template by myself.
-

1.1. The log-likelihood of the softmax regression model can be written as

$$\begin{aligned} l &= \sum_{i=1}^m \log \frac{\exp(\boldsymbol{\theta}_{y^{(i)}}^T \mathbf{x}^{(i)} + b_{y^{(i)}})}{\sum_{j=1}^k \exp(\boldsymbol{\theta}_j^T \mathbf{x}^{(i)} + b_j)} \\ &= \sum_{i=1}^m [\boldsymbol{\theta}_{y^{(i)}}^T \mathbf{x}^{(i)} + b_{y^{(i)}} - \log(\sum_{j=1}^k \exp(\boldsymbol{\theta}_j^T \mathbf{x}^{(i)} + b_j))] \end{aligned}$$

(a) Evaluate the derivation of b_l :

$$\frac{\partial l}{\partial b_l} = \sum_i^m [\mathbb{1}(y^{(i)} = l) - \frac{\exp(\boldsymbol{\theta}_l^T \mathbf{x}^{(i)} + b_l)}{\sum_{j=1}^k \exp(\boldsymbol{\theta}_j^T \mathbf{x}^{(i)} + b_j)}]$$

The $\mathbb{1}(y^{(i)} = l)$ function is defined as:

$$\mathbb{1}(x = l) = \begin{cases} 1, & \text{if } x = l, \\ 0, & \text{if } x \neq l. \end{cases}$$

(b) If we have set the biases to their optimal values, there exists $\frac{\partial l}{\partial b_l} = 0$. Based on (a):

$$\sum_i^m \mathbb{1}(y^{(i)} = l) = \frac{\exp(\boldsymbol{\theta}_l^T \mathbf{x}^{(i)} + b_l)}{\sum_{j=1}^k \exp(\boldsymbol{\theta}_j^T \mathbf{x}^{(i)} + b_j)}$$

Based on the definition of $\hat{P}_y(l)$:

$$\begin{aligned} \hat{P}_y(l) &= \frac{1}{m} \sum_i^m \frac{\exp(\boldsymbol{\theta}_l^T \mathbf{x}^{(i)} + b_l)}{\sum_{j=1}^k \exp(\boldsymbol{\theta}_j^T \mathbf{x}^{(i)} + b_j)} \\ &= \frac{1}{m} \sum_i^m \sum_{\mathbf{x} \in \mathbf{X}} \frac{\exp(\boldsymbol{\theta}_l^T \mathbf{x} + b_l)}{\sum_{j=1}^k \exp(\boldsymbol{\theta}_j^T \mathbf{x}^{(i)} + b_j)} \mathbb{1}(\mathbf{x}^{(i)} = \mathbf{x}) \\ &= \frac{1}{m} \sum_i^m \sum_{\mathbf{x} \in \mathbf{X}} P_{(y|\mathbf{x})}(l | \mathbf{x}) \mathbb{1}(\mathbf{x}^{(i)} = \mathbf{x}) \\ &= \sum_{\mathbf{x} \in \mathbf{X}} P_{(y|\mathbf{x})}(l | \mathbf{x}) \frac{1}{m} \sum_i^m \mathbb{1}(\mathbf{x}^{(i)} = \mathbf{x}) \\ &= \sum_{\mathbf{x} \in \mathbf{X}} P_{(y|\mathbf{x})} \hat{P}_{\mathbf{x}}(\mathbf{x}) \end{aligned}$$

1.2. (a) The MLE is to solve

$$\underset{\mu}{\text{maximize}} \prod_i^n P(x_i | \mu)$$

which is equal to:

$$\begin{aligned} & \underset{\mu}{\text{maximize}} \sum_i^n \log P(x_i | \mu) \\ \Rightarrow & \underset{\mu}{\text{maximize}} \sum_i^n \left(\log \frac{1}{\sqrt{2\pi}\delta^2} - \frac{(x_i - \mu)^2}{2\delta^2} \right) \\ \Rightarrow & \underset{\mu}{\text{minimize}} \sum_i^n (x_i - \mu)^2 \end{aligned}$$

Here we define $f(\mu) = \sum_i^n (x_i - \mu)^2$ to find the μ^* . There exists $\frac{\partial f}{\partial \mu} = 2 \sum_i^n (\mu - x_i) = 0$ when $\mu = \mu^*$.

Then $\mu^* = \frac{\sum_i^n x_i}{n}$.

(b) The MAP problem can be written as:

$$\underset{\mu}{\text{maximize}} P(\mu | x_1, \dots, x_n)$$

which is equal to:

$$\begin{aligned} & \underset{\mu}{\text{maximize}} P(\mu | x_1, \dots, x_n) \\ \Rightarrow & \underset{\mu}{\text{maximize}} P(x_1, \dots, x_n | \mu) P(\mu) \\ \Rightarrow & \underset{\mu}{\text{maximize}} \prod_i^n P(x_i | \mu) P(\mu) \\ \Rightarrow & \underset{\mu}{\text{maximize}} \sum_i^n \log P(x_i | \mu) + \log P(\mu) \\ \Rightarrow & \underset{\mu}{\text{maximize}} \sum_i^n -\frac{(x_i - \mu)^2}{2\delta^2} - \log \sqrt{2\pi\theta^2} - \frac{(\mu - \nu)^2}{2\theta^2} \\ \Rightarrow & \underset{\mu}{\text{minimize}} \sum_i^n \frac{(x_i - \mu)^2}{2\delta^2} + \frac{(\mu - \nu)^2}{2\theta^2} \end{aligned}$$

Here we define $g(\mu) = \sum_i^n \frac{(x_i - \mu)^2}{2\delta^2} + \frac{(\mu - \nu)^2}{2\theta^2}$ to find the μ^* . There exists $\frac{\partial g}{\partial \mu} = \sum_i^n \frac{\mu(\mu - x_i)}{\delta^2} + \frac{\mu(\mu - \nu)}{\theta^2} = 0$ when $\mu = \mu^*$.

Then $\mu^* = \frac{\sum_i^n x_i + \delta^2 \nu}{\theta^2 n + \delta^2}$.

When $n \rightarrow \infty$, μ^* for MLE and MAP is equal.

1.3. First, the square error can be written as

$$\begin{aligned} J(\Theta) &= \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^l ((\Theta^T \mathbf{x}^{(i)})_j - \mathbf{y}_j^{(i)})^2 \\ &= \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^l \left(\sum_{k=1}^n \Theta_{kj} \mathbf{x}_k^{(i)} - \mathbf{y}_j^{(i)} \right)^2 \end{aligned}$$

In order to compute the solution, it needs to find the minimum $J(\Theta)$ where $\frac{\partial J}{\partial \Theta} = 0$. To find the solution, the derivative is:

$$\begin{aligned} \frac{\partial J}{\partial \Theta_{\alpha\beta}} &= \sum_{i=1}^m (\mathbf{x}_\alpha^{(i)} (\sum_{k=1}^n \Theta_{k\beta} \mathbf{x}_k^{(i)} - \mathbf{y}_\beta^{(i)})) \\ &= \sum_{i=1}^m (\mathbf{x}_\alpha^{(i)} (\Theta_\beta^T \mathbf{x}^{(i)} - \mathbf{y}_\beta^{(i)})) \end{aligned}$$

Then we have $\frac{\partial J}{\partial \Theta} = 0$, which means:

$$\begin{aligned} \sum_{i=1}^m \mathbf{x}_\alpha^{(i)} (\Theta_\beta^T \mathbf{x}^{(i)}) &= \sum_{i=1}^m \mathbf{x}_\alpha^{(i)} \mathbf{y}_\beta^{(i)} \\ \Rightarrow \mathbf{X}_\alpha^T \mathbf{X} \Theta_\beta &= \mathbf{X}_\alpha^T \mathbf{Y}_\beta \\ \Rightarrow \mathbf{X}^T \mathbf{X} \Theta_\beta &= \mathbf{X}^T \mathbf{Y}_\beta \end{aligned}$$

So, for $\beta \in (1, 2, \dots, l)$ we have

$$\Theta_\beta = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}_\beta$$

Also,

$$\Theta = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

1.4. Based on Σ is symmetrical as well as the properties of matrix trace, the multivariate normal distribution can be written as

$$\begin{aligned} P_{\mathbf{y}}(\mathbf{y}; \boldsymbol{\mu}, \boldsymbol{\Sigma}) &= \frac{1}{(2\pi)^{\frac{n}{2}} |\boldsymbol{\Sigma}|^{\frac{1}{2}}} \exp\left(-\frac{1}{2}(\mathbf{y} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{y} - \boldsymbol{\mu})\right) \\ &= (2\pi)^{-\frac{n}{2}} |\boldsymbol{\Sigma}|^{-\frac{1}{2}} \exp\left(-\frac{1}{2}(\mathbf{y}^T \boldsymbol{\Sigma}^{-1} - \boldsymbol{\mu}^T \boldsymbol{\Sigma}^{-1})(\mathbf{y} - \boldsymbol{\mu})\right) \\ &= (2\pi)^{-\frac{n}{2}} |\boldsymbol{\Sigma}|^{-\frac{1}{2}} \exp\left(-\frac{1}{2}(\mathbf{y}^T \boldsymbol{\Sigma}^{-1} \mathbf{y} - \boldsymbol{\mu}^T \boldsymbol{\Sigma}^{-1} \mathbf{y} - \mathbf{y}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu} + \boldsymbol{\mu}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu})\right) \\ &= (2\pi)^{-\frac{n}{2}} |\boldsymbol{\Sigma}|^{-\frac{1}{2}} \exp\left(-\frac{1}{2}tr(\mathbf{y}^T \boldsymbol{\Sigma}^{-1} \mathbf{y} - \boldsymbol{\mu}^T \boldsymbol{\Sigma}^{-1} \mathbf{y} - \mathbf{y}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu} + \boldsymbol{\mu}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu})\right) \\ &= (2\pi)^{-\frac{n}{2}} |\boldsymbol{\Sigma}|^{-\frac{1}{2}} \exp\left(-\frac{1}{2}(tr(\mathbf{y}^T \boldsymbol{\Sigma}^{-1} \mathbf{y}) - tr(\boldsymbol{\mu}^T \boldsymbol{\Sigma}^{-1} \mathbf{y}) - tr(\mathbf{y}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}) + tr(\boldsymbol{\mu}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}))\right) \\ &= (2\pi)^{-\frac{n}{2}} |\boldsymbol{\Sigma}|^{-\frac{1}{2}} \exp\left(tr(\boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}) \mathbf{y} - \frac{1}{2}tr(\boldsymbol{\Sigma}^{-1} \mathbf{y} \mathbf{y}^T) - \frac{1}{2}\boldsymbol{\mu}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}\right) \end{aligned}$$

Thus, we can see that multivariate normal distribution is an exponential

family with:

$$\eta = \begin{pmatrix} \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu} \\ -\frac{1}{2} \boldsymbol{\Sigma}^{-1} \end{pmatrix}$$

$$b(\mathbf{y}) = (2\pi)^{-\frac{n}{2}} |\boldsymbol{\Sigma}|^{-\frac{1}{2}}$$

$$T(\mathbf{y}) = \begin{pmatrix} \mathbf{y} \\ \mathbf{y} \mathbf{y}^T \end{pmatrix}$$

$$a(\eta) = \frac{1}{2} \boldsymbol{\mu}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}$$