

**Homework 2**

Zi-Feng WANG

October 28, 2018

- 
- **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 my homework all by myself.
- 

2.1. **Solution:** let  $\nabla J(\theta) = 0$  and we can get the normal equation

$$X^T X \theta = X^T \mathbf{y} \quad (1)$$

from  $x \in \mathbb{R}^n$  we can get that the  $X$  is a  $m \times n$  matrix, hence  $X^T X$  is a  $n \times n$  matrix,  $\theta$  is  $n \times 1$  vector and  $X^T \mathbf{y}$  is  $n \times 1$  vector as well. the (1) can be represented as the following format

$$A \theta = \mathbf{b} \quad (2)$$

in (2), the  $A = X^T X$ ,  $\mathbf{b} = X^T \mathbf{y}$ . if  $A$  is a singular and square matrix, the  $\text{rank}(A) < n$ . there are two possible cases:

- (a) when  $\text{rank}(A) = \text{rank}(A, \mathbf{b}) < n$ , there are infinite number of solutions of  $\theta$
- (b) when  $\text{rank}(A) < \text{rank}(A, \mathbf{b}) < n$ , there is none solution of  $\theta$

2.2. **Solution:**

- (a) before getting  $\nabla_{b_l} \ell$ , unfold  $\ell$  at first

$$\begin{aligned} \ell &= \sum_{i=1}^m \log P_{\mathbf{y}|\mathbf{x}}(y^{(i)}|x^{(i)}) \\ &= \sum_{i=1}^m \log \frac{\exp(\theta_l^T x^{(i)} + b_l)}{\sum_{j=1}^k \exp(\theta_j^T x^{(i)} + b_j)} \\ &= \sum_{i=1}^m (\theta_l^T x^{(i)} + b_l) - \sum_{i=1}^m \log \sum_{j=1}^k \exp(\theta_j^T x^{(i)} + b_j) \end{aligned} \quad (3)$$

from (3), the  $\nabla_{b_l} \ell$  is

$$\begin{aligned}
\nabla_{b_l} \ell &= \nabla_{b_l} \left\{ \sum_{i=1}^m (\theta_l^T x^{(i)} + b_l) - \sum_{i=1}^m \log \sum_{j=1}^k \exp(\theta_j^T x^{(i)} + b_j) \right\} \\
&= \sum_{i=1}^m \mathbb{1}\{y^{(i)} = l\} - \sum_{i=1}^m \frac{1}{\sum_{j=1}^k \exp(\theta_j^T x^{(i)} + b_j)} \frac{\partial (\sum_{j=1}^k \exp(\theta_j^T x^{(i)} + b_j))}{\partial b_l} \\
&= \sum_{i=1}^m \mathbb{1}\{y^{(i)} = l\} - \sum_{i=1}^m \frac{\theta_l^T x^{(i)} + b_l}{\sum_{j=1}^k \exp(\theta_j^T x^{(i)} + b_j)} \\
&= \sum_{i=1}^m (\mathbb{1}\{y^{(i)} = l\} - P(y^{(i)} = l | x^{(i)}; \theta, b))
\end{aligned} \tag{4}$$

(b) *Proof.* when  $\nabla_{b_l} \ell = 0$ , the (4) = 0, that is  
 $\sum_{i=1}^m (\mathbb{1}\{y^{(i)} = l\} - P(y^{(i)} = l | x^{(i)}; \theta, b)) = 0$ .  
so that

$$\begin{aligned}
\sum_{i=1}^m \mathbb{1}\{y^{(i)} = l\} &= \sum_{i=1}^m P(y^{(i)} = l | x^{(i)}; \theta, b) \\
m \frac{\sum_{i=1}^m \mathbb{1}\{y^{(i)} = l\}}{m} &= \sum_{i=1}^m \sum_{j=1}^m P(y^{(i)} = l | x^{(i)}; \theta, b) \mathbb{1}\{x^{(i)} = x^{(j)}\} \\
\hat{P}_y(l) &= \frac{1}{m} \sum_{i=1}^m P(y^{(i)} = l | x^{(i)}; \theta, b) \sum_{j=1}^m \mathbb{1}\{x^{(i)} = x^{(j)}\} \\
\hat{P}_y(l) &= \sum_{x \in \mathcal{X}} P_{y|x}(l|x) \hat{P}_x(x)
\end{aligned}$$

□

2.3. the multivariate normal distribution can be written as

$$\begin{aligned}
p_y(y; \mu, \Sigma) &= \frac{1}{(2\pi)^{\frac{n}{2}} |\Sigma|^{\frac{1}{2}}} \exp\left(-\frac{1}{2}(y - \mu)^T \Sigma^{-1} (y - \mu)\right) \\
&= \frac{1}{(2\pi)^{\frac{n}{2}}} \exp\left\{-\frac{1}{2} \log |\Sigma| - \frac{1}{2}(y - \mu)^T \Sigma^{-1} (y - \mu)\right\}
\end{aligned} \tag{5}$$

the part in  $\exp(\cdot)$  of (5) is

$$\begin{aligned}
&-\frac{1}{2} \log |\Sigma| - \frac{1}{2}(y - \mu)^T \Sigma^{-1} (y - \mu) \\
&= -\frac{1}{2} \log |\Sigma| - \frac{1}{2}(y^T \Sigma^{-1} y - y^T \Sigma^{-1} \mu - \mu^T \Sigma^{-1} y + \mu^T \Sigma^{-1} \mu) \\
&= -\frac{1}{2}(tr(y^T \Sigma^{-1} y) - tr(y^T \Sigma^{-1} \mu) - tr(\mu^T \Sigma^{-1} y)) - \frac{1}{2}(\log |\Sigma| + \mu^T \Sigma^{-1} \mu) \\
&= -\frac{1}{2}(tr(\Sigma^{-1} y y^T) - tr(\Sigma^{-1} \mu y^T) - tr(\mu^T \Sigma^{-1} y)) - \frac{1}{2}(\log |\Sigma| + \mu^T \Sigma^{-1} \mu) \\
&= -\frac{1}{2}tr(\Sigma^{-1} y y^T - \Sigma^{-1} \mu y^T - \mu^T \Sigma^{-1} y) - \frac{1}{2}(\log |\Sigma| + \mu^T \Sigma^{-1} \mu)
\end{aligned} \tag{6}$$

let

$$\begin{aligned}\boldsymbol{\eta} &= -\frac{1}{2}(\Sigma^{-1}, \Sigma^{-1}\mu, \mu^T \Sigma^{-1})^T \\ T(y) &= (yy^T, y^T, y)^T\end{aligned}\tag{7}$$

therefore the (6) can be written as

$$tr(\boldsymbol{\eta}^T T(y)) - \frac{1}{2}(\log |\Sigma| + \mu^T \Sigma^{-1} \mu) = \langle \boldsymbol{\eta}, T(y) \rangle_F - a(\boldsymbol{\eta})\tag{8}$$

combine the (5), (6), (7) and (8), we can show that the multivariate normal distribution is an exponential family

$$P_y(y; \boldsymbol{\eta}) = b(y) \exp(\langle \boldsymbol{\eta}, T(y) \rangle_F - a(\boldsymbol{\eta}))\tag{9}$$

where

$$\begin{aligned}b(y) &= \frac{1}{(2\pi)^2} \\ \boldsymbol{\eta} &= -\frac{1}{2}(\Sigma^{-1}, \Sigma^{-1}\mu, \mu^T \Sigma^{-1})^T \\ T(y) &= (yy^T, y^T, y)^T \\ a(\boldsymbol{\eta}) &= \frac{1}{2}(\log |\Sigma| + \mu^T \Sigma^{-1} \mu)\end{aligned}$$

---

<sup>1</sup>the  $\langle A, B \rangle_F$  is the *Frobenius inner product* used to define the inner product between two matrices  $A$  and  $B$ , which is represented as the trace of their products i.e.  $tr(A^T B)$ .