

# DS 5230 Unsupervised Machine Learning and Data Mining

## Homework 2 By Yanchi Li

### Exercise 1

- a. We directly derive the target equation from the closed form solution of the weights  $\mathbf{w}_*$  that minimize the loss

$$\begin{aligned}\mathbf{w}_* &= (\mathbf{X}^T \mathbf{X} + \lambda I)^{-1} \mathbf{X}^T \mathbf{y} \\ (\mathbf{X}^T \mathbf{X} + \lambda I) \mathbf{w}_* &= \mathbf{X}^T \mathbf{y} \\ \mathbf{X}^T \mathbf{X} \mathbf{w}_* + \lambda \mathbf{w}_* &= \mathbf{X}^T \mathbf{y} \\ \mathbf{X}^T \mathbf{X} \mathbf{w}_* + \lambda \mathbf{w}_* - \mathbf{X}^T \mathbf{y} &= 0 \\ \mathbf{X}^T (\mathbf{X} \mathbf{w}_* - \mathbf{y}) + \lambda \mathbf{w}_* &= 0 \\ 2\mathbf{X}^T (\mathbf{X} \mathbf{w}_* - \mathbf{y}) + 2\lambda \mathbf{w}_* &= 0\end{aligned}\tag{1}$$

This is exactly the same equation with the target equation below

$$\nabla_{\mathbf{w}} L(\mathbf{w})|_{\mathbf{w}=\mathbf{w}^*} = 0\tag{2}$$

- b. By definition,

$$p(\mathbf{w}|\mathbf{y}) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2} \cdot \frac{(\mathbf{X}\mathbf{w} - \mathbf{y})^2}{\sigma^2}\right)\tag{3}$$

and

$$p(\mathbf{w}) = \frac{1}{\sqrt{2\pi s^2}} \exp\left(-\frac{1}{2} \cdot \frac{\mathbf{w}_n^2}{s^2}\right)\tag{4}$$

Then we break  $p(\mathbf{w}|\mathbf{y})$  into  $p(\mathbf{y}|\mathbf{w})p(\mathbf{w})$  and use logarithm to calculate the arg max

value.

$$\begin{aligned}
\mathbf{w}_* &= \arg \max_{\mathbf{w}} p(\mathbf{w}|\mathbf{y}) \\
&= \arg \max_{\mathbf{w}} \log(p(\mathbf{w}|\mathbf{y})) \\
&= \arg \max_{\mathbf{w}} \log[p(\mathbf{y}|\mathbf{w})p(\mathbf{w})] \\
&= \arg \max_{\mathbf{w}} (\log p(\mathbf{y}|\mathbf{w}) + \log p(\mathbf{w})) \\
&= \arg \max_{\mathbf{w}} [\log \frac{1}{\sqrt{2\pi\sigma^2}} \exp(-\frac{1}{2} \cdot \frac{(\mathbf{X}\mathbf{w} - \mathbf{y})^2}{\sigma^2}) + \log \frac{1}{\sqrt{2\pi s^2}} \exp(-\frac{1}{2} \cdot \frac{\mathbf{w}^2}{s^2})] \\
&= \arg \max_{\mathbf{w}} [-\frac{1}{2\sigma^2}(\mathbf{X}\mathbf{w} - \mathbf{y})^2 - \frac{1}{2s^2}\mathbf{w}^2 + \log \frac{1}{\sqrt{2\pi\sigma^2}} + \log \frac{1}{\sqrt{2\pi s^2}}] \\
&= \arg \max_{\mathbf{w}} [-\frac{1}{2\sigma^2}(\mathbf{X}\mathbf{w} - \mathbf{y})^2 - \frac{1}{2s^2}\mathbf{w}^2] \\
&= \arg \min_{\mathbf{w}} [-\frac{1}{2\sigma^2}E(\mathbf{w}) - \frac{1}{2s^2}|\mathbf{w}|^2] \\
&= \arg \min_{\mathbf{w}} [-\frac{1}{2\sigma^2}(E(\mathbf{w}) + \frac{\sigma^2}{s^2}|\mathbf{w}|^2)] \\
&= \arg \min_{\mathbf{w}} [-\frac{1}{2\sigma^2}(E(\mathbf{w}) + \lambda|\mathbf{w}|^2)]
\end{aligned} \tag{5}$$

Where

$$\lambda = \frac{\sigma^2}{s^2} \tag{6}$$

- c. When applying linear regression, we have a equation of  $\mathbf{y} = \mathbf{f} + \boldsymbol{\epsilon}$ , where  $\mathbf{f} = \mathbf{X}\mathbf{w}$ . Then we want to apply the function of Expectation and Covariance to  $\mathbf{y}$ .

$$\begin{aligned}
\mathbb{E}(\mathbf{y}) &= \mathbb{E}(\mathbf{X}\mathbf{w} + \boldsymbol{\epsilon}) \\
&= \mathbb{E}(\mathbf{X}\mathbf{w}) + \mathbb{E}(\boldsymbol{\epsilon}) \\
&= \mathbf{X}\mathbb{E}(\mathbf{w}) \\
&= \mathbf{X}m_0
\end{aligned} \tag{7}$$

and

$$\begin{aligned}
Cov(\mathbf{y}) &= Cov(\mathbf{f}) + Cov(\boldsymbol{\epsilon}) \\
&= Cov(\mathbf{X}\mathbf{w}) + Cov(\boldsymbol{\epsilon}) \\
&= \mathbf{X}Cov(\mathbf{w})\mathbf{X}^T + Cov(\boldsymbol{\epsilon}) \\
&= \mathbf{X}S_0\mathbf{X}^T + \sigma^2 I
\end{aligned} \tag{8}$$

Thus, we may conclude that  $\boldsymbol{\mu} = \mathbf{X}m_0$  and  $\boldsymbol{\Sigma} = \mathbf{X}S_0\mathbf{X}^T + \sigma^2 I$ . And

$$p(\mathbf{y}) \sim Normal(\mathbf{X}m_0, \mathbf{X}S_0\mathbf{X}^T + \sigma^2 I) \tag{9}$$

d. According to the standard identity on Gaussians, we have

$$\boldsymbol{\alpha}|\boldsymbol{\beta} \sim \text{Normal}(\boldsymbol{\alpha} + CB^{-1}(\boldsymbol{\beta} - \mathbf{b}), A - CB^{-1}C^T) \quad (10)$$

when

$$\begin{bmatrix} \boldsymbol{\alpha} \\ \boldsymbol{\beta} \end{bmatrix} \sim \text{Normal}\left(\begin{bmatrix} \mathbf{a} \\ \mathbf{b} \end{bmatrix}, \begin{bmatrix} A & C \\ C^T & B \end{bmatrix}\right) \quad (11)$$

Here we want to substitute  $\boldsymbol{\alpha}$  with  $\mathbf{y}_*$  and  $\boldsymbol{\beta}$  with  $\mathbf{y}$ . Since we already have the mean value and covariance matrix of  $\mathbf{y}$  in part d, we only need to determine the covariance matrix  $C$  between  $\mathbf{y}$  and  $\mathbf{y}_*$  and the mean value  $\boldsymbol{\mu}_*$  of  $\mathbf{y}_*$ . Here we generalize the relationship we derived in the previous part to get

$$\mathbb{E}(\mathbf{y}_*) = \mathbf{x}_*m_0 \quad (12)$$

$$\text{Cov}(\mathbf{y}_*) = \mathbf{x}_*S_0\mathbf{x}_*^T + \sigma^2I \quad (13)$$

$$\text{Cov}(\mathbf{y}_*, \mathbf{y}) = \text{Cov}(\mathbf{y}, \mathbf{y}_*)^T = \mathbf{x}_*S_0\mathbf{X}^T \quad (14)$$

Hence

$$\mathbf{f}_* = \mathbf{x}_*m_0 + \mathbf{x}_*S_0\mathbf{X}^T(\mathbf{X}S_0\mathbf{X}^T + \sigma^2I)^{-1}(\mathbf{y} - \mathbf{X}m_0) \quad (15)$$

$$\boldsymbol{\sigma}_* = \mathbf{x}_*S_0\mathbf{x}_*^T + \sigma^2I - \mathbf{x}_*S_0\mathbf{X}^T(\mathbf{X}S_0\mathbf{X}^T + \sigma^2I)^{-1}(\mathbf{x}_*S_0\mathbf{X}^T)^T \quad (16)$$

e. When we substitute  $m_0$  by 0 and  $S_0$  by  $s^2$  in the equation of  $\mathbf{f}_*$ , we have

$$\begin{aligned} \mathbf{f}_* &= \mathbf{x}_*s^2\mathbf{X}^T(\mathbf{X}s^2\mathbf{X}^T + \sigma^2I)^{-1}\mathbf{y} \\ &= \mathbf{x}_*\mathbf{X}^T(\mathbf{X}\mathbf{X}^T + \lambda I)^{-1}\mathbf{y} \end{aligned} \quad (17)$$

Then we apply a formula from matrix cookbook, which is

$$(P^{-1} + B^TR^{-1}B)^{-1}B^TR^{-1} = PB^T(BPB^T + R)^{-1} \quad (18)$$

In our context here,  $P$  is just the identity  $I$ ,  $B$  is the matrix  $\mathbf{X}$  and  $R$  is  $\lambda I$ . We substitute everything into equation(18) to get

$$\begin{aligned} \mathbf{f}_* &= \mathbf{x}_*\mathbf{X}^T(\mathbf{X}\mathbf{X}^T + \lambda I)^{-1}\mathbf{y} \\ &= \mathbf{x}_*(I + \mathbf{X}^T\frac{1}{\lambda}\mathbf{X})^{-1}\mathbf{X}^T\frac{1}{\lambda}I\mathbf{y} \\ &= \mathbf{x}_*(\lambda I + \mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{y} \end{aligned} \quad (19)$$

Where  $(\lambda I + \mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{y} = \mathbf{w}_*$  which is the closed form solution of ridge regression and the original form in line three of equation(19) is just kernel ridge regression. Here I really don't know how to solve the equation given in item e,

but my intuition tells me that the result should have some relationship with ridge regression and kernel ridge regression. When I derive through equation and just find it similar to the solution of ridge regression.

Then we determine  $\mathbb{E}[f(\mathbf{x}_*)]$

$$\begin{aligned}
 \mathbb{E}[f(\mathbf{x}_*)] &= \mathbb{E}[\mathbf{y}_* - \boldsymbol{\epsilon}] \\
 &= \mathbb{E}[\mathbf{y}_*] \\
 &= \mathbf{X}m_0 \\
 &= 0
 \end{aligned} \tag{20}$$