

## Lecture 12 — September 3

Lecturer: Purnamrita Sarkar

Scribe: Maria Shovgun

**Disclaimer:** These scribe notes have been slightly proofread and may have typos etc.

## 12.1 How collinearity effect ridge and lasso

Let us us assume that we have  $m$  copies of the same feature (similar to add  $m - 1$  copies).

### 12.1.1 Ridge

$$\begin{aligned} \min & \left[ (y - X\beta)^T (y - X\beta) + \lambda\beta^T\beta \right] \\ \hat{\beta} &= (X^T X + \lambda\mathbf{I})^{-1} X^T y \\ a &= \frac{X^T y}{X^T X + \lambda} \end{aligned} \tag{12.1}$$

The same for  $m$  copies.

$$\hat{\beta}' = \left( x^T x \cdot J + \lambda\mathbf{I} \right)^{-1} x^T y \tag{12.2}$$

Using

$$\begin{aligned} X^T X &= (x^T x)\mathbf{J} \\ X^T y &= (x^T y)\mathbf{1}_m \end{aligned} \tag{12.3}$$

Eq. ?? becomes

$$\hat{\beta}' = (x^T y)(c\mathbf{J} + \lambda\mathbf{I})_{1_m}^{-1} \tag{12.4}$$

Now, using the following expressions

$$\begin{aligned} \mathbf{J} &= \mathbf{1}_m \cdot \mathbf{1}_m^T = (m, \mathbf{1}_m) \\ c\mathbf{J} &= (cm, \mathbf{1}_m) \\ c\mathbf{J} + \lambda\mathbf{I} &= (cm + \lambda, \mathbf{1}_m) \end{aligned} \tag{12.5}$$

Eq. ?? becomes

$$\hat{\beta}' = (x^T y) \frac{\mathbf{1}_m}{cm + \lambda} \tag{12.6}$$

Therefore

$$\hat{\beta}' = ax \frac{c + \lambda}{cm + \lambda} = ax \frac{1 + \frac{\lambda}{X^T X}}{m + \frac{\lambda}{X^T X}} \approx \frac{a}{m} \tag{12.7}$$

The last equality holds for small  $\lambda$ .

Now let's solve the HW question.

$$x_{i1} = x_{i2} = x_{i3} = \dots = x_{im} = x \quad (12.8)$$

Then

$$\begin{aligned} & \sum_i \left( y_i - x_i \beta_1 - \sum_{j \neq i} X_{ij} \beta_j \right)^2 + \lambda \sum_j \beta_j^2 \\ &= \sum_i \left( y_i - \sum_j X_{ij} (\beta'_{11} + \beta'_{12} + \dots + \beta'_{1m}) - \sum_j X_{ij} \beta_1 \right)^2 + \lambda \sum_{j=1}^m (\beta'_{1j})^2 + \sum_{j=2}^p (\beta'_j)^2 \\ &= \left\{ \sum_i \left( y_i - \sum_j X_{ij} (\beta'_{11} + \beta'_{12} + \dots + \beta'_{1m}) - \sum_j X_{ij} \beta_1 \right)^2 + \lambda \sum_{j=1}^p (\beta'_j)^2 \right\} + \sum_{j=1}^m (\beta'_{j1})^2 - \beta_1'^2 \end{aligned} \quad (12.9)$$

When  $n$  is large, the first part of the optimization dominates, and so we can make the following approximate argument. Say we keep the first part pinned at its original optimal value,

$$\beta'_{11} + \beta'_{12} + \dots + \beta'_{1m} = a \quad (12.10)$$

Minimizing

$$\min \sum_{k=1}^m (\beta'_{kj})^2 \quad (12.11)$$

Consequently, all  $\beta'_{ij}$  should be the same, such that

$$\sum_j \beta'_{ij} = a \quad (12.12)$$

## 12.2 Logistic regression

### 12.2.1 Multiclass logistic regression

$b$  for every class.

$$P(y = k | x_j) \propto e^{\beta_k^T x_i} = \frac{e^{\beta_k^T x}}{\sum_i e^{\beta_i^T x_i}} \quad (12.13)$$

Assuming

$$\beta'_k = 0_p \quad (12.14)$$

and applying the transformation

$$\beta'_j = \beta_j - \beta_k, j \neq k \quad (12.15)$$

Eq. ?? becomes

$$P(y = k|x_j) = \frac{1}{\sum_c e^{(\beta_c - \beta_k)^T x_i}} \quad (12.16)$$

Discriminative, because does not care about  $x$  distribution (generated). It cares only about how  $(Y|x)$  was generated.

### 12.2.2 Generative analogued (Linear discriminant analysis)

$$X_i|y_i = K \sim N(\mu_k, \Sigma_k), \quad (12.17)$$

where  $\mu$  and  $\Sigma$  are known. how to decide which class involves any particular points (Point classification).

$$P(y_i = k) = \pi_k \quad (12.18)$$

$$\begin{aligned} P(Y = 1|x_\xi) &\propto P(x|Y = 1) P(Y = 1) \\ &= \frac{1}{(2\pi)^{p/2} |\Sigma|} \exp\left(-\frac{(x - \mu_1)^T \Sigma_1^{-1} (x - \mu_1)}{2}\right) \pi_1 \end{aligned} \quad (12.19)$$

Decision rule. Classify as 1 if

$$P(y_i = 1|x) \geq P(Y = 0|x) \quad (12.20)$$

Similarly,

$$\log P(y_i = 1|x) \geq \log P(Y = 0|x) \quad (12.21)$$

$$\begin{aligned} -\frac{1}{2} \log |\Sigma_1| - \frac{(x - \mu_1)^T \Sigma_1^{-1} (x - \mu_1)}{2} + \log \pi_1 &\geq \\ -\frac{1}{2} \log |\Sigma_2| - \frac{(x - \mu_2)^T \Sigma_2^{-1} (x - \mu_2)}{2} + \log \pi_2 \end{aligned} \quad (12.22)$$

If

$$\Sigma_1 = \Sigma_2, \quad (12.23)$$

then we get a linear decision boundary.

$$-\frac{(x - \mu_1)^T \Sigma^{-1} (x - \mu_1) + (x - \mu_2)^T \Sigma^{-1} (x - \mu_2)}{2} \geq \log \frac{\pi_2}{\pi_1} \quad (12.24)$$

$$2\mu_1^T \Sigma^{-1} (x - \mu_1) - 2\mu_2^T \Sigma^{-1} (x - \mu_2) + \frac{\mu_1^T \Sigma^{-1} \mu_1 - \mu_2^T \Sigma^{-1} \mu_2}{2} \geq \log \frac{\pi_2}{\pi_1} \quad (12.25)$$

$$-\left(x - \frac{\mu_1 + \mu_2}{2}\right)^T \Sigma^{-1} (\mu_2 - \mu_1) \geq \log \frac{\pi_2}{\pi_1} \quad (12.26)$$