

Solution to Homework 3b

Name: Chen Shen

NetID: cs5236

1.

The problem with (a): There is no case for test, so that we would never know the performance of the model we trained. There may also be some risk of overfitting.

The problem with (b): The test error could vary significantly depending on samples selected. Only use limited number of samples for training. Problems particularly bad for data with limited number of samples.

The problem with (c): Need more computing capacity. An accurate result requires K fits of parameters.

I would use method (c) when N is relatively small. But when N becomes very large, I still prefer method (c) with a smaller K .

I will use One Standard Error Rule to decide the optimal parameters to do prediction.

2.

$$\begin{aligned} y^k &= \beta^k \mathbf{x}^T \\ \Rightarrow \frac{1}{K} \sum_{k=1}^K y^k &= \frac{1}{K} \sum_{k=1}^K (\beta^k \mathbf{x}^T) \\ &= \frac{1}{K} \sum_{k=1}^K (\beta^k) \mathbf{x}^T \\ &= \bar{\beta} \mathbf{x}^T \end{aligned}$$

3.

(1) When N is very large.

First, we choose a certain K based on the computing power. Then we can get the nonzeros along with the optimal α by K-fold cross validation. Finally, drop the zero ones and sort the coefficients with decreasing magnitude. After that, we can reach the optimal subset.

(2) When N is relatively small.

The steps are same with the previous case, but we can use a larger K to reach a more accurate result. Even $K = N$, so one sample is left out.

4.

For example, filtering method, wrapper method, embedded method, forward stepwise algorithm.

5.

(a) The mean is

$$\begin{aligned}
\frac{1}{N} \sum_i y_i &= \frac{1}{N} \sum_i \frac{y_i^r - \bar{y}}{\sigma_y} \\
&= \frac{1}{N} \frac{\sum_i (y_i^r - \bar{y})}{\sigma_y} \\
&= \frac{1}{N} \frac{\sum_i y_i^r - \sum_i \bar{y}}{\sigma_y} \\
&= 0
\end{aligned}$$

The variance is

$$\begin{aligned}
\frac{1}{N} \sum_i (y_i - 0)^2 &= \frac{1}{N} \sum_i y_i^2 \\
&= \frac{1}{N} \sum_i \left[\frac{y_i^r - \bar{y}}{\sigma_y} \right]^2 \\
&= \frac{1}{N} \sum_i \frac{(y_i^r - \bar{y})^2}{\sigma_y^2} \\
&= \frac{1}{N} \frac{\sum_i (y_i^r - \bar{y})^2}{\sigma_y^2} \\
&= \frac{1}{N} \frac{J \cdot \sigma_y^2}{\sigma_y^2} \\
&= 1
\end{aligned}$$

For $x_{i,j}$, just follow the steps above and replace $y_i, y_i^r, \bar{y}, \sigma_y$ with $x_{i,j}, x_{i,j}^r, \bar{x}_j, \sigma_j$.

(b) Since both the target y and the features x_j are normalized, we have

$$\sum \hat{y} = 0.$$

So

$$\sum \beta_0 + \sum \beta_1 x_1 + \cdots + \sum \beta_J x_J = 0.$$

Because

$$\sum x_j = 0,$$

there should be

$$\sum \beta_0 = 0,$$

i.e.

$$\beta_0 = 0.$$

(c)

$$\begin{aligned}
\hat{y}_i &= \sum_{j=1}^J \beta_j x_{i,j} \\
&= \sum_{j=1}^J \beta_j \frac{x_{i,j}^r - \bar{x}_j}{\sigma_j} \\
&= \sum_{j=1}^J \frac{\beta_j}{\sigma_j} x_{i,j}^r - \sum_{j=1}^J \frac{\beta_j}{\sigma_j} \bar{x}_j \\
&= - \sum_{j=1}^J \frac{\beta_j}{\sigma_j} \bar{x}_j + \sum_{j=1}^J \frac{\beta_j}{\sigma_j} x_{i,j}^r \\
&= \beta_0^r + \sum_{j=1}^J \beta_j^r x_{i,j}^r
\end{aligned}$$

So we have

$$\beta_0^r = - \sum_{j=1}^J \frac{\beta_j}{\sigma_j} \bar{x}_j$$

and

$$\beta_j^r = \frac{\beta_j}{\sigma_j}$$

6.

Without regularization, large positive and negative coefficients cancel each other for correlated features, resulting in high variance of the resulting models. Besides, the regularization can remove the intercept.

7.

Ridge Regression is raised to solve multicollinearity. It uses L2 norm to simplify the calculation. However, it cannot shrink parameters to zero. So it can not be used to do feature selection.

LASSO Regression is raised to do feature selection. It uses L1 norm. So the calculation is more complex and there is no analytical solution. However, it can shrink some parameters to zero so as to select some certain features.

8.

First, we do differentiating with respect to β . Then set the result to be zero.

$$\begin{aligned}
\frac{\partial J(\beta)}{\partial \beta} &= 2\mathbf{A}^T \mathbf{A} \beta - 2\mathbf{A}^T \mathbf{y} + 2\alpha \beta \\
&= \mathbf{A}^T \mathbf{A} \beta - \mathbf{A}^T \mathbf{y} + \alpha \beta \\
&= 0 \\
\Rightarrow (\mathbf{A}^T \mathbf{A} + \alpha \mathbf{I}) \beta_{\text{opt}} &= \mathbf{A}^T \mathbf{y} \\
\Rightarrow \beta_{\text{opt}} &= (\mathbf{A}^T \mathbf{A} + \alpha \mathbf{I})^{-1} \mathbf{A}^T \mathbf{y}
\end{aligned}$$

9.

First, we define an artificial data set by

$$\mathbf{A}^* = (1 + \alpha\lambda)^{-\frac{1}{2}} \begin{pmatrix} \mathbf{A} \\ \sqrt{\alpha\lambda}\mathbf{I} \end{pmatrix},$$

$$\mathbf{y}^* = \begin{pmatrix} \mathbf{y} \\ 0 \end{pmatrix}$$

Let $\alpha^* = \frac{\alpha(1-\lambda)}{\sqrt{1+\alpha\lambda}}$ and $\boldsymbol{\beta}^* = \sqrt{1+\alpha\lambda}\boldsymbol{\beta}$. Then the elastic-net can be rewritten as

$$J(\boldsymbol{\beta}) = \|\mathbf{y}^* - \mathbf{A}^*\boldsymbol{\beta}^*\|^2 + \alpha^*\|\boldsymbol{\beta}^*\|_1$$

which is turned into a LASSO problem.