

**Homework 1: Due Monday, October 3, 2022 by 4:59PM**

**Please read these instructions to ensure you receive full credit on your homework.** Submit the written portion of your homework as a *single* PDF file through Courseworks (less than 5MB). In addition to your PDF write-up, submit all code written by you in their original extensions through Courseworks (e.g., .m, .r, .py, .c). Any coding language is acceptable, but do not submit notebooks. Also, do not wrap your files in .rar, .zip, .tar and do not submit your write-up in .doc or other file type. When resubmitting homeworks, please be sure to resubmit *all files*. Your grade will be based on the contents of one PDF file and the original source code. Additional files will be ignored. We will not run your code, so everything you are asked to show should be put in the PDF file. Show all work for full credit.

**Problem 1 (written)**

Consider a regression problem, in which we observe pairs  $(x_1, y_1), \dots, (x_n, y_n)$ , with  $x_i \in \mathbb{R}^d$  and  $y_i \in \mathbb{R}$ . We seek  $w$  such that

$$x_i w \approx y_i,$$

by *least squares*:

$$w_{\text{LS}} = \arg \min_w \mathcal{L}_{\text{LS}}(w) \equiv \sum_{i=1}^n (x_i^T w - y_i)^2. \quad (1)$$

In lecture, we argued that

$$\nabla \mathcal{L}_{\text{LS}}(w) = 2X^T(Xw - y), \quad (2)$$

where

$$X = \begin{bmatrix} x_1^T \\ x_2^T \\ \vdots \\ x_n^T \end{bmatrix} \in \mathbb{R}^{n \times d}, \quad \text{and} \quad y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} \in \mathbb{R}^n. \quad (3)$$

We used this to show that whenever  $X^T X$  has full rank  $d$

$$w_{\text{LS}} = (X^T X)^{-1} X^T y. \quad (4)$$

**Please answer the following questions:**

**1.1. Number of samples for uniqueness:** How large does  $n$  need to be compared to  $d$  for it to be possible for  $X^T X$  to have full rank  $d$ ? When this condition is not satisfied, how should we modify the formulation (1)?

$$\begin{aligned} X^T X &= \begin{bmatrix} 1 & 1 & \dots & 1 \\ x_1 & x_2 & \dots & x_n \end{bmatrix} \begin{bmatrix} 1 \\ x_1 \\ \vdots \\ x_n \end{bmatrix} = \begin{bmatrix} n & \sum x_i \\ \sum x_i & \sum x_i^2 \end{bmatrix} \\ \text{D-dim. } & \begin{bmatrix} d \times n \end{bmatrix} \begin{bmatrix} n \times d \end{bmatrix} \\ & \downarrow \\ & \begin{bmatrix} d \times d \end{bmatrix} \end{aligned}$$

$X^T X$  is full rank if  $\text{rank} = d$ , is invertible, and  $\neq 0$

if  $X \in \mathbb{R}^{n \times d}$  then  $X^T \in \mathbb{R}^{d \times n}$

$$X^T X w_{LS} = X^T y$$

$$X^T X = \sum_{i=1}^n x_i x_i^T \quad \text{rank of } 1 \text{ if } 1 \text{ data point.}$$

Have  $n$  data points

$X^T X$  bounded by number of data points

$$X^T X = \sum_{i=1}^n x_i x_i^T$$

$$\text{rank}(X^T X) \leq n$$

Need  $n \geq d$  for  $X^T X$  to be full rank

Condition not satisfied: if two inputs are perfectly correlated (eg.  $x_2 = 3x_1$ ).

Then  $X^T X$  is singular and the least squares coefficients  $\tilde{w}$  are not uniquely defined.

If  $n$  is not  $\gg d$  (say,  $n > 10d$ ) then solution is not unique and need to add regularization to constrain the model parameters.

$$w = \arg\min \|y - Xw\|^2 + \lambda g(w)$$

One regularizer is  $g(w) = \|w\|^2$ , called ridge regression, which penalizes large  $w$ .

The regularization parameter  $\lambda$  generates a tradeoff between the first and second terms (loss + regularizer).

Section 5.2.3 Hasties textbook.

**1.2. Unbiased estimates:** Suppose that  $y = Xw_{\text{true}} + z$ , with  $z \sim \mu$  a random (noise) vector with probability distribution  $\mu$ . Under what conditions on  $z$ , is the least squares solution *unbiased*, i.e.,

$$\mathbb{E}[w_{\text{LS}}] = w_{\text{true}}? \quad (5)$$

Note: in lecture we showed that if  $z \sim_{\text{iid}} \mathcal{N}(0, \sigma^2)$ ,  $w_{\text{LS}}$  is unbiased. However, this holds more broadly – please make your conditions as broad as possible for full credit!

**Lemma:** If  $y = Xw + z$  where  $X$  is non random and  $X^T X$  is invertible and  $\mathbb{E}(z) = 0$ , then  $\mathbb{E}(w_{\text{LS}}) = w_{\text{true}}$

$$\text{Proof: } w_{\text{LS}} = (X^T X)^{-1} X^T y = (X^T X)^{-1} (X^T X w + z) = w + (X^T X)^{-1} X^T z$$

$$\begin{aligned} \mathbb{E}(w_{\text{LS}}) &= w + (X^T X)^{-1} X^T \mathbb{E}(z) \\ &= w \end{aligned}$$

The least squares estimates  $w_{\text{LS}}$  are unbiased for  $w$  as long as  $z$  has mean zero. Does not require normally distributed errors. Does not even make assumptions about  $\text{var}(z)$ . To study  $\text{var}(w_{\text{LS}})$  we will need assumptions on  $\text{var}(w)$  but not its mean

**1.3. Least squares by gradient descent:** We can also compute  $w_{\text{LS}}$  iteratively, using gradient descent – namely, by choosing some initial guess  $w^0$  and updating  $w$  as

$$w^{k+1} = w^k - t \nabla \mathcal{L}_{\text{LS}}(w^k). \quad (6)$$

**Part A.** Show that

$$w^{k+1} - w_{\text{LS}} = (I - 2tX^T X)(w^k - w_{\text{LS}}) \quad (7)$$

Hint: use our expressions for  $\nabla \mathcal{L}_{\text{LS}}$  and  $w_{\text{LS}}$ .

$$w^{k+1} = w^k - t \nabla \mathcal{L}_{\text{LS}}(w^k)$$

$$\text{substitute } \nabla \mathcal{L}_{\text{LS}}(w^k) = 2X^T(Xw^k - y)$$

$$w^{k+1} = w^k - 2tX^T(Xw^k - y)$$

$$w^{k+1} = w^k + 2tX^T(y - Xw^k)$$

$$w^{k+1} = w^k + 2t(X^T y - X^T X w^k)$$

substitute  $X^T X w_{LS} = X^T y$

$$\begin{aligned} w^{k+1} &= w^k + 2t (X^T X w_{LS} - X^T X w^k) \\ &= w^k + 2t (X^T X w_{LS}) - 2t (X^T X w^k) \\ &= w^k - 2t (X^T X w^k) + 2t (X^T X w_{LS}) \\ &= w^k \underbrace{(I - 2t X^T X)}_{\downarrow I} + 2t X^T X w_{LS} \end{aligned}$$

then subtract  $w_{LS}$  from both sides

$$\begin{aligned} w^{k+1} - w_{LS} &= w^k (I - 2t X^T X) + 2t X^T X w_{LS} - w_{LS} \\ w^{k+1} - w_{LS} &= w^k (I - 2t X^T X) - w_{LS} (I - 2t X^T X) \end{aligned}$$

$$w^{k+1} - w_{LS} = (w^k - w_{LS}) (I - 2t X^T X)$$

**Part B.** Let  $\lambda_1 \geq \dots \geq \lambda_n$  denote the eigenvalues of the symmetric matrix  $I - 2tX^T X$ . Argue that if  $t$  is chosen such that  $|\lambda_i| < 1$  for all  $i$ , then as  $k \rightarrow \infty$ ,  $w^k \rightarrow w_{LS}$ .

Hint: You can use the following inequality. For a symmetric matrix  $M$  with eigenvalues  $\lambda_1, \dots, \lambda_n$ , and a vector  $x$ ,

$$\|Mx\|_2 \leq \left( \max_i |\lambda_i| \right) \|x\|_2.$$

$t$  chosen such that  $|\lambda_i| < 1$  for all  $i \rightarrow$  full column rank assumption

to ensure that it is strictly better than -1 we need  $t < \frac{1}{\lambda_{\max}(X^T X)}$

$$\begin{aligned} \|X^T X\|_2 &\leq \lambda_{\max} X^T X \\ \|(I - 2t X^T X)X\|_2 &\leq \lambda_{\max} X^T X \end{aligned}$$

$$w_{LS} = (X^T X)^{-1} X^T y$$

$$w_{LS} = w_{RR} = (X^T X + \lambda I)^{-1} X^T y$$

### Problem 2 (written)

As in Problem 1, you have data  $(x_i, y_i)$  for  $i = 1, \dots, n$ , where  $x \in \mathbb{R}^d$  and  $y \in \mathbb{R}$ . Here, we apply ridge regression, setting

$$w_{\text{RR}} = \arg \min_w \sum_{i=1}^n (x_i^T w - y_i)^2 + \lambda w^T w. \quad (8)$$

Please answer the following questions:

2.1. Demonstrate that the ridge regression solution  $w_{\text{RR}}$  is given by

$$w_{\text{RR}} = (\lambda I + X^T X)^{-1} X^T y. \quad (9)$$

Hint: differentiate the objective function in (8) and set the derivative equal to zero.

$$\arg \min_w \underbrace{\sum_{i=1}^n (x_i^T w - y_i)^2}_{\text{least squares } \|y - Xw\|^2} + \underbrace{\lambda w^T w}_{\lambda \|w\|^2}$$

$$\mathcal{L} = (y - Xw)^T (y - Xw) + \lambda w^T w$$

take gradient of  $\mathcal{L}$  w.r.t  $w$  and set to zero

$$\nabla \mathcal{L} = 2X^T X w - 2X^T y + 2\lambda w = 0$$

solve for  $w$

$$w(2\lambda + 2X^T X) - 2X^T y = 0$$

$$w(2\lambda + 2X^T X) = 2X^T y$$

$$w_{\text{RR}} = 2X^T y (2\lambda + 2X^T X)^{-1}$$

2.2. Consider two potential values of  $\lambda$ :  $\lambda = 0.1$  and  $\lambda = 1,000$ . For which value of  $\lambda$  do you expect the ridge regression solution  $w_{\text{RR}}$  to exhibit greater *bias*? For which value of  $\lambda$  do you expect it to exhibit greater *variance*?

The regularization term  $\lambda w^T w$  is imposed to penalize values in  $w$  that are large. This reduces potential high-variance predictions from least squares. As  $\lambda$  increases, the flexibility of the ridge regression fit decreases, leading to decreased variance but increased bias.

So  $\lambda = 1,000$  leads to more bias since the regularizing term is large.  $\lambda = 0.01$  exhibits more variance, since the regularizing term is made very small and basically follows the least squares solution.

The least squares solution has increased variance but decreased bias.

2.3. Describe a procedure for choosing  $\lambda$  in practice.

In practice, bias is related to a model failing to fit the training dataset and variance is related to a model failing to fit the testing dataset.

An optimal value of  $\lambda$  finds the minimum bias and variance despite the tradeoff. One way to do this is through cross validation.

1. Split data into  $K$  roughly equal groups
  2. Train the model on  $K-1$  groups and predict the held-out  $K^{\text{th}}$  group.
  3. Repeat  $K$  times, holding out each group once. Calculate prediction error of the fitted model when predicting the  $K^{\text{th}}$  part of the data. Typical choices of  $K$  are 5 or 10. The case  $K=N$  is known as leave-one-out cross val.
- The cross validation estimate of prediction error is:

$$CV(\hat{f}) = \frac{1}{N} \sum_{i=1}^N L(y_i, \hat{f}^{-x(i)}(x_i))$$

$$\text{where } K\{1, \dots, N\} \mapsto \{1, \dots, K\}$$

4. After all the iterations are done, the model will be trained each time using a different fold. To get the cross-validation score, take the average of all the individual hold out scores.
5. This cross-validation score will have used a specific  $\lambda$  on the data. Steps 1-4 should be repeated for various  $\lambda$ . The optimal value of  $\lambda$  will be that with the highest cross-validation score.

### Problem 3 (coding)

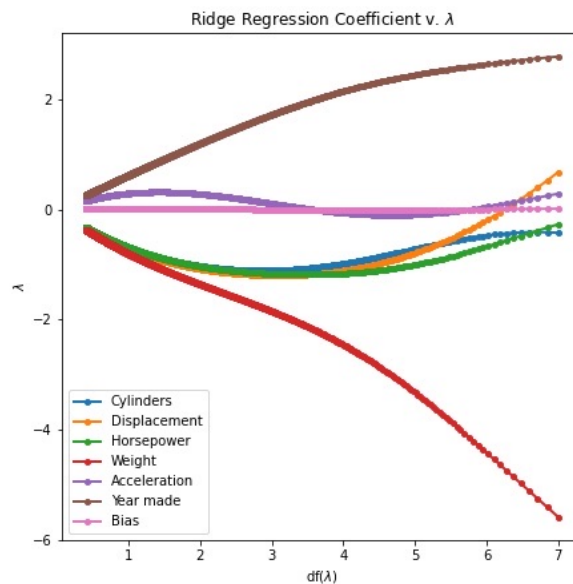
In this problem you will analyze data using the linear regression techniques we have discussed. The goal of the problem is to predict the miles per gallon a car will get using six quantities (features) about that car. The zip file containing the data can be found on Courseworks.<sup>1</sup> The data is broken into training and testing sets. Each row in both “X” files contain six features for a single car (plus a 1 in the 7th dimension) and the same row in the corresponding “y” file contains the miles per gallon for that car.

Remember to submit all original source code with your homework. Put everything you are asked to show below in the PDF file.

Part I. Using the training data only, write code to solve the ridge regression problem

$$\mathcal{L} = \lambda \|w\|^2 + \sum_{i=1}^{350} \|y_i - x_i^T w\|^2.$$

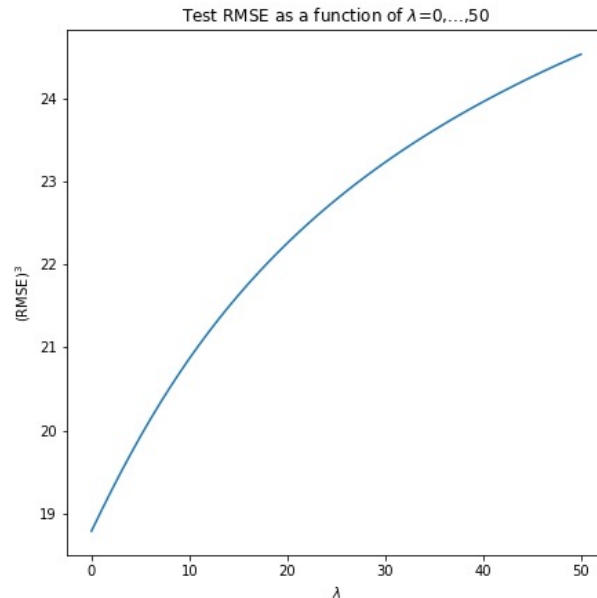
- (a) For  $\lambda = 0, 1, 2, 3, \dots, 5000$ , solve for  $w_{RR}$ . (Notice that when  $\lambda = 0$ ,  $w_{RR} = w_{LS}$ .) In one figure, plot the 7 values in  $w_{RR}$  as a function of  $df(\lambda)$ . You will need to call a built in SVD function to do this as discussed in the slides. Be sure to label your 7 curves by their dimension in  $x$ .<sup>2</sup>



- (b) Two dimensions clearly stand out over the others. Which ones are they and what information can we get from this?

"Year made" and "weight" stand out which means these two features are most predictive in determining the value of  $y$  (miles per gallon) for the car. This is also intuitive since newer cars are more fuel efficient and increasing weight will require more gas/energy to move the car.

- (c) For  $\lambda = 0, \dots, 50$ , predict all 42 test cases. Plot the root mean squared error (RMSE)<sup>3</sup> on the test set as a function of  $\lambda$ —*not* as a function of  $df(\lambda)$ . What does this figure tell you when choosing  $\lambda$  for this problem (and when choosing between ridge regression and least squares)?



With increasing  $\lambda$ , RMSE also increases. Referring to the ridge regression equation:

$$\mathcal{L} = \sum_{i=1}^n \|y_i - x_i^T w\|^2 + \lambda \|w\|^2$$

$\lambda$  controls the weight of the penalty term  $\|w\|^2$ . When  $\lambda$  is zero, the equation becomes a least squares estimate:  $\mathcal{L} = \sum \|y_i - x_i^T w\|^2$ .

Adding a penalty term reduces the variance and increases the bias.

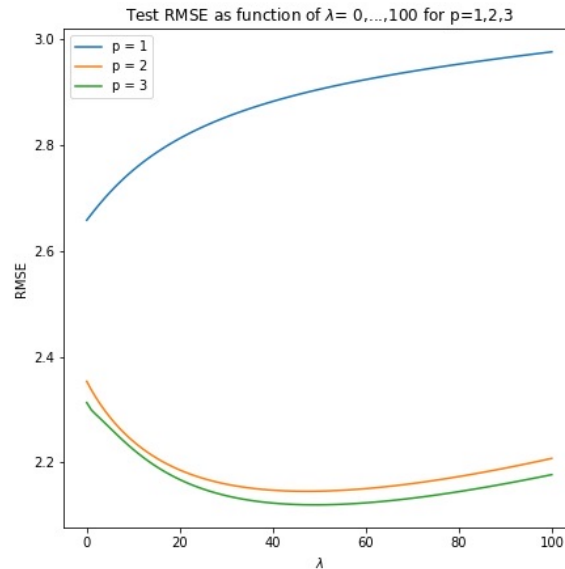
From the plot in part (b) the features most correlated with the predicted variable are most apparent with high bias/increasing  $\lambda$ . Therefore ridge regression/high  $\lambda$  is favored.

However, this plot (c) shows that a  $\lambda = 0$  has the lowest error, and therefore a low  $\lambda$  or least squares should be used.



Part 2. Modify your code to learn a  $p$ th-order polynomial regression model for  $p = 1, 2, 3$ . (You've already done  $p = 1$  above.) For this implementation use the method discussed in the slides. Also, be sure to standardize each additional dimension of your data.

- (d) In one figure, plot the test RMSE as a function of  $\lambda = 0, \dots, 100$  for  $p = 1, 2, 3$ . Based on this plot, which value of  $p$  should you choose and why? How does your assessment of the ideal value of  $\lambda$  change for this problem?



As shown also in plot from (c), a  $\lambda$  value of 0 should be chosen for  $p=1$  since this is the lowest error. However when  $p=2$ , the error decreases with increasing values of  $\lambda$ , with a valley in the RMSE graph at around  $\lambda=40$ . For  $p=3$ , the lowest error is at around  $\lambda=50$ . Therefore, the ideal value for  $\lambda$  depends on the  $p$ th-order of the polynomial regression model. It appears from this graph (d) that increasing  $p$  will decrease the RMSE. But from the data given,  $p=3$  should be used since it has the lowest error, and a  $\lambda$  value of  $\sim 50$ .