

Julio Soldevilla
EECS 545 Winter 2018 — Problem Set 1

Problem 1 THis problem has parts b,c,d,e that we present solutions in this document. The code for the solutions is in the zip file *jsolde_hw1_code.zip* under the name prob1.py.

1. For problem 1b, we used Stochastic Gradient descent to train the regression model. For this problem, The weight for the stochastic grad desc is

$$\begin{bmatrix} 22.94022368 & -0.93943062 & 1.1925544 & 0.21191992 & 0.65780365 & -2.11200674 \\ 2.74910346 & 0.30622184 & -3.12390266 & 2.95424284 & -2.46344951 & -2.00963041 \\ 0.91701884 & -4.05450806 \end{bmatrix}$$

where this last expression is just one 14×1 vector (but for the sake of space we wrote it as above).

From this weight vector, we can see that the **bias term** is 22.94022368.

The **training errors** (the MSE for training data) is 23.1929613454, while the **test errors** (the MSE for test data) is 10.9671896913.

Finally, we can see the plot of the training error at the end of each epoch in figure 1 below.

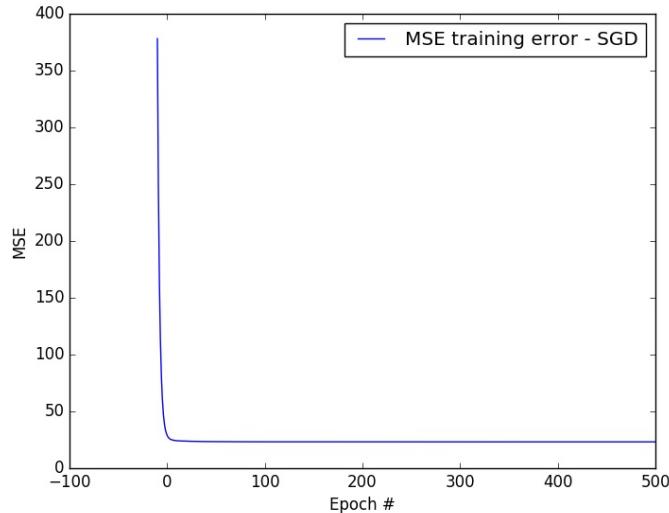


Figure 1: Graph showing the MSE training error when using SGD

2. For part c) we compute the batch gradient descent to train our model. The learned weight vector is

$$\begin{bmatrix} 22.94100877 & -0.936336221.18953456 & 0.21723452 & 0.66970526 & -2.10547565 \\ 2.75126807 & 0.30756426 & -3.12381189 & 2.95817231 & -2.4511536 \\ -2.00722057 & 0.90553088 & -4.05729245 \end{bmatrix}$$

where this last expression is just one 14×1 vector.

From this weight vector, we can see that the **bias term** is 22.94100877.

The **training errors** (the MSE for training data) is 23.191557895, while the **test errors** (the MSE for test data) is 10.9611791412.

Finally, we can see the plot of the training error at the end of each epoch in figure 2 below.

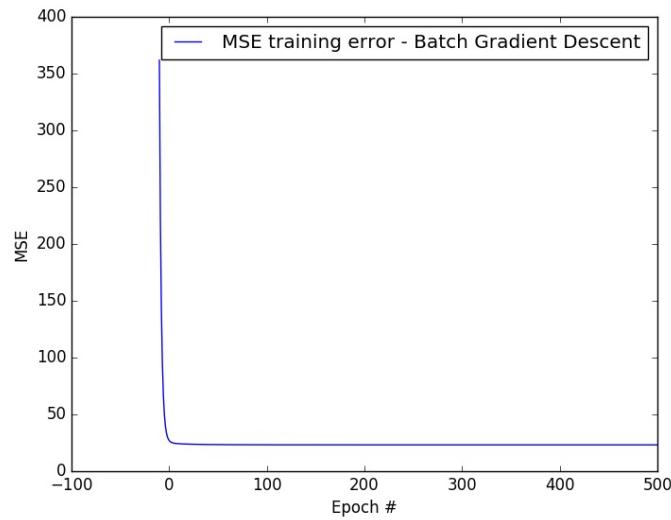


Figure 2: Graph showing the MSE training error when using BGD

3. In this part of the problem, we use the close solution to identify the weights and bias that minimize the MSE. The weight vector we obtain from using the closed form solution is

$$\begin{bmatrix} 22.94100877 & -0.936527281.18983479 & 0.2180906 & 0.66954197 & -2.10545149 \\ 2.75102471 & 0.30777503 & -3.12356704 & 2.96148512 & -2.45469868 \\ -2.00737039 & 0.90552685 & -4.05749492 \end{bmatrix}$$

(which again is a 14×1 vector but for the sake of space we wrote it as above) and from this expression we can see that the **bias that minimizes** the MSE is 22.94100877. Furthermore, from running the code, we see that The **training errors** (the MSE for training data) is 23.1915564692, while the **test errors** (the MSE for test data) is 10.9665431668.

From these computations, we can see that the closed form parameters agree very closely with the ones we computed in the previous part.

4. After computing the random train and test splits we obtained the following mean training error 21.8527415202 and the mean test error 23.3165328596

Problem 2 1. In the first of the problem, we get the following RMSE errors:

- 0th training RMSE: 9.48269466477
- 0th test RMSE: 5.99400813094
- 1-th training RMSE: 4.81576125542
- 1-th test RMSE: 3.3115771419
- 2-th training RMSE: 3.76629417811
- 2-th test RMSE: 7.14746413377
- 3-th training RMSE: 8.18681202323
- 3-th test RMSE: 12.1042413766
- 4-th training RMSE: 8.50491214305
- 4-th test RMSE: 141.55345743

Finally, we can also see the plot of these values in the followin graph (Figure 3):

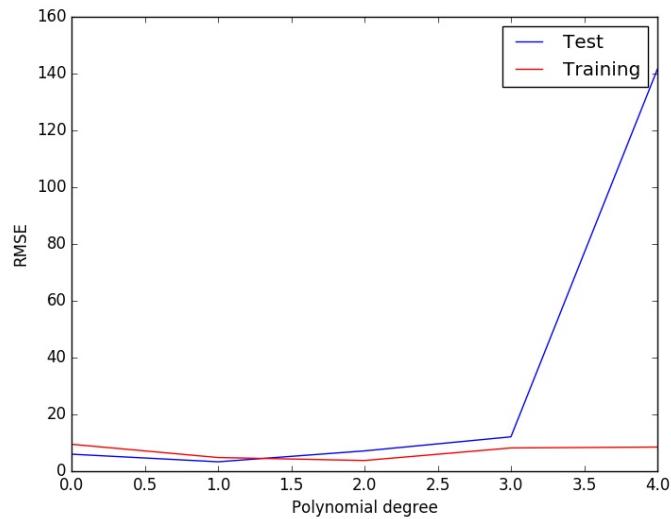


Figure 3: Graph showing the MSE error when doing polynomial features of different degrees

Since the last part of the graph is not very logical, to appreciate better the graph, we can just graph all elements except the last one, which we can see in Figure 4:

2. Finally, we get the following results for when we restrict the training set and start growing:

- This is the RMSE for training with 20% data 11.4601626769
- This is the RMSE for test with 20% data 38.918275945
- This is the RMSE for training with 40% data 2.96802006698
- This is the RMSE for test with 40% data 9.32135308145

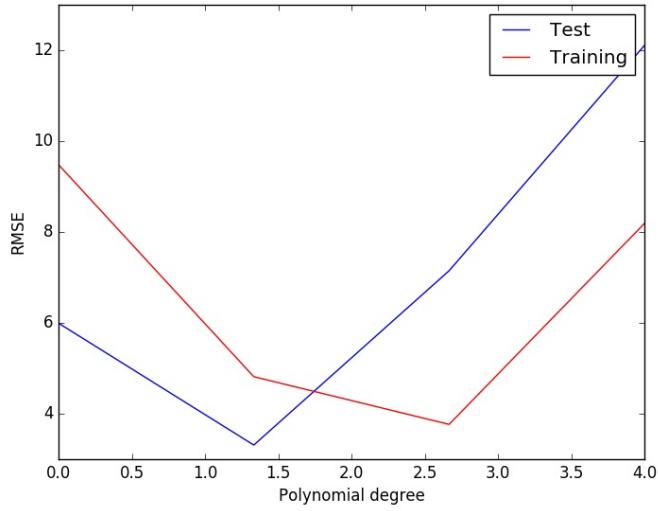


Figure 4: Graph showing the MSE error when doing polynomial features of different degrees, just removing the last value

- This is the RMSE for training with 60% data 3.14261785825
- This is the RMSE for test with 60% data 7.57749760854
- This is the RMSE for training with 80% data 3.25304667758
- This is the RMSE for test with 80% data 4.52585085547
- This is the RMSE for training with 100% data 4.81576125542
- This is the RMSE for test with 100% data 3.3115771419

Again, we can visualize this information as shown in Figure 5:

Problem 3 1. For this first problem, we show the work done in Figure 6:

2. For this problem we report all the possible RMSE, namely validation, training and test set for the possible lambdas:

- The RMSE error of training set with $\lambda = 0$ is 4.75899856011
- The RMSE error of validation set with $\lambda = 0$ is 7.84564027897
- The RMSE error of test set with $\lambda = 0$ is 4.22655365439
- The RMSE error of training set with $\lambda = 0.1$ is 5.31569487929
- The RMSE error of validation set with $\lambda = 0.1$ is 5.49507024934
- The RMSE error of test set with $\lambda = 0.1$ is 3.18380156819
- The RMSE error of training set with $\lambda = 0.2$ is 6.37532816836

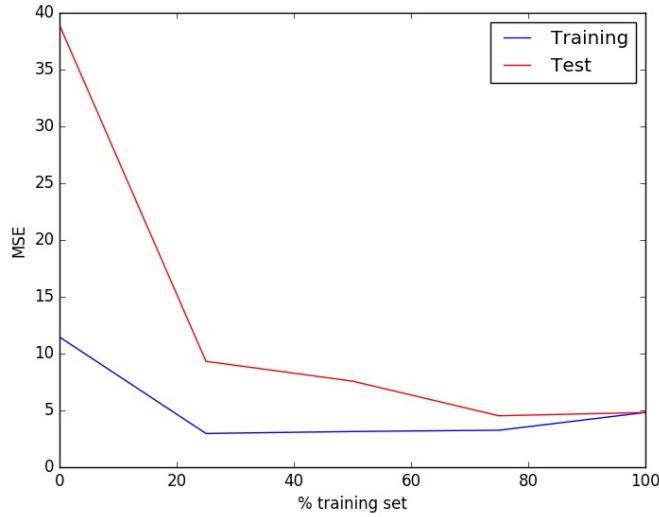


Figure 5: Graph showing the MSE error for training and test data when changing the available percentage of data.

- The RMSE error of validation set with $\lambda = 0.2$ is 4.54395472308
- The RMSE error of test set with $\lambda = 0.2$ is 3.92160042074
- The RMSE error of training set with $\lambda = 0.3$ is 7.51722511924
- The RMSE error of validation set with $\lambda = 0.3$ is 4.48444830133
- The RMSE error of test set with $\lambda = 0.3$ is 5.11613124943
- The RMSE error of training set with $\lambda = 0.4$ is 8.60568912038
- The RMSE error of validation set with $\lambda = 0.4$ is 4.93033974684
- The RMSE error of test set with $\lambda = 0.4$ is 6.29716462016
- The RMSE error of training set with $\lambda = 0.5$ is 9.60402790097
- The RMSE error of validation set with $\lambda = 0.5$ is 5.57508121459
- The RMSE error of test set with $\lambda = 0.5$ is 7.36843458243

Furthermore, we see that the lowest RMSE error on the validation set is 4.48444830133 happening on $\lambda = 0.3$. The corresponding test error for this $\lambda = 0.3$ is 5.11613124943.

Finally, we can see a picture of how the test, validation and training MSE behave for different λ in the table below, this picture corresponds to figure 7.

Problem 4 Finally, we show the theoretical in the following scanned pictures, corresponding to Figure 8, 9, 10 and 11.

$$\begin{aligned}
 & \text{② ③} \quad (\omega^\top \phi - t)^\top (\omega^\top \phi - t) = \omega^\top \phi^\top \phi \omega = \|\phi \omega\|^2 = \|\omega\|^2 \\
 & \text{④} \quad E(\omega) = \frac{1}{2N} \sum_{i=1}^N (\omega^\top \phi(x_i) - t_i)^2 + \frac{\lambda}{2} \|\omega\|^2 \\
 & \quad \left(\frac{1}{2N} (\omega^\top \phi(x_i) - t_i)^2 + \frac{\lambda}{2} \omega^\top \omega \right) \\
 & \quad = \frac{1}{2N} (\phi \omega - t)^\top (\phi \omega - t) + \frac{\lambda}{2} \omega^\top \omega \\
 & \quad = \frac{1}{2N} (\omega^\top \phi^\top \phi \omega - 2 \omega^\top \phi \omega + \omega^\top t) + \frac{\lambda}{2} \omega^\top \omega \\
 & \quad = \frac{1}{2N} \omega^\top \phi^\top \phi \omega - \frac{1}{N} \omega^\top \phi \omega + \frac{\omega^\top t}{2N} + \frac{\lambda}{2} \omega^\top \omega \\
 & \nabla_{\omega} E(\omega) = \frac{\phi^\top \phi \omega}{N} - \frac{\phi^\top t}{N} + \lambda \omega \\
 & \quad = \left(\frac{\phi^\top \phi}{N} + \lambda \mathbb{I} \right) \omega - \frac{\phi^\top t}{N} = 0 \\
 & \Rightarrow \left(\frac{\phi^\top \phi}{N} + \lambda \mathbb{I} \right) \omega = \frac{\phi^\top t}{N} \\
 & \boxed{\omega = \left(\frac{\phi^\top \phi}{N} + \lambda \mathbb{I} \right)^{-1} \frac{\phi^\top t}{N}}
 \end{aligned}$$

Figure 6: Image showing how we computed the closed form solution of the regularization problem in problem 3a.

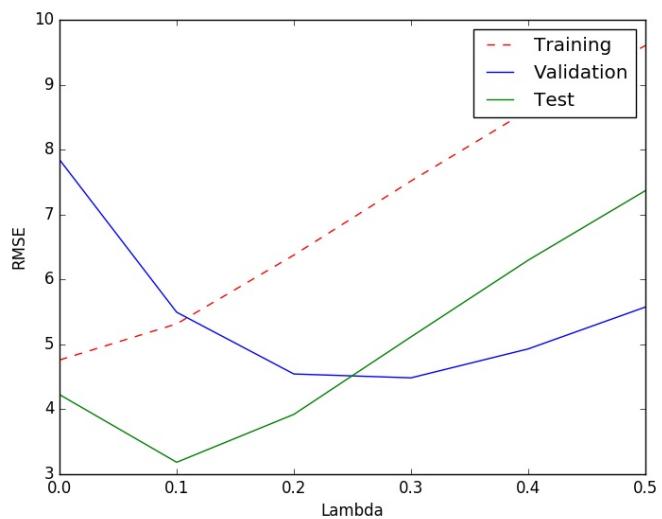


Figure 7: Image showing the behavior of RMSE when we have different regularization parameters.

(4)

$$a) \text{Want to show } E(w) = \frac{1}{2} \sum_{n=1}^N r_n (w^T x_n - t_n)^2$$

$$= \sum_{n=1}^N \frac{1}{2} r_n (w^T x_n - t_n)^2 \quad \text{let } s_n = (w^T x_n - t_n) \text{ then}$$

$$\text{we know } \underbrace{[s_1 \dots s_N]}_{1 \times N \text{ vector}} \underbrace{\begin{bmatrix} s_1 \\ \vdots \\ s_N \end{bmatrix}}_{N \times 1 \text{ vector}} = s_1^2 + \dots + s_N^2 = \sum_{i=1}^N s_i^2 \rightarrow 1 \times 1 \text{ output}$$

Considering the r_n term, we can rewrite $E(w)$ as

$$E(w) = \underbrace{[s_1 \dots s_N]}_{\text{plugging back the value of } s_n} \underbrace{\begin{bmatrix} \frac{r_1}{2} & 0 \\ 0 & \frac{r_N}{2} \end{bmatrix} \begin{bmatrix} s_1 \\ \vdots \\ s_N \end{bmatrix}}_{=}$$

$$= \underbrace{[w^T x_1 - t_1 \dots w^T x_N - t_N]}_{= \begin{bmatrix} w^T x_1 - t_1 \\ \vdots \\ w^T x_N - t_N \end{bmatrix}} \underbrace{\begin{bmatrix} \frac{r_1}{2} & 0 \\ 0 & \frac{r_N}{2} \end{bmatrix}}_{= \begin{bmatrix} \frac{r_1}{2} & 0 \\ 0 & \frac{r_N}{2} \end{bmatrix}} \underbrace{\begin{bmatrix} w^T x_1 - t_1 \\ \vdots \\ w^T x_N - t_N \end{bmatrix}}_{=}$$

$$= \underbrace{[x_w - t]}_{= \begin{bmatrix} x_w - t \\ \vdots \\ x_w - t \end{bmatrix}}^T R \underbrace{[x_w - t]}_{= \begin{bmatrix} x_w - t \\ \vdots \\ x_w - t \end{bmatrix}} \quad \text{where we know}$$

$$R = \underbrace{\begin{bmatrix} \frac{r_1}{2} & 0 \\ 0 & \frac{r_N}{2} \end{bmatrix}}_{N \times N \text{ matrix of scalars } r_i \in \mathbb{R}} \quad x_w - t = \underbrace{\begin{bmatrix} w^T x_1 - t_1 \\ \vdots \\ w^T x_N - t_N \end{bmatrix}}_{\text{mean simplify}} \quad \text{& } w \text{ is an } M \times 1 \text{ vector \& so}$$

$$N \times 1 \text{ vector of scalars } t = \begin{bmatrix} t_1 \\ \vdots \\ t_N \end{bmatrix} \quad \text{and } \begin{bmatrix} x_1 \\ \vdots \\ x_N \end{bmatrix} = X \quad \text{is an } N \times M \text{ vector of vectors where each } x_i \text{ is an } M \times 1 \text{ dimensional vector of features for data } x_i.$$

Figure 8: Image showing part a of problem 4.

We can also go the other way around. Considering the matrices X, R and W defined as before, we have

$$\begin{aligned}
 E(w) &= (Xw - t)^T R (Xw - t) = [Xw - t]^T \begin{bmatrix} r_1 & 0 \\ 0 & \ddots & 0 \\ 0 & 0 & r_N \end{bmatrix} [Xw - t] = \\
 &= \left[w^T \begin{bmatrix} x_1 \\ \vdots \\ x_N \end{bmatrix} - \begin{bmatrix} t_1 \\ \vdots \\ t_N \end{bmatrix} \right]^T \begin{bmatrix} r_1 & 0 & \dots & 0 \\ 0 & \ddots & \vdots & 0 \\ 0 & 0 & \ddots & r_N \\ 0 & 0 & \dots & 0 \end{bmatrix} \left[w^T \begin{bmatrix} x_1 \\ \vdots \\ x_N \end{bmatrix} - \begin{bmatrix} t_1 \\ \vdots \\ t_N \end{bmatrix} \right] = \\
 &= \left[w^T x_1 - t_1 \dots w^T x_N - t_N \right]^T \begin{bmatrix} r_1 & 0 & \dots & 0 \\ 0 & \ddots & \vdots & 0 \\ 0 & 0 & \ddots & r_N \\ 0 & 0 & \dots & 0 \end{bmatrix} \begin{bmatrix} w^T x_1 - t_1 \\ \vdots \\ w^T x_N - t_N \end{bmatrix} = \\
 &= \left[w^T x_1 - t_1 \dots w^T x_N - t_N \right]^T \begin{bmatrix} \frac{r_1}{2} (w^T x_1 - t_1) \\ \vdots \\ \frac{r_N}{2} (w^T x_N - t_N) \end{bmatrix} = \\
 &= \frac{r_1}{2} (w^T x_1 - t_1)^2 + \dots + \frac{r_N}{2} (w^T x_N - t_N)^2 = \frac{1}{2} \sum_{n=1}^N r_n (w^T x_n - t_n)^2 \quad \text{as desired}
 \end{aligned}$$

Figure 9: Image showing a continuation to part a of problem 4a.

(b) From the matrix expression of $E(w)$ we know

$$E(w) = (xw - t)^T R (xw - t)$$

$$= (w^T x^T - t^T) R (xw - t) = (w^T x^T - t^T) (R xw - R t)$$

$$= w^T x^T R x w - t^T R x w - w^T x^T R t - t^T R t$$

Since
 $R x w$ is 1×1 , it
 is equal to its
 transpose

$$\stackrel{\downarrow}{=} w^T x^T R x w - 2w^T x^T R t - t^T R t$$

↑
 also since R is
 diagonal, $R^T = R$

$$\text{Now finding } \nabla_w E = 2x^T R x w - 2x^T R t = 0$$

$$\text{which implies } w = \underline{\underline{(x^T R x)^{-1} x^T R t}}$$

$$\text{So } \underline{\underline{w^* = (x^T R x)^{-1} x^T R t}}$$

Figure 10: Image showing the work done to solve problem 4b.

③ We know from lecture that after applying ln to the likelihood gives:

$$\ln(p(t|x; w)) = \sum_{n=1}^N \ln p(t_n|x_n; w) =$$

$$= \sum_{n=1}^N \left[-\frac{1}{2} \ln 2\pi \sigma_n^2 - \frac{1}{2\sigma_n^2} (w^T x_n - t_n)^2 \right]$$

$$\text{So } w_{ML} = \arg \max_w \ln(p(t|x; w)) =$$

$$= \arg \max_w \sum_{n=1}^N \left[-\frac{1}{2} \ln 2\pi \sigma_n^2 - \frac{1}{2\sigma_n^2} (w^T x_n - t_n)^2 \right]$$

✓ can take out
 $-\frac{1}{2} \ln 2\pi \sigma_n^2$ term
 because maximizing
 w.r.t w

$$= \arg \max_w \sum_{n=1}^N \left[-\frac{1}{2\sigma_n^2} (w^T x_n - t_n)^2 \right]$$

$$= \arg \min_w \sum_{n=1}^N \underbrace{\frac{1}{2\sigma_n^2} (w^T x_n - t_n)^2}$$

So to find the Max.

Likelihood estimate w_{ML}

- we need to minimize a weighted least squares.

thru (i) El(w) from part a)

Hence to find the w_{ML} we need to solve a weighted linear regression problem. In the weighted linear regression problem, the r_i 's are just the terms $\frac{1}{\sigma_i^2}$

Figure 11: Image showing the work done to find the solution to problem 4c.