

# CSE546 - Homework # 1 - Solutions

Cheng-Yen Yang

April 23, 2021

## Problem A.0

**a.**

**Bias** is the expected difference between the predictions or estimates of our model and the ground-truth values from the data, or in my own words, how **far** it is off the target. **Variance** is the expected value of the squared difference between the estimates of our model and the expected value of the estimate, or in my own words, how **much** the model will change when we retrain the model using a new dataset with the same distribution. **Bias-Variance Trade-off** is the effect that when a model has a lower bias it is often found with a higher variance or vice versa. This is because that the error can be derived into bias squared and variance error through bias-variance decomposition which implies that there is a trade-off effect between bias and variance in the models.

**b. Complexity  $\uparrow$  Bias  $\downarrow$  Variance  $\uparrow$  / Complexity  $\downarrow$  Bias  $\uparrow$  Variance  $\downarrow$**

Typically, as the model complexity increases, the bias tends to decrease and the variance tends to increase; as the model complexity decreases, the bias tends to increase and the variance tends to decrease.

**c. False.**

The bias is not likely to be affected by the number of data we use in training, the number of data will directly influence towards variance.

**d. True.**

The variance of a model decreases as the amount of training data available increases.

**e. True.**

A learning algorithm will always generalize better if we use fewer features to represent our data.

**f. Train set.**

One should not use the test set for tuning at any circumstances during the training stage. Making a proper validation set from the original train set is also one of the options if we want to obtain better performance on unseen data.

**g. False.**

The training error of a function on the training set usually provides an underestimate of the true error of that function.

## Problem A.1

**a.**

Since the log of the likelihood function  $\log L(\lambda)$  has the same maximum as the likelihood function  $L(\lambda)$  itself, we have:

$$L(\lambda) = \prod_{i=1}^n \text{Poisson}(x_i|\lambda) = \prod_{i=1}^n e^{-\lambda} \frac{\lambda^{x_i}}{x_i!} \quad (1)$$

and,

$$\log(L(\lambda)) = \log\left(\prod_{i=1}^n e^{-\lambda} \frac{\lambda^{x_i}}{x_i!}\right) \quad (2)$$

$$= \sum_{i=1}^n \log\left(e^{-\lambda} \frac{\lambda^{x_i}}{x_i!}\right) \quad (3)$$

$$= \sum_{i=1}^n (-\lambda + x_i \log(\lambda) - \log(x_i!)) . \quad (4)$$

By making the derivative of the log of likelihood function 0:

$$\frac{d}{d\lambda} \log(L(\lambda)) = \sum_{i=1}^n \left(-1 + \frac{x_i}{\lambda}\right) = 0, \quad (5)$$

we have:

$$\lambda = \sum_{i=1}^n \frac{x_i}{n}. \quad (6)$$

As for the case in A.1.a, the maximum-likelihood estimate of  $\lambda$  is:

$$\lambda = \sum_{i=1}^5 \frac{x_i}{5} = \frac{1}{5}(x_1 + x_2 + x_3 + x_4 + x_5) \quad (7)$$

**b.**

Following similar procedure in A.1.a, as for the case in A.1.b, the maximum-likelihood estimate of  $\lambda$  is:

$$\lambda = \sum_{i=1}^6 \frac{x_i}{6} = \frac{1}{6}(x_1 + x_2 + x_3 + x_4 + x_5 + x_6) \quad (8)$$

**c.**

By using the  $x_i$  given in the problem statement, we have  $\lambda_{5+} = 13/5$  and  $\lambda_{6+} = 16/6 = 8/3$ .

## Problem A.2

The maximum likelihood estimate for  $\theta$  can be express as:

$$\theta_{MLE} = \operatorname{argmax} \prod_{i=1}^n \frac{1}{\theta} \mathbf{1}\{x \in [0, \theta]\} \quad (9)$$

$$= \frac{1}{\theta^n} \quad (10)$$

so we need to pick the smallest possible  $\theta$  which satisfy  $x_i \in [0, \theta]$  for all  $i$ :

$$\theta_{MLE} = \max(x_i). \quad (11)$$

### Problem A.3

**a.**

To prove that  $\mathbb{E}_{\text{train}}[\hat{\epsilon}_{\text{train}}(f)] = \mathbb{E}_{\text{test}}[\hat{\epsilon}_{\text{test}}(f)] = \epsilon(f)$ , we begin by showing that the train error:

$$\mathbb{E}_{\text{train}}[\hat{\epsilon}_{\text{train}}(\hat{f})] = \mathbb{E} \left[ \frac{1}{N_{\text{train}}} \sum_{(x,y) \in S_{\text{train}}} ((f(x) - y)^2) \right] \quad (12)$$

$$= \frac{1}{N_{\text{train}}} \sum_{(x,y) \in S_{\text{train}}} \mathbb{E} [(f(x) - y)^2], \quad (13)$$

$$(14)$$

supposed that the training set  $S_{\text{train}}$  and the test set  $S_{\text{test}}$  are drawn i.i.d. from the underlying distribution  $\mathcal{D}$ . We can actually see the training set as sampling  $N_{\text{train}}$  times from the distribution  $\mathcal{D}$  which gives us:

$$\mathbb{E}_{\text{train}}[\hat{\epsilon}_{\text{train}}(f)] = \frac{1}{N_{\text{train}}} \sum_{(x,y) \in S_{\text{train}}} \mathbb{E} [(f(x) - y)^2] \quad (15)$$

$$= \frac{1}{N_{\text{train}}} N_{\text{train}} \sum_{(x,y) \sim D} \mathbb{E} [(f(x) - y)^2] \quad (16)$$

$$= \epsilon(f). \quad (17)$$

And the same logic is used for the test error:

$$\mathbb{E}_{\text{test}}[\hat{\epsilon}_{\text{test}}(f)] = \mathbb{E} \left[ \frac{1}{N_{\text{test}}} \sum_{(x,y) \in S_{\text{test}}} ((f(x) - y)^2) \right] \quad (18)$$

$$= \frac{1}{N_{\text{test}}} \sum_{(x,y) \in S_{\text{test}}} \mathbb{E} [(f(x) - y)^2] \quad (19)$$

$$= \frac{1}{N_{\text{test}}} N_{\text{test}} \sum_{(x,y) \sim D} \mathbb{E} [(f(x) - y)^2] \quad (20)$$

$$= \epsilon(f). \quad (21)$$

Finally from the above derivation, we let  $f = \hat{f}$ , we have:

$$\mathbb{E}_{\text{test}}[\hat{\epsilon}_{\text{test}}(\hat{f})] = \epsilon(\hat{f}) \quad (22)$$

which the test error is an unbiased estimate of our true error for  $\hat{f}$ .

**b.**

**No**, the above equation will not be true with regards to the training set, which means  $\mathbb{E}_{\text{train}}[\hat{\epsilon}_{\text{train}}(\hat{f})] \neq \epsilon(\hat{f})$ . To be more specific, Eq. 16 will not hold as  $(f(x_i) - y_i)^2$ , or  $(x_i, y_i)$ , are not guaranteed as they are no longer independent to  $\hat{f}$  when  $(x_i, y_i) \in S_{\text{train}}$ .

**c.**

From Problem A.2.a we have:

$$\mathbb{E}_{\text{train}}[\hat{\epsilon}_{\text{train}}(f)] = \mathbb{E}_{\text{test}}[\hat{\epsilon}_{\text{test}}(f)], \forall f \in \mathcal{F}, \quad (23)$$

and from Problem A.2.c's hint we have:

$$\mathbb{E}_{\text{train,test}}[\hat{\epsilon}_{\text{test}}(\hat{f}_{\text{train}})] = \sum_{f \in \mathcal{F}} \mathbb{E}_{\text{test}}[\hat{\epsilon}_{\text{test}}(f)] \mathbb{P}_{\text{train}}(\hat{f}_{\text{train}} = f). \quad (24)$$

Together we can derive:

$$\mathbb{E}_{\text{train,test}}[\hat{\epsilon}_{\text{test}}(\hat{f}_{\text{train}})] = \sum_{f \in \mathcal{F}} \mathbb{E}_{\text{train}}[\hat{\epsilon}_{\text{train}}(f)] \mathbb{P}_{\text{train}}(\hat{f}_{\text{train}} = f) \quad (25)$$

$$= \mathbb{E}_{\text{train}}[\hat{\epsilon}_{\text{train}}(f)], \quad (26)$$

and in general we know that  $\mathbb{E}_{\text{train}}[\hat{\epsilon}_{\text{train}}(f)] \geq \mathbb{E}_{\text{train}}[\hat{\epsilon}_{\text{train}}(\hat{f})]$  as  $\hat{f}_{\text{train}}$  minimize the training error such that  $\hat{\epsilon}_{\text{train}}(\hat{f}_{\text{train}}) \leq \hat{\epsilon}_{\text{train}}(f)$  for all  $f \in \mathcal{F}$  which complete the last portion of the inequality equation.

## Problem B.1

**a.**

For some  $m \leq n$  such that  $n/m$  is an integer define the step function estimator:

$$\hat{f}_m(x) = \sum_{j=1}^{n/m} c_j \mathbf{1}\left\{x \in \left(\frac{(j-1)m}{n}, \frac{jm}{n}\right]\right\} \quad \text{where} \quad c_j = \frac{1}{m} \sum_{i=(j-1)m+1}^{jm} y_i. \quad (27)$$

By the bias-variance decomposition at some  $x_i$  we have

$$\mathbb{E} \left[ (\hat{f}_m(x_i) - f(x_i))^2 \right] = \underbrace{(\mathbb{E}[\hat{f}_m(x_i)] - f(x_i))^2}_{\text{Bias}^2(x_i)} + \underbrace{\mathbb{E} \left[ (\hat{f}_m(x_i) - \mathbb{E}[\hat{f}_m(x_i)])^2 \right]}_{\text{Variance}(x_i)} \quad (28)$$

So intuitively, for a very large  $m$ , it is expected to have high bias and low variance (e.g. we choose  $m = n$ , the model will predict the exact average of all samples). For a very small  $m$ , it is expected to have low bias and high variance (e.g. we choose  $m = 1$ , the model will predict the exact value if falls into the interval of the single sample).

**b.**

We first expand the summation of average bias-squared into:

$$\frac{1}{n} \sum_{i=1}^n (\mathbb{E}[\hat{f}_m(x_i)] - f(x_i))^2 = \frac{1}{n} \sum_{j=1}^{n/m} \sum_{i=jm-m+1}^{jm} (\mathbb{E}[\hat{f}_m(x_i)] - f(x_i))^2, \quad (29)$$

then from Eq. (27) we know that the estimator will predicts the average of the observations within each interval partitioned by  $m$  and  $j$  is the index of the corresponding partition which gives:

$$\mathbb{E}[\hat{f}_m(x_i)] = \mathbb{E}[c_j] = \mathbb{E} \left[ \frac{1}{m} \sum_{i=jm-m+1}^{jm} y_i \right] = \frac{1}{m} \sum_{i=jm-m+1}^{jm} \mathbb{E}[y_i] = \frac{1}{m} \sum_{i=jm-m+1}^{jm} f(x_i) = \bar{f}^{(j)}. \quad (30)$$

Combining Eq. (26) and (27), we have:

$$\frac{1}{n} \sum_{i=1}^n (\mathbb{E}[\hat{f}_m(x_i)] - f(x_i))^2 = \frac{1}{n} \sum_{j=1}^{n/m} \sum_{i=(j-1)m+1}^{jm} (\bar{f}^{(j)} - f(x_i))^2 \quad (31)$$

c.

We first expand the summation of average variance into:

$$\mathbb{E} \left[ \frac{1}{n} \sum_{i=1}^n (\hat{f}_m(x_i) - \mathbb{E}[\hat{f}_m(x_i)])^2 \right] = \mathbb{E} \left[ \frac{1}{n} \sum_{j=1}^{n/m} \sum_{i=jm-m+1}^{jm} (\hat{f}_m(x_i) - \mathbb{E}[\hat{f}_m(x_i)])^2 \right] \quad (32)$$

then from Eq. (30) we have  $\mathbb{E}[\hat{f}_m(x_i)] = \mathbb{E}[c_j]$  and the fact that  $\hat{f}_m(x_i)$  will be identical across all if the  $x_i$  terms in this given  $j$ -th partition:

$$\mathbb{E} \left[ \frac{1}{n} \sum_{j=1}^{n/m} \sum_{i=jm-m+1}^{jm} (\hat{f}_m(x_i) - \mathbb{E}[\hat{f}_m(x_i)])^2 \right] = \frac{1}{n} \sum_{j=1}^{n/m} \sum_{i=jm-m+1}^{jm} \mathbb{E} \left[ (\hat{f}_m(x_i) - \mathbb{E}[\hat{f}_m(x_i)])^2 \right] \quad (33)$$

$$= \frac{1}{n} \sum_{j=1}^{n/m} \sum_{i=jm-m+1}^{jm} \mathbb{E} [(c_j - \mathbb{E}[c_j])^2] \quad (34)$$

$$= \frac{1}{n} \sum_{j=1}^{n/m} m \mathbb{E} [(c_j - \bar{f}^{(j)})^2] \quad (35)$$

$$= \frac{\sigma^2}{m}. \quad (36)$$

Then

$$\mathbb{E} \left[ \frac{1}{n} \sum_{j=1}^{n/m} \sum_{i=jm-m+1}^{jm} (\hat{f}_m(x_i) - \mathbb{E}[\hat{f}_m(x_i)])^2 \right] = \frac{1}{n} \sum_{j=1}^{n/m} \sum_{i=jm-m+1}^{jm} \mathbb{E} \left[ (\hat{f}_m(x_i) - \mathbb{E}[\hat{f}_m(x_i)])^2 \right] \quad (37)$$

$$(38)$$

d.

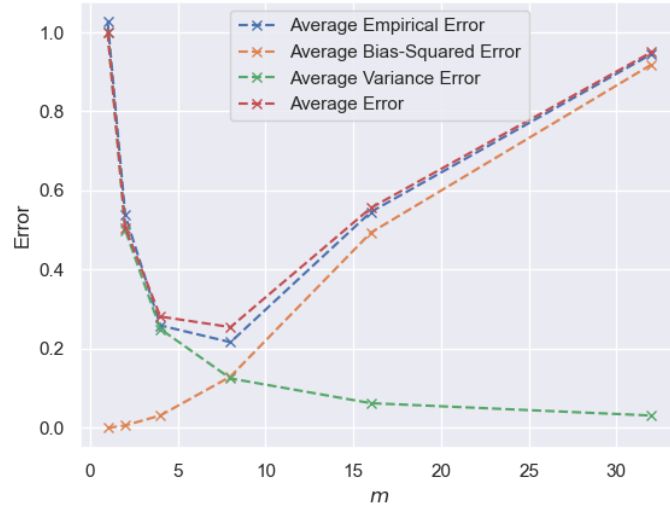


Figure 1: Plot of different error under different  $m$  value.

```

1 import numpy as np
2 import matplotlib.pyplot as plt
3 import seaborn as sns
4
5 sns.set()
6
7 n = 256
8 sigma = 1
9 M = [1, 2, 4, 8, 16, 32]
10 eps = np.random.normal(0, sigma, n)
11
12 X = np.arange(1, n+1) / n
13 Y = 4 * np.sin(np.pi * X) * np.cos(6 * np.pi * X * X) + eps
14 Y_true = 4 * np.sin(np.pi * X) * np.cos(6 * np.pi * X * X)
15
16 # average empirical error
17 average_empirical_error = []
18 for m in M:
19     f_m = np.array([np.mean(Y[j*m-m:j*m]) for j in np.arange(1, n//m + 1)])
20     f_hat = np.array([f_m[int((i-1)//m)] for i in np.arange(1, n+1)])
21     err = np.mean((f_hat - Y_true)**2)
22     average_empirical_error.append(err)
23
24 # average bias-squared error
25 average_bias_squared_error = []
26 for m in M:
27     f_j = np.array([np.mean(Y_true[j * m - m:j * m]) for j in np.arange(1, n // m + 1)])
28     f_hat = np.array([f_j[int((i-1) // m)] for i in np.arange(1, n + 1)])
29     err = np.mean((f_hat - Y_true) ** 2)
30     average_bias_squared_error.append(err)
31
32 # average variance error
33 average_variance_error = []
34 for m in M:
35     average_variance_error.append(sigma**2/m)
36
37 plt.plot(M, average_empirical_error, 'x--', label='Average Empirical Error')
38 plt.plot(M, average_bias_squared_error, 'x--', label='Average Bias-Squared Error')
39 plt.plot(M, average_variance_error, 'x--', label='Average Variance Error')
40 plt.plot(M, np.array(average_bias_squared_error) + np.array(average_variance_error), 'x--', label='Average Error')
41 plt.legend()
42 plt.xlabel('$x$')
43 plt.ylabel('Error')
44 plt.savefig('B-1-4.png')

```

Figure 2: Code used to generate Fig.8.

e.

From B.1.b, we have the average bias-squared and we can derived:

$$\frac{1}{n} \sum_{i=1}^n (\hat{f}^{(j)} - f(x_i))^2 \leq \frac{1}{n} \sum_{i=1}^n (\min f(x_i) - f(x_i))^2 \quad (39)$$

$$\leq \left(\frac{L}{n}\right)^2 (\operatorname{argmin} f(x_i) - i)^2 \quad (40)$$

$$\leq \frac{L^2}{n^2} m^2 \sim O\left(\frac{L^2 m^2}{n^2}\right) \quad (41)$$



## Problem A.4

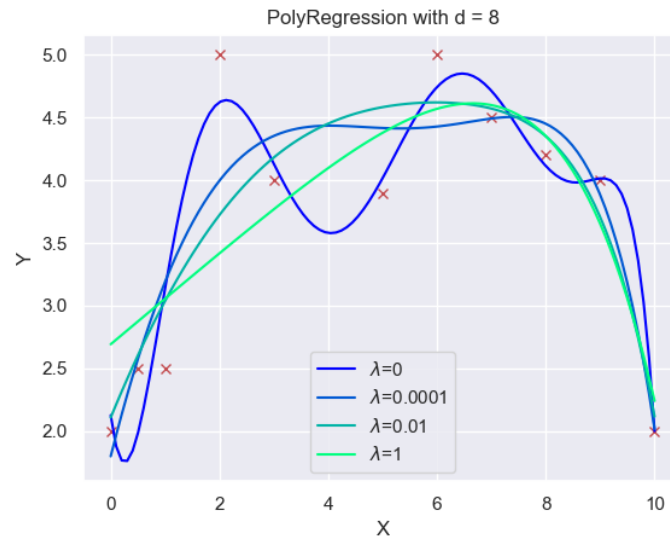


Figure 3: Plot of PolyRegression with  $d = 8$  under different  $\lambda$ .

```

13 class PolynomialRegression:
14
15     def __init__(self, degree=1, reg_lambda=1E-8):
16         """
17         """
18         self.degree = degree
19         self.reg_lambda = reg_lambda
20         self.theta = None
21         self.mean = None
22         self.std = None
23
24
25     def polyfeatures(self, X, degree):
26         """
27         """
28         poly_feat_X = X
29         for i in range(1, self.degree):
30             poly_feat_X = np.c_[poly_feat_X, X**i]
31         assert poly_feat_X.shape == (X.shape[0], self.degree)
32         return poly_feat_X
33
34
35     def fit(self, X, y):
36         """
37         """
38         n = len(X)
39
40         X_ = self.polyfeatures(X, self.degree)
41
42         # standardize
43         self.mean = X_.mean(axis=0)
44         self.std = X_.std(axis=0)
45         X_ = (X_ - self.mean) / self.std
46         X_ = np.c_[np.ones([n, 1]), X_]
47
48         n, d = X_.shape
49
50         reg_matrix = self.reg_lambda * np.eye(d)
51         reg_matrix[0, 0] = 0
52
53         self.theta = np.linalg.pinv(X_.T.dot(X_) + reg_matrix).dot(X_.T).dot(y)
54
55
56     def predict(self, X):
57         """
58         """
59         n, _ = X.shape
60         X_ = self.polyfeatures(X, self.degree)
61         X_ = (X_ - self.mean) / self.std
62         X_ = np.c_[np.ones([n, 1]), X_]
63
64         return X_.dot(self.theta)

```

Figure 4: Screenshot of the polyreg.py.

```

17 if __name__ == "__main__":
18     # Load the data
19     # Load the data
20     # Load the data
21     # Load the data
22     # Load the data
23     filePath = "data/polydata.dat"
24     file = open(filePath, 'r')
25     allData = np.loadtxt(file, delimiter=',')
26     X = allData[:, 0]
27     y = allData[:, 1]
28
29     # regression with degree = d
30     d = 8
31     model = PolynomialRegression(degree=d, reg_lambda=0)
32     model.fit(X, y)
33
34     # output predictions
35     xpoints = np.linspace(np.min(X), np.max(X), 100).reshape(-1, 1)
36     ypoints = model.predict(xpoints)
37
38     # plot curve
39     plt.figure()
40     plt.plot(X, y, 'rx')
41     plt.plot(xpoints, ypoints, color=cmap(0), label='$lambda$=0')
42
43     model = PolynomialRegression(degree=d, reg_lambda=1e-4)
44     model.fit(X, y)
45     ypoints = model.predict(xpoints)
46     plt.plot(xpoints, ypoints, color=cmap(0.35), label='$lambda$=0.0001')
47
48     model = PolynomialRegression(degree=d, reg_lambda=1e-2)
49     model.fit(X, y)
50     ypoints = model.predict(xpoints)
51     plt.plot(xpoints, ypoints, color=cmap(0.7), label='$lambda$=0.01')
52
53     model = PolynomialRegression(degree=d, reg_lambda=1)
54     model.fit(X, y)
55     ypoints = model.predict(xpoints)
56     plt.plot(xpoints, ypoints, color=cmap(1.2), label='$lambda$=1')
57
58     plt.title('PolyRegression with d = ' + str(d))
59     plt.xlabel('X')
60     plt.ylabel('Y')
61     plt.legend()
62     plt.savefig('A_4.png')

```

Figure 5: Screenshot of the test\_polyreg\_univariate.py.

## Problem A.5

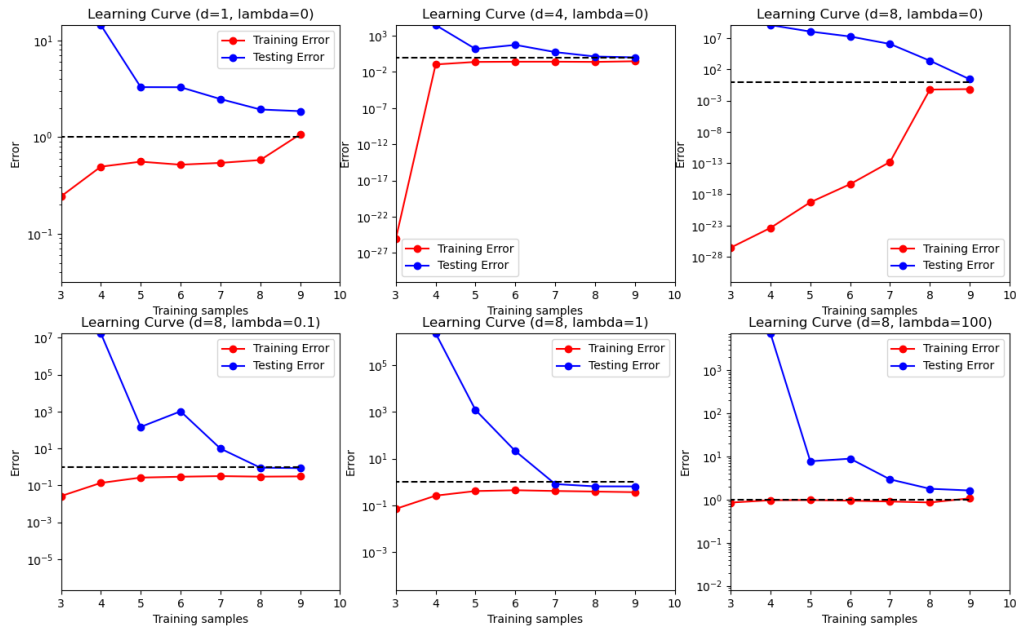


Figure 6: Plot of PolyRegression with  $d = 8$  under different  $\lambda$ .

```

96 def learningCurve(Xtrain, Ytrain, Xtest, Ytest, reg_lambda, degree):
97     """Compute learning curve..."""
117
118     n = len(Xtrain)
119
120     errorTrain = np.zeros(n)
121     errorTest = np.zeros(n)
122
123     #TODO -- complete rest of method; errorTrain and errorTest are already the correct shape
124     assert Xtrain.shape == Ytrain.shape
125     assert Xtest.shape == Ytest.shape
126
127     model = PolynomialRegression(degree=degree, reg_lambda=reg_lambda)
128     for i in range(1, n):
129         X = Xtrain[0:i+1]
130         y = Ytrain[0:i+1]
131
132         model.fit(X, y)
133
134         ytrain = model.predict(X)
135         ytest = model.predict(Xtest)
136
137         errorTrain[i] = np.mean((ytrain-y)**2)
138         errorTest[i] = np.mean((ytest-Ytest)**2)
139
140     return errorTrain, errorTest

```

Figure 7: Screenshot of the polyreg.py.

## Problem A.6

a.

In this problem we will choose a linear classifier to minimize the regularized least squares objective:

$$\widehat{W} = \operatorname{argmin}_{W \in \mathbb{R}^{d \times k}} \sum_{i=1}^n \|W^T x_i - y_i\|_2^2 + \lambda \|W\|_F^2, \quad (42)$$

and from the problem we can deduce the term we ought to minimize:

$$\sum_{i=1}^n \|W^T x_i - y_i\|_2^2 + \lambda \|W\|_F^2 = \sum_{j=1}^k \left[ \sum_{i=1}^n (e_j^T W^T x_i - e_j^T y_i)^2 + \lambda \|W e_j\|^2 \right] \quad (43)$$

$$= \sum_{j=1}^k \left[ \sum_{i=1}^n (w_j^T x_i - e_j^T y_i)^2 + \lambda \|w_j\|^2 \right] \quad (44)$$

$$= \sum_{j=1}^k [\|X w_j - Y e_j\|^2 + \lambda \|w_j\|^2], \quad (45)$$

where  $X = [x_1 \ \dots \ x_n]^T \in \mathbb{R}^{n \times d}$  and  $Y = [y_1 \ \dots \ y_n]^T \in \mathbb{R}^{n \times k}$ . So we take the derivative of the term and set the value to zero to find the minimizer of the regularized least squares objective:

$$\frac{d}{dw_j} \sum_{j=1}^k [\|X w_j - Y e_j\|^2 + \lambda \|w_j\|^2] = 0 \quad (46)$$

$$\sum_{j=1}^k [2X^T (X w_j - Y e_j) + 2\lambda w_j] = 0 \quad (47)$$

$$\sum_{j=1}^k [X^T X w_j - X^T Y e_j + \lambda w_j] = 0 \quad (48)$$

$$\sum_{j=1}^k X^T X w_j + \lambda w_j = \sum_{j=1}^k X^T Y e_j \quad (49)$$

$$\sum_{j=1}^k w_j = \sum_{j=1}^k (X^T X + \lambda I)^{-1} X^T Y e_j \quad (50)$$

Finally we know that  $\widehat{W}$  can be view as a matrix formed by  $w_j$  and the fact that  $Y e_j = y_j$ , so in the last equation, the summation sum over  $k$  making that the minimizer for this problem:

$$\widehat{W} = \operatorname{argmin}_{W \in \mathbb{R}^{d \times k}} \sum_{i=1}^n \|W^T x_i - y_i\|_2^2 + \lambda \|W\|_F^2 = (X^T X + \lambda I)^{-1} X^T Y \quad (51)$$

b.

After we train  $\widehat{W}$  on the MNIST training data with  $\lambda = 0.0001$ , the training error is 14.81% and testing error is 14.66%.

```
1  import ...
2
3
4
5  def load_dataset():
6      mndata = MNIST('./data/')
7      X_train, Y_train = map(np.array, mndata.load_training())
8      X_test, Y_test = map(np.array, mndata.load_testing())
9      X_train = X_train / 255.0
10     X_test = X_test / 255.0
11
12     return X_train, Y_train, X_test, Y_test
13
14
15  def train(X, Y, reg_lambda=1e-4):
16      I = np.eye(X.shape[1])
17      # https://numpy.org/doc/stable/reference/generated/numpy.linalg.solve.html
18      W_hat = np.linalg.solve(X.T.dot(X) + reg_lambda * I, X.T.dot(Y))
19      assert W_hat.shape == (X.shape[1], Y.shape[1])
20      return W_hat
21
22
23  def predict(W, X):
24      preds = np.argmax(W.T.dot(X.T), axis=0)
25      assert preds.shape == (X.shape[0],)
26      return preds
27
28
29  def main():
30      print("Loading MNIST dataset from source")
31      X_train, Y_train, X_test, Y_test = load_dataset()
32      print("Finished loading MNIST dataset from source")
33
34      # one-hot encoding the 10 classes in MNIST
35      # Y_train = np.eye(10)[Y_train]
36      # Y_test = np.eye(10)[Y_test]
37
38      W_hat = train(X_train, np.eye(10)[Y_train], 1e-4)
39
40      Y_train_preds = predict(W_hat, X_train)
41      Y_test_preds = predict(W_hat, X_test)
42
43      train_acc = np.count_nonzero(Y_train == Y_train_preds) / Y_train.shape[0]
44      print("Training Error:", 1 - train_acc)
45      test_acc = np.count_nonzero(Y_test == Y_test_preds) / Y_test.shape[0]
46      print("Testing Error:", 1 - test_acc)
```

Figure 8: Screenshot of the `train_MNIST.py`.

## Problem B.2

a.



Figure 9: Plot of training and validation error under different  $p$  value.

```

52  # Problem B.2.a
53  def main():
54      print("Loading MNIST dataset from source")
55      X_train, Y_train, X_test, Y_test = load_dataset()
56      print("Finished loading MNIST dataset from source")
57
58      n = X_train.shape[0]
59
60      X, Y = shuffle(X_train, Y_train)
61      X_train = X[:int(0.8*n)]
62      Y_train = Y[:int(0.8*n)]
63      X_val = X[int(0.8*n):]
64      Y_val = Y[int(0.8*n):]
65
66      p_list = np.arange(200, 4000, 200)
67
68      train_err = []
69      val_err = []
70
71      for p in p_list:
72
73          G = np.random.normal(0, 0.1, (p, X_train.shape[1]))
74          b = np.random.uniform(0, 2*np.pi, (p,))
75          X_train_transform = np.cos(X_train.dot(G.T) + b)
76          W_hat = train(X_train_transform, np.eye(10)[Y_train], 1e-4)
77          X_val_transform = np.cos(X_val.dot(G.T) + b)
78          Y_train_preds = predict(W_hat, X_train_transform)
79          Y_val_preds = predict(W_hat, X_val_transform)
80
81          print("p=", p)
82          train_acc = np.count_nonzero(Y_train == Y_train_preds) / Y_train.shape[0]
83          train_err.append(1-train_acc)
84          print("Training Error:", 1-train_acc)
85          val_acc = np.count_nonzero(Y_val == Y_val_preds) / Y_val.shape[0]
86          val_err.append(1-val_acc)
87          print("Validation Error:", 1-val_acc)
88
89      plt.figure()
90      plt.plot(p_list, train_err, 'b--x', label='Training Error')
91      plt.plot(p_list, val_err, 'r--x', label='Validation Error')
92      plt.xlabel('$p$')
93      plt.ylabel('Error')
94      plt.legend()
95      plt.savefig('B-2-a.png')

```

Figure 10: Code used to generate Fig.11.

b.

We choose  $p = 4000$  with the training error being 0.017833, validation error being 0.0340, and testing error  $\hat{\epsilon}_{\text{test}}(\hat{f})$  being 0.03020. From the Lemma, we have:

$$\mathbb{P}\left(\left|\left(\frac{1}{m}\sum_{i=1}^m X_i\right) - \mu\right| \geq \sqrt{\frac{(b-a)^2 \log(2/\delta)}{2m}}\right) \leq \delta,$$

where  $a = 0$ ,  $b = 1$ ,  $\sigma = 0.05$  and  $m = 10000$  which let  $\sqrt{\frac{(b-a)^2 \log(2/\delta)}{2m}} \sim 0.0136$  so that the 95% confidence interval of the testing error will be  $0.0302 \pm 0.0136$ .

```

120 g = np.random.normal(0, 0.1, (p, X_train.shape[1]))
121 b = np.random.uniform(0, 2*np.pi, (p,))
122
123 X_train_transform = np.cos(X_train.dot(g.T) + b)
124 X_val_transform = np.cos(X_val.dot(g.T) + b)
125 X_test_transform = np.cos(X_test.dot(g.T) + b)
126
127 W_hat = train(X_train_transform, np.eye(10)[Y_train], 1e-4)
128
129 Y_train_preds = predict(W_hat, X_train_transform)
130 Y_val_preds = predict(W_hat, X_val_transform)
131 Y_test_preds = predict(W_hat, X_test_transform)
132
133 print("p=", p)
134 train_acc = np.count_nonzero(Y_train == Y_train_preds) / Y_train.shape[0]
135 train_err.append(1-train_acc)
136 print("Training Error:", 1-train_acc)
137 val_acc = np.count_nonzero(Y_val == Y_val_preds) / Y_val.shape[0]
138 val_err.append(1-val_acc)
139 print("Validation Error:", 1-val_acc)
140 test_acc = np.count_nonzero(Y_test == Y_test_preds) / Y_test.shape[0]
141 test_err.append(1 - test_acc)
142 print("Test Error:", 1 - test_acc)

```

Figure 11: Code used to get the test error.