

CSE446 HW1

Brian Kang[‡]

Conceptual Questions

A.0

- In your own words, describe what bias and variance are? What is bias-variance tradeoff?
 - Bias: How far we are from fitting out data to our ideal estimator. Or, the error introduced by approximating a real model by a much simpler model.
 - Variance: How much our estimator will change after redrawing a sample and fitting the data again to get the estimator. Or, how much \hat{f} would change if we estimated a model using a different training set.
 - Bias-Variance Tradeoff: The tradeoff between bias and variance within the reducible error part of a model (test MSE) with respect to model flexibility/complexity.
- What happens to bias and variance when the model complexity increases/decreases?
 - When the model complexity increases, bias decreases (our model may be able to capture the required complexity of the true model) and variance increases (noise is introduced by fitting this complex model to sample data). When the model complexity decreases, bias increases and variance decreases. (The irreducible error stays the same regardless of complexity.)
- True or False: The bias of a model increases as the amount of training data available increases.
 - **False.** $Bias[\hat{f}(x)] = \mathbb{E}[\hat{f}(x)] - f(x)$ so training set size will not affect bias.
- True or False: The variance of a model decreases as the amount of training data available increases.
 - **True.** As training data increases, the model will be able to generalize better and variance will decrease.
- True or False: A learning algorithm will generalize better if we use less features to represent our data.
 - **Yes and No.** When facing a high variance issue, decreasing the number of features helps, but for high bias issues increasing them will help.
- To get better generalization, should we use the train set or the test set to tune our hyperparameters?
 - **Training set.** The training set will be split into training and validation sets.
- True or False: The training error of a function on the training set provides an overestimate of the true error of that function.
 - **False.** Underestimation. The function is chosen to minimize the training error. So, fitting this function on true data will return larger error because the function was not trained to minimize the true error.

*Collaborated with Cindy Wu

[†]References: An Introduction to Statistical Learning (James, Witten, Hastie, Tibshirani)

A.1

a.

$$\begin{aligned}
\hat{\theta}_{MLE} &= \operatorname{argmax}_{\theta} P(x_1, x_2, x_3, x_4, x_5 \mid \lambda = \theta) \\
&= \operatorname{argmax}_{\lambda} \prod_{i=1}^5 \frac{e^{-\lambda} \lambda^{x_i}}{x_i!} \\
&= \operatorname{argmax}_{\lambda} (e^{-\lambda})^5 \prod_{i=1}^5 \frac{\lambda^{x_i}}{x_i!} \\
&= \operatorname{argmax}_{\lambda} \log(\bullet) \\
&= \operatorname{argmax}_{\lambda} 5 \log(e^{-\lambda}) + \sum_{i=1}^5 (\log(\lambda^{x_i}) - \log(x_i!)) \\
&= \operatorname{argmax}_{\lambda} -5\lambda + \sum_{i=1}^5 (x_i \log(\lambda) - \log(x_i!)) , \text{ then} \\
\frac{d}{d\lambda}[\bullet] &= -5 + \sum_{i=1}^5 \left(\frac{x_i}{\lambda} - 0\right) \\
&= -5 + \frac{x_1 + \dots + x_5}{\lambda} = 0 , \text{ then} \\
5\lambda &= x_1 + \dots + x_5 , \\
\boxed{\lambda = \hat{\theta}_{MLE} = \frac{x_1 + \dots + x_5}{5}}
\end{aligned}$$

b. Based on the above work, by symmetry,

$$\boxed{\hat{\theta}_{MLE} = \frac{x_1 + \dots + x_6}{6}}$$

c.

$$\begin{aligned}
\lambda \text{ after 5 games} &= \frac{2 + 0 + 1 + 1 + 2}{5} = \boxed{\frac{6}{5}} \\
\lambda \text{ after 6 games} &= \frac{2 + 0 + 1 + 1 + 2 + 4}{6} = \frac{10}{6} = \boxed{\frac{5}{3}}
\end{aligned}$$

A.2

$$\begin{aligned}
\hat{\theta}_{MLE} &= \operatorname{argmax}_{\theta} P(x_1, \dots, x_n \mid \theta) \\
&= \operatorname{argmax}_{\theta} \prod_{i=1}^n \frac{1}{\theta} = \operatorname{argmax}_{\theta} \theta^{-n}
\end{aligned}$$

MLE is the value of θ that maximizes $\frac{1}{\theta^n}$. Since $\frac{1}{\theta^n}$ monotonically decreases in θ , we want the smallest θ such that $\theta \geq x_i$ for $i = 1, \dots, n$. Therefore, it has to be that $\theta = \max(x_1, \dots, x_n)$. Then,

$$\hat{\theta}_{MLE} = \max(X_1, \dots, X_n)$$

Do note that this MLE is not a great estimator for θ because the MLE will ALWAYS underestimate θ , since we know that $Prob(\max(X_1, \dots, X_n) < \theta) = 1$.

A.3

a. *Proof.*

$$\begin{aligned}
\mathbb{E}_{train}[\hat{\epsilon}_{train}(f)] &= \mathbb{E}_{train} \left[\frac{1}{N_{train}} \sum_{(x,y) \in S_{train}} (f(x) - y)^2 \right] \\
&= \frac{1}{N_{train}} \sum_{(x,y) \in S_{train}} \mathbb{E}_{train} [(f(x) - y)^2] \\
&= \frac{1}{N_{train}} * N_{train} * \mathbb{E}_{train} [(f(x) - y)^2] \\
&= \mathbb{E}_{(x,y) \in S_{test} \subset D} [(f(x) - y)^2] \\
&= \epsilon(f)
\end{aligned}$$

The last step is justified because the i.i.d. random samples are drawn from the underlying distribution D . Likewise,

$$\begin{aligned}
\mathbb{E}_{test}[\hat{\epsilon}_{test}(f)] &= \frac{1}{N_{test}} * N_{test} * \mathbb{E}_{test} [(f(x) - y)^2] \\
&= \mathbb{E}_{(x,y) \in S_{test} \subset D} [(f(x) - y)^2] \\
&= \epsilon(f)
\end{aligned}$$

The last step is also justified because the i.i.d. random samples are drawn from the underlying true distribution D .

Therefore, $\mathbb{E}_{train}[\hat{\epsilon}_{train}(f)] = \mathbb{E}_{test}[\hat{\epsilon}_{test}(f)] = \epsilon(f)$.

Using similar reasoning, if $\epsilon(\hat{f}) = \mathbb{E}_D[(\hat{f}(x) - y)^2]$, where \hat{f} is the function trained using the training set, then,

$$\begin{aligned}
\mathbb{E}_{test}[\hat{\epsilon}_{test}(\hat{f})] &= \frac{1}{N_{test}} * N_{test} * \mathbb{E}_{test} [(\hat{f}(x) - y)^2] \\
&= \mathbb{E}_{(x,y) \in S_{test} \subset D} [(\hat{f}(x) - y)^2] \\
&= \epsilon(\hat{f})
\end{aligned}$$

The last step is justified because the i.i.d. samples are drawn from the underlying distribution D , and the training and test sets are mutually exclusive. Keeping the test set as random, the expectation over the test error (using the trained function) will be an unbiased estimator for the true error for \hat{f} , i.e., $\epsilon(\hat{f})$. \square

b. *Proof.* Unlike above, $\mathbb{E}_{train}[\hat{\epsilon}_{train}(\hat{f})] \neq \mathbb{E}_{train}[\epsilon(\hat{f})]$. This is the fact that train error is optimistically biased because it is evaluated on the data it trained on. The \hat{f} function is already trained on the test data such that $\hat{\epsilon}_{train}$ is minimized. Thus, the expectation train of this minimized error cannot be equal to the expectation train of the true error calculated from the trained function. \square

c. *Proof.* First, we are given that $\exists \hat{f}_{train}$ such that, $\hat{\epsilon}_{train}(\hat{f}_{train}) \leq \hat{\epsilon}_{train}(f) \forall f \in F$. Then,

$$\begin{aligned}
\mathbb{E}_{train}[\hat{\epsilon}_{train}(\hat{f}_{train})] &\leq \mathbb{E}_{train}[\hat{\epsilon}_{train}(f)] \forall f \in F \\
&= \epsilon(f) \text{ by (a).}
\end{aligned}$$

Next, we are given information about $\mathbb{E}_{train, test}[\hat{\epsilon}_{test}(\hat{f}_{train})]$,

$$\begin{aligned}
\mathbb{E}_{train, test}[\hat{\epsilon}_{test}(\hat{f}_{train})] &= \sum_{f \in F} \mathbb{E}_{train, test}[\hat{\epsilon}_{test}(f) * \vec{1}\{\hat{f}_{train} = f\}] \\
&= \sum_{f \in F} \mathbb{E}_{test}[\hat{\epsilon}_{test}(f)] * \mathbb{E}_{train}[\vec{1}\{\hat{f}_{train} = f\}] \text{ by indep. between } S_{train} \text{ and } S_{test} \\
&= \sum_{f \in F} \mathbb{E}_{test}[\hat{\epsilon}_{test}(f)] * Prob_{train}(\hat{f}_{train} = f) \\
&= \sum_{f \in F} \epsilon(f) * Prob_{train}(\hat{f}_{train} = f) \text{ by (a)} \\
&= \epsilon(\hat{f}_{train}) * 1
\end{aligned}$$

We know that $\forall f \in F, \epsilon(f) \leq \epsilon(\hat{f}_{train})$.

Therefore, $\mathbb{E}_{train}[\hat{\epsilon}_{train}(\hat{f}_{train})] \leq \mathbb{E}_{train}[\hat{\epsilon}_{train}(f)] = \epsilon(f) \leq \epsilon(\hat{f}_{train}) = \mathbb{E}_{train, test}[\hat{\epsilon}_{test}(\hat{f}_{train})]$. \square

Technical Questions

A.4

```

1  '''
2      Template for polynomial regression
3      AUTHOR Eric Eaton, Xiaoxiang Hu
4  '''
5
6  import numpy as np
7
8
9  # -----
10 # Class PolynomialRegression
11 # -----
12
13 class PolynomialRegression:
14
15     def __init__(self, degree=1, reg_lambda=1E-8):
16         """
17         Constructor
18         """
19         self.regLambda = reg_lambda
20         self.degree = degree
21         self.theta = None
22         self.mu = None
23         self.sd = None
24
25     def polyfeatures(self, X, degree):
26         """
27         Expands the given X into an n * d array of polynomial features of
28         degree d.
29
30         Returns:
31             A n-by-d numpy array, with each row comprising of
32             X, X * X, X ** 3, ... up to the dth power of X.
33             Note that the returned matrix will not include the zero-th power.
34
35         Arguments:
36             X is an n-by-1 column numpy array
37             degree is a positive integer
38         """
39         n, d = X.shape # get the row and column dimension of X
40         X_ = X # save X in different variable "X_"
41         for i in range(2, degree+1): # because last index exclusive and we want the d'th
degree

```

```

42         expX = np.power(X, i) # make d'th power of X
43         X_ = np.concatenate((X_, expX), axis=1) # add d'th power column to the right
44     return X_
45
46 def fit(self, X, y):
47     """
48     Trains the model
49     Arguments:
50         X is a n-by-1 array
51         y is an n-by-1 array
52     Returns:
53         No return value
54     Note:
55         You need to apply polynomial expansion and scaling
56         at first
57     """
58     n = len(X)
59     X = self.polyfeatures(X, self.degree) # poly expand
60     xbar = np.mean(X, axis=0) # get mean of columns
61     sdhat = np.std(X, axis=0) # get std dev of columns
62     X = (X-xbar)/sdhat # standardize
63     X_ = np.c_[np.ones([n, 1]), X] # add 1s column
64     n, d = X_.shape # get dim of new data matrix
65     d = d-1 # remove 1 for the extra column of ones we added to get the original num
        features
66
67     # construct reg matrix
68     reg_matrix = self.regLambda * np.eye(d + 1)
69     reg_matrix[0, 0] = 0
70     # analytical solution (X'X + regMatrix)^-1 X' y
71     self.theta = np.linalg.pinv(X_.T.dot(X_) + reg_matrix).dot(X_.T).dot(y)
72     self.mu = xbar
73     self.sd = sdhat
74
75 def predict(self, X):
76     """
77     Use the trained model to predict values for each instance in X
78     Arguments:
79         X is a n-by-1 numpy array
80     Returns:
81         an n-by-1 numpy array of the predictions
82     """
83     n = len(X)
84     X = self.polyfeatures(X, self.degree) # poly expand
85     X = (X - self.mu) / self.sd # standardize again
86     X_ = np.c_[np.ones([n, 1]), X] # add 1s column
87     return X_.dot(self.theta) # predict the fitted values
88
89
90 #-----
91 # End of Class PolynomialRegression
92 #-----
93
94
95
96 def learningCurve(Xtrain, Ytrain, Xtest, Ytest, reg_lambda, degree):
97     """
98     Compute learning curve
99
100    Arguments:
101        Xtrain — Training X, n-by-1 matrix
102        Ytrain — Training y, n-by-1 matrix
103        Xtest — Testing X, m-by-1 matrix
104        Ytest — Testing Y, m-by-1 matrix
105        regLambda — regularization factor
106        degree — polynomial degree
107
108    Returns:
109        errorTrain — errorTrain[i] is the training accuracy using

```

```

110     model trained by Xtrain[0:(i+1)]
111     errorTest — errorTrain[i] is the testing accuracy using
112     model trained by Xtrain[0:(i+1)]
113
114 Note:
115     errorTrain[0:1] and errorTest[0:1] won't actually matter, since we start displaying
116     the learning curve at n = 2 (or higher)
117     """
118
119 n = len(Xtrain)
120 errorTrain = np.zeros(n)
121 errorTest = np.zeros(n)
122 glm = PolynomialRegression(degree, reg_lambda) # get polynom model
123 for i in range(2, n+1):
124     glm.fit(Xtrain[0:(i+1)], Ytrain[0:(i+1)]) # fit the training data
125     errorTrain[i-1] = np.mean((glm.predict(Xtrain[0:(i+1)]) - Ytrain[0:(i+1)])**2)
126     errorTest[i-1] = np.mean((glm.predict(Xtest[0:(i+1)]) - Ytest[0:(i+1)])**2)
127 return errorTrain, errorTest

```

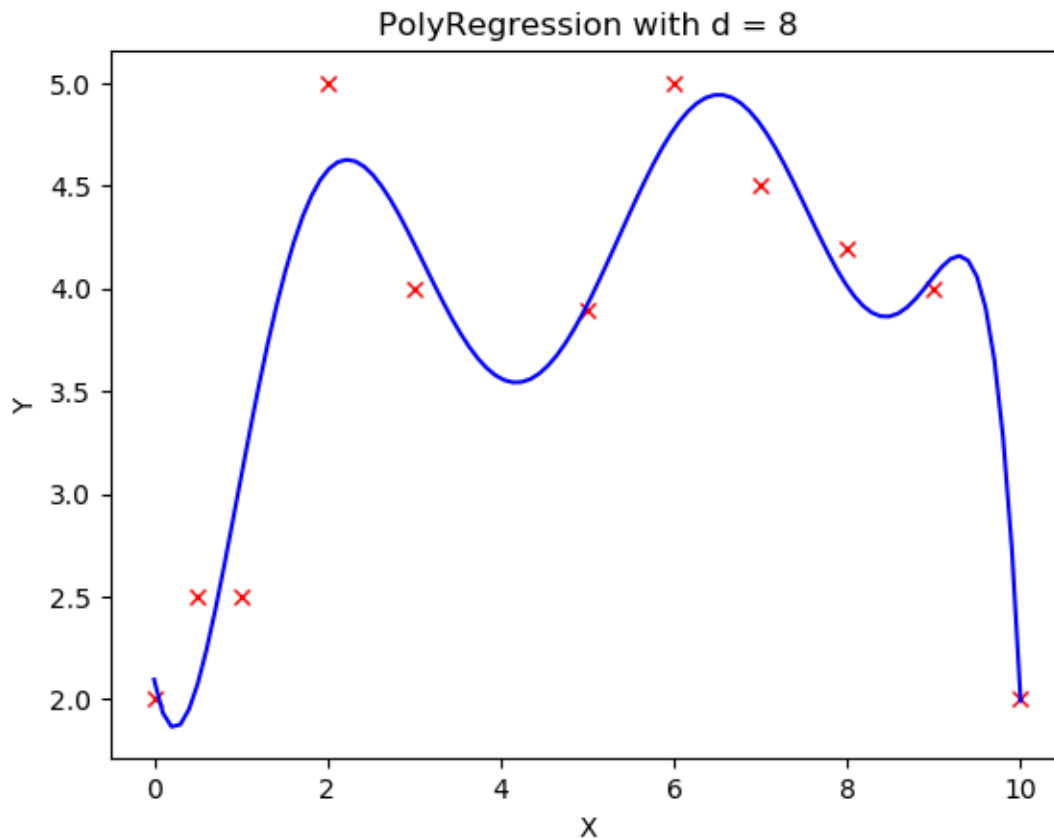


Figure 1: PolyReg Fit

A.5

Code is above with polyreg.py

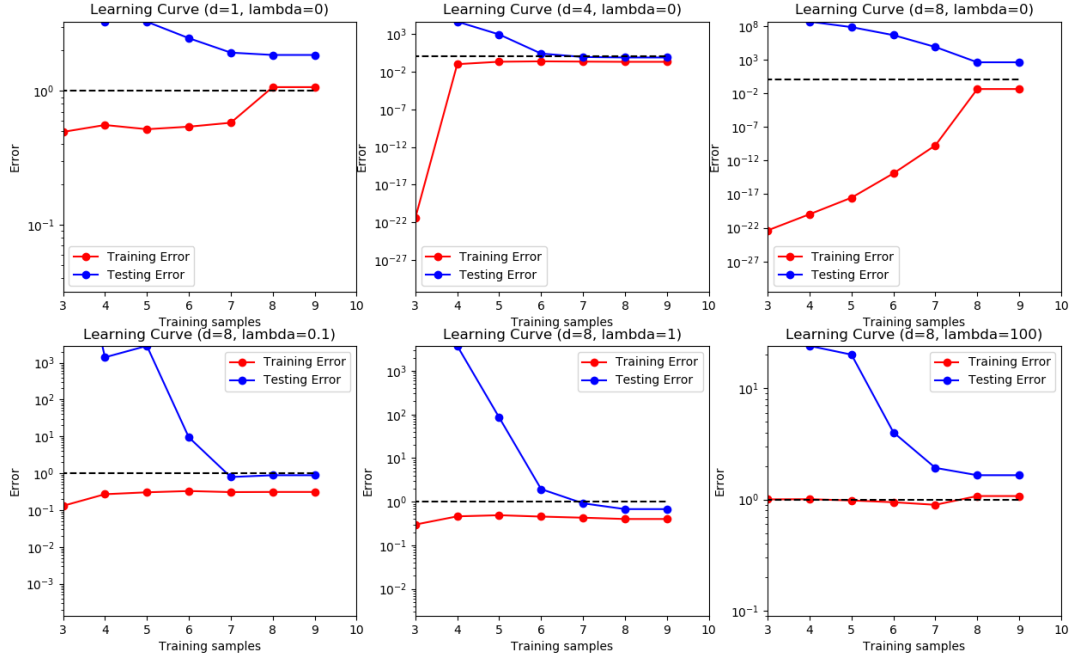


Figure 2: Learning Curves

A.6

a. *Proof.*

$$\begin{aligned}
 \hat{w} &= \underset{w}{\operatorname{argmin}} \sum_{j=0}^k [||Xw_j - Ye_j||^2 + \lambda ||w_j||^2] \\
 &= \underset{w}{\operatorname{argmin}} \sum_{j=0}^k ||Xw_j - Ye_j||^2 + \lambda \sum_{j=0}^k ||w_j||^2 \\
 &= \underset{w}{\operatorname{argmin}} \sum_{j=0}^k (Xw_j - Ye_j)^2 + \lambda \sum_{j=0}^k ||w_j||^2
 \end{aligned}$$

$$\begin{aligned}
 \nabla_w(\bullet) &= \begin{bmatrix} 2(Xw_0 - Ye_0)^T X + 2\lambda w_0 \\ \vdots \\ 2(Xw_k - Ye_k)^T X + 2\lambda w_k \end{bmatrix} \\
 &= \begin{bmatrix} 2X^T X w_0 - 2X^T Y e_0 + 2\lambda w_0 \\ \vdots \\ 2X^T X w_k - 2X^T Y e_k + 2\lambda w_k \end{bmatrix} = \begin{bmatrix} 0 \\ \vdots \\ 0 \end{bmatrix}
 \end{aligned}$$

Now rearrange terms and solve for \hat{w} .

$$\begin{bmatrix} 2X^T X \hat{w}_0 + 2\lambda \hat{w}_0 \\ \vdots \\ 2X^T X \hat{w}_k + 2\lambda \hat{w}_k \end{bmatrix} = \begin{bmatrix} 2X^T Y e_0 \\ \vdots \\ 2X^T Y e_k \end{bmatrix}$$

$$\begin{bmatrix} X^T X \hat{w}_0 + \lambda \hat{w}_0 \\ \vdots \\ X^T X \hat{w}_k + \lambda \hat{w}_k \end{bmatrix} = \begin{bmatrix} X^T Y e_0 \\ \vdots \\ X^T Y e_k \end{bmatrix}$$

$$\begin{bmatrix} (X^T X + \lambda * 1) \hat{w}_0 \\ \vdots \\ (X^T X + \lambda * 1) \hat{w}_k \end{bmatrix} = \begin{bmatrix} X^T Y e_0 \\ \vdots \\ X^T Y e_k \end{bmatrix}$$

$$\begin{bmatrix} \hat{w}_0 \\ \vdots \\ \hat{w}_k \end{bmatrix} = \begin{bmatrix} (X^T X + \lambda * 1)^{-1} X^T Y e_0 \\ \vdots \\ (X^T X + \lambda * 1)^{-1} X^T Y e_k \end{bmatrix}$$

In matrix format, this translates to, $\hat{W} = (X^T X + \lambda I)^{-1} X^T Y$. □

b. Final output:

train error: 0.14805

test error: 0.1466

```

1 import numpy as np
2 from mnist import MNIST
3
4
5 def load_dataset():
6     mndata = MNIST('./data')
7     mndata.gz = True
8     X_train, labels_train = map(np.array, mndata.load_training())
9     X_test, labels_test = map(np.array, mndata.load_testing())
10    X_train = X_train / 255.0
11    X_test = X_test / 255.0
12    return X_train, labels_train, X_test, labels_test
13
14
15 x_train, label_train, x_test, label_test = load_dataset() # load in data
16
17
18 # print("x_train: ")
19 # print(x_train, x_train.shape)
20 # label_train = np.asarray(label_train)
21
22 # print("\noriginal label: ")
23 # print(label_train, label_train.shape)
24 # print(x_test, x_test.shape)
25 # print(label_test, label_test.shape)
26
27 def train(X, Y, reg_lambda):
28     d = X.shape[1]
29     # solve for \hat{w}: (XtX+lambda*I) * w_hat = XtY
30     lhs = np.dot(np.transpose(X), X) + reg_lambda * np.eye(d)
31     rhs = np.dot(np.transpose(X), Y)
32     w_hat = np.linalg.solve(lhs, rhs)
33     return w_hat
34

```



```

35
36 def predict(w_hat, X): # y values we are predicting are labels
37     Y_hat = np.dot(X, w_hat) # equal to  $e_j$ 'th column of  $I * w\_hat^T * x_i$ 'th column from
    X
38     # print("\nDim of y hat: ", Y_hat.shape)
39     return Y_hat.argmax(axis=1) # get argmax_j [ $e_j * w\_hat * x_i$ ]
40
41
42 num_class = 10 # this is k
43 onehot_label_train = np.eye(num_class)[label_train]
44 onehot_label_test = np.eye(num_class)[label_test]
45 # print("\none hot coded label: ")
46 # print(onehot_label_train, onehot_label_train.shape)
47
48 What = train(x_train, onehot_label_train, 1E-4)
49 # print("\ndim of w hat: ", What.shape)
50
51 # on training set
52 pred_train_label = predict(What, x_train)
53 # print("\npredicted label: ")
54 # print(pred_train_label, pred_train_label.shape)
55
56 labelDiff_train = np.equal(label_train, pred_train_label)
57 errorTrain = np.mean(1 - labelDiff_train)
58
59 # on test set
60 pred_test_label = predict(What, x_test)
61 labelDiff_test = np.equal(label_test, pred_test_label)
62 errorTest = np.mean(1 - labelDiff_test)
63
64 print("train error: ", errorTrain)
65 print("test error: ", errorTest)

```