# CSE546 Machine Learning HW1 - A

Bobby Deng — 1663039 — dengy7

15 April 2020

# 1   A.0

## 1.1   In your own words, describe what bias and variance are? What is bias-variance trade-off?

In class, we have derived that the total variability is composed by three parts: irreducible error, bias and variance. Irreducible error is the lowest error that our model can get. Bias is how much our model is deviated from the true model. When the model is two simple, it is more likely to deviate largely from the model. When the model is too complex and fits all training data points. There would not be any bias at all. Variance is basically the variance of the model, when model is a horizontal line, there is no variance at all. But if the model is too complex, the geometry representation of the model becoming too zigzag, and this is where the variance tends to be very big. The bias-variance trade-off is saying that we need to find a balance point where testing error is at the lowest point. Testing error is presented as an u shape, so ideally we want to find the model that produces the testing error at the bottom of the u shape curve.

## 1.2   What happens to bias and variance when the model complexity increases/decreases?

When complexity increase, bias decrease, variance increase. When complexity decrease, bias increase, variance decrease. Bias will be monotonically decrease and finally fits all training data points. Variance will be monotonically increase.

## 1.3   True or False: The bias of a model increases as the amount of training data available increases.

False, bias dose not depend on the amount of training data points.

## 1.4   True or False: The variance of a model decreases as the amount of training data available increases.

True, the variance do depend on the amount of training data.

## 1.5   True or False: A learning algorithm will generalize better if we use less features to represent our data.

False, we want enough useful and good predictors in the modes. The quality of predictor and coefficients determined how good our model might be. And less feature will not be better to present the ground truth model.

## 1.6   To get better generalization, should we use the train set or the test set to tune our hyper-parameters?

We should use train set to train the model and get hyper-parameters. We must not use any testing data in getting the hyper-parameters.

## 1.7 True or False: The training error of a function on the training set provides an overestimate of the true error of that function.

False, since we saw in both practice and theory that the training error is always lower than the testing error, in expectation. And when we train out model, our based model are generated by minimizing the error on training data. So it is an optimistic/underestimated estimation of the true function.

## A.1

Use MLE to estimate $\lambda$ in Poisson Distribution.

### a.

$$P(x|\lambda) = \frac{e^{-\lambda}\lambda^x}{x!}$$

$$P(x_1, x_2, x_3, x_4, x_5|\lambda) = \prod_{i=1}^{n} \frac{e^{-\lambda}\lambda^{x_i}}{x_i!} \text{ (Where n = 5.)}$$

$$\hat{\lambda}_{MLE} = \arg\max_{\lambda}(P(x_1, x_2, x_3, x_4, x_5|\lambda))$$

$$= \arg\max_{\lambda}(\prod_{i=1}^{n} \frac{e^{-\lambda}\lambda^{x_i}}{x_i!})$$

$$= \arg\max_{\lambda} ln(\prod_{i=1}^{n} \frac{e^{-\lambda}\lambda^{x_i}}{x_i!})$$

$$= \arg\max_{\lambda} \sum_{i=1}^{n} ln(\frac{e^{-\lambda}\lambda^{x_i}}{x_i!})$$

$$= \arg\max_{\lambda} \sum_{i=1}^{n} (lne^{-\lambda} + ln\lambda^{xi} - ln(xi!))$$

$$= \arg\max_{\lambda}(-n\lambda + ln\lambda * \sum_{i=1}^{n} xi + \sum_{i=1}^{n} ln(xi!))$$

Then set the derivative to 0.

$$0 = \frac{d}{d\lambda}(-n\lambda + ln\lambda * \sum_{i=1}^{n} xi + \sum_{i=1}^{n} ln(xi!))$$

$$= -n + \frac{1}{\lambda} \sum_{i=1}^{n} xi + 0$$

$$n = \frac{1}{\lambda} \sum_{i=1}^{n} xi$$

$$\hat{\lambda}_{MLE} = \frac{1}{n} \sum_{i=1}^{n} xi, n = 5.$$

### b.

$$\hat{\lambda}_{MLE} = \arg\max_{\lambda}(P(\mathcal{D}|\lambda))$$

$$= \arg\max_{\lambda}(P(x_1, x_2, x_3, x_4, x_5, x_6|\lambda))$$

$$= \frac{1}{n} \sum_{i=1}^{n} xi, n = 6.$$

**c.**

For $\lambda$ after 5, $\hat{\lambda}_{MLE} = (2 + 0 + 1 + 1 + 2)/5 = \frac{6}{5}$.
For $\lambda$ after 6, $\hat{\lambda}_{MLE} = (2 + 0 + 1 + 1 + 2 + 4)/6 = \frac{5}{3}$.

So, we can see the result that the outlier has a big effect on the estimation.

## A.2 German tank problem

From the problem, we know that the serial of number came from a uniform distribution with $[0, \theta]$.

$$f(x_i) = \begin{cases} \frac{1}{\theta} & \text{for } 0 \leq x_i \leq \theta. \\ 0 & \text{Otherwise.} \end{cases}$$

We can derive the MLE with uniform distribution.

$$\hat{\theta}_{MLE} = \arg\max_{\theta}(P(\mathcal{D}|\theta))$$

$$= \arg\max_{\theta}(\prod^n \frac{1}{\theta})$$

$$= \arg\max_{\theta} ln(\prod^n \frac{1}{\theta})$$

$$= \arg\max_{\theta} ln(\frac{1}{\theta^n})$$

$$= \arg\max_{\theta} -nln(\theta)$$

$$\frac{d(ln(L(\theta)))}{d\theta} = -n/\theta$$

We can not set the derivative to zero in this particular problem since $\theta$ is impossible to be infinity. So here, the only thing we can do is let $\theta$ to be as large as it could be. We can not arbitrages assign value to $\theta$. What we can do is to assign the biggest X we have observed to $\theta$.

$$\hat{\theta} = max(x_1, \ldots, x_n)$$

## A.3

**a.**

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

$$E_{train}[\hat{\epsilon}_{train}(f)] = \frac{N_{train}}{N_{train}} E[(f(x) - y)^2]$$

$$E_{train}[\hat{\epsilon}_{train}(f)] = E[(f(x) - y)^2] = \epsilon(f)$$

Since data from train and test are draw iid, the expectation for $(f(x) - y)^2$ should be the same.

$$E_{test}[\hat{\epsilon}_{test}(f)] = \frac{N_{test}}{N_{test}} E[(f(x) - y)^2] = E[(f(x) - y)^2] = \epsilon(f)$$

So,

$$E_{train}[\hat{\epsilon}_{train}(f)] = E_{test}[\hat{\epsilon}_{test}(f)] = \epsilon(f)$$

**b.**

No, the equation is not generally true, because the estimator we got from training is by minimizing the training error. It is not able to be general enough to describe the true model.

**c.**

From the hint we know that,

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

From a., we know that $E_{train}[\hat{\epsilon}_{train}(f)] = E_{test}[\hat{\epsilon}_{test}(f)]$.
So,

$$= \sum_{f \in \mathcal{F}} E_{train}[\hat{\epsilon}_{train}(f)]\mathbb{P}_{train}(\hat{f}_{train} = f)$$

Since $\hat{\epsilon}_{train}(\hat{f}) \leq \hat{\epsilon}_{train}(f)$:

$$\sum_{f \in \mathcal{F}} E_{train}[\hat{\epsilon}_{train}(\hat{f})]\mathbb{P}_{train}(\hat{f}_{train} = f) \leq \sum_{f \in \mathcal{F}} E_{train}[\hat{\epsilon}_{train}(f)]\mathbb{P}_{train}(\hat{f}_{train} = f)$$

$$\sum_{f \in \mathcal{F}} E_{train}[\hat{\epsilon}_{train}(\hat{f})]\mathbb{P}_{train}(\hat{f}_{train} = f) \leq E_{train,test}[\hat{\epsilon}_{test}(\hat{f}_{train})]$$

Because $\sum_{f \in \mathcal{F}} \mathbb{P}_{train}(\hat{f}_{train} = f) = 1$:

$$E_{train}[\hat{\epsilon}_{train}(\hat{f})] \sum_{f \in \mathcal{F}} \mathbb{P}_{train}(\hat{f}_{train} = f) \leq E_{train,test}[\hat{\epsilon}_{test}(\hat{f}_{train})]$$

$$E_{train}[\hat{\epsilon}_{train}(\hat{f}_{train})] \leq E_{train,test}[\hat{\epsilon}_{test}(\hat{f}_{train})]$$

# A.4

```python
import numpy as np


#---------------------------------------------------------------
#  Class PolynomialRegression
#---------------------------------------------------------------

class PolynomialRegression:

    def __init__(self, degree=1, reg_lambda=1E-8):
        """
        Constructor
        """
        self.theta = None
        self.regLambda = reg_lambda
        self.degree = degree
        self.mean = None
        self.std = None

    def polyfeatures(self, X, degree):
        """
        Expands the given X into an n * d array of polynomial features of
        degree d.

        Returns:
        A n-by-d numpy array, with each row comprising of
        X, X * X, X ** 3, ... up to the dth power of X.
        Note that the returned matrix will not include the zero-th power.

        Arguments:
```

```python
31      X is an n-by-1 column numpy array
32      degree is a positive integer
33      """
34      outputX = X[:]
35      for i in range(2, degree + 1):
36        outputX = np.hstack((outputX,X**i))
37      return outputX
38
39
40  def fit(self, X, y):
41      """
42      Trains the model
43      Arguments:
44      X is a n-by-1 array
45      y is an n-by-1 array
46      Returns:
47      No return value
48      Note:
49      You need to apply polynomial expansion and scaling
50      at first
51      """
52
53      X = self.polyfeatures(X, self.degree)
54      # standardization
55      self.mean = np.mean(X, axis=0)
56      self.std = np.std(X, axis=0)
57      X = (X - self.mean) / self.std
58      n = len(X)
59      # add 1s column
60      X_ = np.c_[np.ones([n, 1]), X]
61      n, d = X_.shape
62      reg_matrix = self.regLambda * np.identity(d )
63      reg_matrix[0, 0] = 0
64      self.theta = np.linalg.pinv((X_.T @ X_) + reg_matrix) @ (X_.T @ y)
65      print(self.theta)
66
67  def predict(self, X):
68      """
69      Use the trained model to predict values for each instance in X
70      Arguments:
71      X is a n-by-1 numpy array
72      Returns:
73      an n-by-1 numpy array of the predictions
74      """
75
76      n = len(X)
77      X = self.polyfeatures(X, self.degree)
78      X = (X - self.mean) / self.std
79      # add 1s column
80      X_ = np.c_[np.ones([n, 1]), X]
81      # predict
82      return X_ @ self.theta
83
84
85  #-------------------------------------------------------------
86  #  End of Class PolynomialRegression
87  #-------------------------------------------------------------
```
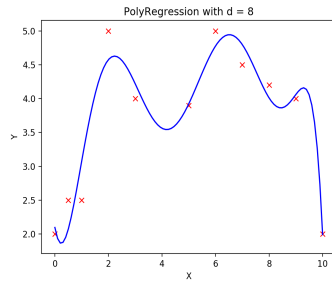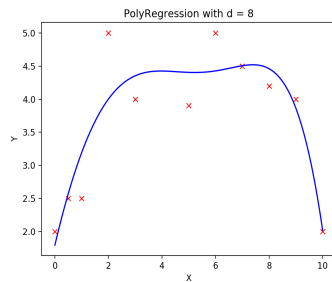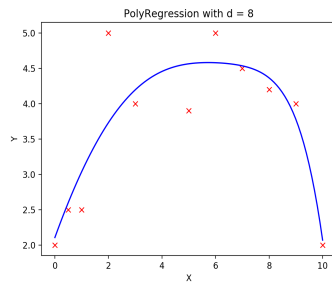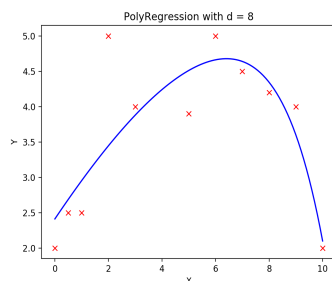
When lambda is 0:

When lambda is 0.0001:



When lambda is 0.01:



When lambda is 0.1:



From the graph, we can see that when lambda increase out model gets smoother. And the model tends to be simpler. And the weight tends to be smaller.

# A.5

```python
def learningCurve(Xtrain, Ytrain, Xtest, Ytest, reg_lambda, degree):
    """
    Compute learning curve

    Arguments:
    Xtrain -- Training X, n-by-1 matrix
    Ytrain -- Training y, n-by-1 matrix
    Xtest -- Testing X, m-by-1 matrix
    Ytest -- Testing Y, m-by-1 matrix
    regLambda -- regularization factor
    degree -- polynomial degree
```

```
12
13    Returns:
14    errorTrain -- errorTrain[i] is the training accuracy using
15    model trained by Xtrain[0:(i+1)]
16    errorTest -- errorTrain[i] is the testing accuracy using
17    model trained by Xtrain[0:(i+1)]
18
19    Note:
20    errorTrain[0:1] and errorTest[0:1] won't actually matter, since we start displaying the learning cur
21    """
22    n = len(Xtrain)
23
24    errorTrain = np.zeros(n)
25    errorTest = np.zeros(n)
26
27    for i in range(3, n+1):
28    print("i = ", i, " Degree= ", degree, " lambda: ", reg_lambda)
29    model = PolynomialRegression(degree, reg_lambda)
30    model.fit(Xtrain[:i], Ytrain[:i])
31
32    trainPredicted = model.predict(Xtrain[:i])
33    singleErrorFromTrain = np.mean((trainPredicted- Ytrain[:i])**2)
34    errorTrain[i-1] = singleErrorFromTrain
35
36    test_predicted = model.predict(Xtest[:i])
37    singleErrorFromTest = np.mean((test_predicted - Ytest[:i])**2)
38    errorTest[i-1] = singleErrorFromTest
39
40    return errorTrain, errorTest
```
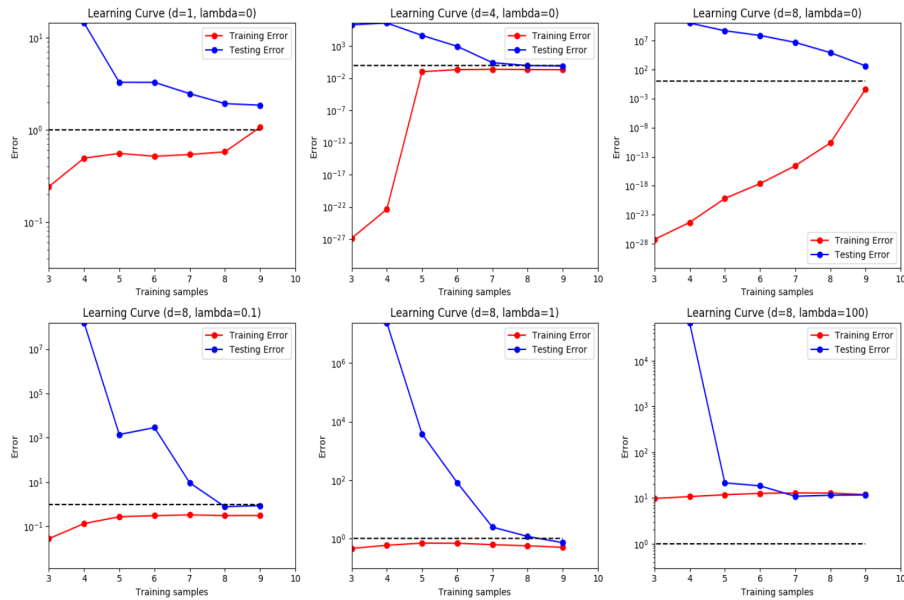


## A.6

**a.**

From the problem we know that:

$$\sum_{i=0}^{n} ||W^T x_i - y_i||_2^2 + \lambda ||W||_F^2 = \sum_{j=0}^{k} \left[ ||X w_j - Y e_j||^2 + \lambda ||w_j||^2 \right]$$

7

Let:
$$A = \sum_{j=0}^{k} \left[ ||Xw_j - Ye_j||^2 + \lambda ||w_j||^2 \right]$$

Take derivative:
$$\frac{\partial A}{\partial w_j} = \sum_{j=0}^{k} 2X^T \left[ Xw_j - Ye_j \right] + 2\lambda w_j$$

Set derivative to 0 to find minimum,

$$0 = \sum_{j=0}^{k} 2X^T \left[ Xw_j - Ye_j \right] + 2\lambda w_j$$

$$0 = \sum_{j=0}^{k} X^T Xw_j + \lambda w_j - X^T Ye_j$$

$$\sum_{j=0}^{k} (X^T X + \lambda I)w_j = \sum_{j=0}^{k} X^T Ye_j$$

Here, $X^T \in \mathbb{R}^{n \times d}$, $Y^T \in \mathbb{R}^{n \times k}$, $w_j \in \mathbb{R}^{d \times 1}$, $x_i \in \mathbb{R}^{d \times 1}$, $y_i \in \mathbb{R}^{k \times 1}$, $\lambda \in \mathbb{R}$, $e_j \in \mathbb{R}^{k \times 1}$, let's write it in matrix form:

$$(X^T X + \lambda I)\hat{W} = X^T Y$$
$$\hat{W} = (X^T X + \lambda I)^{-1} X^T Y$$

**b.**

```
1   import mnist
2   import numpy as np
3
4   mndata = mnist.MNIST("./python-mnist/data/")
5   X_train, labels_train = map(np.array, mndata.load_training())
6   X_test, labels_test = map(np.array, mndata.load_testing())
7   X_train = X_train/255.0
8   X_test = X_test/255.0
9
10  # Train function
11  def train(X,y,lam):
12    n, d = X.shape
13    y = np.eye(10)[y] # put y into a one hot encoding matrix
14      reg_matrix = self.regLambda * np.identity(d )
15      reg_matrix[0, 0] = 0
16      self.theta = np.linalg.pinv((X_.T @ X_) + reg_matrix) @ (X_.T @ y)
17    return W
18
19  def predict(W, X_new):
20    return (X_new @ W_hat).argmax(axis=1)
21
22  # Compute weights
23  W_hat = train(X_train, labels_train, lam = 0.00001)
24  # Computed predicted training label
25  predicted_train = predict(W=W_hat, X_new=X_train)
26  # Compute predicted testing label
27  predicted = predict(W=W_hat, X_new=X_test)
28  # Compute error rate for both training and testing
29  train_error_rate = 1 -( sum(predicted_train == labels_train) / len(labels_train))
30  test_error_rate = 1 - (sum(predicted == labels_test) / len(labels_test))
31  print("Training Accuracy is: ", train_error_rate)
```

```
32  print("Testing Error Rate is: ", test_error_rate)
33  #Training Error is:  0.14806666666666668
34  #Testing Error Rate is:  0.14659999999999995
```