# CSE546 Machine Learning HW1 - B

Bobby Deng — 1663039 — dengy7

15 April 2020

## B.1

### a.

Intuitively, when the function curve is really smooth and flat, meaning overly simple model, it ideally has low variance and high bias. And when the function curve is overly zigzag, meaning too complex model, it ideally has high variance, and low bias. In this problem, when step size, m, is small, the curve tends to be smoother. And when step size, m, is big, each interval is big, so the curve tends to be more zigzag.

### b.

Since:

$$\hat{f}_m(x) = \sum_{j=1}^{n/m} \frac{1}{m} \sum_{i=(j-1)m+1}^{jm} (f(x_i) + \epsilon_i) * 1\{x \in (\frac{(j-1)m}{n}, \frac{jm}{n}]\}$$

Let:

$$A = \frac{1}{n} \sum_{i=1}^{n} (E[\hat{f}_m(x_i)] - f(x_i))^2$$

So,

$$A = \frac{1}{n} \sum_{i=1}^{n} (E[\sum_{j=1}^{n/m} \frac{1}{m} \sum_{i=(j-1)m+1}^{jm} (f(x_i) + \epsilon_i) * 1\{x \in (\frac{(j-1)m}{n}, \frac{jm}{n}]\}] - f(x_i))^2$$

$$= \frac{1}{n} \sum_{i=1}^{n} (\sum_{j=1}^{n/m} (\frac{1}{m} \sum_{i=(j-1)m+1}^{jm} E[(f(x_i) + \epsilon_i)] * 1\{x \in (\frac{(j-1)m}{n}, \frac{jm}{n}]\} - f(x_i))^2$$

Recall that $\epsilon_i$ is centered around 0 and $f(x_i)$ is fixed.

$$= \frac{1}{n} \sum_{i=1}^{n} (\sum_{j=1}^{n/m} \frac{1}{m} \sum_{i=(j-1)m+1}^{jm} f(x_i) * 1\{x \in \frac{(j-1)m}{n}, \frac{jm}{n}]\} - f(x_i))^2$$

$$= \frac{1}{n} \sum_{i=1}^{n} (\sum_{j=1}^{n/m} \bar{f}^{(j)} * 1\{x \in \frac{(j-1)m}{n}, \frac{jm}{n}]\} - f(x_i))^2$$

Now, for all $x_i$, we only need x in certain j interval, which from (j-1)m+1 to jm.

$$= \frac{1}{n} \sum_{j=1}^{n/m} \sum_{i=(j-1)m+1}^{jm} (\bar{f}^{(j)} - f(x_i))^2$$

**c.**

1.

From previous question, we know that:

$$E[\hat{f}_m(x)] = \sum_{j=1}^{n/m} \bar{f}^{(j)} * 1\{x \in (\frac{(j-1)m}{n}, \frac{jm}{n}]\}$$

Let

$$A = E[\frac{1}{n}\sum_{I=1}^{N}(\hat{f}_m(x_i) - E[\hat{f}_m(x_i)])^2]$$

Then let's compute something:

$$E[c_j] = E[\frac{1}{m}\sum_{i=(j-1)m+1}^{jm} y_i]$$

$$= E[\frac{1}{m}\sum_{i=(j-1)m+1}^{jm} f(x_i) + \epsilon_i]$$

$$= \frac{1}{m}\sum_{i=(j-1)m+1}^{jm} E[f(x_i) + \epsilon_i]$$

Recall that $\epsilon_i$ is centered around 0 and $f(x_i)$ is fixed.

$$= \frac{1}{m}\sum_{i=(j-1)m+1}^{jm} f(x_i)$$

$$= \bar{f}^{(j)}$$

Now, with $E[c_j] = \bar{f}^{(j)}$, we have:

$$A = \frac{1}{n}\sum_{I=1}^{N} E[(\hat{f}_m(x_i) - E[\hat{f}_m(x_i)])^2]$$

$$= \frac{1}{n}\sum_{I=1}^{N} E[(\sum_{j=1}^{n/m} c_j * 1\{x \in (\frac{(j-1)m}{n}, \frac{jm}{n}]\} - E[\sum_{j=1}^{n/m} c_j * 1\{x \in (\frac{(j-1)m}{n}, \frac{jm}{n}]\}])^2]$$

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

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

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

2.

$$A = \frac{1}{n}\sum_{I=1}^{N} E[(\hat{f}_m(x_i) - E[\hat{f}_m(x_i)])^2]$$

$$= \frac{1}{n}\sum_{I=1}^{N} E[(\sum_{j=1}^{n/m} c_j * 1\{x \in (\frac{(j-1)m}{n}, \frac{jm}{n}]\} - \sum_{j=1}^{n/m} E[c_j] * 1\{x \in (\frac{(j-1)m}{n}, \frac{jm}{n}]\})^2]$$

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

2

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

$$= \frac{1}{n} \sum_{j=1}^{n/m} m \, var(c_j)$$

$$= \frac{1}{n} \sum_{j=1}^{n/m} m \, var(\frac{1}{m} \sum_{i=(j-1)m+1}^{jm} (f(x_i) + \epsilon_i))$$

$$= \frac{1}{n} \sum_{j=1}^{n/m} m * \frac{1}{m^2} \sum_{i=(j-1)m+1}^{jm} var(f(x_i) + \epsilon_i)$$

$$= \frac{1}{n} \sum_{j=1}^{n/m} \frac{1}{m} \sum_{i=(j-1)m+1}^{jm} \sigma^2$$

$$= \frac{1}{n} \sum_{i=1}^{n} \frac{\sigma^2}{m}$$
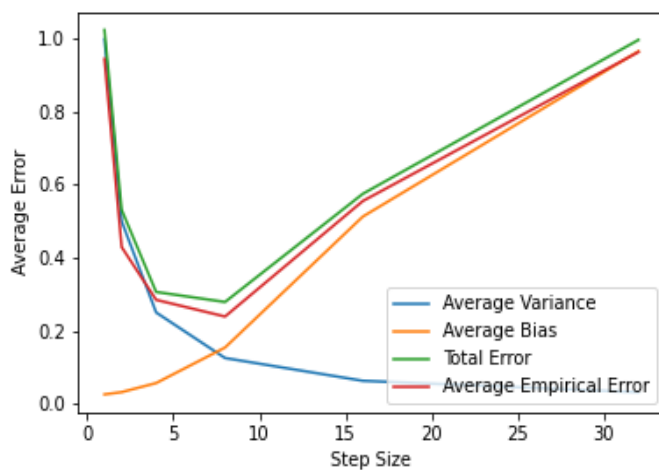
$$= \frac{\sigma^2}{m}$$

**d.**

```
1   import pandas as pd
2   import numpy as np
3   import matplotlib.pyplot as plt
4   import random
5
6   np.random.seed(10)
7   n = 256
8   sigma2 = 1
9   mean = 0
10  m_list = [1,2,4,8,16,32]
11  xList = np.array(range(1,n+1))/n
12
13  def f(x):
14      return 4 * np.sin(np.pi * x) * np.cos(6 * np.pi * x ** 2)
15
16  def y_i(x):
17      return f(x) +  np.random.normal(0,1)
18
19  def cj(m,j):
20      sum = 0.0
21      for i in range((j-1) * m + 1, j * m + 1):
22          sum += y_i(i/n)
23      return sum/m
24
25  def f_hat(x, m):
26      sum = 0.0
27      for j in range(1, int(n/m) + 1):
28          if ((x > (j -1) *1.0* m /n) and (x <= j*m*1.0/n)):
29              sum += cj(m,j)
30      return sum
31
32  def f_bar(j,m):
33      sum = 0.0
34      for i in range((j-1)*m+1, j*m + 1):
35          sum += f(i/n)
```

```python
36        return sum/m
37
38   def bias(m,n):
39        output = 0.0
40        for j in range(1, int(n/m)+1):
41            for i in range((j-1)*m, j*m):
42                output += (f_bar(j, m) - f(i/n))**2
43        return output/n
44
45   def variance(sigma2, m):
46        return sigma2 / m
47
48   #Initialize
49   empirical_error = []
50   bias_list = []
51   bias_sum = 0.0
52   variance_list = []
53   total_error = []
54
55   #Start iteration
56   for m in m_list:
57        empirical_error_sum = 0.0
58        # empirical_error
59        for i in range(1,n+1):
60            empirical_error_sum += (f_hat(i/n, m) - f(i/n))**2
61        empirical_error.append(empirical_error_sum/n)
62        # variance
63        bias_list.append(bias(m,n))
64        # bias
65        variance_list.append(variance(sigma2, m))
66        #total error
67        total_error = np.array(bias_list) + np.array(variance_list)
68
69   plt.plot(m_list,variance_list, label="Average Variance")
70   plt.plot(m_list,bias_list, label="Average Bias")
71   plt.plot(m_list,total_error, label="Total Error")
72   plt.plot(m_list,empirical_error, label="Average Empirical Error")
73   plt.xlabel("Step Size")
74   plt.ylabel("Average Error")
75   plt.legend()
```

**e.**

In the average bias squared expression $\frac{1}{n}\sum_{j=1}^{n/m}\sum_{i=(j-1)m+1}^{jm}(\bar{f}^{(j)}-f(x_i))^2$, we can notice that $\frac{1}{n}\sum_{j=1}^{n/m}\sum_{i=(j-1)m+1}^{jm}$ does not make the total expression much bigger or smaller. It is only doing the summation and then do the average. I will just ignore it for a moment.

From the question, we know that:

$$min_{i=(j-1)m+1,...,jm}f(x_i) \leq \bar{f}_{(j)} \leq max_{i=(j-1)m+1,...,jm}f(x_i)$$

Based on the L-Lipschitz rule, we got:

$$|\bar{f}_{(j)} - f(x_i)| \leq |max_{i=(j-1)m+1,...,jm}f(x_i) - min_{i=(j-1)m+1,...,jm}f(x_i)|$$

$$|\bar{f}_{(j)} - f(x_i)| \leq \frac{L}{n}|max_{i=(j-1)m+1} - min_{i=(j-1)m+1}|$$

$$|\bar{f}_{(j)} - f(x_i)| \leq \frac{L}{n}|m|$$

Recall that each interval has m elements

$$(\bar{f}_{(j)} - f(x_i))^2 \leq (\frac{L}{n}|m|)^2$$

$$\leq \frac{L^2 m^2}{n^2}$$

As for the total error, and recall $\sigma^2 = 1$,

$$(\bar{f}_{(j)} - f(x_i))^2 + \frac{\sigma^2}{m} \leq \frac{L^2(m)^2}{n^2} + \frac{1}{m}$$

Minimizing the total error with respect to m.

$$\frac{dO}{dm} = \frac{2L^2 m}{n^2} - m^{-2}$$

$$0 = \frac{2L^2 m}{n^2} - m^{-2}$$

$$m = (\frac{n^2}{2L^2})^{\frac{1}{3}}$$

From above, we know that m can not be 0, can not be negative. And it is really intuitive that when m increases, the bias term decrease, and variance increase. It make sense that it need to find a balance point according to the data size. Also we can see that there is a positive relationship between m and n, and there is a negative relationship between n and L.

For this particular scenario, n is 256, and we can see the lowest valley is when m is 8.

$$8 = (\frac{256^2}{2L^2})^{\frac{1}{3}}$$

$$L = 8$$

# B.2

**a.**

```python
import numpy as np
import matplotlib.pyplot as plt
import mnist

mndata = mnist.MNIST('./python-mnist/data/')
X_train, labels_train = map(np.array, mndata.load_training())
X_test, labels_test = map(np.array, mndata.load_testing())
X_train = X_train/255.0
```
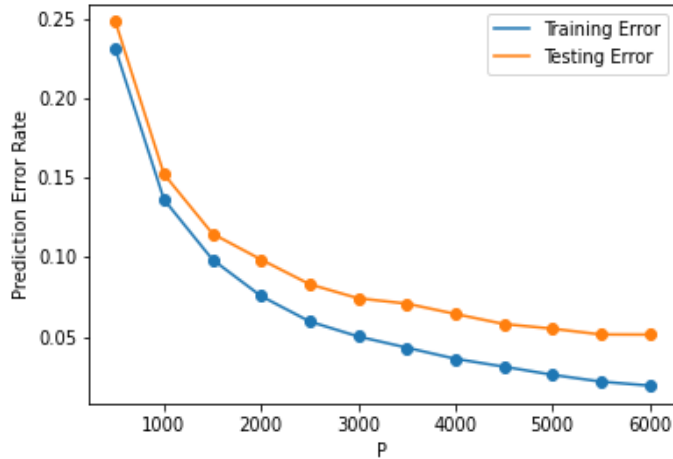
```python
 9   X_test = X_test/255.0
10   # This function trains the model the return the weights
11   def train(X, Y):
12     lambda_ = 0.0001
13     n, d = np.shape(X)
14     W = np.linalg.solve(X.T @ X + lambda_ * np.identity(d), X.T @ Y)
15     return W
16
17   # This function do the prediction
18   def predict(W,X):
19     return (X @ W).argmax(axis = 1)
20
21   # This function apply the transformation to data
22   def h1(X_train, X_test, p):
23     n, d = X_train.shape
24     sigma = np.sqrt(0.1)
25     G = np.random.normal(0, sigma, p * d).reshape(p, d)
26     b = np.random.uniform(0, 2 * np.pi, p).reshape(p, 1)
27     h_train = np.cos(np.dot(X_train, G.T) + b.T)
28     h_test = np.cos(np.dot(X_test, G.T) + b.T)
29
30     return h_train, h_test, G, b
31
32   n, d = X_train.shape
33   training_error_all = []
34   validing_error_all = []
35   W_list = []
36   Gb_list = []
37   train_index = np.random.choice(np.arange(n), int(X_train.shape[0] * 0.8), replace=False)
38   valid_index = np.setdiff1d(np.arange(n), train_index)
39   # loop from p=500 to p=6000, step=500
40   for p in list(range(500, 6001, 500)):
41     h_train, h_test, G, b = h1(X_train[train_index, :], X_train[valid_index, :], p)
42     # h = h1(X_train, p)
43     # Train test split with 80%-20%
44     Gb_list.append((G,b))
45     train_data = h_train
46     valid_data = h_test
47     y_train = np.eye(10)[labels_train[train_index]]
48
49     # Compute weights
50     W_hat = train(train_data, y_train)
51     W_list.append(W_hat)
52     # Compute train predicted
53     predict_train = predict(W_hat, train_data)
54     predict_train = labels_train[train_index] - predict_train
55     train_error_single = np.count_nonzero(predict_train) / len(predict_train) #train_size
56     training_error_all.append(train_error_single)
57     # Compute test predicted
58     predicted_test = predict(W_hat, valid_data)
59     predicted_test = labels_train[valid_index] - predicted_test
60     valid_error_single = np.count_nonzero(predicted_test) / len(predicted_test)
61     validing_error_all.append(valid_error_single)
62
63     print("p: ", p,", train_err: ", train_error_single, ", test_err: ", valid_error_single)
64
65   x_index = list(range(500, 6001, 500))
66   plt.plot(x_index, training_error_all, label="Training Error")
```

```
67  plt.scatter(x_index, training_error_all)
68  plt.plot(x_index, validing_error_all, label="Testing Error")
69  plt.scatter(x_index, validing_error_all)
70  plt.xlabel("P")
71  plt.ylabel("Prediction Error Rate")
72  plt.legend()
```



**b.**

From the question, we could have:

$$\frac{1}{m}\sum_{i=1}^{m} x_i - \sqrt{\frac{\log(\frac{2}{\delta})}{2m}} \le \mu \le \frac{1}{m}\sum_{i=1}^{m} x_i + \sqrt{\frac{\log(\frac{2}{\delta})}{2m}}$$

And $\frac{1}{m}\sum_{i=1}^{m} x_i$ is basically the mean test error.

```
1   n, d = X_train.shape
2   p = 6000
3   dt = 0.05
4   sigma = np.sqrt(0.1)
5
6   h_train, h_test, G, b = h1(X_train, X_test, p)
7   y_train = np.eye(10)[labels_train]
8
9   # Compute weights
10  W_hat = train(h_train, y_train)
11  # Compute test predicted
12  predicted_test = predict(W_hat, h_test)
13  valid_error_single = 1 - (sum(labels_test == predicted_test) / len(labels_test))
14
15  H = np.sqrt(np.log(2/dt)/(2*len(labels_test)))
16  print(f'The test_error is {valid_error_single}')
17  print(f'Confidence Interval:[{valid_error_single - H} : {valid_error_single + H}]')
18  # The test_error is 0.0460000000000004
19  # Confidence Interval:[0.03241898484259385 : 0.059581015157406235]
```