# CSE546 Machine Learning HW2

Bobby Deng — 1663039 — dengy7

3 May 2020

## A.0

**a. [2 points] Suppose that your estimated model for predicting house prices has a large positive weight on 'number of bathrooms'. Does it implies that if we remove the feature "number of bathrooms" and refit the model, the new predictions will be strictly worse than before? Why?**

Answer: Yes, when there is a large positive weight comparatively to other weights. It means that feature is important and may have a strong positive linear relationship with the target class.

**b. [2 points] Compared to L2 norm penalty, explain why a L1 norm penalty is more likely to result in a larger number of 0s in the weight vector or not?**

Answer: The gradient for l1 is always the same, however l2 norm produce a gradient with diminishing return as weights move to zero. If we remove that important feature it spread the effect to other features which produces more bias. So we want to keep important feature and remove less important features.

**c. [2 points] In at most one sentence each, state one possible upside and one possible downside of using the following regularizer: $(\sum_i |w_i|^{0.5})$**

Answer: The downside it this norm tends to choose bigger coefficients compared to L2 and might lead to over-fitting.One possible upside could be that it only produce positive gradient and may converge faster.

**d. [1 points] True or False: If the step-size for gradient descent is too large, it may not converge.**

Answer: True

**e. [2 points] In your own words, describe why SGD works.**

Answer: SGD refers to Stochastic Gradient Descent. And it is computing the gradients of one sample and update the weights. And do this process over and over again until some converge conditions are met.

**f. [2 points] In at most one sentence each, state one possible advantage of SGD (stochastic gradient descent) over GD (gradient descent) and one possible disadvantage of SGD relative to GD.**

Answer: SGD might be a little bit computationally efficient. The disadvantage is that SGD produce unstable gradients and it is likely to bouncing back and forth.

## A.1

**a.**

To prove it is a norm, we want to prove:

$$||x + y|| \leq ||x|| + ||y||$$

First:

$$|a + b| - > |a + b|^2 = (a + b)^2 = a^2 + 2ab + b^2$$
$$|a| + |b| - > (|a| + |b|)^2 = a^2 + 2|a||b| + b^2$$

And $2|a||b| \geq 2ab$, So:

$$|a + b| \leq |a| + |b|$$

Then from above we can get:

$$\sum_{i=1}^{n} |a + b| \leq \sum_{i=1}^{n} |a| + \sum_{i=1}^{n} |b|$$

So it satisfy the triangle inequality:

$$||a + b|| \leq ||a|| + ||b||$$

So it is a norm.

**b.**

This is not a norm. For points x(1,4), y(4,1), it does not satisfy the triangle inequality.

**B.1**

$$\sqrt{\sum_{i=1}^{n} x_i^2} \leq \sum_{i=1}^{n} |x_i|$$

$$\sum_{i=1}^{n} x_i^2 \leq (\sum_{i=1}^{n} |x_i|)^2$$

Let $|x_i| = Z_i$

$$\sum_{i=1}^{n} Z_i^2 \leq (\sum_{i=1}^{n} Z_i)^2$$

$$\sum_{i=1}^{n} Z_i^2 \leq (\sum_{j=1}^{n} Z_j)(\sum_{i=1}^{n} Z_i)$$

$$\sum_{i=1}^{n} Z_i^2 \leq (Z_1 + Z_2 + \cdots + Z_n)(Z_1 + Z_2 + \cdots + Z_n)$$

$$\sum_{i=1}^{n} Z_i^2 \leq \sum_{i=1}^{n} Z_j^2 + \sum_{j=1}^{n} \sum_{j \neq i}^{n} Z_j Z_i$$

With this equality, we can conclude that with higher order, there will be more terms on the right side of the equation. So:

$$||x||\infty \leq ||x||_2 \leq ||x||_1$$

**A.2**

- 1 is not convex, the line segment from b to c it not convex set.
- 2 is convex, every points inside is convex set.
- 3 is not convex, the line segment from a to d is not convex set.

**A.3**

- a. Yes
- b. No, the segment from a to c is not convex set.
- c. No, the line segment from a to d is not convex set.
- d. Yes

# B.2

## a.

The goal it to achieve triangle inequality that:

$$f(ax + (1-a)y) \leq af(x) + (1-a)f(y)$$

Firstly, lets compute the left part:

$$f(ax + (1-a)y) = |ax + (1-a)y|$$

Assume:

$$|ax + (1-a)y| \leq a|x| + (1-a)|y|$$

Square left and right:

$$a^2x^2 + 2a(1-a)xy + (1-a)^2y^2 \leq a^2x^2 + 2a(1-a)|x||y| + (1-a)^2y^2$$

Since x and y could positive and negative so:

$$2a(1-a)xy \leq 2a(1-a)|x||y|$$

So in fact the assumption is right:

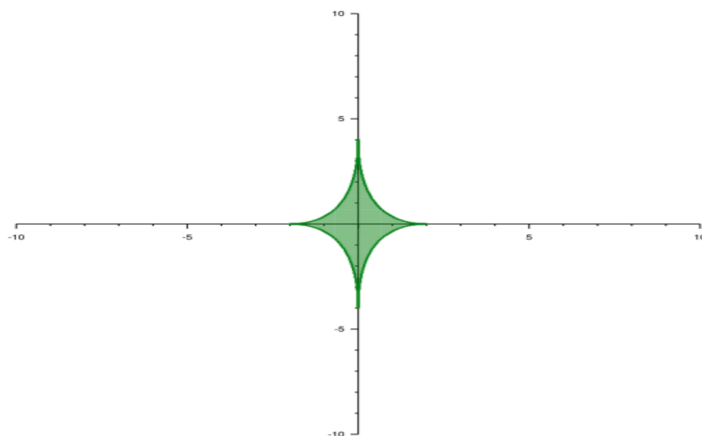$$|ax + (1-a)y| \leq a|x| + (1-a)|y|$$

So it is convex function.

## b.

The same as previous question, we need to prove the triangle inequality to show it is convex function.

$$||ax + (1-a)y|| \leq a||x|| + (1-a)||y|| \leq a + (1-a) = 1$$

So from above, we proved that it satisfy the triangle inequality, and so it is convex function.

## c.

No, it is not convex function and we can see it from the graph that if we draw two points on the graph and it could be above the function area. So the set is not a convex set.



$$\left( |x \cdot 1|^{\frac{1}{2}} + |x \cdot 2|^{\frac{1}{2}} \right)^2 \leq 4$$

## B.3

### a.

We want to prove the triangle inequality:

$$f(ax_1 + (1-a)x_2) \le af(x_1) + (1-a)f(x-2)$$

Let:

$$w_i = ax_i + (1-a)y_i$$

We know these two functions are convex function separately, now we compute 2 parts of the function separately:

$$\sum_{i=1}^{n} \ell_i(w_i) \le a\sum_{i=1}^{n} \ell(x_i) + (1-a)\sum_{i=1}^{n} \ell(y_i)$$

$$\sum_{i=1}^{n} \lambda||w_i|| \le a\sum_{i=1}^{n} \lambda||x_i|| + (1-a)\sum_{i=1}^{n} \lambda||y_i||$$

So, we combine the two equation into one:

$$\sum_{i=1}^{n} \ell_i(w_i) + \sum_{i=1}^{n} \lambda||w_i|| \le a\sum_{i=1}^{n} \ell(x_i) + (1-a)\sum_{i=1}^{n} \ell(y_i) + a\sum_{i=1}^{n} \lambda||x_i|| + (1-a)\sum_{i=1}^{n} \lambda||y_i||$$

Move terms around:

$$\sum_{i=1}^{n} \ell_i(w_i) + \sum_{i=1}^{n} \lambda||w_i|| \le a\sum_{i=1}^{n} \ell(x_i) + a\sum_{i=1}^{n} \lambda||x_i|| + (1-a)\sum_{i=1}^{n} \ell(y_i) + (1-a)\sum_{i=1}^{n} \lambda||y_i||$$

Simplify and we get the triangle inequality:

$$\sum_{i=1}^{n} \ell_i(w_i) + \sum_{i=1}^{n} \lambda||w_i|| \le a\sum_{i=1}^{n} \left(\ell(x_i) + \lambda||x_i||\right) + (1-a)\sum_{i=1}^{n} \left(\ell(y_i) + \lambda||y_i||\right)$$

So, the sum of the two convex function is still a convex function.

### b.

We could easily find the local minimum in convex function and in convex function local minimum is global minimum.
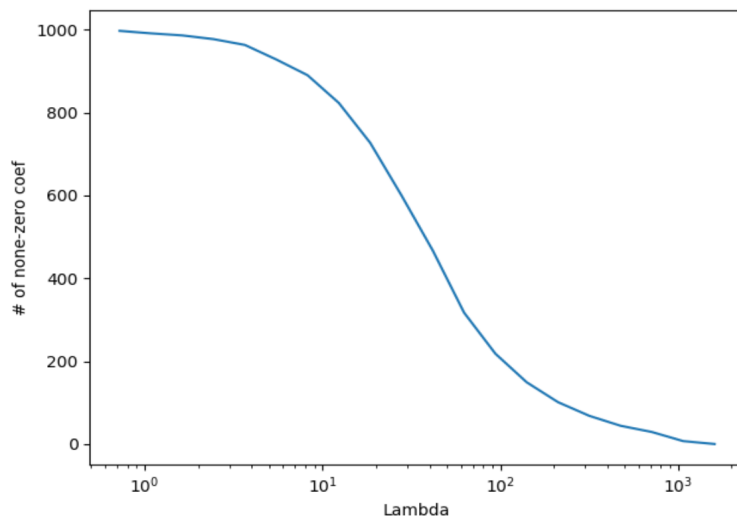
## A.4

```
1  import numpy as np
2  import matplotlib.pyplot as plt
3  import pandas as pd
4
5  def generate_x(n,d):
6    return np.random.standard_normal((n,d))
7
8  def generate_y(n, d, k, X):
9    w = np.zeros((d,1))
10   for i in range(0, d):
11     if (0 <= i and i < k):
12       w[i] = (i + 1)/k
13     else:
14       w[i] = 0
15
16   y = np.zeros((n, 1))
17   for i in range(n):
```

```
18         y[i] = w.T @ X[i] + np.random.standard_normal()
19     return y, w
20
21
22 def compute_initial_lamb(x, y):
23     n, d = x.shape
24     lamb_array = np.zeros((d, 1))
25     for k in range(d):
26         lamb_temp = 2 * abs(x[:, k].T @ (y - np.mean(y)))
27         lamb_array[k] = lamb_temp
28     return max(lamb_array)
29
30 class Lasso:
31
32     def __init__(self, lamb=0.001, delta=0.05):
33         self.lamb = lamb
34         self.last_w = None
35         self.b = 0.0
36         self.delta = delta
37         self.loss_list = []
38         self.last_selected_coef = []
39         self.selected_feature_index = [1, 3, 5, 7, 12]
40
41     def coordinate_descent(self, X, y, initial_w):
42         n, d = X.shape
43         W = initial_w
44         a = 2 * np.sum(np.power(X, 2), axis=0)
45
46         not_converge = True
47         while not converged:
48             self.b = np.average(y - X.dot(W))
49             loss_prev = loss
50             prev_w = np.copy(W)
51             for k in range(d):
52                 X_k = X[:, k]
53                 prev_wk = np.copy(W[k])
54                 W[k] = 0
55                 c_k = np.dot((y - (self.b * np.ones((n, 1)) + X.dot(W))).T, X_k)
56                 if 2 * c_k + self.lamb < 0:
57                     W[k] = (c_k * 2 + self.lamb) / a[k]
58                 elif 2 * c_k - self.lamb > 0:
59                     W[k] = (c_k * 2 - self.lamb) / a[k]
60                 else:
61                     W[k] = 0
62
63             if sum(abs(W - prev_w)) <= sum(abs(self.delta * prev_w)):
64                 not_converge = False
65                 print(W)
66
67             loss = np.sum(np.power((self.b * np.ones((n, 1)) + X.dot(W) - y), 2)) + self.lamb * np.sum(abs
68             self.loss_list.append(loss)
69             self.last_w = W
70             self.last_selected_coef = self.last_w.T[0][self.selected_feature_index]
71
72     def predict(self, X):
73     return X.dot(self.last_w)
```
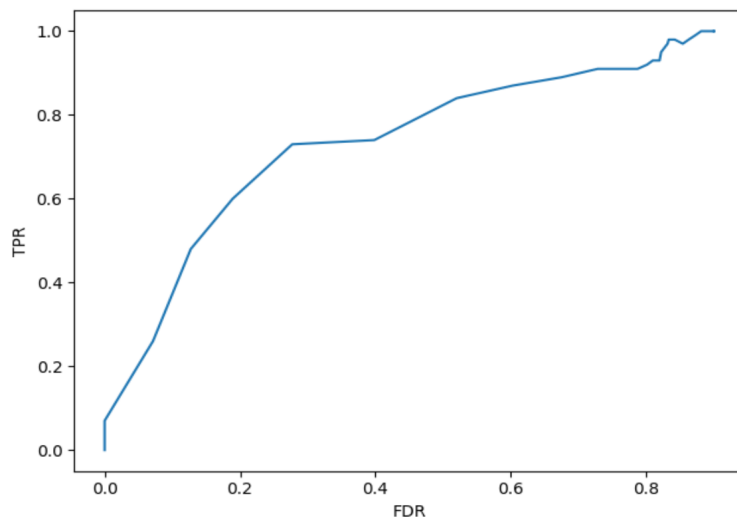
**a.**



```python
1   from A4_A5_starter import *
2
3   if __name__ == '__main__':
4     n = 500
5     d = 1000
6     k = 100
7
8     X_train = generate_x(n, d)
9     y_train, W_init = generate_y(n, d, k, X_train)
10    lam = compute_initial_lamb(X_train, y_train)
11    number_of_nonezero_feature = []
12    FDR_list = []
13    TPR_list = []
14
15    lam_list = lam * (1/1.5) ** np.arange(0, 20)
16
17    for lam in lam_list:
18
19      print("lam", lam)
20      lasso = Lasso(lam, delta=0.4)
21      lasso.coordinate_descent(X_train, y_train, np.zeros((d,1)))
22      last_w = lasso.last_w.copy()
23      print("Number of coe > 0:", sum(abs(last_w) > 0))
24      number_nonezero = sum(last_w != 0)
25      number_of_nonezero_feature.append(number_nonezero)
26
27
28    plt.plot(lam_list, number_of_nonezero_feature)
29    plt.xscale('log')
30    plt.xlabel("Lambda")
31    plt.ylabel("# of none-zero coef")
32    plt.show()
```

**b.**



```
1   from A4_A5_starter import *
2
3   if __name__ == '__main__':
4       n = 500
5       d = 1000
6       k = 100
7
8       X_train = generate_x(n, d)
9       y_train, W_init = generate_y(n, d, k, X_train)
10      lam = compute_initial_lamb(X_train, y_train)[0]
11      number_of_nonezero_feature = []
12      FDR_list = []
13      TPR_list = []
14
15      lam_list = lam * (1/1.5) ** np.arange(0, 40)
16
17      for lam in lam_list:
18
19          print("lam", lam)
20          lasso = Lasso(lam, delta=0.001)
21          lasso.coordinate_descent(X_train, y_train, np.zeros((d,1)))
22          last_w = lasso.last_w
23          print("Number of coe > 0:", sum(abs(last_w) > 0))
24          number_nonezero = sum(last_w != 0)
25          number_of_nonezero_feature.append(number_nonezero)
26
27          incorrect_none_zero = sum(last_w[W_init == 0] != 0)
28          number_correct_none_zero = sum(last_w[W_init != 0] != 0)
29          if incorrect_none_zero == 0:
30              FDR = 0
31              FDR_list.append(0)
32          else:
33              FDR = incorrect_none_zero / number_nonezero
34              FDR_list.append(FDR)
35          TPR = number_correct_none_zero / k
36          TPR_list.append(TPR)
37
38          print("FDR: ", FDR, " TPR: ", TPR)
39
```

```
40
41
42   plt.plot(FDR_list, TPR_list)
43   plt.xlabel("FDR")
44   plt.ylabel("TPR")
45   plt.show()
```
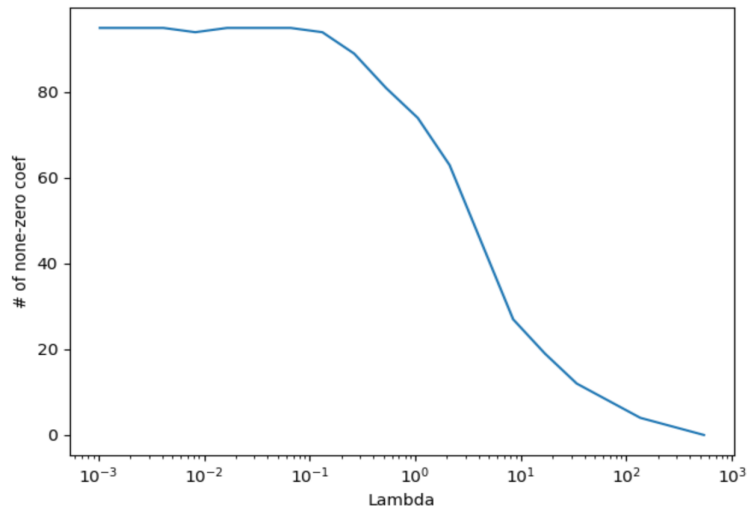
**c.**

When lambda gets bigger, number of none-zero term decrease until 0. When lambda is small enough, there is no zero term. I can see that it is important to chooes an proper lambda.
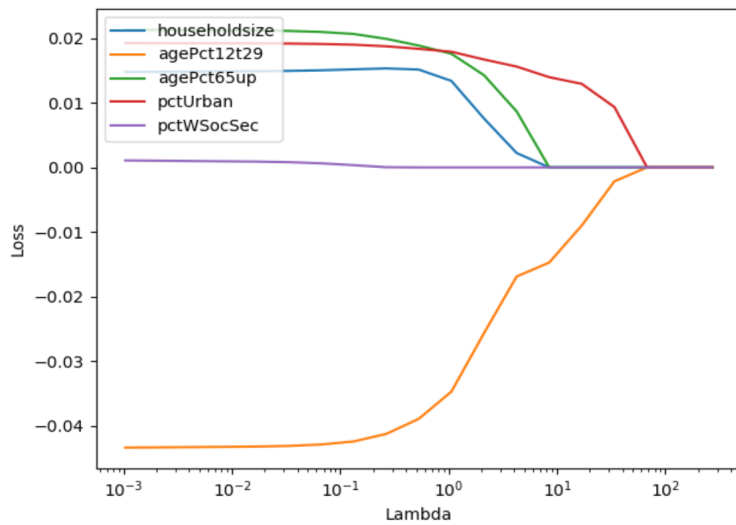
In the TPR-FDR relationship plot, they tends to have a positive relationship. When lambda is at max, all feature coeff are 0, so TPR and FDR are both 0. When lambda decrease, TPR and FDR both increase.
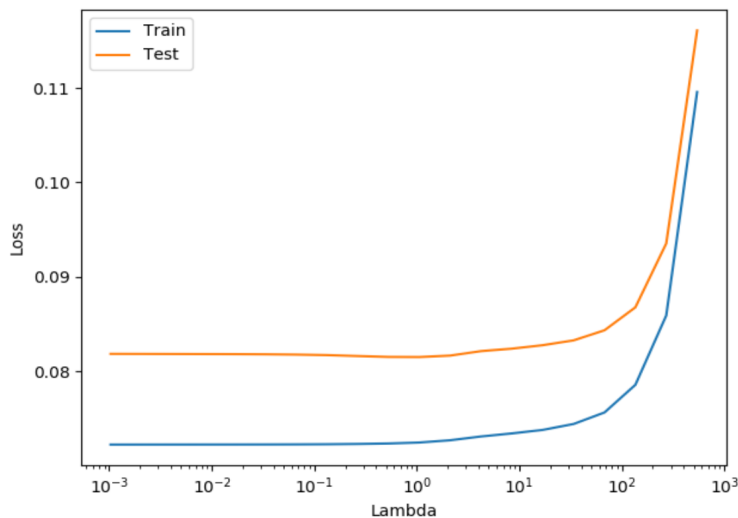
# A.5

**a.**



**b.**

**c.**



**d.**

- Coeff Name

- 0.068025 PctIlleg

- -0.070175 PctKids2Par

- For highest positive feature: PctIlleg(percentage of kids born to never married (numeric - decimal)). A reasonable explanation for this is that in poor area people had less sense of contraception, can lead to unwanted kids burn. Also may become one of the reason. And poor districts may have higher crime rate as well.

- For highest negative feature: PctKids2Par( percentage of kids in family housing with two parents (numeric - decimal)). This is really reasonable. With two parents probabaly means a good family, and higher such rate means the district are more good families. And this district might be rich and well educated and has good measure for preventing crime to take place.

**e.**

This does not make sense, because mathematically, when where are more older people, a lot of other things will change as well. Meaning coefficients of other variable may change and then the model would not be the same. Then the crime rate may not decrease. On the other hand, crime is not produced by older people. It does not make sense to introduce more older people just in order to lower the crime rate. We should focused on what caused crime and reduce that factor.

```python
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt


class Lasso:

    def __init__(self, lamb=0.001, delta=0.05):
        self.lamb = lamb
        self.last_w = None
        self.b = 0.0
        self.delta = delta
        self.loss_list = []
        self.last_selected_coef = []
        self.selected_feature_index = [1, 3, 5, 7, 12]
```

```python
    def coordinate_descent(self, X, y, initial_w):
      n, d = X.shape
      W = initial_w
      a = 2 * np.sum(np.power(X, 2), axis=0)
      loss = np.sum(np.power((self.b * np.ones((n, 1)) + X.dot(W) - y), 2)) +
       self.lamb * np.sum(abs(W))

      not_converge = True
      while not_converge:
        self.b = np.average(y - X.dot(W))
        loss_prev = loss
        prev_w = np.copy(W)
        for k in range(d):
          X_k = X[:, k]

          W[k] = 0
          c_k = np.dot((y - (self.b * np.ones((n, 1)) + X.dot(W))).T, X_k)
          if 2 * c_k + self.lamb < 0:
            W[k] = (c_k * 2 + self.lamb) / a[k]
          elif 2 * c_k - self.lamb > 0:
            W[k] = (c_k * 2 - self.lamb) / a[k]
          else:
            W[k] = 0

        if sum(abs(W - prev_w)) <= sum(abs(self.delta * prev_w)):
          not_converge = False
          print(W)

        loss = np.sum(np.power((self.b * np.ones((n, 1)) + X.dot(W) - y), 2)) + self.lamb * np.sum(abs
        self.loss_list.append(loss)
        self.last_w = W
        self.last_selected_coef = self.last_w.T[0][self.selected_feature_index]

    def predict(self, X):
      return X.dot(self.last_w)


def compute_initial_lamb(x, y):
  n, d = x.shape
  lamb_array = np.zeros((d, 1))
  for k in range(d):
    lamb_temp = 2 * abs(x[:, k].T @ (y - np.mean(y)))
    lamb_array[k] = lamb_temp
  return max(lamb_array)

if __name__ == "__main__":
  df_train = pd.read_table("crime-train.txt")
  df_test = pd.read_table("crime-test.txt")
  y_train = df_train["ViolentCrimesPerPop"].values.reshape(df_train.shape[0], 1)
  x_train = df_train.drop("ViolentCrimesPerPop", axis=1).values
  x_test = df_test.drop("ViolentCrimesPerPop", axis=1).values
  y_test = df_test["ViolentCrimesPerPop"].values.reshape(df_test.shape[0], 1)


  n ,d = x_train.shape
  lambda_max = compute_initial_lamb(x_train, y_train)
  initial_model = Lasso(lambda_max)
```

```python
74      initial_model.coordinate_descent(x_train, y_train, np.zeros((d, 1)))
75      initial_w = initial_model.last_w
76
77      lam_list = (lambda_max) * (1 / 2) ** np.arange(0, 20)
78      householdsize_list = []
79      agePct12t29_list = []
80      agePct65up_list = []
81      pctUrban_list = []
82      pctWSocSec_list = []
83      trained_w = []
84      initial_loss_train = np.mean((initial_model.predict(x_train) - y_train)**2)
85      initial_loss_test = np.mean((initial_model.predict(x_test) - y_test)**2)
86
87      loss_train_list = [initial_loss_train]
88      loss_test_list = [initial_loss_test]
89      number_of_nonezero_feature = []
90
91      # lam_list = [30,30]
92      for lam in lam_list[1:]:
93        model = Lasso(lam)
94        model.coordinate_descent(x_train, y_train, initial_w)
95        w_new = model.last_w.copy()
96        number_of_nonezero_feature.append(np.count_nonzero(w_new))
97        householdsize_list.append(w_new[1])
98        agePct12t29_list.append(w_new[3])
99        agePct65up_list.append(w_new[5])
100       pctUrban_list.append(w_new[7])
101       pctWSocSec_list.append(w_new[12])
102
103       loss_train = np.mean((model.predict(x_train) - y_train)**2)
104       loss_train_list.append(loss_train)
105       loss_test = np.mean((model.predict(x_test) - y_test)**2)
106       loss_test_list.append(loss_test)
107
108     selected_coef_history = [householdsize_list, agePct12t29_list,
109     agePct65up_list, pctUrban_list, pctWSocSec_list]
110
111     # A.5 a
112     # Plot lambda against number_of_nonezero_feature
113     plt.plot(lam_list[1:], number_of_nonezero_feature)
114     plt.xscale('log')
115     plt.xlabel("Lambda")
116     plt.ylabel("# of none-zero coef")
117     plt.show()
118
119     # A.5 b
120     # plot 5 different feature change with different lambda
121     features = ['householdsize', 'agePct12t29', 'agePct65up', 'pctUrban', 'pctWSocSec']
122     for i, feature in enumerate(features):
123       plt.plot(lam_list[1:], selected_coef_history[i], label=feature)
124       plt.xscale('log')
125       plt.xlabel("Lambda")
126       plt.ylabel("Loss")
127       plt.legend(loc="upper left")
128     plt.show()
129
130     # A.5 c
131     plt.plot(lam_list, loss_train_list, label="Train")
```

11

```
132    plt.plot(lam_list, loss_test_list, label="Test")
133    plt.xscale('log')
134    plt.xlabel("Lambda")
135    plt.ylabel("Loss")
136    plt.legend(loc="upper left")
137    plt.show()
138
139    coeff_data = pd.DataFrame({"Coeff": (list(w_new.T[0])),
140    "Name": df_train.columns[1:]}).sort_values(["Coeff"], ascending=False)
141    coeff_data[:10]
142    coeff_data[-10:]
```

# A.6

### a.

We know:

$$\mu_i(w, b) = \frac{1}{1 + exp(-y_i(b + x_i^T w))}$$

Rewrite the above equiation:

$$1 + exp(-y_i(b + x_i^T w)) = \frac{1}{\mu_i(w, b)}$$

$$exp(-y_i(b + x_i^T w)) = \frac{1}{\mu_i(w, b)} - 1$$

$$exp(-y_i(b + x_i^T w)) = \frac{1}{\mu_i(w, b)} - \frac{\mu_i(w, b)}{\mu_i(w, b)}$$

$$exp(-y_i(b + x_i^T w)) = \frac{1 - \mu_i(w, b)}{\mu_i(w, b)}$$

Now we have: $exp(-y_i(b + x_i^T w)) = \frac{1 - \mu_i(w,b)}{\mu_i(w,b)}$ Then we compute gradient of L2 logistic regression:
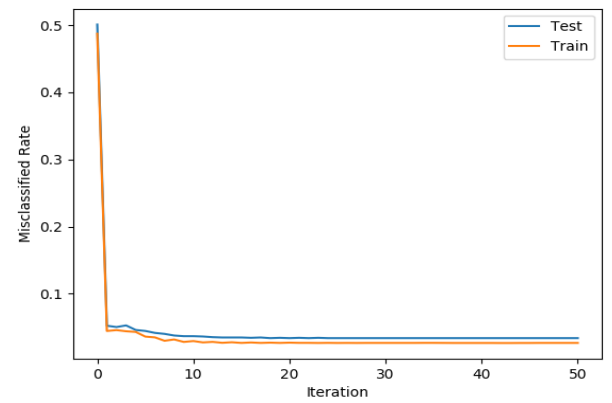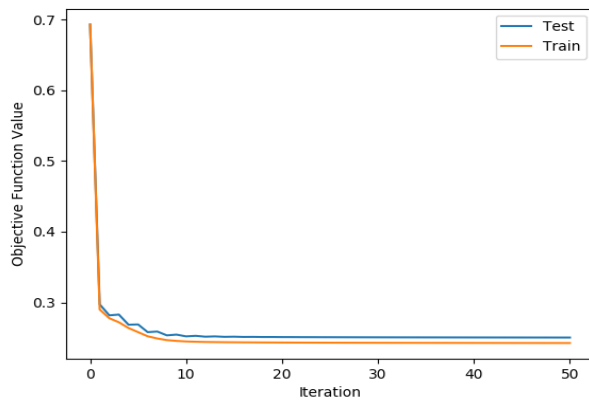To derive the gradient of w:

$$\nabla_w J(w, b) = \frac{1}{n} \sum_{i=1}^n \frac{exp(-y_i(b + x_i^T w))}{1 + exp(-y_i(b + x_i^T w))} * (-y_i x_i^T) + 2\lambda w$$

$$= \frac{1}{n} \sum_{i=1}^n \frac{1 - \mu_i(w, b)}{\mu_i(w, b)} * \mu_i(w, b) * (-y_i x_i^T) + 2\lambda w$$

$$= \frac{1}{n} \sum_{i=1}^n \left(1 - \mu_i(w, b)\right) * (-y_i x_i^T) + 2\lambda w$$
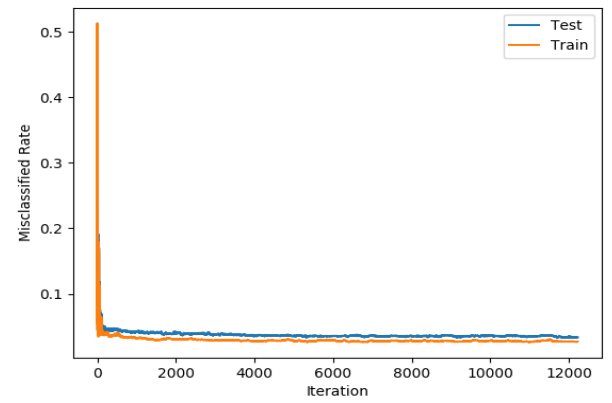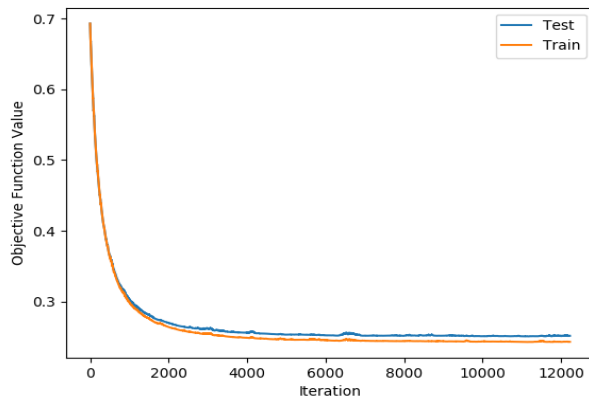
To derive the gradient of b:

$$\nabla_b J(w, b) \frac{1}{n} \sum_{i=1}^n \frac{exp(-y_i(b + x_i^T w))}{1 + exp(-y_i(b + x_i^T w))} * (-y_i)$$

$$= \frac{1}{n} \sum_{i=1}^n \frac{1 - \mu_i(w, b)}{\mu_i(w, b)} * \mu_i(w, b) * (-y_i)$$

$$\frac{1}{n} \sum_{i=1}^n \left(1 - \mu_i(w, b)\right) * (-y_i)$$
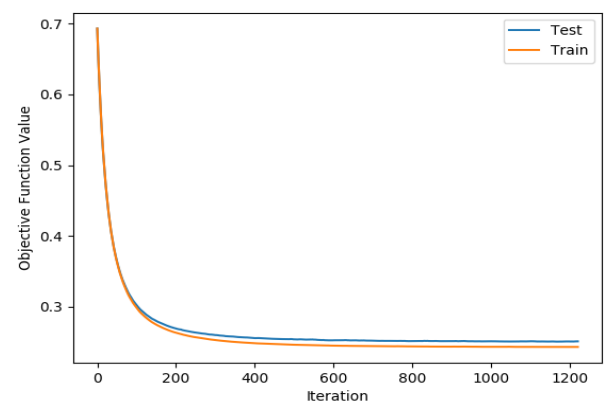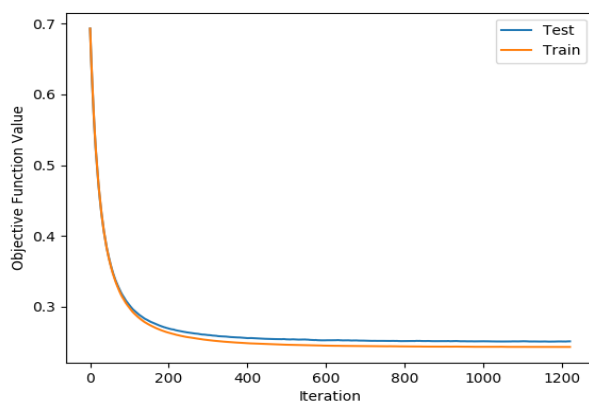
12

**b.**





**c.**





**d.**





```python
1  import numpy as np
2  import scipy.linalg as la
3  import matplotlib.pyplot as plt
4  from mnist import MNIST
5
6  np.random.seed(666)
7  lam = 0.1
```

```python
8
9
10  def load_data():
11    mndata = MNIST('python-mnist/data/')
12    X_train, labels_train = map(np.array, mndata.load_training())
13    X_test, labels_test = map(np.array, mndata.load_testing())
14    X_train = X_train / 255.0
15    X_test = X_test / 255.0
16
17    labels_train = labels_train.astype(np.int16)
18    labels_test = labels_test.astype(np.int16)
19    X_train = np.vstack( (X_train[labels_train == 2], X_train[labels_train == 7]))
20    y_train = np.hstack((labels_train[labels_train == 2], labels_train[labels_train == 7]))
21    y_train[y_train == 2] = -1
22    y_train[y_train == 7] = 1
23
24    X_test = np.vstack((X_test[labels_test == 2], X_test[labels_test == 7]))
25    y_test = np.hstack((labels_test[labels_test == 2], labels_test[labels_test == 7]))
26    y_test[y_test == 2] = -1
27    y_test[y_test == 7] = 1
28
29    print(X_train.shape, y_train.shape, y_train.sum())
30    return (X_train, y_train, X_test, y_test)
31
32
33  def descent(X, Y, w, b, learning_rate=0.01):
34
35    u = 1.0 / (1.0 + np.exp(-Y * (b + X.dot(w))))
36    gradient_b = (-Y * (1 - u)).mean()
37    b -= learning_rate * gradient_b
38
39    u = 1.0 / (1.0 + np.exp(-Y * (b + X.dot(w))))
40    xy = np.multiply(X.T, Y)
41    gradient_w = (- xy * (1 - u)).mean(axis=1) + 2 * lam * w
42    w -= learning_rate * gradient_w
43
44    return (w, b)
45
46
47  def objective_function_value(X, Y, w, b):
48    log_error = np.log(1.0 + np.exp(-Y * (b + X.dot(w))))
49    obj_value = log_error.mean() + L * np.linalg.norm(w, 2)
50
51    predicted = b + X.dot(w)
52    predicted[predicted < 0] = -1
53    predicted[predicted >= 0] = 1
54    correct = np.sum(predicted == Y)
55
56    error = 1.0 - float(correct) / float(X.shape[0])
57
58    return (obj_value, error)
59
60
61  def binary_logistic_regression(X_train, Y_train, X_test, Y_test, learning_rate,
62  epochs, batch=0, save_plt_name="A6"):
63    n, d = X_train.shape
64    w = np.zeros(d)
65    b = 0
```

14

```
66
67    iterations = []
68    test_objective_list = []
69    train_objective_list = []
70    test_error_list = []
71    train_error_list = []
72    obj_train, error = objective_function_value(X_train, Y_train, w, b)
73    test_obj, test_error = objective_function_value(X_test, Y_test, w, b)
74    test_objective_list.append(test_obj)
75    train_objective_list.append(obj_train)
76    test_error_list.append(test_error)
77    train_error_list.append(error)
78    iterations.append(0)
79
80    i = 1
81    for epoch in range(epochs):
82        n, d = X_train.shape
83        ramdom_index = np.random.permutation(n)
84        X_train = X_train[ramdom_index]
85        Y_train = Y_train[ramdom_index]
86        X_list = np.array_split(X_train, n / batch)
87        Y_list = np.array_split(Y_train, n / batch)
88
89
90        for X_split, Y_split in zip(X_list, Y_list):
91            w, b = descent(X_split, Y_split, w, b, learning_rate)
92            obj_train, error = objective_function_value(X_train, Y_train, w, b)
93            obj_test, test_error = objective_function_value(X_test, Y_test, w, b)
94
95            test_objective_list.append(obj_test)
96            train_objective_list.append(obj_train)
97            test_error_list.append(test_error)
98            train_error_list.append(error)
99            iterations.append(i)
100           i += 1
101
102   plt.plot(iterations, test_objective_list, label="Test")
103   plt.plot(iterations, train_objective_list, label="Train")
104   plt.xlabel("Iteration")
105   plt.ylabel("Objective Function Value")
106   plt.legend()
107   plt.savefig("/Users/yinruideng/Desktop/senior_spring/cse546/hw/hw2/latex/"+
108   save_plt_name + "_1.png")
109   plt.show()
110
111   plt.plot(iterations, test_error_list, label="Test")
112   plt.plot(iterations, train_error_list, label="Train")
113   plt.xlabel("Iteration")
114   plt.ylabel("Misclassified Rate")
115   plt.legend()
116   plt.savefig("/Users/yinruideng/Desktop/senior_spring/cse546/hw/hw2/latex/"+
117   save_plt_name + "_2.png")
118   plt.show()
119
120 if __name__ == '__main__':
121   X_train, Y_train, X_test, Y_test = load_data()
122   n, d = X_train.shape
123   print("######  Gradient Descent ######")
```

```
124    binary_logistic_regression(X_train, Y_train, X_test, Y_test, 0.5, 50, batch=n,
125    save_plt_name="A6_b")
126    print("#######  Stochastic Gradient Descent  #######")
127    binary_logistic_regression( X_train, Y_train, X_test, Y_test, 0.001, 1, batch=1,
128    save_plt_name="A6_c")
129    print("#######  Mini Batch Gradient Descent  #######")
130    binary_logistic_regression( X_train, Y_train, X_test, Y_test, 0.01, 10, batch=100,
131    save_plt_name="A6_d")
```

# B.4

## a.

In this problem, in order to compute the gradient vector of the weights, we could compute individual gradient separately and then put them into an vector.

From the question we can rewrite the gradient equation and add softmax to it.

$$\nabla_w \mathcal{L}(W) = -\sum_{i=1}^{n} x_i (y_i - \frac{exp(W^T x_i)}{\sum_{j=i}^{k} exp(W^T x_i)})^T$$

To prove the above equation, we compute the gradients individually:

$$\mathcal{L}(W) = -\sum_{i=1}^{n} \sum_{\ell=i}^{k} \mathbb{1}\{y_i = \ell\} log(\frac{exp(W^T x_i)}{\sum_{j=i}^{k} exp(W^T x_i)})$$

$$\nabla_w^t \mathcal{L}(W) = -\sum_{i=1}^{n} \nabla_w^t \Big( \mathbb{1}\{y_i = t\} log(\frac{exp(W^T x_i)}{\sum_{j=i}^{k} exp(W^T x_i)}) \Big)$$

Then we take derivative, and note the following because derivative other than $w^j = w^t$ equals to 0:

$$\frac{\partial(\sum_{j=1}^{k} exp(w^{(j)} x_i)))}{\partial w^t} = exp(w^{(t)} x_i)) x_i$$

Then with this in mind:

$$\nabla_{w^t} \mathcal{L}(W) = -\sum_{i=1}^{n} \Big( \mathbb{1}\{y_i = t\}(\frac{\sum_{j=i}^{k} exp(w^{(t)} x_i)}{exp(w^{(t)} x_i)}) \frac{exp(w^{(t)} x_i) * \sum_{j=i}^{k} exp(w^{(t)} x_i) x_i - exp(w^{(t)} x_i) exp(W^T x_i) x_i}{\Big( \sum_{j=i}^{k} exp(w^{(t)} x_i) \Big)^2} \Big)$$

Then simplify the equation:

$$\nabla_{w^t} \mathcal{L}(W) = -\sum_{i=1}^{n} \Big( \mathbb{1}\{y_i = t\} x_i \frac{\sum_{j=i}^{k} exp(w^{(t)} x_i) - exp(w^{(t)} x_i)}{\sum_{j=i}^{k} exp(w^{(t)} x_i)} \Big)$$

$$\nabla_{w^t} \mathcal{L}(W) = -\sum_{i=1}^{n} x_i \Big( \frac{\mathbb{1}\{y_i = t\} \sum_{j=i}^{k} exp(w^{(t)} x_i) - \mathbb{1}\{y_i = t\} exp(w^{(t)} x_i)}{\sum_{j=i}^{k} exp(w^{(t)} x_i)} \Big)$$

$$\nabla_{w^t} \mathcal{L}(W) = -\sum_{i=1}^{n} x_i \Big( \mathbb{1}\{y_i = t\} - \frac{exp(w^{(t)} x_i)}{\sum_{j=i}^{k} exp(w^{(t)} x_i)} \Big)$$

Then put individual gradient together.

$$\nabla_W \mathcal{L}(W) = \Big[ \sum_{i=1}^{n} x_i \Big( \mathbb{1}\{y_i = 1\} - \frac{exp(w^{(t)} x_i)}{\sum_{j=i}^{k} exp(w^{(1)} x_i)} \Big), \dots, \sum_{i=1}^{n} x_i \Big( \mathbb{1}\{y_i = n\} - \frac{exp(w^{(t)} x_i)}{\sum_{j=i}^{k} exp(w^{(1)} x_i)} \Big) \Big]$$

$$\nabla_W \mathcal{L}(W) = \sum_{i=1}^{n} x_i \Big[ \mathbb{1}\{y_i = 1\} - \frac{exp(w^{(1)} x_i)}{\sum_{j=i}^{k} exp(w^{(j)} x_i)}, \dots, \mathbb{1}\{y_i = n\} - \frac{exp(w^{(n)} x_i)}{\sum_{j=i}^{k} exp(w^{(j)} x_i)} \Big]$$

$$\nabla_W \mathcal{L}(W) = \sum_{i=1}^{n} x_i (\mathbf{y_i} - \hat{\mathbf{y}}_i^{(w)})^T$$

16

**b.**

$$\mathbf{J}(W) = \frac{1}{2} \sum_{i=1}^{n} ||\mathbf{y}_i - W^T x_i||_2^2$$

And because of:

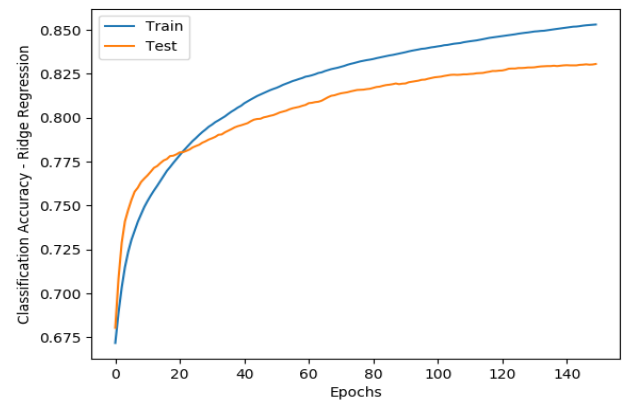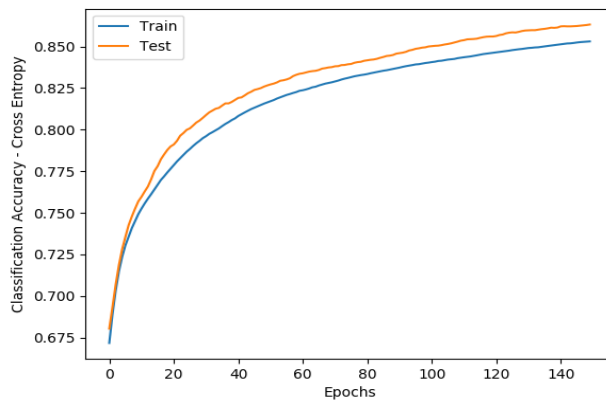$$\tilde{\mathbf{y}}_i^{(w)} = W^T x_i$$

So:

$$\nabla \mathbf{J}(W) = -\sum_{j=1}^{n} x_i (\mathbf{y}_i - W^T x_i)^T$$

$$\nabla \mathbf{J}(W) = -\sum_{j=1}^{n} x_i (\mathbf{y}_i - \tilde{\mathbf{y}}_i^{(w)})^T$$

**c.**

When epoch goes up, accuracy goes up as well. I tries different step size and find 0.05 works fine for me. When using 0.01 step size, I got lower accuracy.



```python
import torch
import numpy as np
import matplotlib.pyplot as plt
from mnist import MNIST

def load_data():
    mndata = 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
    X_test = X_test / 255.0
    return X_train, labels_train, X_test, labels_test
X_train, y_train, X_test, y_test = load_data()

# One hot encoding
def one_hot(y_train_, m):
    n = len(y_train_)
    reformed_tensor = torch.zeros(n, m)
    for i in range(n):
        index = y_train_[i]
        reformed_tensor[i][index] = 1
    return reformed_tensor


# convert to tensor
```

```
26    X_train_ = torch.tensor(X_train, dtype=torch.double)
27    y_train_ = torch.tensor(y_train, dtype=torch.int64)
28
29    X_test_ = torch.tensor(X_test, dtype=torch.double)
30    y_test_ = torch.tensor(y_test, dtype=torch.int64)
31
32    W = torch.zeros(784, 10, requires_grad=True, dtype=torch.double)
33    W_mse = torch.zeros(784, 10, requires_grad=True, dtype=torch.double)
34    step_size = 0.01
35    epochs = 100
36    train_accuracy_list = []
37    test_accuracy_list = []
38    train_accuracy_list_mse = []
39    test_accuracy_list_mse = []
40
41    epochs = list(range(epochs))
42    for epoch in epochs:
43      print("Epoch: ", epoch)
44      y_hat = torch.matmul(X_train_, W)
45      y_hat_mse = torch.matmul(X_train_, W_mse)
46      # cross entropy combines softmax calculation with NLLLoss
47      loss = torch.nn.functional.cross_entropy(y_hat, y_train_)
48      loss_mse = torch.nn.functional.mse_loss(y_hat_mse, one_hot(y_train_, 10).double())
49      # computes derivatives of the loss with respect to W
50      loss.backward()
51      loss_mse.backward()
52      # gradient descent update
53      W.data = W.data - step_size * W.grad
54      W_mse.data = W_mse.data - step_size * W_mse.grad
55      # .backward() accumulates gradients into W.grad instead
56      # of overwriting, so we need to zero out the weights
57
58      # Cross Entropy
59      max_index_train = torch.max((torch.matmul(X_train_, W)), dim=1).indices.numpy()
60      num_corrected_prediction_train = sum(max_index_train == y_train)
61      train_accu = num_corrected_prediction_train / len(y_train)
62      train_accuracy_list.append(train_accu)
63
64      max_index_test = torch.max((torch.matmul(X_test_, W)), dim=1).indices.numpy()
65      num_corrected_prediction_test = sum(max_index_test == y_test)
66      test_accu = num_corrected_prediction_test / len(y_test)
67      test_accuracy_list.append(test_accu)
68
69      # MSE
70      max_index_train_mse = torch.max((torch.matmul(X_train_, W_mse)), dim=1).indices.numpy()
71      num_corrected_prediction_train_mse = sum(max_index_train_mse == y_train)
72      train_accu_mse = num_corrected_prediction_train / len(y_train)
73      train_accuracy_list_mse.append(train_accu_mse)
74
75      max_index_test_mse = torch.max(torch.matmul(X_test_, W_mse), dim=1).indices.numpy()
76      num_corrected_prediction_test_mse = sum(max_index_test_mse == y_test)
77      test_accu_mse = num_corrected_prediction_test_mse / len(y_test)
78      test_accuracy_list_mse.append(test_accu_mse)
79
80
81
82      print("Train Accuracy: ", train_accu)
83      print("Test Accuracy: ", test_accu)
```

18

```
84
85    W.grad.zero_()
86    W_mse.grad.zero_()
87
88  plt.plot(epochs, train_accuracy_list, label="Train")
89  plt.plot(epochs, test_accuracy_list, label="Test")
90  plt.xlabel("Epochs")
91  plt.ylabel("Classification Accuracy - Cross Entropy")
92  plt.legend()
93  plt.savefig("/Users/yinruideng/Desktop/senior_spring/cse546/hw/hw2/latex/B4_c_1.png")
94  plt.show()
95
96  plt.plot(epochs, train_accuracy_list_mse, label="Train")
97  plt.plot(epochs, test_accuracy_list_mse, label="Test")
98  plt.xlabel("Epochs")
99  plt.ylabel("Classification Accuracy - Ridge Regression")
100 plt.legend()
101 plt.savefig("/Users/yinruideng/Desktop/senior_spring/cse546/hw/hw2/latex/B4_c_2.png")
102 plt.show()
```