

CSE446 HW2

Brian Kang[‡]

Conceptual Questions

A.0

- a. Suppose that your estimated model for predicting house prices has a large positive weight on 'number of bathrooms'. Does it imply that if we remove the feature "number of bathrooms" and refit the model, the new predictions will be strictly worse than before? Why?
 - Not necessarily. There is a possibility that removing that feature would decrease complexity of the model where we already had serious overfitting. This can lead to better predictions. But, given that we are predicting house prices, keeping "number of bathrooms" sounds like a good idea to me.
- b. 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?
 - The L2 norm penalizes "more smoothly" in a sense. Recall the graph with a diamond constraint region for L1 and circle for L2 and a SSR contour plot for each. For L1 norm, we keep the y-axis β_2 value while the x-axis $\beta_1 = 0$ with the penalization (leading to many 0s in the weight vector). But, for the L2 norm, β_1 is close to zero and β_2 is close to the apex of the circle.
- c. In at most one sentence each, state one possible upside and one possible downside of using the following regularizer: $\sum_i |w_i|^{0.5}$
 - Like the lasso, the $L(1/2)$ norm can help with feature selection by shrinking some weights to zero; however, the downside could be pushing too many weights down to zero since the penalty constraint region is smaller than L1.
- d. True or False: If the step-size for gradient descent is too large, it may not converge.
 - True. By keeping the step-size too large, we may head in the direction of the optimal point, but just largely bypass it and keep repeating it.
- e. In your own words, describe why SGD works.
 - SGD works in the sense that we don't need all of the data points to get an idea of the direction towards the minimum point. Getting a general direction from one point each step at a time is sufficient.

*Collaborated with Cindy Wu

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

- f. 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.
- Benefit: Randomly choosing one point to get a general direction and repeating that many times to approach the minimum is very cost efficient compared to GD.
 - Cost: SGD naturally involves more noise and so may require much more steps/updates than GD to approach the optimal value.

A.1

a. *Proof.*

- i $\forall i \in \{1, \dots, n\}, |x_i| \geq 0$, so $f(x) = \sum_{i=1}^n |x_i| \geq 0$. If $|x_i| = 0 \forall i \in \{1, \dots, n\}$, then $f(x) = 0$.
- ii $f(ax) = \sum_{i=1}^n |ax_i| = |a| \sum_{i=1}^n |x_i| = |a|f(x)$
- iii $f(x+y) = \sum_{i=1}^n |x_i + y_i| \leq \sum_{i=1}^n |x_i| + \sum_{i=1}^n |y_i| = f(x) + f(y)$

Therefore, $f(x) = \sum_{i=1}^n |x_i|$ is a norm. □

b. *Proof.* Let $x, y \in \mathbb{R}^2, x, y \geq 0$. Then,

$$\begin{aligned}
 g(x+y) &= (|x_1 + y_1|^{1/2} + |x_2 + y_2|^{1/2})^2 \\
 &= |x_1 + y_1| + |x_2 + y_2| + 2 * |x_1 + y_1|^{1/2} * |x_2 + y_2|^{1/2} \\
 &= |x_1 + x_2| + |y_1 + y_2| + 2 * |x_1 + y_1|^{1/2} * |x_2 + y_2|^{1/2} \\
 &= (|x_1 + x_2|^{1/2})^2 + (|y_1 + y_2|^{1/2})^2 + 2 * |x_1 + y_1|^{1/2} * |x_2 + y_2|^{1/2} \\
 &\not\leq (|x_1 + x_2|^{1/2})^2 + (|y_1 + y_2|^{1/2})^2 = g(x) + g(y)
 \end{aligned}$$

Thus, $\exists x, y \geq 0$ such that $g(x+y) \not\leq g(x) + g(y)$, i.e., $g(x) = \sum_{i=1}^n (|x_i|^{1/2})^2$ is not a norm. □

A.2

I is not convex because $\lambda b + (1 - \lambda)c \notin I$.

II is convex.

III is not convex because $\lambda a + (1 - \lambda)d \notin III$.

A.3

- a. I is convex.
- b. II is not convex because $f(\lambda a + (1 - \lambda)b) \not\leq \lambda f(a) + (1 - \lambda)f(b)$.
- c. III is not convex on $[a, d]$ because $f(\lambda a + (1 - \lambda)c) \not\leq \lambda f(a) + (1 - \lambda)f(c)$.
- d. III is convex on $[c, d]$.

Programming Problems

LASSO: Code is all together

```
1 # Brian Kang
2 import numpy as np
3 import matplotlib.pyplot as plt
4 import pandas as pd
5
6 # make synthetic data
7 n = 500 # row
8 d = 1000 # column
9 k = 100 # ???
10 sd = 1 # standard normal sd
11 X = np.random.randn(n, d) # sample from standard normal dist
12 error = np.random.randn(n) # sample from standard normal dist
13 wj = np.empty(d) # made up weights
14 for j in range(d):
15     if j < k:
16         wj[j] = j / k
17     else:
18         wj[j] = 0
19 Y = np.dot(X, wj) + error # synthetic data
20
21
22 # get smallest lambda that shrinks all w to 0
23 def lambdaMax(X, Y):
24     ybar = Y - np.mean(Y)
25     return np.amax(2 * np.absolute(np.dot(X.T, ybar)))
26
27
28 # do coordinate descent algorithm for lasso
29 def coordDescentLasso(X, Y, w0, lmbda, tol):
30     n, d = X.shape
31     a = np.zeros(d)
32     c = np.zeros(d)
33     wNew = np.copy(w0) # copy of initial weights, changing one vars doesn't
34     # affect the other
35     doCoordDescent = True # while not converged, do...
36     while doCoordDescent:
37         wOld = np.copy(wNew) # after first loop, wNew is now wOld
38         b = np.sum(Y - np.dot(X, wNew)) / n # get b value
39         for k in range(d):
40             a[k] = 2 * np.sum(np.power(X[:, k], 2))
41             resid = Y - (b + np.dot(X, wNew) - X[:, k] * wNew[k]) # subtract
42             # since we sum_{j!=k}
43             c[k] = 2 * np.sum(X[:, k] * resid)
44             # update weights
45             if c[k] < -lmbda:
46                 wNew[k] = (c[k] + lmbda) / a[k]
47             elif c[k] > lmbda:
48                 wNew[k] = (c[k] - lmbda) / a[k]
49             else:
50                 wNew[k] = 0
```

```

49         doCoordDescent = (np.amax(np.absolute(wNew - wOld)) > tol) # false if
error < tolerance
50     return wNew, b
51
52
53 def problem4():
54     w0 = np.random.randn(d) # make up starting weights
55     # keep weights, b values
56     keepW = [] # list, not np.array
57     keepB = []
58
59     # Compute lasso along regularization path
60     lambdaMax = lambdaMax(X, Y) # Largest penalty
61     # print(lambdaMax)
62     numlasso = 20 # number of lasso's to run. increase until regPath[numlasso]
close to 0
63     regPath = [lambdaMax] * numlasso
64     for i in range(numlasso):
65         regPath[i] = regPath[i] / np.power(1.5, i)
66     # print(regPath)
67     for lambda in regPath:
68         w, b = coordDescentLasso(X, Y, w0, lambda, 0.1)
69         keepW.append(w)
70         keepB.append(b)
71
72     # a)
73     keepFeat = []
74     for w in keepW:
75         keepFeat.append(np.sum(w > 0)) # count how many features kept
76
77     plt.figure(0)
78     plt.plot(regPath, keepFeat)
79     plt.xlabel("$\lambda$")
80     plt.ylabel("Number of Non-Zero Features")
81     plt.xscale('log')
82     plt.savefig('hw2_1.png')
83
84     # b)
85     FDR = [] # false discovery rate
86     TPR = [] # true positive rate
87     for w in keepW:
88         FDR.append(np.sum(w[k:] > 0) / (d - k))
89         TPR.append(np.sum(w[:k] > 0) / k)
90
91     plt.figure(1)
92     plt.scatter(FDR, TPR)
93     plt.xlabel("False Discovery Rate")
94     plt.ylabel("True Positive Rate")
95     plt.savefig('hw2_2.png')
96
97
98 def problem5():
99     df_train = pd.read_table("crime-train.txt") # load datasets
100     df_test = pd.read_table("crime-test.txt")
101

```

```

102 d = 95 # number of features (columns)
103 Xbody = df_train.iloc[:, 1:96].copy()
104 X_train = Xbody.values # get a copy
105 # print(X_train)
106 X_test = df_test.iloc[:, 1:96].copy().values
107 Y_train = df_train['ViolentCrimesPerPop'].values # convert to numpy array
108 # print(Y_train)
109 Y_test = df_test['ViolentCrimesPerPop'].values
110
111 w0 = np.zeros(d) # initial weights
112 # keep weights, b values
113 keepW = [] # list, not np.array
114 keepB = []
115
116 # Compute lasso along regularization path
117 lambdaMax = lambdaMax(X_train, Y_train) # Largest penalty
118 lambda = 1 # temporary arbitrary number
119 regPath = []
120 i = 0
121 while lambda > 0.01: # keep doing lasso until lambda <= 0.01
122     lambda = lambdaMax / np.power(2, i)
123     # print(lambda, ", ", ", ", i)
124     regPath.append(lambda)
125     w, b = coordDescentLasso(X_train, Y_train, w0, lambda, 0.1)
126     w0 = w # use new weights as the weights for next lasso
127     keepW.append(w)
128     keepB.append(b)
129     i += 1
130
131 # a)
132 keepFeat = []
133 for w in keepW:
134     keepFeat.append(np.sum(w > 0)) # count how many features kept
135
136 plt.figure(2)
137 plt.plot(regPath, keepFeat)
138 plt.xlabel("$\lambda$")
139 plt.ylabel("Number of Non-Zero Features")
140 plt.xscale('log')
141 plt.savefig('hw2_3.png')
142
143 # b)
144 # get coefficients of specific vars from the lassos
145 vars = ["agePct12t29", "pctWSocSec", "pctUrban", "agePct65up", "householdsize"]
146 for i in vars:
147     indx = Xbody.columns.get_loc(i)
148     getCoeff = [row[indx] for row in keepW]
149     plt.figure(3)
150     line, = plt.plot(regPath, getCoeff)
151     line.set_label(i)
152
153 plt.xlabel("$\lambda$")
154 plt.ylabel("Coefficients")
155 plt.xscale('log')

```

```

156 plt.legend()
157 plt.savefig('hw2_4.png')
158
159 # c)
160 # get squared errors from the lassos
161 errorTrain = [np.mean(np.square(np.dot(X_train, w) + b - Y_train)) for w, b
in zip(keepW, keepB)]
162 errorTest = [np.mean(np.square(np.dot(X_test, w) + b - Y_test)) for w, b in
zip(keepW, keepB)]
163
164 plt.figure(4)
165 plt.plot(regPath, errorTrain, label="train")
166 # line1.set_label("Train")
167 plt.plot(regPath, errorTest, label="test")
168 # line2.set_label("Test")
169 plt.xlabel("$\lambda$")
170 plt.ylabel("MSE")
171 plt.xscale('log')
172 plt.legend()
173 plt.savefig('hw2_5.png')
174
175 # d)
176 # print(regPath)
177 indx = 4 # index where lambda approx = 30
178 getCoeff = [row[indx] for row in keepW] # the weights
179 print("Most Positive Lasso Coefficient: ", np.max(getCoeff))
180 # Out: "Most Positive Lasso Coefficient: 0.005256729455865754"
181 maxindx = np.argmax(getCoeff)
182 print("Corresponding Variable:", Xbody.columns[maxindx])
183 # Out: "perCapInc"
184 print("Most Negative Lasso Coefficient: ", np.min(getCoeff))
185 # Out: "Most Negative Lasso Coefficient: 0.0"
186 minindx = np.argmin(getCoeff)
187 print("Corresponding Variable:", Xbody.columns[minindx])
188 # Out: "population"
189
190
191 problem4()
192 problem5()

```

Plots and solutions are below!

A.4

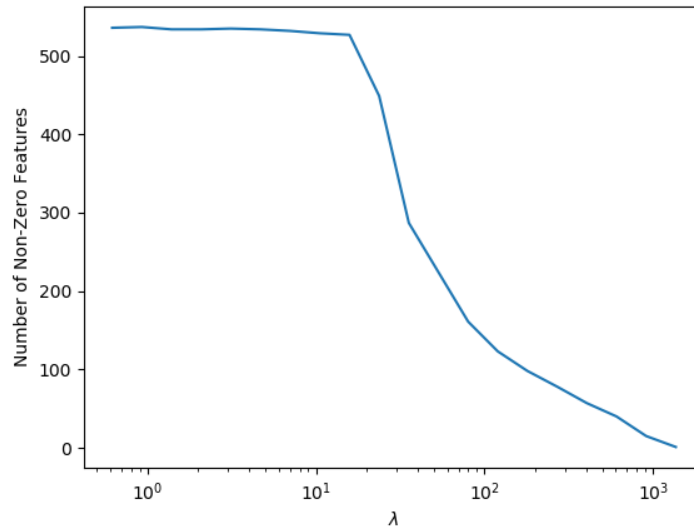


Figure 1: Part a code is above.

a.

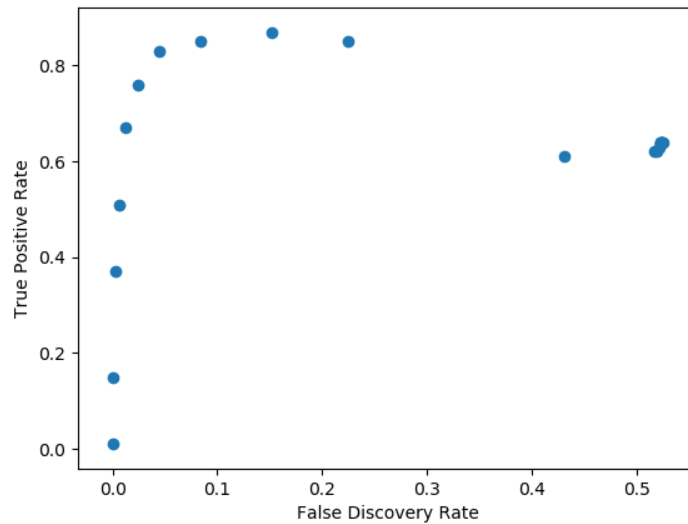


Figure 2: Part b code is above.

b.

- c. We observe that as lambda decreases by linear fractions (from right to left in figure 1), we move from left to right in FDR v.s. TPR (figure 2). The sharp decrease in figure 1 and sharp increase in figure 2 towards ideal (upper left-hand corner) are correlated to each other, as the lasso is at

work with feature selection with respect to various regularization lambda values. But as we go further, FDR increases rapidly and TPR even decreases, which is reasonable since the number of non-zero features kept increase in a decreasing rate.

A.5

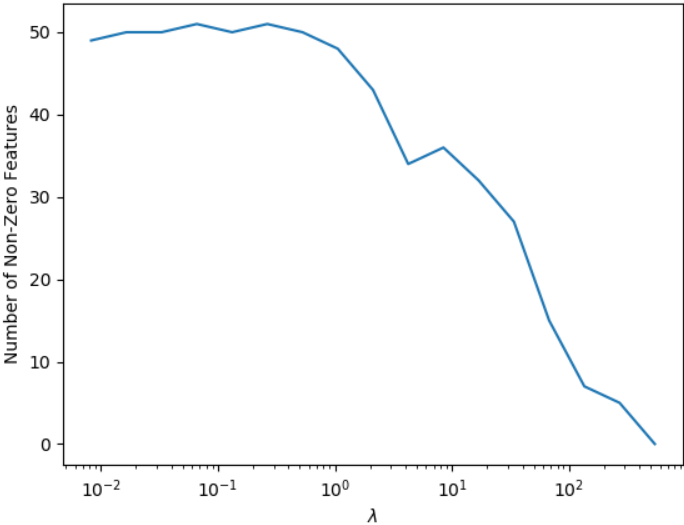


Figure 3: Part a code is above.

a.

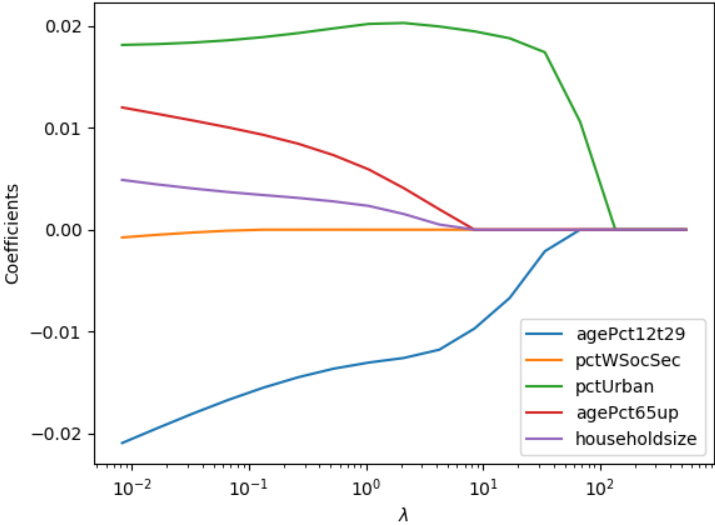


Figure 4: Part b code is above.

b.

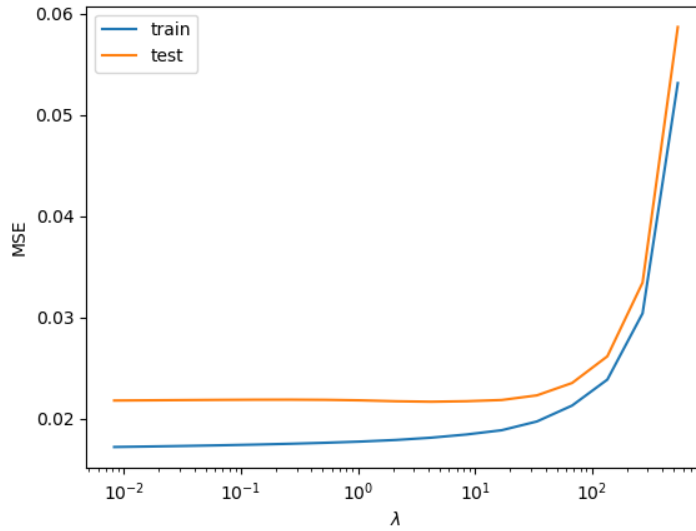


Figure 5: Part c code is above.

- c.
- d. We found out that the variable with the most positive lasso coefficient = 0.00525673 is **per-CapInc**, and the variable with the most "negative" lasso coefficient = 0.0 is **population**. Since we are predicting Per Capita Violent Crimes, it is reasonable to say that Per Capita Income is a important feature because usually there is a correlation between crime rates and income. But when it comes to population, the lasso shrunk the coefficient to zero because Per Capita Violent Crimes is probably calculated by dividing population. The machine doesn't know this so population is treated as an insignificant feature for prediction.
- e. This method of reasoning is incorrect and it is known as Reverse Causality. Just because there is a correlation, we cannot conclude with a solution resulting in "if X, then Y" because there is always the case of "if Y, then X" for correlations. And in cases like this (and the fire truck/burning building example) this is what is happening. For causal inference we need a clear experiment design and domain knowledge to know that "fire trucks don't cause fire."

LOGIT: Code is all together

```
1 # Brian Kang
2 import numpy as np
3 import matplotlib.pyplot as plt
4 from mnist import MNIST
5
6
7 def load_dataset():
8     mndata = MNIST('./data')
9     mndata.gz = True
10    X_train, labels_train = map(np.array, mndata.load_training())
11    X_test, labels_test = map(np.array, mndata.load_testing())
12    X_train = X_train / 255.0
13    X_test = X_test / 255.0
14
15    # get index array of 0 or 1
16    keep_train = np.bitwise_or(labels_train == 2, labels_train == 7)
17    keep_test = np.bitwise_or(labels_test == 2, labels_test == 7)
18
19    # delete columns that don't give 2 or 7 response
20    X_train = X_train[keep_train]
21    labels_train = labels_train[keep_train].astype(int)
22    X_test = X_test[keep_test]
23    labels_test = labels_test[keep_test].astype(int)
24
25    # binary classify to -1, 1
26    labels_train[labels_train == 2] = -1
27    labels_train[labels_train == 7] = 1
28    labels_test[labels_test == 2] = -1
29    labels_test[labels_test == 7] = 1
30    return X_train, labels_train, X_test, labels_test
31
32
33 x_train, bin_train, x_test, bin_test = load_dataset() # load in data
34
35
36 # calculate mu
37 def mu(w, b, X, Y):
38     return 1 / (1 + np.exp(-Y * (b + X.T @ w)))
39
40
41 # get the cost function J
42 def J(w, b, X, Y, lmbda):
43     _, n = X.shape
44     return np.sum(np.log(1 / mu(w, b, X, Y))) / n + lmbda * np.linalg.norm(w) **
45     2
46
47 # gradient of w on J (dell w J)
48 def dwJ(w, b, X, Y, lmbda):
49     d, n = X.shape
50     mu_i = mu(w, b, X, Y)
51     # do matrix algebra much faster
52     # https://stackoverflow.com/questions/26089893/understanding-numpys-einsum
```

```

53     return np.einsum("i, i, di -> d", (1 - mu_i), -Y, X) / n + 2 * lambda * w
54
55
56 # gradient of b on J (dell b J)
57 def dbJ(w, b, X, Y):
58     d, n = X.shape
59     mu_i = mu(w, b, X, Y)
60     # do element wise multiplication faster
61     return np.einsum("i, i -> ", (1 - mu_i), -Y) / n
62
63
64 # gradient descent
65 def gd(w0, b0, X, Y, lambda, stepSize, iter):
66     # step size = learning rate
67     keepJ = []
68     # start will all zeros
69     keepW = np.zeros((iter, len(w0)))
70     keepB = np.zeros(iter)
71     keepW[0] = w0 # 0th row is the starting weights
72     keepB[0] = b0 # 0th row is the starting intercept
73     keepJ.append(J(w0, b0, X, Y, lambda)) # keep the 0th cost J
74     for i in range(1, iter):
75         # new = old - step size * gradient
76         keepW[i] = keepW[i - 1] - stepSize * dwJ(keepW[i - 1], keepB[i - 1], X,
Y, lambda)
77         # print(i,":")
78         # print("keepW: ", keepW)
79         keepB[i] = keepB[i - 1] - stepSize * dbJ(keepW[i - 1], keepB[i - 1], X,
Y)
80         # print("keepB: ", keepB)
81         keepJ.append(J(keepW[i], keepB[i], X, Y, lambda))
82         # print("keepJ: ", keepJ)
83         i += 1
84     return keepW, keepB, keepJ
85
86
87 # classification error
88 def classifier(w, b, X, Y):
89     d, n = X.shape
90     levelx = np.sign(b + X.T @ w) > 0 # classify to 1 or -1
91     levely = Y == 1
92     return 1 - np.sum(np.equal(levelx, levely) / n)
93
94
95 def sgd(w0, b0, X, Y, lambda, stepSize, iter, batchSize):
96     keepJ = []
97     # start will all zeros
98     keepW = np.zeros((iter, len(w0)))
99     keepB = np.zeros(iter)
100    keepW[0] = w0 # 0th row is the starting weights
101    keepB[0] = b0 # 0th row is the starting intercept
102    keepJ.append(J(w0, b0, X, Y, lambda)) # keep the 0th cost J
103    for i in range(1, iter):
104        batch = np.random.randint(0, len(Y), batchSize) # random index
selection

```

```

105         # new = old - step size * gradient
106         keepW[i] = keepW[i - 1] - stepSize * dwJ(keepW[i - 1], keepB[i - 1], X
[: , batch], Y[batch], lambda)
107         keepB[i] = keepB[i - 1] - stepSize * dbJ(keepW[i - 1], keepB[i - 1], X
[: , batch], Y[batch])
108         keepJ.append(J(keepW[i], keepB[i], X, Y, lambda))
109         i += 1
110     return keepW, keepB, keepJ
111
112
113 def problem6b():
114     d, _ = x_train.T.shape
115     w0 = np.zeros(d)
116     b0 = 0
117     lambda = 0.1
118     iter = 100
119     trainW, trainB, trainJ = gd(w0, b0, x_train.T, bin_train, lambda, 0.05, iter)
120     testW, testB, testJ = gd(w0, b0, x_test.T, bin_test, lambda, 0.05, iter)
121
122     plt.figure(0)
123     plt.plot(range(0, iter), trainJ, label="train")
124     plt.plot(range(0, iter), testJ, label="test")
125     plt.xlabel("Iteration")
126     plt.ylabel("Cost J(w,b)")
127     plt.legend()
128     plt.savefig('hw2_6.png')
129
130     train_error = [classifier(w, b, x_train.T, bin_train) for w, b in zip(trainW
, trainB)]
131     test_error = [classifier(w, b, x_test.T, bin_test) for w, b in zip(testW,
testB)]
132     plt.figure(1)
133     plt.plot(range(0, iter), train_error, label="train")
134     plt.plot(range(0, iter), test_error, label="test")
135     plt.xlabel("Iteration")
136     plt.ylabel("Misclassification Error")
137     plt.legend()
138     plt.savefig('hw2_7.png')
139
140
141 def problem6c():
142     d, _ = x_train.T.shape
143     w0 = np.zeros(d)
144     b0 = 0
145     lambda = 0.1
146     iter = 100
147     batch = 1
148     trainW, trainB, trainJ = sgd(w0, b0, x_train.T, bin_train, lambda, 0.05, iter
, batch)
149     testW, testB, testJ = sgd(w0, b0, x_test.T, bin_test, lambda, 0.05, iter,
batch)
150
151     plt.figure(2)
152     plt.plot(range(0, iter), trainJ, label="train")
153     plt.plot(range(0, iter), testJ, label="test")

```

```

154     plt.xlabel("Iteration")
155     plt.ylabel("Cost J(w,b)")
156     plt.legend()
157     plt.savefig('hw2_8.png')
158
159     train_error = [classifier(w, b, x_train.T, bin_train) for w, b in zip(trainW
, trainB)]
160     test_error = [classifier(w, b, x_test.T, bin_test) for w, b in zip(testW,
testB)]
161     plt.figure(3)
162     plt.plot(range(0, iter), train_error, label="train")
163     plt.plot(range(0, iter), test_error, label="test")
164     plt.xlabel("Iteration")
165     plt.ylabel("Misclassification Error")
166     plt.legend()
167     plt.savefig('hw2_9.png')
168
169
170 def problem6d():
171     d, _ = x_train.T.shape
172     w0 = np.zeros(d)
173     b0 = 0
174     lambda = 0.1
175     iter = 100
176     batch = 100
177     trainW, trainB, trainJ = sgd(w0, b0, x_train.T, bin_train, lambda, 0.05, iter
, batch)
178     testW, testB, testJ = sgd(w0, b0, x_test.T, bin_test, lambda, 0.05, iter,
batch)
179
180     plt.figure(4)
181     plt.plot(range(0, iter), trainJ, label="train")
182     plt.plot(range(0, iter), testJ, label="test")
183     plt.xlabel("Iteration")
184     plt.ylabel("Cost J(w,b)")
185     plt.legend()
186     plt.savefig('hw2_10.png')
187
188     train_error = [classifier(w, b, x_train.T, bin_train) for w, b in zip(trainW
, trainB)]
189     test_error = [classifier(w, b, x_test.T, bin_test) for w, b in zip(testW,
testB)]
190     plt.figure(5)
191     plt.plot(range(0, iter), train_error, label="train")
192     plt.plot(range(0, iter), test_error, label="test")
193     plt.xlabel("Iteration")
194     plt.ylabel("Misclassification Error")
195     plt.legend()
196     plt.savefig('hw2_11.png')
197
198
199 problem6b()
200 problem6c()
201 problem6d()

```

Plots and solutions are below!

A.6

a. *Proof.*

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

Or equivalently, we can get,

$$\begin{aligned}
\nabla_w J(w, b) &= \frac{1}{n} \sum_{i=1}^n \frac{\exp(-y_i(b + x_i^T w))}{1 + \exp(-y_i(b + x_i^T w))} * -y_i x_i + 2\lambda w \\
&= \frac{1}{n} \sum_{i=1}^n \frac{1 + \exp(-y_i(b + x_i^T w)) - 1}{1 + \exp(-y_i(b + x_i^T w))} * -y_i x_i + 2\lambda w \\
&= \frac{1}{n} \sum_{i=1}^n \left(1 - \frac{1}{1 + \exp(-y_i(b + x_i^T w))}\right) * -y_i x_i + 2\lambda w \\
&= \frac{1}{n} \sum_{i=1}^n (1 - \mu_i(w, b)) * -y_i x_i + 2\lambda w
\end{aligned}$$

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

□

b. Part b plots for gradient descent.

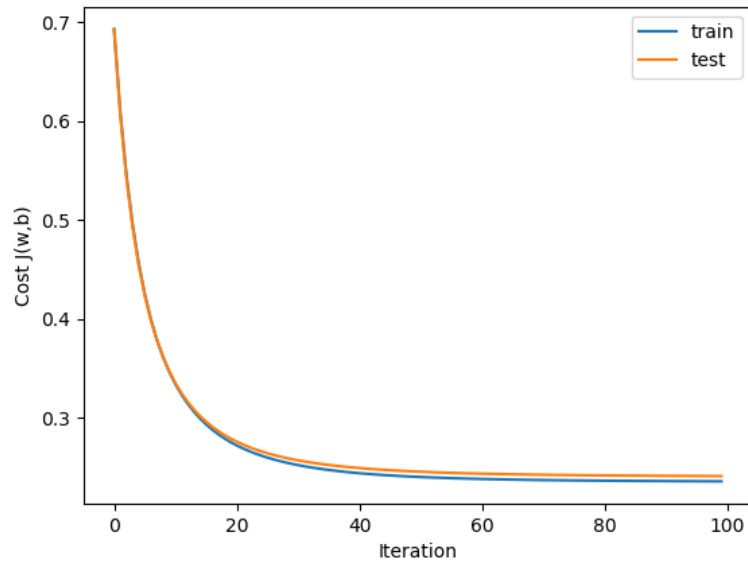


Figure 6: Part i ($J(w,b)$ Cost) code is above.

i

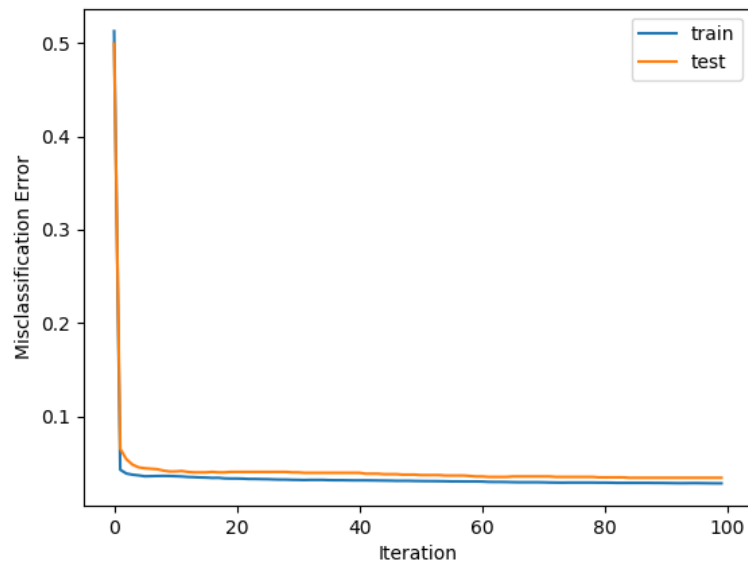


Figure 7: Part ii (Misclassification Error) code is above.

ii

c. Part c plots for stochastic gradient descent with batch size=1.

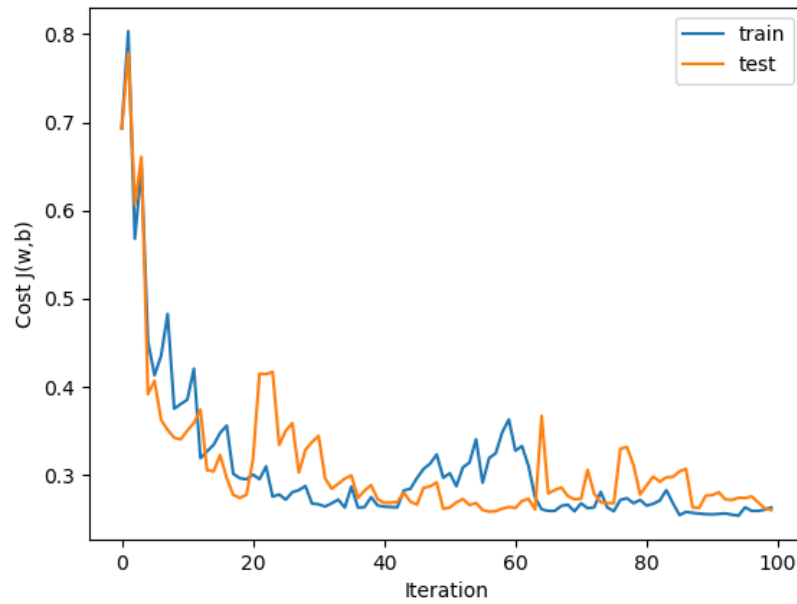


Figure 8: $J(w,b)$ Cost code is above.

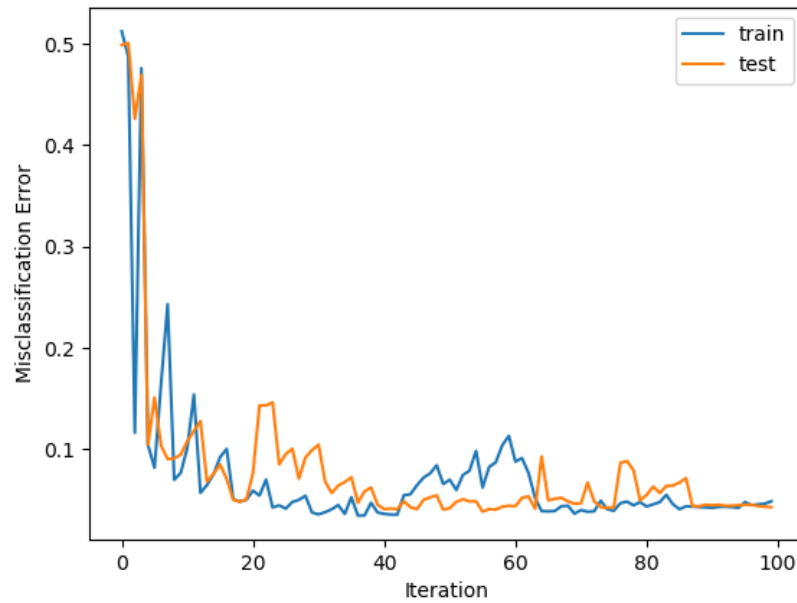


Figure 9: Misclassification Error code is above.

d. Part d plots for stochastic gradient descent with batch size=100.

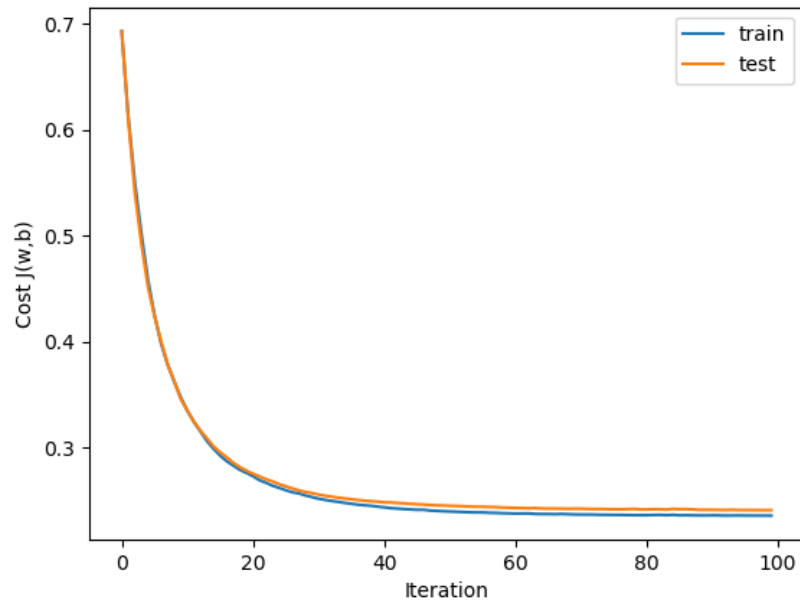


Figure 10: $J(w,b)$ Cost code is above.

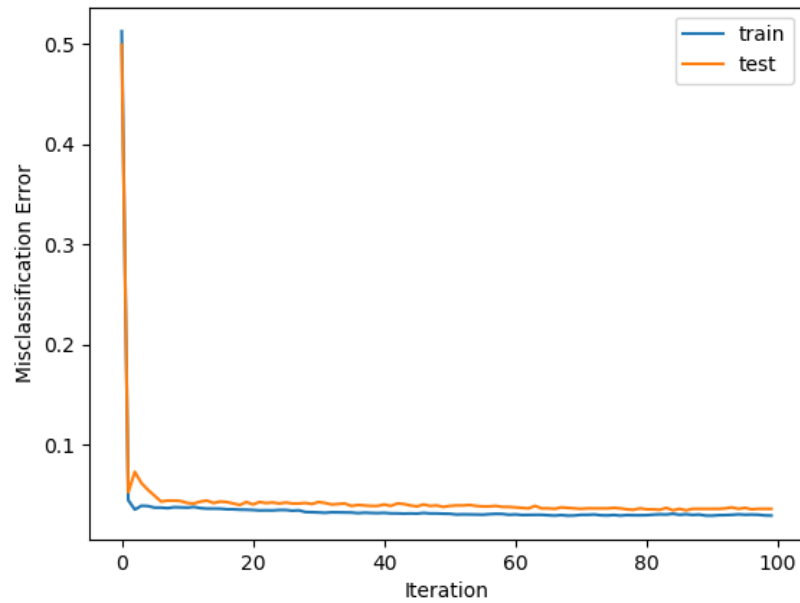


Figure 11: Misclassification Error code is above.