

Homework #2

CSE 446: Machine Learning

Eric Boris: 1976637

Conceptual Questions

Problem 0

- a. **No**, because it's possible that there are other highly correlated features that also have large positive weights. In which case, any one could be removed without affecting the quality of the predictions.
- b. The L1 norm penalty results in a more sparse features than the L2 norm because the **L1 ball is diamond shaped with vertices on the axes** while the **L2 norm is spherical**. This means that an intersection between the function's contour lines with the L1 ball are more likely to occur on an axis, resulting in zeros for the other axes, than with the L2 norm.
- c. Advantage: Better **weight sparsity** than L1. Disadvantage: **Non-Convex**.
- d. k-fold cross validation works by dividing the training data into k subsets. For k iterations, one of the k subsets is used as a test set and the model is trained over the remaining k-1 subsets. The resulting error is found by averaging the error over the k iterations. **Because the training set using k-fold cross is smaller**, k-1 than for example n-1 as in LOOCV, **the bias is higher**. But training time is reduced for k-fold cross from LOOCV. In short, increasing k increases the size of the training set and therefore, increases training time and reduces bias. **k=10** could be a reasonable choice because k is sufficiently large to have **low bias but fast training times**.
- e. **True** because the gradient update values can get stuck stepping over the minimum.
- f. Because **on average it goes in the direction of the true gradient**.
- g. Advantage: SGD is **less computationally intensive** than GD. Disadvantage: SGD requires **more training** than GD to achieve the same error.

Convexity and Norms

Problem 1

- a. Prove that $\|x\| = f(x) = \left(\sum_{i=1}^n x_i^2\right)^{\frac{1}{2}}$ is a norm.
 - Non-negativity: $\|x\| \geq 0$ for all $x_i \in \mathbb{R}^n$ since $x_i^2 \geq 0$ and $\|x\| = 0$ if and only if $x = 0$.

- Absolute scalability: $\|ax\| = |a|\|x\|$ for all $a \in \mathbb{R}$ and $x \in \mathbb{R}^n$.

$$\begin{aligned}
\|ax\| &= \left(\sum_{i=1}^n (ax_i)^2 \right)^{\frac{1}{2}} \\
&= \left(\sum_{i=1}^n a^2 x_i^2 \right)^{\frac{1}{2}} \\
&= \left(a^2 \sum_{i=1}^n x_i^2 \right)^{\frac{1}{2}} \\
&= (a^2)^{\frac{1}{2}} \left(\sum_{i=1}^n x_i^2 \right)^{\frac{1}{2}} \\
&= |a| \left(\sum_{i=1}^n x_i^2 \right)^{\frac{1}{2}} && \text{since } (a^2)^{\frac{1}{2}} > 0 \text{ for all } a \in \mathbb{R} \setminus 0 \\
&= |a|\|x\|
\end{aligned}$$

- Triangle inequality: $\|x + y\| \leq \|x\| + \|y\|$ for all $x, y \in \mathbb{R}^n$. Let $\langle x, y \rangle$ be the inner product $x \cdot y = \sum_{i=1}^n x_i y_i$. Prove using Cauchy-Schwarz inequality, i.e. $|\langle x, y \rangle|^2 \leq \langle x, x \rangle \cdot \langle y, y \rangle$.

$$\begin{aligned}
(\|x + y\|)^2 &= \|x\|^2 + 2\langle x, y \rangle + \|y\|^2 \\
&\leq \|x\|^2 + 2(\|x\| \cdot \|y\|) + \|y\|^2 \\
&= (\|x\| + \|y\|)^2
\end{aligned}$$

Thus $\|x + y\| \leq \|x\| + \|y\|$ since $(\|x + y\|)^2 = (\|x\| + \|y\|)^2$ and all are positive by non-negativity.

- b. Prove by contradiction that $g(x) = \left(\sum_{i=1}^n |x_i|^{\frac{1}{3}} \right)^3$ is not a norm.

Assume that $g(x)$ is a norm.

By the triangle inequality $g(x + y) = g(x) + g(y)$ for all $x, y \in \mathbb{R}^n$.

Let $x = [1, 0]^T$ and $y = [0, 1]^T$.

Trivially, $x, y \in \mathbb{R}^n$.

$$g(x + y) = (1 + 1)^3 = 8$$

$$g(x) + g(y) = 1^3 + 1^3 = 2$$

$$g(x + y) \neq g(x) + g(y)$$

Thus, since $g(x)$ is a norm and the triangle inequality doesn't hold is a contradiction, $g(x)$ is not a norm.

Problem 2

- **I is convex.**
- **II is not convex** because points b and c are in the set but the line segment \overline{bc} is not.
- **III is not convex** because points a and d are in the set but the line segment \overline{ad} is not.

Problem 3

- **I is convex** on $[a, c]$.
- **II is not convex** on $[a, d]$ because $f(\lambda c + (1 - \lambda)d) \geq \lambda f(c) + (1 - \lambda)f(d)$ on $[c, d]$ for an arbitrary λ .

- **II is convex** on $[a, b]$.
- **III is not convex** on $[a, c]$ because $f(\lambda a + (1 - \lambda)b) \geq \lambda f(a) + (1 - \lambda)f(b)$ on $[a, b]$ and $f(\lambda b + (1 - \lambda)c) \geq \lambda f(b) + (1 - \lambda)f(c)$ on $[b, c]$ for an arbitrary λ .

Problem 4

Prove that $f(x) = (x^T A x)^{\frac{1}{2}}$ is convex given that $A \in \mathbb{R}^{n \times n}$ is a symmetric positive definite matrix.

Prove that $f(x)$ is a norm. Since A is a real symmetric matrix, A is orthogonally diagonalizable. By spectral decomposition, there exist matrices $V, \Lambda \in \mathbb{R}^{n \times n}$ such that V is orthogonal and $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_n)$ such that $A = V \Lambda V^T$. Because A is positive definite, for all $i \in [1, \dots, n]$, $\lambda_i > 0$. Then $\Lambda^{\frac{1}{2}} = \text{diag}(\lambda_1^{\frac{1}{2}}, \dots, \lambda_n^{\frac{1}{2}})$ holds. So we can write $f(x) = \|\Lambda V^T x\|_2$. Since ΛV^T is invertible, $\lambda_n \|x\|_2^2 \leq \langle x, A x \rangle \leq \lambda_1 \|x\|_2^2$. Which implies that $\sqrt{\lambda_n} \|x\|_2 \leq \|x\|_A \leq \sqrt{\lambda_1} \|x\|_2$. Hence, $f(x) = \|\Lambda^{\frac{1}{2}} V^T x\|_2$. Thus $f(x)$ is a norm. Prove that all norms are convex. By the definition of convexity, an arbitrary function $g(x)$ is convex if and only if sub-additivity holds. Since $g(x)$ is a norm, the triangle inequality $\|\lambda x + (1 - \lambda)y\| \leq \lambda \|x\| + (1 - \lambda)\|y\|$ holds for all $x, y \in \mathbb{R}^n$ and $\lambda \in [1, 0]$ by problem 1a. By homogeneity then, $g(\lambda x + (1 - \lambda)y) = \lambda g(x) + (1 - \lambda)g(y)$. Therefore, all norms are convex. And since $f(x)$ is a norm, $f(x)$ is convex.

Lasso on a real dataset

Problem 5

- a. See Figure 1.

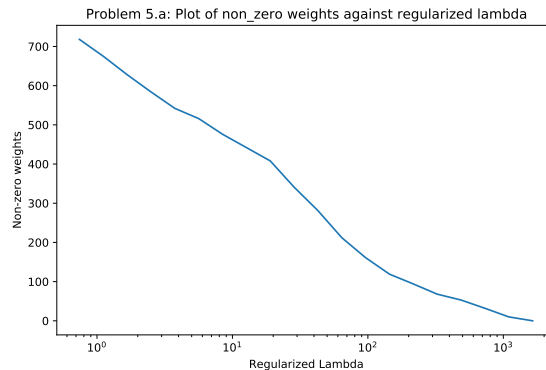


Figure 1

- b. See Figure 2.
- c. We see that higher λ values increase the number of zeros in the weights leading to a more sparse solution. But we encounter problems when λ is too small or too large. When λ is too small, we have a high false discovery rate and when λ is too large the weights become too sparse and the model has poor performance. These data suggest that we should look for values of lambda that balance these trade-offs.

```

1 # Lasso Part 1 - Problem 5
2
3 import sys
4 import numpy as np

```

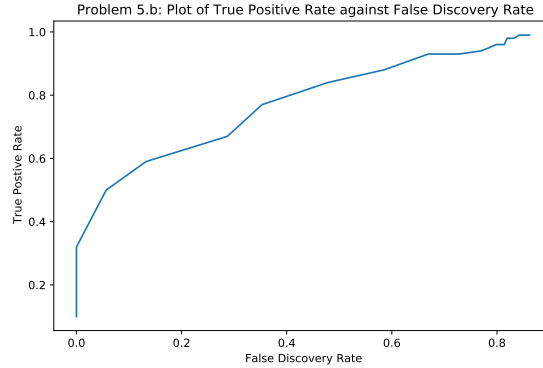


Figure 2

```

5 import matplotlib.pyplot as plt
6 from Lasso import Lasso
7 import Performance as pf
8 import Supplemental as sp
9 import DataManagement as dm
10
11 def part_a(lambdas, non_zeros):
12     plt.figure(figsize=(8, 5))
13     plt.plot(lambdas, non_zeros)
14     plt.title('Problem 5.a: Plot of non-zero weights against regularized lambda')
15     plt.xscale('log')
16     plt.ylabel('Non-zero weights')
17     plt.xlabel('Regularized Lambda')
18     plt.savefig('P5_a.pdf')
19     plt.show()
20
21 def part_b(FDR, TPR):
22     plt.figure(figsize=(8, 5))
23     plt.plot(FDR, TPR)
24     plt.title('Problem 5.b: Plot of True Positive Rate against False Discovery Rate')
25     plt.ylabel('True Postive Rate')
26     plt.xlabel('False Discovery Rate')
27     plt.savefig('P5_b.pdf')
28     plt.show()
29
30 def main(args):
31     # Training data values.
32     if len(args) == 5:
33         n = args[1]
34         d = args[2]
35         k = args[3]
36         sigma = args[4]
37     else:
38         n = 500
39         d = 1000
40         k = 100
41         sigma = 1
42
43     # Get the synthetic data.
44     X, y, w_actual = dm.synthetic_data(n, d, k, sigma)
45
46     # False Discovery Rate (FDR) and True Positive Rate (TPR).
47     FDR = []
48     TPR = []
49
50     # Count of nonzero weights.
51     non_zeros = []
52
53     # Let begin as the max lambda value

```

```

54 # and be reduced by constant ratio 1.5.
55 max_lam = sp.max_lambda(X, y)
56
57 # Let lambdas hold the precomputed regularized lambdas.
58 regularized_lambdas = sp.regularized_lambdas(max_lam)
59
60 # Passing in the values of w_pred into the model is faster
61 # than initializing with 0 weights each time.
62 w_pred = None
63
64 for rl in regularized_lambdas:
65     print(f'Lambda: {rl}')
66
67     # Train the model.
68     model = Lasso(rl)
69     model.train(X, y, w_pred, delta=1E-3, verbose=True)
70     w_pred = model.w
71
72     # Count non-zeros for part a.
73     non_zeros.append(np.sum(abs(w_pred) > 1E-14))
74
75     FDR.append(pf.fdr(w_actual, w_pred))
76     TPR.append(pf.tpr(w_actual, w_pred))
77
78 # Plot the graphs for part a and part b.
79 part_a(regularized_lambdas, non_zeros)
80 part_b(FDR, TPR)
81
82 if __name__ == '__main__':
83     main(sys.argv)

```

Problem 6

- a.
 - Percentage of population that is of hispanic heritage (racePctHispanic): Hispanic is a race term with a meaning that has historically been changed as a result of policy decisions.
 - Number of people living in areas classified as urban (numbUrban): Urban is a term with a meaning that has changed as a result of policy decisions.
 - Percentage of people under the poverty level (PctPopUnderPov): Poverty level has changed as a result of policy decisions.
- b.
 - Percentage of males who are divorced (MalePctDivorce): It might be assumed that higher rates of divorce cause higher crime rates but it could also be the case that higher crime rates cause in an area cause higher divorce rates.
 - Number of vacant households (HousVacant): It might be assumed that more vacant households cause higher crime rates but it could also be the case that higher crime rates cause there to be more vacant houses.
 - Number of homeless people counted in the street (NumStreet): It might be assumed that more homeless people in the street cause higher crime rates but it could also be the case that higher crime rates cause there to be more homeless people in the street.
- a. See Figure 3.
- b. See Figure 4.
- c. See Figure 5.
- d. See Figure 6.
 - Minimum weight: index 39, PctKids2Par with value -0.1247.
 - Maximum weight: index 45, PctIlleg with value 0.0687.

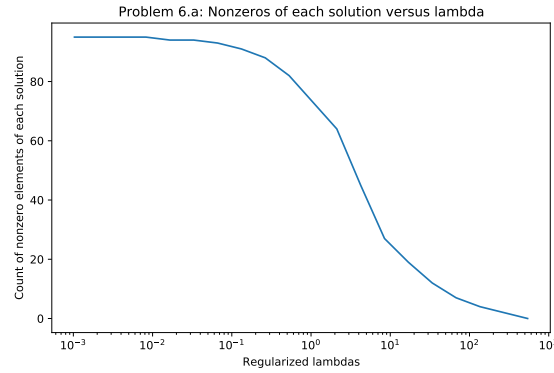


Figure 3

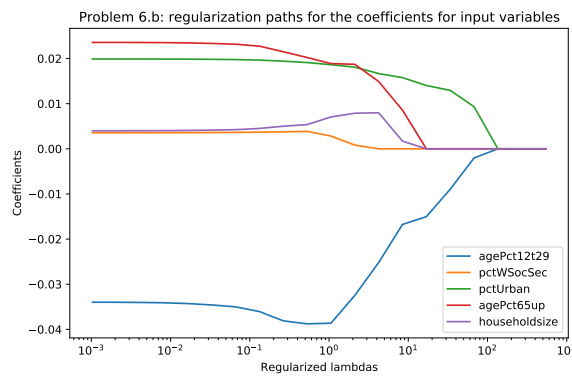


Figure 4

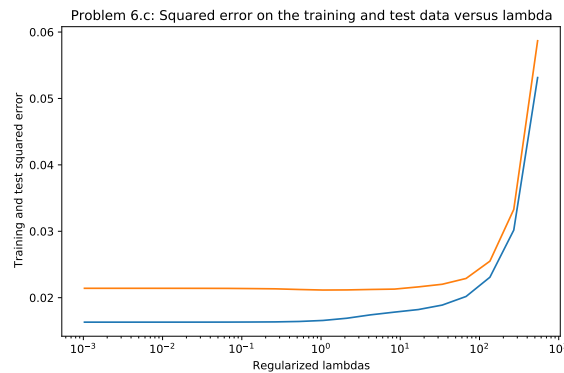


Figure 5

Crime rates correlate most positively with “percentage of kids born to never married” (PctIlleg) and correlate most negatively with “percentage of kids in family housing with 2 parents” (PctKids2Par). These data suggest that crime rates increase with increases in percentage of kids born to never married and decrease with increases in percentage of kids in family housing with 2 parents.

- e. The dictum that “Correlation does not equal Causation” holds. More over-65 years old people living in an area doesn’t cause lower crime rates. There are many confounding factors. One could be that older people move away from high crime rate areas. So rather than their presence indicating low crime rates

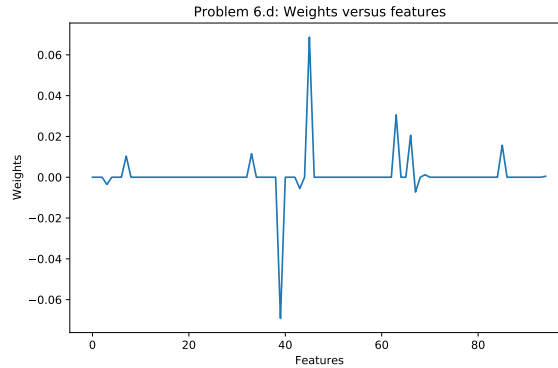


Figure 6

their absence can indicate high crime rates. And, should older people return, the crime rates could remain the same, or even increase, if we assume that older people are more likely to be targeted by criminals than younger people.

```

1 # Lasso Part 2 – Problem 6
2
3 import pandas as pd
4 import numpy as np
5 import matplotlib.pyplot as plt
6 from Lasso import Lasso
7 import DataManagement as dm
8 import Supplemental as sp
9 import Performance as pf
10
11 def part_a(regularized_lambdas, non_zeros):
12     ''' Plot the graph defined in part a. '''
13     plt.figure(figsize=(8, 5))
14     plt.plot(regularized_lambdas, non_zeros)
15     plt.title('Problem 6.a: Nonzeros of each solution versus lambda')
16     plt.xlabel('Regularized lambdas')
17     plt.ylabel('Count of nonzero elements of each solution')
18     plt.xscale('log')
19     plt.savefig('P6.a.pdf')
20     plt.show()
21
22 def part_b(regularization_paths, regularized_lambdas, coefficient_names, coefficient_indices):
23     ''' Plot the graph defined in part b. '''
24     plt.figure(figsize=(8, 5))
25     for coeff, label in zip(np.array(regularization_paths)[: , coefficient_indices].T,
26                             coefficient_names):
27         plt.plot(regularized_lambdas, coeff, label=label)
28     plt.title('Problem 6.b: regularization paths for the coefficients for input variables')
29     plt.xlabel('Regularized lambdas')
30     plt.ylabel('Coefficients')
31     plt.xscale('log')
32     plt.legend()
33     plt.savefig('P6.b.pdf')
34     plt.show()
35
36 def part_c(regularized_lambdas, train_mse, test_mse):
37     ''' Plot the graph defined in part c. '''
38     plt.figure(figsize=(8, 5))
39     plt.plot(regularized_lambdas, train_mse, label='train-mse')
40     plt.plot(regularized_lambdas, test_mse, label='test-mse')
41     plt.title('Problem 6.c: Squared error on the training and test data versus lambda')
42     plt.xlabel('Regularized lambdas')
43     plt.ylabel('Training and test squared error')
44     plt.xscale('log')

```

```

44 plt.savefig('P6_c.pdf')
45 plt.show()
46
47 def part_d(weights):
48     ''' Plot the graph defined in part d. '''
49     plt.figure(figsize=(8, 5))
50     plt.plot(weights)
51     plt.title('Problem 6.d: Weights versus features')
52     plt.xlabel('Features')
53     plt.ylabel('Weights')
54     plt.savefig('P6_d.pdf')
55     plt.show()
56
57 def main():
58     # Load the data frames.
59     df_train = pd.read_table('data/crime-train.txt')
60     df_test = pd.read_table('data/crime-test.txt')
61
62     # Split the data and labels.
63     # Let column be the column to split the data frames on.
64     column = 'ViolentCrimesPerPop'
65     X_train, y_train = dm.split_on_column(df_train, column)
66     X_test, y_test = dm.split_on_column(df_test, column)
67
68     # Let the following be the list of regularized lambdas to train over.
69     max_lambda = sp.max_lambda(X_train, y_train)
70     regularized_lambdas = sp.regularized_lambdas(max_lambda, n=20, constant=2)
71
72     # Let the following hold data relevant to each training iteration for plotting graphs.
73     non_zeros = []
74     train_mse = []
75     test_mse = []
76     paths = []
77
78     # Passing in the values of w_pred into the model is faster
79     # than initializing with 0 weights each time.
80     w_pred = None
81
82     for rl in regularized_lambdas:
83         print(f'Lambda: {rl}')
84
85         # Train the model.
86         model = Lasso(rl)
87         model.train(X_train.values, y_train.values, w_pred, delta=1E-4, verbose=True)
88
89         # Use as the initialization weights on the next iteration.
90         # Must copy to prevent side effects.
91         w_pred = np.copy(model.w)
92         paths.append(w_pred)
93
94         # Run the model.
95         y_hat_train = model.predict(X_train)
96         y_hat_test = model.predict(X_test)
97
98         # Record model performance.
99         non_zeros.append(pf.non_zeros(w_pred))
100         train_mse.append(pf.mse(y_train, y_hat_train))
101         test_mse.append(pf.mse(y_test, y_hat_test))
102
103     # Plot part a.
104     part_a(regularized_lambdas, non_zeros)
105
106     # Plot part b.
107     coefficient_names = ['agePct12t29', 'pctWSocSec', 'pctUrban', 'agePct65up', 'householdsize']
108     coefficient_indices = [X_train.columns.get_loc(name) for name in coefficient_names]
109     part_b(paths, regularized_lambdas, coefficient_names, coefficient_indices)
110
111     # Plot part c.
112     part_c(regularized_lambdas, train_mse, test_mse)

```



```

113
114 # Plot part d.
115 model_d = Lasso(30)
116 model_d.train(X_train.values, y_train.values, w_pred, delta=1E-4, verbose=True)
117 part_d(model_d.w)
118
119 min_w_val = float('inf')
120 min_w_idx = None
121 max_w_val = float('-inf')
122 max_w_idx = None
123 for i, w in enumerate(model_d.w):
124     if w > max_w_val:
125         max_w_val = w
126         max_w_idx = i
127     if w < min_w_val:
128         min_w_val = w
129         min_w_idx = i
130 print(f'Min Weight: index={min_w_idx} val={min_w_val}')
131 print(f'Max Weight: index={max_w_idx} val={max_w_val}')
132
133 for i, weight in enumerate(model_d.w):
134     if weight != 0:
135         print(f'{i}: {weight}')
136
137 for i, label in enumerate(X_train):
138     print(f'{i}: {label}')
139
140 if __name__ == '__main__':
141     main()

```

Logistic Regression

Problem 7

a. Derive the gradient $\nabla_w J(w, b)$.

$$\begin{aligned}
 J(w, b) &= \frac{1}{n} \sum_{i=1}^n \log(1 + \exp(-y_i(b + x_i^T w))) + \lambda \|w\|_2^2 \\
 &= \frac{1}{n} \sum_{i=1}^n \log\left(1 + \left(\frac{1}{\mu_i(w, b)} - 1\right)\right) + \lambda \|w\|_2^2 \\
 &= \frac{1}{n} \sum_{i=1}^n \log\left(\frac{1}{\mu_i(w, b)}\right) + \lambda \|w\|_2^2 \\
 \nabla_w J(w, b) &= \nabla_w \frac{1}{n} \sum_{i=1}^n \log\left(\frac{1}{\mu_i(w, b)}\right) + \nabla_w \lambda \|w\|_2^2 \\
 &= \frac{1}{n} \sum_{i=1}^n \mu_i(w, b) \left(\frac{1}{\mu_i(w, b)} - 1\right) (-y_i)(x_i) + 2\lambda w \\
 &= \frac{1}{n} \sum_{i=1}^n (\mu_i(w, b) - 1) (y_i) x_i + 2\lambda w
 \end{aligned}$$

Thus,

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

Derive the gradient $\nabla_b J(w, b)$.

$$\begin{aligned}
 J(w, b) &= \frac{1}{n} \sum_{i=1}^n \log \left(\frac{1}{\mu_i(w, b)} \right) + \lambda \|w\|_2^2 \\
 \nabla_b J(w, b) &= \nabla_b \frac{1}{n} \sum_{i=1}^n \log \left(\frac{1}{\mu_i(w, b)} \right) + \nabla_b \lambda \|w\|_2^2 \\
 &= \frac{1}{n} \sum_{i=1}^n \mu_i(w, b) \left(\frac{1}{\mu_i(w, b)} - 1 \right) (-y_i) + 0 \\
 &= \frac{1}{n} \sum_{i=1}^n (\mu_i(w, b) - 1) y_i
 \end{aligned}$$

Thus,

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

- b. • See Figure 7.

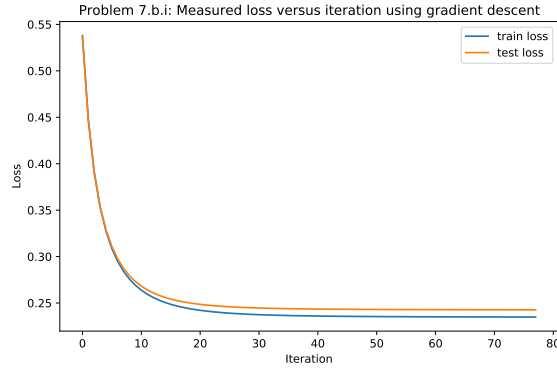


Figure 7

- See Figure 8.

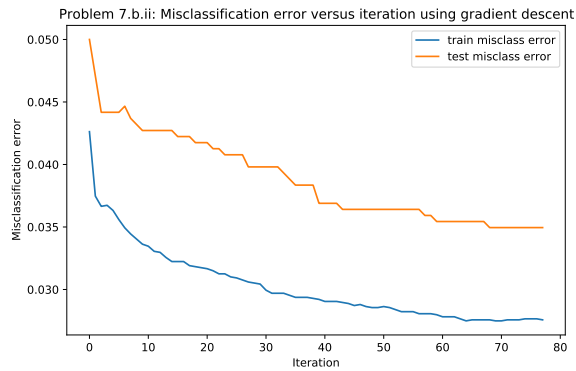


Figure 8

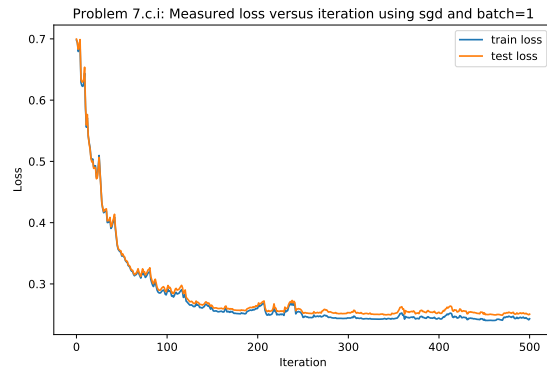


Figure 9

- c.
- See Figure 9.
 - See Figure 10.

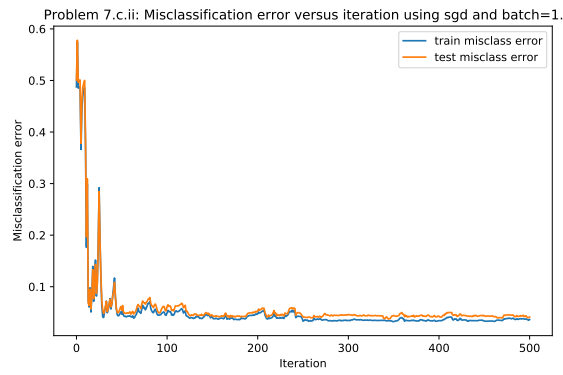


Figure 10

- d.
- See Figure 11.

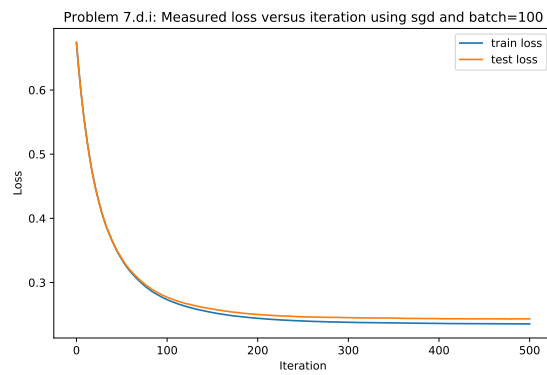


Figure 11

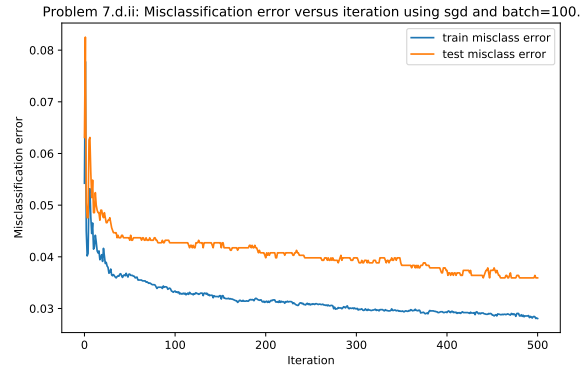


Figure 12

- See Figure 12.

```

1  # Logistic Regression – Problem 7
2
3  import matplotlib.pyplot as plt
4  import DataManagement as dm
5  import Supplemental as sp
6  import Performance as pf
7  from GradientDescent import gradient_descent
8  from StochasticGradientDescent import stochastic_gradient_descent as sgd
9
10 def two-part-plot(x1, x1_label, x2, x2_label, title, x_label, y_label, file_name):
11     plt.figure(figsize=(8, 5))
12     plt.plot(x1, label=x1_label)
13     plt.plot(x2, label=x2_label)
14     plt.title(title)
15     plt.xlabel(x_label)
16     plt.ylabel(y_label)
17     plt.legend()
18     plt.savefig(file_name)
19     plt.show()
20
21 def main():
22     # Load the mnist data with all columns.
23     print(f'Loading data')
24     X_train_all, y_train_all, X_test_all, y_test_all = dm.load_mnist()
25
26     # Strip from each data set all columns except columns 2 and 7.
27     print(f'Stripping columns')
28     keep_cols = (2, 7)
29
30     i_train = dm.indices(y_train_all, keep_cols)
31     X_train = dm.strip_cols(X_train_all, i_train)
32     y_train = dm.strip_cols(y_train_all, i_train)
33
34     i_test = dm.indices(y_test_all, keep_cols)
35     X_test = dm.strip_cols(X_test_all, i_test)
36     y_test = dm.strip_cols(y_test_all, i_test)
37
38     # Set the y label values.
39     print('Setting label values')
40     pos, neg = 1, -1
41     y_train[y_train == 7] = pos
42     y_train[y_train == 2] = neg
43     y_test[y_test == 7] = pos
44     y_test[y_test == 2] = neg
45
46     # Run gradient descent.
47     print('Running gradient descent')

```

```

48 w, b, w_history, b_history, train_loss = gradient_descent(X_train, y_train, alpha=0.1,
49 verbose=True)
50 print('Plotting part b')
51 # Plot measured loss versus iteration.
52 test_loss = [sp.gradient_loss(X_test, y_test, w, b) for w, b in zip(w_history, b_history)]
53
54 #part b.i(train_loss, test_loss)
55 two_part_plot(x1 = train_loss,
56 x1_label = 'train loss',
57 x2 = test_loss,
58 x2_label = 'test loss',
59 title = 'Problem 7.b.i: Measured loss versus iteration using gradient descent',
60 x_label = 'Iteration',
61 y_label = 'Loss',
62 file_name = 'P7_b.i.pdf')
63
64 # Plot misclassification error versus iteration.
65 y_hat_train = [sp.gradient_predict(X_train, w, b) for w, b in zip(w_history, b_history)]
66 y_train_error = [pf.misclass_error(y_train, y_hat) for y_hat in y_hat_train]
67
68 y_hat_test = [sp.gradient_predict(X_test, w, b) for w, b in zip(w_history, b_history)]
69 y_test_error = [pf.misclass_error(y_test, y_hat) for y_hat in y_hat_test]
70
71 two_part_plot(x1 = y_train_error,
72 x1_label = 'train misclass error',
73 x2 = y_test_error,
74 x2_label = 'test misclass error',
75 title = 'Problem 7.b.ii: Misclassification error versus iteration using gradient descent',
76 x_label = 'Iteration',
77 y_label = 'Misclassification error',
78 file_name = 'P7_b.ii.pdf')
79
80 # Run stochastic gradient descent with batch = 1.
81 print('Running Stochastic Gradient Descent with batch=1')
82 w, b, w_history, b_history, train_loss = sgd(X_train, y_train, alpha=0.01, batch_size=1,
83 max_iterations=500, verbose=True)
84
85 print('Plotting part c')
86 # Plot measured loss versus iteration using sgd and batch = 1.
87 test_loss = [sp.gradient_loss(X_test, y_test, w, b) for w, b in zip(w_history, b_history)]
88
89 two_part_plot(x1 = train_loss,
90 x1_label = 'train loss',
91 x2 = test_loss,
92 x2_label = 'test loss',
93 title = 'Problem 7.c.i: Measured loss versus iteration using sgd and batch=1',
94 x_label = 'Iteration',
95 y_label = 'Loss',
96 file_name = 'P7_c.i.pdf')
97
98 # Plot misclassification error versus iteration using sgd and batch = 1.
99 y_hat_train = [sp.gradient_predict(X_train, w, b) for w, b in zip(w_history, b_history)]
100 y_train_error = [pf.misclass_error(y_train, y_hat) for y_hat in y_hat_train]
101
102 y_hat_test = [sp.gradient_predict(X_test, w, b) for w, b in zip(w_history, b_history)]
103 y_test_error = [pf.misclass_error(y_test, y_hat) for y_hat in y_hat_test]
104
105 two_part_plot(x1 = y_train_error,
106 x1_label = 'train misclass error',
107 x2 = y_test_error,
108 x2_label = 'test misclass error',
109 title = 'Problem 7.c.ii: Misclassification error versus iteration using sgd and batch=1.',
110 x_label = 'Iteration',
111 y_label = 'Misclassification error',
112 file_name = 'P7_c.ii.pdf')
113
114 # Run stochastic gradient descent with batch = 1.
115 print('Running Stochastic Gradient Descent with batch=100')

```

```

115 w, b, w_history, b_history, train_loss = sgd(X_train, y_train, alpha=0.01, batch_size=100,
116         max_iterations=500, verbose=True)
117 print('Plotting part c')
118 # Plot measured loss versus iteration using sgd and batch = 100.
119 test_loss = [sp.gradient_loss(X_test, y_test, w, b) for w, b in zip(w_history, b_history)]
120
121 two_part_plot(x1 = train_loss,
122     x1_label = 'train loss',
123     x2 = test_loss,
124     x2_label = 'test loss',
125     title = 'Problem 7.d.i: Measured loss versus iteration using sgd and batch=100',
126     x_label = 'Iteration',
127     y_label = 'Loss',
128     file_name = 'P7-d.i.pdf')
129
130 # Plot misclassification error versus iteration using sgd and batch = 100.
131 y_hat_train = [sp.gradient_predict(X_train, w, b) for w, b in zip(w_history, b_history)]
132 y_train_error = [pf.misclass_error(y_train, y_hat) for y_hat in y_hat_train]
133
134 y_hat_test = [sp.gradient_predict(X_test, w, b) for w, b in zip(w_history, b_history)]
135 y_test_error = [pf.misclass_error(y_test, y_hat) for y_hat in y_hat_test]
136
137 two_part_plot(x1 = y_train_error,
138     x1_label = 'train misclass error',
139     x2 = y_test_error,
140     x2_label = 'test misclass error',
141     title = 'Problem 7.d.ii: Misclassification error versus iteration using sgd and batch=100.',
142     x_label = 'Iteration',
143     y_label = 'Misclassification error',
144     file_name = 'P7-d.ii.pdf')
145
146 if __name__ == '__main__':
147     main()

```

Additional files used for Problems 5-7

```

1 # Define Data Management functions for Problems 5-7.
2
3 import numpy as np
4 from mnist import MNIST
5
6 def synthetic_data(n, d, k, sigma):
7     ''' Generate synthetic training data. '''
8     # Let X be an (n, d) matrix with values drawn from a normal distribution.
9     X = np.random.normal(0.0, 1.0, (n, d))
10
11     # Let w be a (d, ) array with values d/k if d in range {1, ..., k} else 0.
12     w = np.arange(d) / k
13     w[k + 1:] = 0
14
15     # Let y be a (n, ) array with values y_i = w^T x_i + epsilon_i
16     # where epsilon is a (n, ) array with values drawn from a normal distribution.
17     epsilon = np.random.normal(0.0, sigma, (n, ))
18     y = X.dot(w) + epsilon
19
20     return X, y, w
21
22 def split_on_column(df, column):
23     ''' Return split the given data frame df on column s.t.
24     X consists of all data except column and
25     y consists of only the data in column. '''
26     return df.drop(column, axis=1), df[column]
27
28 def load_mnist():
29     ''' Load and return the mnist training and testing datasets. '''
30     data = MNIST('data/python-mnist/data/')
31
32     # Load the data.

```

```

33 X_train, y_train = map(np.array, data.load_training())
34 X_test, y_test = map(np.array, data.load_testing())
35
36 # Reduce the data.
37 X_train = X_train / 255.0
38 X_test = X_test / 255.0
39
40 return X_train, y_train, X_test, y_test
41
42 def indices(data, cols):
43     indices = 0
44     for c in cols:
45         indices += (data == c).astype('int')
46     return indices
47
48 def strip_cols(data, indices):
49     ''' Return data with all columns not in cols removed. '''
50     return data[indices.astype('bool')].astype('float')

1 # Implementation of Gradient Descent for Problem 7.
2
3 import numpy as np
4 import Supplemental as sp
5 def gradient_descent(X, y, alpha, regularized_lambda=1E-1, w_init=None, b_init=0, delta=1E-4,
6                     max_iterations=1E4, verbose=False):
7     # Initialize the given weights if none given.
8     d = X.shape[1]
9     if w_init is None:
10         w_init = np.zeros(d)
11
12     # Initialize the working weights and bias.
13     w_curr, b_curr = w_init, b_init
14     w_prev = w_init + float('inf')
15
16     # Store the current iteration of the function.
17     i = 0
18
19     # Store the states of each iteration.
20     w_history = []
21     b_history = []
22     loss_history = []
23
24     # While not converged.
25     dw = float('inf')
26     while dw >= delta and i <= max_iterations:
27         # Store the previous iteration for step comparison.
28         # Copy to prevent side effects.
29         w_prev = np.copy(w_curr)
30
31         # Step down the gradient.
32         w_curr = w_curr - alpha * sp.gradient_w(X, y, w_curr, b_curr, regularized_lambda)
33         b_curr = b_curr - alpha * sp.gradient_b(X, y, w_curr, b_curr)
34
35         # Compute the loss.
36         loss = sp.gradient_loss(X, y, w_curr, b_curr, regularized_lambda)
37
38         # Store results of current iterations.
39         w_history.append(w_curr)
40         b_history.append(b_curr)
41         loss_history.append(loss)
42
43         # Output results of this iteration.
44         if verbose and i % 10 == 0:
45             print(f'{i}\tLoss: {loss}')
46
47         # Update for next iteration.
48         dw = np.linalg.norm(w_curr - w_prev, np.inf)
49         i += 1
50     return w_curr, b_curr, w_history, b_history, loss_history

```

```

1  # Implement Lasso for Problems 5–6.
2
3  import numpy as np
4
5  class Lasso:
6      def __init__(self, regularized_lambda):
7          self.rl = regularized_lambda
8
9      def train(self, X, y, w, delta=1E-3, verbose=False):
10         ''' Train the lasso using coordinate descent. '''
11         n, d = X.shape
12         self.w = np.zeros(d) if w is None else w
13
14         # Let history hold the history of loss measures.
15         self.history = []
16
17         # Precompute a to speed things up considerably.
18         a = 2 * np.sum(X ** 2, axis=0)
19
20         # Let i be the current iteration.
21         i = 0
22
23         # while not converged
24         w_change = float('inf')
25         while w_change >= delta:
26             # For measuring iteration differences.
27             w_prev = np.copy(self.w)
28
29             #  $b \leftarrow \frac{1}{n} \sum_{i=1}^n (y_i - \sum_{j=1}^d w_j * x_{i,j})$ 
30             self.b = np.mean(y - X.dot(self.w))
31
32             # for k in {1, 2, ..., d} do
33             for k in range(d):
34                 # Access precomputed a.
35                 a_k = a[k]
36
37                 # Used in computing c_k
38                 not_k = np.arange(d) != k
39
40                 #  $c_k \leftarrow 2 * \sum_{i=1}^n x_{i,k} * (y_i - (b + \sum_{j \neq k} w_j * x_{i,j}))$ 
41                 c_k = 2 * np.sum(X[:, k] * (y - (self.b + X[:, not_k].dot(self.w[not_k]))),
42                                axis=0)
43
44                 # w_k is a piecewise assignment.
45                 #  $w_k \leftarrow (c_k + \lambda) / a_k$  if  $c_k < -\lambda$ 
46                 #  $w_k \leftarrow (c_k - \lambda) / a_k$  if  $c_k > \lambda$ 
47                 #  $w_k \leftarrow 0$  if  $c_k \in [-\lambda, \lambda]$ 
48                 condlist = [c_k < -self.rl, c_k > self.rl, ]
49                 funclist = [(c_k + self.rl) / a_k, (c_k - self.rl) / a_k, 0]
50                 self.w[k] = np.float(np.piecewise(c_k, condlist, funclist))
51
52             # Compute the loss.
53             loss = self.loss(X, y)
54             self.history.append(loss)
55
56             # Output progress.
57             if verbose:
58                 print(f'\t{i}\tLoss: {loss}')
59
60             # Update the change in w with the new values.
61             w_change = np.linalg.norm(self.w - w_prev, ord=np.inf)
62
63             # Increment the current iteration.
64             i += 1
65
66         return self.history
67
68 def loss(self, X, y):

```



```

68         ''' Compute the lasso loss. '''
69         return (np.linalg.norm(X.dot(self.w) + self.b - y)) ** 2 + self.rl * np.linalg.norm(
            self.w, ord=1)
70
71     def predict(self, X):
72         ''' Predict y_hat using the trained model. '''
73         return X.dot(self.w) + self.b

```



```

1  # Define Performance functions for Problems 5-7.
2
3  import numpy as np
4
5  # Let the following be the threshold, below which a value
6  # is considered to be zero.
7  ZERO_THRESHOLD = 1E-14
8
9  def true_positive(actual, predicted, threshold=ZERO_THRESHOLD):
10     ''' Return the count of true positives. '''
11     return np.sum(np.logical_and(abs(actual)>threshold, abs(predicted)>threshold))
12
13  def true_negative(actual, predicted, threshold=ZERO_THRESHOLD):
14     ''' Return the count of true negative. '''
15     return np.sum(np.logical_and(abs(actual)<=threshold, abs(predicted)<=threshold))
16
17  def false_positive(actual, predicted, threshold=ZERO_THRESHOLD):
18     ''' Return the count of false positives. '''
19     return np.sum(np.logical_and(abs(actual)<=threshold, abs(predicted)>threshold))
20
21  def false_negative(actual, predicted, threshold=ZERO_THRESHOLD):
22     ''' Return the count of false negatives. '''
23     return np.sum(np.logical_and(abs(actual)>threshold, abs(predicted)<=threshold))
24
25  def fdr(actual, predicted):
26     ''' Return the False Discovery Rate of the given data. '''
27     tp = true_positive(actual, predicted)
28     fp = false_positive(actual, predicted)
29     return (fp / (tp + fp))
30
31  def tpr(actual, predicted):
32     ''' Return the True Positive Rate of the given data. '''
33     tp = true_positive(actual, predicted)
34     fn = false_negative(actual, predicted)
35     return (tp / (tp + fn))
36
37  def non_zeros(array, threshold=ZERO_THRESHOLD):
38     ''' Return the count of non-zero values in the given np array. '''
39     return np.sum(abs(array) > threshold)
40
41  def mse(actual, predicted):
42     ''' Return the Mean Squared Error of the given data. '''
43     return np.mean((actual - predicted) ** 2)
44
45  def misclass_error(actual, predicted):
46     ''' Return the misclassification error. '''
47     return 1 - np.mean(actual == predicted)

```



```

1  # Implementation of Stochastic Gradient Descent for Problem 7.
2
3  import numpy as np
4  import Supplemental as sp
5
6  def stochastic_gradient_descent(X, y, alpha, batch_size, regularized_lambda=1E-1, w_init=None,
    b_init=0, delta=1E-4, max_iterations=1E4, verbose=False):
7      # Initialize the given weights if none given.
8      n, d = X.shape
9      if w_init is None:
10         w_init = np.zeros(d)
11
12     # Initialize the working weights and bias.

```

```

13 w_curr, b_curr = w_init, b_init
14 w_prev = w_init + float('inf')
15
16 # Store the current iteration of the function.
17 i = 0
18
19 # Store the states of each iteration.
20 w_history = []
21 b_history = []
22 loss_history = []
23
24 # While not converged.
25 dw = float('inf')
26 while dw >= delta and i <= max_iterations:
27     # Batch the descent.
28     batch_index = np.random.choice(n, batch_size)
29     X_batch = X[batch_index]
30     y_batch = y[batch_index]
31
32     # Store the previous iteration for alpha comparison.
33     # Copy to prevent side effects.
34     w_prev = np.copy(w_curr)
35
36     # Step down the gradient.
37     w_curr = w_curr - alpha * sp.gradient_w(X_batch, y_batch, w_curr, b_curr, regularized_lambda)
38     b_curr = b_curr - alpha * sp.gradient_b(X_batch, y_batch, w_curr, b_curr)
39
40     # Compute the loss.
41     loss = sp.gradient_loss(X, y, w_curr, b_curr, regularized_lambda)
42
43     # Store results of current iterations.
44     w_history.append(w_curr)
45     b_history.append(b_curr)
46     loss_history.append(loss)
47
48     # Output results of this iteration.
49     if verbose:
50         print(f'{i}\tLoss: {loss}')
51
52     # Update for next iteration.
53     dw = np.linalg.norm(w_curr - w_prev, np.inf)
54     i += 1
55
56 return w_curr, b_curr, w_history, b_history, loss_history

1 # Define Supplemental functions for Problems 5–7.
2
3 import numpy as np
4
5 def max_lambda(X, y):
6     ''' Return the smallest value of lambda for which the solution w_hat is entirely zero. '''
7     return np.max(np.sum(2 * X * (y - np.mean(y))[:, None], axis=0))
8
9 def regularized_lambdas(max_lambda, n=20, constant=1.5):
10     ''' Return a list of n regularized lambdas starting from max_lambda and
11         decreasing by a constant ratio. '''
12     return [(max_lambda / (constant ** i)) for i in range(n)]
13
14 def _mu(X, y, w, b):
15     ''' Compute the inverse exponential. '''
16     return 1 / (1 + np.exp(-y * (b + X.dot(w))))
17
18 def gradient_w(X, y, w, b, regularized_lambda=1E-1):
19     ''' Return the gradient J loss function with respect to w. '''
20     return np.mean((( _mu(X, y, w, b) - 1) * y)[:, None] * X, axis=0) + (2 * regularized_lambda *
21 w)
22
23 def gradient_b(X, y, w, b):
24     ''' Return the gradient J loss function with respect to b. '''

```

```

24     return np.mean(((mu(X, y, w, b) - 1) * y), axis=0)
25
26 def gradient_loss(X, y, w, b, regularized_lambda=1E-1):
27     ''' Return the loss function J. '''
28     return np.mean(np.log(1 + np.exp(-y * (b + X.dot(w))))) + (regularized_lambda * w.dot(w))
29
30 def gradient_predict(X, w, b):
31     ''' Perform a prediction using the gradient weights. '''
32     return np.sign(X.dot(w) + b)

```