

APPLIED MACHINE LEARNING

HOMEWORK 6

Daniel Gonzalez
Colton Piper
25th of October, 2018

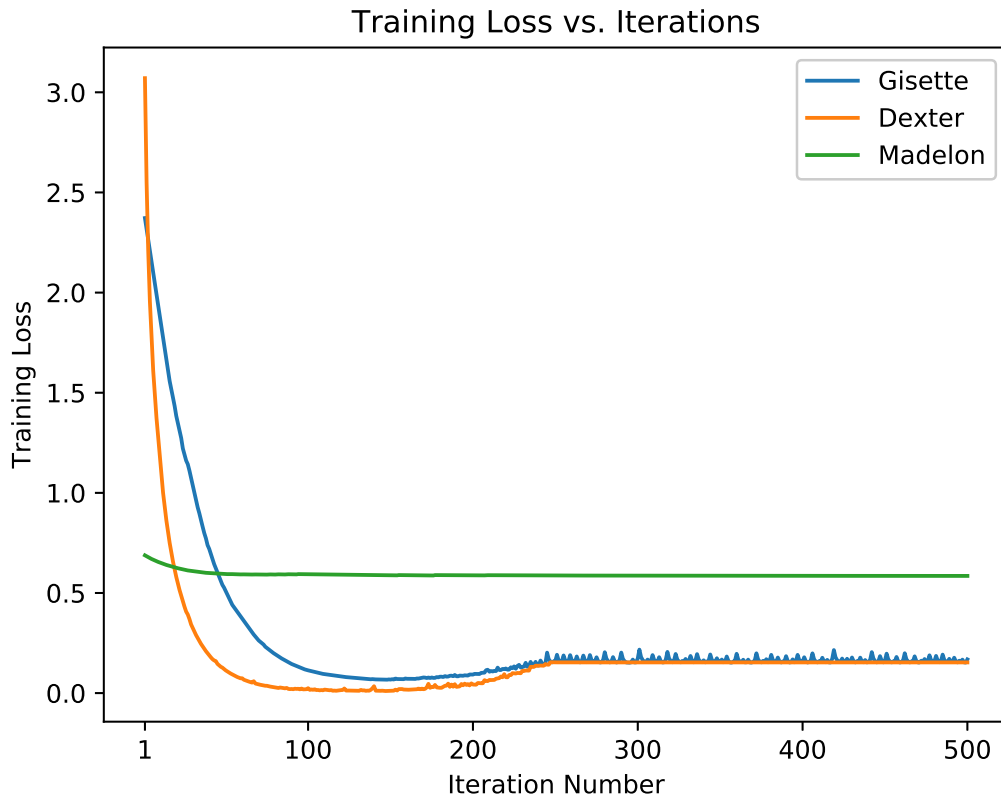
1 Results

1.1 Table

Table 1: Summary of Results: Original Learning Rates

Data Set	Learning Rate η	Features Selected k	Training Error	Test Error
gisette	5	10	7.65%	8.30%
		30	3.28%	4.20%
		100	2.60%	3.00%
		300	2.87%	3.90%
dexter	10	10	7.00%	12.67%
		30	0.67%	14.33%
		100	0.00%	9.00%
		300	0.00%	11.00%
madelon	0.01	10	38.00%	39.50%
		30	35.80%	40.83%
		100	31.20%	40.83%
		300	25.90%	42.83%

1.2 Figures



2 Appendix: Code

If the code looks too small, please zoom in on the pdf. The screenshots are .png images, so you should be able to zoom in and read at whatever is a comfortable size for you. The first two screenshots are of the actual FSA code, while the last screenshot is the python code for normalizing and pre-processing the data.

```
1 # Daniel Gonzalez, FSU Mathematics PhD
2 # Colton Piper, FSU Mathematics PhD
3 # Applied Machine Learning Assignment 6
4
5 import re
6 import numpy as np
7 import matplotlib.pyplot as plt
8
9 training_data = ["/data_norm/gisette/gisette_train.data.npy", "/data_norm/dexter/dexter_train.csv.npy", "/data_norm/madelon/madelon_train.data.npy"]
10 test_data = ["/data_norm/gisette/gisette_valid.data.npy", "/data_norm/dexter/dexter_valid.csv.npy", "/data_norm/madelon/madelon_valid.data.npy"]
11
12 training_labels = ["/data_norm/gisette/gisette_train.labels.npy", "/data_norm/dexter/dexter_train.labels.npy", "/data_norm/madelon/madelon_train.labels.npy"]
13 test_labels = ["/data_norm/gisette/gisette_valid.labels.npy", "/data_norm/dexter/dexter_valid.labels.npy", "/data_norm/madelon/madelon_valid.labels.npy"]
14
15 #SHRINKAGE FACTOR
16 s = 0.001
17 #ANNEALING PARAMETER
18 M = 100
19 #NUMBER OF ITERATIONS
20 N_itr = 500
21 #LEARNING RATES FOR DIFFERENT DATA SETS
22 η_master = [5, 10, 0.01]
23 #VARIABLE SELECTION PARAMETER
24 K = [10, 30, 100, 300]
25
26 #LOSS FUNCTION
27 def lorenz_loss(X, Y, β):
28     Z = Y * np.dot(X, β)
29     loss = np.log(1 + (Z - 1)**2)
30     loss[Z > 1] = 0
31     return np.sum(loss)/X.shape[0] + s * (np.linalg.norm(β)**2)
32
33 #GRADIENT OF LORENZ LOSS
34 def gradient(X, Y, β):
35     ybx = Y * np.dot(X, β)
36     f = (2*(ybx - 1)/(1 + (ybx - 1)**2))
37     f[ybx > 1] = 0
38     return np.dot(X.T, f * Y)/X.shape[0] + 2*s*β
39
40 #VARIABLE ELIMINATION: RETAIN ONLY THE M1 FEATURES WITH THE LARGEST L2 NORM
41 def feature_selection(β, M1):
42     sorted_features = np.flip(np.sort(β**2, kind="mergesort", axis=None))
43     threshold = sorted_features[M1-1]
44     count = len([f for f in sorted_features[M1:] if f == threshold])
45
46     β[β**2 < threshold] = 0
47     for i in range(0, β.shape[0]):
48         if β[i]**2 == threshold and count > 0:
49             β[i] = 0
50             count -= 1
51         elif count <= 0:
52             break
53     return β
54
55 #FEATURE SELECTION WITH ANNEALING: GRADIENT DESCENT FOLLOWED BY VARIABLE ELIMINATION
56 def fsa(X, Y, validation, η, k):
57     N = X.shape[0]
58     M = X.shape[1]
59     β = np.zeros(M);
60     loss = []
61     for i in range(1, N_itr+1):
62         β = β - η * gradient(X, Y, β)
63
64         #FIND THE NUMBER OF FEATURES TO RETAIN
65         M1 = int(k + (M - k)*max(0, (N_itr - 2*i)/(2*i*η + N_itr)))
66         #M1 = k + (M - k)*max(0, int((N_itr - 2*i)/(2*i*η + N_itr)))
67         β = feature_selection(β, M1)
68         indices = np.where(β == 0)
69
70         #DELETE THE IRRELEVANT FEATURES FROM THE WEIGHTS AND DATA
71         β = np.delete(β, indices)
```

1.

```

54
55 #FEATURE SELECTION WITH ANNEALING: GRADIENT DESCENT FOLLOWED BY VARIABLE ELIMINATION
56 def fsa(X, Y, validation, η, k):
57     N = X.shape[0]
58     M = X.shape[1]
59     β = np.zeros(M);
60     loss = []
61     for i in range(1, N_itr+1):
62         β = β - η * gradient(X, Y, β)
63
64         #FIND THE NUMBER OF FEATURES TO RETAIN
65         M1 = int(k + (M - k)*max(0, (N_itr - 2*i)/(2*i*M + N_itr)))
66         #M1 = k + (M - k)*max(0, int((N_itr - 2*i)/(2*i*M + N_itr)))
67         β = feature_selection(β, M1)
68         indices = np.where(β == 0)
69
70         #DELETE THE IRRELEVANT FEATURES FROM THE WEIGHTS AND DATA
71         β = np.delete(β, indices)
72         X = np.delete(X, indices, 1)
73         validation = np.delete(validation, indices, 1)
74         loss.append(lorezn_loss(X, Y, β))
75     return β, X, validation, loss
76
77 #PREDICTION
78 def predict(X, Y, β):
79     p = 1/(1 + np.exp(-np.dot(X, β)))
80     p[p > 0.5] = 1
81     p[p < 0.5] = -1
82     error = np.count_nonzero(Y - p)/Y.shape[0] * 100
83     return error
84
85 #MAIN PROGRAM
86 f = open("report/results/results.txt", "w")
87 l = open("report/results/losses.txt", "w")
88
89 losses = []
90 for n, f_train_data, f_test_data, f_train_labels, f_test_labels in zip(n_master, training_data, test_data, training_labels, test_labels):
91     print("Processing: " + f_train_data.split("/")[-2], file=f)
92     print("Processing: " + f_train_data.split("/")[-2], file=l)
93     for k in K:
94         print("\tk = " + str(k), file=f)
95         print("\tk = " + str(k), file=l)
96         Xtrain = np.load(f_train_data)
97         Ytrain = np.load(f_train_labels)
98         Xtest = np.load(f_test_data)
99         Ytest = np.load(f_test_labels)
100
101         #TRAINING THE MODEL
102         β, Xtrain, Xtest, loss = fsa(Xtrain, Ytrain, Xtest, η, k)
103         if k == 10:
104             losses.append(loss)
105
106         #PREDICT NEW LABELS AND COMPUTE THE MISCLASSIFICATION ERROR
107         error_train = predict(Xtrain, Ytrain, β)
108         print("\t\tTraining error: " + str(error_train), file=f)
109         error_test = predict(Xtest, Ytest, β)
110         print("\t\tValidation error: " + str(error_test), file=f)
111
112 x_axis = range(1, 501)
113 fig, ax = plt.subplots()
114
115 ax.plot(x_axis, losses[0], label="Gisette")
116 ax.plot(x_axis, losses[1], label="Dexter")
117 ax.plot(x_axis, losses[2], label="Madelon")
118
119 legend = ax.legend(loc='best')
120 plt.title(r'Learning Rate eta = [5, 10, 0.01]')
121 plt.xticks([1, 100, 200, 300, 400, 500], [1, 100, 200, 300, 400, 500])
122 plt.xlabel("Iteration Number")
123 plt.ylabel("Training Loss")
124 plt.savefig("../report/figures/graph.eps", format="eps", dpi=1000, bbox_inches="tight")

```

12

```

1 # Daniel Gonzalez, FSU Mathematics PhD
2 # Colton Piper, FSU Mathematics PhD
3 # Applied Machine Learning Assignment 6
4
5 import re
6 import numpy as np
7
8 training_data = ["/data/gisette/gisette_train.data", "/data/dexter/dexter_train.csv", "/data/madelon/madelon_train.data"]
9 test_data = ["/data/gisette/gisette_valid.data", "/data/dexter/dexter_valid.csv", "/data/madelon/madelon_valid.data"]
10
11 training_labels = ["/data/gisette/gisette_train.labels", "/data/dexter/dexter_train.labels", "/data/madelon/madelon_train.labels"]
12 test_labels = ["/data/gisette/gisette_valid.labels", "/data/dexter/dexter_valid.labels", "/data/madelon/madelon_valid.labels"]
13
14 #MAIN DATA PROCESSING
15 for f_train_data, f_test_data, f_train_labels, f_test_labels in zip(training_data, test_data, training_labels, test_labels):
16     with open(f_train_data) as train, open(f_test_data) as test, open(f_train_labels) as train_labels, open(f_test_labels) as test_labels:
17         print("Processing: " + str(f_train_data.split("/")[-1]))
18         data = []
19         valid = []
20         data_labels = []
21         valid_labels = []
22
23         #INPUT DATA
24         for line in train:
25             data.append([float(x) for x in re.split(r'[ , ]', line.strip().strip("\n"))])
26         for line in test:
27             valid.append([float(x) for x in re.split(r'[ , ]', line.strip().strip("\n"))])
28         for line in train_labels:
29             data_labels.append([float(x) for x in re.split(r'[ , ]', line.strip().strip("\n"))])
30         for line in test_labels:
31             valid_labels.append([float(x) for x in re.split(r'[ , ]', line.strip().strip("\n"))])
32
33         data = np.asarray(data)
34         valid = np.asarray(valid)
35         data_labels = np.asarray(data_labels)
36         valid_labels = np.asarray(valid_labels)
37
38         data_labels = data_labels.reshape(data_labels.shape[0])
39         valid_labels = valid_labels.reshape(valid_labels.shape[0])
40
41         #NORMALIZE
42         avg = np.mean(data, axis=0)
43         std = np.std(data, axis=0)
44
45         data = (data - avg)/std
46         data = data[:, (np.isfinite(data)).any(axis=0)]
47
48         valid = (valid - avg)/std
49         valid = valid[:, (np.isfinite(valid)).any(axis=0)]
50
51         #RELABEL
52         data_labels[data_labels != 1] = -1
53         valid_labels[valid_labels != 1] = -1
54
55         #OUTPUT DATA
56         print("\t Writing normalized training data...")
57         np.save("/data_norm/" + f_train_data.split("/")[-2] + ".n" + f_train_data.split("/")[-1], data)
58
59         print("\t Writing normalized validation data...")
60         np.save("/data_norm/" + f_test_data.split("/")[-2] + ".n" + f_test_data.split("/")[-1], valid)
61
62         print("\t Writing training labels...")
63         np.save("/data_norm/" + f_train_labels.split("/")[-2] + ".n" + f_train_labels.split("/")[-1], data_labels)
64
65         print("\t Writing validation labels...")
66         np.save("/data_norm/" + f_test_labels.split("/")[-2] + ".n" + f_test_labels.split("/")[-1], valid_labels)

```

"normalize.py" 66L, 2952C

1.