

CSE 847 Home Assignment 4

Submitted by: Ritam Guha (MSU ID: guharita)

Date: March 31, 2021

1 Logistic Regression: Experiment

As a part of this experiment, we are supposed to build a classifier based on logistic regression. The experimentation is performed on **Spam Email Detection** dataset. The dataset is publicly available at: [Spam Email Dataset](#). The dataset consists of 4601 samples classified as 1 (Spam) and 0 (non-Spam).

First, the labels were converted to $+1/-1$ scheme from the existing $+1/0$ scheme. Multiple training sizes were used to track the classification accuracies for varying training load. The test data was always fixed and consisted of the last 2601 entries of the data. On the other hand the training size was increased as: $n = 200, 500, 800, 1000, 1500, 2000$.

1.1 Code

The code is uploaded to a public Github repo at: [CSE 847 HA 4](#). For reference, the code is also provided below:

```
1 data = importdata('Data/Spam Email Detection/data.xlsx'); % import data
2 labels = importdata('Data/Spam Email Detection/labels.xlsx'); % import labels
3 labels(labels==0) = -1; % transforming 0/1 labels to -1/+1
4 data(:, size(data,2)+1) = ones(1, size(data,1))';
5
6
7 global cur_train_size;
8
9 % initializing training and test data
10 train_size = [200, 500, 800, 1000, 1500, 2000];
11 test_size = 2601;
12 [num_samples, num_features] = size(data);
13 test_data = data(num_samples - test_size + 1 : num_samples, :);
14 test_labels = labels(num_samples - test_size + 1 : num_samples);
15 accuracy = zeros(1, size(train_size, 2));
16
17 for i = 1: size(train_size, 2)
18     cur_train_size = train_size(i);
19
20     train_data = data(1 : cur_train_size, :);
21     train_labels = labels(1 : cur_train_size);
22
```

```

23     weights = logRegression(train_data , train_labels , 1e-5, 1000);
24     accuracy(1, i) = compute_accuracy(test_data , test_labels , weights);
25     fprintf('Training Size = %d: %f\n', cur_train_size , accuracy(1,i));
26 end
27
28 % Plot the variation of accuracy with training size
29 figure;
30 X = train_size;
31 Y = accuracy;
32 fig = plot(X, Y);
33 xlabel('Training Size');
34 ylabel('Classification Accuracy (in %)');
35 title('Variation of Classification Accuracy with Training Size')
36 saveas(fig , strcat('Results/Logistic Train/Accuracy_Variance.jpg'));
37
38
39 function [accuracy] = compute_accuracy(data, labels , weights)
40     global cur_train_size
41
42     % getting the predicted labels and computing accuracy
43     predicted_labels = sigmoid(data * weights);
44     predicted_labels(predicted_labels > 0.5) = 1;
45     predicted_labels(predicted_labels <= 0.5) = -1;
46     correct_predictions = sum(predicted_labels == labels);
47     accuracy = (correct_predictions/size(data,1) * 100);
48
49     % create the confusion matrix
50     fig = confusionchart(labels , predicted_labels);
51     title(strcat('Confusion Chart for Training Size: ', int2str(cur_train_size))
52           );
53     saveas(fig , strcat('Results/Logistic Train/Conf_Chart_Train-Size-', int2str(
54         cur_train_size), '.jpg'));
55 end
56
57 function [weights] = logRegression(data, labels , epsilon , maxiter)
58
59     % setting the default parameter values
60     if nargin < 4
61         if ~exist('epsilon')
62             epsilon = 1e-6;
63         end
64         if ~exist('maxiter')
65             maxiter=1000;
66         end
67     end
68
69     %
70     % code to train a logistic regression classifier

```

```

68 %
69 % INPUTS:
70 %     data      = n * (d+1) matrix with n samples and d features, where
71 %                column d+1 is all ones (corresponding to the intercept term)
72 %     labels    = n * 1 vector of class labels (taking values 0 or 1)
73 %     epsilon   = optional argument specifying the convergence
74 %                criterion - if the change in the absolute difference in
75 %                predictions, from one iteration to the next, averaged across
76 %                input features, is less than epsilon, then halt
77 %                (if unspecified, use a default value of 1e-5)
78 %     maxiter    = optional argument that specifies the maximum number of
79 %                iterations to execute (default=1000)
80 %     train_size = number of samples used for training
81 %     test_size  = number of samples used for testing
82 %
83 % OUTPUT:
84 %     weights = (d+1) * 1 vector of weights where the weights correspond to
85 %                the columns of "data"
86
87
88 % initializing parameters
89 num_features = size(data, 2);
90 w = zeros(num_features, 1);
91 iter = 1;
92 eta = 0.00001;
93 prev_error = Inf;
94 cur_error = -Inf;
95
96 % loop running till convergence
97 while(iter <= maxiter && (abs(cur_error - prev_error) >= epsilon))
98     % compute train error
99     z = -labels .* (data * w);
100    error(1, iter) = mean(log(1 + exp(z)));
101
102    prev_error = cur_error;
103    cur_error = error(1, iter);
104
105    % use the gradient of the loss function wrt the training data to
106    % update the weight
107    dw = (mean(-exp(-z)./(1+exp(-z)) .* (data .* labels)))';
108    w = w - (eta * dw);
109    iter = iter+1;
110 end
111
112 iter = iter-1;
113 weights = w;
114

```

```

115 % plot the train error over the iterations
116 figure;
117 hold on;
118 x = linspace(1,iter,iter);
119 plot(x, error);
120 legend('Train Error');
121 hold off;
122 end
123
124 function [val] = sigmoid(input)
125 % sigmoid function implementation
126 val = 1./(1 + exp(-input));
127 end

```

1.2 Experimental Outcome

From the experimentation, it was observed that the classification accuracy varied over the increasing number of training samples as shown in [Figure 1](#). From the Figure, it is visible that at first, increasing the number of samples in the training schedule was improving the performance. But when the training size exceeded 1000, the classification accuracy started dropping beyond that point. The best classification accuracy obtained by this classifier was 91.080354 for training size of 1000.

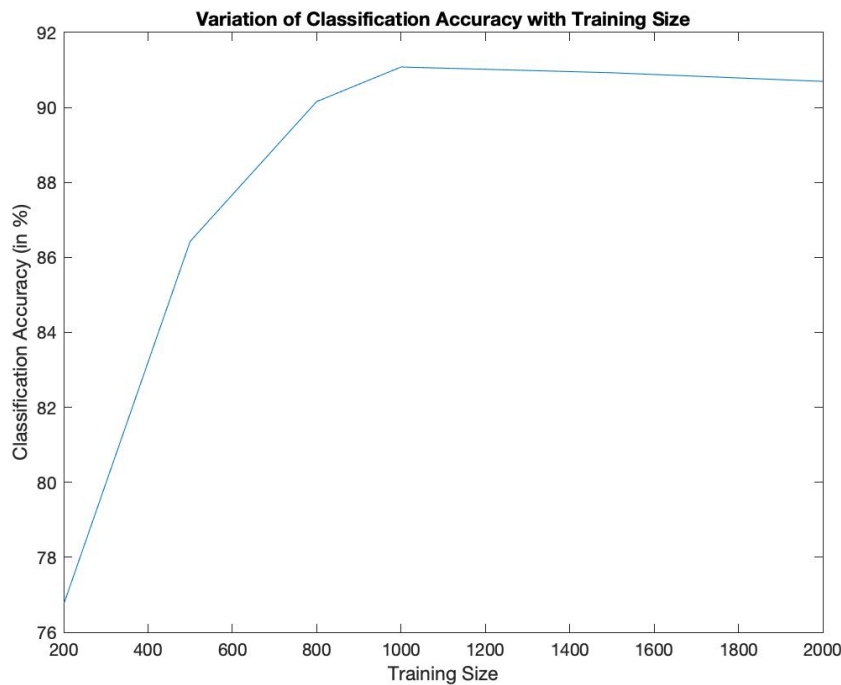


Figure 1: Variation of Classification Accuracy with the Number of Samples available for Training

In addition to this, the confusion matrices for these experimentation were also plotted. The matrices are provided in [Figure 2](#). False positives are really crucial for spam email detection services. If an important mail

gets marked as spam, it may be fatal for the users. So, classification accuracy cannot be always used as an evaluation metric. The reason is that classification accuracy works on the total false cases and does not differentiate between false positives and false negatives. But, in some cases, false positives may be worse than false negatives. In case of spam email detection, that is the scenario. So, we should always check the confusion matrices to arrive at the final model based on the business aspect of the service. For this reason, the confusion matrices for the given experiments were checked and it was observed that $n = 1000$ gives the least number of false positives as well. So, it can be considered as the best model provided by the experimentation.

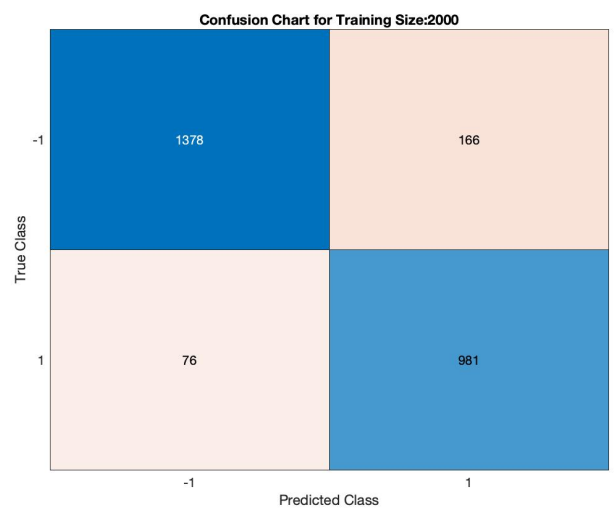
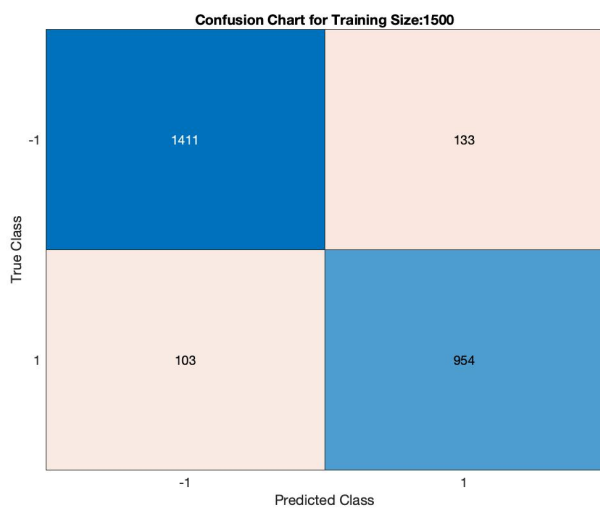
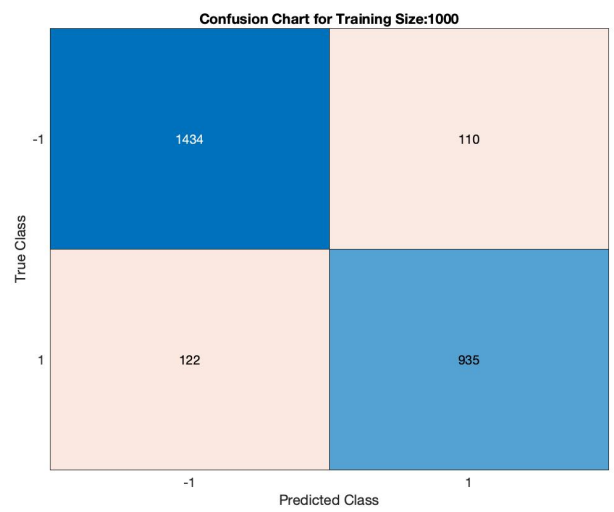
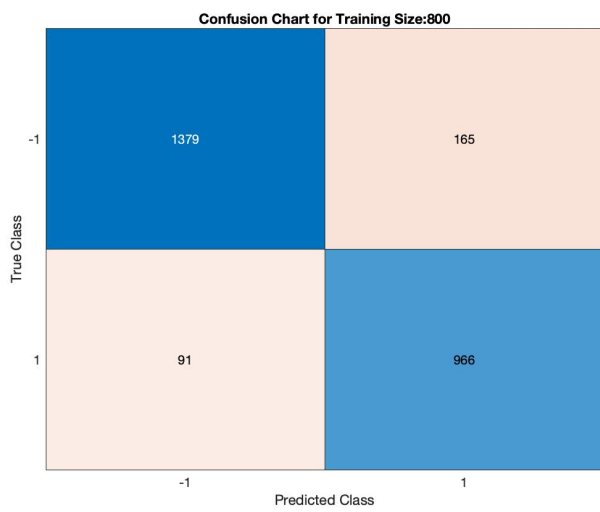
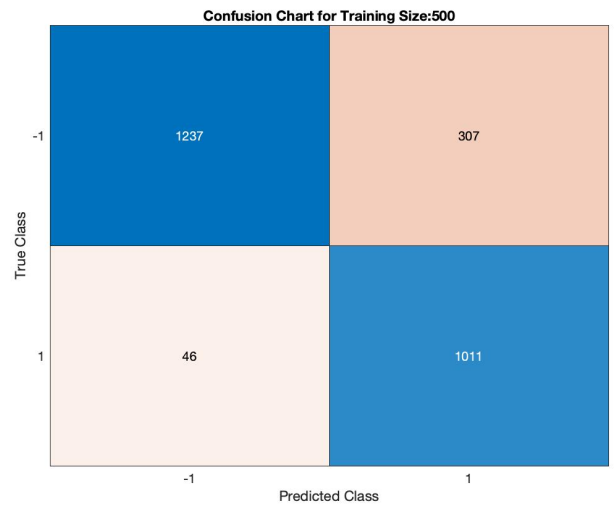
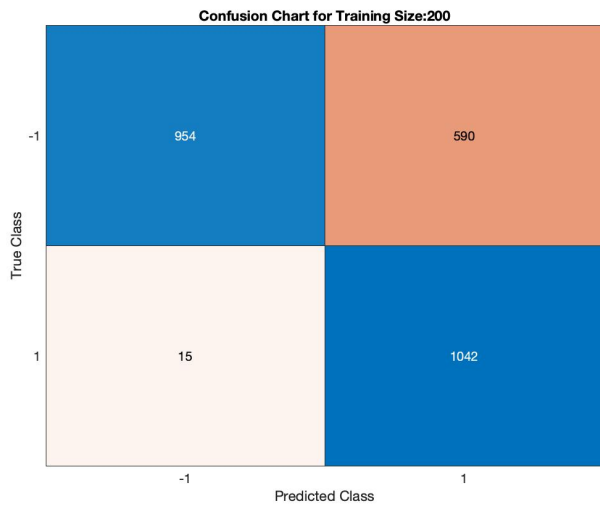


Figure 2: Plot of Confusion Matrices obtained for different training sizes

2 Sparse Logistic Regression: Experiment

The goal of this experiment is to add sparse regularization to the logistic regression framework and observe how it helps to reduce the number of features used for the classification (feature selection). For the simulation of the experimentation, *Alzheimers* dataset is used as an application which is publicly available at: [Alzheimers Dataset](#).

The dataset is pre-divided into training and test samples. There are 172 training samples and 74 test samples in the dataset. Each sample has an associated label of either +1 (Alzheimer's Disease patient) or -1 (Mild Cognitive Impairment patient).

2.1 Code

The code is uploaded to a public Github repo at: [CSE 847 HA 4](#). For reference, the code used to implement the sparse logistic regression is also presented below:

```
1 % import train and test data
2 data = importdata('Data/Alzheimers/ad_data.mat');
3 train_data = data.X_train;
4 train_labels = data.y_train;
5 test_data = data.X_test;
6 test_labels = data.y_test;
7 features = importdata('Data/Alzheimers/feature_name.mat');
8
9 % possible values for the regularization parameter
10 par = [0.01, 0.1, 0.2, 0.3, 0.4, 0.5, 0.6, 0.7, 0.8, 0.9, 1];
11
12 global cur_par
13
14 % perform sparse regularized logistic regression for all possible
15 % parameters
16 for i=1:size(par, 2)
17     cur_par = par(i);
18     [w, c] = logL1Regression(train_data, train_labels, cur_par);
19     scores = sigmoid(test_data * w + c);
20
21     [~, ~, ~, AUC(1,i)] = perfcurve(test_labels, scores, 1); % get the AUC
22     accuracy(1, i) = compute_accuracy(test_data, test_labels, w, c);
23     num_nz_weights(1, i) = nnz(w);
24 %     fprintf('Par= %d: Non-zero weights: %d, Accuracy: %f, AUC: %f\n', cur_par,
25 %         num_nz_weights(1, i), accuracy(1, i), AUC(1,i));
26     fprintf('%d\t %f\t %f\n', num_nz_weights(1, i), accuracy(1, i), AUC(1,i));
27 end
28
29 % Plot the variation of accuracy, AUC and No. of non-zero weights with training
30 % size
31 figure;
32 hold on;
33 X = par;
```

```

32 Y1 = normalize(accuracy, 'norm');
33 Y2 = normalize(AUC, 'norm');
34 Y3 = normalize(num_nz_weights, 'norm');
35 fig = plot(X, Y1);
36 plot(X, Y2)
37 plot(X, Y3);
38 legend('Accuracy', 'AUC', 'No. of Non-zero Weights')
39 xlabel('Regularization Parameter');
40 ylabel('Normalized Metric Scores');
41 title({'Variation in Classification Accuracy, AUC and Number of non-zero weights', 'with Regularization Parameter'})
42 saveas(fig, strcat('Results/L1 Logistic Train/Metric-Variance.jpg'));
43 hold off
44
45
46 function [accuracy] = compute_accuracy(data, labels, weights, c)
47     global cur_par
48
49     % getting the predicted labels and computing accuracy
50     predicted_labels = sigmoid(data * weights + c);
51     predicted_labels(predicted_labels > 0.5) = 1;
52     predicted_labels(predicted_labels <= 0.5) = -1;
53     correct_predictions = sum(predicted_labels == labels);
54     accuracy = (correct_predictions/size(data,1) * 100);
55
56     % create the confusion matrix
57     fig = confusionchart(labels, predicted_labels);
58
59     title(strcat('Confusion Chart for Regularization Parameter: ', string(
        cur_par)));
60     saveas(fig, strcat('Results/L1 Logistic Train/Conf-Chart-Reg-Par-', string(
        cur_par), '.jpg'));
61 end
62
63
64 function [w, c] = logL1Regression(data, labels, par)
65 % OUTPUT w is equivalent to the first d dimension of weights in logistic train
66 % c is the bias term, equivalent to the last dimension in weights in logistic
    train.
67 % Specify the options (use without modification).
68 opts.rFlag = 1; % range of par within [0, 1].
69 opts.tol = 1e-6; % optimization precision
70 opts.tFlag = 4; % termination options.
71 opts.maxIter = 5000; % maximum iterations
72
73 [w, c] = LogisticR(data, labels, par, opts);
74 end

```



```

75
76
77 function [val] = sigmoid(input)
78     % sigmoid function implementation
79     val = 1./(1 + exp(-input));
80 end

```

2.2 Experimental Outcome

As a part of the experiment, the regularization parameters was varied as:

$par = 0.01, 0.1, 0.2, 0.3, 0.4, 0.5, 0.6, 0.7, 0.8, 0.9, 1$

Every different model was judged by 3 imporant metrics: Classification Accuracy, AUC and Number of non-zero weights. $L1$ -regularization leads to sparse models, so it is really interesting to see how change in the regularization parameter can lead to a change in the number of non-zero weights. As number of zero-weights increase in the model, the number of features used to build the classification model reduces. Thus $L1$ -regularization helps to perform Feature Selection.

The result of the experimentation is presented in [Figure 3](#) and [Table 1](#). From the Figure, it can be observed that the number of non-zero weights decreases as the value of regularization parameter is increased. This means that the model is becoming more sparse with increase in the value of the parameter. On the other hand, it can be seen that accuracy is increasing as opposite to the variation in number of non-zero weights. This observation was pretty interesting. Particularly, in the later stages of the experimentation, the model was using only one non-zero weight and it was still giving almost 75% accuracy. A careful look at the dataset made the idea clear. The test data consisted of imbalanced samples where approx. 75% samples were having label -1 . So, the model was simply classifying every sample to -1 and that is the reason why it was getting such a high accuracy using only 1 feature. But it did not produce our intended model.

In order to get some more insights to the problem, in the next stage, the confusion matrices for the experiments were plotted. The matrices are presented in [Figure 4](#). After the parameter value reaches 0.7, it starts assigning all the labels to -1 , so there's no false negatives or true positives. Although it is able to achieve good classification accuracy, it is not a good model. For this reason, in case of imbalanced data, classification accuracy cannot serve as an appropriate metric. AUC (Area Under Curve) is a better metric compared to accuracy in such a scenario. Intuitioanlly, AUC denotes the probability of a model assigning higher value to a random positive example in comparison to a random negative example. In terms of AUC, the second model of the experimentation, i.e. parameter value of 0.1, can be considered to be the best model out of all as it achieved the highest AUC score.

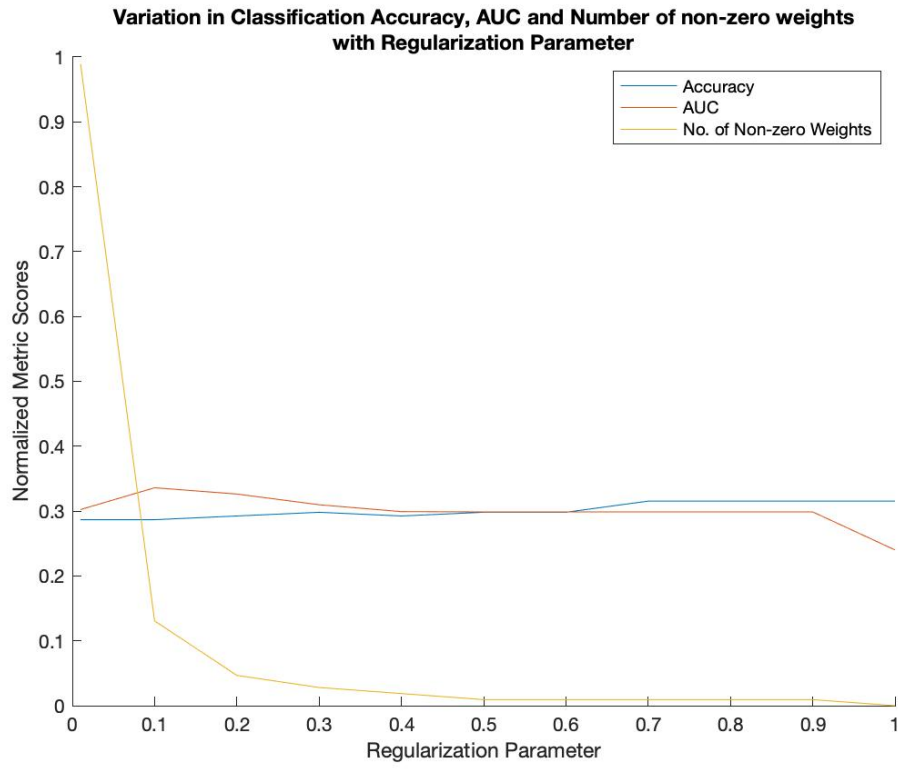


Figure 3: Plot of different metrics obtained for different values of regularized parameters

Par	No. of Non-zero Weights	Accuracy	AUC
0.01	106	67.567568	0.629665
0.1	14	67.567568	0.699522
0.2	5	68.918919	0.679426
0.3	3	70.27027	0.644976
0.4	2	68.918919	0.622967
0.5	1	70.27027	0.62201
0.6	1	70.27027	0.62201
0.7	1	74.324324	0.62201
0.8	1	74.324324	0.62201
0.9	1	74.324324	0.62201
1	0	74.324324	0.5

Table 1: Different Metric Scores for different values of the regularization parameter

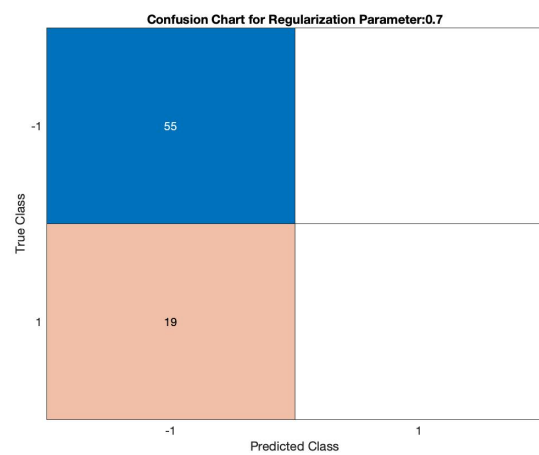
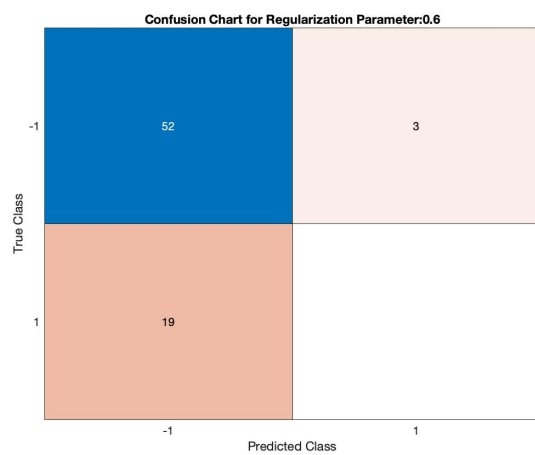
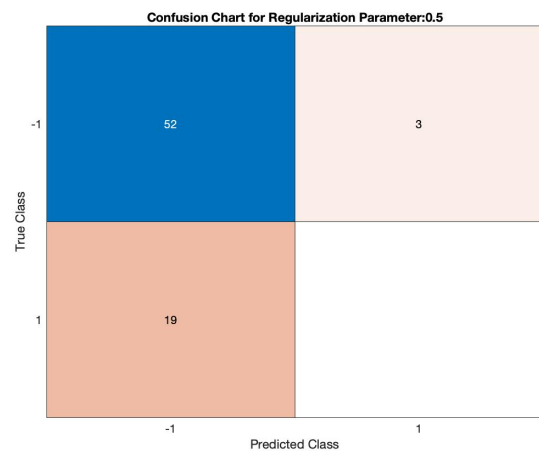
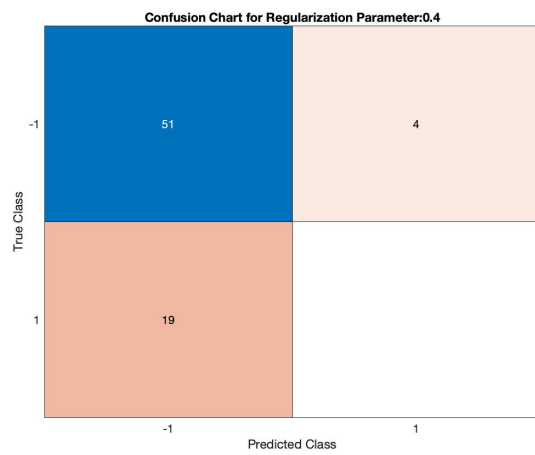
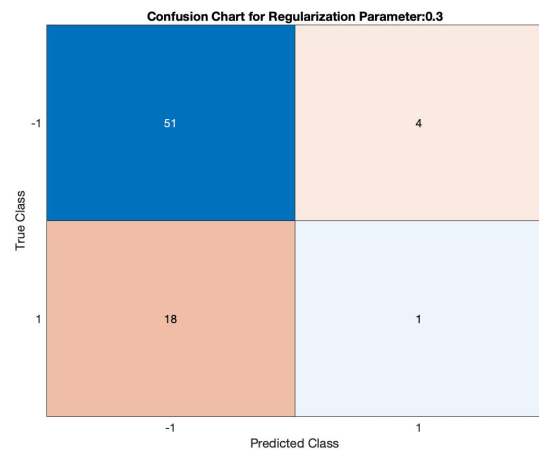
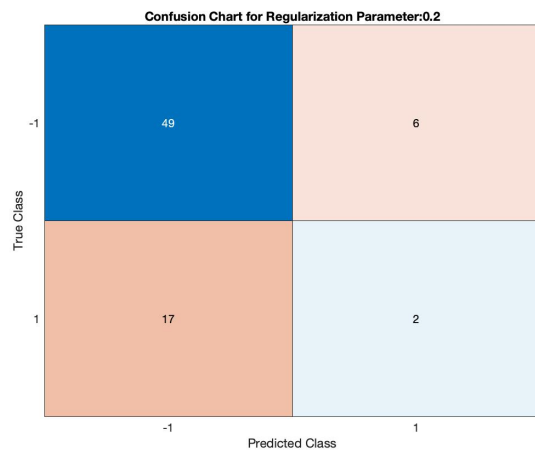
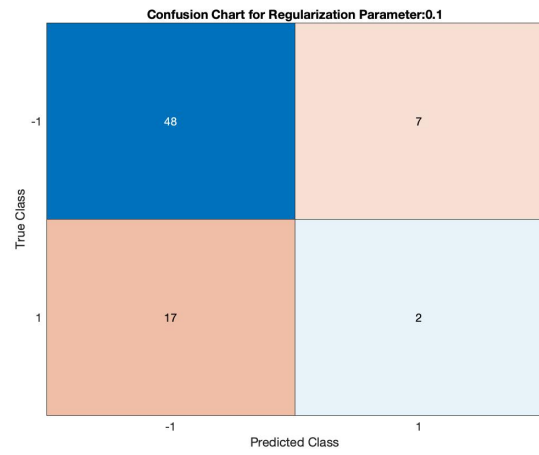
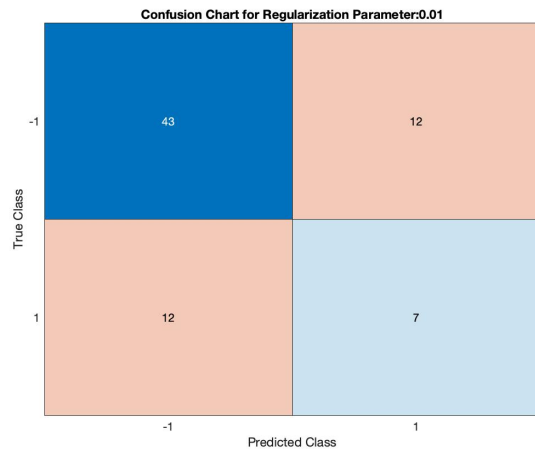


Figure 4: Plot of Confusion Matrices obtained for different values of the Regularization Parameter