

UNIVERSITY OF VICTORIA  
Department of Electrical and Computer  
Engineering  
ECE 403/503 Optimization for Machine  
Learning  
LABORATORY REPORT

Experiment No: 4

Title: Breast Cancer Diagnosis via Logistic Regression

Date of Experiment: 11/21/2019

Report Submitted on: 11/28/2019

To: TA

Laboratory Group No.:

Name(s): Kurt Elliott & Raj Deepak S. N.

## Objectives:

In this experiment, we investigate using a logistic regression to build a model for breast cancer patients that can classify new samples as either benign or malignant. Then a gradient descent algorithm with backtracking line search will be used to optimize the model for best results.

## Introduction:

The logistic regression model is a regression technique that uses statistical probability to classify a model. The model classifies based on probability, if the values of the model is  $< 0.5$  the model is classified as -1 (benign), if the model is  $> 0.5$  the model is classified as 1(malignant). A Gradient descent algorithm is then applied to the model to find the optimum solution. A backtracking line search algorithm is also applied to the gradient descent to insure that it takes the most optimal path to the optimized solution of the logistic regression model.

The dataset to be used has 30 feature values from each of 569 patients of which 357 patients were diagnosed as benign while the other 212 patients were diagnosed as malignant.

## Results:

The following is the MATLAB result from the script used in this lab, and a brief description of how the program is implemented. The script was sliced into sections for readability, therefore refer to appendix for full script.

The script used in this lab consists of three primary sections. The first section of the code is used to load and extract the training and testing datasets. The datasets consists of input vectors  $\mathbf{x}$ , which consist of 30 components

(features), and  $y$  output vectors, which consist of 2 components (malignant & benign).

```
1 %% Step 1
2 % load datasets into matlab
3 - load D_bc_te
4 - load D_bc_tr
5
6 %% Step 2
7 % format testing and training datasets and normalization dataset
8
9 - NumFeatures = 30;
10 - K = 75;
11
12 - Xtrain = zeros(30,480);
13 - for i = 1:30
14 -     xi = D_bc_tr(i,:);
15 -     mi = mean(xi);
16 -     vi = sqrt(var(xi));
17 -     Xtrain(i,:) = (xi - mi)/vi;
18 - end
19 - Xtest = zeros(30,89);
20 - for i = 1:30
21 -     xi = D_bc_te(i,:);
22 -     mi = mean(xi);
23 -     vi = sqrt(var(xi));
24 -     Xtest(i,:) = (xi - mi)/vi;
25 - end
26
27 - ytrain = D_bc_tr(31,:);
28 - ytest = D_bc_te(31,:);
29
30 - Dtrain = [Xtrain;ytrain];
31 - Dtest = [Xtest;ytest];
```

The next step is to build the model, this is done through two functions that were made in matlab. First function is the “objective function” (refer to lab manual E4.2).

```

1  function f = f_wcdc(w,D)
2  -   X = D(1:30,:);
3  -   y = D(31,:);
4  -   N = length(y);
5  -   f = 0;
6  -   wt = w(:)';
7  -   for i = 1:N
8  -       xi = [X(:,i);1];
9  -       fi = log(1+exp(-y(i)*(wt*xi)));
10 -       f = f+fi;
11 -   end
12 -   f = f/N;

```

The second function is the gradient of the objective function (refer to lab manual E4.3).

```

1  function g = g_wcdc(w,D)
2  -   X = D(1:30,:);
3  -   y = D(31,:);
4  -   N = length(y);
5  -   g = zeros(31,1);
6  -   wt = w(:)';
7  -   for i = 1:N
8  -       xi = [X(:,i);1];
9  -       ei = exp(-y(i)*(wt*xi));
10 -       ci = -y(i)*ei/(1+ei);
11 -       gi = ci*xi;
12 -       g = g+gi;
13 -   end
14 -   g = g/N;

```

Next is the model training section. The data is formatted into matrices and the learning rate and the direction of the gradient descent are obtained. Then the model is trained to a episodal threshold.

```

36 %% Step 4
37 % Minimize f(w)
38 - eps = 10^-9;
39 - w = zeros(NumFeatures+1,1);
40 - f = zeros(1,K);
41 - k = 1;
42
43 % step 2: find d
44 - gk = g_wcdc(w,Dtrain);
45 - dk = -gk;
46
47 % step 3: find alpha
48 - alpha = bt_lsearch(w,dk,'f_wcdc','g_wcdc',Dtrain);
49
50 %% Training
51 - while ((norm(alpha*dk) >= eps) && (k<+K))
52
53 -     w = w + alpha*dk;
54
55 -     % step 2: find d
56 -     gk = g_wcdc(w,Dtrain);
57 -     dk = -gk;
58
59 -     % step 3: find alpha
60 -     alpha = bt_lsearch(w,dk,'f_wcdc','g_wcdc',Dtrain);
61
62 -     %check if f is decreasing.
63 -     f(k) = f_wcdc(w,Dtrain);
64
65 -     k = k + 1;
66 - end

```

The final part is testing the logistic regression model by comparing prediction values with the test model prediction in the test dataset. This was done by checking the predictions label, than checking it will the tests label and counting the false positive and false negatives. A misclassification ratio was than obtained by adding the false positives and negatives together and dividing by 100.

```

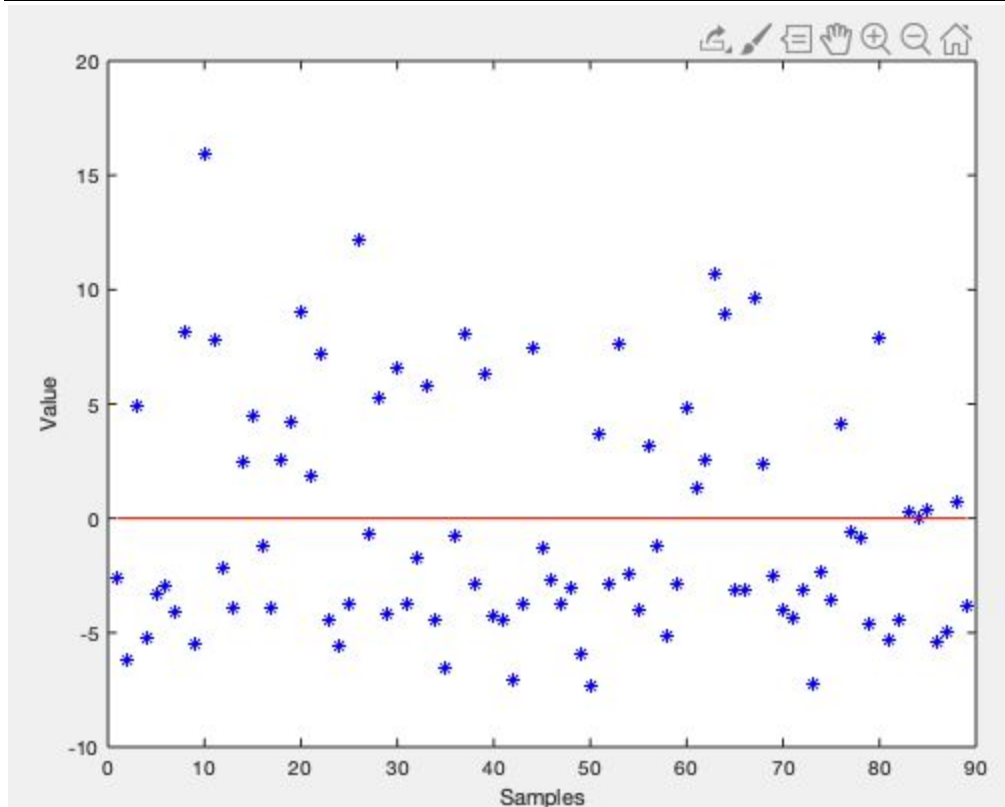
69 - %% Testing
70 - Dt = [Xtest; ones(1,89)];
71 - Result = w'*Dt;
72 - TestLabel = zeros(1,length(Result));
73 - FalsePos = 0;
74 - FalseNeg = 0;
75 -
76 - if PLOT == 1
77 -     AfterTr = Result;
78 -     Z = zeros(1,length(AfterTr));
79 -
80 -     figure, plot(AfterTr,'b*');
81 -     hold on
82 -     plot(Z,'r-');
83 -     xlabel('Samples');
84 -     ylabel('Value');
85 - end
86 -
87 - for ii = 1:length(Result)
88 -     if Result(ii) < 0
89 -         TestLabel(ii) = -1;
90 -         if ytest(ii) > 0
91 -             FalseNeg = FalseNeg + 1;
92 -         end
93 -     else
94 -         TestLabel(ii) = 1;
95 -         if ytest(ii) > 0
96 -             FalsePos = FalsePos + 1;
97 -         end
98 -     end
99 - end

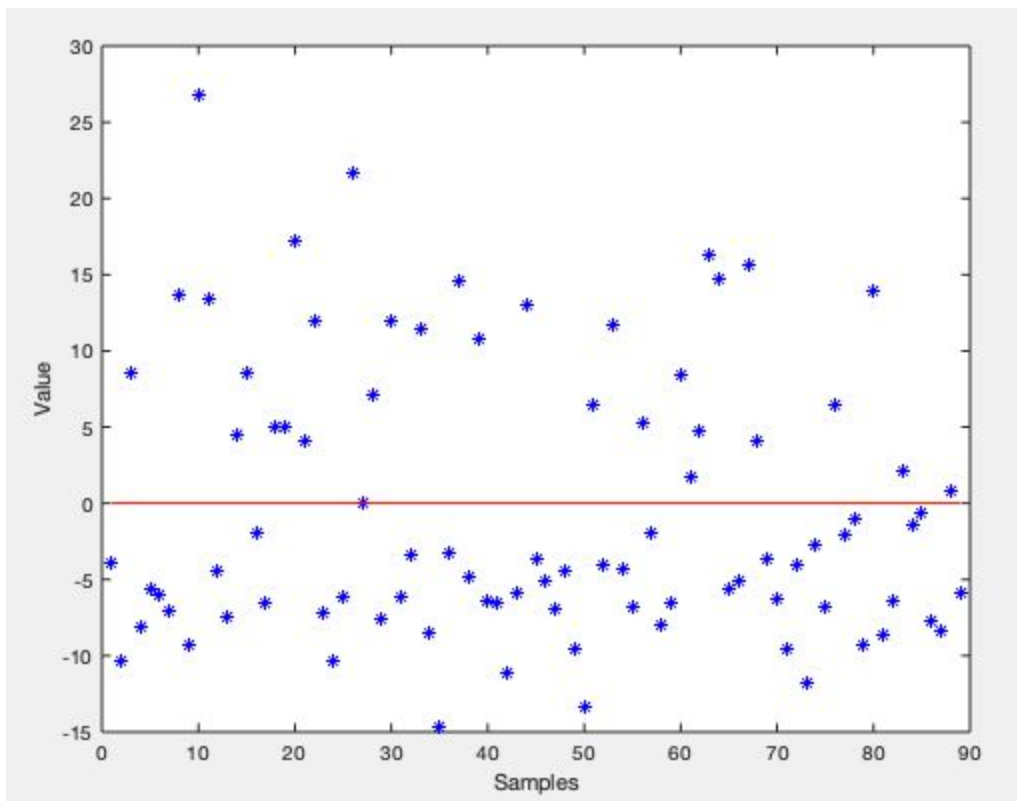
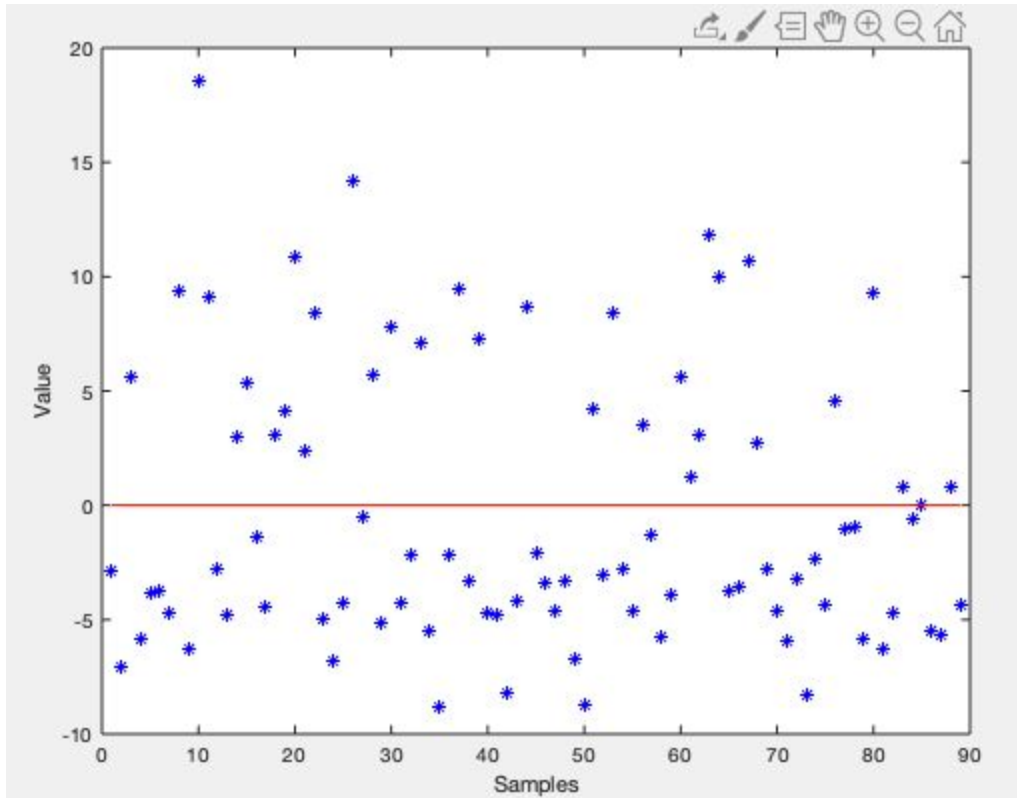
```

## Discussion:

After training the model, analysis was done on the testing dataset and the results from the true values and the model predicted values were plotted. This procedure was done 3 times to evaluate different iteration values  $K=5$ ,  $K=12$ ,  $K=75$ . The different iteration values gave the following results for false positives and false negatives and are plotted below.

K-Iterations	K=5	K=12	K=75
False positive	31	31	32
False negative	1	1	0
mis-classification Ratio	0.32	0.32	0.32







## Conclusions:

The logistic regression classification model with gradient descent and backtracking line search for optimization was investigated for breast cancer diagnosis. The results show that all the points above and below the line are cases of benign and malignant respectively, with misclassified points on or in very close proximity to the line. The model is 68% efficient in separating benign and malignant cases. It has a high number of false positive results which means that a high number of patients are diagnosed with breast cancer when in actuality it is not true.

# Appendix:

---

## Table of Contents

Step 1 .....	1
Step 2 .....	1
Step 3 .....	1
Step 4 .....	2
Training .....	2
Testing .....	2

## Step 1

load datasets into matlab

```
load D_bc_te
load D_bc_tr
```

*Error using evalin  
Unrecognized function or variable 'experiment3'.*

## Step 2

format testing and training datasets and normalization dataset

```
NumFeatures = 30;
K = 75;

Xtrain = zeros(30,480);
for i = 1:30
    xi = D_bc_tr(i,:);
    mi = mean(xi);
    vi = sqrt(var(xi));
    Xtrain(i,:) = (xi - mi)/vi;
end
Xtest = zeros(30,89);
for i = 1:30
    xi = D_bc_te(i,:);
    mi = mean(xi);
    vi = sqrt(var(xi));
    Xtest(i,:) = (xi - mi)/vi;
end

ytrain = D_bc_tr(31,:);
ytest = D_bc_te(31,:);

Dtrain = [Xtrain;ytrain];
Dtest = [Xtest;ytest];
```

## Step 3

Objective function and gradient