

## CS171 Assn2 Report

### Part 0:

Downloaded the dataset and loaded it into the workspace by using the load function and the import data feature on MatLab. Use the load(breastcancerwisconsin) to load the mat file and then you can access breastcancerwisconsin1 which is the matrix of data points.

Missing Values: I took care of the missing values by replacing them by 0 so that they would not affect the distance function.

I also took out the serial number of each data point because it was not useful.

### Part 1:

```
9 function y_pred = knn_classifier(x_test, x_train, y_train, k, p)
10     [sizeTest, ~] = size(x_test); %get size of x_test
11     [sizeTrain, ~] = size(x_train);
12
13     y_pred = zeros([sizeTest 1]);
14     %calculate distance
15     for i=1:sizeTest
16         distanceBetweenAllPoints = zeros(sizeTrain, 2);
17         distanceBetweenAllPoints(:,2) = y_train ;
18         %calculate distance between x_test data point and all x_train
19         for j=1:sizeTrain
20             distanceBetweenAllPoints(j,1) = distance(x_test(i,:), x_train(j,:), p);
21         end
22         %sort distances
23         distanceBetweenAllPoints = sortrows(distanceBetweenAllPoints);
24
25         %get k nearest neighbors by grabbing first k data points, then use
26         %mode on the vector for most frequent value
27         distanceBetweenAllPoints = distanceBetweenAllPoints(1:k,2);
28         y_pred(i,1) = mode(distanceBetweenAllPoints(1:k,1));
29     end
30 end
31
32
33
34
```

Command Window

```
> 0.985714>>
```

I got a 98% accuracy when using my knn algorithm with an 80/20 split. I split the first 80% of the data as the training data and the last 20% as the test data.

The distance function used is the one from Assignment 1.

This function is implemented by calculating the distance between the test data point and every training data point. After calculating the distance I sorted the list and chose the first k datapoints or neighbors. Then I calculated the mode of those k neighbors to determine the label of the test datapoint.

```

7 %get xtrain, xtest and ytrain
8 %80percent train 20% is test
9 [n,~] = size(xdata);
10 eightyCutOff = floor(n*.8);
11 eightyCutOffPlus1 = eightyCutOff+1;
12
13 xtrain = xdata(1:eightyCutOff,:);
14 xtest = xdata(eightyCutOffPlus1:n,:);
15 ytrain = ydata(1:eightyCutOff,:);
16
17 %test with k = 1 and p = 2
18 k = 1;
19 p = 2;
20 ypred = knn_classifier(xtest, xtrain, ytrain, k, p);
21 sizeTest = n - eightyCutOff;
22
23 %compute accuracy
24 count = 0;
25
26 for i=1:sizeTest
27     if(ypred(i,1) == ydata(eightyCutOff+i,1))
28         count = count +1;
29     end
30     fprintf('%i , %i\n', ypred(i,1), ydata(eightyCutOff+i,1));
31 end
32 fprintf('%f',count/sizeTest);
33

```

//Test file used to get 98% accuracy. 'Question1.m'  
 I take floor of 80% of the size of the xdata for the xtrain set  
 Then I take that 80% cut off plus 1 for xtest

Part2.1:

```

4
5 %shuffle rows
6 data = data(randperm(end),:);
7 xdata = data(:,(2:10));
8 ydata = data(:,11);
9
10 numFolds = 10;
11 [n,~] = size(xdata);
12 k = 1;
13 p = 2;
14 sizeTest = floor(n/numFolds);
15 ypred = zeros(sizeTest,1);
16 a=0;
17 accuracy = zeros(numFolds, 1);
18 sensitivity = zeros(numFolds, 1);
19 specificity = zeros(numFolds, 1);

```

I shuffle before cross validating so that there is diverse data in my partitions.

```

21 - for w=1:numFolds
22 -     if(w==1)
23 -         xtest = xdata(1:sizeTest,:);
24 -         xtrain = xdata(sizeTest+1:n,:);
25 -         ytrain = ydata(sizeTest+1:n,:);
26 -         ypred = knn_classifier(xtest, xtrain, ytrain, k, p);
27 -     else
28 -         a = (w-1)*(sizeTest+1);
29 -         xtest = xdata(a+1:a+sizeTest,:);
30 -         xtrain = xdata(1:a,:);
31 -         xtrain = [xtrain; xdata(a+sizeTest+1:n,:)];
32 -         ytrain = ydata(1:a,:);
33 -         ytrain = [ytrain; ydata(a+sizeTest+1:n,:)];
34 -         fprintf('test range: %i -%i, train trange: 1-%i, %i-n', a+1,a+sizeTest,a,a+sizeTest+1);
35 -         ypred = knn_classifier(xtest, xtrain, ytrain, k, p);
36 -     end
37 -     %compute statistics here
38 -     %compute accuracy
39 -     count = 0;
40 -     truePositive = 0;
41 -     trueNegative = 0;
42 -     falsePositive = 0;
43 -     falseNegative = 0;
44 -     for i=1:sizeTest
45 -         if(ypred(i,1) == ydata(a+i,1))
46 -             count = count +1;
47 -         end

```

For Cross Validation I randomly shuffled all the rows of the data. Then I partitioned it into 10 equal pieces.

```

3 -         if(ypred(i,1) == 4 && ydata(a+i,1) == 4)
4 -             truePositive = truePositive + 1;
5 -         end
6 -         if(ypred(i,1) == 2 && ydata(a+i,1) == 2)
7 -             trueNegative = trueNegative + 1;
8 -         end
9 -         if(ypred(i,1) == 4 && ydata(a+i,1) == 2)
10 -             falsePositive = falsePositive + 1;
11 -         end
12 -         if(ypred(i,1) == 2 && ydata(a+i,1) == 4)
13 -             falseNegative = falseNegative + 1;
14 -         end
15 -     end
16 -     accuracy(w,1) = count/sizeTest;
17 -     sensitivity(w,1) = truePositive/(truePositive+falseNegative);
18 -     specificity(w,1) = trueNegative/(trueNegative+falsePositive);
19 -     fprintf('fold %i: accuracy - %f, sensitivity - %f, specificity - %f\n',w, accuracy(w,1),sensitivity(w,1),specificity(w,1));
20 - end
21 - fprintf('Mean of accuracy: %f, Std: %f\n', mean(accuracy), std(accuracy));
22 - fprintf('Mean of sensitivity: %f, Std: %f\n', mean(sensitivity), std(sensitivity));
23 - fprintf('Mean of specificity: %f, Std: %f\n', mean(specificity), std(specificity));

```

```

Command Window

fold 1: accuracy - 0.985507, sensitivity - 1.000000, specificity - 0.977778
fold 2: accuracy - 0.942029, sensitivity - 0.789474, specificity - 1.000000
fold 3: accuracy - 0.942029, sensitivity - 0.937500, specificity - 0.945946
fold 4: accuracy - 0.971014, sensitivity - 1.000000, specificity - 0.955556
fold 5: accuracy - 0.956522, sensitivity - 0.962963, specificity - 0.952381
fold 6: accuracy - 0.971014, sensitivity - 0.913043, specificity - 1.000000
fold 7: accuracy - 0.956522, sensitivity - 0.900000, specificity - 0.979592
fold 8: accuracy - 0.956522, sensitivity - 0.913043, specificity - 0.978261
fold 9: accuracy - 0.985507, sensitivity - 0.952381, specificity - 1.000000
fold 10: accuracy - 0.927536, sensitivity - 0.888889, specificity - 0.952381
Mean of accuracy: 0.959420, Std: 0.019081
Mean of sensitivity: 0.925729, Std: 0.061621
Mean of specificity: 0.974189, Std: 0.021482

```

I got around 96% accuracy for cross validation with a std of .019.

The function will iterate between all partitions as the test set. So there is 1 test set and 9 training at each iteration.

## **Part2.2:**

For k=1:10 and p=1 and p=2 cross validated performance I got these metrics(with graphs below).

k:1, p:1

Mean of accuracy: 0.955072, Std: 0.018648  
Mean of sensitivity: 0.913699, Std: 0.071672  
Mean of specificity: 0.978242, Std: 0.020230

k:1, p:2

Mean of accuracy: 0.900000, Std: 0.160434  
Mean of sensitivity: 0.836362, Std: 0.236064  
Mean of specificity: 0.932327, Std: 0.132157

k:2, p:1

Mean of accuracy: 0.885507, Std: 0.150992  
Mean of sensitivity: 0.780124, Std: 0.233288  
Mean of specificity: 0.940881, Std: 0.127053

k:2, p:2

Mean of accuracy: 0.892754, Std: 0.153255  
Mean of sensitivity: 0.799852, Std: 0.227857  
Mean of specificity: 0.940706, Std: 0.127073

k:3, p:1

Mean of accuracy: 0.910145, Std: 0.163082  
Mean of sensitivity: 0.864192, Std: 0.242342  
Mean of specificity: 0.934214, Std: 0.132467

k:3, p:2

Mean of accuracy: 0.908696, Std: 0.162688  
Mean of sensitivity: 0.864192, Std: 0.244836  
Mean of specificity: 0.932131, Std: 0.131477

k:4, p:1

Mean of accuracy: 0.900000, Std: 0.164882  
Mean of sensitivity: 0.828537, Std: 0.250379  
Mean of specificity: 0.936101, Std: 0.133044

k:4, p:2

Mean of accuracy: 0.905797, Std: 0.166544  
Mean of sensitivity: 0.849180, Std: 0.254851  
Mean of specificity: 0.934214, Std: 0.132467

k:5, p:1

Mean of accuracy: 0.913043, Std: 0.163824  
Mean of sensitivity: 0.869963, Std: 0.242439  
Mean of specificity: 0.934214, Std: 0.132467

k:5, p:2

Mean of accuracy: 0.911594, Std: 0.163746  
Mean of sensitivity: 0.872954, Std: 0.245974  
Mean of specificity: 0.931889, Std: 0.132295

k:6, p:1

Mean of accuracy: 0.907246, Std: 0.157313  
Mean of sensitivity: 0.844747, Std: 0.239601  
Mean of specificity: 0.938184, Std: 0.126548

k:6, p:2

Mean of accuracy: 0.905797, Std: 0.162717  
Mean of sensitivity: 0.846612, Std: 0.240121  
Mean of specificity: 0.936101, Std: 0.133044

k:7, p:1

Mean of accuracy: 0.902899, Std: 0.160522  
Mean of sensitivity: 0.844747, Std: 0.239601  
Mean of specificity: 0.931889, Std: 0.132295

k:7, p:2

Mean of accuracy: 0.914493, Std: 0.164457  
Mean of sensitivity: 0.875487, Std: 0.244369  
Mean of specificity: 0.933775, Std: 0.132909

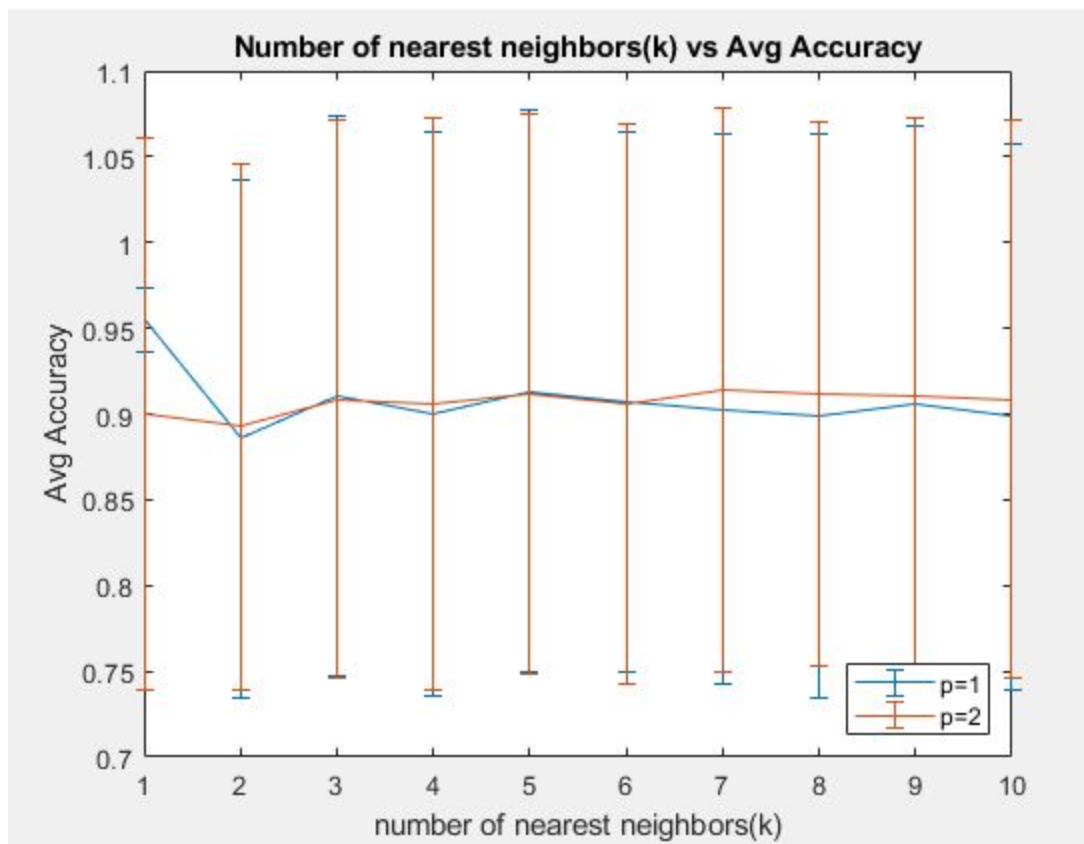
k:8, p:1

Mean of accuracy: 0.898551, Std: 0.164535  
Mean of sensitivity: 0.826440, Std: 0.249754  
Mean of specificity: 0.933775, Std: 0.132909

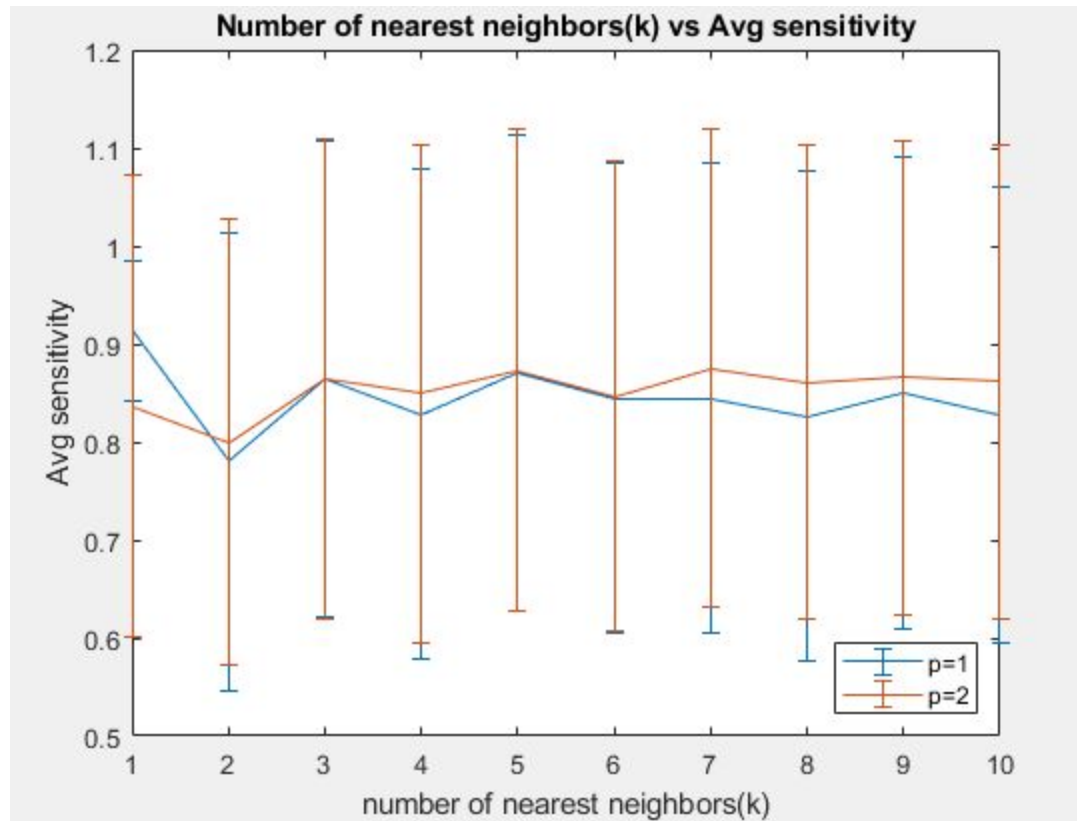
k:8, p:2

Mean of accuracy: 0.911594, Std: 0.158679

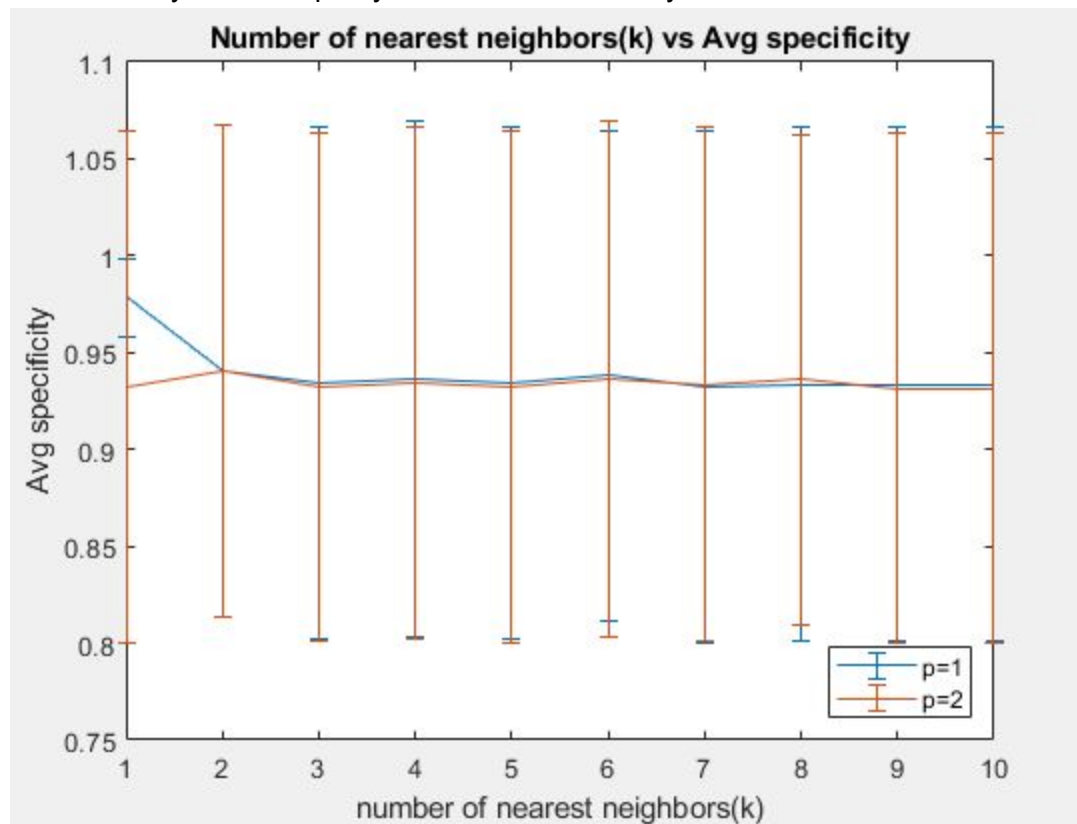
Mean of sensitivity: 0.861201, Std: 0.241471  
 Mean of specificity: 0.935859, Std: 0.126449  
 k:9, p:1  
 Mean of accuracy: 0.905797, Std: 0.162142  
 Mean of sensitivity: 0.850060, Std: 0.240065  
 Mean of specificity: 0.933775, Std: 0.132909  
 k:9, p:2  
 Mean of accuracy: 0.910145, Std: 0.162939  
 Mean of sensitivity: 0.865963, Std: 0.241851  
 Mean of specificity: 0.931692, Std: 0.131916  
 k:10, p:1  
 Mean of accuracy: 0.898551, Std: 0.159201  
 Mean of sensitivity: 0.827202, Std: 0.233207  
 Mean of specificity: 0.933775, Std: 0.132909  
 k:10, p:2  
 Mean of accuracy: 0.908696, Std: 0.162544  
 Mean of sensitivity: 0.861963, Std: 0.241924  
 Mean of specificity: 0.931692, Std: 0.131916



For Accuracy p=1 and k=1 yield the best accuracy.



For sensitivity,  $k=1$  and  $p=1$  yield the best sensitivity.



For specificity,  $k=1$  and  $p=1$  yield the best specificity

```
for i=1:sizeTest
    if(ypred(i,1) == ydata(a+i,1))
        count = count +1;
    end
    if(ypred(i,1) == 4 && ydata(a+i,1) == 4)
        truePositive = truePositive + 1;
    end
    if(ypred(i,1) == 2 && ydata(a+i,1) == 2)
        trueNegative = trueNegative + 1;
    end
    if(ypred(i,1) == 4 && ydata(a+i,1) == 2)
        falsePositive = falsePositive + 1;
    end
    if(ypred(i,1) == 2 && ydata(a+i,1) == 4)
        falseNegative = falseNegative + 1;
    end
end

accuracy(w,1) = count/sizeTest;
sensitivity(w,1) = truePositive/(truePositive+falseNegative);
specificity(w,1) = trueNegative/(trueNegative+falsePositive);
```

I measured sensitivity by using the formula  $\text{truePositive}/(\text{truePositive}+\text{falseNegative})$

And I measured specificity by  $\text{trueNegative}/(\text{trueNegative}+\text{falsePositive})$

I used these rules to figure out which is positive or negative for specificity and sensitivity:

positive = 4

negative = 2

true positive = got 4 and actual is 4

false positive = got 4 and actual is 2

true negative = got 2 and actual is 2

false negative = got 2 and actual is 4



### Part3.1

I used

Train\_perceptron function

```
1 function w = train_perceptron(input_x,input_y, w_init)
2     w = w_init;
3     numData = size(input_x(:,1),1);
4     attributeSize = size(input_x(1,:),2);
5     %set learning rate alpha
6     alpha = .8;
7
8     %for each data point
9     for i=1:numData
10         net = 0;
11         %compute net
12         for j=1:attributeSize
13             net = net + input_x(i,j)*w(j,1);
14         end
15         %compute sigmoid of function
16         net = 1.0/(1+exp(-net));
17         %calculate error
18         change = input_y(i, 1) - net;
19         %adjust weight
20         for j=1:attributeSize
21             w(j,1) = w(j,1) + alpha*change*input_x(i,j);
22         end
23     end
24
25 end
```

For train\_perceptron I went through each data point and implemented this formula from the slides:

**For each training data point  $x$**

1. Compute  $y = \text{sigmoid}(\sum_{i=1}^n w_i x_i)$
2. Compute error =  $(d_j - y)$
3. "Correct" the weights:  $w_i \leftarrow w_i + \text{error} * x_i$

For the error I used gradient descent with an alpha of .8 .



```

1 function y_pred = classify_perceptron(input_x,w)
2 %CLASSIFY_PERCEPTRON Summary of this function goes here
3 % Detailed explanation goes here
4 sizeX = size(input_x(:,1),1);
5 y_pred = zeros(sizeX,1);
6 for m=1:sizeX
7     sum = 0;
8     for i=1:size(input_x(1,:),2)
9         sum = sum + w(i,1)*input_x(m,i);
10    end
11    if(sign(sum) < 0)
12        y_pred(m,1) = 4;
13    else
14        y_pred(m,1) = 2;
15    end
16 end
17 end
18
19

```

For classify\_perceptron function for each data point I multiplied each weight by the attribute. Then I used sign as the activation function. I defined -1 as '4' class and +1 as '2' class.

#### Evaluating Perceptron:

First I shuffled the rows to create sufficient diversity between partitions of cross validation.

```

6 %shuffle rows
7 data = data(randperm(end),:);
8 xdata = data(:,(2:10));
9 ydata = data(:,11);
10

```

```

26 - for d=1:numFolds
27 -     %when q=1 initialize w to 0 else to random
28 -     for q=1:2
29 -         %q=1;
30 -         w = zeros(size(xdata(1,:), 2),1);
31 -
32 -         %get xtrain and ytrain
33 -         a = ((d-1)*(sizeTest+1));
34 -         xtest = xdata(a+1:a+sizeTest,:);
35 -         xtrain = xdata(1:a,:);
36 -         xtrain = [xtrain; xdata(a+sizeTest+1:n,:)];
37 -         ytrain = ydata(1:a,:);
38 -         ytrain = [ytrain; ydata(a+sizeTest+1:n,:)];
39 -
40 -         %convert y values
41 -         y_input = zeros(size(ytrain,1),1);
42 -
43 -         for i=1:n-sizeTest
44 -             if(ytrain(i,1) == 2)
45 -                 y_input(i,1) = 1;
46 -             else
47 -                 y_input(i,1) = 0;
48 -             end
49 -         end

```

For each fold I calculated averages based on if I was initializing the weights to zero or initializing the weights to random values. This part obtains the xtrain and ytrain for the current partition. I convert the ytraining data to +1 if it is '2' class or 0 if it is '4' class

```

50 -
51 -     %used for when q=2 so that we can get 10 independent runs of
52 -     %the training so that randomness does not affect it
53 -     avgWeight = zeros(9, 10);
54 -     avgWeight(:,1) = w;
55 -     if(q==2)
56 -         %initialize w to random values from [-25,25]
57 -         for u=2:10
58 -             w = -25 + 50*rand(size(xdata,2),1);
59 -             avgWeight(:,u) = train_perceptron(xtrain,y_input,w);
60 -         end
61 -         for u=1:9
62 -             actualWeight(u,1) = mean(avgWeight(u,:),2);
63 -         end
64 -         w = actualWeight;
65 -     end
66 -     w = train_perceptron(xtrain, y_input, w);
67 -
68 -     %compute y label
69 -     ypred = classify_perceptron(xtest,w);

```

If I am initializing the weights to random values then i do 10 iterations of training and find the weight of each iteration then I average it to find the weights to give to the perceptron classifier- this way I can reduce randomness from affecting the results.

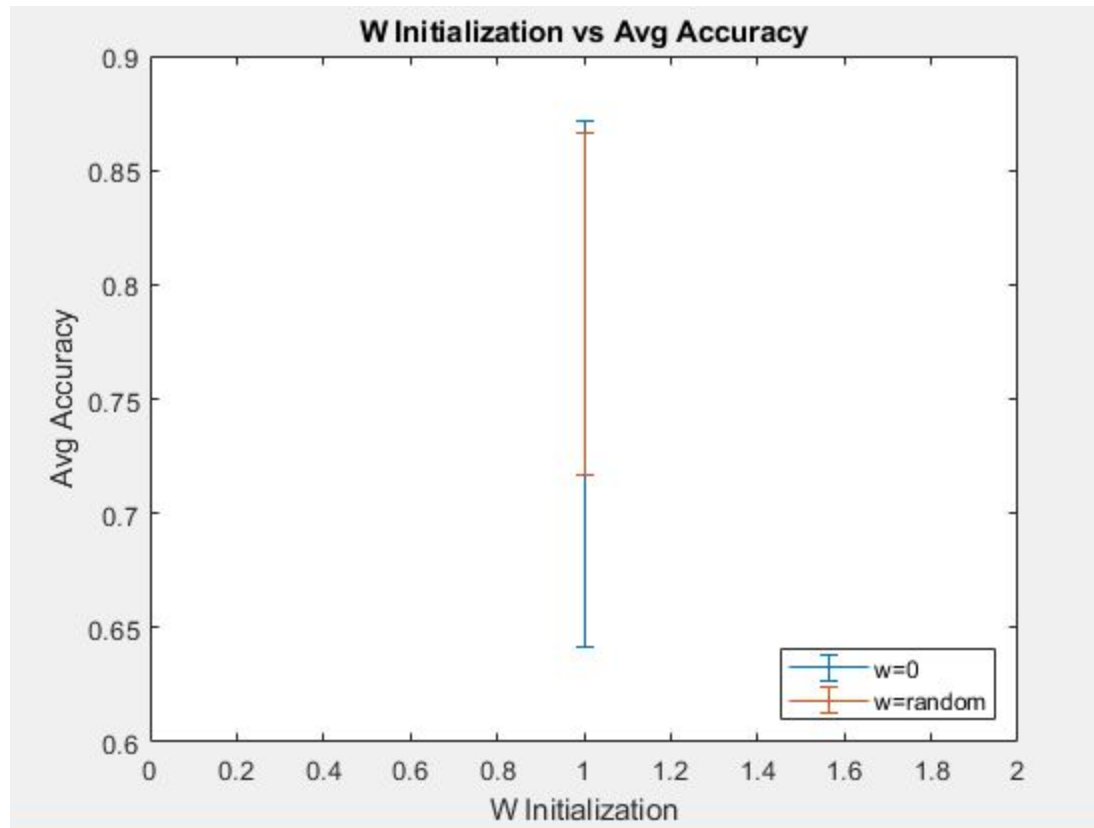
```

72 - count = 0;
73 - truePositive = 0;
74 - trueNegative = 0;
75 - falsePositive = 0;
76 - falseNegative = 0;
77 - for i=1:sizeTest
78 -     if(ypred(i,1) == ydata(a+i,1))
79 -         count = count +1;
80 -     end
81 -     if(ypred(i,1) == 4 && ydata(a+i,1) == 4)
82 -         truePositive = truePositive + 1;
83 -     end
84 -     if(ypred(i,1) == 2 && ydata(a+i,1) == 2)
85 -         trueNegative = trueNegative + 1;
86 -     end
87 -     if(ypred(i,1) == 4 && ydata(a+i,1) == 2)
88 -         falsePositive = falsePositive + 1;
89 -     end
90 -     if(ypred(i,1) == 2 && ydata(a+i,1) == 4)
91 -         falseNegative = falseNegative + 1;
92 -     end
93 -
94 - end
95 - accuracy(d,q) = count/sizeTest;
96 - sensitivity(d,q) = truePositive/(truePositive+falseNegative);
97 - specificity(d,q) = trueNegative/(trueNegative+falsePositive);

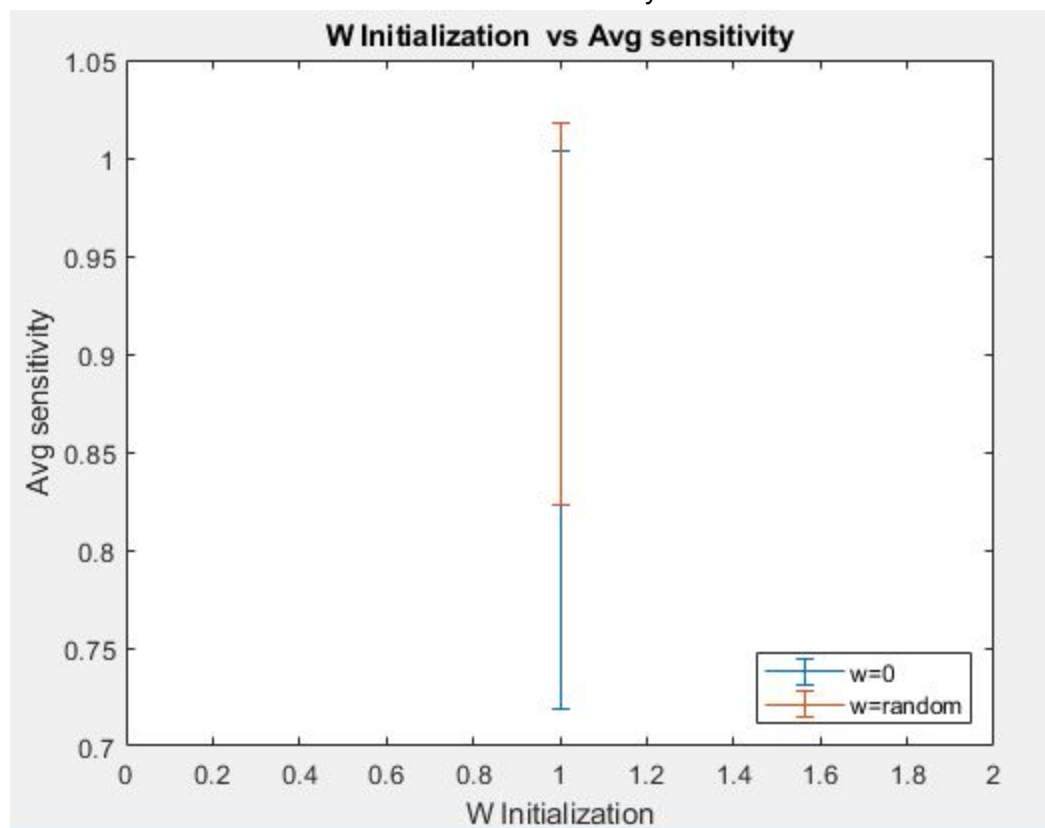
```

Then I calculate the accuracy, sensitivity, and specificity of each fold.

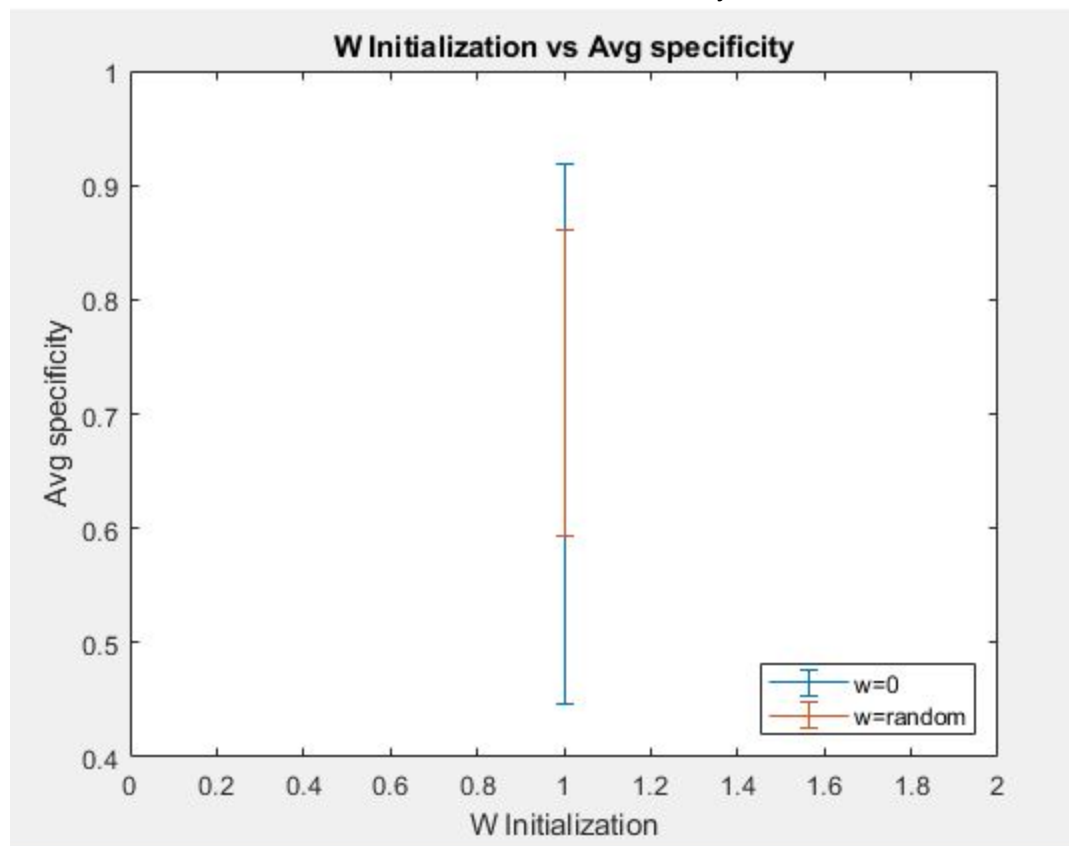
After that I take the average of the performance measures and plot them with respective standard deviations.



The random initialization works better for accuracy.



The random initialization also works better for sensitivity



The random initialization works best for specificity.