CPT_S 570

Machine Learning, Fall 2020

Homework # 2

Tayyab Munir

11716089

## Question 1 – Analytical Part

1. Suppose we have $n_+$ positive training examples and $n-$ negative training examples. Let $C_+$ be the center of the positive examples and $C_-$ be the center of the negative examples, i.e., $C_+ = \frac{1}{n_+} \sum_{i: y_i=+1} x_i$ and $C_- = \frac{1}{n_-} \sum_{i: y_i=-1} x_i$. Consider a simple classifier called CLOSE that classifies a test example x by assigning it to the class whose center is closest.

   • Show that the decision boundary of the CLOSE classifier is a linear hyperplane of the form sign(w· x + b). Compute the values of w and b in terms of $C_+$ and $C_-$.

**Solution:**

We have two points $C_+$ be center of positive examples and $C_-$ be center of negative examples respectively. We need to show that the boundary of the CLOSE classifier is a linear hyperplane of the form sign(w· x + b).

Here w is the weight vector, which is the distance between the 2 centroids $C_+$ and $C_-$ which can be written as $w = |C_+ - C_-|$.

By substituting the value of w in sign(w· x + b) we get

$$(C_+ - C_-) . x + b = 0$$

Consider an example x which has to be classified, let suppose it lies between two clusters such that

$$x = \frac{(C_+ - C_-)}{2}$$

After substituting the value of x we got

$$(C_+ - C_-) . \frac{(C_+ - C_-)}{2} + b = 0$$

$$\frac{1}{2} ((C_+ - C_-) . (C_+ - C_-)) + b = 0$$

$$\frac{1}{2} (C_+{}^2 - C_-{}^2) + b = 0$$

Finally, we can write w and b in terms of $C_+$ and $C_-$

$$w = (C_+ - C_-)$$
$$b = -\frac{1}{2} (C_+{}^2 - C_-{}^2)$$

   • Recall that the weight vector can be written as a linear combination of all the training examples: $w = \sum_{i=1}^{n_+ + n_-} \alpha_i\, y_i\, x_i$. Compute the dual weights (α's). How many of the training examples are support vectors?

**Solution**:

In the above answer we got our weight vector as

$$w = (C_+ - C_-)$$

by substituting values, we get

$$w = (\frac{1}{n_+} \sum_{i:\, y_i=+1} x_i \; - \; \frac{1}{n_-} \sum_{i:\, y_i=-1} x_i)$$

- $$w = (\sum_{i:\, y_i=+1} \frac{1}{n_+} x_i \; - \; \sum_{i:\, y_i=-1} \frac{1}{n_-} x_i)$$

Now, we have the weight vector

$$w = \sum_{i=1}^{n_+ n_-} \propto_i y_i x_i$$

$$w = \sum_{i=1}^{n_+} \propto_i y_i x_i \sum_{j=1}^{n_-} \propto_j y_j x_j$$

For positive values lets insert $y_i = 1$ and for negative values $y_j = -1$

- $$w = \sum_{i=1}^{n_+} \propto_i x_i \sum_{j=1}^{n_-} \propto_j x_j$$

By comparing both the weights we can get $\propto_i$

$$\propto_i = \frac{1}{n_+} \text{ for positive class points}$$

$$\propto_i = \frac{1}{n_-} \text{ for negative class points}$$

2. Suppose we use the following radial basis function (RBF) kernel: $K(x_i, x_j) = \exp\left(-\frac{1}{2}\|x_i - x_j\|^2\right)$, which has some implicit unknown mapping $\varphi(x)$.

- Prove that the mapping $\varphi(x)$ corresponding to RBF kernel has infinite dimensions.

**Solution**:

Radial basis function (RBF) kernel: $K(x_i, x_j) = \exp\left(-\frac{1}{2}\|x_i - x_j\|^2\right)$

By expanding using $(a + b)^2$

$$= \exp\left(\frac{-1}{2}(\|x_i\|^2 + \|x_j\|^2 - 2\|x_i\|\,\|x_j\|)\right)$$

$$= \exp(-(\|x_i\|^2))\,\exp\left(-\left(\|x_j\|^2\right)\right)\exp\left(2\|x_i\|\,\|x_j\|\right)$$

By Taylor Series expansion we got

$$= \exp(-(\|x_i\|^2))\,\exp\left(-\left(\|x_j\|^2\right)\right)\underbrace{\sum_{n=0}^{\infty} \frac{2^n\,\|x_i\|^n\,\|x_j\|^n}{n!}}_{\exp\,(2\|x_i\|\,\|x_j\|)}$$

This equation is the upper bound, hence the RBF has infinite dimensions.

- Prove that for any two input examples xi and xj , the squared Euclidean distance of their corresponding points in the higher-dimensional space defined by $\varphi$ is less than 2, i.e., $\|\varphi(x_i) - \varphi(x_j)\|^2 \leq 2$.

**Solution**:

$$\|\varphi(x_i) - \varphi(x_j)\|^2$$

$$= (\varphi(x_i) - \varphi(x_j)) \cdot (\varphi(x_i) - \varphi(x_j))$$

$$= \varphi(x_i) \cdot \varphi(x_i) + \varphi(x_j) \cdot \varphi(x_j) - 2 \cdot \varphi(x_i) \cdot \varphi(x_j)$$

$$= 2 - 2\exp(-\frac{1}{2}\| x_i - x_j \|^2)\,2$$

$$\|\varphi(x_i) - \varphi(x_j)\|^2 < 2$$

Hence, we can see that the squared Euclidean distance of their corresponding points in the higher dimensional space is less than 2.

3. The decision boundary of a SVM with a kernel function (via implicit feature mapping $\varphi(.)$) is defined as follows: $w \cdot \varphi(x) + b = \sum_{i \in y_i \alpha_i} K(x_i, x) + b = f(x; \alpha, b)$, where w and b are parameters of the decision boundary in the feature space phi defined by the kernel function K, SV is the set of support vectors, and $\alpha i$ is the dual weight of the ith support vector.

Let us assume that we use the radial basis function (RBF) kernel $K(x_i, x_j) = \exp\left(-\frac{1}{2}\|x_i - x_j\|^2\right)$; also assume that the training examples are linearly separable in the feature space $\varphi$ and SVM finds a decision boundary that perfectly separates the training examples. If we choose a testing example $x_{far}$ that is far away from any training instance xi (distance here is measured in the original feature space < d). Prove that $f(x_{far}; \alpha, b) \approx b$.

**Solution**:
We need to prove that $f(x_{far}; \alpha, b) \approx b$
Whereas,

$$\|x_{far} - x_i\| \gg 0, \qquad \forall i \in sv$$

Here SV is the set of support vectors.

$$K(x_{far} - x_i) \approx 0, \qquad \forall i \in sv$$

$$\sum_{i \in SV} y_i \alpha_i K(x_{far}, x_i) \approx 0, \forall i \in sv$$

So, first part of the function is evaluated to a '0' and we are only left with the other part of the function i.e. 'b'.

$$f(x_{far}; \alpha, b) \approx b$$

Hence proved.

4. The function $K(x_i, x_j) = -\langle x_i, x_j \rangle$ is a valid kernel. Prove or Disprove it.

**Solution**:

The function $K(x_i, x_j) = -\langle x_i, x_j \rangle$ is a valid kernel is a valid kernel if and only if it satisfies the mercer's theorem. In order to make the given function satisfy the mercer's theorem we have to satisfy following condition that

$$K(x_i, x_j) = \sum_{i=1}^{n} \sum_{j=1}^{n} K(x_i, x_j) c_i c_j \geq 0$$

If this condition is true and the value is greater than 0 then the given function is a valid kernel.

After substituting the values we get

$$K(x_i, x_j) = -\langle x_i, x_j \rangle$$

$$K(x_i, x_j) = \sum_{i=1}^{n} \sum_{j=1}^{n} -\langle x_i, x_j \rangle c_i c_j$$

$$K(x_i, x_j) = -\sum_{i=1}^{n} \sum_{j=1}^{n} \langle x_i, x_j \rangle c_i c_j$$

Since we get a negative symbol and a negative is always less than 0 we can say that the kernel $K(x_i, x_j) = -\langle x_i, x_j \rangle$ is not a valid kernel.

5. You are provided with n training examples: $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$, where xi is the input example, yi is the class label (+1 or -1). The teacher gave you some additional information by specifying the costs for different mistakes $C_+$ and $C_-$, where $C_+$ and $C_-$ stand for the cost of misclassifying a positive and negative example respectively.

   a. How will you modify the Soft-margin SVM formulation to be able to leverage this extra information? Please justify your answer.

**Solution**:

We have $S = \{(x_i, y_i) \mid i = 1, 2, \dots, l\} \subset R^d \times \{+1, -1\}$ be our training sample and our cost of misclassification is $C_+$ and $C_-$ for making a mistake for positive and negative examples. $C_+$ is for limiting the influence of positive outliers whereas, $C_-$ is for limiting the influence of negative outliers. For learning from a finite sample, Lets divide S into subsets $S \pm 1$ which contain the indices of all positive and negative examples respectively. The optimization problem can now be formulated as

$$Min \ \frac{1}{2}||w||2 \ + \ C \sum_{i \in s+1} C + 1 \ x_i \ \xi_{i\dots k}$$

$$+ \ C \sum_{i \in s+1} C - 1 \ x_i \ \xi_{i\dots k}$$

$$s.t \ \ y'(w.\varphi(x_i) + b) >= 1 - \xi_{i\dots k}. \quad \xi i \geq 0$$

6. You are provided with a set of n training examples: $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$, where xi is the input example, yi is the class label (+1 or -1). Suppose n is very large (say in the order of millions). In this case, standard SVM training algorithms will not scale due to large training set. Tom wants to devise a solution based on "Coarse-to-Fine" framework of problem solving. The basic idea is to cluster the training data; train a SVM classifier based on the clusters (coarse problem); refine the clusters as needed (fine problem); perform training on the finer problem; and repeat until convergence. Suppose we start with $k_+$ positive clusters and $k_-$ negative clusters to begin with (a cluster is defined as a set of examples). Please specify the mathematical formulation (define all the variables used in your formulation) and concrete algorithm for each of the following steps to instantiate this idea:

   a. How to define the SVM training formulation for a given level of coarseness: a set of k+ positive clusters and a set of k− negative clusters?

**Solution**:

Given are k+ positive clusters and k- negative clusters to begin with and we need to convert this to a finer approximation to control the complexity of learning. Let us consider the nearest centroid classification so we assume a centroid from each of the clusters. Which means we have k+ positive centroids and k- negative centroids to reduce the complexity we can consider these points as the training examples for the model. After formulating an SVM optimization problem we have

$$Minimize \ \frac{1}{2} \ ||w||^2 \ + \ C \sum_{i=1}^{k_+ + k_-} \xi_i$$

$$s.t. \ y_i(x_i \ \times \ w \ + b) \geq \ 1 - \ \xi_i$$

   b. How to refine the clusters based on the resulting SVM classifier?

**Solution**:

One way to refine the SVM is to perform training on the clusters that are closer to the decision boundary. The idea is to find the clusters that are close to the decision boundary and break those clusters into smaller clusters and train on those clusters. This will result in a finer classification over the support vectors.

   c. What is the stopping criteria?

**Solution**:
Accuracy score of the SVM can be stopping criteria. We can compare the accuracies for each iteration and if the difference is low then we can stop the training.

   d. Optional question: For what kind of problems will this solution fail?

**Solution**:
This solution can only be failed if there are unnecessary or not related training examples.

7. You are provided with a set of n training examples: $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$, where xi is the input example, yi is the class label (+1 or -1). Suppose n is very large (say in the order of millions). In this case, online kernelized Perceptron algorithms will not scale if the number of allowed support vectors are unbounded.

   a. Suppose you have trained using kernelized Perceptron algorithm (without any bounds on support vectors) and got a set of support vectors SV. Tom wants to use this classifier for real-time prediction and cannot afford more than B kernel evaluations for each classification decision. Please give an algorithm to select B support vectors from SV. You need to motivate your design choices in order to convince Tom to use your solution.

**Solution**:

We are assuming that we have trained using kernelized perceptron algorithm without any bounds which means we have used all the training examples. Now let us suppose that we have X support vectors by training on the given n training examples. Since Tom can only use B kernel evaluations, we need to select B support vectors from X that are useful. By performing a greedy search approach, we can select B support vectors by removing the support vectors from X that contributes less towards the quality of function. Perform the same operation until the size of X is equal to that of B.

   b. Tom wants to train using kernelized Perceptron algorithm but wants to use at most B support vectors during the training process. Please modify the standard kernelized Perceptron training algorithm (from class slides) for this new setting. You need to motivate your design choices in order to convince Tom to use your solution.

**Solution**:

Since Tom wants to use at most B kernels, so for each training example we will draw a circle in such a way that it comprises only those examples that are of the same class. Count the number of examples in each circle and store it as N. If a training example has a large value of N then ignore the example as it will not be close to the decision boundary otherwise, select only the examples which have the smallest values for N which depicts that they are closer to the decision boundary. Run Kernelized Perceptron over the smaller dataset for only B support vectors. We run the algorithm over only those examples that are closer to the decision boundary and perform only B kernel evaluations which reduces time complexity.

# Question 2 – Programming and Empirical Analysis Part

1. Empirical analysis question. You can use a publicly available SVM classifier implementation

(e.g., scikit-learn) for SVM related experiments. scikit-learn (http://scikit-learn.org/stable/modules/svm.html).

You will use the Fashion MNIST data (https://github.com/zalandoresearch/fashion-mnist). There is a fixed training and testing set. From training data, use first 80 percent for \training" and last 20 percent as \validation data".

Each example is a 28x28 grayscale image, associated with a label from 10 classes: 0 Tshirt/top, 1 Trouser, 2 Pullover, 3 Dress, 4 Coat, 5 Sandal, 6 Shirt, 7 Sneaker, 8 Bag, 9 Ankle boot.

You will use ONLY training data for training and testing data for evaluation.

```
6     # Data and Packages Loading
7
8     import pandas as pd
9     import numpy as np
10    import matplotlib.pyplot as plt
11    from sklearn import svm
12
13
14    Data=pd.read_csv('C:/Users/tayya/OneDrive/Desktop/Fall 2020/Machine Learning/Assignment 2/Data/fashion-mnist
15    y_train=np.zeros(len(Data))
16    X_train=np.zeros([len(Data),784])
17    for i in range(len(Data)):
18        y_train[i]=Data.loc[i,'label']
19        X_train[i]=Data.loc[i,'pixel1':'pixel784']
20
21
22    Data_test=pd.read_csv('C:/Users/tayya/OneDrive/Desktop/Fall 2020/Machine Learning/Assignment 2/Data/fashion-
23    y_test=np.zeros(len(Data_test))
24    X_test=np.zeros([len(Data_test),784])
25    for i in range(len(Data_test)):
26        y_test[i]=Data_test.loc[i,'label']
27        X_test[i]=Data_test.loc[i,'pixel1':'pixel784']
28
```

(a) Using a linear kernel, train the SVM on the training data for different values of *C* parameter: 10-4; 10-3; 10-2; 10-1; 100; 101; 102; 103; 104. Compute the training accuracy, validation accuracy, and testing accuracy for the SVM obtained with different values of the C parameter. Plot the training accuracy, validation accuracy, and testing accuracy as a function of *C* (*C* value on x-axis and Accuracy on y-axis) { one curve each for training, validation, and testing data. Also, plot the number of support vectors (if applicable for the SVM toolkit you are using) as a function of *C*. List your observations.
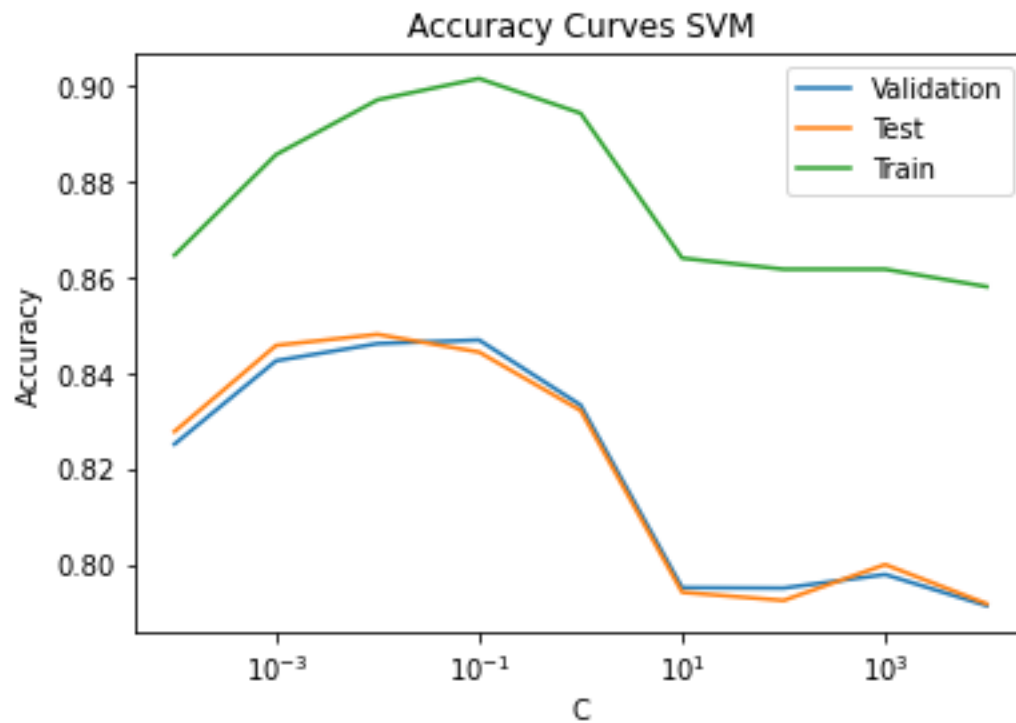
**Solution:**

```
29
30    # Q1 (a)
31
32    # Training for different Values of C
33
34    from sklearn import preprocessing
35    X_train=preprocessing.scale(X_train)
36    X_test=preprocessing.scale(X_test)
37    from sklearn.svm import LinearSVC
38    clf = LinearSVC(C=1e-4,random_state=0, tol=1e-5)
39    clf.fit(X_train[0:48000], y_train[0:48000])
40
41
42    clf1 = LinearSVC(C=1e-3,random_state=0, tol=1e-5)
43    clf1.fit(X_train[0:48000], y_train[0:48000])
44
45
46    clf2 = LinearSVC(C=1e-2,random_state=0, tol=1e-5)
47    clf2.fit(X_train[0:48000], y_train[0:48000])
48
49
50    clf3 = LinearSVC(C=1e-1,random_state=0, tol=1e-5)
51    clf3.fit(X_train[0:48000], y_train[0:48000])
52
53
54    clf4 = LinearSVC(C=1,random_state=0, tol=1e-5)
55    clf4.fit(X_train[0:48000], y_train[0:48000])
56
57
58    clf5 = LinearSVC(C=1e1,random_state=0, tol=1e-5)
59    clf5.fit(X_train[1:48000], y_train[1:48000])
60
61
62    clf6 = LinearSVC(C=1e2,random_state=0, tol=1e-5)
63    clf6.fit(X_train[0:48000], y_train[0:48000])
64
65
66    clf7 = LinearSVC(C=1e3,random_state=0, tol=1e-5)
67    clf7.fit(X_train[0:48000], y_train[0:48000])
68
69
70    clf8 = LinearSVC(C=1e4,random_state=0, tol=1e-5)
71    clf8.fit(X_train[0:48000], y_train[0:48000])
72
```

```python
# Computation of Testing, Validation and Training Accuracies

CL=[clf,clf1,clf2,clf3,clf4,clf5,clf6,clf7,clf8]
C_mat=np.array([1e-4,1e-3,1e-2,1e-1,1,1e1,1e2,1e3,1e4])
Test_Accuracy_test=np.zeros(len(C_mat))
Test_Accuracy_valid=np.zeros(len(C_mat))
Test_Accuracy_Train=np.zeros(len(C_mat))
X_valid=X_train[48000:60000]
y_valid=y_train[48000:60000]
X_T=X_train[0:48000]
y_T=y_train[0:48000]
for i in range(C_mat.shape[0]):
    Mistake=0
    for k in range(X_test.shape[0]):
        yp=CL[i].predict([X_test[k]])
        if(yp!=y_test[k]):
            Mistake=Mistake+1
    Test_Accuracy_test[i]=(len(X_test)-Mistake)/len(X_test)
    Mistake=0
    for k in range(X_T.shape[0]):
        yp=CL[i].predict([X_train[k]])
        if(yp!=y_T[k]):
            Mistake=Mistake+1
    Test_Accuracy_Train[i]=(len(X_train)-Mistake)/len(X_train)
    Mistake=0
    for k in range(X_valid.shape[0]):
        yp=CL[i].predict([X_valid[k]])
        if(yp!=y_valid[k]):
            Mistake=Mistake+1
    Test_Accuracy_valid[i]=(len(X_valid)-Mistake)/len(X_valid)


plt.plot(C_mat,Test_Accuracy_valid,label="Validation")
plt.plot(C_mat,Test_Accuracy_test,label="Test")
plt.plot(C_mat,Test_Accuracy_Train,label="Train")
plt.xscale('log',basex=10)
plt.xlabel('C')
plt.ylabel('Accuracy')
plt.title('Accuracy Curves SVM')
plt.legend()
plt.show()
```

Accuracy Curves SVM

As can be seen from graph above, the best value of C based on Validation accuracy is C=1e-1. Also, it is observed that validation accuracy is good approximation of test accuracy or the true error.

(b) Select the best value of hyper-parameter $C$ based on the accuracy on validation set and train a linear SVM on the *combined* set of training and validation examples. Compute the testing accuracy and the corresponding confusion matrix: a $10 \times 10$ matrix.

**Solution:**

```
115
116    # Q1 (b)
117
118    clfbest = LinearSVC(C=1e-1,random_state=0, tol=1e-5)
119    clfbest.fit(X_train, y_train)
120
121    X_train.shape
122
123
124    from sklearn.metrics import confusion_matrix
125    Mistake=0;
126    yp=np.zeros(X_test.shape[0])
127    for k in range(X_test.shape[0]):
128        yp[k]=clfbest.predict([X_test[k]])
129        if(yp[k]!=y_test[k]):
130            Mistake=Mistake+1
131    Test_Accuracy_best=(len(X_test)-Mistake)/len(X_test)
132    print(Test_Accuracy_best)
133    ## CONFUSION MATRIX
134    Conf_Mat=confusion_matrix(y_test, yp)
135    import pandas as pd
136    Conf_Table=(pd.DataFrame(Conf_Mat,columns=["Tshirt/Top", "Trouser","Pullover","Dress","Coat","Sandal","Shirt","Sneaker","Bag","Ankle-Boot"],
137                             index=["Tshirt/Top", "Trouser","Pullover","Dress","Coat","Sandal","Shirt","Sneaker","Bag","Ankle-Boot"]))
138    print(Conf_Table)
139
```

In [96]: X_train.shape
Out [96]: (60000, 784)

Conf_Table - DataFrame

| Index | Tshirt/Top | Trouser | Pullover | Dress | Coat | Sandal | Shirt | Sneaker | Bag | Ankle-Boot |
|---|---|---|---|---|---|---|---|---|---|---|
| Tshirt/Top | 815 | 4 | 21 | 46 | 1 | 6 | 90 | 0 | 17 | 0 |
| Trouser | 5 | 968 | 4 | 15 | 1 | 1 | 5 | 1 | 0 | 0 |
| Pullover | 19 | 7 | 765 | 9 | 118 | 0 | 72 | 1 | 9 | 0 |
| Dress | 30 | 19 | 15 | 877 | 27 | 2 | 25 | 0 | 5 | 0 |
| Coat | 1 | 2 | 73 | 30 | 804 | 1 | 85 | 2 | 2 | 0 |
| Sandal | 1 | 3 | 0 | 0 | 0 | 907 | 0 | 57 | 11 | 21 |
| Shirt | 169 | 7 | 111 | 43 | 87 | 1 | 558 | 4 | 20 | 0 |
| Sneaker | 0 | 0 | 0 | 0 | 0 | 40 | 0 | 912 | 1 | 47 |
| Bag | 7 | 2 | 9 | 13 | 2 | 18 | 18 | 5 | 925 | 1 |
| Ankle-Boot | 0 | 0 | 1 | 0 | 1 | 18 | 0 | 38 | 2 | 940 |

(c) Repeat the experiment (a) with the best $C$ value from (a) with polynomial kernel of degree 2, 3, and 4. Compare the training, validation, testing accuracies, and the number of support vectors for different kernels (linear, polynomial kernel of degree 2, polynomial kernel of degree 3, and polynomial kernel of degree 4). List your observations.

**Solution:**

```python
140
141    # Q1 (c)
142
143    # Training for different degrees corresponding to best C value (1e-1).
144
145    clfpol1 = svm.SVC(C=1e-1,kernel='poly',degree=2)
146    clfpol1.fit(X_train[0:48000], y_train[0:48000])
147
148
149    clfpol2 = svm.SVC(C=1e-1,kernel='poly',degree=3)
150    clfpol2.fit(X_train[0:48000], y_train[0:48000])
151
152
153    clfpol3 = svm.SVC(C=1e-1,kernel='poly',degree=4)
154    clfpol3.fit(X_train[0:48000], y_train[0:48000])
155
156
157    CL=[clfpol1,clfpol2,clfpol3]
158    Degrees=np.array([2,3,4])
159    Test_Accuracy_test_poly=np.zeros(len(Degrees))
160    Test_Accuracy_valid_poly=np.zeros(len(Degrees))
161    Test_Accuracy_Train_poly=np.zeros(len(Degrees))
162    X_valid=X_train[48000:60000]
163    y_valid=y_train[48000:60000]
164    X_T=X_train[0:48000]
165    y_T=y_train[0:48000]
166    for i in range(Degrees.shape[0]):
167        Mistake=0
168        for k in range(X_test.shape[0]):
169            yp=CL[i].predict([X_test[k]])
170            print(k)
171            if(yp!=y_test[k]):
172                Mistake=Mistake+1
173        Test_Accuracy_test_poly[i]=(len(X_test)-Mistake)/len(X_test)
174        print("Done1")
175        Mistake=0
176        for k in range(X_valid.shape[0]):
177            print(k)
178            yp=CL[i].predict([X_valid[k]])
179            if(yp!=y_valid[k]):
180                Mistake=Mistake+1
181        Test_Accuracy_valid_poly[i]=(len(X_valid)-Mistake)/len(X_valid)
182        print('Done2')
183        Mistake=0
184        for k in range(X_T.shape[0]):
185            print(k)
186            yp=CL[i].predict([X_train[k]])
187            if(yp!=y_T[k]):
188                Mistake=Mistake+1
189        Test_Accuracy_Train_poly[i]=(len(X_train)-Mistake)/len(X_train)
190        print('Done3')
191
```

```
192
193    print(Test_Accuracy_valid_poly)
194    print(Test_Accuracy_Train_poly)
195    print(Test_Accuracy_test_poly)
196    Degre=np.array([2,3,4])
197    plt.plot(Degre,Test_Accuracy_valid_poly,label="Validation")
198    plt.plot(Degre,Test_Accuracy_test_poly,label="Test")
199    plt.plot(Degre,Test_Accuracy_Train_poly,label="Train")
200    plt.xlabel('Degree')
201    plt.ylabel('Accuracy')
202    plt.title('Accuracy Curves SVM- Polynomial Kernel')
203    plt.legend()
204    plt.show()
205
206
207    NSV=np.zeros(3)
208    NSV[0]=len(CL[0].support_vectors_)
209    NSV[1]=len(CL[1].support_vectors_)
210    NSV[2]=len(CL[2].support_vectors_)
211    plt.plot(Degre,NSV)
212    plt.xlabel('Degree')
213    plt.ylabel('# of SVs')
214    plt.title('Number of Support Vectors- Polynomial Kernel')
215    plt.show()
216
217
```
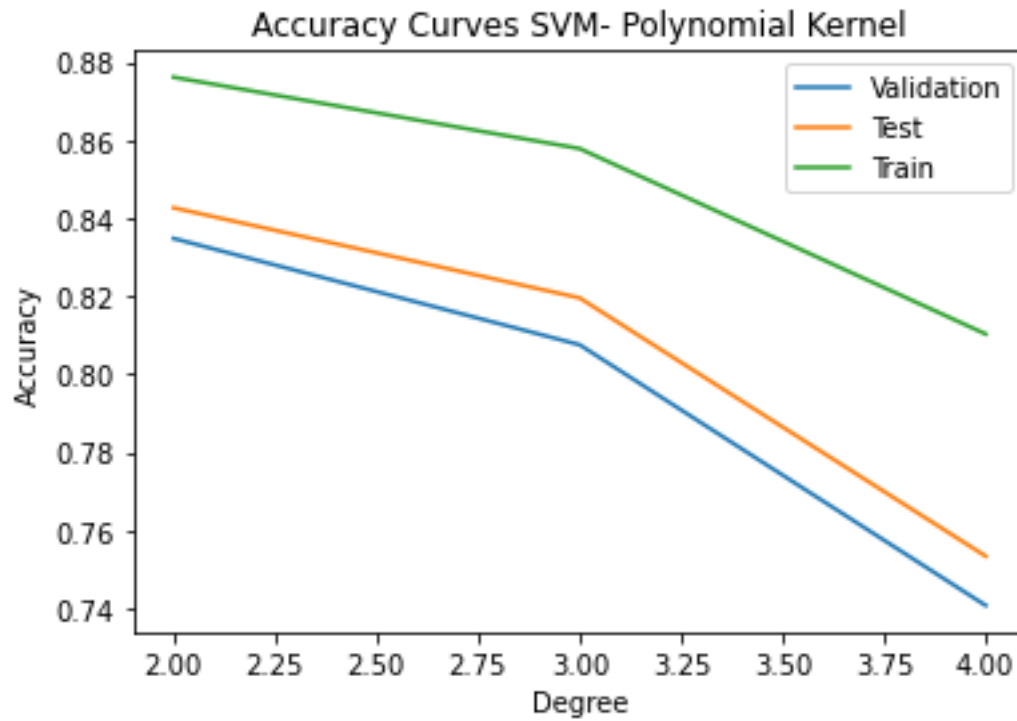
Training for different degrees corresponding to best C value (1e-1).
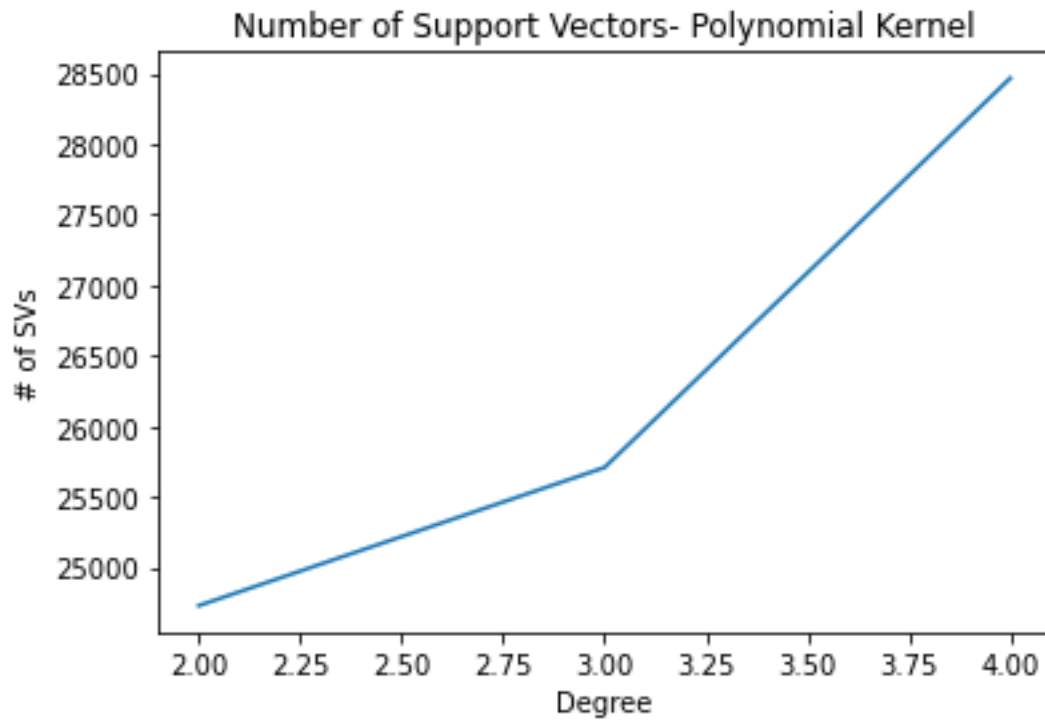
[0.83475    0.8075      0.74066667]

[0.87608333    0.85778333      0.81025    ]

[0.8426    0.8195    0.7533]

Accuracy Curves SVM- Polynomial Kernel

It is observed from the above graph that best degree based on validation accuracy is 2.00.



Number of Support Vectors- Polynomial Kernel

Degree 2 corresponds to least number of support vectors.

2. Programming question. You will implement the kernelized Perceptron training algorithm (discussed in the class) for *multi-class* classification.

(a) You will use the Fashion MNIST data. Train the kernelized Perceptron classifier for 5 iterations with polynomial kernel (pick the best degree out of 2, 3, and 4 from the above experiment). Plot the number of mistakes as a function of training iterations. Compute the training, validation, and testing accuracy at the end of 5 iterations.

```python
1    # # Q 2  [Kernelized Perceptron]
2
3    import pandas as pd
4    import numpy as np
5    import matplotlib.pyplot as plt
6    import time
7
8
9    Data=pd.read_csv('C:/Users/tayya/OneDrive/Desktop/Fall 2020/Machine Learning/Assignment 2/Data/fashion-mnist_train.csv')
10   y_train=np.zeros(len(Data))
11   X_train=np.zeros([len(Data),784])
12   for i in range(len(Data)):
13       y_train[i]=Data.loc[i,'label']
14       X_train[i]=Data.loc[i,'pixel1':'pixel784']
15
16
17   Data_test=pd.read_csv('C:/Users/tayya/OneDrive/Desktop/Fall 2020/Machine Learning/Assignment 2/Data/fashion-mnist_test.csv')
18   y_test=np.zeros(len(Data_test))
19   X_test=np.zeros([len(Data_test),784])
20   for i in range(len(Data_test)):
21       y_test[i]=Data_test.loc[i,'label']
22       X_test[i]=Data_test.loc[i,'pixel1':'pixel784']
23
```

This function below takes whole X_train as input and computes Kernel of test feature vector with each training example feature. Output is a row vector of length 1xN where N is the number of training examples.

```python
24
25   def poly_Kernel(X_train,xj,p):
26       K=np.zeros(len(X_train))
27       K=(1+np.matmul(X_train,np.reshape(xj,(784,1))))**p
28       return K
29
30
```

This function below takes Alpha matrix, X_train, number of classes (k) and degree of Polynomial kernel as input and predicts the output for a test feature.

```python
31   def Predict_Y(Alpha_M,x,X_train,k,p):
32       pred_vect=np.zeros(k)
33       pred_vect=np.matmul(Alpha_M,poly_Kernel(X_train,x,p))
34       pred=np.argmax(pred_vect)
35       return pred
36
```

This function below trains the kernelized perceptron for a polynomial kernel of degree "p". User also specifies the number of classes (k) and number of training passes (T)

```python
def Train_Kernelized_Perceptron(X_train,y_train,k,T,p):
    Alpha_M=np.zeros([k,len(X_train)])
    Errors=np.array([])
    for i in range(T):
        Count=0
        for j in range(len(X_train)):
            t0 =time.time()
            yhat=Predict_Y(Alpha_M,X_train[j],X_train,k,p)
            t1= time.time()
            if(yhat!=y_train[j]):
                Alpha_M[int(yhat),j]=Alpha_M[int(yhat),j]-1
                Alpha_M[int(y_train[j]),j]=Alpha_M[int(y_train[j]),j]+1
                Count=Count+1
                if(j%1000==0):
                    print(t1-t0,j)
                #print(j)
        Errors=np.append(Errors,Count)
        print("Pass Comp")
    return Alpha_M, Errors
```

```python
# Training


from sklearn import preprocessing
X_train=preprocessing.scale(X_train)
X_test=preprocessing.scale(X_test)


X_train_set=X_train[0:48000]
y_train_set=y_train[0:48000]
X_train_val=X_train[48000:60000]
y_train_val=y_train[48000:60000]

```

```
72
73
74      [Alpha_M,Errors]=Train_Kernelized_Perceptron(X_train_set,y_train_set,10,5,2)
75
76
77      Training_IT=[1,2,3,4,5]
78      plt.plot(Training_IT,(48000-Errors)/48000)
79      plt.xlabel('Training Iterations')
80      plt.ylabel('Accuracy')
81      plt.title('Online Learning Curve-Kernelized Perceptron')
82      plt.show()
83
84
85      Training_IT=[1,2,3,4,5]
86      plt.plot(Training_IT,Errors)
87      plt.xlabel('Training Iterations')
88      plt.ylabel('Mistakes')
89      plt.title('Online Learning Curve-Kernelized Perceptron')
90      plt.show()
91
```

[Alpha_M,Errors]=Train_Kernelized_Perceptron(X_train_set,y_train_set,10,5,2)

Output:

0.0625145435333252 0

0.015633344650268555 4000

0.01561427116394043 9000

0.01563096046447754 12000

0.015628576278686523 17000

0.030728578567504883 22000

0.015624046325683594 25000

0.01562666893005371 26000

0.0159912109375 47000

Pass Comp

0.015986919403076172 5000

0.015638351440429688 9000

0.014977693557739258 12000

0.013988733291625977 17000

0.0062906742095947266 22000

0.014012813568115234 25000

0.01497650146484375 42000

0.015624761581420898 47000

Pass Comp

0.01564311981201172 0

0.014961719512939453 8000

0.015614986419677734 17000

0.01599264144897461 25000

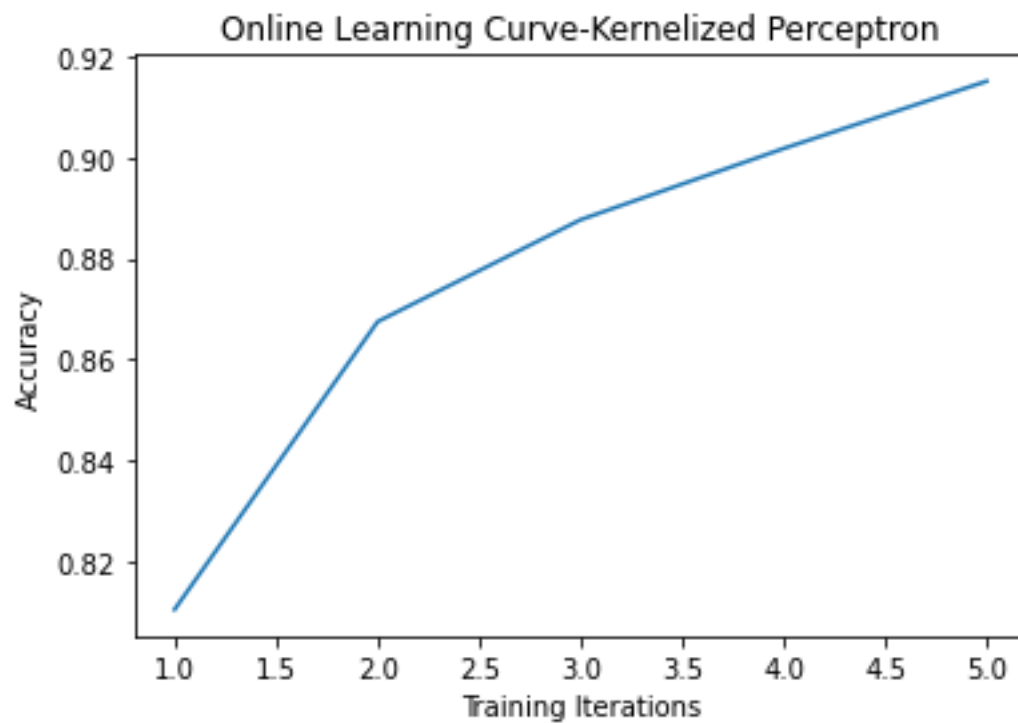0.01499319076538086 47000

Pass Comp

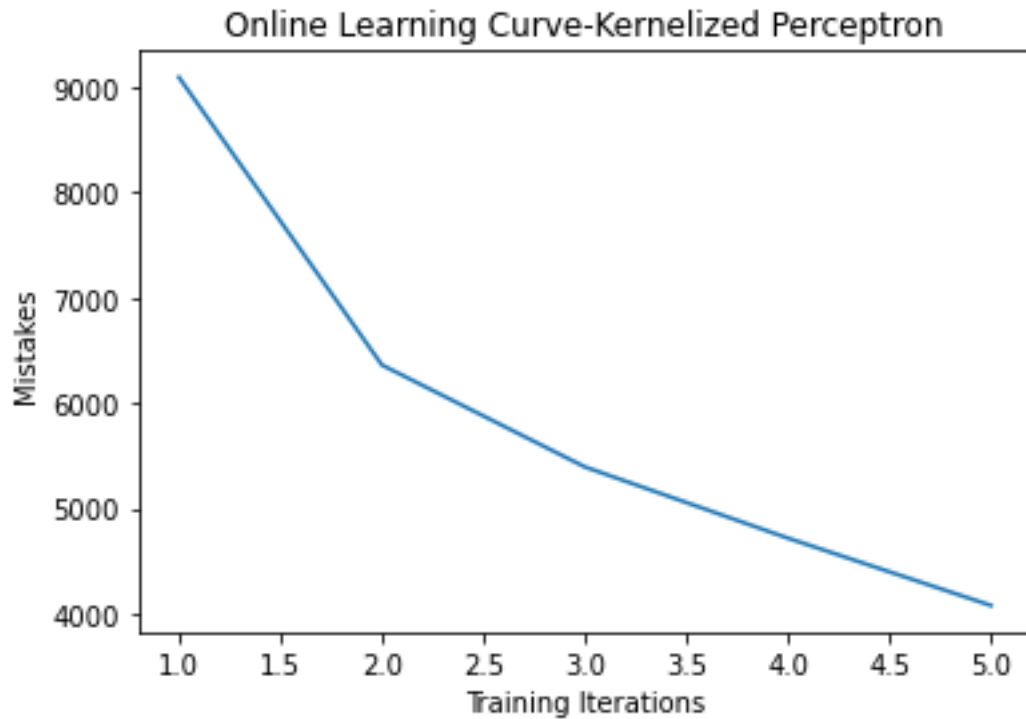0.015989303588867188 17000

0.014985084533691406 25000

Pass Comp

0.01599287986755371 12000

0.015989303588867188 25000

Pass Comp

Online Learning Curve-Kernelized Perceptron

Comparison of the Training, Test, and validation Accuracy

```
94
95
96      Mistake=0
97      for i in range(len(X_train_val)):
98          yhat=Predict_Y(Alpha_M,X_train_val[i],X_train_set,10,2)
99          print(i)
100         if(yhat!=y_train_val[i]):
101             Mistake=Mistake+1
102
103
104     Mistake_test=0
105     for i in range(len(X_test)):
106         yhat=Predict_Y(Alpha_M,X_test[i],X_train_set,10,2)
107         print(i)
108         if(yhat!=y_test[i]):
109             Mistake_test=Mistake_test+1
110
111
112     print("Training Accuracy    :",(48000-Errors[4])/48000)
113     print("Validation Accuracy :",(12000-Mistake)/12000)
114     print("Test Accuracy        :",(10000-Mistake_test)/10000)
115
```

Training Accuracy: 0.9150208333333333

Validation Accuracy: 0.8583333333333333

Test Accuracy: 0.8639

3. Programming question. \Breast Cancer" Classifier using Decision Trees. You will use the following dataset for this question:
https://archive.ics.uci.edu/ml/datasets/Breast+Cancer+Wisconsin+%28Diagnostic%29. You will use the first 70 percent examples for training, next 10 percent examples for validation, and last 20 percent examples for testing.

```
3
4    import pandas as pd
5    import numpy as np
6
7    missing_values=['?']
8    Data=pd.read_csv('C:/Users/tayya/OneDrive/Desktop/Fall 2020/Machine Learning/Assignment 2
9    Data=Data.replace(np.nan,0)
10   y_train=np.zeros(len(Data))
11   X_train=np.zeros([len(Data),9])
12   for i in range(len(Data)):
13       if(Data.loc[i,10]==2):
14           y_train[i]=1
15       else:
16           y_train[i]=0
17       X_train[i]=Data.loc[i,1:9]
18
19
20   import math
21   def Entropy(y_train):
22       Count1=0
23       Pr1=0;
24       Pr2=0;
25       for i in range(len(y_train)):
26           if y_train[i]==1:
27               Count1=Count1+1
28       Count0=len(y_train)-Count1
29       if(len(y_train)!=0):
30           Pr1=(Count1)/len(y_train)
31           Pr2=1-Pr1;
32       else:
33           H=0
34       if((Pr1!=0)&(Pr2!=0)):
35           H=-Pr1*math.log(Pr1)-Pr2*math.log(Pr2)
36       else:
37           H=0
38       return H
39
```

(a) Implement the ID3 decision tree learning algorithm that we discussed in the class. The key step in the decision tree learning is choosing the next feature to split on. Implement the information gain heuristic for selecting the next feature. Please see lecture notes or https://en.wikipedia.org/wiki/ID3_algorithm for more details. Do the following to select candidate thresholds for continuous features: Sort all candidate values for feature $f$ from training data. Suppose $f1; f2; \cdots ; fn$ is the sorted list. The candidate thresholds are chosen as $fi + (fi+1 - fi)=2$ for $i=1$ to $n$.

**Solution:**

> • Here the function "Split" choses the best feature and corresponding threshold based on Maximum Information gain. Output of this function are "Feat_Index (Best Feature)" and corresponding best threshold "Threshold".
> • Function split divides the data into two parts by taking data, threshold and index of the feature as input

```
40
41    # Part a.
42
43    def subset(X_train,y_train,TH,i):
44        X_T1=[]
45        y_T1=np.array([])
46        X_T2=[]
47        y_T2=np.array([])
48        for k in range(len(X_train)):
49            if(X_train[k,i]>=TH):
50                X_T1.append(X_train[k,:])
51                y_T1=np.append(y_T1,y_train[k])
52            if(X_train[k,i]<TH):
53                X_T2.append(X_train[k])#np.append(X_T2,X_train[k,i])
54                y_T2=np.append(y_T2,y_train[k])
55        return X_T1,y_T1,X_T2,y_T2
56
57
58
59    def Split(X_train,y_train):
60        H=Entropy(y_train)
61        H_max=np.zeros(X_train.shape[1])
62        T_max=np.zeros(X_train.shape[1])
63        for i in range(X_train.shape[1]):
64            X_T=np.sort(X_train[:,i])
65            H1=np.zeros(len(X_T)-1)
66            T1=np.zeros(len(X_T)-1)
67            for j in range(len(X_T)-1):
68                TH=X_T[j]+(X_T[j+1]-X_T[j])/2
69                [X_T1,y_T1,X_T2,y_T2]=subset(X_train,y_train,TH,i)
70                #print(len(y_T1),len(y_T2))
71                #if((len(y_T1)!=0)|(len(y_T2)!=0)):
72                H1[j]=Entropy(y_T1)*len(y_T1)/(len(y_T1)+len(y_T2))+Entropy(y_T2)*len(y_T2)/(len(y_T1)+len(y_T2))
73                T1[j]=TH
74            H_max[i]=max(H-H1)
75            T_max[i]=T1[np.argmax(H-H1)]
76        Feat_index=np.argmax(H_max)
77        Threshold=T_max[Feat_index]
78        return Feat_index,Threshold
79
```

(b) Run the decision tree construction algorithm on the training examples. Compute the accuracy on validation examples and testing examples.

**Solution:**

"Train-ID3" trains the decision tree by recursive approach, the format of the decision tree node is "Feat_Index, Threshold, Child_1, Child_2". Child nodes could be impure or terminal nodes. Each Child also contains same information ("Feat_Index, Threshold, Child_1, Child_2") if it's not a leaf node. If Child node is a leaf node it contains "[0]" or "[1]".

```
80
81      # Part b.
82
83
84      def Train_ID3(X_train,y_train):
85          Child_1=[]
86          Child_2=[]
87          [Feat_index,Threshold]=Split(X_train,y_train)
88          [X_T1,y_T1,X_T2,y_T2]=subset(X_train,y_train,Threshold,Feat_index)
89          if((Entropy(y_T1)==0)&(len(y_T1)>0)):
90              Child_1=[int(np.mean(y_T1))]
91          elif(len(y_T1)==0):
92              Child_1=[]
93          if((Entropy(y_T2)==0)&(len(y_T2)>0)):
94              Child_2=[int(np.mean(y_T2))]
95          elif(len(y_T2)==0):
96              Child_2=[]
97          if((Entropy(y_T1))!=0):
98              subtree1=Train_ID3(np.array(X_T1),y_T1)
99              Child_1.extend(subtree1)
100         if((Entropy(y_T2))!=0):
101             subtree2=Train_ID3(np.array(X_T2),y_T2)
102             Child_2.extend(subtree2)
103         print('Training Decision Tree ...')
104         return Feat_index,Threshold,Child_1,Child_2
105
106
107     Tree=Train_ID3(X_train[0:490],y_train[0:490])
108
```

In [7]: Tree=Train_ID3(X_train[0:490],y_train[0:490])

Training Decision Tree ...

Training Decision Tree ...

Training Decision Tree ...

Training Decision Tree ...

Training Decision Tree ...

Training Decision Tree ...

Training Decision Tree ...

Training Decision Tree ...

Training Decision Tree ...

Training Decision Tree ...

Training Decision Tree ...

Training Decision Tree ...

Training Decision Tree ...

Training Decision Tree ...

Training Decision Tree ...

Training Decision Tree ...

Training Decision Tree ...

Training Decision Tree ...

Training Decision Tree ...

Training Decision Tree ...

Training Decision Tree ...

Training Decision Tree ...

Training Decision Tree ...

Training Decision Tree ...

Training Decision Tree ...

Training Decision Tree ...

Tree

Out[8]:

(2,

2.5,

[5,

2.5,

[0,

6.5,

[2, 4.5, [0], [7, 7.0, [1, 5.0, [0], [1]], [0]]],

[3,

5.5,

[0],

[7,

7.5,

[0],

[7,

6.5,

[1],

[1,

3.5,

[5, 6.0, [6, 5.5, [0], [1, 6.5, [1], [2, 4.5, [0], [1]]]], [1]],

[0, 5.5, [1], [0]]]]]]],

[2,

3.5,

[0,

6.5,

[1, 4.5, [0], [2, 5.5, [0, 7.5, [1], [0]], [0]]],

[4, 4.5, [5, 0.5, [0], [1]], [1]]],

[1]]],

[0, 5.5, [5, 1.5, [0], [1]], [5, 4.5, [0, 3.5, [0], [1]], [1]]])

```
112
113    def pred(X_test,C):
114        F=C[0]
115        T=C[1]
116        if(X_test[F]>=T):
117            if ((C[2]==[1])|(C[2]==[0])):
118                y_pred=C[2]
119                C=C[2]
120            else:
121                C=C[2]
122                y_pred=pred(X_test,C)
123        else:
124            if ((C[3]==[1])|(C[3]==[0])):
125                y_pred=C[3]
126                C=C[3]
127            else:
128                C=C[3]
129                y_pred=pred(X_test,C)
130        return y_pred
131
132
```

"Pred" function predicts the output "ypred" for a given value of feature vector.

```
132
133
134     def Testing(X_test,y_test,Tree):
135         Mistakes=0
136         for i in range(len(X_test)):
137             ypred=pred(X_test[i],Tree)
138             if(ypred[0]!=y_test[i]):
139                 Mistakes=Mistakes+1
140         Accuracy=(len(X_test)-Mistakes)/len(X_test)
141         return Accuracy
142
143
144     X_trains=X_train[0:490]
145     y_trains=y_train[0:490]
146     Accuracy=Testing(X_trains,y_trains,Tree)
147     print('Training Accuracy:',Accuracy)
148
149     X_trainv=X_train[490:560]
150     y_trainv=y_train[490:560]
151     Accuracy=Testing(X_trainv,y_trainv,Tree)
152     print('Validation Accuracy:',Accuracy)
153
154     X_test=X_train[560:699]
155     y_test=y_train[560:699]
156     Accuracy=Testing(X_test,y_test,Tree)
157     print('Test Accuracy:',Accuracy)
158
```

"Testing" function checks the testing accuracy using decision tree.

70% of Data is used for training, 10% for Validation and 20% for testing.

Training Accuracy: 1.0

Validation Accuracy: 0.9571428571428572

Test Accuracy: 0.9424460431654677

(c) Implement the decision tree pruning algorithm discussed in the class (via validation data).

**Solution:**

In this part pruning algorithm is designed for decision tree.

> • "Majority" fuction finds the majority classes of the children in a decision tree node
> • "Prune_Tree" prunes the decision tree nodes if validation accuracy increases by placing a majority class at an impure node. Algorithm runs recursively until it has pruned all the nodes if validation accuracy improves

```
159
160    # Part c.
161
162    def Majority(Tree,X_train,y_train):
163        [X_T1,y_T1,X_T2,y_T2]=subset(X_train,y_train,Tree[1],Tree[0])
164        One=0
165        Zero=0
166        for i in range(len(y_T1)):
167            if(y_T1[i]==0):
168                Zero=Zero+1
169            else:
170                One=One+1
171        if(Zero>=One):
172            Mp1=0
173        else:
174            Mp1=1
175        One=0
176        Zero=0
177        for i in range(len(y_T2)):
178            if(y_T2[i]==0):
179                Zero=Zero+1
180            else:
181                One=One+1
182        if(Zero>=One):
183            Mp2=0
184        else:
185            Mp2=1
186        return [Mp1],[Mp2]
187
188
189
190    def New_Tree(Tree,Mp1,Mp2):
191        Tr_P=[Tree[0],Tree[1],Mp1,Mp2]
192        return Tr_P
193
194
```

```python
def Prune_Tree(Tree,X_val,y_val):
    T=[]
    [Mp1,Mp2]=Majority(Tree,X_train,y_train)
    Tr_P1=New_Tree(Tree,Mp1,Tree[3])
    Accuracy_p=Testing(X_val,y_val,Tr_P1)
    Accuracy=Testing(X_val,y_val,Tree)
    if(Accuracy_p>Accuracy):
        T.extend(Tr_P1)
    else:
        if((Tree[2]!=[0])&(Tree[2]!=[1])):
            Tr_P1=Prune_Tree(Tree[2],X_val,y_val)
    Tr_P1=New_Tree(Tree,Tree[2],Mp2)
    Accuracy_p=Testing(X_val,y_val,Tr_P1)
    Accuracy=Testing(X_val,y_val,Tree)
    if(Accuracy_p>Accuracy):
        T.extend(Tr_P1)
    else:
        if((Tree[3]!=[0])&(Tree[3]!=[1])):
            Tr_P1=Prune_Tree(Tree[3],X_val,y_val)
    return T


Tp=Prune_Tree(Tree,X_trainv,y_trainv)
```

(d) Run the pruning algorithm on the decision tree constructed using training examples. Compute the accuracy on validation examples and testing examples. List your observations by comparing the performance of decision tree with and without pruning.

To debug and test your implementation, you can employ scikit-learn (http://scikit-learn.org/stable/modules/tree.html).

**Solution:**

• The validation accuracy improves by pruning but testing accuracy decreases slightly. We can try with a slightly larger validation set to improve testing accuracy.

```
221
222     # Part d.
223
224     X_trains=X_train[0:490]
225     y_trains=y_train[0:490]
226     Accuracy=Testing(X_trains,y_trains,Tp)
227     print('Training Accuracy:',Accuracy)
228
229     X_trainv=X_train[490:560]
230     y_trainv=y_train[490:560]
231     Accuracy=Testing(X_trainv,y_trainv,Tp)
232     print('Validation Accuracy:',Accuracy)
233
234     X_test=X_train[560:699]
235     y_test=y_train[560:699]
236     Accuracy=Testing(X_test,y_test,Tp)
237     print('Test Accuracy:',Accuracy)
238
239
```

Training Accuracy: 0.9836734693877551

Validation Accuracy: 0.9857142857142858

Test Accuracy: 0.935251798561151