

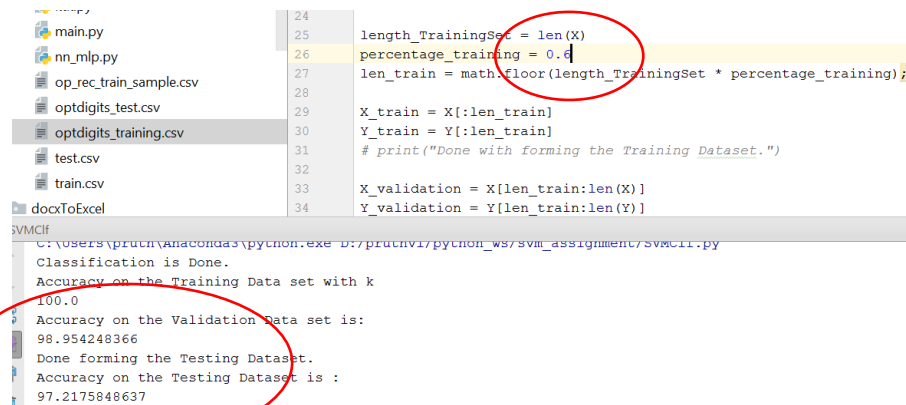
(Worked together with Chaitanya Sri Krishna Lolla)

Digit Recognition dataset

- RESULTS:

Splitting the given data into 2 parts to perform cross validation (as a step for pruning the data).

Training on anything below 70% for cross validation which is resulting in a 98% accuracy on training set and 97.21% accuracy on the validation data.



```

24
25 length_TrainingSet = len(X)
26 percentage_training = 0.7
27 len_train = math.floor(length_TrainingSet * percentage_training);
28
29 X_train = X[:len_train]
30 Y_train = Y[:len_train]
31 # print("Done with forming the Training Dataset.")
32
33 X_validation = X[len_train:len(X)]
34 Y_validation = Y[len_train:len(Y)]

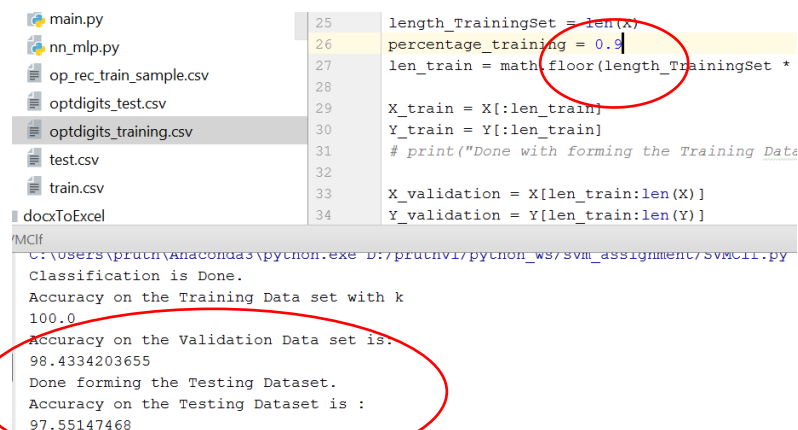
```

```

C:\Users\pruthi\anaconda3\python.exe D:/pruthi/python_ws/svm_assignment/svmcl1.py
Classification is Done.
Accuracy on the Training Data set with k
100.0
Accuracy on the Validation Data set is:
97.2175848637
Done forming the Testing Dataset.
Accuracy on the Testing Dataset is :
97.2175848637

```

But on cross validating it with a 90-10 split of training data it is resulting in a better fit as it can be observed in the accuracy on testing data:



```

25 length_TrainingSet = len(X)
26 percentage_training = 0.9
27 len_train = math.floor(length_TrainingSet *
28
29 X_train = X[:len_train]
30 Y_train = Y[:len_train]
31 # print("Done with forming the Training Data
32
33 X_validation = X[len_train:len(X)]
34 Y_validation = Y[len_train:len(Y)]

```

```

C:\Users\pruthi\anaconda3\python.exe D:/pruthi/python_ws/svm_assignment/svmcl1.py
Classification is Done.
Accuracy on the Training Data set with k
100.0
Accuracy on the Validation Data set is:
97.55147468
Done forming the Testing Dataset.
Accuracy on the Testing Dataset is :
97.55147468

```

Above accuracy is a resultant of the classifier with all the default values provided by the sklearn classifier, With kernel="poly".

With default SVM parameters given by the classifier this model it gives a poor classification.

```

optdigits_training.csv 40 clf = SVC()
test.csv               41 clf = clf.fit(X_train, Y_train)
train.csv              42 print("Classification is Done.")
docxToExcel            43
                        44 output_Predicted = clf.predict(X_train);
/Clf
C:\Users\pruthi\Anaconda3\python.exe D:/pruthiVI/python_ws/svm_assignment/svmClf.py
Classification is Done.
Accuracy on the Training Data set with k
100.0
Accuracy on the Validation Data set is:
68.1462140992
Done forming the Testing Dataset.
Accuracy on the Testing Dataset is :
57.5403450195

```

Amazon review dataset

- Result:

Using nltk library services to lemmatize, tokenize the reviews along with removing the punctuation from the review texts as a part of preprocessing steps as in done in the previous assignment.

Pre-processing the text and making it into two sets of positive and negative reviews.

Picking top 5000 of the most commonly occurring words and assigning them numerical values to pass as parameters for the decision tree classifier method.

Used TF-IDF to build a matrix of numerical representation of the words in each sentence.

And the same procedure is repeated on the test dataset and the accuracy is measured.

- Used a normal SVM () method and checked for the accuracy of model :

```

(0, 2158)      1
(0, 2896)      1
(0, 2764)      1
(0, 1749)      1
(0, 4414)      1
(0, 2264)      1
(0, 4463)      1
(0, 254)       1
(0, 2271)      1
(0, 4385)      1
Accuracy :  0.764279399974
27948
[('the', 152974), ('it', 99372), ('i', 924
Accuracy of test :  76.2240447651 %

```

- Now I 've added the kernel="poly" as I've observed that this is giving a same fit:

```
(0, 2765) 1
(0, 1748) 1
(0, 4415) 1
(0, 2265) 1
(0, 4464) 1
(0, 254) 1
(0, 2272) 1
(0, 4386) 1
accuracy : 0.764279399974
.7948
('the', 152974), ('it', 99372), ('i'
accuracy of test : 76.2240447651 %
```