

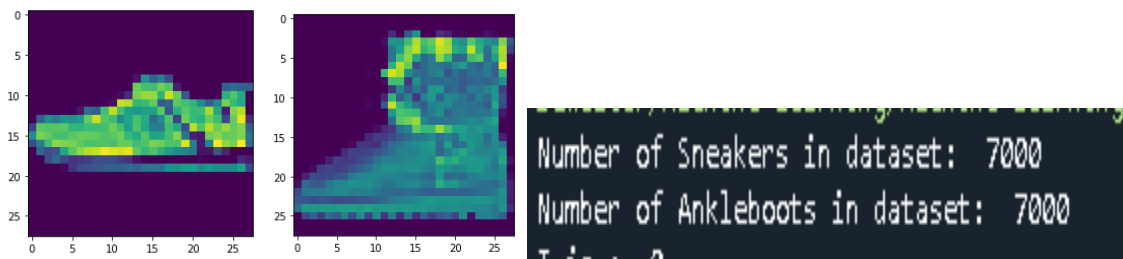
ASSIGNMENT 2 R00122087:

Task 1:

For task 1 I read in the data frame and create 2 subframes of ankle boots and sneakers separating by label. Next, I create 2 sub frames to store the labels before dropping the labels in the data sets. With that done I get the number of samples in each sub frame using the `len()` function as each row represents a sample in the form of 28x28 pixel values.

Finally I plot out the first row of each dataset to display the image of both an ankle boot and sneaker and reshaping the image to 28x28 to display each sample correctly.

These steps can be seen in the code from line 21 to 52.



Task 2:

Task 3 - Perceptron:

The perceptron classifier is created in the `perceptron_classifier` function which can be viewed in line 159 to 171.

Passing in the train and test data and label variables created in task 2 and is rather straightforward and I keep track of the time it takes for both the training and predicting of the perceptron using the [timeit](#) module.

The training and prediction times are passed back to task2 for later use as well as the accuracy score of the perceptron prediction.

```
def perceptron_classifier(train_data , train_labels , test_data, test_labels):  
  
    clf = linear_model.Perceptron()  
    starttime = timeit.default_timer()  
    clf.fit(train_data , train_labels)  
    training_time = timeit.default_timer() - starttime  
    print("Time taken for perceptron training: ", training_time)  
    starttime = timeit.default_timer()  
    prediction = clf.predict(test_data)  
    prediction_time = timeit.default_timer() - starttime  
    print("Time taken for perceptron prediction: " , prediction_time)  
    score = metrics.accuracy_score(test_labels, prediction)  
    return score , training_time , prediction_time
```

Task 4 Support Vector Machine:

The code for the svm classifier is viewable from line 173 to 193.

I begin by declaring an array to hold the gamma values which a for loop will iterate through and produce 5 different sets of results representing a svm classifier using each gamma value from 1e1 to 1e5.

I track the scores of these individual prediction values for each gamma and also track the training and prediction times once again for the svm which are passed back to the task 2 function for future use. The accuracy score is also passed back.

```
def svm_classifier(train_data , train_labels , test_data , test_labels):
    gamma_ranges = [1e1,1e2,1e3,1e4,1e5]
    best_scores=[]

    for gamma_range in gamma_ranges:
        scores = []
        clf = svm.SVC(kernel = 'rbf' ,gamma=gamma_range )
        starttime = timeit.default_timer()
        clf.fit(train_data,train_labels)
        training_time = timeit.default_timer() - starttime
        #print("Time taken for svm training: ", timeit.default_timer() - training_time)

        starttime = timeit.default_timer()
        prediction = clf.predict(test_data)
        prediction_time = timeit.default_timer() - starttime
        #print("Time taken for svm prediction: ", timeit.default_timer() - prediction_time)
        svm_score = metrics.accuracy_score(test_labels , prediction)
        scores.append(svm_score)
        #print('svm score is: ', svm_score)
    best_scores.append(max(scores))
    print("Best Score is : " , max(best_scores) , " At Gamma: " , gamma_ranges[best_scores.index(max(best_sc
    return svm_score , training_time , prediction_time
```

Task 5 KNN:

The knn_classifier function is viewable between lines 194 and 221. I create a Kfold with 4 splits and use the variable k value to vary the amount of k's I use. The best k is determined by assessing the accuracy scores created by the loop and selecting the best one. I then create a new classifier to train and predict using that k value. I then record that accuracy and return it as well as the best k value. Once again, the time for training the model and predicting are recorded, however in this case I record the times for the best k only.

```
def knn_classifier(train_data , train_labels , test_data , test_labels):

    kf = model_selection.KFold(n_splits=4)
    knn_scores = []
    k = 20
    for i in range(1,k):
        scores = []
        for train_index, test_index in kf.split(train_data):
            clf = neighbors.KNeighborsClassifier(n_neighbors=k)
            clf.fit(train_data.iloc[train_index].values, train_labels.iloc[train_index].values)
            prediction = clf.predict(train_data.iloc[test_index].values)
            score = metrics.accuracy_score(train_labels.iloc[test_index].values , prediction)
            scores.append(score)
        knn_scores.append(np.mean(scores))
    print("best value: ",np.argmax(knn_scores))
    best_k = np.argmax(knn_scores)+1
    print("Best K: ", best_k)

    clf = neighbors.KNeighborsClassifier(n_neighbors = best_k)
    starttime = timeit.default_timer()
    clf.fit(train_data,train_labels)
    train_time = timeit.default_timer() - starttime
    starttime = timeit.default_timer()
    prediction = clf.predict(test_data)
    prediction_time = timeit.default_timer() - starttime
    knn_score = metrics.accuracy_score(test_labels , prediction)
    print(" Best knn score: ", knn_score, "At K: " , best_k)
    return best_k, knn_score, train_time , prediction_time
```

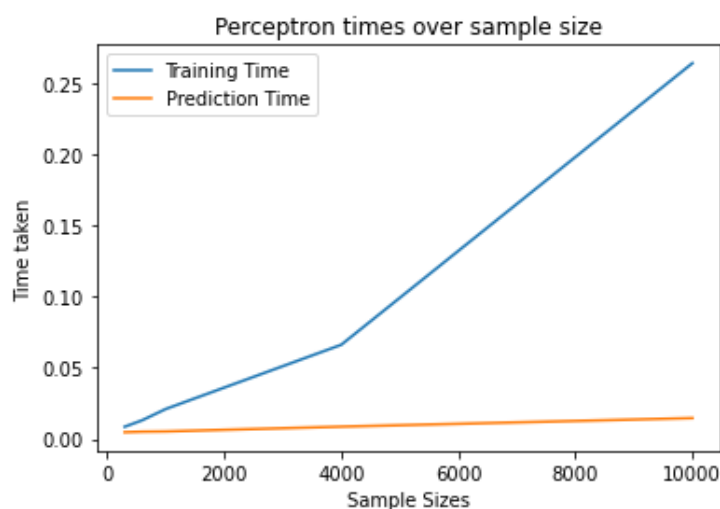
Task 6 Decision Tree's:

Task 7- Comparison:

The trends of each classifier we can extrapolate from the training and prediction times are:

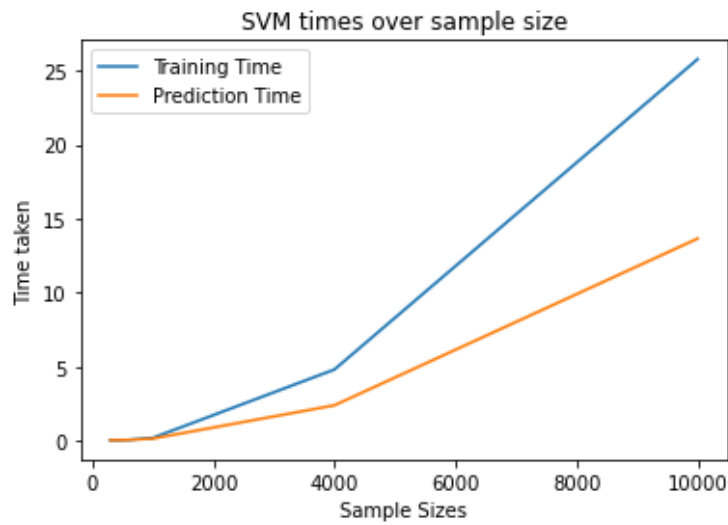
1. Perceptron:

As the sample size increases the prediction time also increases at a noticeable level, while prediction time is only impacted at an arguably insignificant level.



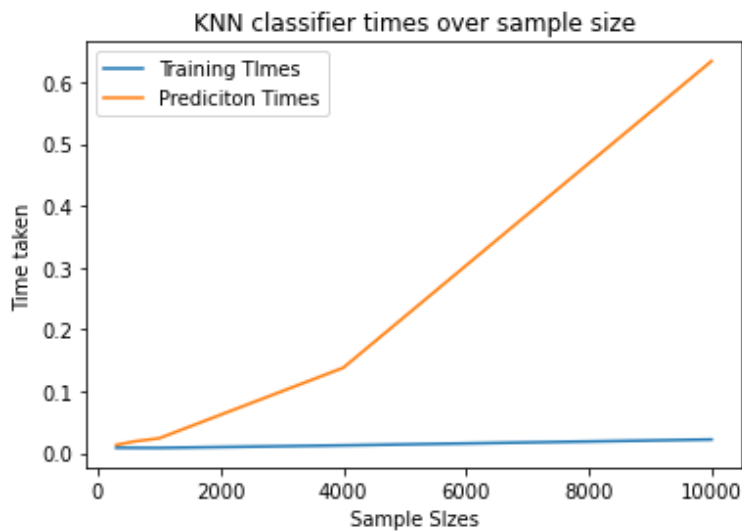
2. Support Vector Machine:

As the sample size increases in the support vector machine the training time also increases similarly to the perceptron, the prediction time also increases.



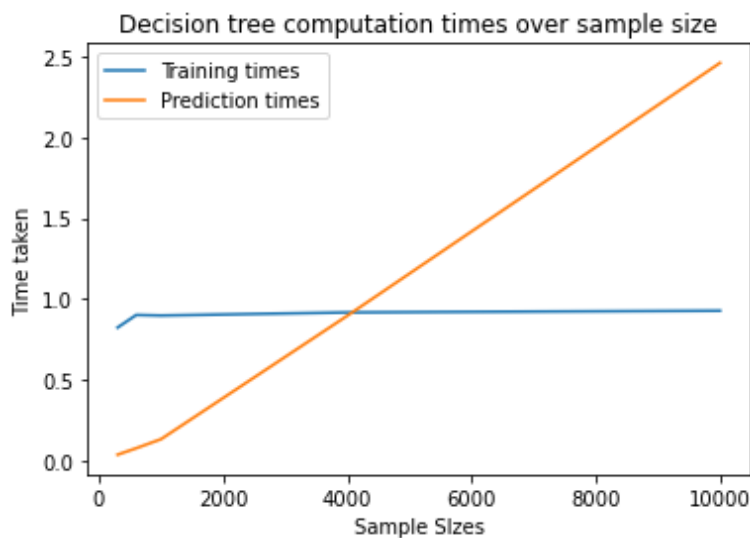
3. K-Nearest Neighbours:

In the k-nearest Neighbours data we can see that the training time remains the same as the optimal classifier has already been found before this operation. Prediction time increases which we can logically find reasonable due to the fact that as the sample size increases the prediction time also increases.



4. Decision Tree's:

I believe that there is an error in my code for the decision tree's however I will press on in my conclusion with them. In this case the training times for the decision tree's stayed almost stagnant while the prediction time saw a large increase in the amount of time taken to predict. I am not sure why this is happening but I believe it is due to the increased amount of sample sizes but that leads to the question why the training time is not also effected.



5. Conclusion:

In my opinion for the sample sizes and the nature of the samples I would recommend the use of the perceptron in this case. As we can see from the charts the training and prediction time is much smaller than other classifiers and therefore would scale well with a growing dataset.

