

## **Final Report: Time Series Assignment**

Prepared For: Jennifer Winikus  
Professor of CSE 454  
Department of Computer Science and Engineering  
University at Buffalo

Prepared By: Diakov, Alan  
Department of Computer Science Engineering  
School of Engineering and Applied Sciences  
University at Buffalo

## Introduction

This project was carried out by a student attending the University at Buffalo to examine the representation and classification of data in the form of time series.

This was done by using a synthetic control chart time series data set from the University of Irvine. It came in the form of a 600 by 60 table, where each row is a time series. Every 100 rows represented a different classification of time series (Normal, Cyclic, Increasing trend, Decreasing trend, Upward shift, Downward shift).

First, the number of data points had to be reduced by using PAA and SAX. Then an “expert” approach of the KNN algorithm was used to perform classification using manhattan and euclidian distance on the original data set and the PAA of it. Comparisons in performance were observed.

## Methods

This section will detail the methods gone about representing and classifying the data.

### Task 1: Create a PAA of the time series.

PAA is the Piecewise Aggregate Approximation of the data. In other words, it takes segments of size  $N$  of the data and uses the average of those  $N$  values to reduce the amount of data, while accurately representing it.

To implement the PAA, the data is taken in as a matrix and the int  $c$  as inputs. Data is the data set, where each row represents a time series. The int  $c$  is the size of the desired segments.

Initialize a matrix of size  $[600, 60/c]$  to store the PAA. Since each segment is of size  $c$ , there will be  $60/c$  segments. Iterate through all the times series in data. For each time series, get the average of every segment of size  $c$  and store that value into the PAA.

### Task 2: Create a SAX of the time series.

SAX is the Symbolic Aggregate Approximation of the data. In other words, it is the PAA of the data, but then each value of the PAA is given a letter representing the region it falls into. The regions are determined through the normal distribution of the data.

To implement the SAX, take in the data as a matrix the int  $L$  and the int  $c$  as inputs. Data is the data set, where each row represents a time series. The int  $c$  is the size of the desired segments for the PAA. The int  $L$  is the number of letters to be used to represent the data.

First, use the previous function to create a PAA of the data. Then find the standard deviation and mean of the data. Use the inverse norm of ( $\%$ , mean, std) to find which region each value in the PAA falls into, then assign a letter representing that region.

### Task 3: The processing process to establish training and testing data generation.

It is generally good practice to take 10% of the data and put it aside for testing. My “splitData” function takes the first 90 data points of each set of data and puts it in DataTrain, while the last 10 data points are put in DataTest. Therefore DataTrain has 540 rows of data, and DataTest has 60 rows of data. Also for the Labels, they correspond to:

1. Normal
2. Cyclic
3. Increasing trend
4. Decreasing trend

5. Upward shift
6. Downward shift

#### **Task 4: Explanation of my classification process and results.**

There are 4 different classification scenarios I had to consider.

1. The normal data set trained using euclidian distance.
2. The normal data set trained using manhattan distance.
3. The PAA of the data trained using euclidian distance.
4. The PAA of the data trained using manhattan distance.

For each classification, the method to train the data is almost identical. For the data set in question (original data or PAA of it), and the choice of distance, I train the data using:

```
“mdl=fitcknn(dataTrain, trainLabel, 'distance', @(x,y)dist_func(x,y), 'NumNeighbors', 1);”
```

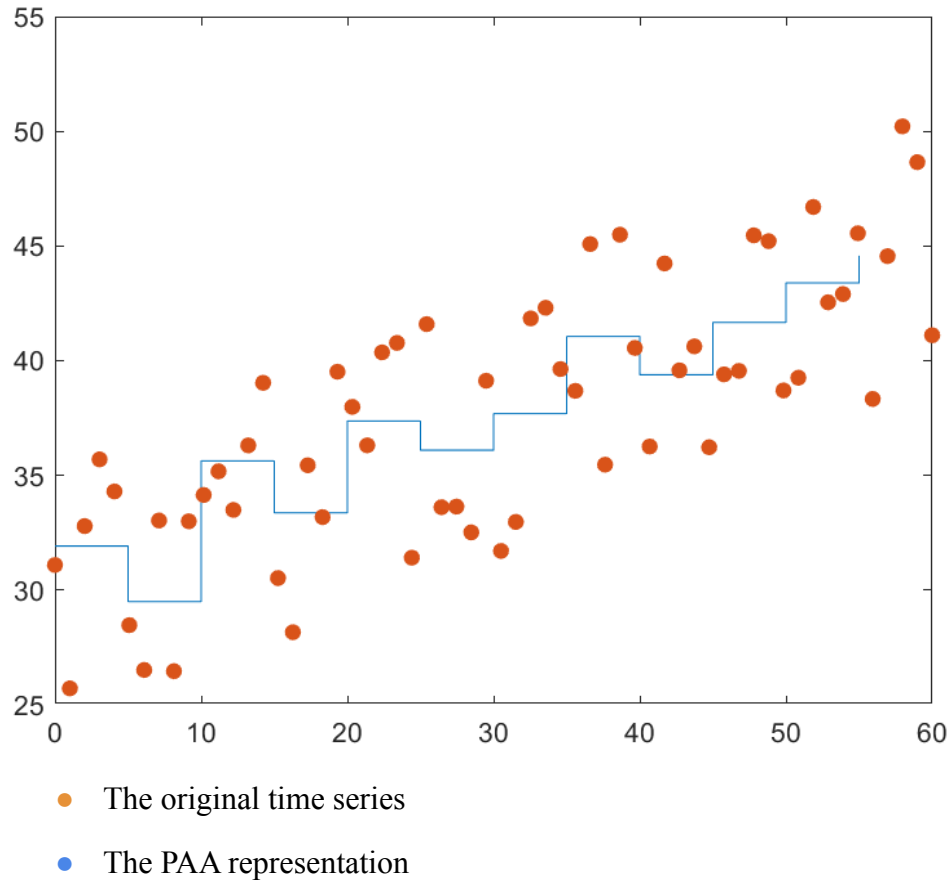
This trains the data through an “expert” approach. mdl makes a prediction using the nearest neighbor via the current distance function and chooses the class label of that neighbor as the prediction for the current test time series.

## Results

This section will show the results of the project.

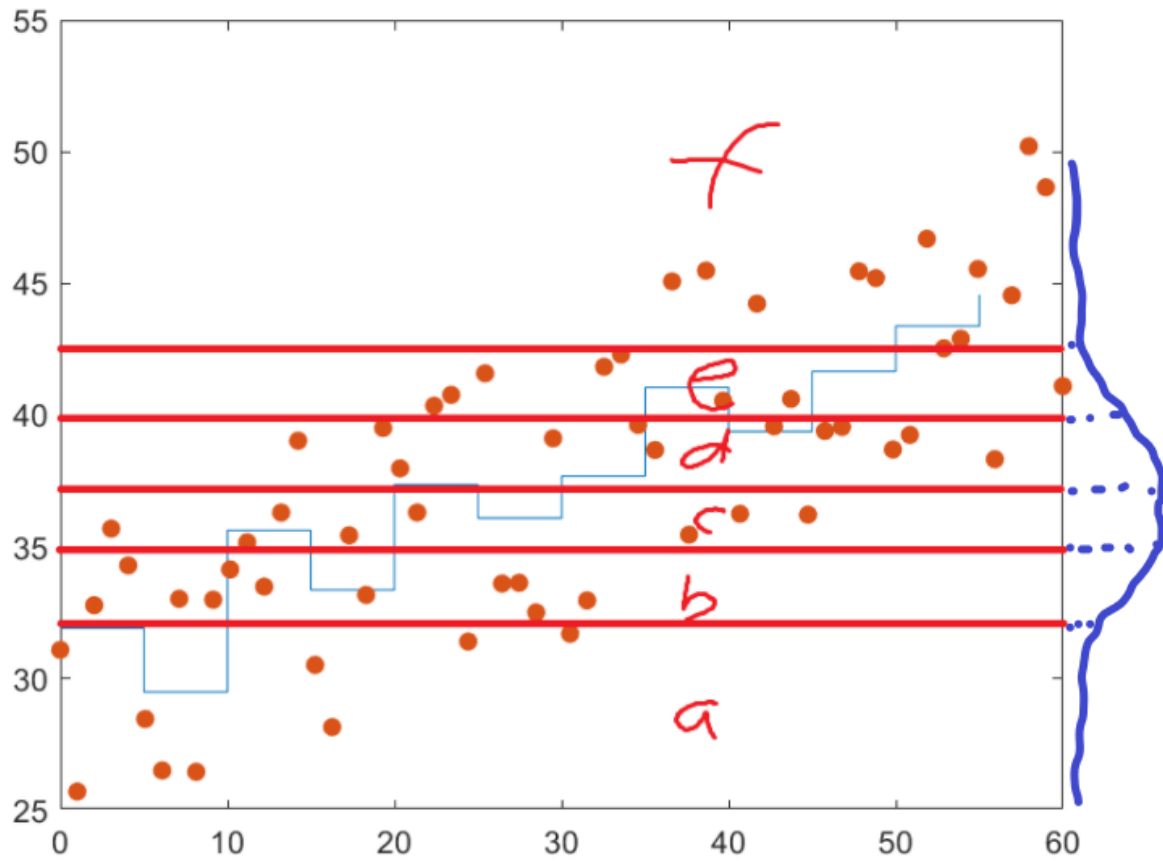
### The plot of a time series and the PAA representation:

Here is an example of a time series that is classified as increasing to show the PAA represents the data accurately.



### The plot of a time series and the SAX representation:

Here is an example of a time series that is classified as increasing to show the PAA represents the data accurately.



The SAX representation of the data becomes: “aacbdcdedeff”. By looking at the letters representing the data it is easy to see that this time series is increasing.

Results of the classification process. A confusion matrix will be displayed for each.

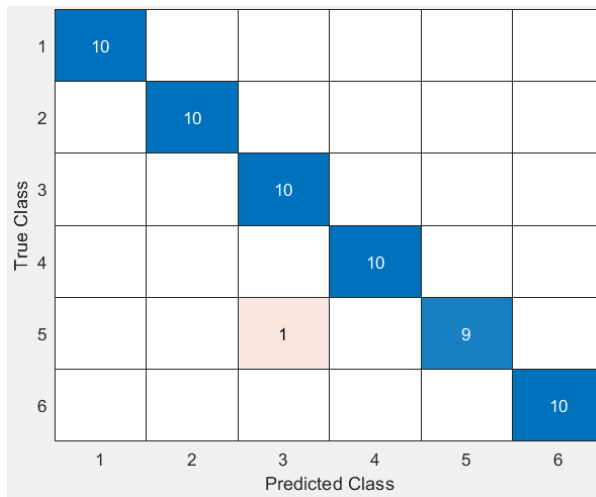
1. The normal data set trained using euclidian distance.

True Class	1	10					
	2		10				
	3			10			
	4				10		
	5					10	
	6						10
		1	2	3	4	5	6
		Predicted Class					

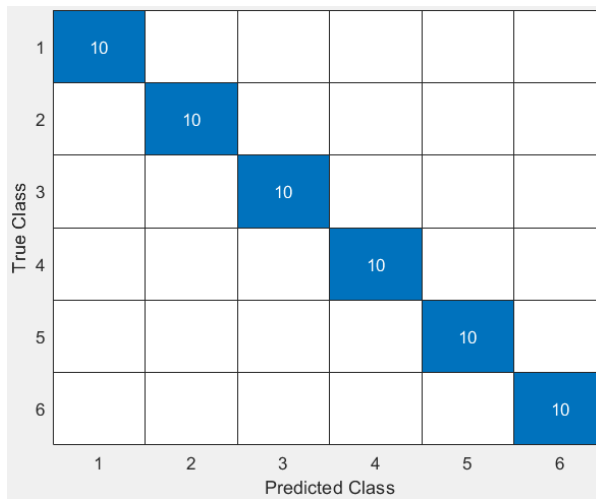
2. The normal data set trained using manhattan distance.

True Class	1	10					
	2		10				
	3			10			
	4				10		
	5					10	
	6						10
		1	2	3	4	5	6
		Predicted Class					

3. The PAA of the data trained using euclidian distance.



4. The PAA of the data trained using manhattan distance.





## Conclusion

This section will compare the results of the effectiveness of euclidian and manhattan distance.

### Comparisons and conclusion of the classification:

- There is no difference between the accuracy of the model on the original data using euclidian distance and manhattan distance. The reason is that both models have 100% accuracy on the test data.
- The model for using manhattan distance on the PAA still has 100% accuracy, but the euclidian distance on the PAA has only 98% accuracy. It mistook a time series that was categorized as an upward shift as an increasing time series, which is a fairly reasonable mistake. Therefore we can say that using the manhattan distance is more accurate for the PAA set.