

Improving the discovery of Motifs in high-dimensional sequences of varying length

M. Adnan Haider

Master of Science
School of Informatics
University of Edinburgh

2013

Abstract

Acknowledgements

Many thanks to my mummy for the numerous packed lunches; and of course to Igor, my faithful lab assistant.

Declaration

I declare that this thesis was composed by myself, that the work contained herein is my own except where explicitly stated otherwise in the text, and that this work has not been submitted for any other degree or professional qualification except as specified.

(M. Adnan Haider)

Table of Contents

| | | |
|----------|---|-----------|
| 1 | Introduction | 1 |
| 2 | Datasets | 5 |
| 3 | DTW-Background | 7 |
| 4 | Improving DTW | 11 |
| 4.1 | Feature Selection | 14 |
| 4.1.1 | Signal Filter | 15 |
| 4.1.2 | Downsampling | 19 |
| 4.2 | Domain Dependent Feature extraction | 21 |
| 4.2.1 | Mel Cepstrum Cepstral Coefficients | 22 |
| 5 | Extending DTW | 28 |
| 5.1 | Domain-independent feature extraction | 28 |
| 5.2 | Adapting DTW | 34 |
| 5.2.1 | Testing the methodology | 39 |
| | Bibliography | 46 |

Chapter 1

Introduction

Over the course of the last decade, the mining of time-series data have received considerable attention within the data mining and machine learning community. The term 'time series' denotes a set of observations concurring any activity against different periods of time. The duration of time period may be either in the order of milliseconds or monthly or even annually depending on the domain.

Mathematically, a time series is defined by the values $y_1, y_2 \dots y_n$ at times $t_1, t_2 \dots t_n$ where $y = f(t)$. The time t_i acts as an independent variable to estimate dependent variables y_i . The dimensionality of the series is denoted as \mathbf{n} where ' \mathbf{n} ' denotes the length of the sequence.

Time series analysis is used in many applications ranging from sales forecasting, budgetary analysis to stock market analysis and many more. One particular domain where the application of time series analysis is currently very popular is *motif* discovery- the problem of efficiently locating frequent/interesting sub-patterns in the data. The knowledge of motifs has been seen to have important applications in various aspects of data mining tasks. For instance motifs can use applied :

- to discover association rules [1]. the reflect information of 'primitive shapes.
- to specify the number of clusters for unsupervised clustering algorithms. Clustering is one of the most frequently used data mining tasks. It in-

volves an unsupervised process for partitioning a dataset into a specified number of meaningful groups. The knowledge of motifs give a good approximation on the number of meaningful groups that are present in the data. [2].

- to identify important sub-patterns in DNA and gene sequences [3]

In the analysis of speech data, motifs also play a very important role. Recent research have shown that detecting and isolating motifs in speech utterances is equivalent to extracting frequently spoken words spoken by the speaker(s) [4, 5]. The methodology proposed in these papers is based on constructing a framework that is entirely data driven i.e there is no intermediate recognition stage that maps an audio signal to a symbolic representation such as the 'phone' states in the HMM model. This results in the word acquisition process to be unsupervised which presents a completely different approach to the current speech recognition systems that are built using a supervised training methodology employing manually transcribed speech to model the underlying speech process.

To identify and extract motifs from time series data, various clustering algorithms have been proposed. The two most commonly used approaches are:

1. Dynamic time warping algorithm(DTW) [6, 7, 8, 9, 10, 11] that clusters similar sequences separated by time shifts and/or scale.
2. Single value decomposition(SVD)[12]. The entire time series data is approximated by a low-rank approximation matrix achieved through mapping the data onto a low dimensional orthogonal feature space.

Both these algorithms suffer from severe setbacks when applied directly on the 'raw' time series data. The DTW algorithm for example, suffers from large run-times when the length of the time series sequences are very long. This is because the time complexity of the algorithm is quadratic and is dependent on the dimensionality of the sequences i.e the length of the sequences. To address this issue window constraints (Itakura parallelogram[13], Sakoe-Chiba band[7]) are imposed that reduces size of the search space by forcing the algorithm to look for optimal path that are along the diagonal. Although introducing the window constraint does improve the time complexity but the reduction the search space leads to a substantial decrease in the accuracy of

the algorithm[11].

The SVD on the other hand also suffers from high computational cost when the number of samples \ll the size of the dimension of the data. For long time sequences, the high computational cost incurred by the SVD is therefore quite pronounced as each point in the time series corresponds to a feature/attribute. However, the main drawback of SVD is that it cannot be applied to time-series datasets where the length of the sequences vary. This constraint greatly reduces the type of time series domains to which SVD can be applied to. The speech corpus is an prime example of one such domain where SVD cannot be applied directly. Data sets comprised of speech utterances are a good example where recorded utterances do not share the same dimensionality (i.e the same length) as acoustically similar signals may be contracted/expanded versions of each other due to speaker variations, context etc.

The purpose of this project is to employ machine learning techniques to tackle and resolve the drawbacks incurred by both algorithms and thus improve their performance in handling long time series sequences that vary in length. In the first half of the project, I will be solely concentrating on improving the performance of the DTW algorithm in problem domains where the minimising the time complexity is a high priority. In the second half of the project, I will investigating ways to successfully adapt the SVD to work on long sequences that vary in length.

For this project, I will be using 3 time series data sets (details in chapter 2):

- TIGITS
- INLINESKATE
- CINC_ECG_TORSO

The prime objective here is to improve the accuracy of the DTW and SVD algorithm in tackling high dimensional time series sequences that vary in length while minimising the run-time to a minimum. To evaluate and compare the merits of different proposed changes, I combined both algorithms with the K nearest neighbour classifier and partitioned each of the above data sets into two disjoints subsets: one to be used as a test set and the other as a train-

ing set. The reason for choosing the nearest neighbour classifier is because it shares almost the same methodology as the algorithms used to detect motifs. Motif detection algorithms are memory based i.e they rely on comparing each sequence with other sequences in the data set. K nearest neighbours differs from motif detection on two aspects: each sequence in the test set is compared with sequences that belong only to the training data set and the entire process is supervised i.e correct label information are available to check the accuracy of the algorithm. The availability of correct labels allow accuracy and run-time scores of the nearest neighbour methods to serve as a useful evaluation metric to compare and evaluate various adaptations of the DTW and SVD algorithm. Improvements in SVD and DTW that lead to greater classification accuracy scores and low run-times are mostly likely to result in better extraction of motifs in the unsupervised context.

The dissertation is organised as follows: Chapter 2 gives a description of the 3 time-series datasets used for this project. Chapter 4 provides a detailed background description of the DTW algorithm. Chapter 3 and 4 investigates methods to improve the performance of the DTW algorithm in terms of both accuracy and speed. Chapter4...

Chapter 2

Datasets

The primary dataset that I have used for this project is the 'TIGITS' corpus. (I need to give more description here)

For Training : The entire training data is used –To contain the computational complexity, I am using samples from production 'a'

For the training set : To reduce the average mean time , I am using half of the training data set by choosing samples from one production: I have chosen:

225 samples from the boy category

234 samples from the girl category

495 samples from the men category

513 samples from the women category

Note : the size of the training set is half of the original training set but contains examples of all classes [1-9]

For the test set : Due to the high computational complexity, I am using only 1/3 of the test set I have chosen 162 random samples from boys

162 random samples from girls

326 random samples from men

326 random samples from women

Apart from the TIGITS, I have used two datasets from the UCR database:

The description of the data sets used for the next set of experiments are as follows:

1. CinC_ECG_torso

- Length of the time series:1639
- Size of test set:1380
- Size of training set:40
- Number of classes:4

2. InLineSkate

- Length of the time series:1882
- Size of test set:550
- Size of training set:100
- Number of classes:7

Chapter 3

DTW-Background

The Dynamic Time Warping algorithm measures the similarity between sequences varying in both time and speed. Formally, the problem formulation of the algorithm is stated as follows: Given two time series X , and Y , of lengths $|X|$ and $|Y|$,

$$X = x_1, x_2, \dots, x_{|X|} \quad (3.1)$$

$$Y = y_1, y_2, \dots, y_{|Y|} \quad (3.2)$$

construct a warping path W

$$W = w_1, w_2, \dots, w_k \text{ where } \max(|X|, |Y|) \leq k \leq |X| + |Y|$$

- Here k denotes the length of the warping path and the m th element of the warping path is $w_l = (n_l, m_l) \in [1 : N] \times [1 : M]$ for $l \in [1 : k]$ where n_l is an index from the time series X and m_l is an index from the time series Y .

To properly understand the mechanism of the DTW algorithm, the definition of some key terminologies must first be stated:

1. Warping path: An (N, M) -warping path (or simply referred to as warping path if N and M are clear from the context) is a sequence $w = (w_1, \dots, w_k)$ with $w_l = (n_l, m_l) \in [1 : N] \times [1 : M]$ for $l \in [1 : k]$ satisfying the following three conditions.

- (a) Boundary condition: $p_1 = (1, 1)$ and $p_k = (N, M)$. The boundary condition enforces that the first elements of X and Y as well as the last elements of X and Y to be aligned with each other. In other words, the alignment refers to the entire sequences X and Y .
- (b) Monotonicity condition requires that the path will not turn back on itself, both the i and j indexes either stay the same or increase, they never decrease.
- (c) Step-size condition: $p_{l+1} - p_l \in \{(1, 0), (0, 1), (1, 1)\}$ for $l \in [1:k-1]$. The step size condition expresses a kind of continuity condition: no element in X and Y can be omitted and there are no replications in the alignment

Intuitively speaking, the (N, M) warping path $p = (p_1, \dots, p_k)$ defines an alignment between two sequences $X = (x_1, x_2, \dots, x_N)$ and $Y = (y_1, y_2, \dots, y_M)$ by assigning the element x_i of X to the element y_j of Y .

2. Optimum Warping Path :

The optimal warp path corresponds to the minimum-distance warp path, where the distance of a warp path W is given as

$$Dist(W) = \sum_{i=1}^K dist(X, Y)_{|(w_i)}$$

$dist(X, Y)_{|(w_i)}$ represents the distance computed using an appropriate cost function between the time series points of x_{ni} of sequence X and y_{mi} of sequence Y .

$$dist(X, Y)_{|(w_i)} = dist(x_{ni}, y_{mi})$$

The goal of the DTW algorithm is to compute the distance of the optimal warping path between two time series sequences. Instead of attempting to solve the entire problem all at once, the algorithm utilises the technique of dynamic programming to find an optimum alignment between two sequences through the computation of local distances between the points in the temporal sequences. The algorithm proceeds by iteratively filling in values for each cell (i, j) in the $|X|$ by $|Y|$ cost matrix D . The value of the cell (i, j) is given by

$D(x_{ni}, y_{mj})$ which corresponds to the minimum- distance warp path :

$$D(i, j) = \text{Dist}(i, j) + \min(D(i-1, j), D(i-1, j-1), D(i, j-1))$$

An outline of the baseline DTW algorithm is given below:

Algorithm 1 Value-Based DTW

```

1: procedure VALUE-BASED(seq1, seq2) ▷ two raw sequences
2:   DTW= zeros(length(seq1)+1,length(seq2)+1)
3:   for i=1: to length(seq1) do ▷ Initialise the DTW cost matrix
4:     DTW(i,0) = ∞
5:   end for
6:   for i=1 to length(seq2) do
7:     DTW(0,i) = ∞
8:   end for
9:   for i=2 to length(seq1) do
10:    for j=2 to length(seq2) do ▷ cost(a,b)≡euclid(a,b)
11:      DTW(i,j) = cost(seq1(i),seq2(j)) + min{ DTW(i-1,j)+DTW(i,j-1)+DTW(i-1,j-1)}
12:    end for
13:  end for
14:  return result =  $\frac{\text{DTW}(n,m)}{nm}$  ▷ n=length(seq1), m=length(seq2)
15: end procedure

```

The figure below gives an example of the optimal path found by the algorithm.

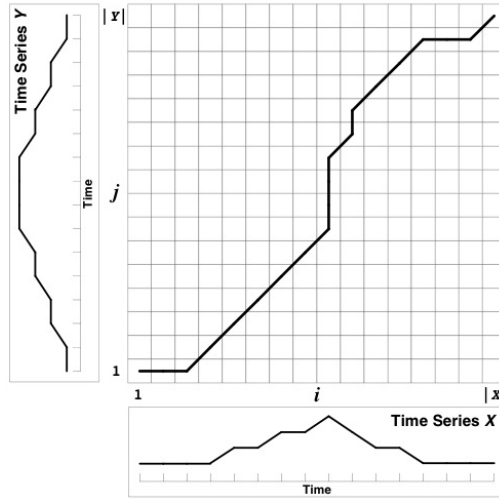


Figure 3.1: A cost matrix with the minimum-distance warp path traced through it.

The computational complexity of the DTW algorithm is $O(n^2)$ where n denotes the length of the sequences that are being compared. Thus for time series domains having high dimensions (long sequences), the time and computational costs incurred by the algorithm are quite high. To address this issue, two well-known global window constraints are employed: the Sakoe-Chiba band[18] and the Itakura parallelogram[19]. Figure 3.2 gives an illustration of the use of both constraints:

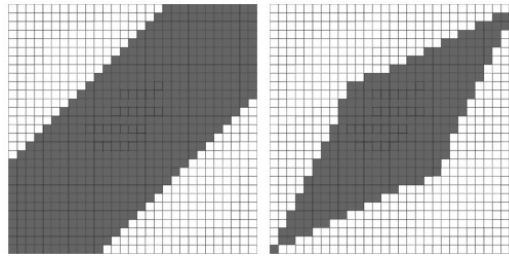


Figure 3.2: Two constraints: Sakoe-Chiba Band (left) and an Itakura Parallelogram (right), both have a width of 5.

The Sakoe-Chiba band runs along the main diagonal and has a fixed (horizontal and vertical) width. The Itakura parallelogram on the other hand describes a region that constrains the slope of a warping path. To constraint the time complexity to a minimum, the vast majority of the data mining researchers use a Sakoe-Chiba Band with a 10% width.

Chapter 4

Improving DTW

The Dynamic Time Warping(DTW) algorithm is the one of the oldest algorithm that is used to compare and cluster sequences varying in time, scale and speed. Given two temporal sequences, the algorithm utilises the technique of dynamic programming to compute the cost of the optimum alignment path between them. The computed cost gives an indication of the degree of similarity. The smaller the cost, the more similar the sequences are. Intuitively speaking, DTW can be seen as a clustering algorithm that clusters patterns that share roughly the same shape. As I have stated in the previous chapter, the time and computational complexity of this algorithm is $O(n^2)$ where n denotes the length of the sequences that are being compared. Thus for time series domains having high dimensions i.e long sequences, DTW becomes a very unattractive choice as it suffers from high computational and time costs.

To address the issue of the curse of dimensionality, DTW algorithms employ a window constraint to reduce the search space. The most commonly used are Sakoe-Chuba Band[7] and the Itakura window constraint [13]. Figure[3.2] gives an illustration of the nature of these window constraints. Such constraints determine the allowable shapes that a warping path can take by restricting the DTW to construct optimal warping paths only through a restricted number of cells around the diagonal region of the cost matrix. As the dimensionality(length) of the sequences increases, the size of the window is adjusted accordingly. Rigid window constraints impose a more rigid alignment that prevent an overly temporal skew between two sequences, by

keeping frames of one sequence from getting too far from the other. The vast majority of the data mining researchers use a Sakoe-Chiba Band with a 10% width for the global constraint [14] to constraint the time complexity of DTW to a minimum. For finding motifs in data sets comprising of speech utterances, the use of rigid window constraints is highly undesirable. Utterances corresponding to the same lexical identity may suffer from large variations in speed, scale and time due to the word(s) being spoken in difference contexts, by different speakers, in different environments etc. Thus it is necessary to explore alternative techniques to window constraints that can reduce the time complexity of DTW without degrading its accuracy.

Before investigating methods to improve the DTW algorithm itself, it is highly necessary to first understand the nature of the data sequences that the DTW is presented with. Achieving a thorough understanding of the data can result in the extraction of a smaller set of relevant features that can be used to achieve better discrimination between different classes/motifs. In this chapter, I investigate domain-dependent preprocessing techniques to improve the performance of the baseline DTW algorithms in clustering patterns that have long lengths.

There are presently two groups of preprocessing techniques commonly used to address this issue:

- Feature Selection
- Feature Extraction

Feature selection techniques involve selecting only a subset of attributes from the original data. With respect to the time series data, the technique is analogous to sub-sampling the sequence. To remove redundant features, one of the most popular feature selection approach is the exploratory data analysis(EDA). EDA is an approach to data analysis that postpones the usual assumptions about what kind of model the data follows with the more direct approach of allowing the data itself to reveal its underlying structure and models. The particular techniques employed in EDA are often quite simple, consisting of various techniques of:

1. Plotting the raw data such as data traces, histograms, histograms, probability plots, lag plots, block plots, and Youden plots.

2. Plotting simple statistics such as mean plots, standard deviation plots, box plots, and main effects plots of the raw data.
3. Positioning such plots so as to maximise our natural pattern-recognition abilities, such as using multiple plots per page.

Apart from removing redundant features, we can also construct more useful features from the existing ones. Feature extraction processes are concerned with the range of techniques that apply an appropriate functional mapping to the original attributes to extract new features. The intuition behind feature extraction is that the data vectors $\{x_n\}$ typically lie close to a non-linear manifold whose intrinsic dimensionality is smaller than that of the input space as a result of strong correlations between the input features. Hence by using appropriate functional mapping, we obtain a smaller set of features that capture the intrinsic correlation between the input features. By doing so, we move from working in high dimensional spaces to working in low dimensional spaces. The choice of appropriate functional mapping can also improve the clustering of data. For example, let's consider figure 4.1: The left-hand plot represents the locations of two dimensional data points in the original input space. The colours red and blue denote the classes to which the data points belong to. To cluster the data with respect to their classes, it will be ideal if we can partition the input space into disjoint regions where intraclass variation is small and interclass separation is large. For this example, this is achieved by mapping the points to a feature space spanned by two gaussian basis functions (shown on the right). Now, we can partition the feature space into two disjoint regions, one of each cluster.

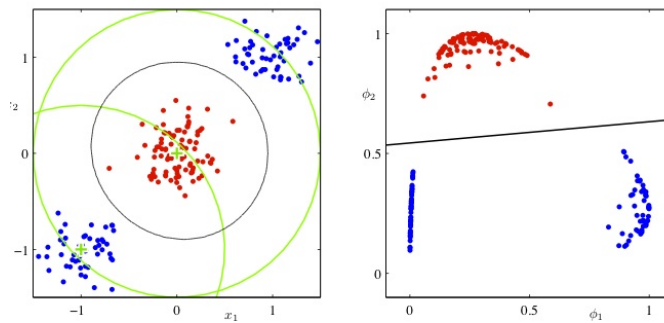


Figure 4.1: The figure on the right corresponds to location of the data points in the feature space spanned by gaussian basis functions $\phi_1(x)$ and $\phi_2(x)$

In the following sections, I investigate whether the application of exploratory data analysis and the construction of features that integrate metadata mainly knowledge of the domain can aid the baseline DTW algorithm in achieving low run times and good accuracy on long time series sequences. The primary data set that I will be using is the TIDIGITS corpus as the time series sequences of this data set have lengths that are in the order of 10^4 .

4.1 Feature Selection

The computational and time complexity associated with the DTW algorithm is governed by the length of the time series. Removing redundant features can result in a great reduction in the time complexity without any negative impact on the accuracy. To get an idea about the underlying structure of the data, I studied the plots of the time series sequences along with listening to the individual samples. Figure 4.2 shows the plot of raw signal corresponding to the word '8' by a speaker from the *boy* category.

From the visual and auditory analysis, I have made the following observations:

- Long durations of silence occupy the beginning and end of each utterance. The durations of the interesting regions that actually contain information about the spoken digit is quite small in comparison to the durations of silence regions. Removing these silence segments allow reduction in the dimensionality of the time series sequences with minimal loss of information.
- The recordings are highly distorted when played in *matlab*. The distorted signals fail to provide any type of auditory clue about category of the speaker i.e whether the speaker belongs to { boy,girl, men,women} the signal has to be played multiple time to correctly identify the spoken word.

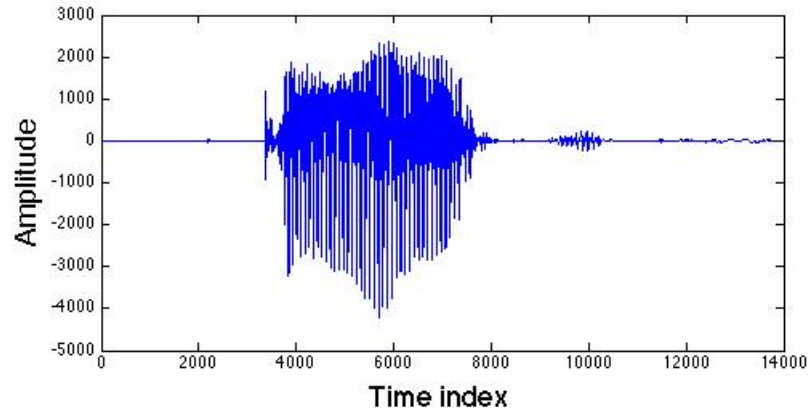


Figure 4.2: 'Raw 'signal

4.1.1 Signal Filter

To remove the redundant attributes from the time series utterances, I have constructed the following algorithm: 'SignalFilter' that performs feature selection by removing segments corresponding to durations of silence. An outline of the algorithm is as follows:

Algorithm 2 SignalFilter

```

1: procedure SIGNALFILTER(signal) ▷ raw signal
2:   threshold = 0
3:   maxAmplitude = max(rawSignal)
4:   Adapt the threshold based on the value taken by the maximum amplitude
5:   output ← removeSilence(rawSignal, threshold)
6:   return output
7: end procedure

```

The algorithm employs an adaptive threshold to select and remove redundant attributes. The value of the threshold is dependent on the maximum amplitude of the sequence. The value set is comparatively higher for utterances of speakers having a loud and deep voice and lower for utterances for speakers having gentle and low voice.

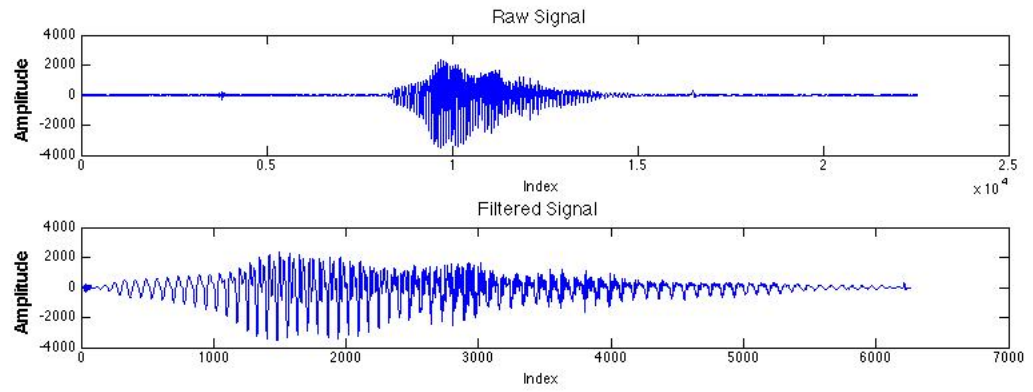


Figure 4.3: shows the raw acoustic signal corresponding to the utterances of the digit '8' alongside with the filtered signal at the bottom.

From figure 4.3, it can be observed that the filter preserves the interesting patterns associated with the utterance while succeeding in reducing the dimensionality of the data. To investigate the effect of introducing this prior feature selection step on the performance of the baseline DTW algorithm, I conducted the following experiment:

SETUP:

- Dataset used: TIDIGITS Test-set size : 976 samples

| category | sample size |
|----------|-------------|
| boy | 162 |
| girl | 162 |
| men | 326 |
| women | 326 |

Training data set size: 1467

| category | sample size |
|----------|-------------|
| boy | 225 |
| girl | 234 |
| men | 495 |
| women | 513 |

- An outline of DTW used algorithm used for this experiment is given below.

Algorithm 3 Value-Based DTW

```

1: procedure VALUE-BASED(seq1, seq2) ▷ two raw sequences
2:   DTW = zeros(length(seq1)+1, length(seq2)+1)
3:   w = max( $\lceil 0.1 * \max(n, m) \rceil$ , abs(n-m)) ▷ Window constraint
4:   for i=1: to length(seq1) do ▷ Initialise the DTW cost matrix
5:     DTW(i,0) =  $\infty$ 
6:   end for
7:   for i=1 to length(seq2) do
8:     DTW(0,i) =  $\infty$ 
9:   end for
10:  for i=2 to length(seq1) do
11:    for j=max(2, i-w) to min(length(seq2), i+w) do ▷
12:      cost(a,b)  $\equiv$  euclid(a,b)
13:      DTW(i,j) = cost(seq1(i), seq2(j)) + min{ DTW(i-1,j)+DTW(i,j-1)+DTW(i-1,j-1) }
14:    end for
15:  end for
16:  return result =  $\frac{DTW(n,m)}{nm}$  ▷ n=length(seq1), m=length(seq2)
17: end procedure

```

Note : The main objective of my research is to improve the accuracy of the DTW algorithm while reducing the time and computational cost to a **minimum**. Even after applying the feature selection process, from initial experiments on few samples, I have found that the computational cost incurred by the algorithm is still very high. Hence to minimise run time, I employed the Sakoe-Chuba band that has an adaptive window size of : $w = \max(\lceil 0.1 * \max(n, m) \rceil, \text{abs}(n-m))$. The lower bound of the window size is set to 10% of the size of the longest sequence which is the standard size that vast majority of the data mining researchers [14] use to keep the time complexity of DTW to a minimum.

- RESULTS:

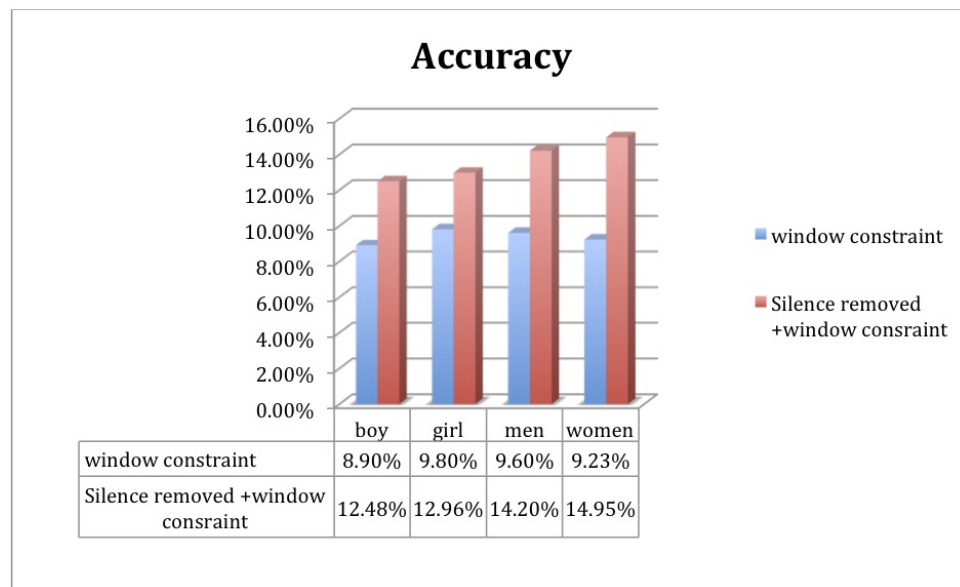


Figure 4.4

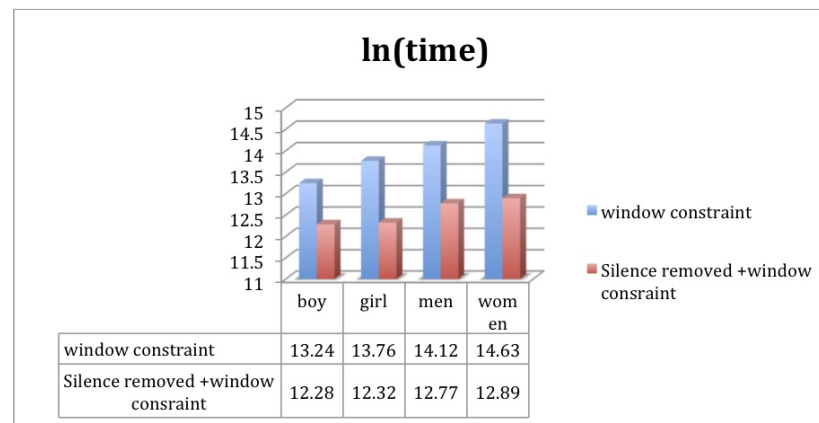


Figure 4.5

Observations:

- The DTW algorithm achieves very poor accuracy. The poor accuracy is a result of one or a combination of the following three factors:
 - the use of raw values as features: the numerical value associated with each time point is not a complete picture of the data point in relation to the rest of the sequence.
 - the use of the rigid window constraint- the low accuracy may be attributed to the optimum warping path between similar patterns lying outside boundaries of the Sakoe-Chuba bands[5, 4].

- not integrating knowledge of the domain: The data set is comprised of speech utterances. It is a widely known fact that for speech data, the MFCC feature vectors capture the information of phones that make up a word. Since different lexical identities are composed of different phones, these use of MFCC vectors can increase the inter class variance of the samples. (details of MFCC to follow)[4, 15, 5, 16, 17].
- Removing ‘silence’ segments improve **both** the accuracy and the run-time of the algorithm. While the decrease in run-time is quite obvious, the increase in accuracy is not. All utterances contains durations of silence. Taking the silence regions into account decreases the inter class variance as they bring in an unwanted notion of similarity in dissimilar patterns.

4.1.2 Downsampling

From conducting further exploratory data analysis, I have observed that if I down-sample the utterances by $\frac{1}{2}$ which in other words corresponds to decreasing the sampling frequency by half, the global shape is still preserved even though some local information is lost as shown in figure 4.6. This results in further reduction in the dimensionality of the sequences i.e the length of the sequences.

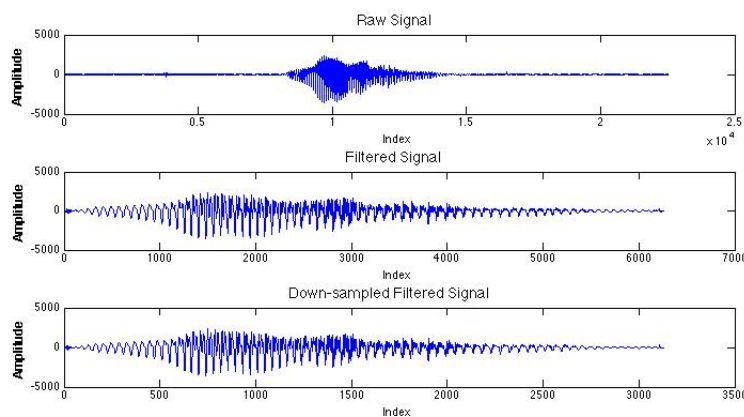


Figure 4.6: shows the raw acoustic signal of the digit ‘8’ (top figure), the silence removed version of the signal(middle) and the silence removed and down sampled version of the acoustic signal (bottom)

To investigate the effect of performing further dimensionality reduction on the time series sequences through down sampling on the accuracy of the DTW, I performed the following experiment:

Dataset: The TIDIGITS training and set used in 4.1

Algorithm : baseline DTW augmented with an adaptive window constraint(4.1)

Approaches being compared: Removing silence utterances VS Removing silence utterance + downsampling VS baseline

Results:

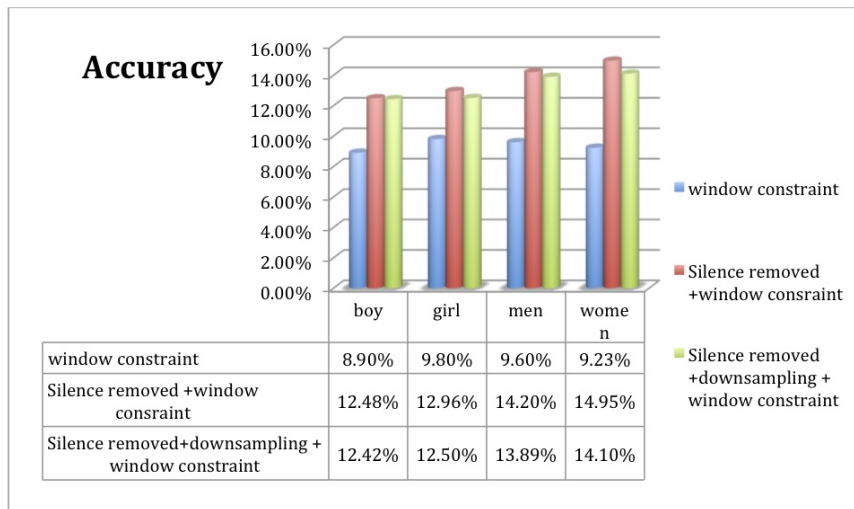


Figure 4.7

Observations:

Performing the two stage feature selection step that involves removing 'silence' utterances followed by downsampling allow DTW to still construct more optimal paths that are aligned along the main diagonal for similar patterns than using the entire 'raw' time series sequences. This procedure results in a 4% increase in accuracy on average. Furthermore, the loss of local information through downsampling the non-silence regions leads to a minimal decrease in the accuracy of the algorithm in comparison to using the entire non silence regions for pattern matching. This supports the claim that the classes are differentiated mainly by their global shape.

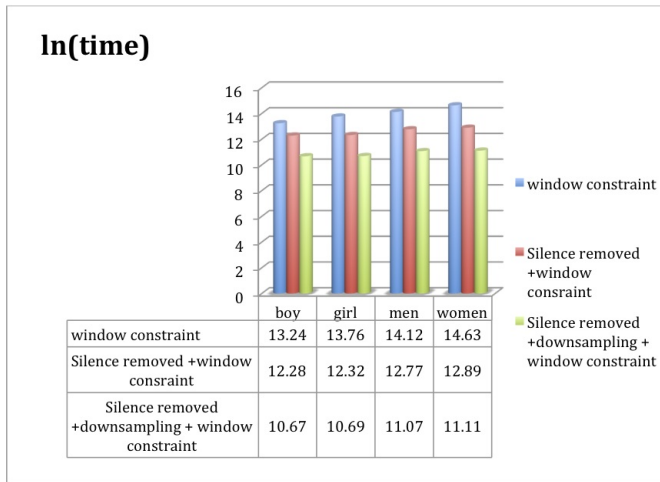


Figure 4.8: Integrating downsampling decreases the ln(time)

Augmenting the down sampling procedure to the feature selection process leads to a further decrease of 1.5 in the ln(run time) of the DTW algorithm in comparison to just removing silence regions as a preprocessing step. This is expected since in the first stage of the preprocessing phase, redundant features are dropped which reduces the dimensionality(i.e length) of the sequences. The length of the sequences is reduced even further through downsampling the the output sequences of stage 1. Since the computational cost of DTW is directly dependent on the length of the sequences, thus decrements in dimensionality leads to a decrease in the computational cost.

4.2 Domain Dependent Feature extraction

From the analysis conducted so far, it can be concluded that heuristically selecting only significant attributes from the time series sequences does **improve** both the accuracy and the speed of the DTW. However, from the observation of the experimental results gathered so far, it is quite clear that the accuracy of the algorithm is very low. So far, I have investigated the effect on using ‘raw’ values of the time series sequences on the performance of the DTW. In this section, I investigate the contribution of using the window constraint and the contribution of using domain dependent features on the performance of the DTW algorithm. There are two motivations behind conducting this analysis:

- The primary motivation is to improve the speed of the DTW algorithm

without degrading the accuracy. Choosing an appropriate functional mapping, can help map the data to a lower dimensional feature space that can capture the intrinsic qualities of the data. Thus constructing appropriate functional mappings can achieve both dimensionality reduction and higher accuracy.

- The other motivation is to investigate to what degree is the low accuracy error credited to not using domain dependent features and to using a rigid window constraint has on the low accuracy(The features that we have considered so far are the raw values indexed by time). Understanding the underlying factors that that influence the performance of the algorithm can provide insights on what aspect of the algorithm needs to be changed to gain better performance.

4.2.1 Mel Cepstrum Cepstral Coefficients

The primary data set that I am working with is the TIGITS corpus which is composed of speech utterances. For speech, the most commonly used features are the MFCC features-mel cepstrum cepstral coefficients. This feature representation is based on the idea of the cepstrum. For human speech, a speech waveform is created when a glottal source waveform of a particular frequency is passed through the vocal tract which because of its shape has a particular filtering characteristic. The exact position of the vocal tract is in fact the key attribute in providing useful information about phones(units of sounds). Cepstrum provides a useful way to separate the information of the vocal tract from the glottal source.

An outline of the MFCC feature extraction is given below:

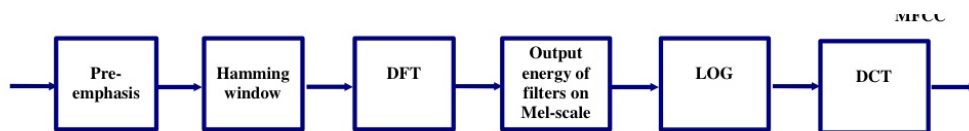


Figure 4.9: MFCC feature extraction

- Pre-emphasis: boosts the energy of the signal at high frequencies to improve phone detection

- ii Windowing: partitions the time series sequence into frames using a hamming window
- iii DFT: extracts spectral information at different frequency bands
- iv Mel scale : transforming to mel scale improves recognition performance by allowing the model to take into account the property of human hearing
- v Log : makes the feature less sensitive to variations in input such as power variations on the speakers mouth.
- vi Cepstrum : separate the information of the vocal tract from the glottal source. The first 12 cepstral values from spectrum of the log of the spectrum are used

The windowing process allows dimensionality reduction i.e the length of the time-series sequences is reduced. Each sequence is segmented into frames of length 20 to 30 ms. These frames are then mapped using the procedure above to MFCC feature vectors. Since the result sequence of vectors is much smaller than the length of the original sequence resulting in the size of the DTW cost matrix is much smaller than before which in turn lowers the time and computation cost incurred by the algorithm.

The experiments conducted in section 4.1.1 and 4.1.2, have shown that the DTW algorithm performs very poorly in terms of accuracy on the TIDIGITS test data when it employs the narrow Sakoe-Chuba band to minimise the time complexity. The reason for this low accuracy was credited to the influence of one or a combination of the following factors : using a narrow window constraint, using raw attribute values and not using domain-dependent features. Having already investigated the influence of using raw values(4.1), for this part of the project, I constructed the following 3 models in order to investigate the relative isolated impacts of the other two factors. The models were tested on the TIDIGITS data set, thus allowing the results to be compared with the results of the previous experiments so far :

- i Model 1 employs the MFCC feature extraction mapping as a preprocessing step. The extracted features are used by the valued-based DTW algorithm augmented with the adaptive Sakoe-Chuba band constraint

described in 4.1.1 to cluster similar patterns. The performance of this model can be compared to the performance of the base line model to measure the relative impact on the accuracy and run-time of the DTW of replacing each numerical value associated with each time index with a domain dependent feature vector.

- ii Model 2 employs a two stage preprocessing step. The feature selection procedure discussed in 4.1.1 is first applied to remove redundant features followed by MFCC feature extraction that achieves further dimensionality reduction. In this model dimensionality reduction occurs at both stages of the pre-processing step. The downsampling method discussed in 4.1.2 was deemed not necessary as a feature selection step since the feature extraction phase allows further reduction in dimensionality without any loss of information. The sequence of vectors are then used by the value-based DTW algorithm augmented with the window constraint to find optimal warping along the main diagonal of the DTW cost matrix.

By comparing the results of this model with the results of model 1 and the model discussed in 4.1, the optimum preprocessing strategy can be determined.

- iii Model 3 is identical to the version 2 with the exception that this version does not employ the window constraint. The main purpose of using this model is to check the relative impact of using window constraints. The difference in the performance between using model 2 and model 3 can be used to check to what degree is employing such a rigid window constraint affect the accuracy and run-time of the algorithm.

Experimental setup:

Dataset : The TIDIGITS training and test set (Chapter 2, 4.1.1)

RESULTS:

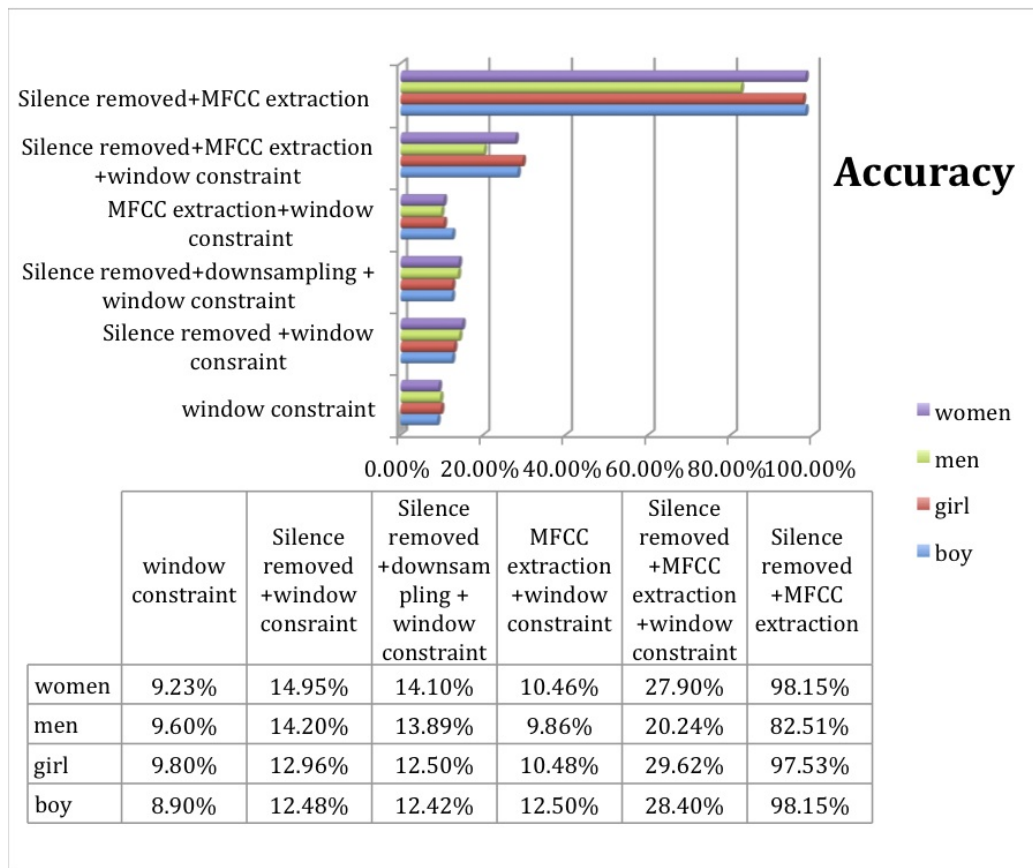


Figure 4.10: Accuracy achieved by all models

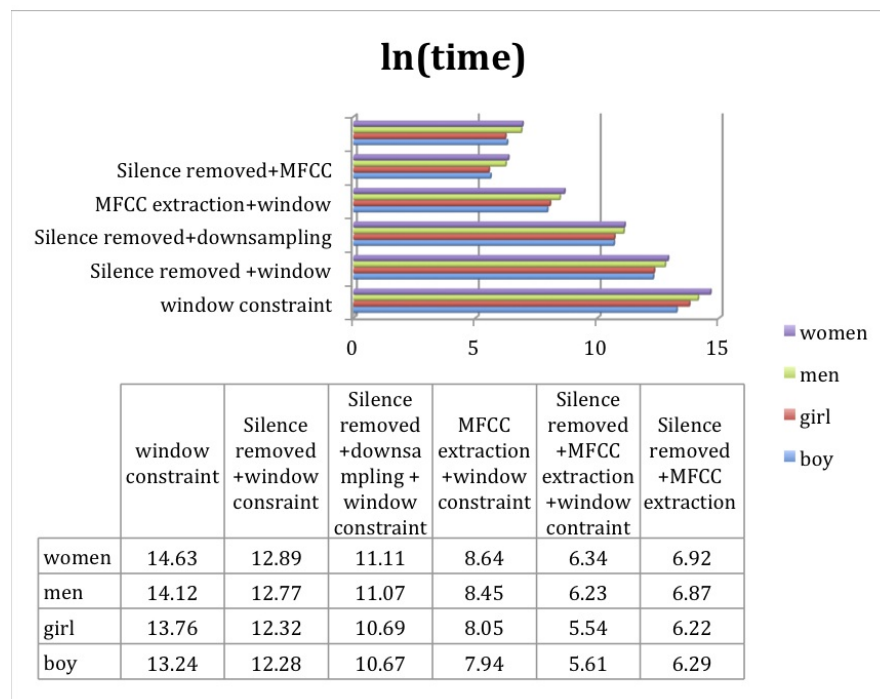


Figure 4.11: The run time incurred by all models

Observations:

- Partitioning the raw acoustic signal into frames and replacing each frame with an MFCC feature vector surprisingly has minimal effect on the accuracy of the DTW algorithm equipped with the adaptive Sakoe-Chuba band constraint.
- In comparison to the base-line model, the accuracy of the window-constrained DTW only increases by 1.3% on average across all 4 categories. The run-time on the other hand has improved considerably which is obvious since the acoustic signal is mapped to smaller sequence of frames. From examining figure 4.12, it can be seen that the windowing process results in the average $\ln(\text{run time})$ has decreased by 5.7.
- Removing redundant features seems to have a greater relative impact on the accuracy than employing domain dependent MFCC features. This is evident if we compare the results of using model 1 with the model discussed in 4.1.1. The presence of silence regions force the optimal warping paths between patterns of the same lexical identity to lie outside the regions in the cost matrix bounded by the Sakoe-Chuba band. This forces the DTW to find an **approximate** optimal path within the allowed regions which in turn leads to an increase in the error rate.
- Combining attribute/feature selection with MFCC feature extraction as a preprocessing step achieves greater improvement in accuracy and speed than using either of these approaches alone. In comparison to just using MFCC feature extraction as a preprocessing step, the accuracy has been boosted up by 15.17% on average while the $\ln(\text{time})$ have been reduced by 2.36. Similarly, in comparison to just using feature selection as a preprocessing step, the accuracy has increased by 12.9% on average while the $\ln(\text{time})$ have been reduced by 6.5.
- From the above observation alone, we can deduce two facts for this data set, one of which is not that obvious: removing redundant features have a greater relative impact on improving the performance of a constrained DTW than employing domain dependent features and the second fact which is quite obvious is the strong dependency between the algorithm's accuracy and the size of the window constraint. It can be concluded that

the use of the rigid window constraint is the main contributing factor for achieving small accuracy rates.

Since the main goal is to improve **both** the accuracy and speed of DTW in handling sequences of high-dimensionality i.e long lengths, removal of silence segments provides an ideal mechanism to improve **both** the time complexity and the accuracy of the algorithm. Dropping the window constraint although will lead to considerable increase in the accuracy by will result in longer run-times. This is not ideal when working on large high dimensional datasets.

- Model 3 achieves almost near perfect accuracy. Dropping the window constraint improves the accuracy by 67.15% on average over model 2. In other words 67.15% of the time on average, the optimum warping path lie outside the regions bounded by the Sakoe-Chuba band constraints, This confirms that patterns belonging to the same lexical identity can have an overly temporal askew between them as result of the different contexts in which the words are spoken and/or as a result of different speakers speaking the same word.
- The run time of model 3 is lower than the run time of all models with the exception of model 2. Model 3 therefore provides the best balance in the fulfilling the two conflicting objectives of high accuracy and high speed that any of the other investigated models.

Conclusion

Although the DTW algorithm is domain independent, from the analysis of the above experiments, it has been observed that augmenting the DTW with a domain dependent preprocessing technique can greatly improve its performance in terms of **both** accuracy and speed with out the need of any global constraints. For example, for the TIDIGITS data set, employing the preprocessing steps of 'silence' removal followed by MFCC feature extraction allows DTW to improve not only its accuracy but also speed even without the use of global window constraints.

Chapter 5

Extending DTW

So far, we have investigated methodologies that integrate meta data (i.e knowledge of the domain) in the pre-processing stage. The problem with such methodologies is that the same algorithm cannot be extended across multiple domains since the feature extraction process is highly domain-dependent. The MFCC feature vectors, for example, that we considered in the previous section can only be employed for data sets that comprise of speech utterances. In the first half of this chapter, I investigate feature extraction methodologies that are entirely data driven so that we can construct a methodology for improving DTW's speed and accuracy that can be extended across multiple domains. In the second half of this chapter, I investigate alternative measures to using window constraints that can improve the performance of the algorithm terms of **both** time and accuracy across all time series domains.

5.1 Domain-independent feature extraction

Ideally, we require features that reflect information about the structure of the data. This allows the DTW to build a complete picture of the data point in relation to the rest of the sequence and hence achieve better optimal alignments between similar sequences. The fundamental problem of baseline (value-based) DTW is that the numerical value of a data point in a time series sequence is not a complete picture of the data point in relation to the rest of the sequence. The context such as the position of the points in relation to their neighbours

is ignored. To fix this issue, an alternative form of DTW known as *derivative* DTW is proposed but it too fails to achieve better performance across all domains as it ignores to take into account the common sub-patterns between two sequences (mainly global trends).

For feature extraction, the methodology that I have used for this setup is based on Xie and Wiltgen's paper[2]. In their paper, the authors highlight a domain-independent feature extraction process where each point in the time series sequence is replaced by a 4 dimensional vector. In this vector, the first two features correspond to information regarding the local trends around a point and the last two features reflect the position of that point with respect to the global shape of the sequence. From experiments conducted on the UCR data sets, they have observed that embedding DTW with this feature extraction process yields greater accuracy across all datasets.

Definition of local feature given in [2] is as follows:

$$f_{\text{local}}(r_i) = (r_i - r_{i-1}, r_i - r_{i+1})$$

The first feature reflects the difference between the values of the current index and the previous index while the second feature reflects the difference between the values in the current index and the succeeding index.

The extraction of global features however, is constrained by two factors: the features must reflect information about global trends and must be in the same scaling order as the local features. Being in the same scale allows them to be combined with local features. In [2], the authors used the following method to extract global features from the time series sequence:

$$f_{\text{global}}(r_i) = (r_i - \sum_{k=1}^{i-1} \frac{r_k}{i-1}, r_i - \sum_{k=i+1}^M \frac{r_k}{M-i})$$

The first feature represents the deviation of the value of the current index from the mean of the values of the sequence that has been seen so far while the second feature represents the deviation of the current value from the mean of the values that is yet to be seen. This formulation allows the detection of significant 'drops' or 'rises' in the series.

Note : The local and global features have no definition for the first and last

points in a sequence. To keep the terminology clear, I refer to the length of the time series as the dimension of the time series. Each dimension i.e each point in the time series can be a value or in this case a 4-d vector.

When working with high dimensional time series (i.e sequences with long lengths) data, the main drawback of employing this feature extraction method is that it does not offer the advantage of dimensionality reduction. The dimensionality of a transformed time series sequence is just two dimensions less than the dimensionality of the original sequence. The DTW algorithm combined with this feature extraction process therefore suffers from the curse of dimensionality as before. To address this issue, the DTW algorithm is subjected to the adaptive Sakoe-Chuba band window constraint (4.1.1) that reduces the search space by restricting the algorithm to look for optimal paths only through limited cells in the DTW cost matrix.

Xie and Witgen[2] have already shown that augmenting this feature extraction methodology to the DTW algorithm does allow the algorithm to achieve better accuracy on datasets from different domains. However, due to the availability of sufficient computing power, they didn't use any window constraints when performing their experiments. For problem scenarios where the speed of the DTW is considered a priority, it will be interesting to investigate whether this methodology can allow DTW to achieve better performance in terms of accuracy over the base line method when subjected to the window constraint.

To investigate the effect of introducing this prior feature selection step on the performance of the DTW algorithm that employs a rigid window constraint, I conducted the following 2 experiment:

Objective Compare the affect of using global and local features to using raw values on the performance of the DTW subjected to a window constraint

- Experiment 1

Datasets: TIDIGITS data set(4.1.1)

When conducting experiments using MFCC feature vectors, it was observed that removing silence segments is a more prominent contributing factor in improving a window constrained DTW's performance than employing the MFCC feature extraction process that integrates meta data

i.e knowledge of domain into the algorithm. to ensure that this finding is conclusive for this dataset, apart from inspecting whether using the new domain-independent features are better than using raw values, I also measure their relative importance in improving the performance of the DTW against the ‘silence’ removal phase.

RESULTS

A summary of the results are given below:

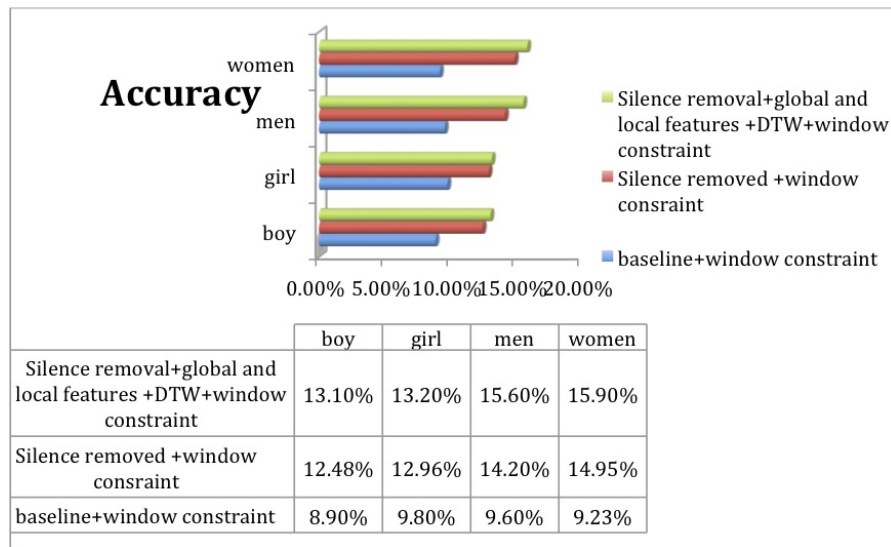


Figure 5.1

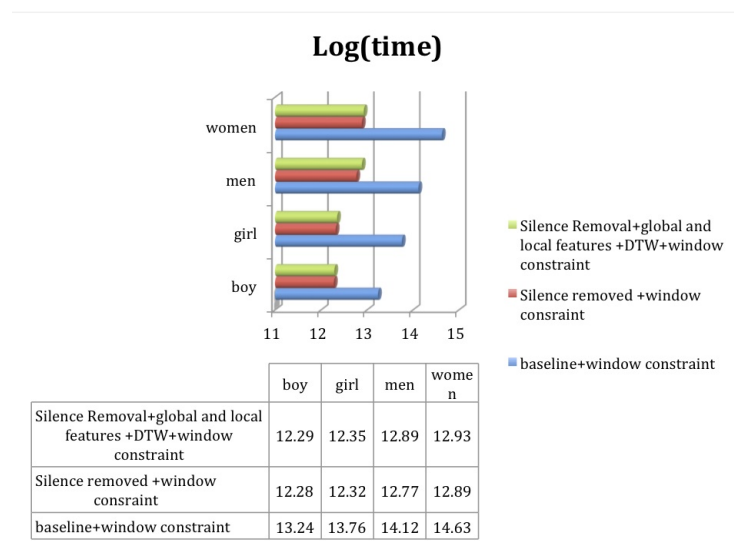


Figure 5.2

Observation:

Surprisingly, by comparing these results with the previous experiments, it can now be concluded that for the TIDIGITS dataset(2.1), removing redundant features i.e removing regions of silence is a greater contributing factor in improving the performance of a window constrained DTW than applying a feature extraction process that integrates meta data(4.2.1) or that captures information about local and global trends. Using local and global features only leads to an average improvement of 0.8% over the model that performs only silence removal as a preprocessing step.

One obvious observation is the poor performance shown by all 3 models in terms of accuracy. From the analysis conducted so far, the reason for this poor performance can be narrowed down to the use of the rigid window constraint imposed to minimise the time complexity. The computational cost incurred by the algorithm is higher than the version used for the model of 4.1.1 One possible explanation is that the cost of applying the euclidean metric on vectors > cost of applying the euclidean metric on points. Since the euclidean metric is applied mn times. The overall computational cost increases.

- Experiment 2 Datasets: UCR datasets: InlineSkate and Cinc_ECG_Torso

The results are as follows:

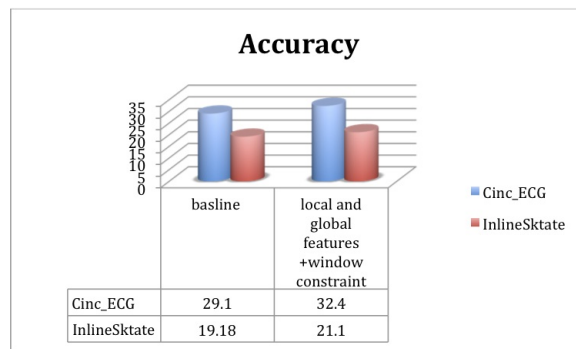


Figure 5.3: Using features that reflect information of trends improves the accuracy of the algorithm

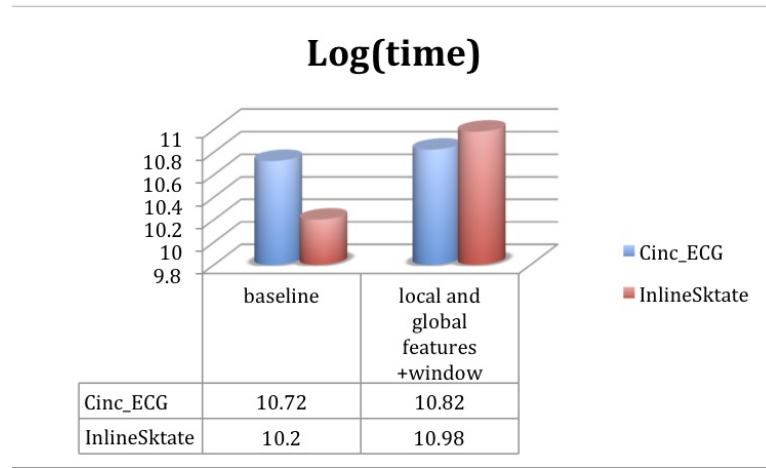


Figure 5.4: The computation time incurred

Observation

- The differences between performances of the two versions of DTW are consistent with the observations made in the previous experiment. For the Cinc_ECG_Torso time series data set, using global and local features improves the accuracy by 3.2% whereas for the InlineSkate data set, the accuracy improves by 1.92%. Thus it is safe to conclude that replacing each point in the time series sequence with a vector that reflect the relation of the point with the respect to the trends increase the number of the optimal warping paths that lie with in the regions bounded by the Sakoe-Chuba band constraint. The time complexity incurred by the algorithm under the new model is worse in comparison to the baseline model. Thus although we achieve an improvement in accuracy, we are loosing performance in speed. The accuracy of both the baseline and the current model is very low. This shows that a majority of the optimal warping between sequences belonging to the same class lie outside the Sakoe-Chuba band.

From the results of the experiments that has been conducted so far, it can be observed that the prominent factor responsible for the low accuracy of the DTW is the window constraint that keeps points(or vectors) of one sequence from getting too far from the other. Increasing the width of the Sakoe-Chuba band will although increase the accuracy of the DTW (as seen in 4.2) but will definitely cause an reduction in speed. One of the primary goals of this

project is to improve the performance of the DTW in handling sequences of high dimensionality(i.e log lengths). Thus **minimising** time complexity is as important as improving the accuracy. This provides the motivation to investigate alternative methods that can a better balance between the two conflicting goals of accuracy and speed.

In the next section, I investigate a self-proposed method that is aimed to help DTW tackle the two conflicting objectives more effectively than using rigid window constraints. The new proposed methodology uses the data-driven feature extraction process that is discussed in the current section. The aim here to use the advantages presented by this domain independent feature extraction process [2] without being subject to the model's drawbacks.

5.2 Adapting DTW

The feature extraction methodology discussed above maps the time series sequence to a time series sequence of vectors whose length is $\|X_n\| - 2$. (where $\|X_n\|$ denotes the length of the original time series sequence). The DTW augmented with these features will still suffer from large time and computational complexity if the dimensionality of the data is high. In the MFCC feature extraction process, the time series sequence is first segmented into series of frames of length 20ms i.e 200 points. Through appropriate functional mapping, each frame is then mapped to a vector. Because the length of the resultant sequence of vectors is much smaller than the length of the original time series, the size of the DTW cost matrix. This reduces the search space and thus decrease the time complexity of the DTW algorithm.

Using the MFCC feature extraction method as motivation, in the proposed model the sequence of 4d vectors extracted using the feature extraction process discussed in 5.1 are segmented using windows of size 50 which in the case of the speech data corresponds to width of 5 ms . The original time series is reduced to series of matrices where the columns of the matrices consist of 4-d feature vectors corresponding to a particular time slice. The length of the series is now 50 times smaller than before. Now if we adapt the cost function of DTW to work on series of frames rather than series of vectors as before

we can achieve a large improvement in both accuracy and computational cost associated than imposing a **window** constraint.

The problem now can be shifted to finding an appropriate kernel that can be used to compute the similarity between matrices composed of feature vectors. Ideally, we want a metric that takes into account the variation of speed and time when comparing two similar subsequences. We will want to compare the global and local properties associated with a point in one subsequence with the global and local properties of points at different regions in the second sub-sequence illustrated by figure 2. Using a euclidean metric in this scenario is inappropriate. The euclidean metric in this context is identical to linear time warping where the two subsequences will be matched based on a linear match of the two temporal dimensions. In our context, we need a kernel that computes the similarity between two sub-sequences by warping the time axis.

The motivation behind the kernel that I propose for aiding DTW to tackle high-dimensional sequences(i.e sequences with long lengths) comes from the polynomial kernel.

Let x and z be two dimensional vectors. Consider the simple polynomial kernel of degree 2 : $k(x, z) = (x^T z)^2$. This kernel can expressed as :

$$\begin{aligned}
 k(x, z) &= (x^T z)^2 \\
 &= (x_1 z_1 + x_2 z_2)^2 \\
 &= x_1^2 z_1^2 + 2x_1 z_1 x_2 z_2 + x_2^2 z_2^2 \\
 &= (x_1^2, 2x_1 x_2, x_2^2)(z_1^2, 2z_1 z_2, z_2^2)^T \\
 &= \phi(x)^T \phi(z)
 \end{aligned}$$

The 2nd order polynomial kernel is equivalent to a corresponding feature mapping $\phi(x)$ that maps a two dimensional vector to $(x_1^2, 2x_1 x_2, x_2^2)$ where each attribute is monomial of order 2 . Generalising this notion to order M then $k(x, z) = (x^T z)^M$ contains all monomials of order M. Now, if we imagine x and z to be two images, then the polynomial kernel represents a particular weighted sum of all possible products of M pixels in the first image with M pixels in the second image.

Using this as motivation I propose the following kernel:

$$k(x, z) = \left\langle \sum_{i=1}^n x_i, \sum_{j=1}^n z_j \right\rangle$$

where n denotes the length of the window and x_i and z_j represents the 4-dimensional features indexed by the points in two sub-sequences.

To motivate the reasoning behind the construction of this particular kernel lets consider the following signals:

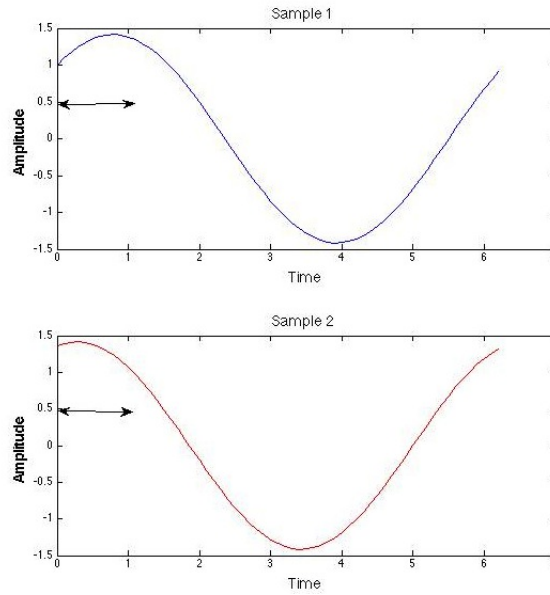


Figure 5.5: Two signals separated by translation

The signal denoted by the 'red' color is a 'slower' version of the signal denoted by the 'blue' color. In the above example, if we are comparing the similarity between the time slices spanned by the arrows, an ideal kernel must be invariant to the time offsets of the signals and thus should consider all possible pairings between the vectors in the subsequences. Intuitively speaking, the kernel must behave like a DTW algorithm.

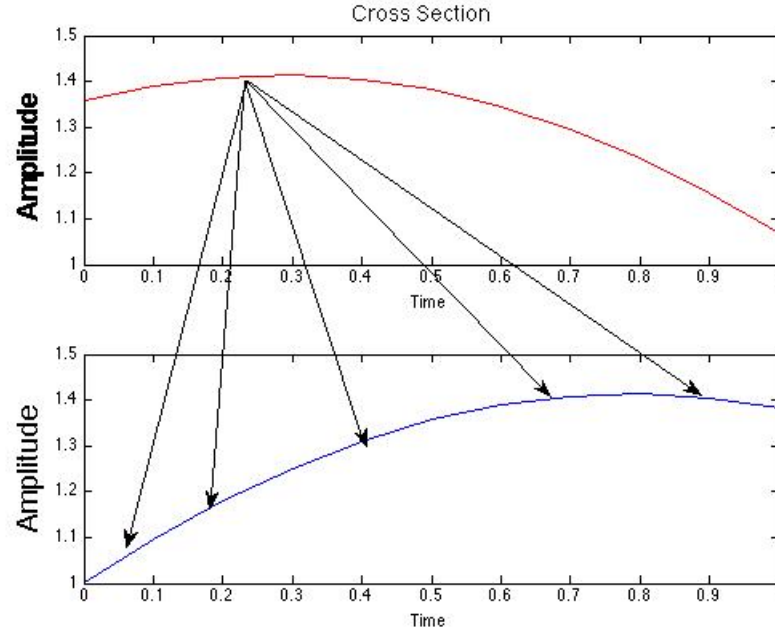


Figure 5.6: Two identical subsequences varying in time

For time slices of width n , the kernel metric can be expanded and expressed as :

$$\begin{aligned}
 k(x, z) &= \left\langle \sum_{i=1}^n x_i, \sum_{j=1}^n z_j \right\rangle \\
 &= \langle (x_1 + x_2 + x_3 + \dots), (z_1 + z_2 + z_3 + \dots) \rangle \\
 &= \langle x_1, z_1 \rangle + \langle x_1, z_2 \rangle + \langle x_1, z_3 \rangle + \dots + \langle x_2, z_1 \rangle + \langle x_2, z_2 \rangle + \langle x_2, z_3 \rangle + \dots
 \end{aligned}$$

From above expression, we can see that the proposed kernel corresponds to a sum of all possible dot products of pairs belonging to the set $\{(x_i, z_j) | x_i \in \text{seq1}, z_j \in \text{seq2}\}$. Similar to the polynomial kernel, the proposed kernel allows us to match all possible pairs of vectors belonging to the two sub-sequences given by the matrices. It is easy to check that this proposed kernel is in fact a valid kernel:

- $K(x, z) = K(z, x) \Rightarrow$ the function is symmetric.
- The kernel satisfies Mercer's theorem : $K(x, z) = \phi(x)^T \phi(z)$ where the feature mapping corresponds to a finite summation of vectors $\phi(y) = \sum_{i=1}^n y_i$.

Augmenting the kernel to the DTW algorithm allows DTW to work on high-

dimensional time sequences without using a window constraint. However the accuracy and computational cost of the DTW is now dependent on the size of the time slices used to segment the original sequences in the first place. To use this kernel as an appropriate cost function in the DTW algorithm, we need a functional mapping that:

1. constraints the codomain to be in the range from 0 to ∞ .
2. ensures larger values given by the function signify great degree of dissimilarity and smaller values signify a high degree of similitude.

An ideal cost function that make use of dot products is the *arc-cosine*. Hence I embedded the kernel function in the cosine distance:

$$\theta = \frac{\langle X, Z \rangle}{|X||Z|}$$

where $X = \sum_{i=1}^n x_i$ and $Z = \sum_{j=1}^n z_j$

A formal outline of the algorithm is as follows:

Algorithm 4 Adapted DTW

```

1: procedure VALUE-BASED(seq1, seq2)      ▷ two sequences of feature vectors
2:   seq_1 ← segment(seq1, n)  ▷ Segment the sequences using a window of
   size n
3:   seq_2 ← segment(seq2, n)
4:   for i=1: to length(seq_1) do          ▷ Initialise the DTW cost matrix
5:     DTW(i, 0) =  $\infty$ 
6:   end for
7:   for i=1 to length(seq_2) do
8:     DTW(0, i) =  $\infty$ 
9:   end for
10:  for i=2 to length(seq_1) do
11:    for j=max(2, i-w) to min(length(seq_2), i+w) do
12:      DTW(i, j) =  $\theta = \frac{\langle X, Z \rangle}{|X||Z|} + \min\{ \text{DTW}(i-1, j) + \text{DTW}(i, j-1) + \text{DTW}(i-1, j-1) \}$ 
13:    end for                                ▷  $X = \sum_{i=1}^n x_i$  and  $Z = \sum_{j=1}^n z_j$ 
14:  end for
15:  return result =  $\frac{\text{DTW}(n, m)}{nm}$           ▷ n=length(seq1), m=length(seq2)
16: end procedure

```

5.2.1 Testing the methodology

The changes that I have introduced in the previous section to the 'Dynamic Time Warping' algorithm is aimed to improve the algorithm's performance in handling long time series sequences. In this section, I investigate the performance of my proposed methodology against the versions of the DTW algorithm that employ the Sakoe-Chuba band constraint (discussed in 4.1.1). To reduce the time complexity to a minimum, the window constraint that I have used so far is the most rigid constraint that can be imposed on the DTW algorithm without forcing the algorithm to behave like the euclidean metric. In order to measure the difference in performance between using my proposed method and employing the window constraint, I have conducted the following experiments:

Experiment 1

Data set used : The TIDIGITS test and training data.

Model Used : The model discussed in 5.2. The model framework consists of a two stage preprocessing step that involves feature selection process consisting of 'silence' removal followed the domain independent feature extraction discussed in the previous section.

Variables being compared: employing window constraints vs the new proposed changes

RESULTS: the results are as follows:

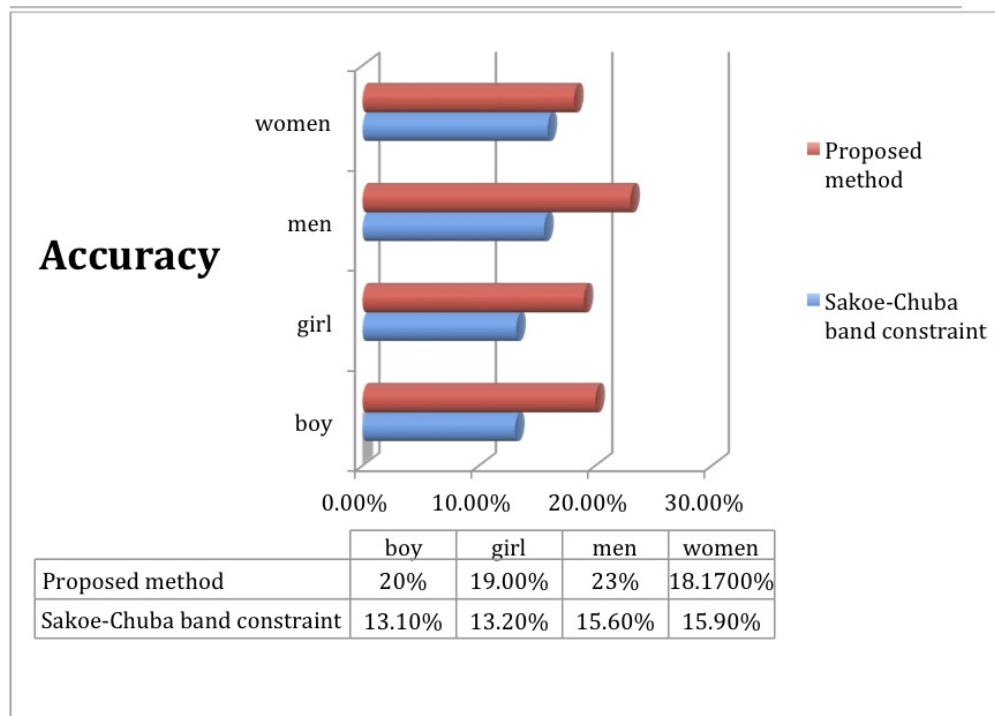


Figure 5.7: Accuracy

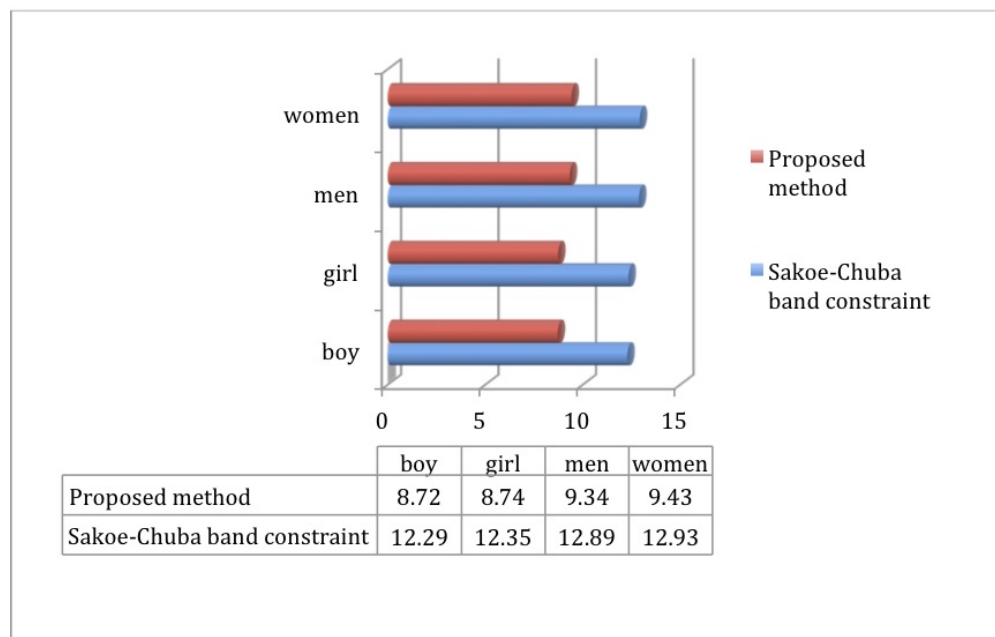


Figure 5.8: Time complexity in log(time)

There are quite number of interesting observations that can made from the graphs and the tables given by figures 5.7 and 5.8.

- The proposed changes to DTW allow the algorithm to achieve better ac-

curacy on test samples across all categories than employing the rigid Sakoe-Chuba band constraint of $w = \max(\lceil 0.1 * \max(n, m) \rceil, \text{abs}(n - m))$. The most interesting result is that the new algorithm incurs a lower computational cost than before. Thus introducing these changes have improved both the accuracy and time complexity of the algorithm. From the results noted in the tables, it can be seen that the accuracy of DTW has increased by 6.54% on average and the average log(time) has decreased by 3.1. The reduction of the time complexity is mainly due to the partitioning of the sequence into time slices of width 5 ms. The reduction in the length of the sequences by an order of 50 results in the shrinkage of the search space of DTW thus causing the algorithm to improve its speed. As we have seen so far, increasing the speed of DTW negatively impact the accuracy, in this scenario, the new methodology actually provides an exception. The use of the kernel function improves the accuracy of the DTW which implies that matching frames using the new cost function is a better alternative than employing the euclidean distance between points/vectors confined by the window constraint.

In the previous chapter, we have seen that for the TIDIGITS data set, constructing a preprocessing methodology that involves ‘silence’ removal followed by MFCC feature extraction allows DTW to minimise its time complexity and at the same achieve high accuracy without the use of the global window constraint. However for situations where the need for a global window constraint is deemed necessary, it will be interesting to investigate whether the new changes do in fact provide a better alternative to using rigid window constraint for very long sequences. To check that the new changes allow DTW to perform equally well using domain dependent features, I have repeated the same experiment again but this time, as a preprocessing step I have just used the MFCC feature extraction

A summary of the results are as follows:

Experiment 2

Data set used : The TIDIGITS test and training data.

Model : The model framework consists of single preprocessing step that in-

volves the extraction of MFCC features

Variables being compared: employing window constraints vs the new proposed changes

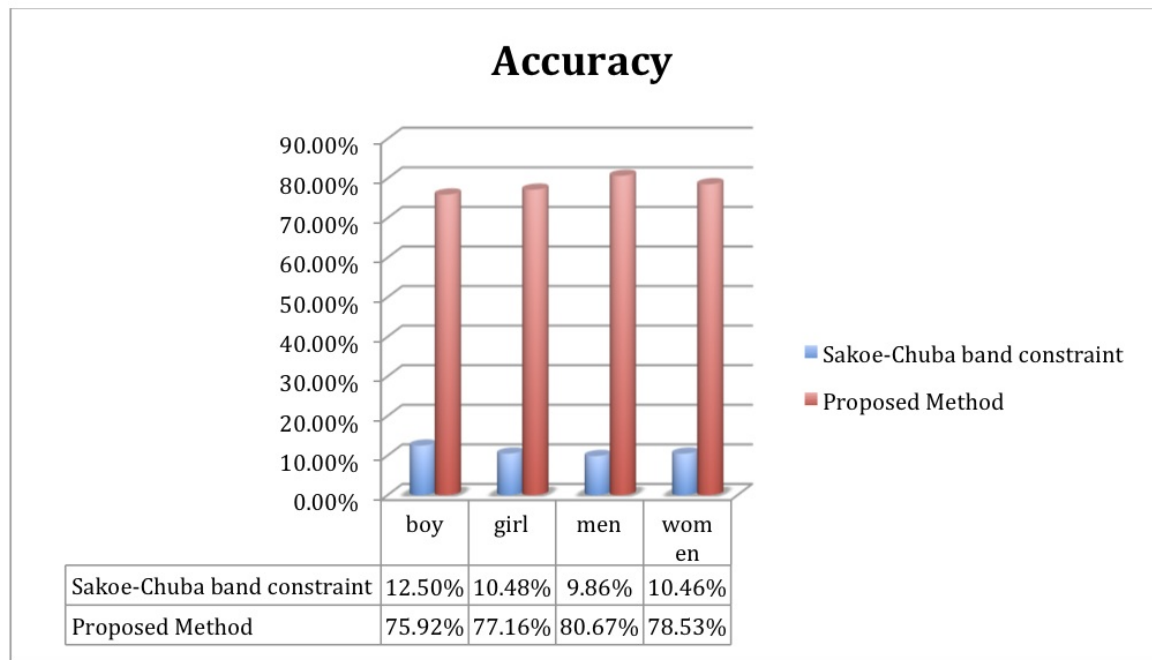


Figure 5.9: Accuracy

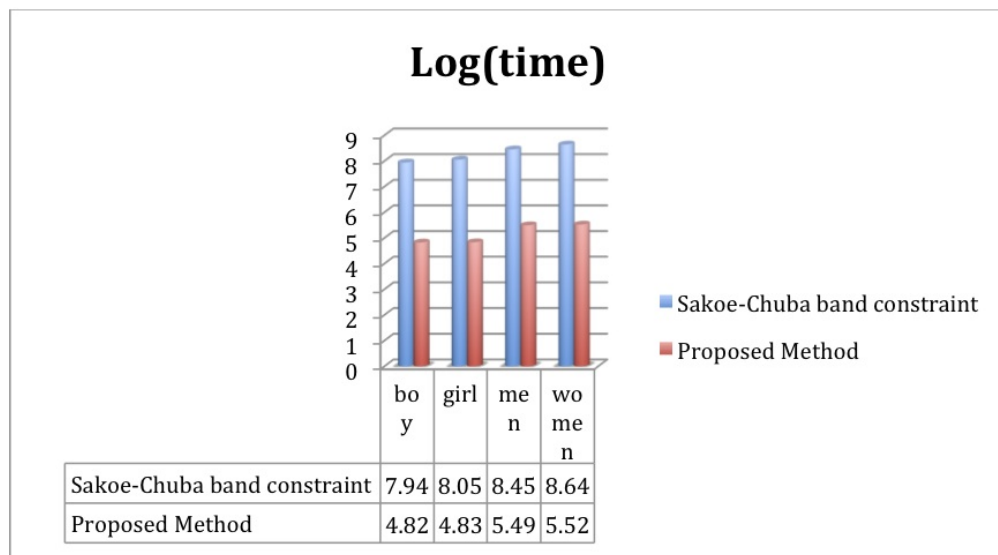


Figure 5.10: Time complexity in log(time)

Observation

From the above results, it can be observed that the new changes does provide an a **better** alternative approach to using a very rigid constraint for scenarios where the time complexity is of high priority and the use of window constraints can not be avoided. For this data set, using the new approach improves the accuracy of the DTW by over 60% and the $\log(\text{time})$ is reduced by order. This shows that that not only the time complexity is exponentially reduced but the accuracy has also been boosted. However, since the tests so far has been conducted on the TIDIGITS test set, it is possible that this new approach is only tailored for this particular data set. Thus to confirm that the performance of the new algorithm is not tailored for this particular time-series data set, I have the test the performance of the DTW augmented with the new changes on the UCR data sets: InlineSkate and CINC_ECG_TORSO.

Experiment 3

Data sets used : InlineSkate and CINC_ECG_TORSO

Setup : One of the main goals of this analysis is to construct a methodology that can achieve good performance in both speed and accuracy across multiple domains. The model framework that I am using for this experiment involves the use of the feature extraction process discussed in 5.1 and [2]. From the experiments conducted on the UCR datasets, Xie and Witgen[2] have shown at the cost of higher computational complexity m the accuracy of the DTW algorithm does improve by a significant amount when the algorithm substitutes each raw value with a 4-d vector that reflects information about local and global trends. However, from the analysis I conducted in 5.1, it was observed that under the window constraint, the use of these features made little contribution in improving the performance of the algorithm. If replacing the window constraint with the new changes does improve the performance of the DTW algorithm then we will have a model where the use of global and local features does allow the algorithm to achieve greater performance even when the decreasing the time complexity is a major priority.

In the proposed approach, the width of the time slices has been kept fixed at a default value of 50 so far. In this experiment, I also investigate the influence of this parameter on the performance of the DTW. Decreasing its value reduces the size of the time slices which in principal should increase both accuracy

and time-complexity . The core kernel used by the new algorithm is based on the function:

$$k(x,z) = < \sum_{i=1}^n x_i, \sum_{j=1}^n z_j >$$

$k(x,z)$ represents the sum of all possible dot-products. Using smaller subsequences allow the similarity measure to be dominated by the dot products of points whose local and global features are most alike. However, this suffers from the drawback of achieving lesser dimensionality reduction. Thus the time and computational complexity suffers.

RESULTS

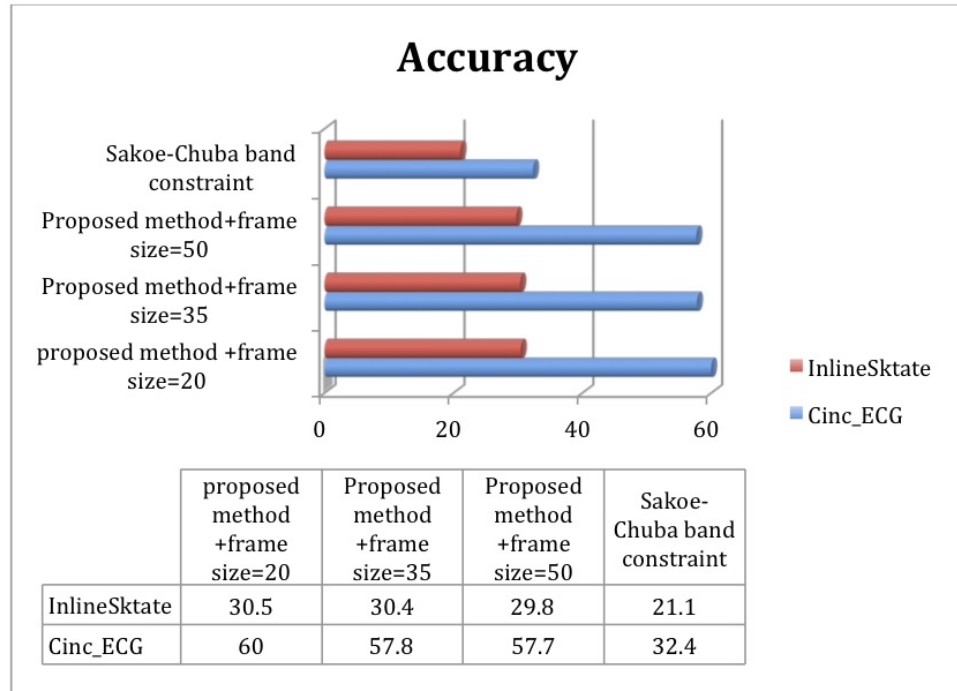


Figure 5.11: Accuracy

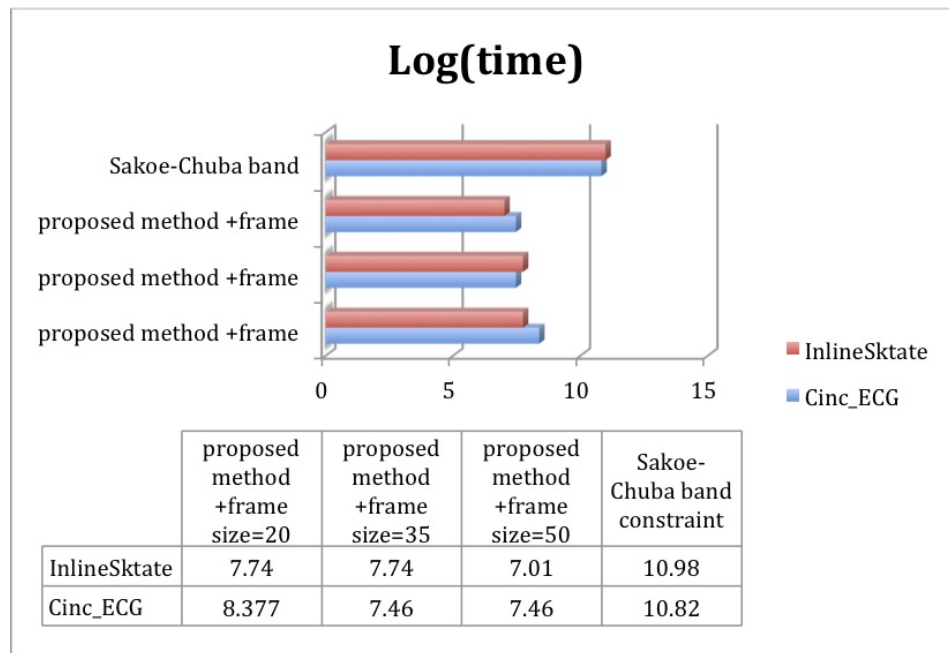


Figure 5.12: Accuracy

Observation

- The proposed changes does provide a **better** alternative approach to using rigid constraints for problem domains where the time complexity is of high priority and the use of window constraints can not be avoided. The accuracy of the DTW algorithm under the new methodology is significantly higher then employing a window constrained DTW on sequences of extracted local and global features. In my analysis in 5.1, I have already shown that the latter model achieves better accuracy than the baseline model. Thus the new methodology improves the performance of the DTW in both speed and time and also can be **applied** to time series sequences belonging to different domains. Hence, combining this approach with the feature extraction method discussed in [2] results in the construction of a model that can applied for **different types** of time series data sets.
- Decreasing the size of the time slices only leads to a minimal increase in accuracy. Thus using the default value of 50 seems to be safe option as the algorithm better accuracy and speed than using the rigid Sakoe-Chuba band constraint .

Bibliography

- [1] R. Gautam Das, King ip Lin, Mannila, H., "Rule Discovery From Time Series," *4th Int'l Conference on Knowledge Discovery and Data*, 1998.
- [2] U. Fayyad, C. Reina, and P. S. Bradley, "Initialization of Iterative Refinement Clustering Algorithms," 1998.
- [3] A. Bazma, I. JONASSEN, I. EIDHAMMER, and D. GILBERT, "Approaches to the Automatic Discovery of Patterns in Biosequences," *Journal of Computational Biology*, vol. 5, pp. 279–305, Jan. 1998.
- [4] A. S. Park and J. R. Glass, "Unsupervised Pattern Discovery in Speech," *IEEE Transactions on Audio, Speech, and Language Processing*, vol. 16, pp. 186–197, Jan. 2008.
- [5] Y. Zhang and J. R. Glass, "Towards multi-speaker unsupervised speech pattern discovery," in *2010 IEEE International Conference on Acoustics, Speech and Signal Processing*, pp. 4366–4369, IEEE, 2010.
- [6] S. Salvador and P. Chan, "FastDTW : Toward Accurate Dynamic Time Warping in Linear Time and Space,"
- [7] H. Sakoe and S. Chiba, "Dynamic programming algorithm optimization for spoken word recognition," pp. 159–165, May 1990.
- [8] L. Rabiner, A. Rosenberg, and S. Levinson, "Considerations in dynamic time warping algorithms for discrete word recognition," *IEEE Transactions on Acoustics, Speech, and Signal Processing*, vol. 26, pp. 575–582, Dec. 1978.
- [9] Y. Xie and B. Wiltgen, "Adaptive Feature Based Dynamic Time Warping," vol. 10, no. 1, pp. 264–273, 2010.

- [10] X. Xi, E. Keogh, C. Shelton, L. Wei, and C. A. Ratanamahatana, "Fast time series classification using numerosity reduction," in *Proceedings of the 23rd international conference on Machine learning - ICML '06*, (New York, New York, USA), pp. 1033–1040, ACM Press, June 2006.
- [11] A. W.-C. Fu, E. Keogh, L. Y. H. Lau, C. A. Ratanamahatana, and R. C.-W. Wong, "Scaling and time warping in time series querying," *The VLDB Journal*, vol. 17, pp. 899–921, Mar. 2007.
- [12] F. Korn, H. V. Jagadish, and C. Faloutsos, "Efficiently supporting ad hoc queries in large datasets of time sequences," in *Proceedings of the 1997 ACM SIGMOD international conference on Management of data - SIGMOD '97*, vol. 26, (New York, New York, USA), pp. 289–300, ACM Press, June 1997.
- [13] F. Itakura, "Minimum prediction residual principle applied to speech recognition," *IEEE Transactions on Acoustics, Speech, and Signal Processing*, vol. 23, pp. 67–72, Feb. 1975.
- [14] E. K. Chotirat Ann Ratanamahatana, "Three Myths about Dynamic Time Warping Data," *Mining, in the Proceedings of SIAM International Conference on Data Mining*, 2005.
- [15] M. De Wachter, M. Matton, K. Demuynck, P. Wambacq, R. Cools, and D. Van Compernelle, "Template-Based Continuous Speech Recognition," *IEEE Transactions on Audio, Speech and Language Processing*, vol. 15, pp. 1377–1390, May 2007.
- [16] M. De Wachter, K. Demuynck, P. Wambacq, and D. Van Compernelle, "A locally weighted distance measure for example based speech recognition," in *2004 IEEE International Conference on Acoustics, Speech, and Signal Processing*, vol. 1, pp. I–181–4, IEEE, 2004.
- [17] M. A. Carlin, S. Thomas, A. Jansen, and H. Hermansky, "Rapid Evaluation of Speech Representations for Spoken Term Discovery,"