

DSP Lab Report Draft

Introduction

Sound is very common in the nature. Animals use sound to communicate with each other, such as birds, elephants and monkeys. Also, humans use sounds to communicate with each other. The sound human use follows particular rules, which is called language. It is nature for humans to communicate with each other using sound. People can understand with each other quickly using sound, or in other words, talking.

For a long time, the way people interact with computers is through input devices such as keyboard and mouse. People need to type the word they want to express into the input text field. As the technology develops, computer is becoming more powerful. The power of modern computer enables us to realize a system to recognize the word human speaks. This system is called speech recognition system. Nowadays there are a lot of speech recognition systems such as Siri, Cortana, Google Assistant etc. You can awake Siri by saying “Hey Siri!” by default. The sound wave you make will first be transformed into digital signal by your phone, and then by applying algorithms, the digital signal will be transformed into words or some commands.

There are many algorithms to implement the speech recognition system, and our paper mainly focuses on a basic but important method, Dynamic Time Warping (DTW). It is the fundamental of many modern speech recognition systems. It compares the input signal with a standard signal pattern to get the most similar word.

Dynamic Time Warping (DTW)

Dynamic Time Warping algorithm can calculate the similarity between two digital signals in time domain. Its powerful strength of Dynamic Time Warping algorithm is that although the two digital signals may vary in time, for example, two people may have different speaking speed, they can still have very high-level similarity. This is the difficulty in comparing two speech signals as everyone has a different speaking speed.

DTW algorithm accomplishes this difficulty by calculating an “optimal warping path” between the two digital signals. Rather than directly compare two signals, DTW will find a best match for the input signal. Figure below shows how “optimal warping path” works. For figure (a), it just compares the corresponding points of the two signals. For figure (b), the input signal will find a best match from the signal pattern to compare. By applying this method, DTW can calculate the similarity between two different signals in time domain.

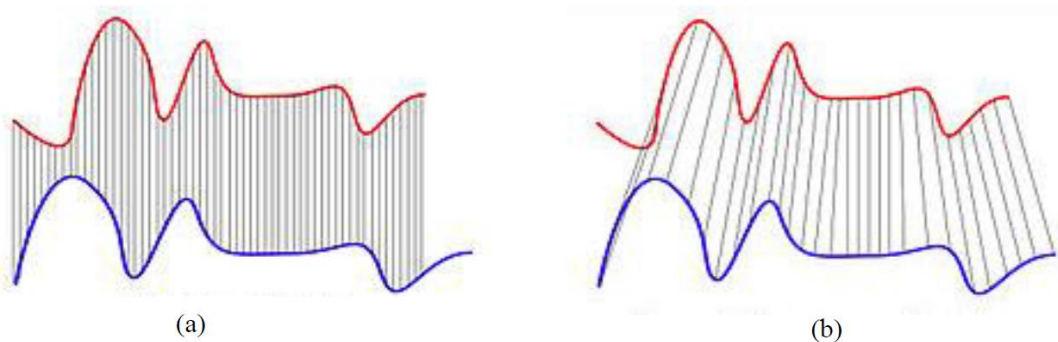


Figure 1. (a) The original alignment of two sequences (b) alligments with DTW

Euclidean Distance

$$Dist(x, y) = |x - y| = [(x_1 - y_1)^2 + (x_2 - y_2)^2 + \dots + (x_n - y_n)^2]^{\frac{1}{2}}$$

Implementation

DTW used a programming method called Dynamic Programming. We will use a two-dimensional array $DTW[m+1][n+1]$. $DTW[m+1][n+1]$ represents the optimal distance

between two input signal sequences lets say $audioIn[m]$ and $sampleAudio[n]$.

$DTW[m+1][n+1]$ can be calculated by the following recursive formula:

$$DTW[m + 1][n + 1] = Dist(m, n) + \min \begin{pmatrix} DTW[m + 1][n] \\ DTW[m][n + 1] \\ DTW[m][n] \end{pmatrix}$$

And here is the pseudo-code:

```
1 function DTWDistance(s[1:n], t[1:m]) {
2     initialize array DTW[1:n][1:m]
3     for i = 0 : n
4         for j = 0 : m
5             DTW[i][j] = infinity
6     DTW[0][0] = 0
7     for i = 1 : n
8         for j = 1 : m
9             min = minimum(DTW[i-1][j],
10                          DTW[i][j-1],
11                          DTW[i-1][j-1])
12             DTW[i][j] = dist(s[i], t[j]) + min
13     return DTW[n][m]
14 }
```

By applying DTW, we can calculate the similarity and output an optimal word.

Future Application

DTW algorithm for now is widely used in speech recognition. Many projects is built on it for example the ESP system (<https://github.com/damellis/ESP>).

In my point of view, I think it can also be used in safety area. It can be used to unlock some system. It can be used to compare the similarity between two sound pattern and return lock or unlock.