# Digit Recognition Using Dynamic Time Warping <span>spate_235</span>

## Project Overview

This project focuses on recognizing spoken digits using the Dynamic Time Warping (DTW) algorithm. DTW is a powerful technique for measuring similarity between two sequences, which may vary in time or speed. The dataset provided includes 300 spoken digit recordings (30 utterances for each digit from three speakers: Jackson, Nicolas, and Theo). Each recording is stored as a matrix with dimensions $M \times 13$, where $M$ denotes the number of frames and 13 represents the feature dimensions per frame.

The primary objective is to take each test utterance and compute its alignment cost with all other recordings to identify the digit being spoken accurately.

## Methodology

Step 1: Dissimilarity Matrix Computation

For each test utterance $T$ and reference utterance $R$, we compute an $M \times N$ dissimilarity matrix $S$, where $M$ and $N$ represent the number of frames in $T$ and $R$, respectively. The elements of $S$ are calculated using the Euclidean distance between corresponding frames of $T$ and $R$.

Step 2: Accumulated Cost Matrix and Optimal Path Calculation

Using the dissimilarity matrix $S$, we construct an accumulated DTW cost matrix $D$. Each element $D(m,n)$ is computed by considering the minimum of three neighboring elements, which represent different potential paths:

- $D(m-1,n)$: Insertion

- $D(m,n-1)$: Deletion

- $D(m-1,n-1)$: Match

The formula used is:

$$D(m,n)=S(m,n)+\min(D(m-1,n),D(m,n-1),D(m-1,n-1))$$

Once $D(M,N)$ is calculated, the backtracking process starts from this position, tracing back to $(0,0)$ to identify the optimal alignment path.

### Step 3: Normalized Alignment Cost

The alignment cost between $T$ and $R$ is determined by summing the dissimilarity values along the optimal path and normalizing it by the path length. This normalized cost reflects how similar the two utterances are, with lower costs indicating higher similarity.

### Step 4: Confusion Matrix Generation

For each test utterance, the alignment costs with all reference utterances are sorted, and the closest 29 reference recordings are selected. The confusion matrix is then computed based on how often the correct digit is identified among these nearest neighbors.

## Results

### DTW Cost Matrix and Optimal Path

The first image below illustrates a sample DTW cost matrix between two recordings, along with the optimal path (highlighted in red). The optimal path represents the lowest-cost alignment between the test and reference utterances, indicating the areas of highest similarity.

### Improved Confusion Matrix

The confusion matrix below demonstrates the model's performance in distinguishing different digits. Each row represents the actual digit spoken, and each column represents the predicted digit. High values along the diagonal indicate successful recognition, while off-diagonal values represent misclassifications.