Contextualized Spoken Word Representations using

Convolutional Autoencoders

Prakamya Mishra (1610110256) Pranav Mathur (1610110257)

INTRODUCTION

- ★ This paper first depicts the process of splitting long audio files into individual words by cutting the silences from the file and calculating the Mahalanobis distance to identify these silences.
- ★ The second and main aim is to generate a deep learning model to visualise audio files in the form of unique vector representations.

RELATED WORK

- ★ Word2Vec has bestowed a prominent impact on the NLP field. With words being represented as vectors that provide context with reference to the nearby words, the vector representations can be used for sentiment analysis
- ★ There have been a few attempts by embedding word segments as fixed length vectors, which also proved to have a great application in Spoken Term Detection (STD).

WORD2VEC

- ★ Word2Vec is one of the most popular techniques to learn word embeddings using shallow neural network.
- ★ It's objective is to have words with similar context occupy close spatial positions in a graph.
- ★ There are two types of models available :
- 1. Common Bag Of Words (CBOW) Model
- 2. Skip-Gram Model:

AUDIO SPLITTING

- ★ Audio Splitting was achieved by using PyDub, a simple python library that can be used for audio manipulation at a high level interface.
- ★ This paper deals with the basic classification into the following three categories -
 - (i) silence (S), where no speech is produced,
 - (ii) unvoiced (U), where the vocal cords are not vibrating
 - (iii) voiced (V), in which the vocal cords get tensed and allow air flow to generate speech.

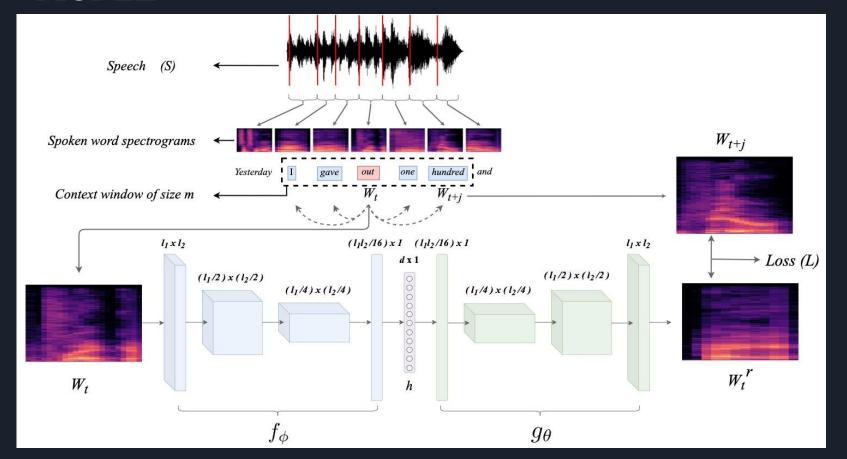
DATASET PREPARATION

- ★ The publicly available Trump speeches (Audio and word transcription) dataset [4] has been used for training the model.
- ★ This dataset contains the complete audio files as well as the JSON files that contain a list of all the words spoken in the speech.
- ★ A spectrogram is a visual representation of the spectrum of frequencies of an audio file as it varies with time. Each audio file was converted to an image spectrogram of size 1004 X 642 pixels.

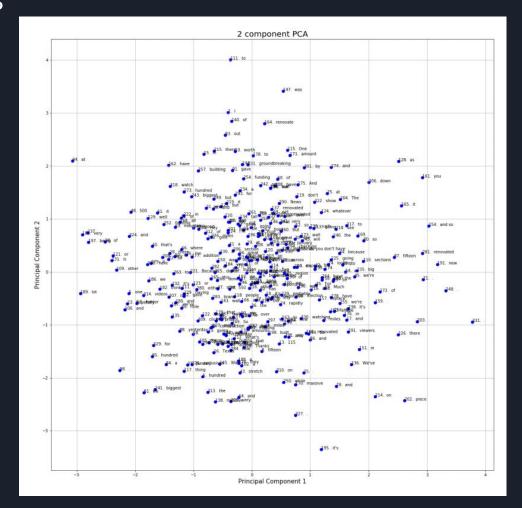
MODEL

- ★ Convolutional Autoencoders have been used to learn a 16 dimensional latent vector representation of each spoken word using the context in which it was spoken.
- **★** This neural network model contains in itself 2 different models.

MODEL



Results



Results

