# Speech Emotion Recognition using Semantic Information

Navya (210657)

## 1   Implementation Details

All the codes given in the original github link by the author were ran other than evaluation. All the training codes(4 files) contained code which had attributes of tensorflow 1.x version and not compatible with tensorflow 2.x version which had to be converted and to work with tensorflow 2.x version and python 3.9

Folders segmentations and tfrecords present in the datasent link were generated by the code by running the first two codes.

For training, all the other codes also have to be copy pasted in the same jupyter notebook before running training as the training files uses those functions.

As aeneas is a python c/c++ library, if you are not able to download manually using pip as it is problematic, you can use the following link to install all its dependencies and library in one quick installer.

aeneas library download link

## 2   Dataset

Dataset provided only had the audio file and transcript file of the participants and did not contain the labels/annotations for emotional(arousal, valence, and liking) data.

The data was generated by me by studying the required format and the necessary fields. The folders(turns and labels) are both created by me, and the labels data contains all zeroes as this was to check whether the code would run.

The data for which labeled data was present had no transcript for as there was no conversation happening in those files and the transcript of conversations was necessary for generating segmentation and training records.

## 3   Results

Running the two codes which generate the word segmentation using the audio and transcript file should give an output of this:

I was only able to train the model for 5 epochs using a small dataset containing only 4 files as otherwise my laptop was showing resource exhausted error/out of memory error.

As all the labels data was zero and not the actual data, the model did not show any improvement during training but rather the same loss was encountered in all five epochs.

(a) speech2word mappings



(b) generating tfrecords of the data

Figure 1: output



(a) Resource Exhausted error



(b) Training output

Figure 2: output