

Speech Emotion Recognition

Utilizing Speech Emotion Recognition to Enhance Human-Computer Interaction



Table of contents

01

Introduction/
Problem
Statement

02

Use Cases

03

Working
Principle

04

Results

05

Lessons
Learned

06

QnA



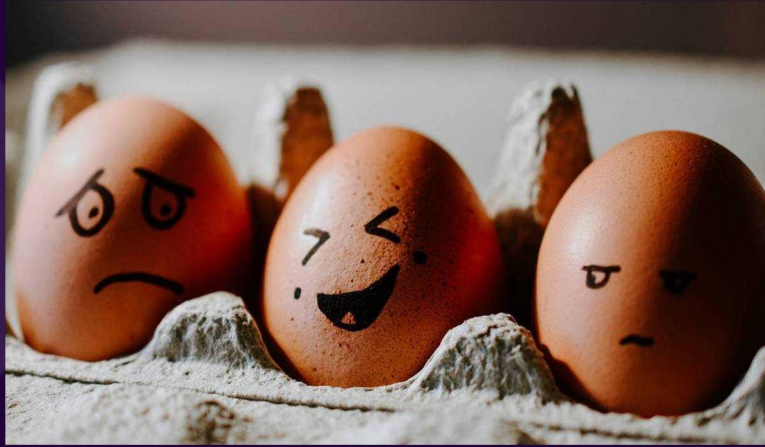


01

Introduction

- Human speech carries complex information beyond the spoken words.
- Utilizing those non-linguistic voice characteristics such as pitch and tone to identify underlying emotions.





Problem Statement

- Near-Real-time prediction of emotions from audio input.
- Aiming for accurate emotion labels based on audio signal analysis.
- The increasing need is driven by the rising demand for advanced communication technologies.
- We are classifying 8 different emotions: Surprise, Angry, Calm, Disgust, Sad, Fear, Neutral, Happy.



Use Cases

Human-Computer Interaction

Improve user engagement and satisfaction by responding empathetically to emotional cues.

Customer Service

Analyze customer calls to identify emotional states and provide appropriate support.

Mental Health Monitoring

Detect signs of distress or emotional dysregulation in speech patterns.

Educational Technology

Assess student engagement and understanding by detecting emotions during lectures or discussions.



Working Principle

1. Gathering Data

RAVDESS, TESS, CREMA-D

2. Pre-Processing

Data augmentation involves creating synthetic data samples by introducing small perturbations to the original training set.

3. Feature Extraction

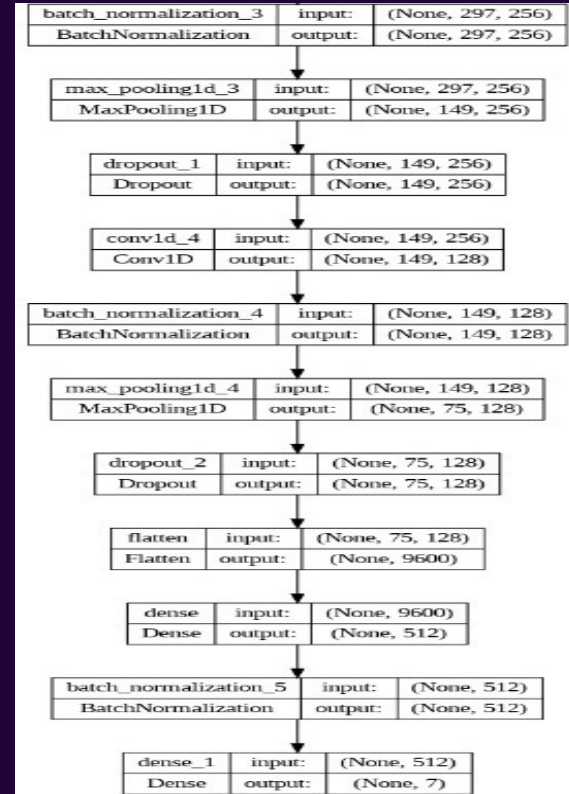
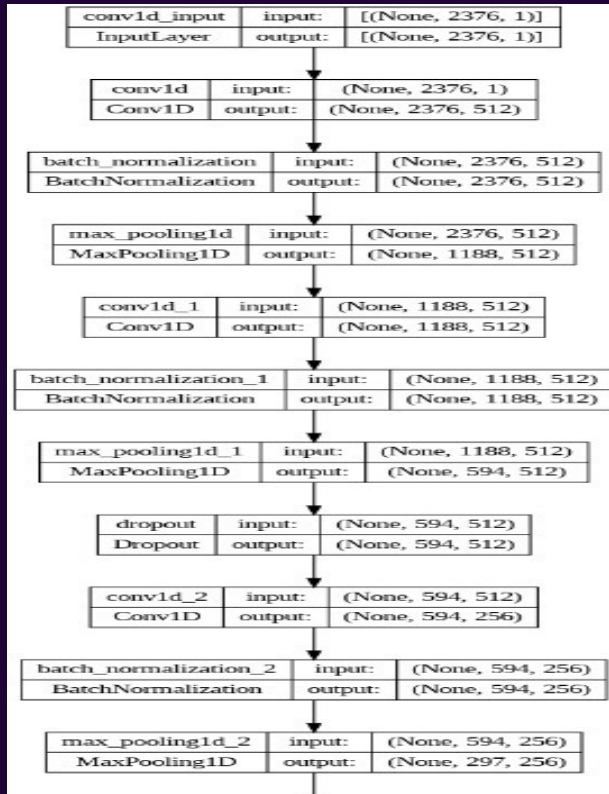
Features such as Zero-Crossing Rate (ZCR), Root Mean Square Error (RMSE), and Mel-frequency Cepstral Coefficients (MFCC) are extracted from the recorded audio using Librosa

4. Data processing and Encoding

Feature vectors are processed, handling missing values, standardized, and emotion labels are one-hot encoded for model training.



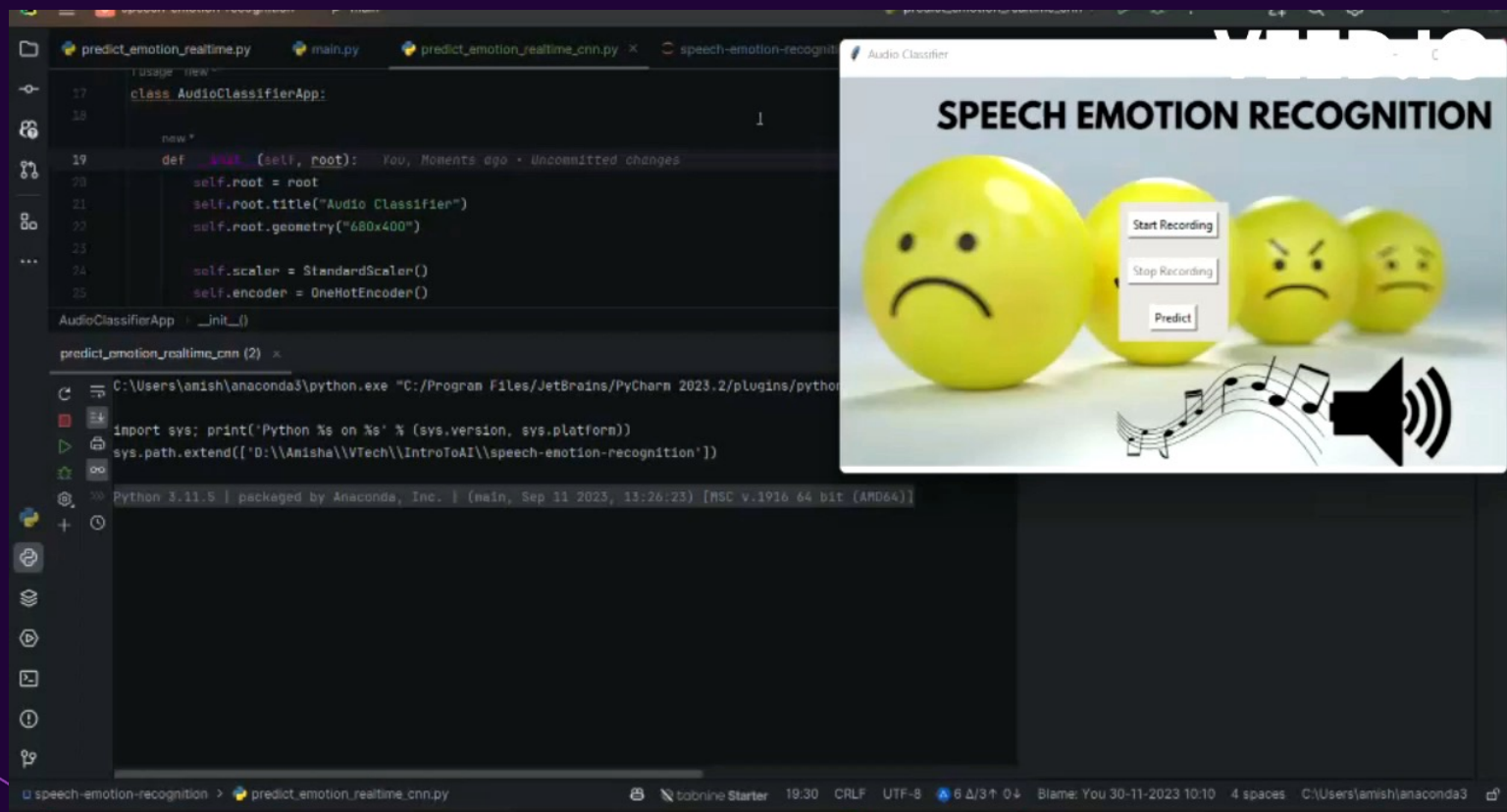
Build CNN Model



Results



Demo



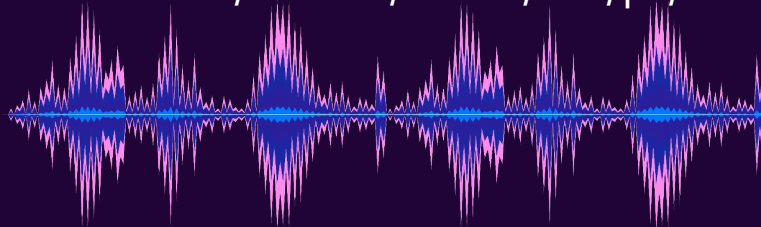
Lessons Learned

- Selection audio features are crucial for emotion detection through speech.
- Training the model with good quality data can help in improving accuracy drastically.
- Tweaking the Epochs to optimum level plays an important role in improving accuracy.



References

- [1]<https://www.kaggle.com/datasets/uwrfkaggler/ravdess-emotional-speech-audio>
- [2]<https://www.kaggle.com/datasets/ejlok1/toronto-emotional-speech-set-tess>
- [3]<https://jonathanbgn.com/speech/2020/10/31/emotion-recognition-transfer-learning-wav2vec.html>
- [4]<https://github.com/CheyneyComputerScience/CREMA-D>
- [5]<https://www.sciencedirect.com/science/article/abs/pii/S1746809420300501>



The background is a deep purple with several large, out-of-focus circles in a slightly lighter shade of purple. Overlaid on this are several glowing, wavy lines in a light purple or magenta color. In the center-left, there is a prominent, bright blue and white waveform that resembles a sound wave or a signal graph, with sharp peaks and valleys.

Thank You

Q&A