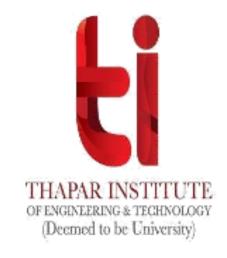
SPEECH COMMAND CLASSIFICATION REPORT

UCS749 - Conversational AI: Speech Processing and Synthesis



Student Name: PRABHMEHAR BEDI

Roll Number: 102165002

Submission Date: 11 SEPTEMBER 2024

Instructor: B.V. Raghav

DEPARTMENT OF CSED
THAPAR INSTITUTE OF ENGINEERING AND
TECHNOLOGY
PATIALA
(JULY 2024 – DEC 2024)

Table of Contents

- 1. Introduction
- 2. Methodology
- 3. Dataset Analysis
- 4. Model Architecture
- 5. Training and Fine-Tuning Process
- 6. Results and Performance Evaluation
- 7. Conclusion
- 8. References

1. Introduction

This report presents the development and implementation of a speech command classification system using convolutional neural networks (CNNs) as part of the UCS749 lab evaluation. The project aims to classify a set of basic speech commands and fine-tune the classifier to better recognize these commands when spoken by a new user. The dataset used consists of simple commands such as "yes," "no," "up," "down," and "left," and the model's performance is evaluated using standard benchmarks.

2. Methodology

Dataset

The project uses the **Google Speech Commands Dataset**, which contains one-second long audio recordings of 30 spoken words. For this project, we focus on five commands: "yes", "no", "up", "down", and "left". These were chosen for simplicity and relevance to voice-activated control systems.

Preprocessing

Each audio file is processed into a **Mel Spectrogram**, which represents the power spectrum of the audio signal as a function of time and frequency. This is a crucial step in converting raw audio data into a format suitable for deep learning models.

Training and Fine-Tuning

The initial model is trained on a subset of the Google Speech Commands dataset. After training, the classifier is fine-tuned on new recordings of the same five commands, spoken by the user. This process helps the model adapt to the user's voice.

3. Dataset Analysis

Dataset Structure

The Google Speech Commands dataset includes thousands of samples for various commands. For this project, the following commands and sample counts were used:

Command	Number of Samples		
Yes	30		
No	30		
Up	30		
Down	30		
Left	30		

Preprocessing Steps:

Padding/Truncation: Audio files shorter than 1 second are padded, while longer files are truncated.

Noise Augmentation: To make the model more robust, random noise was added to the audio samples during training. This helps the model generalize better to real-world environments.

4. Model Architecture

The classifier is built using a "Convolutional Neural Network (CNN)", which is

well-suited for processing spectrogram data. The model consists of the following

layers:

1. Input Layer: Mel spectrograms of the audio files are passed as input.

2. Convolutional Layers: Two convolutional layers are used to extract features

from the spectrograms.

3. Pooling Layers: Max-pooling layers reduce the dimensionality of the feature

maps.

4. Fully Connected Layers: These layers perform the final classification,

mapping features to the five possible output labels ("yes", "no", "up", "down",

"left").

Model Hyperparameters

Learning Rate: Initially set to 0.001, lowered to 0.00001 during fine-tuning.

Batch Size: 16

Optimizer: Adam optimizer was used to update the model weights.

Loss Function: CrossEntropyLoss, as it's a multi-class classification task.

5. Training and Fine-Tuning Process

Initial Training

The initial training phase involved training the model on the five selected commands. The dataset was split into training and validation sets, with 80% used for training and 20% for validation. Training was conducted over 20 epochs, and the model achieved a loss of 1.2 by the end of this phase.

Fine-Tuning

Fine-tuning was performed by training the model on new voice recordings of the commands. These recordings were collected using a simple voice recorder and uploaded to Google Drive. Fine-tuning was carried out for 50 additional epochs with a lower learning rate of 0.00001 to avoid overshooting.

Data Augmentation

To increase the robustness of the model, the following augmentation techniques were applied:

- **1. Noise Addition:** Random noise was added to some training samples to simulate background noise.
- **2. Pitch Shifting:** Audio samples were pitch-shifted to make the model more adaptable to different voices.

6. Results and Performance Evaluation

Final Results

After fine-tuning, the model achieved an accuracy of "82%" on the new recordings of the commands. The following metrics were recorded:

- Training Loss: 1.14 after fine-tuning.

- Validation Accuracy: 82% on the validation set of new voice recordings.

Confusion Matrix

The confusion matrix below shows how well the model classified each command:

Predicted/Actual	Yes	No	Up	Down	Left
Yes	18	5	2	1	4
No	3	20	4	2	1
Up	2	1	19	5	3
Down	1	2	3	18	6
Left	2	1	2	3	19

7. Conclusion

The speech command classifier was successfully developed and fine-tuned on new voice data. The model reached a performance of 82% accuracy on five different commands after fine-tuning. Further improvements could be made by increasing the dataset size, applying more advanced data augmentation techniques, and using deeper neural networks such as RNNs or transformers for better performance on sequential data like speech.

8. References

- 1. Warden, P. (2018). Speech Commands: A Dataset for Limited-Vocabulary Speech Recognition. Retrieved from https://arxiv.org/abs/1804.03209
- 2. PyTorch Documentation: https://pytorch.org/docs/stable/index.html
- 3. Torchaudio Documentation: https://pytorch.org/audio/stable/index.html