# UNIVERSITÀ DI PAVIA

Machine Learning Course

---

# Speech Commands Recognition

---

Andrea Alberti

Department of Computer Engineering - Data Science

University of Pavia, Italy

Email: andrea.alberti01@universitadipavia.it

February 25, 2024

## Abstract

This project aimed to develop a highly accurate speech recognition model capable of distinguishing between 35 different spoken words, trained on a dataset comprising 105,829 audio recordings by numerous speakers.

To achieve this objective, various MLP architectures were evaluated, resulting in a test accuracy of approximately 49%. However, the task's complexity, with numerous classes to classify, posed a significant challenge in developing a robust model. Additionally, some words had fewer recordings, which adversely impacted the overall model performance. The model choice was based on the comparison between different model's architectures and normalization techniques. The project provides insights about the significance of addressing dataset imbalances, the importance of feature normalization, and the performance of an MLP model on the speech recognition task.

# Contents

# 1 Introduction

This project is about the implementation of a Neural Network for the classification of audio files. This problem is known as *Speech Command Recognition* and its importance has been raising in the last years thanks to the development of smart devices and the need of a more natural interaction with them. The main library exploited for the implementation is *pvml*.

## 1.1 Available Data

The data set includes 105829 recordings of 35 words uttered by many different speakers. For the analysis it is divided into three subsets: *train*, *validation* and *test*.

## 1.2 Goal

The goal is to fit on the data a Multi Layer Perceptron (MLP) so that it can be exploited to classify new unseen data on the 35 different classes.

## 1.3 Multi Layer Perceptron

A Multi Layer Perceptron (MLP) is a type of Feed Forward Neural Network that consists of zero or more hidden layers of interconnected neurons between an input layer and an output layer. Each neuron in an MLP receives input signals from all neurons in the previous layer, applies an activation function to the weighted sum of its inputs, and outputs the result to the next layer. The parameters of the MLP (weights and biases) are learned optimizing a chosen Loss Function (Cross Entropy) through a process called Backpropagation.

# 2 Model Building

Prior to start building the model is important to preprocess the data, making them suitable to be passed to the training algorithm. This process is known as *Feature Extraction* and it is a fundamental one in the model creation.

## 2.1 Features Extraction

Starting from the raw audio data, extraction of features consists in the creation of a spectrogram for each word. This latter has been made uniform in size by padding or truncating the original signals.

## 2.2 Data Visualization

To better understand the data, it is useful to visualize some of the spectrograms. In figure 1 are reported some examples about some words. The red zones are those in which the frequencies are more present. It is important to underline that the same word can have different spectrograms depending on the speaker.
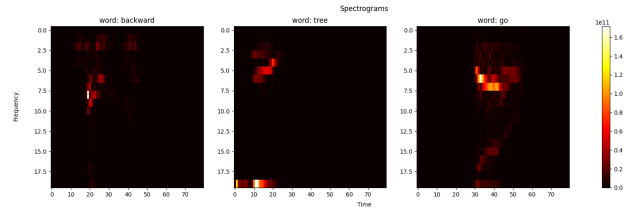


Figure 1: Spectrograms example

## 2.3 Feature Normalization

Another important step in the pre-processing phase is the normalization of the features. Among the several normalization techniques available, the four used in this project are *Mean-Var*, *Min-Max*, *Whitening* and *Max-Abs* normalization. All these methods are compared in figure 2. It is evident that the best one is the Mean-Variance normalization. This is because it helps to reduce the impact of differences in the volume or intensity of the audio recordings, centering the data around zero and scaling them to have unit variance. These normalized data allow the model to capture underlying patterns.
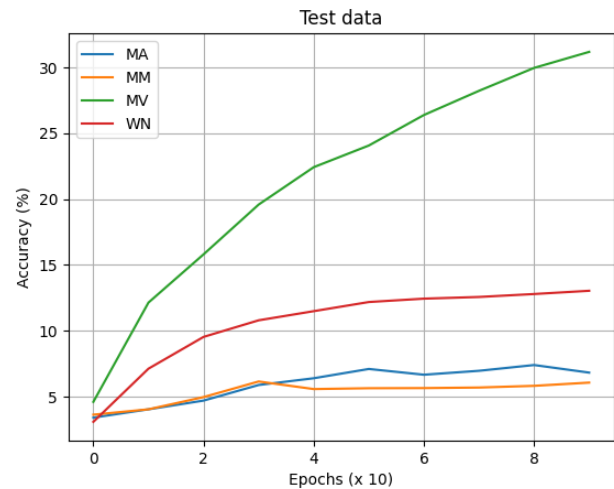


Figure 2: Normalization comparison

## 2.4 Train Classifier

**Batch size**
To optimize the training the data are divided into mini-batches. The size of the batch could impact on the model performance, therefore different values have been tested and the result are reported in figure 3. The used model had not hidden layers, therefore its performance is poor, thus all the sizes converges to 18% accuracy on the test set but requiring different number of epochs. The fastest was the 50 batch size and it is chosen as the best one.
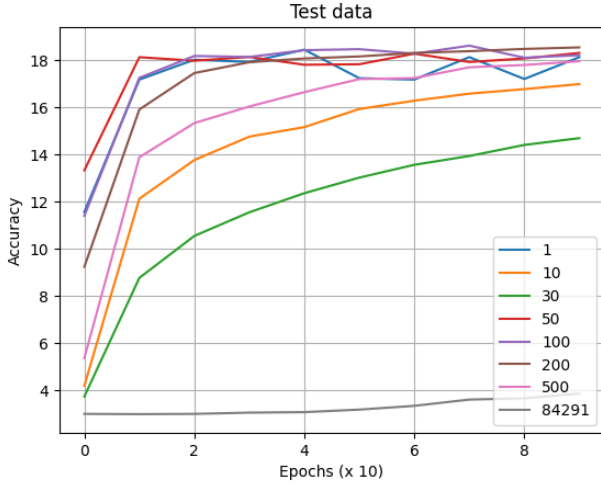
Figure 3: Batch comparison

**Architecture**

The architecture of the network strongly influences its performance therefore different configurations results are compared in figure 4. For time reasons, the number of epochs is limited to 100 disregarding that some networks have more parameters and could take more time to express their full potential. In this greedy approach the best architecture is that with many hidden layers. This type of architecture was trained with 400 epochs, in figure 5 are shown the results. The blue and orange model reached 80% in training accuracy, however, when tested on the independent test set, their performance was not as good, suggesting that they were overfitting to the training data. To address this problem, several regularization values were tested, and a regularization parameter of $\lambda = 0.001$ was found to be the most effective. Despite this effort, overfitting was only slightly reduced, but not entirely resolved. Ultimately, the red model was selected as the best-performing model and used for further analysis.
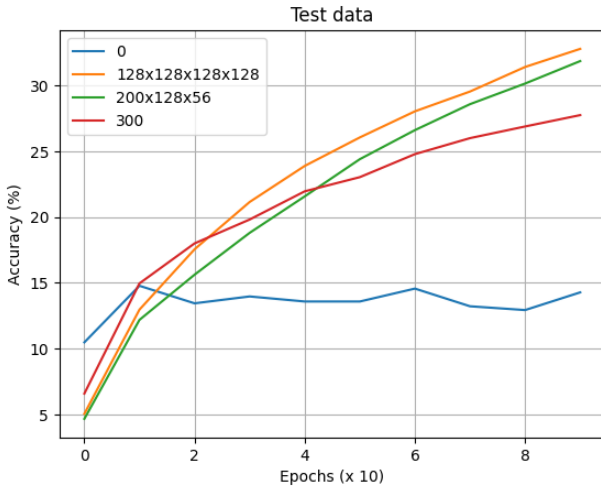


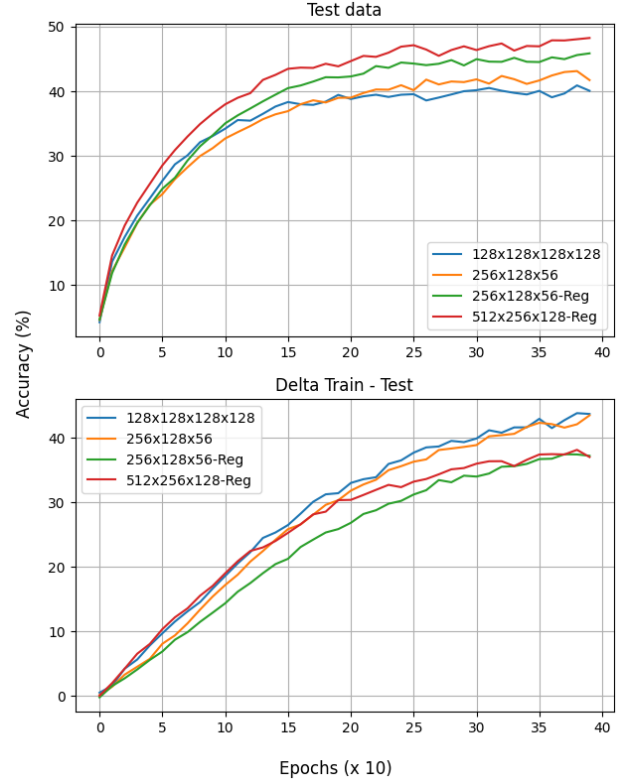Figure 4: Architecture comparison



Figure 5: Deep Networks comparison

# 3 Model Analysis

To analyze the behavior of the chosen model was created a confusion matrix in figure 12 that shows the number of correct and incorrect predictions for each class. The darker the color, the higher the number of samples of class i (row index) classified as class j (column index). As expected, the diagonal is the most populated part of the matrix, meaning that the model correctly classifies a word lots of times.
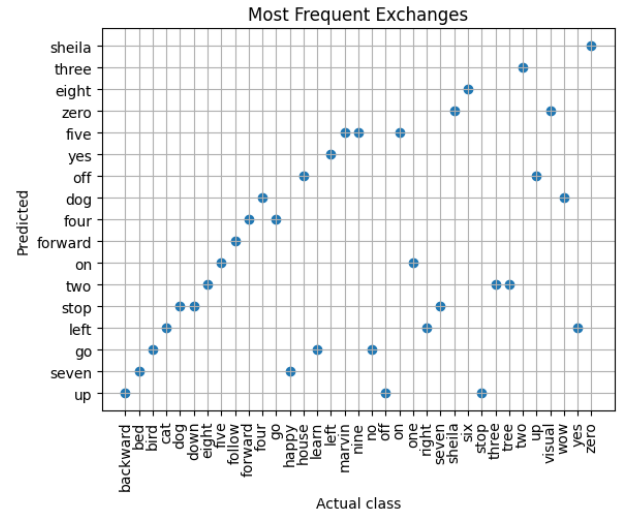


Figure 6: Most exchanged

**Most Exchanged Classes**

The most exchanged classes are those in which the

model has more difficulties to distinguish between. In figure 6 are reported for each word the class to which the model classifies them (excluded correct classifications) most of times. Lots of exchanges are reasonable like *Forward* classified to *Four* and *Three* to *Two* whose spectrograms are very similar (figure 7).
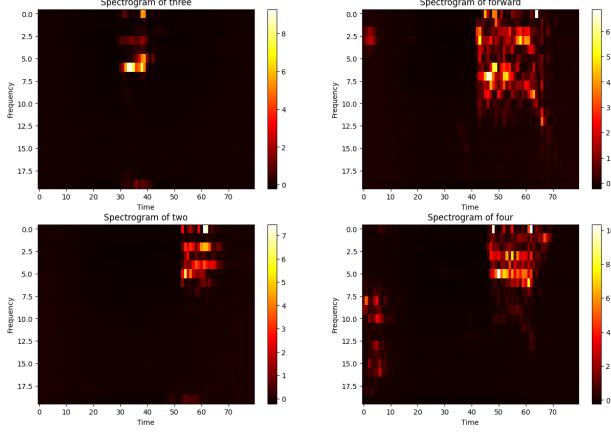


Figure 7: Three VS Two, Forward VS Four

**Most Misclassified Samples**

Some words are more difficult to be classified than others and in figure 8 are reported the most misclassified samples. *Tree* and *Learn* are the most misclassified words, listening to the audio files of *learn* it is possible to notice that lots of time the word is mispronounced and even a human could have difficulties to understand it. However, the actual reason between the misclassified words is explained by figure 9. It is evident that the most misclassified words are those with less samples, therefore the model has less information to learn and this highly affects its overall performance.
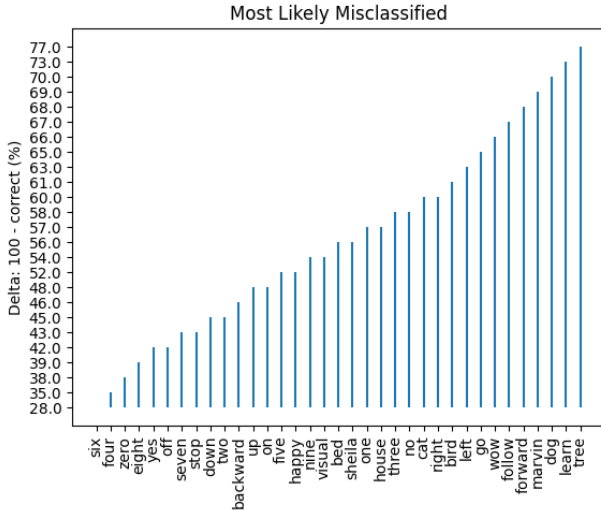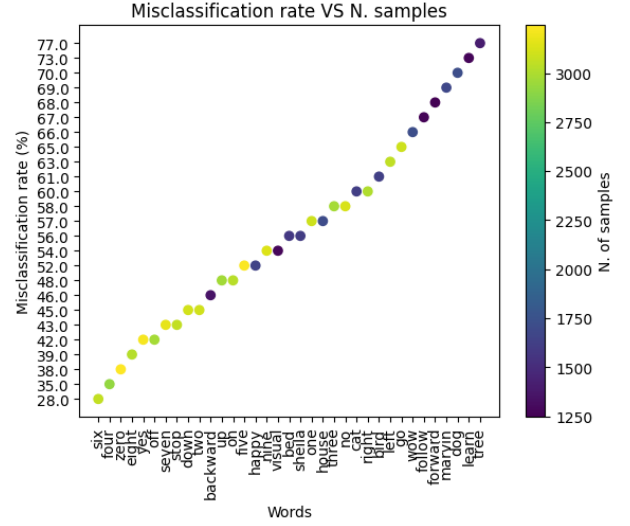


Figure 8: Most misclassified



Figure 9: Misclassified VS N. of samples

**Weights Visualization**

An interesting way to understand how the model works is to visualize the weights of the first layer, shown in figure 11. The red parts votes in favour of the specific class, while the blue parts votes against it. The weights distribution is not uniform and changes from word to word, according to their spectrograms. Remembering that the spectrograms has on y axis the frequency and on x axis the time, it is possible to notice on which frequencies the model focused. For example, the high frequencies strongly penalized words like *Follow* and favours words like *up* and *wow*. In the case of *Forward*, *Four* and *Follow*, due to the similarity of their spectrograms, the model assigns almost the same weights to the same frequencies. Indeed, the model learnt to extract relevant acoustic details, as shown by the overlapping between the assigned weights and the spectrograms of a specific word (figure 10).
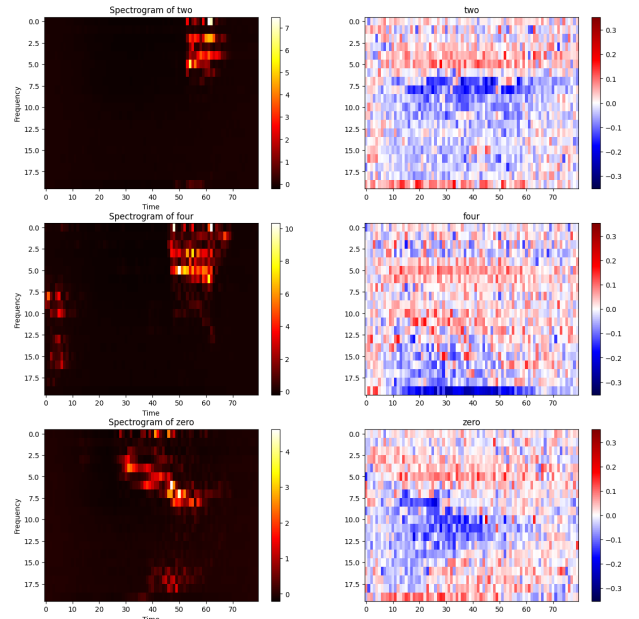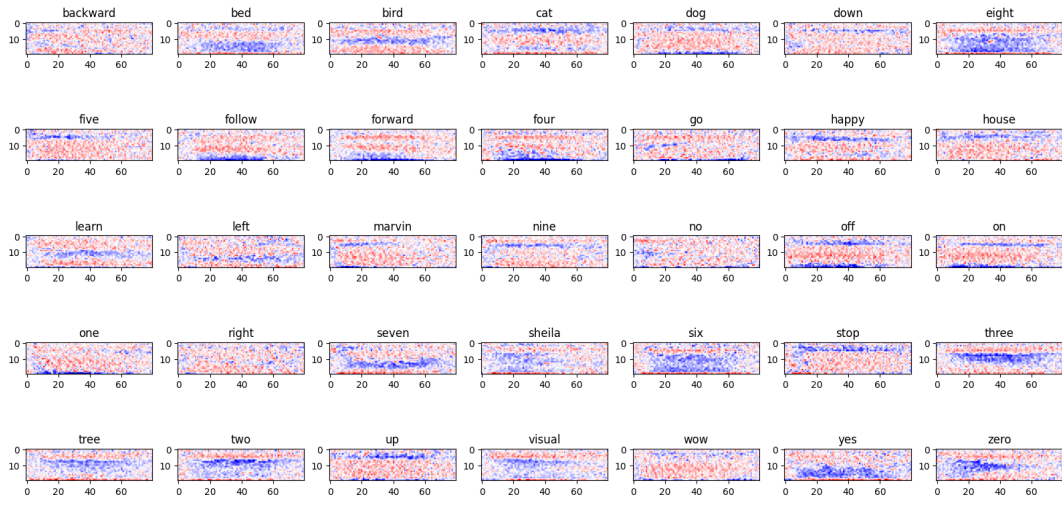


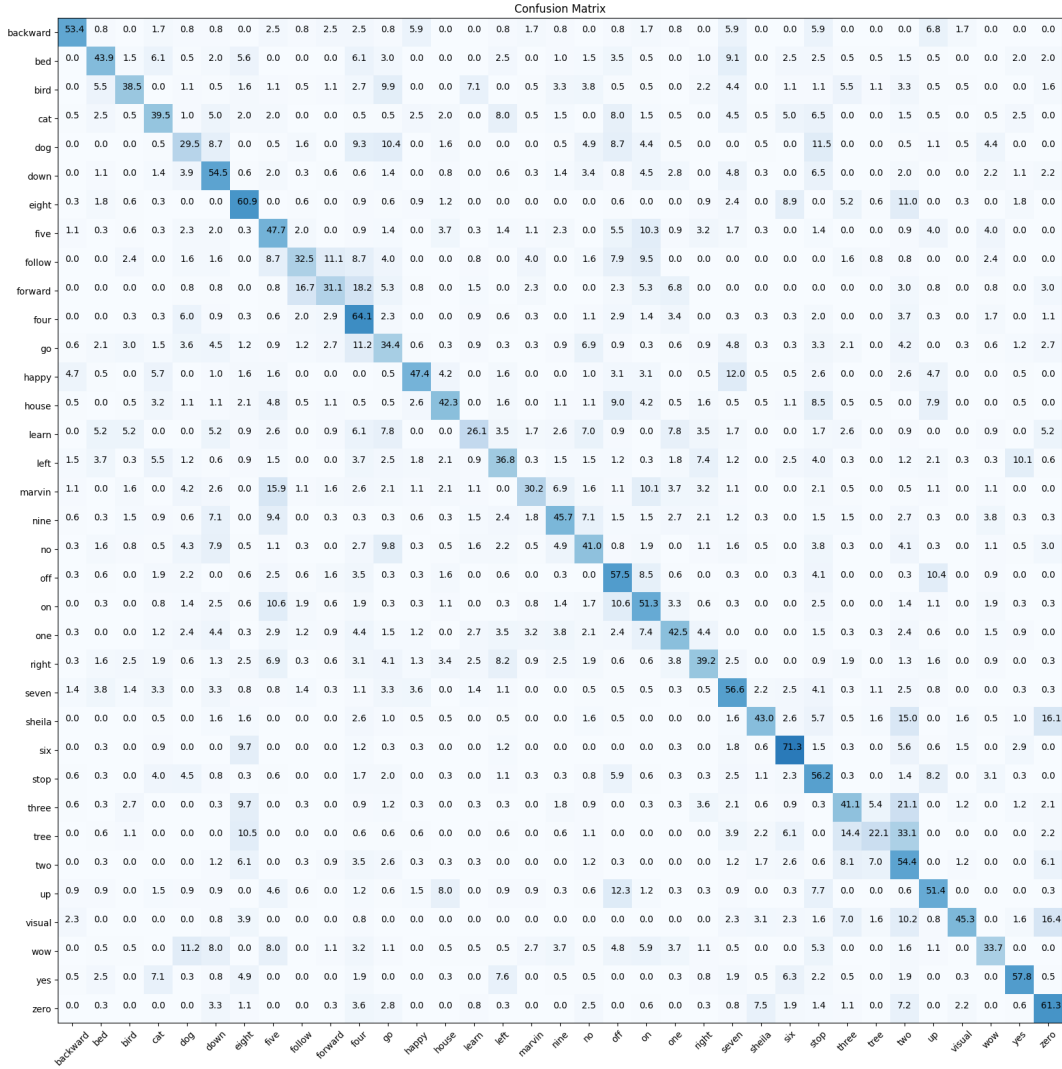Figure 10: Spectrogram VS Weights

Figure 11: Weights visualization

Confusion Matrix

Figure 12: Test accuracies comparison