*Springboard DSC*

# CAPSTONE PROJECT 2

# URBAN SOUND CLASSIFICATION

Milestone Report 2

Vikram Lucky

Dr.vlucky@gmail.com

01 May 2019

## Description and Objectives

We all get exposed to vast variety of urban sounds every day. Unlike speech and audio signals, Urban sounds are usually unstructured sounds. They include various real-life noises generated by human activities ranging from car horns, engine idling to music and children playing. The aim of this work is to extract features from audio files and using those features create Machine Learning (ML) models to accurately label such audio files into their specific categories (*out of 10 categories, in this case*). The idea is to discover relevant patterns in the audio features and use the discovered patterns as characterizing attributes for ML models. In this project I have used Logistic Regression, Decision trees, XGBoost, and Deep Neural Networks. I compared the results of these models, and picked the one that most accurately labeled sounds into categories.

## The Client

The automatic classification of urban sounds is relevant in many areas and has a variety of applications including surveillance, highlight extraction, environmental monitoring, video summarization etc. Most Importantly, It also has the potential of improving the quality of life of city dwellers by providing a data-driven understanding of urban sound and noise patterns.

## The Data

Publicly available dataset called "*UrbanSound-8k*", created by Justin Salamon, Christopher Jacoby, and Juan Pablo Bello is used for this project.

This dataset is a collection of 8,732 short audio clips (upto 4 secs) of urban sound areas. The dataset is divided into 10 classes: *air conditioner, car horn, children playing, dog bark, drilling, engine idling, gun shot, jackhammer, siren and street music.*

## Data Wrangling

As with all unstructured data formats, a data has a couple of preprocessing steps which have to be followed before it is presented for analysis. The first and most important step is to extract features from audio files, so it is ready for machine learning analysis. We used library called "*LibROSA*" for feature extraction. LibROSA is a python package for music and audio analysis. It provides the building blocks necessary for  to create audio information retrieval systems. LibROSA provides several methods to extract different features from sound clips. For this project we used the following methods to extract various features:

- *melspectrogram: Compute a Mel-scaled power spectrogram*
- *mfcc: Mel-frequency cepstral coefficients*
- *Spectral_contrast: Compute spectral contrast*
- *Chroma: Compute a chromagram from a waveform or power spectogram*
- *Tonnetz: Compute the tonal centroid features*

A Python function was written which went over each audio file present in the provided training set and extracted above mentioned features from each one of them and it returned:
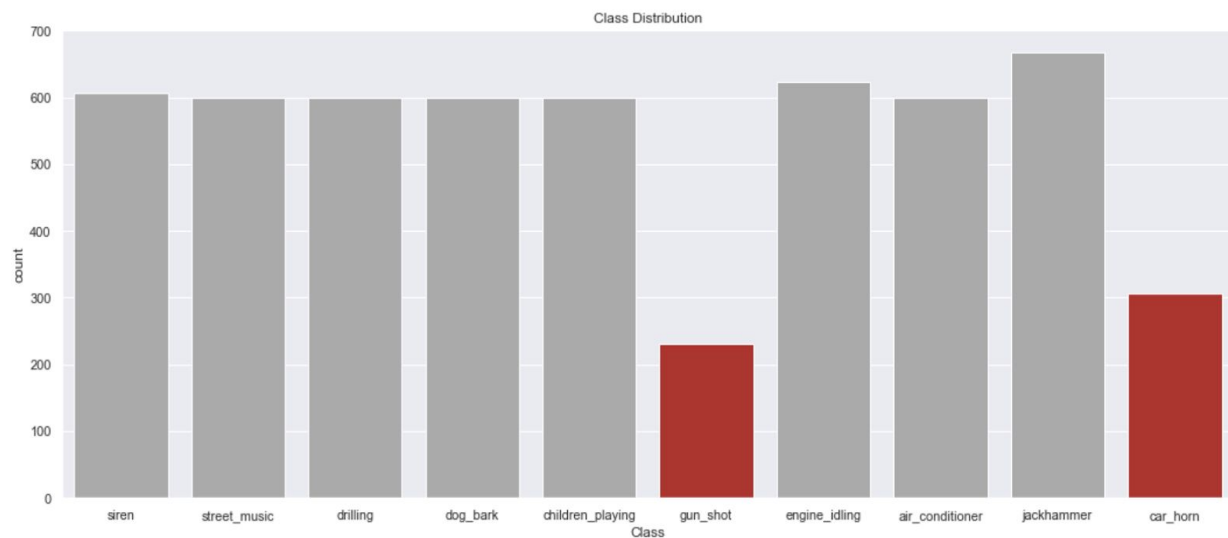
- *Tonnetz* data frame (5435, 7)
- *Mfccs* data frame (5435, 41)
- *Chroma* data frame (5434, 13)
- *Mel Spectrogram* data frame (5435, 129)

- *Contrast* data frame (5435, 8)
- df: A data frame combined of all of the data frames above (5435, 196)

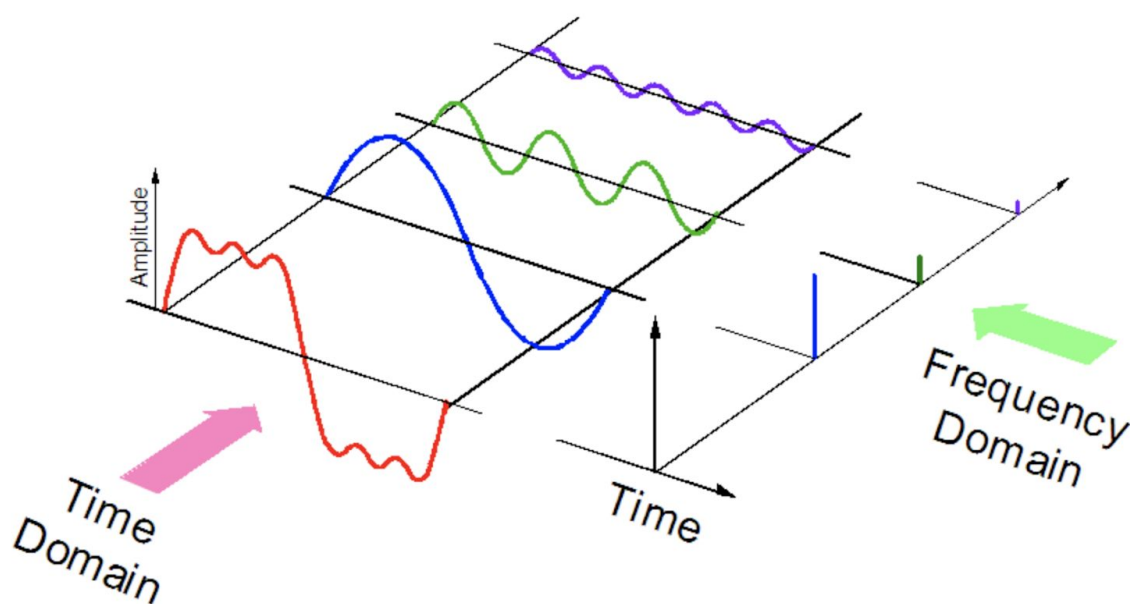# Exploratory Data Analysis and Inferential Statistics

In this section we present a series of questions we posed as part of our storytelling.

We begin by looking into the distribution of Classes:



Bar graph above shows distribution of classes. It can be observed that *car_horn* and *gun_shot* are under represented, when compared to other classes present in dataset. As other classes were almost equally distributed, I decided not to balance minority classes by introducing additional external or synthetic data, until I look at the classification report.
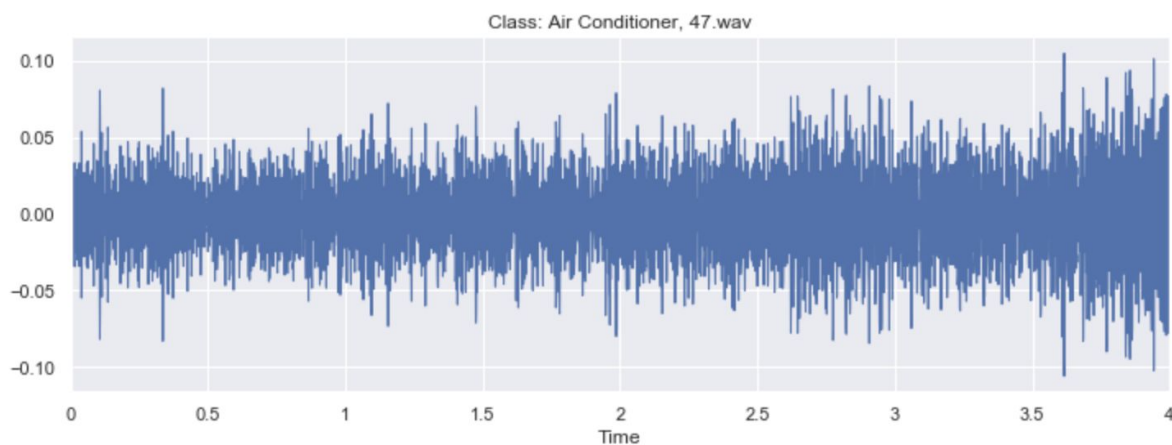
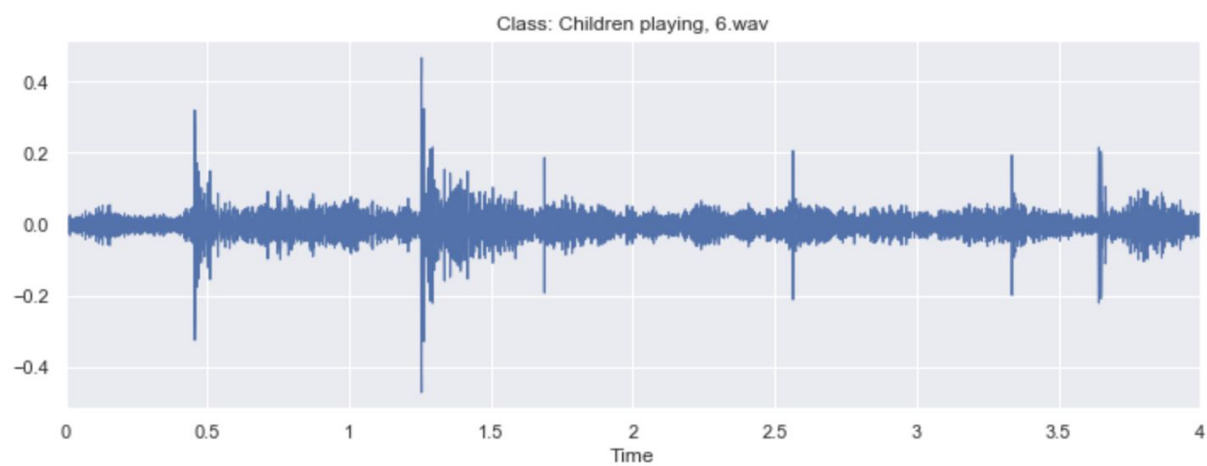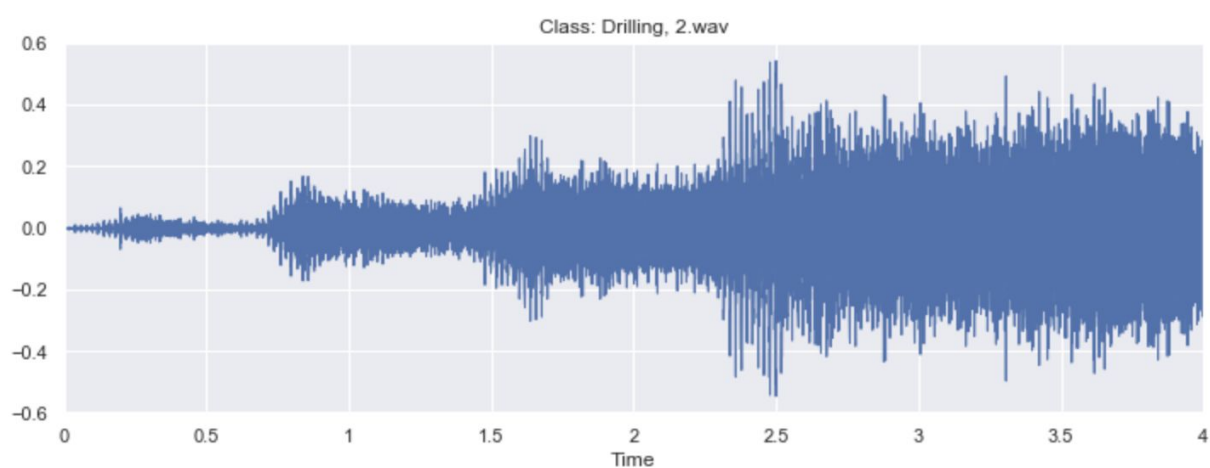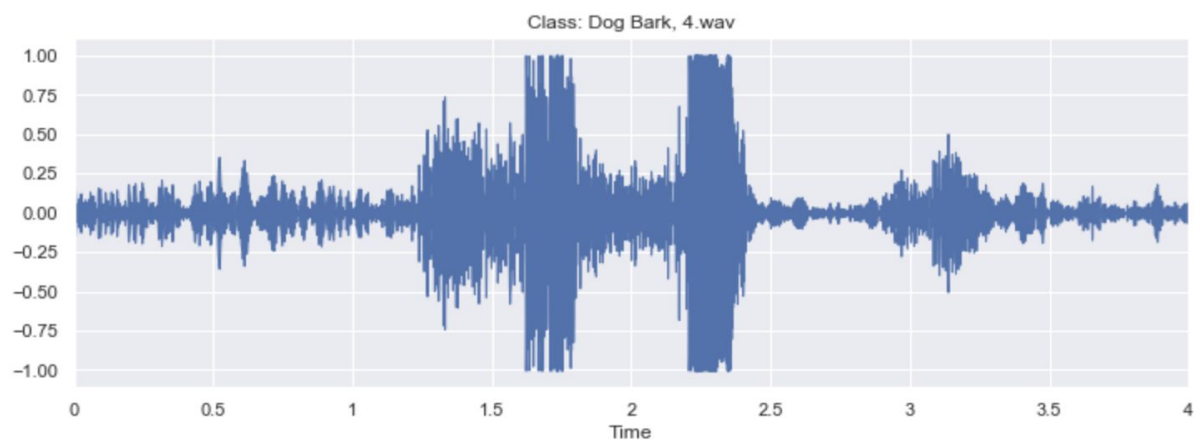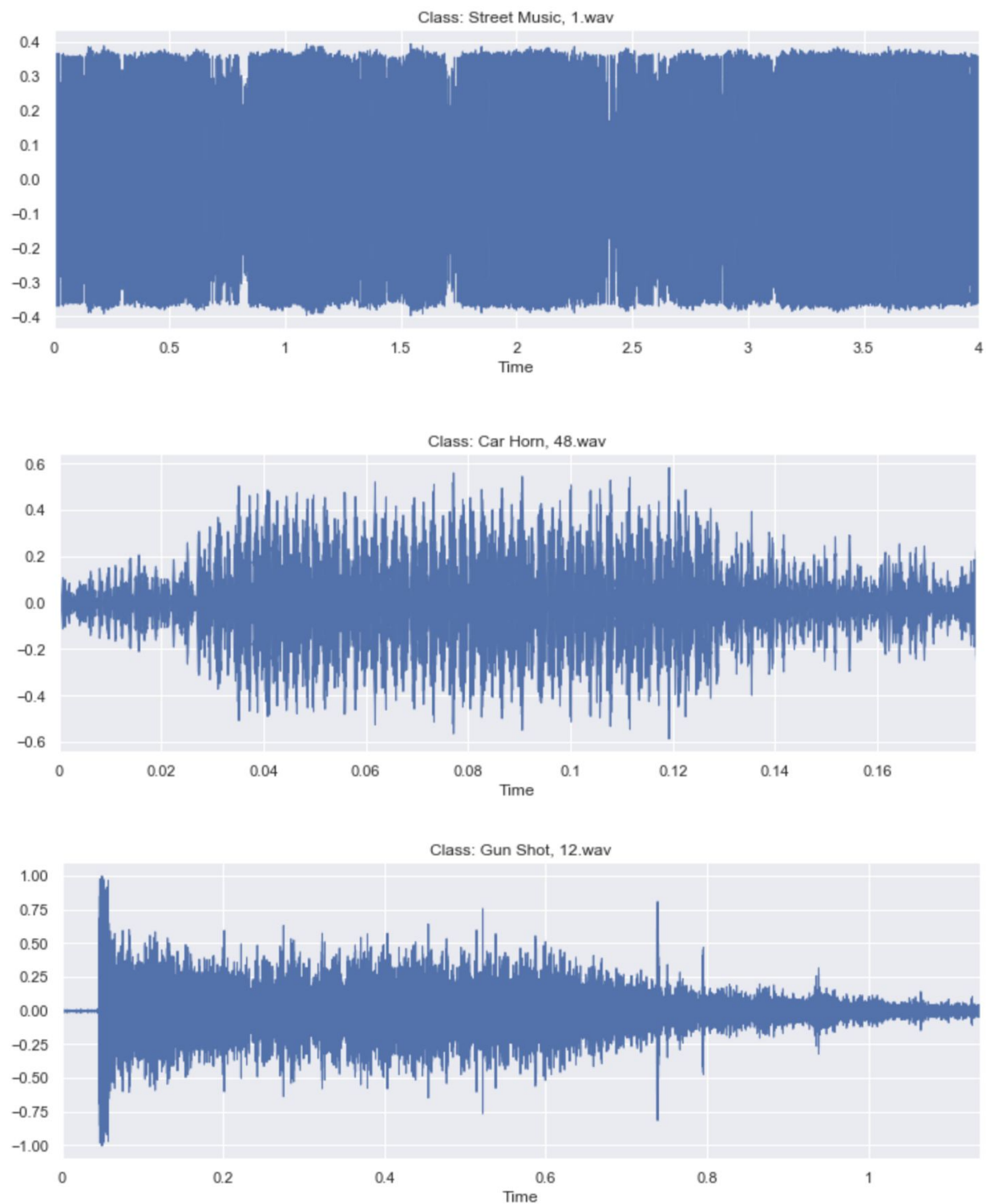There are 2 commony ways to represent sound:



- **Time domain**: each sample represents the variation in air pressure.
- **Frequency domain:** at each time stamp we indicate the amplitude for each frequency.

As shown in the above picture, there are 2 common ways to represent audio sound: Time domain and Frequency domain

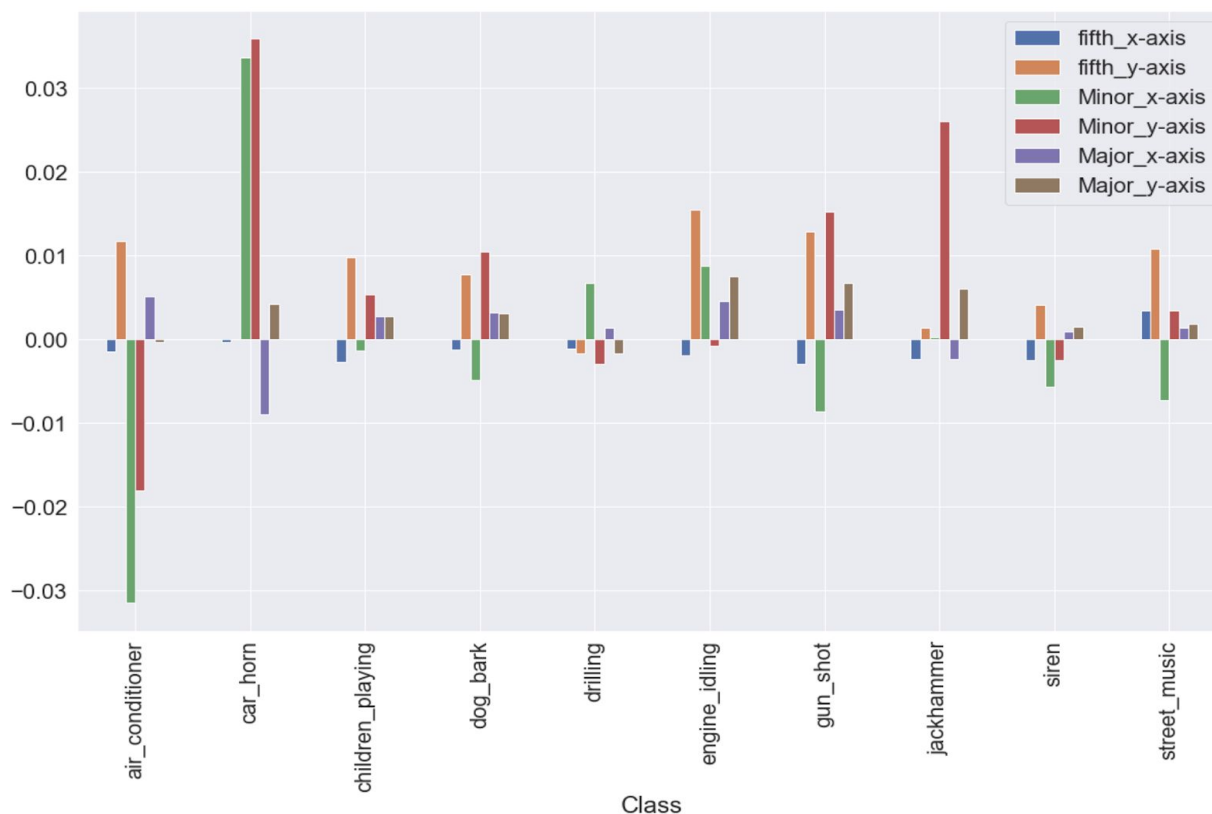## Time domain: wave plots of different class types:

Class: Dog Bark, 4.wav



Class: Drilling, 2.wav



Class: Children playing, 6.wav

Class: Street Music, 1.wav



Class: Car Horn, 48.wav



Class: Gun Shot, 12.wav

The above waveform plots show the differences in the sound waves for each specific class.

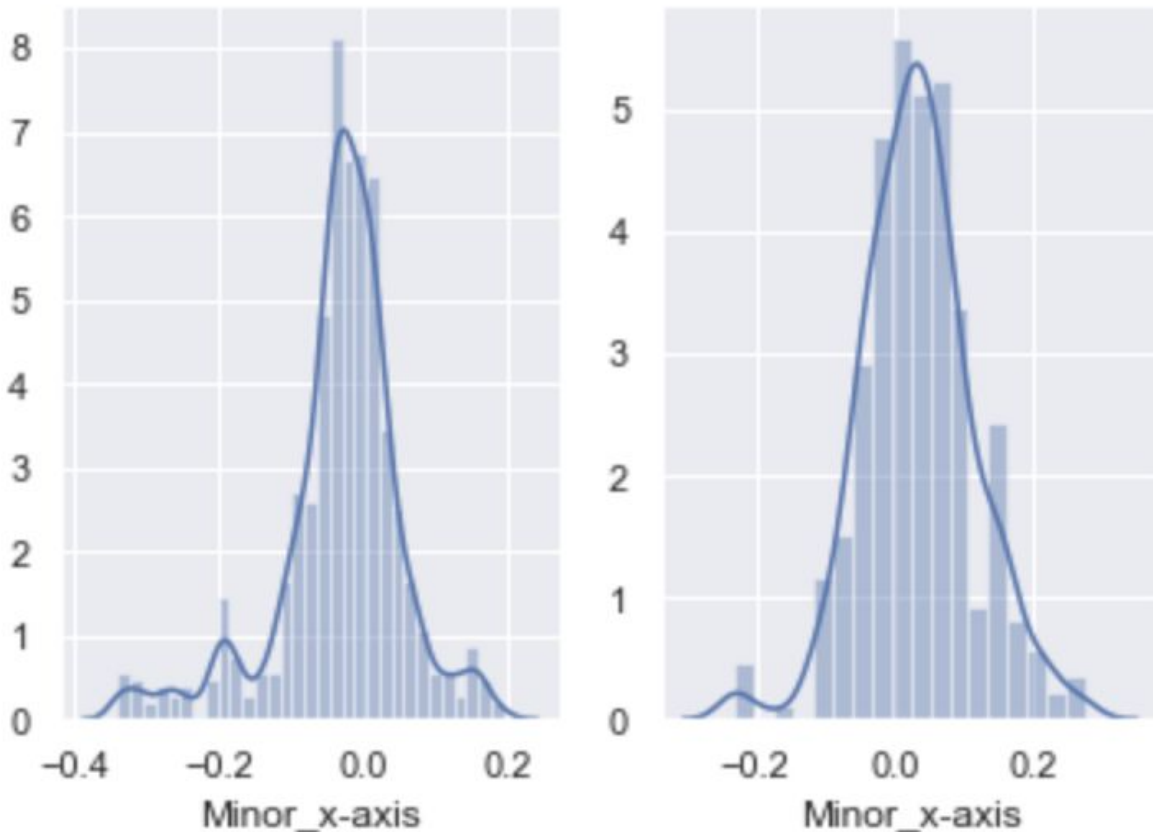**Bar graph of Tonnetz:** Tonnetz computes the tonal centroid features.



Above graph shows average tonal centroid features for each frame for different classes. A noticeable pattern can be observe between *car horn* and *air conditioner* classes:

- *Average Minor_x-axis and Minor_y-axis is higher for **car_horn** class compare to other classes*
- *Average Minor_x-axis and Minor_y-axis is the lowest for **air conditioner** class compare to all of the other classes.*

**Inferential Stats:**

Hypothesis test was conducted to check the relationship between Minor_x-axis of **car_horn** and **air conditioner** classes. Is difference between their mean statistically significant?
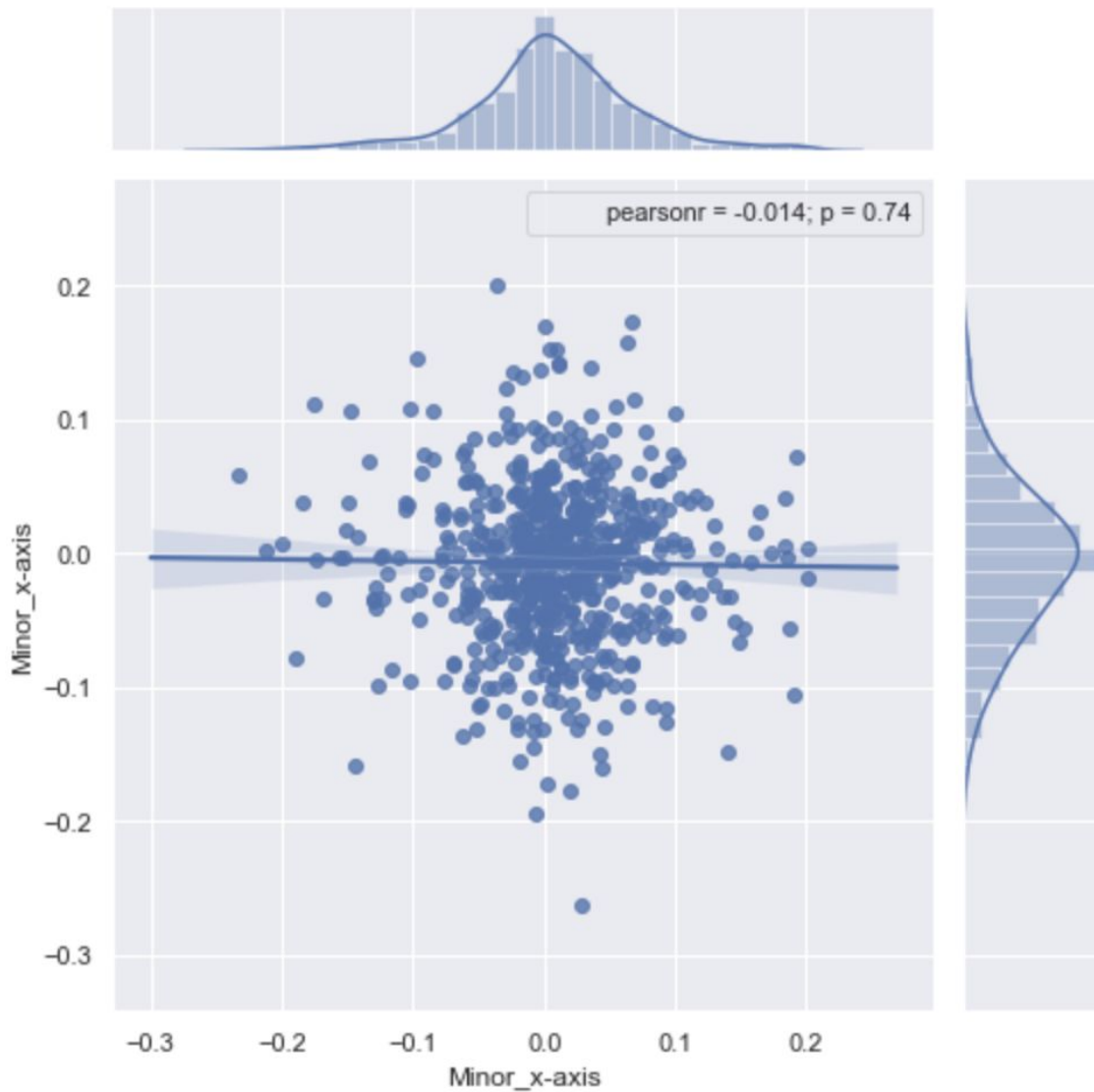


Ho: Difference between car_horn's mean and air_conditioner's mean is not statistically significant

H1: Difference between car horn class's mean and air conditioner class's mean is statistically significant

a: 0.05

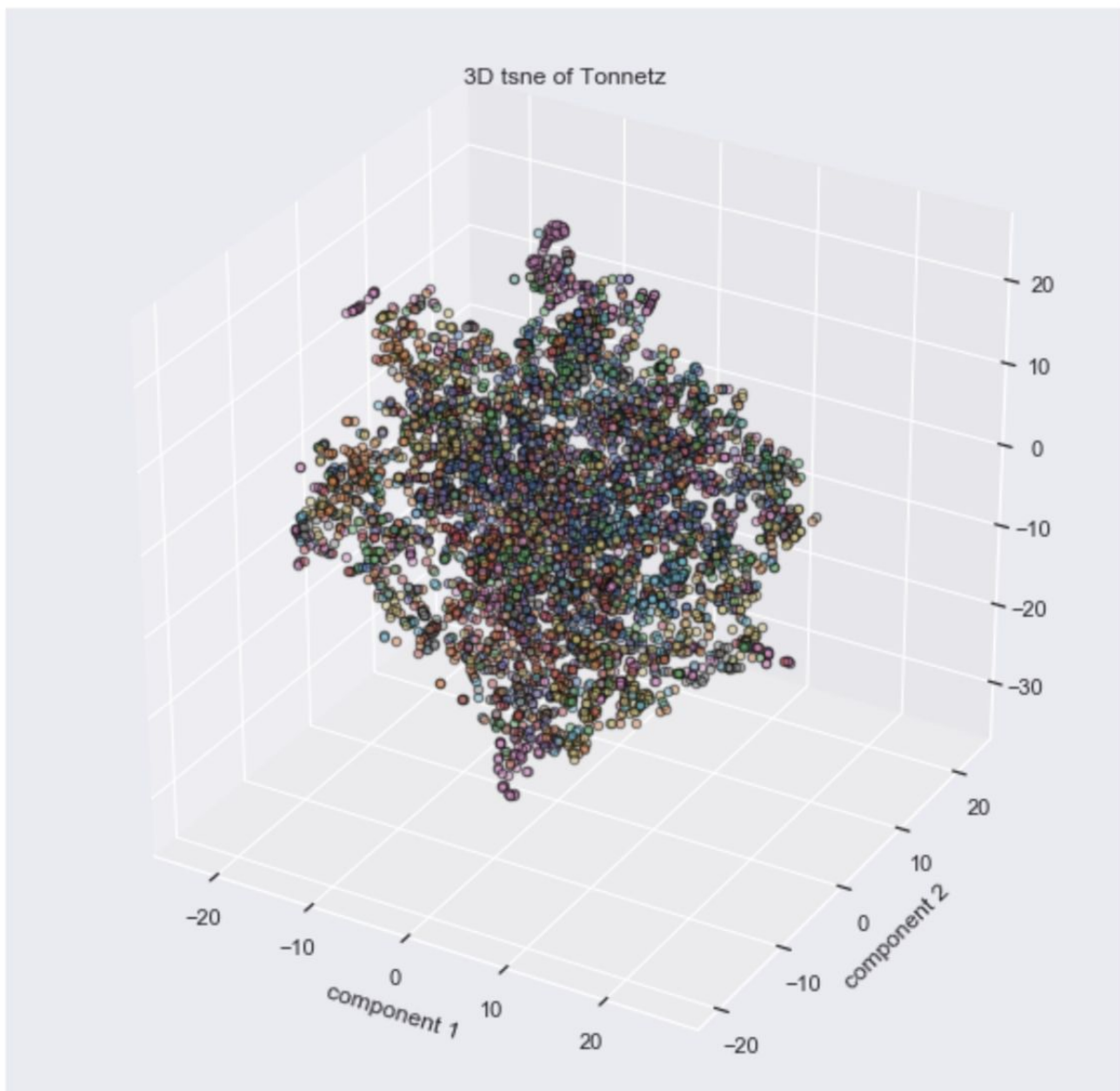- P - value observed were less than the threshold, therefore Ho was rejected.

**Pearson's** correlation coefficient test was conducted to measure the statistical relationship between **Minor_x-axis** of drilling and street music classes.
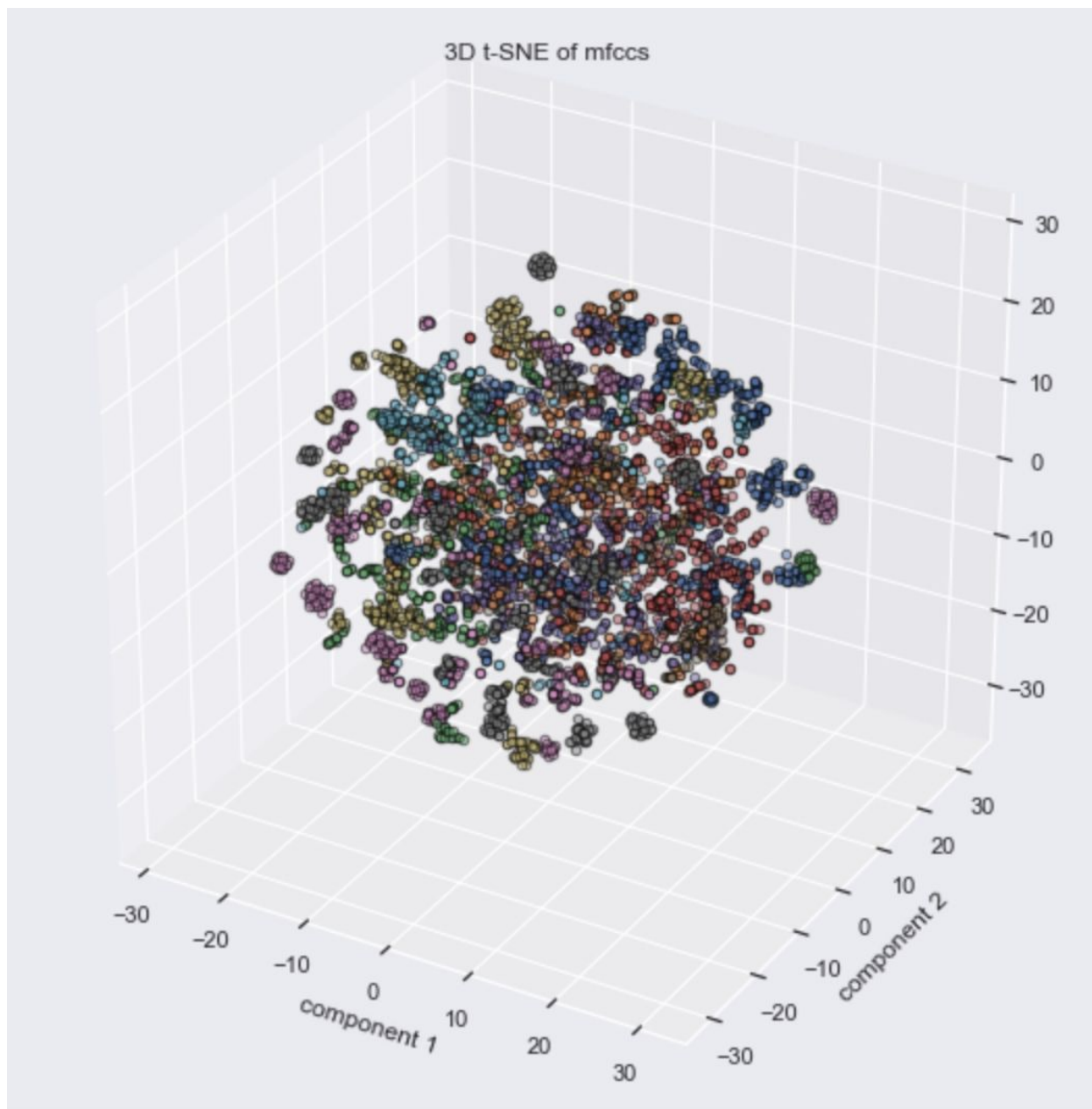


As p value is 0.74, much higher than threshold 0.05. The null hypothesis could not be rejected, therefore we cannot say whether the observed correlation between the two variables is statistically significant or not.

# Data Visualization using 3D Graphs

As the number of column increases it became harder for us to identify any hidden patterns using basic histograms and bar graphs. Therefore visualization technique like **t-SNE**(t-distributed stochastic neighbor embedding) was used to explore data in 3D space

Above graphs shows Tonnetz features in 3D space, different colors represents different classes. It's hard to differentiate visually between classes using Tonnetz.



Above graph shows Mfccs features in 3D space, Much better than Tonnetz representation of data. As some classes grouped together can be easily observed visually.
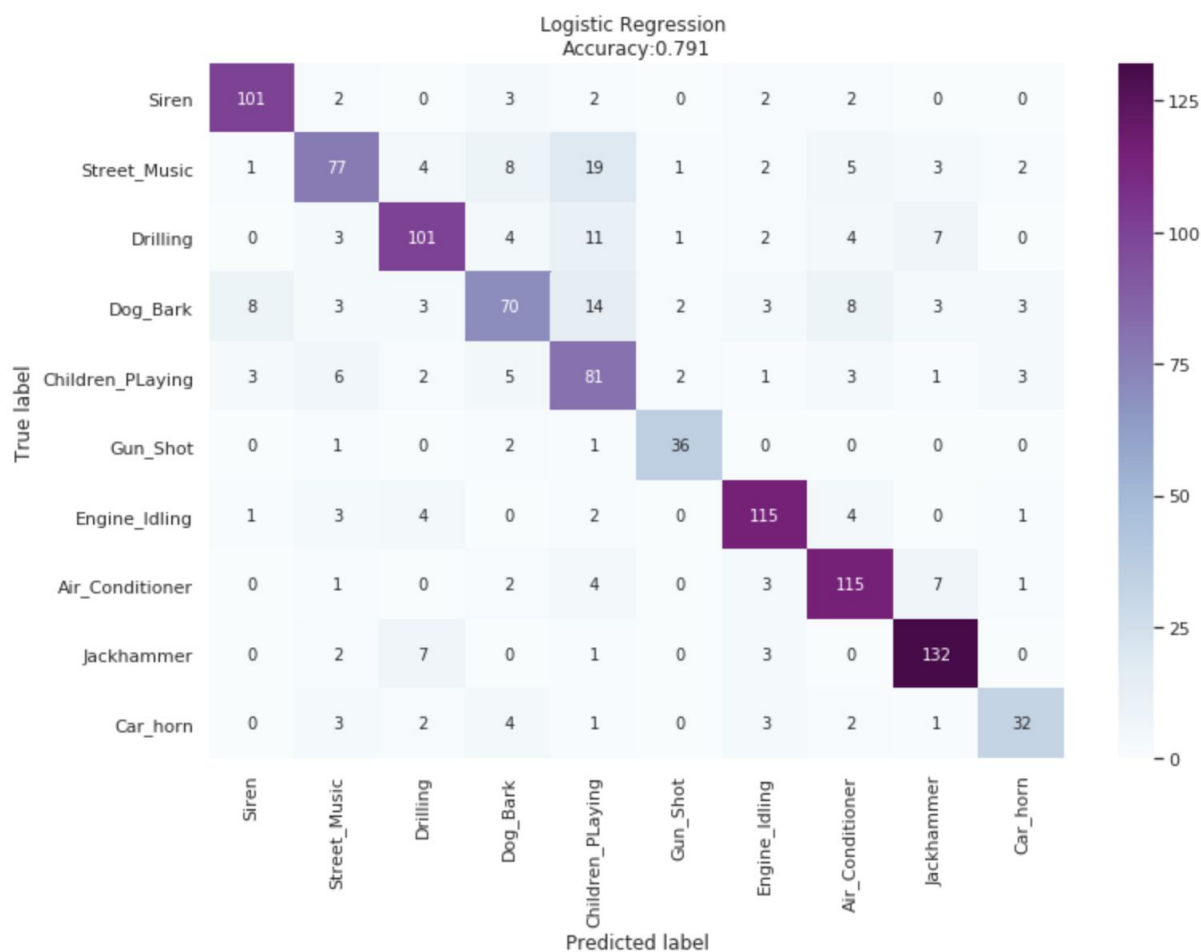
## PREDICTIVE MODELING

# *Baseline modeling*

In this part of the project, our first task was to create a baseline set of classification experiments to compare our results with. The goal of the baseline experiments was not to produce optimal parameters and to maximize accuracy, but to study the problems and characteristics of the dataset itself.

```
{'C': 1, 'penalty': 'l1', 'solver': 'liblinear'}
0.7854645814167434
```

|  | precision | recall | f1-score | support |
|---|---|---|---|---|
| 0 | 0.89 | 0.90 | 0.89 | 112 |
| 1 | 0.76 | 0.63 | 0.69 | 122 |
| 2 | 0.82 | 0.75 | 0.78 | 133 |
| 3 | 0.72 | 0.59 | 0.65 | 117 |
| 4 | 0.59 | 0.76 | 0.66 | 107 |
| 5 | 0.86 | 0.90 | 0.88 | 40 |
| 6 | 0.84 | 0.88 | 0.86 | 130 |
| 7 | 0.80 | 0.86 | 0.83 | 133 |
| 8 | 0.86 | 0.91 | 0.88 | 145 |
| 9 | 0.76 | 0.67 | 0.71 | 48 |
| micro avg | 0.79 | 0.79 | 0.79 | 1087 |
| macro avg | 0.79 | 0.79 | 0.78 | 1087 |
| weighted avg | 0.79 | 0.79 | 0.79 | 1087 |

We used Logistic Regression with various regularization techniques. We ran 5-fold cross validation, and in each fold model was trained on 80% of the data and tested on the rest 20%. Our baseline classifier achieved an average classification accuracy of ~79% averaged across all classes.



Logistic Regression
Accuracy:0.791

## Extended Analysis

Above we created a baseline model using Logistic Regression and we achieved an accuracy of ~79% across all classes. In the above heatmap we can see that the model did not distinguished well between some of the classes for example: street music and children playing. Now that we have baseline model set up, something that we can compare our sophisticated

models to. We are ready to jump into extended analysis part of this project, In extended analysis our goal will be to create a more sophisticated machine learning model using different algorithms and techniques that captures our data well and successfully beats our baseline model's performance.