# Università degli Studi di Milano

## Data Science and Economics

Machine Learning Project

# Urban Sound Classification with Neural Networks

Antonella D'Amico
961150

Academic year 2021/2022

## Declaration

*I declare that this material, which I now submit for assessment, is entirely my own work and has not been taken from the work of others, save and to the extent that such work has been cited and acknowledged within the text of our work. I understand that plagiarism, collusion, and copying are grave and serious offences in the university and accept the penalties that would be imposed should I engage in plagiarism, collusion or copying. This assignment, or any part of it, has not been previously submitted by me or any other person for assessment on this or any other course of study.*

# Contents

**Abstract.** The aim of this project is to implement a Neural Network able to classify urban sound contained in UrbanSound8K dataset. Three different Neural Network architecture was developed with three different features extracted from audio files (MFCCs, Mel-Spectogram and Chromagram) as input.

# Dataset Description

## 1.1 Variables Description

In this experimental project we will use UrbanSound 8k dataset[1]. It contains 8732 labeled sound excerpts ($\leq$ 4 sec) of urban sounds from 10 different classes:

| Class | Conut |
|---|---|
| air_conditioner | 1000 |
| car_horn | 429 |
| children_playing | 1000 |
| dog_bark | 1000 |
| drilling | 1000 |
| engine_idling | 1000 |
| gun_shot | 374 |
| jackhammer | 1000 |
| siren | 929 |
| street_music | 1000 |

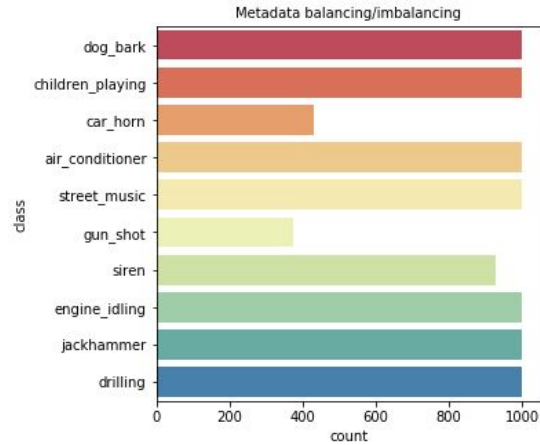Table 1: Classes and their count



Figure 1.1: Classes distribution

Each sound sample is labeled with the class to which it belongs.
After downloading the dataset we can see that contains two main parts:

- **Audio files** stored in *audio* folder. This folder contains a total of 10 sub-folds labeled from fold1 to fold10. In each of them there are stored a variable amount of audio files in .wav format.

- **Metadata** stored in *metadata* folder. It contains a .csv files where main audio information are listed, such as: filename, numeric classID and the corresponding class according to those showed above.

---

[1]https://urbansounddataset.weebly.com/urbansound8k.html

| | slice_file_name | fsID | start | end | salience | fold | classID | class |
|---|---|---|---|---|---|---|---|---|
| 0 | 100032-3-0-0.wav | 100032 | 0.0 | 0.317551 | 1 | 5 | 3 | dog_bark |
| 1 | 100263-2-0-117.wav | 100263 | 58.5 | 62.500000 | 1 | 5 | 2 | children_playing |
| 2 | 100263-2-0-121.wav | 100263 | 60.5 | 64.500000 | 1 | 5 | 2 | children_playing |
| 3 | 100263-2-0-126.wav | 100263 | 63.0 | 67.000000 | 1 | 5 | 2 | children_playing |
| 4 | 100263-2-0-137.wav | 100263 | 68.5 | 72.500000 | 1 | 5 | 2 | children_playing |

Figure 1.2: Metadata file

## 1.2 Features extraction

From a visual inspection we can see, in **Fig.1.3**, how a children playing wave sound look like.
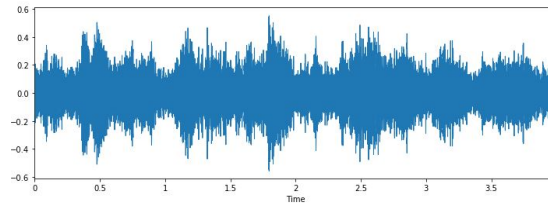


Figure 1.3: Children_playing wave sound

For the purpose of our analysis we extract three main features from each raw time series file. Analyzing the different sound waves of the ten classes contained in the dataset, it is possible to notice how some are very different from each other (the sound of a gunshot is very short if it is compared to a drill); others have repetitive sounds, such as the pneumatic hammer and the drill, and for this reason they are similar in shape. The fact that some sounds may have similar shapes may not be great for a classification analysis point of view. For a more detailed and precise analysis, we decided to extract three of the main features contained in sounds: **MFCCs**, **Mel-Spectogram** and **Chromagram**.

We use librosa library to extract those features by developing three different functions, that were then applied to predefined fold in order to construct our training and test set. For all functions the value for corresponding n_mfcc, n_mel and n_chroma parameters are set eqaul to 50.

MFCCs, Mel-Spectogram and Chromagram can be displayed in plots. In **Fig.1.4** an example of how the raw time series file *17973-2-0-21.wav*, classified as children_playing, can be represented.

4

(a) MFCCs



(b) Mel-Spectogram



(c) Chromagram

Figure 1.4: Children playing features

## 1.3 Data Preparation

For the aim of this project we will use folds 1, 2, 3, 4 and 6 as *training set* and folds 5, 7, 8, 9, 10 as *test set*. We extract features in one time[2] using the three functions developed before. We store the obtained results in lists divided for extracted feature and training/test set.

After each extraction and before starting our analysis, features and classes are combined in a dataframe, where the column **features** contains a list of values which corresponds to statistics of the extracted features and column **class** contains the label associated to that file. In **Fig.1.5** an example of a dataframe created with features and class obtained from the use of the function extract_mfcc:

---

[2]Please note that this operation will take about 30min. It could seems lot of time, but we also performed an experiment by extracting them separately and, at the end, they take the same time. We preferred to put all together so we avoid repetition of code.

Figure 1.5: Training dataframe for MFCCs features

We perform this operation for all the extracted features. Then, we proceed by splitting features and class, again for all the extracted features, both for training and test set. Since class are categorical we need to transform them into binary variable by using One-Hot Encode techniques. Thanks to techniques we are able to create number of dummy variables equal to the number of class, in our case 10, and for each row mark with 1 the dummy that corresponds to the category associated with that file and 0 otherwise.

```
array(['dog_bark', 'dog_bark', 'dog_bark', ..., 'street_music',
       'street_music', 'street_music'], dtype='<U16')
```

(a) before OHE

```
array([[0., 0., 0., ..., 0., 0., 0.],
       [0., 0., 0., ..., 0., 0., 0.],
       [0., 0., 0., ..., 0., 0., 0.],
       ...,
       [0., 0., 0., ..., 0., 0., 1.],
       [0., 0., 0., ..., 0., 0., 1.],
       [0., 0., 0., ..., 0., 0., 1.]], dtype=float32)
```

(b) after OHE

Figure 1.6: y_train data for MFCCs

6

# Modeling

To perform our analysis three different models of Neural Networks were implemented. For all of them **Categorical Cross-entropy** loss function and **Adam** optimizer were used. **Softmax** activation function is used in output layer of all the three model, while for hidden layers **Rectified Linear Unit**, also known as ReLU, activation function is used. The latter is the most communly used activation function in deep learning models. It return 0 if it receives any negative input, but for any positive value x, it returns back that value.

For all models batch size is set to **32** and epochs to **100**.

## 2.1 Model 1

The first model is composed by three hidden layers with a corresponding number of neurons equal to 512, 256 and 127.

```
Layer (type)                Output Shape            Param #
=================================================================
dense_97 (Dense)            (None, 256)             13056

dense_98 (Dense)            (None, 512)             131584

dense_99 (Dense)            (None, 256)             131328

dense_100 (Dense)           (None, 127)             32639

dense_101 (Dense)           (None, 10)              1280
=================================================================
Total params: 309,887
Trainable params: 309,887
Non-trainable params: 0
```

Figure 2.1: Summary of Model 1

## 2.2 Model 2

In the second model we add an hidden layers of 64 neurons with respect to Model 1.

```
Layer (type)                Output Shape            Param #
=================================================================
dense_102 (Dense)           (None, 256)             13056

dense_103 (Dense)           (None, 512)             131584

dense_104 (Dense)           (None, 128)             65664

dense_105 (Dense)           (None, 256)             33024

dense_106 (Dense)           (None, 64)              16448

dense_107 (Dense)           (None, 10)              650
=================================================================
Total params: 260,426
Trainable params: 260,426
Non-trainable params: 0
```

Figure 2.2: Summary of Model 2

## 2.3 Model 3

Third model is composed by the same numbers and corresponding neurons of Model 2, but in order to prevent overfitting a **Dropout** equal to 0.2 is added for each hidden layers.

```
Layer (type)                 Output Shape              Param #
=================================================================
dense_108 (Dense)            (None, 256)               13056

dropout_45 (Dropout)         (None, 256)               0

dense_109 (Dense)            (None, 512)               131584

dropout_46 (Dropout)         (None, 512)               0

dense_110 (Dense)            (None, 128)               65664

dropout_47 (Dropout)         (None, 128)               0

dense_111 (Dense)            (None, 256)               33024

dropout_48 (Dropout)         (None, 256)               0

dense_112 (Dense)            (None, 64)                16448

dropout_49 (Dropout)         (None, 64)                0

dense_113 (Dense)            (None, 10)                650
=================================================================
Total params: 260,426
Trainable params: 260,426
Non-trainable params: 0
_____
```

Figure 2.3: Summary of Model 3

# Results

Three different models where used in this project with three different types of features as input (MFCCs, Mel-Spectogram and Chormagram).

| | Avg Accuracy | Std |
|---|---|---|
| Model 1 | 48.37% | 0.026 |
| Model 2 | 51.25% | 0.0228 |
| Model 3 | 52.67% | 0.0323 |

(a) MFCCs

| | Avg Accuracy | Std |
|---|---|---|
| Model 1 | 39.28% | 0.0157 |
| Model 2 | 43.18% | 0.0215 |
| Model 3 | 45.77% | 0.0202 |

(b) Mel-Spectogram

| | Avg Accuracy | Std |
|---|---|---|
| Model 1 | 37.87% | 0.0455 |
| Model 2 | 35.79% | 0.0335 |
| Model 3 | 40.02% | 0.0591 |

(c) Chromagram

Figure 3.1: Average Accuracy for all the features

**Fig.3.1** reports average accuracy and standard deviation of each model divided for features.

Best model for each feature is **Model 3**, with corresponding accuracy of 52.67% for MFCCs, 45.77% for Mel-Spectogram and 40.02% for Chromagram.

They are very low values, this means that prediction will not be so precise. To go more in depth, if we compare the accuracy and the validation accuracy of those best model for each of the extracted features, we can notice that the models are overfitting. This is due to the fact that there is a strong difference between training accuracy, which is relatively good, and validation accuracy which is very low and so, not good.



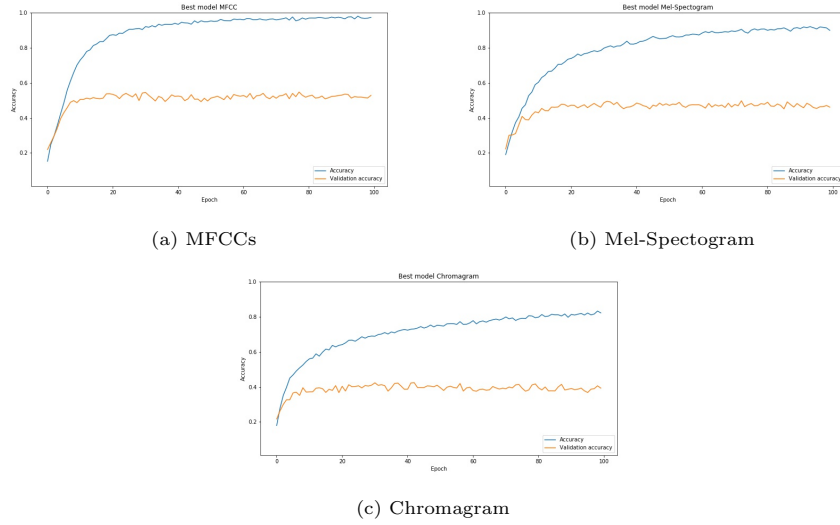(a) MFCCs

(b) Mel-Spectogram

(c) Chromagram

Figure 3.2: History for best model of each feature

This not optimal result can be due to the fact that our models have memorized files contained in the training set, without having learned how to correct map the features given in input to the respective class.

It was earlier stated in data description section, that it is difficult to visualise the difference between some of the classes. In particular, the following sub-groups are similar in shape: repetitive sounds for drilling and jackhammer, spikes for dog bark and gun shot. By plotting the confusion matrix for the best model (Model 3) for each features category we are able to see if also our model struggled to differentiate those classes. In confusion matrix rows represents true labels while columns those predicted by the model.

(a) MFCCs

| True Label \ Predicted | air_conditioner | car_horn | children_playing | dog_bark | drilling | engine_idling | gun_shot | jackhammer | siren | street_music |
|---|---|---|---|---|---|---|---|---|---|---|
| air_conditioner | 157 | 0 | 114 | 42 | 12 | 7 | 9 | 40 | 22 | 97 |
| car_horn | 2 | 159 | 8 | 15 | 5 | 4 | 1 | 4 | 2 | 21 |
| children_playing | 11 | 1 | 308 | 58 | 12 | 12 | 22 | 3 | 24 | 49 |
| dog_bark | 13 | 0 | 79 | 325 | 9 | 8 | 18 | 3 | 26 | 19 |
| drilling | 8 | 2 | 55 | 43 | 236 | 4 | 17 | 75 | 12 | 48 |
| engine_idling | 38 | 0 | 75 | 37 | 12 | 218 | 37 | 7 | 6 | 53 |
| gun_shot | 4 | 1 | 27 | 38 | 0 | 0 | 105 | 0 | 6 | 3 |
| jackhammer | 59 | 0 | 9 | 3 | 178 | 10 | 4 | 163 | 3 | 23 |
| siren | 7 | 0 | 25 | 57 | 3 | 9 | 3 | 1 | 262 | 26 |
| street_music | 25 | 4 | 72 | 29 | 24 | 8 | 10 | 8 | 25 | 295 |

(b) Mel-Spectogram

| True Label \ Predicted | air_conditioner | car_horn | children_playing | dog_bark | drilling | engine_idling | gun_shot | jackhammer | siren | street_music |
|---|---|---|---|---|---|---|---|---|---|---|
| air_conditioner | 165 | 13 | 15 | 34 | 60 | 38 | 11 | 19 | 10 | 135 |
| car_horn | 1 | 144 | 14 | 18 | 15 | 3 | 4 | 5 | 2 | 15 |
| children_playing | 13 | 12 | 228 | 55 | 25 | 37 | 4 | 7 | 67 | 52 |
| dog_bark | 17 | 9 | 57 | 316 | 18 | 6 | 2 | 3 | 37 | 35 |
| drilling | 16 | 4 | 43 | 42 | 203 | 3 | 36 | 53 | 19 | 81 |
| engine_idling | 86 | 7 | 77 | 4 | 9 | 197 | 23 | 2 | 21 | 57 |
| gun_shot | 8 | 0 | 26 | 1 | 15 | 12 | 89 | 8 | 9 | 16 |
| jackhammer | 12 | 0 | 44 | 0 | 218 | 0 | 98 | 54 | 1 | 25 |
| siren | 5 | 3 | 38 | 44 | 15 | 4 | 3 | 6 | 267 | 8 |
| street_music | 30 | 13 | 46 | 43 | 36 | 3 | 11 | 30 | 11 | 277 |

(c) Chromagram

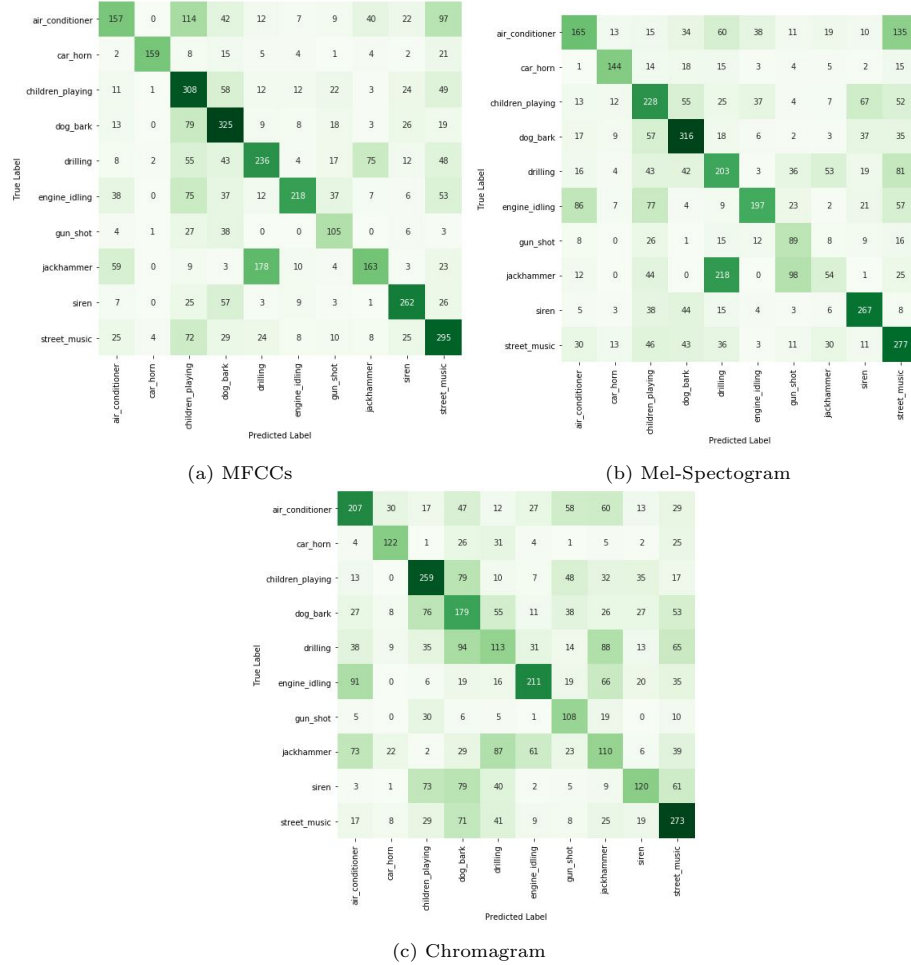| True Label \ Predicted | air_conditioner | car_horn | children_playing | dog_bark | drilling | engine_idling | gun_shot | jackhammer | siren | street_music |
|---|---|---|---|---|---|---|---|---|---|---|
| air_conditioner | 207 | 30 | 17 | 47 | 12 | 27 | 58 | 60 | 13 | 29 |
| car_horn | 4 | 122 | 1 | 26 | 31 | 4 | 1 | 5 | 2 | 25 |
| children_playing | 13 | 0 | 259 | 79 | 10 | 7 | 48 | 32 | 35 | 17 |
| dog_bark | 27 | 8 | 76 | 179 | 55 | 11 | 38 | 26 | 27 | 53 |
| drilling | 38 | 9 | 35 | 94 | 113 | 31 | 14 | 88 | 13 | 65 |
| engine_idling | 91 | 0 | 6 | 19 | 16 | 211 | 19 | 66 | 20 | 35 |
| gun_shot | 5 | 0 | 30 | 6 | 5 | 1 | 108 | 19 | 0 | 10 |
| jackhammer | 73 | 22 | 2 | 29 | 87 | 61 | 23 | 110 | 6 | 39 |
| siren | 3 | 1 | 73 | 79 | 40 | 2 | 5 | 9 | 120 | 61 |
| street_music | 17 | 8 | 29 | 71 | 41 | 9 | 8 | 25 | 19 | 273 |

Figure 3.3: Confusion matrices

In 3.3.(a) we plot the confusion matrix of Model 3 which takes as input Mel-frequency cepstral coefficients. As we can see, air conditioner is mostly confused with children playing and street music. Jackhammer is confused with drilling, but this is not surprising since they have similar shape (repetitive) if we look at their wave sound.

In 3.3.(b) we have the confusion matrix of Model 3 which takes as input Mel-spectogram. Here, air conditioner is largely confused with street music, as before, and also with drilling and engine. Jackhammer is highly confused with drilling, more than in the previous case where MFCCs were used as input.

In 3.3.(c) confusion matrix of Model 3 with Chromagram as input is displayed. Drilling here is largely confused with dog bark and jackhammer. We see an improvement in air conditioner performance with respect to use MFCCs as input for the same model.

# Conclusion

In this experiment three different models were created in order to classify urban sounds contained in the UrbanSound8K dataset. From each audio file three different features were extracted and then given as input to our different models. The performance of models developed in this project are not so high. Despite that, Model 3 is able to classify sound clips to a high degree of accuracy, with respect to the other two models.
MFCCs seems to be the best features that can be used to classify audio data with respect to the other two extracted features.
Overall, the best model in our experimental project for urban sound classification is Model 3 with MFCCs feature as input.

# Appendix

Link to Github Repository:

https://github.com/antonelladamico17/Urban-Sound-Classification