

Urban sound classification

Francesco Tomaselli

August 8, 2021

Contents

1	Introduction	2
2	Feature extraction	3
2.1	Dataset structure	3
2.2	First dataset	3
2.3	Extended dataset	4
3	Model definition	6
3.1	Neural network structure	6
3.2	Initial training set results	7
3.3	Hyperparameter tuning	7
4	Results and final remarks	10
4.1	Test set results	10
4.2	Possible improvements	10

1 Introduction

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 my 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.

The goal of this project is to build a neural network to classify audio files from the *UrbanSound8k* dataset [?].

This dataset contains audio divided in ten classes, each one representing a different type of city sound, for instance, we can find *car horns*, *dogs barking*, *sirens*, etc. A deeper discussion about the dataset is made at Subsection 2.1 on the following page.

The presented methodology is composed of three main parts. The first step is to extract relevant features from audio files, this is discussed on Section 2 on the next page.

The next step consists in composing and refining a neural network to classify the data obtained from the previous step. This part is discussed in Section 3 on page 6.

Lastly, results from the classification, namely accuracy and standard deviation among test sets, are presented in Section 4 on page 10.

Implementation details The project is developed in *Python* and some useful packages were used, in particular *Librosa* and *Dask* to extract features from audio files and *TensorFlow* and *Keras* to build the Neural Network. [?] [?] [?] [?] [?] Other libraries were used to complete some minor tasks, they are named in the following Sections.

The project folder is structured as follows:

- *src*: this contains the source code for the project. Each sub-folder contains code for a specific part of the processing. In particular, the *data* folder holds classes to extract features and to manage the dataset, the *model* folder contains the class to create the Neural Network, then *utils* stores utility functions to measure performances;
- *data*: here data is stored, there is a *processed* and *raw* sub-folders, where the first stores computed datasets and the second original data;
- *models*: trained models are saved here;
- *notebooks*: the folder contains the Jupyter notebooks used in the project, where code from *src* is practically used.

2 Feature extraction

This Section presents the original dataset structure and the steps followed to create training and test sets from it.

Note that the models mentioned in this section are three layers *multilayer perceptron* a reasonable number of neurons, trained for 100 epochs with default parameters for the *Stochastic Gradient Descent optimizer*. [?] [?]

The accuracy on the training is computed with a *cross-validation* approach. [?] Details about models and validation techniques used in the project can be found at Section 3 on page 6.

2.1 Dataset structure

The dataset contains ten folds of audio samples, each one about four seconds long. The samples are divided in ten classes among ten folds. The following table shows the classes and their relative frequency in the dataset.

Class name	Number of samples
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

The training set consists of the first four folds plus the sixth, the other folds create five different test sets. The table shows the numerosity of the different sets.

Dataset	Number of samples
Training set	4499
Test set 5	936
Test set 7	838
Test set 8	806
Test set 9	816
Test set 10	837

2.2 First dataset

Choosing the features to extract was challenging due to the inexperience working on audio files, thus some *research* on what features to choose was necessary [?].

The Librosa library provides many feature to choose from, to keep it simple, for the first try with this dataset, the extracted features are these three ones:

1. *Mel-frequency cepstral coefficients*;
2. *Chromagram*;
3. *Root-mean-square*.

Each feature consists of an array of arrays containing measurements. A series of functions were applied to each sub-array and results were concatenated in a final feature vector [?]. The functions applied are *minimum*, *maximum*, *mean* and *median* from the *Numpy* library [?].

This approach resulted in 132 components feature vectors.

Feature scaling After testing some Neural Networks on the first dataset the results were not promising. One of the reasons is the big difference in ranges among feature vector components.

To mitigate this effect a *StandardScaler* from *scikit learn* was applied [?][?]. The result is a dataset where each feature has more or less a distribution centered in zero with unit variance. This lead to an improvement on the results using the same model as before.

2.3 Extended dataset

To improve results on the training set new features are added, namely:

1. *Zero-crossing rate*;
2. *Roll-off frequency*;
3. *Spectral flux onset strength*.

After applying the same four functions to these three new arrays, a total of 12 new features are added to the dataset [?]. Scaling yield to promising results on the first dataset, so the same approach is applied to the extended dataset.

After testing a network on the new dataset results improved once again.

PCA Adding new features can lead to better results in the end but they all need to be useful to the model, so the extended dataset was subject of some experiments with feature selection, in particular *PCA* algorithm from *scikit learn* was applied [?].

The main idea is to select a reduced number of features from the total, without losing information. This approach often leads to better results, as useless features are discarded.

After some experiments with the number of features to select, 120 out of 144 were selected. This led to a small improvement on the training set, so this final reduced dataset is chosen to perform hyperparameter tuning.

3 Model definition

This Section presents how the networks used on the training set are structured, the second Subsection gives an overview of the performances on the different training sets, lastly, the Hyperparameter tuning phase is described.

3.1 Neural network structure

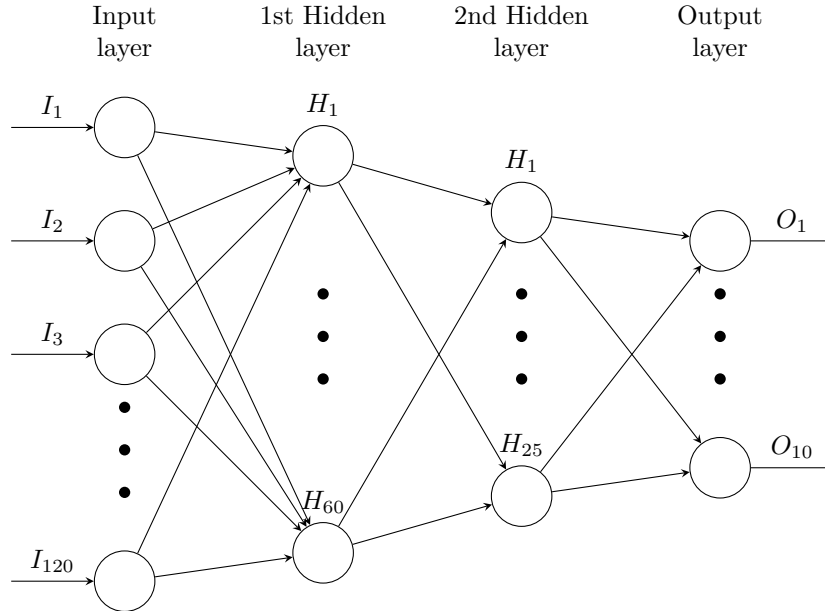
The starting point for the neural network structure was a reasonable network in terms of hidden neurons to prevent over-fitting, indeed a high number of units in the hidden layers would end up in learning too much from the dataset.

The rule of thumb followed to decide hidden neurons quantity is the following:

$$\#hidden\ neurons = \frac{2}{3} \#input\ neurons + \#output\ neurons$$

The next step was to decide the hidden layer number. To respect the number of hidden neurons, two hidden layers were considered. A more complex structure in terms of layers number would mean having a real small number of neurons per layer.

To give reference, this is the model used on the PCA training set, mentioned at the end of Subsection 2.3 on page 4, with 120 input features.



The activation function for the input and hidden layers is a *Relu* and the output one uses a *Softmax* [?] [?]. The loss used is the *Sparse Categorical Crossentropy loss* as it is suited for this kind of problems [?].

As discussed in the previous Section, models with this logic for construction were tested on the four training sets to choose the one to tune the final network on.

3.2 Initial training set results

Section 2 on page 3 talked about the four different training sets obtained from the dataset and, without going into details, stated that there was continuous improvement. This subsection presents the results obtained on the training set.

Note that all the random seeds used by Tensorflow were fixed to make results reproducible.

Cross validation To estimate performance on the training set cross-validation with 5 folds was used. Basically the dataset is divided into five folds, and a model is repeatedly trained on four and tested on one. The mean accuracy on the testing folds gives a hint about the model performance.

Results For each training set presented, a model was defined with the structure presented at the beginning of this Section, and these are the results:

Training set	Mean accuracy	St. deviation
132 features unscaled	0.0669	0.0488
132 features scaled	0.5201	0.0668
144 features scaled	0.5586	0.0842
120 features reduced with PCA	0.5608	0.0998

There is a huge improvement after scaling the training set, after that small refinements were made. As accuracy is the best on the last dataset, this was the one selected to perform the Hyperparameter tuning.

3.3 Hyperparameter tuning

Choosing the training set with PCA applied, led to the best results with cross validation, although the model was reasonable, it can not be the final one, as many parameters were left on their default value, for instance, learning rate and momentum of the optimizer were not touched.

The main goal now is to experiments with different model parameters to find the best one.

Grid vs Random search Two of the most commonly used strategies in Hyperparameter optimization are *Grid* and *Random Search* [?].

In both cases we define ranges of parameters to test different combinations, for instance, fixed the number of neurons, one could try to find the best combination of learning rate and momentum that optimize accuracy on the training set.

While similar, the two methodologies differs in the amount of exploration they do. The Grid search try all the possible combinations of parameters, while the Random approach fixes a number of iterations and picks an arbitrary combination each time.

Obviously the first one is more computationally expensive than the second, if we fix a small amount of possible iterations, but in theory it finds a better result than going the random route. Nonetheless the Grid Search can led to over-fitting, and in practice Random Search in preferred.

To get the best out of the two techniques, a first Random Search is performed on a larger parameter space, then, a Grid Search is used on much smaller ranges to optimize the previously found model.

Random search Starting from the model introduced in the previous Subsection, parameters are tweaked to find a better one. The optimizer used is the Stochastic Gradient Descent. The considered ranges for parameters for this run are:

1. *Neurons*: first and last layers stay the same, while the two hidden layers are tested with a number of neurons respectively equals to:

$$60 + 2i \text{ and } 25 + 2j, \text{ with } i, j \in \{-2, -1, 0, 1, 2\}$$

2. *Learning rate*: 0.001, 0.01, 0.1, 0.5;
3. *Momentum*: 0.0, 0.01, 0.1, 1.

An early stopper is used on the fit with 100 *epochs*, *batch size* is fixed at 32, and cross validation with 5 folds is performed on each combination. The total possible models are 400, but the search was performed with 100 iterations in total.

The best model found was the following:

- *Neurons*: 120 for input, 62 for the first hidden layer, 27 for the second, and 10 for output;
- *Momentum*: 0.1;
- *Learning rate*: 0.1.

Comparing the initial model with the one found now, a small improvement can be seen both in accuracy and standard deviation:

Model	Mean accuracy	St. deviation
Initial model	0.5608	0.0998
Random search result	0.5857	0.0715

Grid search The best model from the random search was then fine-tuned with a Grid Search, where ranges are in the proximity of the parameters found by the previous search.

Parameters this time was:

1. *Neurons*: first and last layer stay the same, while the two hidden layers are tested with a number of neurons respectively equals to:

$$62 + i \text{ and } 27 + j, \text{ with } i, j \in \{-1, 0, 1\}$$

2. *Learning rate*: 0.08, 0.1, 0.12;
3. *Momentum*: 0.08, 0.1, 0.12.

The total possible models was 81, the best found is the following:

- *Neurons*: 120 for input, 62 for the first hidden layer, 28 for the second, and 10 for output;
- *Momentum*: 0.12;
- *Learning rate*: 0.08.

As before, a comparison shows that Grid Search fine tuned the model to have better accuracy, unfortunately standard deviation increased by a small quantity:

Model	Mean accuracy	St. deviation
Initial model	0.5608	0.0998
Random search result	0.5857	0.0715
Grid search result	0.5877	0.0763

This last model is the one selected to evaluate test set performances.

4 Results and final remarks

This last Section shows the results obtained on the different test sets and possible improvements.

4.1 Test set results

The final model selected is the last one found on previous Section, found by the Grid Search. As stated on the first Section, the five test sets are made out of the folds number five, seven, eight, nine and ten.

The following are the results on those folds:

Test set	Accuracy
Fold 5	0.6731
Fold 7	0.6301
Fold 8	0.7184
Fold 9	0.6630
Fold 10	0.6906

Finally, the mean accuracy and standard deviation for the test sets are:

Mean accuracy	Standard deviation
0.6750	0.0293

4.2 Possible improvements

Testing different scalers,

testing different feature selection methods

testing various neural network structures

balancing classes