

# Engine Detection Classification

Rayden R. Dodd, IEEE

**Abstract**—This Document provides the development process of an engine detection model designed as a binary classifier to be able to tell when the sound input is a car engine idling or something else. Through a rigorous approach of iterative training and validation the model achieves the detection of car engines. This paper also touches on the challenges encountered, sound variability and model comparison.

**Index Terms**—Engine Sound Classification, Binary Classification, MFCCs, Machine Learning, ExtraTrees.

## I. INTRODUCTION

For our Major Design Experience project we were tasked

with creating an machine learning engine classifier model. The requirements for this model were that it had to constantly be listening to its surroundings and once it detected a car engine is idling it then has to run this sound clip through a machine learning model that is able to classify within 80% accuracy what brand of car it, as well as its next 2 brand predictions. This document will go over the engine detection part of this project. This includes talking about how the dataset for training was made, choice in model, optimization of the best model, and model testing.

## II. DATA SET CREATION

To create the dataset first we needed sound recordings of stuff that wasn't car engine idling. So first I found a GitHub repo [karolpiczak/ESC-50: ESC-50: Dataset for Environmental Sound Classification \(github.com\)](https://github.com/karolpiczak/ESC-50/) that had a dataset of different environmental sounds. This dataset contained 2000 5 second audio clips of various noises including but not limited to frogs, breathing, can opening, chainsaws, ext.

Task ID	Worker ID	Demographic	Category	Region	City	Worker prediction	Ground truth	Recording	Corner
705407956	1.00	63	POL	74	Italy	Chirping birds	Chirping birds	1:20002-A.ogg	
705407957	1.00	111	POL	51	Poland	Hand saw	Vacuum cleaner	1:000210-A.ogg	
705407958	1.00	120	POL	74	Italy	House cleaner	House cleaner	1:000211-A.ogg	
705407959	1.00	63	ROU	78	Latvia	Thunderstorm	Thunderstorm	1:03196-A.ogg	
705407960	1.00	22	BIN	2	Sarajevo	Car knock	Door knock	1:03136-A.ogg	
705407961	1.00	41	HUN	12	Tatransky	Can opening	Can opening	1:310404-A.ogg	
705407962	0.95	86	RUS	78	Tyumen	Drum	Drum	1:031365-A.ogg	
705407964	0.91	93	ESP	60	Valladolid	Door knock	Door knock	1:03199-A.ogg	
705407965	1.00	133	ESP	1	Sarajevo	Door knock	Door knock	1:031991-A.ogg	
705407966	1.00	76	ITA	7	Rome	Clapping	Clapping	1:030409-A.ogg	
705407967	1.00	81	ESP	29	Madrid	Clapping	Clapping	1:204089-A.ogg	
705407968	1.00	53	POL	74	Latvia	Clapping	Clapping	1:204090-A.ogg	
705407970	1.00	63	POL	74	Latvia	Clapping	Clapping	1:110537-A.ogg	
705407971	1.00	79	ITA	7	Rome	Thunderstorm	Thunderstorm	1:310521-A.ogg	
705407972	1.00	74	BIN	3	Sarajevo	Thunderstorm	Thunderstorm	1:310522-A.ogg	
705407973	0.95	86	RUS	78	Tyumen	Footsteps	Footsteps	1:115145-A.ogg	
705407974	0.95	53	ESP	60	Valladolid	Footsteps	Footsteps	1:115146-A.ogg	
705407975	1.00	76	ITA	7	Rome	Fireworks	Fireworks	1:115147-A.ogg	
705407976	1.00	58	ESP	29	Madrid	Fireworks	Fireworks	1:115148-A.ogg	
705407977	1.00	81	ESP	29	Madrid	Clapping	Clapping	1:11520-A.ogg	
705407978	1.00	91	BGR	54	Shumen	Clapping	Clapping	1:11521-A.ogg	
705407979	1.00	109	ESP	1	Barcelona	Clapping	Clapping	1:11522-A.ogg	
705407980	1.00	58	ROU	7	Rome	Airplane	Airplane	1:11601-A.ogg	
705407981	1.00	76	ITA	7	Rome	Mosquito tick	Mosquito tick	1:11610-A.ogg	
705407982	1.00	81	ESP	74	Latvia	Running water	Running water	1:11611-A.ogg	
705407983	1.00	81	ESP	29	Madrid	Train	Train	1:119125-A.ogg	
705407984	0.95	81	ESP	39	Madrid	Drum	Drum	1:119126-A.ogg	
705407985	0.85	51	IND	2	Hyderabad	Water drops	Water drops	1:12053-A.ogg	
705407986	1.00	104	GRI	HS	Longisland	Water drops	Water drops	1:12054-A.ogg	
705407987	1.00	70	ITA	7	Rome	Water drops	Water drops	1:12055-A.ogg	
705407988	0.94	112	BIN	1	Sarajevo	Water drops	Water drops	1:12056-A.ogg	
705407989	1.00	120	ESP	29	Madrid	Church bells	Church bells	1:13572-A.ogg	
705407990	0.95	86	RUS	78	Tyumen	Church alarm	Church alarm	1:13573-A.ogg	
705407990	0.95	86	RUS	78	Tyumen	Clock alarm	Clock alarm	1:13013-A.ogg	
705407991	1.00	56	SAR	0	NYC	Ring	Ring	1:12057-A.ogg	
705407992	1.00	63	POL	74	Latvia	Keyboard typing	Keyboard typing	1:1377-A.ogg	
705407993	1.00	81	IND	2	Hyderabad	Keyboard typing	Keyboard typing	1:1378-A.ogg	
705407993	1.00	104	GRI	HS	Longisland	Clock alarm	Clock alarm	1:34242-A.ogg	

**Fig. 1.** This is a sample of the raw dataset with different classes shown with the audio file and where the audio was taken from.

Next, I needed to add our car engine recordings into the dataset. First of all, our recordings needed cleaning. For most of our car engine recordings the first 1-2 mins is just background noise until we turn the car on. Thus, using a python script that could detect the sudden change from a low volume recording of nothing to a loud engine starting I was able to cut the audio files into two background noise and engine idling. Because of this we had more environmental noise to add to our data set as well. Now that the files were cut to only include engine idling, we had to cut down each video into 5 second segments as the rest of the data set were also 5 second audio clips. Using another python code, I was able to chop all of the audio clips into 5 second segments leaving us with 2000 unique car engine idling clips. This was perfect as we had 2000 non engine idling clips and thus this made out database balanced which is beneficial in model training.

## III. FEATURE EXTRACTION

To be able to train a model on this database, we first needed to convert the .WAV and .ogg audio files into numerical features that a machine learning algorithm can interpret. To this end, we utilized Mel-Frequency Cepstral Coefficients (MFCCs), which are a type of feature widely used in the processing of audio signals, particularly in the context of speech and music. MFCCs effectively represent the short-term power spectrum of a sound, based on a linear cosine transform of a log power spectrum on a nonlinear mel scale of frequency. While spectrograms also provide a visual and numerical representation of the frequency spectrum over time, MFCCs are used to distill this information into a smaller set of features that capture the most important aspects of the sound for the purpose of auditory perception. This condensation into a more compact representation makes MFCCs more amenable to machine learning models, which can efficiently process and classify audio data based on these reduced features without a significant loss of information critical to distinguishing between different sounds or voices.

By employing MFCCs, we transformed the audio features into structured NumPy arrays, enabling their integration with scikit-learn machine learning models.

## IV. SELECTION OF BINARY CLASSIFICATION MODEL

### A. Testing Methodology

When coming up with a binaray classification model to train our data on our Subject Matter Expert (SME) suggested that we use Random Forest. To test this I used a list of different binary classification models to compare and see which model would give the highest accuracy. These models include:

- 1) Adaptive Boosting Classifier (AdaBoostClassifier)
- 2) Bagging Classifier (BaggingClassifier)
- 3) Decision Tree Classifier (DecisionTreeClassifier)

- 4) Extremely Randomized Trees Classifier  
(ExtraTreesClassifier)
- 5) Gaussian Naive Bayes (GaussianNB)
- 6) Gradient Boosting Classifier  
(GradientBoostingClassifier)
- 7) K-Nearest Neighbors Classifier  
(KNeighborsClassifier)
- 8) Logistic Regression (LogisticRegression)
- 9) Multinomial Naive Bayes (MultinomialNB)
- 10) Passive Aggressive Classifier  
(PassiveAggressiveClassifier)
- 11) Perceptron (Perceptron)
- 12) Quadratic Discriminant Analysis  
(QuadraticDiscriminantAnalysis)
- 13) Random Forest Classifier (RandomForestClassifier)
- 14) Ridge Classifier (RidgeClassifier)
- 15) Stochastic Gradient Descent Classifier  
(SGDClassifier)
- 16) Support Vector Machine (SVM)

To compare these models fairly, I employed two distinct cross-validation techniques: Leave One Out Cross Validation (LOOCV), which splits the dataset such that each instance serves once as the test set while the remainder forms the training set, ensuring every data point is utilized for validation, and K-Folds Cross Validation, dividing the dataset into for our case 10(K) equal parts, or "folds," and then iteratively uses one fold for testing and the remaining for training, with results averaged for thorough evaluation. After each iteration of LOOCV and K-Folds, I recorded the Accuracy, Recall, Precision, and F1 metrics, allowing these performance indicators to be averaged at the end of training to provide a comprehensive assessment of each model's effectiveness.

### B. Results

The results of this showed that the best model out of the 16 different models was ExtraTreesClassifier which had a 99.02% accuracy in the LOOCV and in second place Random Forest with 98.25% accuracy. The worst performing model was MultinomialNB 53.03%. The ExtraTreesClassifier performed better, most likely due to the fact that it makes use of extreme randomness in selecting splits for each node in the decision trees. Unlike Random Forest, which seeks the most optimal split among a random subset of the features, ExtraTreesClassifier randomly selects the split points, resulting in a higher level of model variance and diversity. This approach can be particularly effective in complex classification tasks like distinguishing engine sounds, as it may better capture the nuanced differences between classes. On the other hand, the MultinomialNB performed so poorly because it operates under the assumption that all features are independent of each other—an assumption that doesn't hold true for audio data. Audio signals, especially features extracted for machine learning like MFCCs, exhibit significant correlation between features, which likely led to the algorithm's inability to accurately model the data and resulted in its low accuracy.

Comparing these two models against each other on the same set of data we get the following results:

TABLE I  
F1 METRICS FOR EXTRATREES

Metric	Class 0	Class 1	Macro Avg	Weighted Avg
Precision	0.98	0.98	0.98	0.98
Recall	0.99	0.98	0.98	0.98
F1-Score	0.98	0.98	0.98	0.98
Support	550	448	998	998

TABLE II  
ACCURACY PERCENTAGES FOR EXTRATREES

Overall Accuracy	98.30%
K-Fold Accuracy	0.99 (+/- 0.01)
LOOCV Accuracy	0.99 (+/- 0.20)

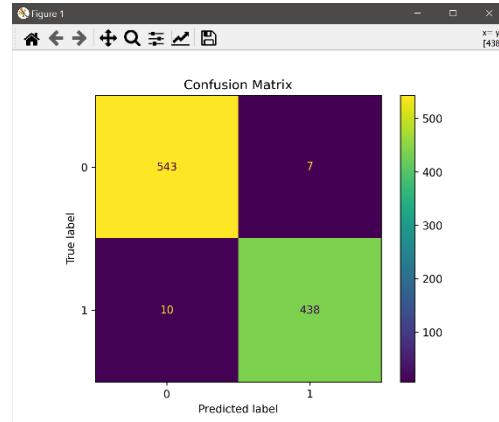


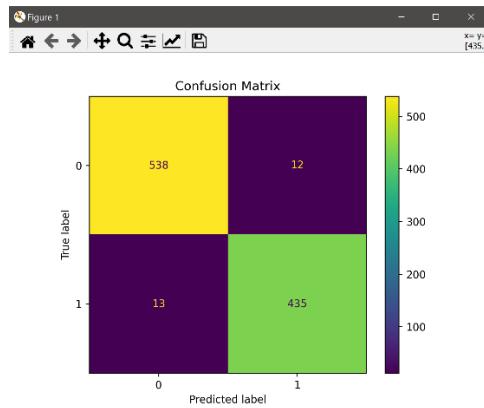
Fig. 2. Confusion matrix for ExtraTrees

TABLE III  
F1 METRICS FOR RANDOM FOREST

Metric	Class 0	Class 1	Macro Avg	Weighted Avg
Precision	0.98	0.97	0.97	0.97
Recall	0.98	0.97	0.97	0.97
F1-Score	0.98	0.97	0.97	0.97
Support	550	448	998	998

TABLE IV  
ACCURACY PERCENTAGES FOR RANDOM FOREST

Overall Accuracy	97.49%
K-Fold Accuracy	0.98 (+/- 0.01)
LOOCV Accuracy	0.98 (+/- 0.27)



**Fig. 3.** Confusion matrix for Random Trees

These results show that overall ExtraTrees is slightly better than Random trees due to the better accuracy and F1-scores.

#### V. MODEL OPTIMIZATION

To optimize the ExtraTree model I used the sklearn documentation page to find all of the available parameters for the model and using GridSearch Cross Validation built into sklearn after some time it showed the best parameters which were : n\_estimators=250, max\_features = 7, min\_samples\_leaf =1, random\_state=42, n\_jobs=-1. This brought the single run accuracy from 98.3% to 98.4% so only a small bump in performance

#### V. CONCLUSION

In conclusion, developing and comparing binary classification models to detect engine sounds has shown us how challenging and intricate this task can be. The ExtraTreesClassifier stood out with its impressive accuracy and ability to avoid overfitting, demonstrating the value of using a more random approach in decision trees. This project didn't just test different machine learning models; it also gave us insights into how to fine-tune these models for better results. Additionally, the poor performance of the Multinomial Naive Bayes model highlighted the need to choose the right machine learning strategy, considering the data's specific traits, like the connection between audio features. Our careful work on preparing the data, picking out features, testing models, and refining them has set a solid foundation for future work in audio classification. This could open up new possibilities in automated engine diagnostics and more. This project was an important part of our major design experience, contributing useful knowledge and techniques to the field of machine learning.