

# Machine Learning: Dog Bark Identification

Bryan Khuu, Pawel Juda. Spring 2020

This research paper introduces and executes approaches to identifying whether a dog has barked in an audio file using machine learning statistical models.

The project repository is located on GitHub at [paweljuda/MLDogBarkIdentification](https://github.com/paweljuda/MLDogBarkIdentification)

## Part A: Introduction

### Background and Rationale

Speech recognition technology has become increasingly advanced within the passing decade. The increase in computing power of our personal devices has allowed for larger clusters of data to be obtained and managed. This is particularly significant as voice recognition is a computationally heavy process. Audio signals contain numerous features which can be extracted for data. Human voice for example, can be analyzed for unique frequencies, bandwidth, and other features which make it distinguishable from other audio or people.

In the following sections, audio samples will be processed and analyzed to determine if similar voice recognition technology can be implemented to recognize dog speech.

### Data Sourcing

To improve reliability and prevent overfitting, the models are trained on multiple breeds. A reliable dataset which was found and used comes from kaggle, provided by the user mmoreaux. The link to the dataset is <https://www.kaggle.com/mmoreaux/audio-cats-and-dogs> and contains 113 dog audio files and 167 cat audio files. Although cat audio files do not fall within the scope of the project's goal, they contribute to a part of the negative train/test set. The second dataset was provided by a handy dandy bark maker, Leo, a personal dog assistant for this project. This dataset provides 14 longer audio files, containing numerous barks. Lastly, a third dataset was created from a mix of audio files from various YouTube videos.

## Part B: Preprocessing and Feature Extraction

### Data Preprocessing

Before any feature can be extracted, the audio files must go through pre-processing. This process sets up and modifies audio samples to make it easier and more useful for analysis.

The first step of pre-processing is noise reduction. This helps to eliminate unwanted background audio. This process is visualized in the spectrograms shown in figure 1 and 2. Spectrograms plot frequency and signal strength with respect to time. As can be seen, this process only preserves

significant events, where a possible dog bark might have occurred.

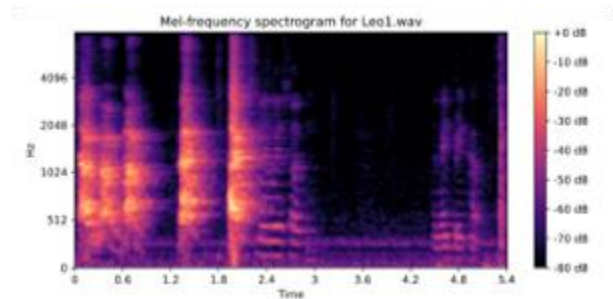


Figure 1: Original mel-spectrogram of a sample audio input.

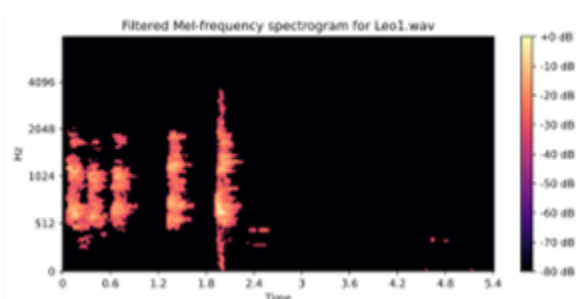


Figure 2: Amplitude filtered mel-spectrogram.

The new spectrogram, seen in figure 2, displays a sharper edge of where a significant audio event happens. These edges are used to help splice long audio files into multiple smaller files which only contain the desired audio. One spliced sample is shown in figure 3.

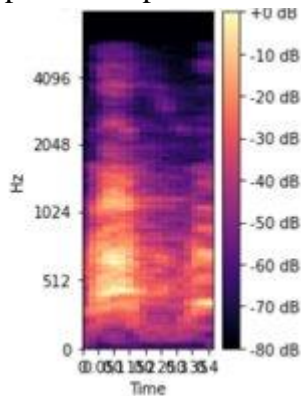


Figure 3: Sample of a spliced region from mel-spectrogram determining a significant event.

The next step in this process is to classify the newly spliced audio files. This is done by individually examining each audio split and labelling it as 1 if it contains a dog audio, or 0, if it contains no dog audio. This is known as machine supervised learning, a process where input and corresponding outputs are manually written into the system. These data points are used to train the model.

1111	Leo12-13.wav	1
1112	Leo12-14.wav	0
1113	Leo13-0.wav	0
1114	Leo13-1.wav	1
1115	Leo13-2.wav	1
1116	Leo13-3.wav	1
1117	Leo13-4.wav	1
1118	Leo13-5.wav	1
1119	Leo13-6.wav	0
1120	Leo13-7.wav	1
1121	Leo13-8.wav	1
1122	Leo13-9.wav	1
1123	Leo13-10.wav	1
1124	Leo13-11.wav	1
1125	Leo14-0.wav	0
1126	Leo14-1.wav	1
1127	Leo14-2.wav	1
2086	monkey-43.wav	0
2087	monkey-44.wav	0
2088	monkey-45.wav	0
2089	monkey-46.wav	0
2090	monkey-47.wav	0
2091	monkey-5.wav	0
2092	monkey-6.wav	0
2093	monkey-7.wav	0
2094	monkey-8.wav	0
2095	monkey-9.wav	0
2096	pig noise-0.wav	0
2097	pig noise-1.wav	0
2098	pig noise-2.wav	0
2099	pig noise-3.wav	0
2100	Recording_1-0.wav	0
2101	Recording_1-1.wav	0
2102	Recording_1-10.wav	0

Figure 4: Samples of classified spliced audio file data.

## Feature Extraction

Feature extraction is performed on all the samples and stored in the data frame using a schema as shown in figure 5. The features used are mel-spectrogram, MFCC, chromagram, zero-crossing rate (ZCR), root-mean square (RMS), spectral centroid, spectral bandwidth, spectral roll-off frequency, spectral flatness, and spectral contrast. The library, LibROSA, was used for extracting these features.

	File Name	Mel-Spectrogram	MFCCS	Chromagram	ZCR	RMS	Centroid	Bandwidth	Roll-off	Flatness	Contrast	Class
0	cat_1_audio_sp118.wav	[[[0.48189516, 0.5245578, 0.37867386, 0.7248...]]	[[[-174.48797, -132.7548, -148.42489, -189.775...]]	[[[0.44264818, 0.1191244, 0.61389227, 0.514214...]]	[[[0.4971288625, 0.4697559375, 0.49328515425,...]]	[[[0.401089317779918, 0.4514866828619934, 0...]]	[[[1686.178392642087, 2853.1176544689984, 2866...]]	[[[1638.459558867188, 1735.163668761884, 1767...]]	[[[3477.4328666878, 3832.49151425, 3865.389968...]]	[[[0.4812842331, 0.461528282, 0.48017851405,...]]	[[[24.7893211459683, 22.4732186049534, 24.53...]]	0.0
1	cat_1_audio_sp115.wav	[[[0.413389145, 0.47323262, 0.512524...]]	[[[-201.84723, -187.47483, -191.96774, -138.59...]]	[[[0.41574832, 0.47988884, 0.50622975, 0.18932...]]	[[[0.4234375, 0.40787171875, 0.4388859375, 0...]]	[[[0.462714176895724, 0.1877893747893248, 0...]]	[[[747.4391348128279, 458.4768418823789, 0...]]	[[[1167.77667251272, 948.179882468818, 927.2...]]	[[[1881.7994453116, 823.46328131, 796.728516...]]	[[[0.4888854-05, 0.4615288-05, 0.2467524-05,...]]	[[[124.11879477932192, 14.41388813292368, 26...]]	0.0
2	cat_1_audio_sp112.wav	[[[0.4338883, 0.2737826, 0.474776135, 0.46258...]]	[[[-213.38127, -232.52164, -248.46269, -245.54...]]	[[[0.24172398, 0.11887581, 0.14831159, 0.4222...]]	[[[0.41953335, 0.42734575, 0.444921875, 0.4008...]]	[[[0.472155118619237, 0.1213867744413888, 0...]]	[[[488.4946163143227, 569.463687734897, 0...]]	[[[1972.481120861489, 839.46319134876, 988.1...]]	[[[1838.1615234375, 1285.4238115, 984.3945325...]]	[[[0.4495164-05, 0.4615288-05, 7.35913984-05,...]]	[[[23.2781678784224, 24.49866468772726, 23.49...]]	0.0
3	cat_1_audio_sp113.wav	[[[0.316785, 0.4078789, 0.404927395, 0.42748...]]	[[[-157.48764, -156.87764, -205.12647, -159.15...]]	[[[0.7562376, 0.4346786, 1.0, 1.0, 0.4222...]]	[[[0.424424825, 0.464921875, 0.40382421875, 0...]]	[[[0.4318574442279888, 0.4212584537042389, 0...]]	[[[147.7377538911472, 942.4728839523967, 1227...]]	[[[1327.788156781118, 1349.7932737683247, 133...]]	[[[1342.6981266875, 1285.4238115, 984.3945325...]]	[[[0.48881864629, 0.4881154394, 3.97286424-05,...]]	[[[122.42182908786583, 11.41528623644518, 13...]]	0.0
4	cat_18_audio_sp118.wav	[[[0.1578544, 0.4527843, 0.49912259, 0.1881...]]	[[[-243.44387, -232.38843, -238.43538, -233.98...]]	[[[0.424899921, 0.4382196475, 0.448979957, 0.47...]]	[[[0.43515625, 0.448212875, 0.46382421875, 0...]]	[[[0.421628742441374445, 0.423884160872464,...]]	[[[3761.4287552268266, 1868.19878835555, 1837...]]	[[[1834.689149683893, 1933.1512673751213, 195...]]	[[[1243.84755426, 3811.77893315, 4184.4748846...]]	[[[0.4828648518, 0.46173266, 7.94916164-05, 0...]]	[[[15.11619317895849, 4.118811418677231, 19.2...]]	0.0
5	cat_18_audio_sp115.wav	[[[0.48888459, 0.41332124, 0.417852475, 0.4...]]	[[[-149.1658, -145.88129, -141.8291, -158.947...]]	[[[0.41648455, 0.44948235, 0.43782557, 0.47789...]]	[[[0.4083128125, 0.143878125, 0.19785557, 0...]]	[[[0.4518481156947181, 0.46151993384382189, 0...]]	[[[2316.1864992181952, 2838.1857942127894, 278...]]	[[[1583.151721138888, 1534.989293773312, 158...]]	[[[4382.8988159375, 4434.396486...]]	[[[0.4838818864, 0.48878511, 6.6585874-05, 7.5...]]	[[[9.451744997836214, 9.77852798887794, 18.28...]]	0.0
6	cat_18_audio_sp112.wav	[[[0.42238536, 0.4527843, 0.48888459, 0.482...]]	[[[-146.1648, -136.27327, -136.42848, -247.78...]]	[[[0.4935965, 0.394241, 1.0, 1.0, 0.44743285, 0.4...]]	[[[0.498828125, 0.143878125, 0.19785557, 0...]]	[[[0.43184895178293888, 0.4394861839988328,...]]	[[[2487.645971488362, 2832.42770284643, 2747...]]	[[[1888.1648388131862, 1738.43473878817, 1553...]]	[[[4888.4971893375, 4749.486421875, 4354.27346...]]	[[[0.415736377, 0.482464646, 6.4045264197, 0...]]	[[[18.43718812815424, 11.39481321429648, 26...]]	0.0
7	cat_18_audio_sp113.wav	[[[0.47843236, 0.4527843, 0.467854525, 0.48...]]	[[[-204.88387, -201.17859, -209.13165, -176.14...]]	[[[0.4347583, 0.5276421, 0.23468718, 0.2278135...]]	[[[0.463426875, 0.143878125, 0.19785557, 0...]]	[[[0.42542899353439546, 0.4615178125, 0.47621615848312455,...]]	[[[1456.158994838998, 1849.4638819889996, 1888...]]	[[[1583.54413573388, 1473.58879936865, 1413...]]	[[[1746.77734378, 3871.421323815, 3892.944335...]]	[[[0.4834687812, 0.48878511, 6.4045264197, 0...]]	[[[26.44831548852647, 18.748118131286779, 19.82...]]	0.0
8	cat_18_audio_sp114.wav	[[[0.42744785, 0.4527843, 0.48888459, 0.48...]]	[[[-219.26818, -214.4248, -214.23637, -208.189...]]	[[[0.2192137, 0.10988254, 0.08819181, 0.4754...]]	[[[0.41662189375, 0.4615178125, 0.478125, 0...]]	[[[0.4348893518473136, 0.438672237353873634, 0.47621615848312455,...]]	[[[1363.7981366177163, 1323.8479844142438, 138...]]	[[[1341.958483963628, 1263.4541188141294, 114...]]	[[[1267.3878, 3871.421323815, 1993.821898825...]]	[[[0.48881739, 0.48878511, 2.94799154-05, 0...]]	[[[10.3312643296612, 18.13813213286779, 19.31...]]	0.0
9	cat_18_audio_sp113.wav	[[[0.486464375, 0.4527843, 0.48888459, 0.48...]]	[[[-207.73629, -154.16131, -168.48153, -145.08...]]	[[[0.4448833, 0.38843426, 0.2378833, 0.146981425...]]	[[[0.4611234375, 0.4615178125, 0.478125, 0...]]	[[[0.42451873892713887, 0.4615178125, 0.47621615848312455,...]]	[[[1363.8896878972524, 1327.363791218279, 138...]]	[[[1354.444943156845, 1482.498126443124, 1508...]]	[[[1293.3861328125, 3432.421323815, 3822.5435...]]	[[[0.48881739, 0.48878511, 6.4045264197, 0...]]	[[[7.294123754828811, 8.41398511648849, 16.84...]]	0.0

Figure 5: First 10 rows of dataset dataframe holding extracted features and file classification

## Part C: Model Setup and Results

### Baseline Method: Logistic Regression

Logistic regression is a basic way of separating and classifying samples. This method can be visualized using the 3D and 4D models generated below. There are areas where similarly classified data points are clumped together. Using logistic regression, these data points can be separated using boundaries like hyperplanes. Data points which land on same side of the boundary get the same determined classification. This logic also holds true for higher dimensional models as more features are included, though it becomes impossible to visualize. The logistic regression model utilized 7 of the features extracted: RMS, ZCR, spectral centroid, bandwidth, roll-off, flatness and contrast. These features were used to calculate corresponding  $\beta$  values and an intercept.

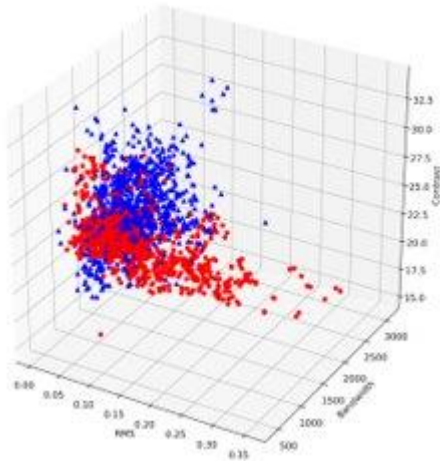


Figure 6: 3D Feature model, RMS, BW, Contrast

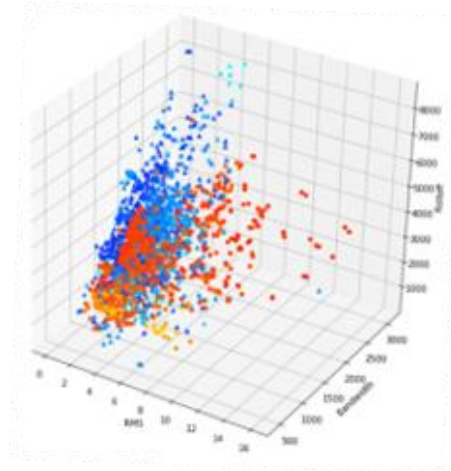


Figure 7: 4D Feature model, RMS, BW, Contrast, Roll-Off

The sklearn library provides a convenient function which calculates beta values with which a hyper plane can be constructed for the dataset. After separating the samples into a train and test set, these were the results: accuracy (76-81%), recall (~73%), precision (~77%). Utilizing an offset from  $-0.5$  to  $0$  for the hyperplane realized even greater accuracy, though it also runs small risks of decreasing accuracy. The following values were found using an offset of  $-0.5$ : accuracy (75-84%).

### Convolutional Neural Network (CNN)

A second method of approach are convolutional neural networks (CNN). The results of the linear regression model were used as a baseline to compare against our second model choice, a CNN, whose architecture is shown in figure 8. The is composed of two layers of 2D convolutions and max pooling + batch normalization. The first 2D convolution utilizes a  $5 \times 5$  kernel whilst the second one is a  $3 \times 3$  kernel. After the two layers the result is flatten and passed through a two hidden layer neural network with a final activation layer utilizing a sigmoid function on the basis that the CNN is solving a binary classification problem.

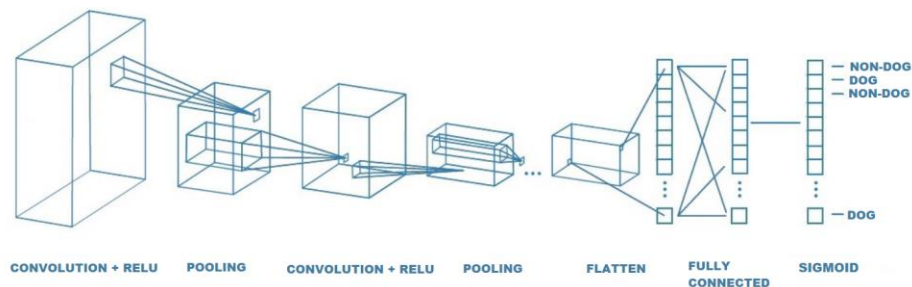


Figure 8: Convolutional Neural Network Architecture

One of the issues with using the mel-spectrogram as an image input for the CNN was that all image inputs need to be of the same size. This was fixed by padding spectrograms of all samples with the global minimum to a uniform maximum width (heights of spectrograms were the same). The padded spectrograms were then normalized to a value between 1 and 0 to improve performance of the CNN. A sample of the resulting padded and normalized spectrograms is shown in figure 9 below.

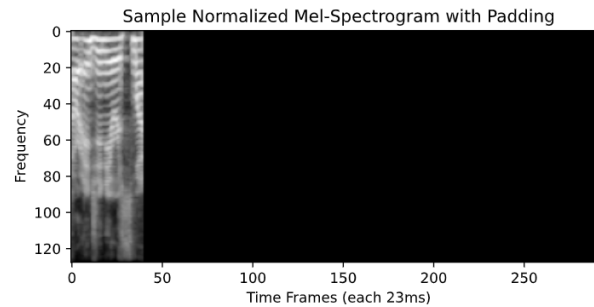


Figure 9: Example normalized and padded mel-spectrogram.

The results of the CNN model predictions are shown in table 1 below where they are compared to the results of the baseline logistic regression.

Method	True Pos	True Neg	False Pos	False Neg	Precision	Recall	Accuracy
<b>LogReg</b>	33%	45%	10%	12%	77%	73%	78%
<b>CNN</b>	40%	57%	0.5%	3%	99%	93%	97%

Table 1: Model Performance Comparisons

## Part D: Conclusions and Future Work

The results of logistics regression and convolutional neural network were strong. As a baseline, logistic regression got a fair 80%. The CNN provided a much stronger approach reaching a whopping 97% accuracy with a 99% precision. However, its validation accuracy was unstable, as seen in figure 10, which will require changes to the architecture, learning rate and decay. Additionally, the environment used (WSL Ubuntu 16.04) for training the CNN did not support CUDA and therefore an NVIDIA GPU was not recognized, therefore setting a limiter on the number of epochs used, which was 30 where it would have otherwise been set to 100. More epochs could stabilize the validation accuracy.

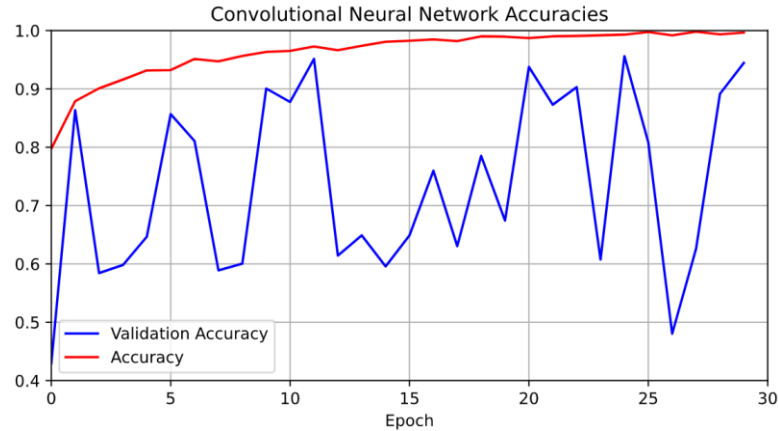


Figure 10: CNN Model Accuracies per Epoch

One improvement that may have improved results is a better splicing of the samples. There are cases where dog barks occur too closely together for the noise reduction method to finely separate the events. The dataset would produce more viable features if they were spliced apart into singular dog barks. This can be better achieved using an edge filtering convoluting kernel. As seen in figure.11 and figure.12 below, the filter makes each peak more distinct. The edge filtering convolutional kernel used was  $[-2, -1, 1, 4, 1, -1, -2]$ .

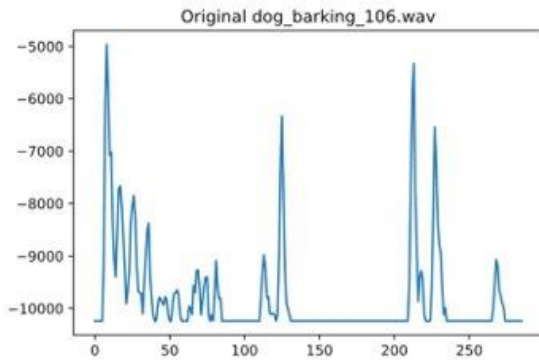


Figure 11: Original audio file's sum of mel-spectrogram columns

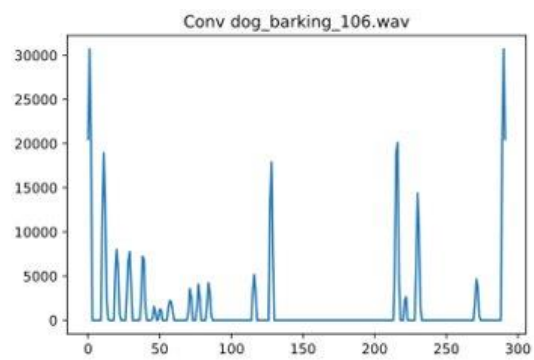


Figure 12: Convolution of audio file's sum of mel-spectrogram columns

For future work, there are many possible tasks. One thing that would be fun to implement is a machine learning model that can identify a specific breed or even a specific dog. Another task could include trying to trick the model into marking a fake bark as a real one. This could be made a person trying to imitate a dog. Lastly, an option for using our trained models on real life data would include using an open microphone to record audio and in real time predict if a bark occurred and take recordings of the number of barks as well as the timings.