# A music detector for hearing aids

Audio explorers - Software challenge
Herring Oats Team
https://github.com/indianatoms/AudioExplorers

Piotr Samir Saffarini          Tomasz Mariusz Koziak          Aleksander Nagaj

**Abstract**

In this project, we are trying to make the most efficient music binary classifier to hearing aids devices. Many solutions were evaluated in terms of their accuracy, memory usage and computational time. During the experiments, we discovered AI algorithms perform better than the conventional machine learning ones. We implemented the Convolutional Neural Network (CNN) that gave 98% accuracy on the validation data. We also examined the possibility of making the input data more compact by extracting different features based on log-mel spectograms as well as performing Principal Component Analysis (PCA). Finally, we presented the ideas of improving the model.

## 1    Introduction

Hearing aids help people around the world. Providing patients with the most advanced and reliable technology is one of the most important goals of this field. Noise muting with respect to relevant sounds like music or conversations still proves to be a current research topic [7]. In this project we focused on recognizing if music is present in an audio snippet. That classifier could be useful to later amplify this sound or to attenuate the noise when music is not played 1. Approach that we decided to take in this project was to classify two seconds samples in a form of 79 log-mel spectrograms. For that purpose we benachmarked many existing state of the art algorithms with the idea of finding the most appropriate one. Additionally, we took an effort of optimizing the classification scheme to fit the memory constraints of hearing aids.
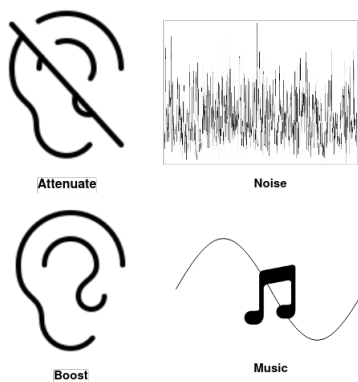


Figure 1: Overview of the desired solution

### 1.1    Problem Overview

The goal of the project was to develop a binary classifier that is able to recognize music to enhance the user experience of hearing aid devices and to make them more autonomous. Low computational power as well as limited memory of the device are the critical aspects that we had to take into consideration while doing this project.

**Data** was already converted from 2s audio samples to 79 log-mel spectograms (each consist of 25ms). The array was as followed: 10 500 samples x 30 frequencies bands x 79 time frames.

## 1.2 Project constraints

**Memory:** According to instructions, the device cannot store 2sec of audio, so more data efficient way is necessary. One 2sec sample is the size of 79x30. We assumed two variants: First is that it is feasable to store whole 2 sec sample with 79 log-mel spectograms, and second that more compact way of storing is necessary. We considered a real time computing variant which is not necessarily frame-by frame (1x30). Grouping e.g. 6 frames with PCA component reduction to 5 would give the same size (6x5 = 1x30) as frame-by-frame but with more time variation to the classification input.

**Computational Power:** On top of memory constraint, the device cannot perform complicated tasks due to CPU and battery limitations. Frame-by-frame computing is in contrary to this constraint, so we decided to also group frames into 5-10 and classify them all at once. This is where first variant (classifying 2 sec samples) has the advantage over the frame-by-frame computing.

# 2 Investigated Solutions

## 2.1 Getting to know the data

Before experimenting with the potential solutions to the problem, we run statistical analysis of the data in the format as one continuous stream of frames (samples*timeframes x frequencies). When comparing music data with other data, basic statistical features of each frequency bands were very close to each other (Figure 6. This indicated the usual frame-by-frame music recognition may not be feasible, especially in conventional machine learning algorithms. Bearing that in mind, instead of pure frame-by frame (1x30) we decided to investigate the whole spectrograms (79x30) or part of it, as it seemed that the detection of music samples is more possible with that approach.
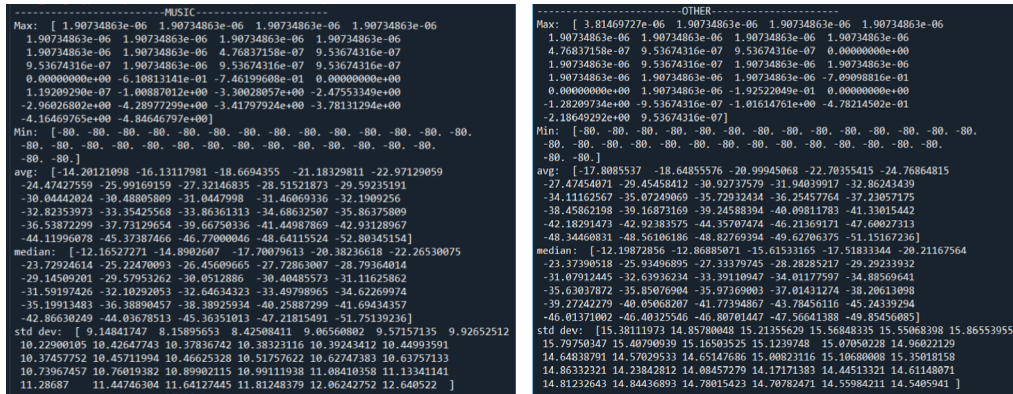


Figure 2: Basic statistical features of each frequency band every frame

## 2.2 Baseline models

We tried implementing many various baseline classification models with different parameters to get the feeling how well they perform on our data. We discovered that in some machine learing models PCA components reduction gave better end results in terms of accuracy and/or size of prediction model.

**Kmeans** We used Kmeans for classifying the data with 2 centroids which would be 'music' and 'other'. Frame-by frame classification accuracy was like a coin-toss probability and PCA dimension reduction did not produce better results. We did not include the results of classifying whole sample as it was too big to compute on our computers, so definitely is not suitable for small hearing aids devices.

**Logistic Regression** Although it performed similar to Kmeans, logistic regression model proved to be very lightweight and provided 67% accuracy using sample by sample approach.

**Random Forest** Among all conventional machine learning models (excluding AI) we have tested [5] Random Forest gave the best results. Unfortunately, the memory size is too big to implement in hearing aids industry. Even after compression the model was over 300MB.

**CNN** Convolutional Neural Network gave the best preliminary results in testing. As there are many different AI solutions besides CNN, we decided to further explore this path.

| Classification Model | frame-by-frame | sample-by-sample |
|---|---|---|
| KMeans | 50% accuracy 6.3MB model | Model too big |
| Logistic Regression | 49% accuracy 1KB model | 67% accuracy 3KB model |
| Random Forest | 88% accuracy 2GB model | Model too big |
| CNN | 95% 10 frames with PCA total 37 features 1MB model | 98% accuracy 3MB model |

Table 1: Baseline models overview

# 3 Final Model

In the final solution we decided to implement the convolutional neural network (CNN) (Figure 3) as it achieved the best values in terms of accuracy (98%) when being ran on all 79 log-mel spectrograms in each sample. The network was trained on 13440 samples and validated on 3360 samples with a *batch size* of 32. It consists of two *convolutional* layers with kernels of size $3 \times 3$; each is followed by $2 \times 2$ *max-pooling* and 0.2 *dropout* layers. The *fully connected* layer is a $1 \times 3456$ vector which is followed by 64 and 32 *dense* layers which produce the output of the last layer, which uses *sigmoid* activation for binary classification. All other hidden layers use *ReLU* as an activator.

Bearing in mind that hearing aids are memory limited, we have also implemented CNN for just 10 time-frames with PCA component reduction, where we achieved the accuracy of 95%. In this approach, we are waiting for ten frames to be accumulated in the buffer and then the recognition will be performed. Thanks to that, we will be able to implement near real time recognition (5FPS) instead of taking 79 frames at once (0.5 FPS). At the same time, thanks to PCA, the input to CNN would be almost the size of one time-frame (37 features instead of 10x30). Thus, hearing aid device would not need to store additional data, therefore the memory constraint is satisfied.
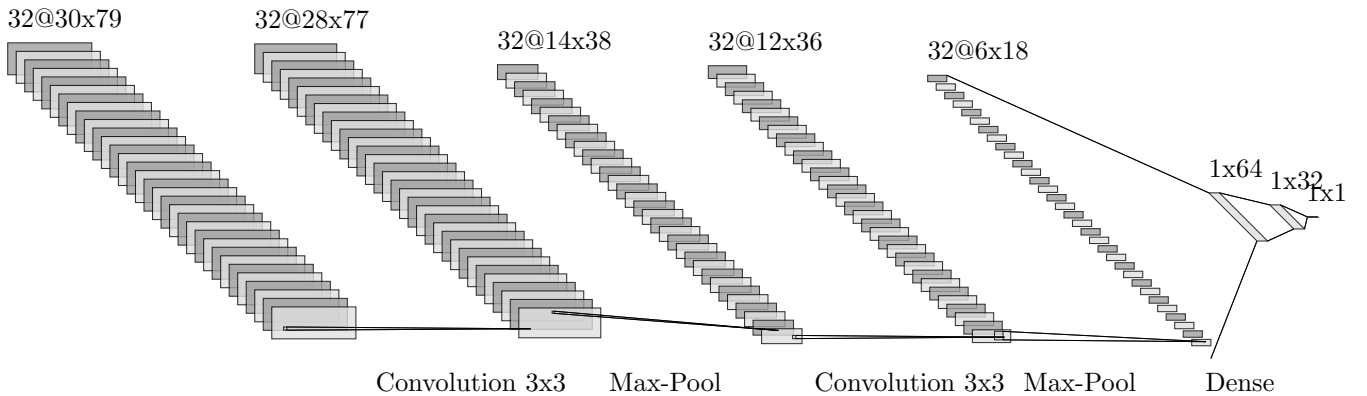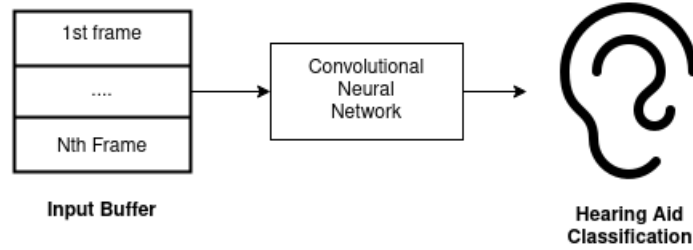


Figure 3: CNN architecture

Figure 4: Final Model Detection

After fine tuning our models, we achieved the following results:

| Classification Model | Accuracy | Precision | F1 score | Recall |
|---|---|---|---|---|
| CNN sample-by-sample | 0.983 | 0.997 | 0.983 | 0.969 |
| CNN 10 frames with PCA | 0.948 | 0.981 | 0.947 | 0.915 |

Last but not least, in order to make the model more lightweight we decided to test it using integers insted of floats. With this approach we achieved the accuracy of 94%. Thanks to that conversion and assuming that the production model will be implemented in C we will be able to save 6 bytes of memory per variable.

# 4 Project Feasibility

The main trade-off that we tried to solve is efficiency vs. memory and power. To increase the accuracy, we took the whole sample at once. As the solution has to be more memory efficient, we tweaked the data with PCA that 10 frames consuming similar amount of memory as one original frame. Moreover, the buffer allows to trigger the classifier less frequently, thus saving the power 10 times compared to frame-by-frame classifier.

This solution may sacrifice a little bit of user experience, as the processing would be computed at 5FPS. This is where the power constraint comes into play - analyzing whole sample would save more energy, but sacrificing the UX.

It was also possible to compute PCA component reduction frame-by-frame, but then computations are 10x more often than doing it once after 10 frames. Both options are equally good and possible with our application, it just depends on the priority. We decided that we can sacrifice a little bit more memory for the sake of saving the CPU power and battery.

On top of taking into consideration model sizes 1 and input for the models we described throughout the 3, we measured the computational time of our classification model. We carried out the tests on Intel(R) Core(TM) i5-10210U CPU @ 1.60GHz. Computing PCA for each sample took 0.28ms in total. CNN classifier takes 0.16ms. The computational times of CNN classifier are very similar when taking the whole sample as well as 10 frames as an input (0.15 ms) We are aware that this will take significantly more time on small hearing aid device, but also the code could be greatly optimized - coding with C instead of python and make operations that are prepared for the device's hardware.

Moreover, we can clearly see, that although PCA may result in saving the memory, it will consume almost 2x more power than the classification itself. Maybe extracting different features described in 5.2 is a path to explore.

# 5 Conclusions and Future Work

## 5.1 Conclusions

Two CNN models were developed, but as the one that takes all 79 frames at once perfoms better (98% accuracy) versus the one with PCA and 10 frames (95%). We are aware that the first model might not be memory-efficient, but as the data was given in this form, we decided to compute the final test results with this neural network. However, with PCA solution we developed, the accuracy is only slightly worse. In a long run, when only shorter samples are available and with the buffer we implemented or a sliding window, this would be the optimal solution.

After performing PCA on one frame, we discovered that first 6 principal components gives 95.2% of variance. We repeated the experiment on 10 frames and results are similar as seen on Figure 5. The fact that only 12% (or 20% in one frame case) of input data holds over 95% of information make us wonder that there might be more compact audio metrics that we could utilize. This is described in Section 5.2
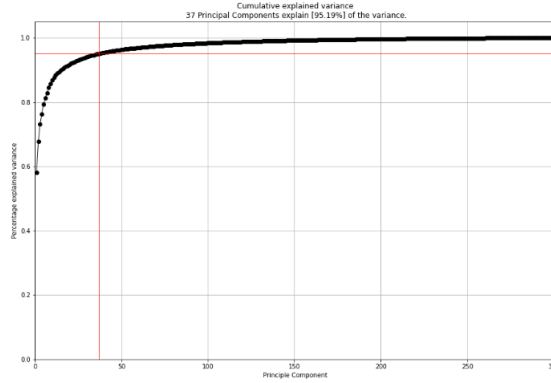


Figure 5: Principal component analysis of 10 frames[6]

## 5.2 Future Work

**Transformer Neural Networks** [2] could be a way to obtain even better results than CNN. They are better suited for sequential data such as text or music, as it makes the connections with the previous inputs. This is an improved version of LSTM [4] network. If we were given more time, this would be the main focus of our work to get even better results.

Besides given log-mel spectograms, we discovered **other metrics** that could be derived from audio signal or based on the spectogram itself. Some features are useful to process speech and singing[6]. Spectral envelope and detail, formants, MFCC [1]) are one of many [3] that could be explored more in depth. During the project we explored some with no improved results, but if raw audio files and more time, we believe this approach may bring valuable results.

On top of our CNN classifier solution, we could also optimize its performance on the whole samples by introducing **PID controller** on top of the buffer/sliding window. If 'music' was detected for a longer period, it would take more outputs classified as 'other' to fully 'convince' the model that music is no longer played. It could be useful if music had some parts that are only with lyrics, without the instruments. We believe that it would greatly improve the user experience. As it was not the essence of this project to implement it, we focus on other aspects of the problem instead.
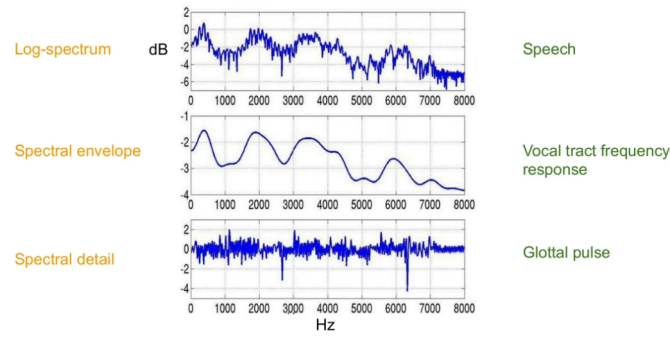
Figure 6: Speech features[6]

# References

[1] Amit Meghanani, Anoop C. S., A.G.R.: An exploration of log-mel spectrogram and mfcc features for alzheimer's dementia recognition from spontaneous speech `https://ieeexplore.ieee.org/document/9383491`

[2] Ankit, U.: Transformer neural network `https://towardsdatascience.com/transformer-neural-network-step-by-step-breakdown-of-the-beast-b3e096dc857f`

[3] Feinberg, D.R.: Voicelab `https://voice-lab.github.io/VoiceLab/`

[4] Olah, C.: Understanding lstm networks `http://colah.github.io/posts/2015-08-Understanding-LSTMs/?fbclid=IwAR0St715-AtEjh6eVlAP7DeMqkCtvHCyUd32-1iynsndQZQp6pdHN5tw39g`

[5] Ortner, A.: Top 10 binary classification algorithm `https://medium.com/thinkport/top-10-binary-classification-algorithms-a-beginners-guide-feeacbd7a3e2`

[6] Velardo, V.: Mel-frequency cepstral coefficients `https://www.youtube.com/watch?t=1294&v=4_SH2nfbQZ8&feature=youtu.be`

[7] Venkatesh, S., Moffat, D., Miranda, E.R.: You only hear once: A yolo-like algorithm for audio segmentation and sound event detection. Applied Sciences **12**(7) (2022), `https://www.mdpi.com/2076-3417/12/7/3293`