# DISCLAIMER

My prediction data in my google drive somehow disappeared so I had to run all the predictions again tonight (28.feb).

So I ran my analysis on not all the data and I will be updating the numbers and discussion of the results if the results change a lot.

Numbers in this paper are generated with;

Musicnn: 434/574 chunks predicted

VGGish: 265/574 chunks predicted

I will submit it not fully finished because I want to submit on time and like I said I will be updating it, when prediction data is finished

# Audio And Music Processing Lab
## Module 1 - Large audio datasets assignment

For this task I will annotate music data by hand. I was given 574 audio chunks to annotate. I have eleven characteristics for annotating. The characteristics are;

- If the song is electric
- If the song is acoustic
- If the song is aggressive
- If the song is relaxed
- If the song is happy
- If the song is sad
- If its a party song
- If the song is tonal or atonal
- Danceability
- If the song is instrumental
- Finally if the song is not instrumental then what gender is the singing voice

One of the main purposes for annotating data by hand is to transform that data into a form, that is suitable for computer-aided analysis. I will take my annotated data and compare it to two different pre-existing models, musicnn and vggish, that predict these characteristics automatically and analyse the accuracy of those models.

My process was pretty straight forward. I used the programming environment Colab because it is the easiest way to use Essentia (I have not figured out how to install Essentia on my macBook computer). I installed all the prediction models using !wget. I iterated over all the sound chunks and for each sound chunk I used each model to predict all the characteristics and saved the results in a json file so it was formatted the same way as the data which was annotated by hand. The predictions from the models gave me an array of probabilities for each characteristic, so I find what index gives me the best probability and that is what the model classified. I then calculate the accuracy for each model given that my annotation are ground truth and calculate confusion matrices for the results. I also noticed that some data was flipped and calculated the accuracies with that in mind.

Accuracy table for musicnn model:

| | mood_acoustic | mood_electronic | mood_aggressive | mood_relaxed | mood_happy | mood_sad | mood_party | tonal_atonal | danceability | voice_instrumental | gender |
|---|---|---|---|---|---|---|---|---|---|---|---|
| **Accuracy** | 0.757505773672 0554 | 0.730414746543 7788 | 0.760368663594 4701 | 0.518433179723 5023 | 0.684331797235 0223 | 0.428571428571 42855 | 0.836405529953 91 7 | 0.578341013824 8848 | 0.688940092165 8986 | 0.797235023041 4746 | 0.983870967741 9355 |

Accuracy table for vggish model

| | mood_acoustic | mood_electronic | mood_aggressive | mood_relaxed | mood_happy | mood_sad | mood_party | tonal_atonal | danceability | voice_instrumental | gender |
|---|---|---|---|---|---|---|---|---|---|---|---|
| **Accuracy** | 0.829545454545 4546 | 0.840909090909 0909 | 0.776515151515 1515 | 0.621212121212 1212 | 0.643939393939 3939 | 0.636363636363 6364 | 0.803030303030 303 | 0.337121212121 21 21 | 0.727272727272 7273 | 0.840909090909 0909 | 0.954545454545 4546 |

Confusion matrices for vggish model

| | | Actual Values | |
|---|---|---|---|
| | mood_acoustic | Positive | Negative |
| **Predicted Values** | Positive | 66 | 38 |
| | Negative | 7 | 153 |

| | | Actual Values | |
|---|---|---|---|
| | mood_electronic | Positive | Negative |
| **Predicted Values** | Positive | 133 | 28 |
| | Negative | 14 | 89 |

| | | Actual Values | |
|---|---|---|---|
| | mood_aggressive | Positive | Negative |
| **Predicted Values** | Positive | 29 | 46 |
| | Negative | 13 | 176 |

| | mood_relaxed | Actual Values | |
|---|---|---|---|
| | | Positive | Negative |
| **Predicted Values** | Positive | 97 | 99 |
| | Negative | 1 | 67 |

| | mood_happy | Actual Values | |
|---|---|---|---|
| | | Positive | Negative |
| **Predicted Values** | Positive | 31 | 57 |
| | Negative | 37 | 139 |

| | mood_sad | Actual Values | |
|---|---|---|---|
| | | Positive | Negative |
| **Predicted Values** | Positive | 136 | 92 |
| | Negative | 4 | 32 |

| | mood_party | Actual Values | |
|---|---|---|---|
| | | Positive | Negative |
| **Predicted Values** | Positive | 172 | 40 |
| | Negative | 12 | 40 |

| | tonal_atonal | Actual Values | |
|---|---|---|---|
| | | Positive | Negative |
| **Predicted Values** | Positive | 86 | 90 |
| | Negative | 85 | 3 |

| | danceability | Actual Values | |
|---|---|---|---|
| | | Positive | Negative |
| **Predicted Values** | Positive | 75 | 10 |
| | Negative | 62 | 117 |

| | voice_instrumental | Actual Values | |
|---|---|---|---|
| | | Positive | Negative |
| **Predicted Values** | Positive | 165 | 33 |
| | Negative | 9 | 57 |

|                     |          | Actual Values |          |
| ------------------- | -------- | ------------- | -------- |
|                     | gender   | Positive      | Negative |
| **Predicted Values** | Positive |           144 |       66 |
|                     | Negative |             7 |       47 |

The vggish model gave me better results. It has better accuracy in 7 out of 11 characteristics. One interesting result is tonality in both models, it looks like I misrepresented what tonality is. You can see when notating objective characteristics like acoustic, electronic and voice/instrumentals both models give a decent score there and for me that is expected. For the more subjective characteristics the accuracy score fluctuates a little bit. I think its a bit tough to get an accurate model for predicting such subjective characteristics, the feeling of sadness and happiness and just mood in general can be pretty different individually I think.