Milan, 10-15 January 2021

# Mood detection analyzing lyrics and audio signal based on deep learning architectures

K. Pyrovolakis, P. Tzouveli, G. Stamou School of Electrical and Computer Engineering National Technical University of Athens Athens, Greece

## The Emotional Model

- Rusel's Circumplex is the emotional model we used in our work
- According to Circumplex all emotions human are distributed in a 2D space with axes of valence and arousal
- Each quadrum represents a mood class

## From Lyrics to Mood

- Each world in lyrics is attributed to pair of valence and arousal values
- The set of values is computed with the help of dictionaries which contain emotional information
- A general pair of valence and arousal values is computed for each song

### Data

Results

- The dataset we used is the MoodyLyrics Dataset
- 2.000 song titles with their corresponding mood label
- Mood labels: {happy, angry, sad, relaxed}
- Audio data (37.989 samples, from web and augmentation)
- Lyrics data (18.115 samples, from web and augmentation)
- Compute BERT Embeddings



Valence (V) and arousal (A) values	Mood
$A > A_t$ and $V > V_t$	Нарру
$A > A_t$ and $V < -V_t$	Angry
$A < -A_t$ and $V < -V_t$	Sad
$A < -A_t$ and $V > V_t$	Relaxed

# From Audio to Mood

Association between structural features of music and emotion

Structural Feature	Definition	Associated Emotion	
Tempo	The speed or pace of a musical piece	Fast tempo:happiness, excitement anger. Slow tempo: sadness, seren- ity.	
Mode	The type of scale	Major tonality: happiness, joy. Mi- nor tonality: sadness.	
Loudness	The physical strength and amplitude of a sound	Intensity, power, or anger	
Melody	The linear succession of musical tones that the listener perceives as a single entity	Complementing harmonies: happi- ness, relaxation, serenity. Clashing harmonies: excitement, anger, un- pleasantness.	
Rhythm	The regularly recur- ring pattern or beat of a song	Smooth/consistent rhythm: happiness, peace. Rough/irregular rhythm: amusement, uneasiness. Varied rhythm: joy.	

Features extracted from audio that we experimented with:

- Spectogram
- Mel Spectogram
- Log-Mel Spectogram
- MFCCs
- Chroma features
- Centroid tonal features
- Spectral contrast

# System Architecture

Prediction

Embedded Audio Features Lyrics ENCODER 2D Convolution ENCODER 2D Max Pooling ENCODER 2D Convolution ENCODER 2D Max Pooling Model  $A_1$ ENCODER 2D Col ENCODER 2D Max Pooling ENCODER Flatte ENCODER Fully-Connected ENCODER Fully-Connected ENCODER Class Prediction ENCODER Audio Embedded ENCODER Lyrics Features ully-Cor Model M Fully-Connected Class Prediction Fully-Conne Class

How the multichannel system (M₁) is developed?

- Train BERT-base uncased model  $(T_2)$  on lyrics
- Train CNN model (A<sub>1</sub>) on audio signal
- System M<sub>1</sub> is implemented as the fusion of  $A_1$  and  $T_2$  with a common classifier of two fully connected layers







AILS

### Lyric Analysis Subsystem

We trained BERT model  $(T_2)$  and compared its results with several text analysis techniques

Model	Embedding Method	Loss	Accuracy %
$T_1'$	BoW	1.287	65.49
$T_1'$	TF-IDF	1.381	67.98
$T_1$	Word2Vec	1.262	41.66
$T_1$	GloVe	1.064	53.33
$T_2$	Bert	1.353	69.11



#### Audio Analysis Subsystem

 $\exists_2$ 

Model

We trained CNN model ( $A_1$ ) and experimented with different possible feature combinations

Feature Combination	Accuracy %
Mel	64.97
Mel, Log-Mel	68.38
Mel, Chroma, Tonnetz, Spectral Contrast	60.86
Log-Mel, Chroma, Tonnetz, Spectral Contrast	58.96
MFCC, Chroma, Tonnetz, Spectral Contrast	65.36
Mel, Log-Mel, MFCC, Chroma, Tonnetz	69.77
Mel, Log-Mel, MFCC, Chroma, Tonnetz, Spec- tral Contrast	70.34

### Fuse Analysis System

We used the already trained subsytems to train our multichannel model (M<sub>4</sub>) and compared its results with the previous models

Model	Loss	Accuracy %	Computational Time
$T_1'$	1.381	67.98	0m 25.391s
$T_2$	1.353	69.11	18m 12.444s
$A_1$	0.743	70.51	80m 13.064s
$M_1$	0.156	94.58	3m 38.551s

