

Second meeting

24/10/2024

PLAYLIST

- title
- Title song 1
- Title song 2
- ...

CORPUS



- Should we use Json file or CSV file ?

Dataset

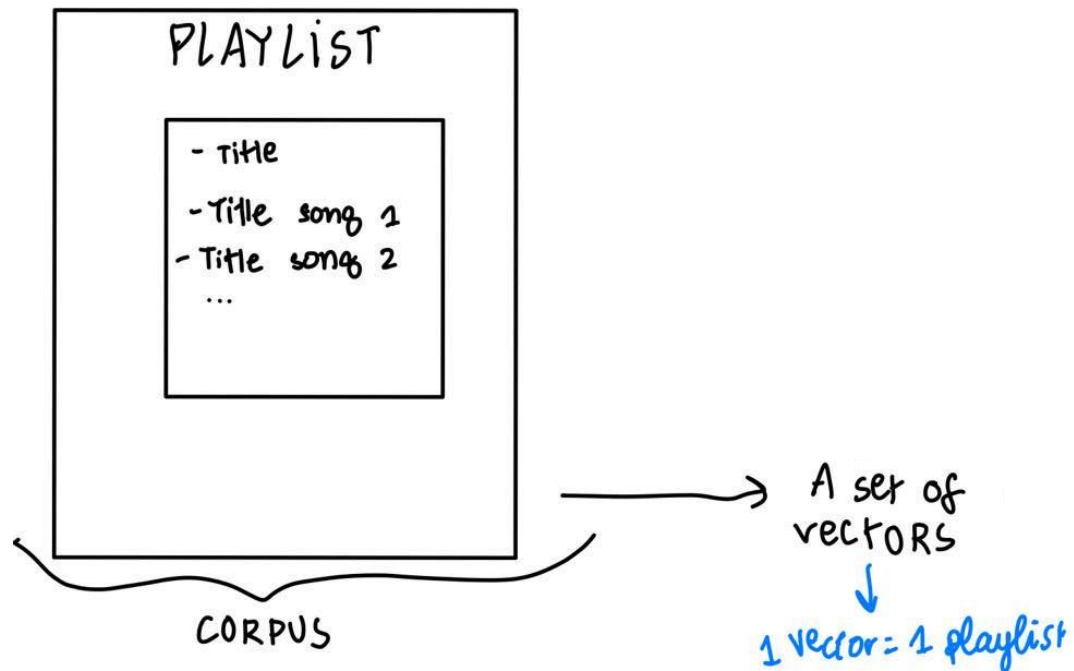
We converted the original JSON files in an equivalent CSV version.

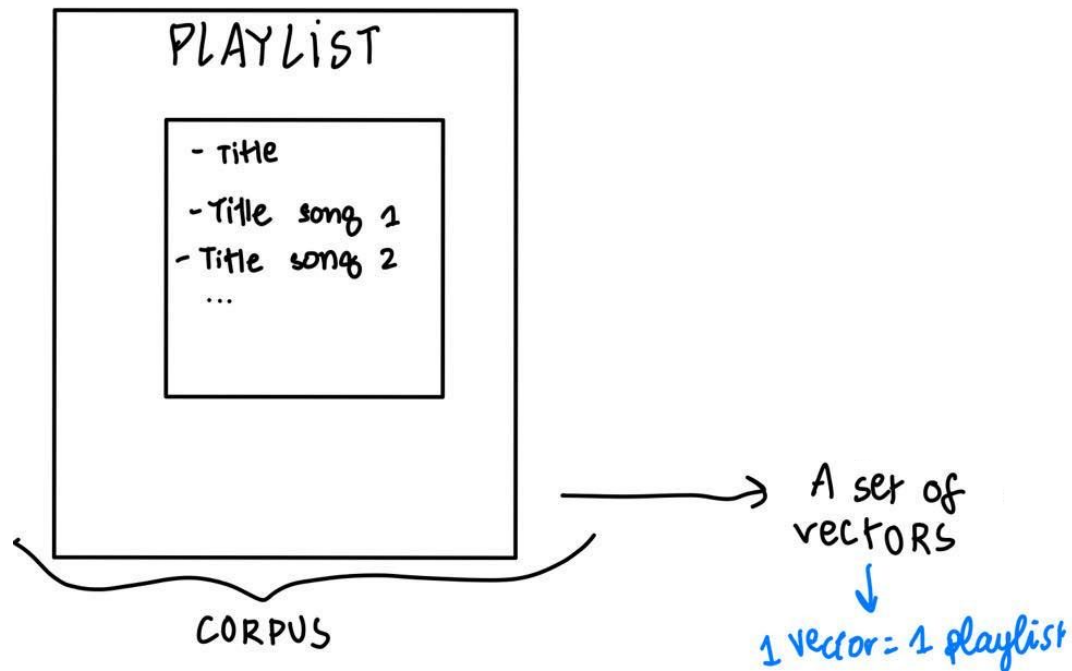
```
python evaluation/mpd2csv.py --mpd_path /path/to/mpd --out_path dataset  
python evaluation/challenge2csv.py --challenge_path /path/to/challenge.json --out_path dataset
```

- Cleaning data

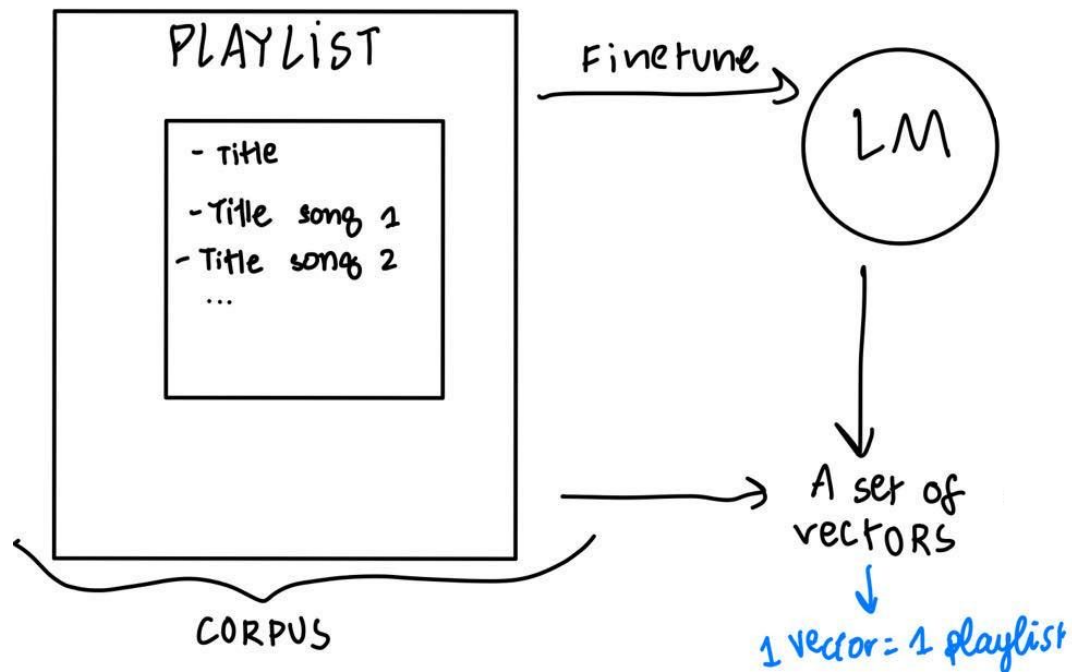
- lowercasing;
- detecting and separating emoji from words;
- separating the skin code from the emoji;
- detecting and separating emoticons from words;
- transforming space-separated single letters into words (e.g. "w o r k o u t" becomes "workout");
- remove '#' from hashtags.

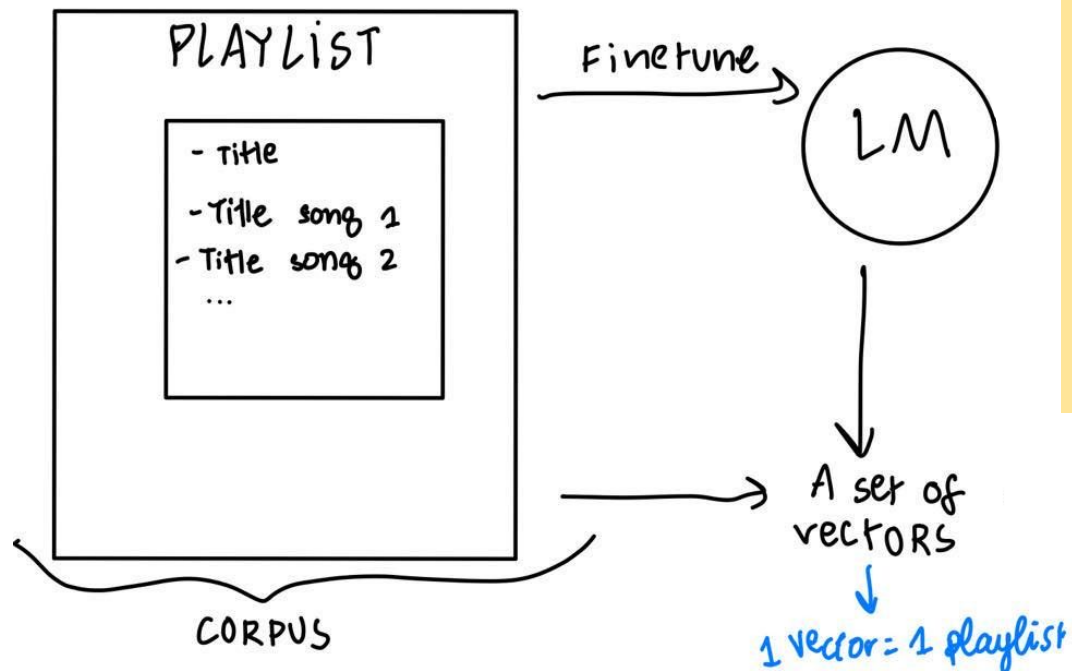
using "BERT un_cased" ???





- Embedding using fastText?
- Word2Vec vs Title2Vec
- Title2Vec: K-means algorithm





- Bert: bidirectional context
- Pre-training: mask words and sentences
- Easy to fine-tune (few additional layers needed)

