# Million Songs

Team *Win or Lose* - Kezhi, Yinchen, Shaoze, Jiache

# 1. Data Preparation

The following codes are added to `~/.bashrc` to automatically mount millionsongs upon booting.

```
1 echo "password" | sudo -S sshfs /home/hadoopuser/ece472
  ↪ -o allow_other -o Port=2223
  ↪ ece472@focs.ji.sjtu.edu.cn: -o
  ↪ IdentityFile=~/.ssh/id_ed25519 1>/dev/null
  ↪ 2>/dev/null
2 echo "password" | sudo -S mount
  ↪ /home/hadoopuser/ece472/millionsong.iso
  ↪ /home/hadoopuser/ece472/ 1>/dev/null 2>/dev/null
```

```
Track.h5
├── analysis
│   ├── bars confidence
│   └── ...
├── metadata
│   ├── artist terms
│   └── ...
└── musicbrainz
    ├── artist mbtags
    └── ...
```

- We used **python h5py** module to extract information in *.h5* files.

- We collect useful fields in all *.h5* files and summarize them in one single *.avro* file

- Divide the dataset into 26 parts

- Create one *.avro* file for all files in each part

- Use **pyspark** to parallelize the process

- Merge 26 small *.avro* files into one single big file

With 8 cores processing in parallel, the total time cost to form one *.avro* file is reduced from 6 hours to 3 hours.

The *.avro* file generated which is consisted of all desired features of one million songs is approximately 150Mb.

Table: Statistic of Raw Data

|        | tempo    | hotness | year    | time_signature | ... |
|--------|----------|---------|---------|----------------|-----|
| count  | 1000000  | 581965  | 1000000 | 1000000        |     |
| mean   | 123.889  | 0.356   | 1030    | 3.59           |     |
| std    | 35.056   | 0.234   | 999     | 1.22           |     |
| min    | 0.000    | 0.000   | 0       | 0              |     |
| 25%    | 97.995   | 0.215   | 0       | 3              |     |
| 50%    | 122.086  | 0.378   | 1969    | 4              |     |
| 75%    | 144.089  | 0.532   | 2002    | 4              |     |
| max    | 302.300  | 1.000   | 2011    | 7              |     |

What features are needed to describe a Song?

1. Duration
   - The duration cannot directly determine the style of songs
   - Adjusting granularity & OneHotEncoder

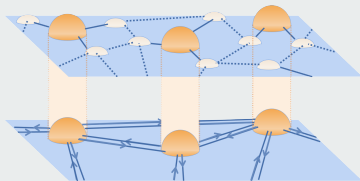2. Segments_loudness_max / Segments_loudness_max_time
   - The most load segment can be regarded as the **"musical climax"**
   - Classify this feature by the occurrence of musical climax.

3. Segments_pitches
   - Analyze **emotion** of song by computing the mean of pitches.
   - Associate with the Segments_loudness_max

**Prepare the data for Model:**

- Missing & Abnormal Value

- Distribution (Left/Right Skew)

- Granularity (Year)

- Encoding (One-hot, Hash...)

- Combination

- Normalize & Standardize

**Issues:**

- Mostly $< 0$ with 359/1000000 abnormal values

- Left skewed



Figure: Distribution of Raw Data

**Methods:**

- Replace exceptions with median

- $loudness\_log = ln(-loudness)$
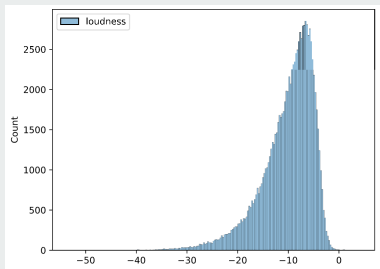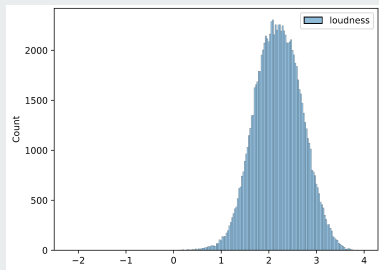


Figure: Distribution of Adjusted Data

```
1  processed_data = process_data_gm(data, (
2      # Customized Column Transformer
3      (log_transform_negative_column, ['loudness'], None),
4      (classify_year_column, ['year'], None),
5      (proportion_fade_out, None, None),
6      (fade_out_time, None, None),
7
8      # Exception Handling
9      (drop_zeros, [List_of_Features], None),
10     (fill_hotness_na_with_0, None,  None),
11
12     # Normalize Selected Features
13     (normalized_selected_columns, [List_of_Features], None),
14
15     # Select Features
16     (select_columns, ['Log_loudness', 'bar_num', 'beat_num',
17                        # Add features here
18                        ], None),
19 ))
```

# 2. Drill Database Query

In this part, we used drill to perform simple database queries, including:

- Find the range of dates covered by the songs in the dataset, i.e. the age of the *oldest* and of the *youngest* songs

- Find the *hottest* song that is the *shortest* and has the *highest energy* with the *lowest tempo*.

- Find the name of the album with the *most tracks*.

- Find the name of the band(artists) who recorded the *longest song*.

Given the avro file, we first created a table in drill

```
1  -- read avro file from local file system (~100M)
2  create table dfs.tmp.`songs` as
3  select * from dfs.`F:/avro/songs.avro`;
4  -- change path
5  use dfs.tmp;
```

# Detailed Solutions

1. the oldest and youngest songs

```
1  select max(year_end) from songs where year_end > 0;
2  +--------+
3  | EXPR 0 |
4  +--------+
5  | 2011   |
6  +--------+
7  1 row selected (0.321 seconds)
8
9  select min(year_end) from songs where year_end > 0;
10 +--------+
11 | EXPR 0 |
12 +--------+
13 | 1922   |
14 +--------+
15 1 row selected (0.323 seconds)
```

Therefore, the dataset covered songs from 1922 to 2011, namely the age of the songs vary from 12 years to 101 years.

2. the hottest, shortest, highest energy, lowest tempo

```
1  select id, title from songs
2  where hotness <> 'NaN'
3  order by hotness desc, duration asc, energy desc, tempo asc
4  limit 5;
5
6  +------------------+-----------------------------------+
7  |        id        |              title                |
8  +------------------+-----------------------------------+
9  | SONASKH12A58A77831 | Jingle Bell Rock                |
10 | SOAVJBU12AAF3B370C | Rockin Around The Christmas Tree |
11 | SOEWAKD12AB01860D5 | Holiday                         |
12 | SOAAXAK12A8C13C030 | Immigrant Song (Album Version)  |
13 | SOAXLDX12AC468DE36 | La Tablada                      |
14 +------------------+-----------------------------------+
15 5 rows selected (0.497 seconds)
```

Therefore, Jingle Bell Rock is the song the hottest song that is the
shortest and shows highest energy with lowest tempo.

3. the album with the most songs

```
1  select album_id, album_name, count(album_id) as numSongs
2  from songs
3  group by album_id, album_name
4  order by numSongs desc
5  limit 1;
6
7  +----------+------------------------------------+----------+
8  | album_id |            album_name              | numSongs |
9  +----------+------------------------------------+----------+
10 | 60509    | First Time In A Long Time:         | 85       |
11 |          | The Reprise Recordings             |          |
12 +----------+------------------------------------+----------+
13 1 row selected (1.076 seconds)
```

First Time In A Long Time: The Reprise Recordings is the album with
most tracks. Indeed, it has 4CDs and 80+ tracks.

4. the band with longest song

```
1  select artist_name, title, duration from songs
2  order by duration desc limit 1;
3
4  +-----------------------------+-----------+----------+
5  |         artist_name         |   title   | duration |
6  +-----------------------------+-----------+----------+
7  | Mystic Revelation of Rastafari | Grounation | 3034.9058 |
8  +-----------------------------+-----------+----------+
9  1 row selected (0.468 seconds)
```

Therefore, the band `Mystic Revelation of Rastafari` has
recorded `Grounation` which has highest duration.

# 3. Big Data Recommendation

Similarity between $Song_A = [a_1, a_2, ..., a_n]$ and $Song_B = [b_1, b_2, ..., b_n]$

**①** $L_1$ Norm

$$d_{L_1} = \sum_{i=1}^{n} |a_i - b_i|$$

**②** Cosine Similarity

$$cos\theta = \frac{\sum_{i=1}^{n}(a_i \times b_i)}{\sqrt{\sum_{i=1}^{n} a_i^2} \times \sqrt{\sum_{i=1}^{n} b_i^2}}$$

**③** Combination

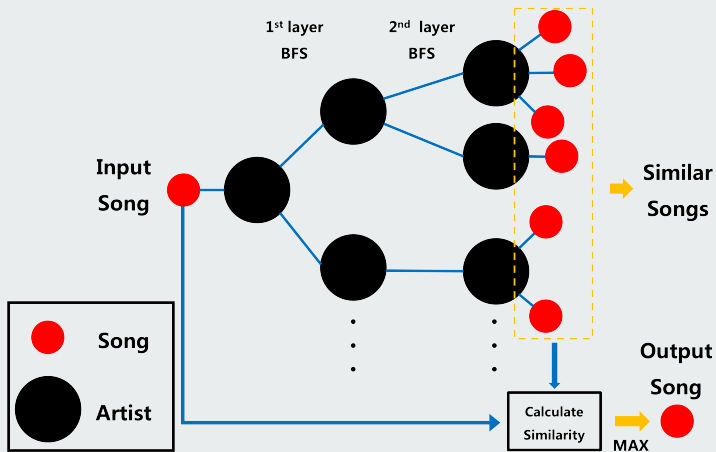$$Similarity(Song_A, Song_B) = \lambda cos\theta - d_{L2}$$
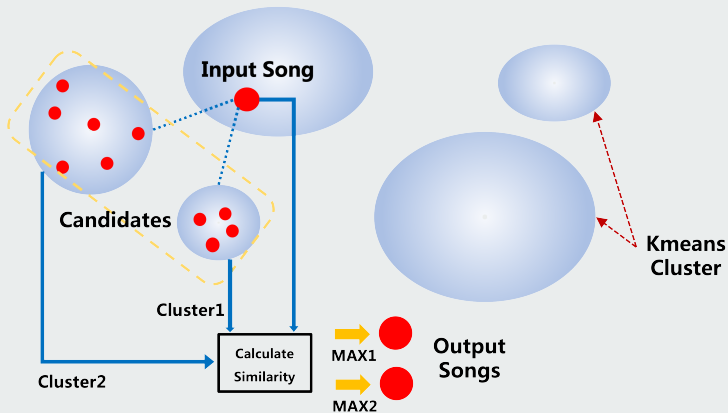
Figure: Two layers BFS

Input song: (Old Man Mose, The Bristols)
Recommended Song:

1. $L_1$ Norm: (The Story of Two, Micragirls)

2. Cosine Similarity: (Kentish (demo), Modwheelmood)

3. Combination:
   When $\lambda \leq 307$, (Kentish (demo), Modwheelmood).
   When $\lambda \geq 308$, (The Story of Two, Micragirls).

# Diverse Recommendation Implementation

*Choose different $\lambda$ to reach diverse recommendations!*

```python
import numpy as np
from numpy import ndarray
def calcDistance(song1: tuple, song2: ndarray) -> float:
    # song1: (feature: ndarray, track_Id: str)
    weight = 325
    return weight * (np.dot(song1[0], song2) / \
            (np.linalg.norm(song1[0]) * \
            np.linalg.norm(song2)) \
            - np.sum(np.abs(song1[0]-song2)), song1[1])
```

```
1     def getKMostRelatedCenter(track_id, data, cluster_centers,
      ↪  k=3)-> list:
2
3     song = data[data["track_id"] == track_id]
4     song = song.loc[:, features_list].to_numpy()
5     dis = []
6     for center in cluster_centers:
7         dis.append(np.dot(song, center) / (np.linalg.norm(song) *
          ↪  np.linalg.norm(center)))
8
9     sorted_np = np.argsort(np.array(dis), axis=0)[:,0] < k
10    index = range(0, np.shape(cluster_centers)[0])
11
12    cluster_centers_withindex = np.hstack((cluster_centers,
      ↪  np.array(index).reshape(-1, 1)))
13    selected_centers = cluster_centers_withindex[sorted_np][:, -1]
14    return list(selected_centers)
15
```

Function of mappers and reducer:

- mapper_artist.py: given an artist id as input, find all the similar artists (according to the database)

- mapper_song.py: given the similar artists, find all of their songs

- mapper_distance.py: given the list of songs, calculate their cosine similarity from the given song

- reducer.py: find the song with largest similarity

We implement a `driver.sh` to run the three map reduce job:

```
1  # How to run
2  cd MapReduce
3  time bash ./driver.sh
4  # Result
5  + hdfs dfs -cat /project_distance/part-00000
6  ('Story Of Two', 'TRMHEEB12903C9F3C9',
   ↪  0.9140447260983122)
7
8  real    4m18.674s
9  user    1m27.983s
10 sys     0m5.830s
```

Use the same strategy, we implement the map and reduce in spark:

```python
1  sc = SparkContext()
2  # find the list of similar artists
3  for i in range(depth):
4      artists += sc.parallelize(artists, 4)\
5                  .map(artistNeighbor).reduce(merge_lists)
6  # find all the songs of similar artists
7  songs: list = sc.parallelize(artists, 12)\
8                  .map(getArtistSongs).reduce(merge_lists)
9  # find the feature of the input song
10 features = sc.parallelize(features, 100)\
11                .map(lambda x: (np.concatenate((x[1:2], x[3:-7]))\
12                .astype(np.float64), (x[-6],x[-5],x[-1]))).collect()
13 # reduce to get the song with largest similarity
14 result = sc.parallelize(features, 100)\
15             .map(lambda x: calcDistance(x, songFeature))\
16             .reduce(lambda x, y: max(x, y))
```

For pyspark, there is an interface for you to play with :).

# Pyspark Performance

```
1   Please enter the name of a song:  Old Man Mose
2   Too many songs have the same name! Please choose a specific author
    ↪   from the list:
3   ['Jesse Fuller', 'George Lewis And His New Orleans Stompers', 'Louis
    ↪   Armstrong', 'Manhattan Transfer', 'Kenny Ball And His Jazzmen',
    ↪   'The Bristols']
4   The author of your song: The Bristols
5   The song you choose:
6   Name: Old Man Mose, Author: The Bristols, Id: TRYESJS12903CDF730
7   Please enter the depth of the BFS: 2
8   Num of similar artists in the 1th layer: 48.
9   Num of similar artists in the 2th layer: 1134.
10  Num of similar songs: 29818.
11  (0.9140447260983123, ("Story Of Two", "TRMHEEB12903C9F3C9", "The
    ↪   Micragirls"))
12
13  real    0m32.016s
14  user    0m2.668s
15  sys     0m2.587s
```

That is **8** times speed up than Mapreduce!

# 4. Year Prediction

We want to train a model to predict the year of a song with its features.

```
1  SELECT COUNT(*) AS year_num FROM songs WHERE year<>0;
2  +----------+
3  | year_num |
4  +----------+
5  | 515576   |
6  +----------+
7  1 row selected (0.402 seconds)
```

Approximately half of the songs have years labeled. We use 80% of them as training set and 20% as validation set.

- Years are grouped by every 5 neighboring years so that total number of classes is reduced from 88 to 18.

- Each feature is normalized to the range of $[0, 1]$

- Very few data are missing (less than 1%). We just fill them as 0 and it would be not a big deal to the model.

# PCA

We fitted a PCA model on the training set, which keeps 6 principal components from 15 original features.

```
1  pca = PCA(k=6, inputCol='features',
   ↪ outputCol='pca_features')
2  pca_model = pca.fit(feature_df)
3  return pca_model
4  pca_df = pca_model.transform(feature_df)
5  res_df =pca_df.select('coarse_classified_class_year',
6  'pca_features').rdd.map(lambda x: Row(tag=x[0],
7  PC1=float(x[1][0]), PC2=float(x[1][1]),
8  PC3=float(x[1][2]), PC4=float(x[1][3]),
9  PC5=float(x[1][4]), PC6=float(x[1][5])))
10 .toDF()
```

We trained a logistic regression model to do the classification task.

```
lrm = LogisticRegressionWithLBFGS.train(
sc.parallelize(data_train, 200), iterations=200,
numClasses=18)
print("Weights:", lrm.weights)
lrm.save(sc,'file:///home/hadoopuser/project
              /predict/model')
```

```
1  (dfy['coarse_classified_class_year']==
2   dfp['year_predict']).value_counts()
```
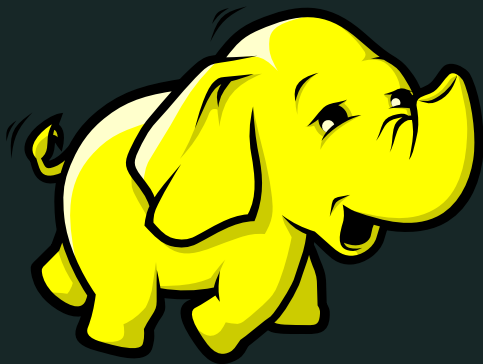
**No PCA**                           **Applying PCA**

```
1  False    79282
2  True     23275
3  Name: count,
4  dtype: int64
```

```
1  False    71696
2  True     30861
3  Name: count,
4  dtype: int64
```

Applying PCA raises prediction accuracy from 22.69% to 30.09%. For a 18 classification task, this result is satisfying enough.

**Thank you!**