**Project Proposal**                                              **Shoaib Ahmed**
CS401                                                              Roll : 1501049
Date: 30/08/18

# Problem Statement

Prediction of the release year of a song from audio features, using the YearPredictionMSD[1], a subset of the Million Songs Dataset[2]
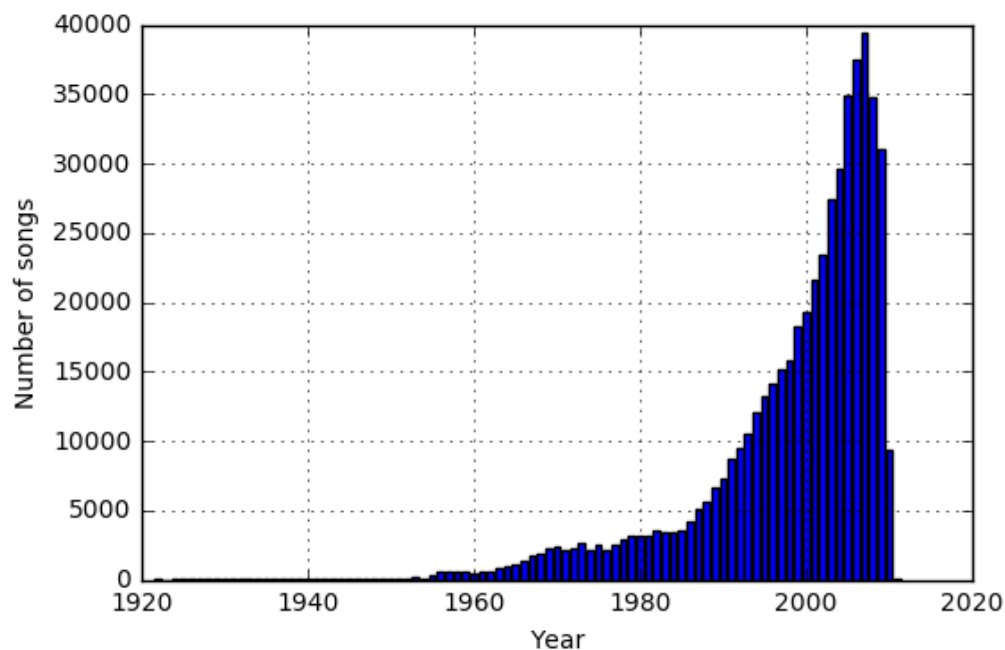
# About the Problem

Year prediction is defined as estimating the year in which a song was released based on its audio features. Million Song Dataset, a famous data set with freely available collection of audio features and metadata for a million contemporary popular music tracks. There are many attractive features in the Million Song Database. Here, what we focus are *Timbre Average* and *Timbre Covariance* features, and then do the release year prediction based on these.

# About the Dataset

The data set used for this project is UCI Machine Learning Respository [1] YearPredictionMSD which is a subset of the Million Song Dataset[2]. The songs here are mostly western, commercial tracks ranging from 1922 to 2011, with a peak in the year 2000s.

Figure 1: Trend of Data in the Dataset through the years.

# Data Collection

Specific characteristics of the dataset[1] used, are:

- Size: 460MB

- File Format: .txt

- Number of songs: 515345

- Attributes: 90 (12 = timbre average, 78 = timbre covariance)

The dataset contains features and labels for each song, where, label is the year in which the song was released. There are 12 timbre-average and 78 timbre-covariance features. Features have been extracted from the 'timbre'[4] features from The Echo Nest API[5]. We take the average and covariance over all 'segments', each segment being described by a 12-dimensional timbre vector. The API divides the song into multiple time segments. These segments were delimited by onsets and (or) other discontinuities in the signal

Figure 2: Screenshot of the data. Each line contains the Year followed by the 90 features.



# Construction/Implementation

I am planning to implement this research paper[3]. It is implements four types of machine learning algorithms on Apache Spark Machine Learning library (MLlib). The algorithms applied on the dataset are Linear Regression, Random Forest and Gradient Boosted Trees. The data preprocessing step includes data rescaling and normalization as the dataset at our disposal is highly skewed towards latest sound tracks, i.e., the number of songs released after 1990 is much greater than the number of songs released before 1990, so, we divide the dataset into two parts - Before 1990 and After 1990

# References

[1] *Source :* https://archive.ics.uci.edu/ml/datasets/yearpredictionmsd

[2] T. Bertin-Mahieux *Original Source :* http://labrosa.ee.columbia.edu/millionsong

[3] Mishra, P., Garg, R., Kumar, A.,2016 *2016 International Conference on Advances in Computing, Communications and Informatics (ICACCI).*
    https://ieeexplore.ieee.org/document/7732275/

[4] *Timbre :* https://en.wikipedia.org/wiki/Timbre

[5] *EchoNest :* http://the.echonest.com/