



Digital Receipt

This receipt acknowledges that Turnitin received your paper. Below you will find the receipt information regarding your submission.

The first page of your submissions is displayed below.

Submission author: Philip Steinke
Assignment title: Assignment 3:
Submission title: Assignment 3 preprocessing music...
File name: MATH2349_Assignment_3_AshBe...
File size: 0
Page count: 21
Word count: 5,655
Character count: 32,958
Submission date: 03-Jun-2018 11:41PM (UTC+1000)
Submission ID: 971740063

MATH2349 Semester 1, 2018

CODE

Assignment 3

Ben Cole s3412349, Ashleigh Olney s3686808, Phil Steinke s3725547, Ellen Toumpas s3708633

Links: GitHub: <https://github.com/turnitindatapreprocessing/assignment3> assignment notes: https://drive.google.com/open?id=1agf9_uD-6ZC-XH5X030C7P0F7P0a25205b0uH0g how to work on 8 MarkDown and R files

Required packages

Executive Summary

We found a dataset with 7 tables that gave us the potential to explore our abilities to merge datasets. However upon investigating the full scope of all the information and variables we had access to we realised playing and chasing with ALL of the data would make this project a lot bigger than it needed to be. Therefore before proceeding with the preprocessing a decision had to be made to select a subset of this data to play with.

Once the data was submitted problems were experienced even from the import, with the multi-index nature of the original data creating issues with column names and data types on import. Actions were required to fix this, correct data types and rename a few factors.

Understanding the data to get preliminary insight into data that may form impossible records and contain missing values was performed. Moving it to making our data able to fit the tidy data concepts, genre, album columns were separated from gathered in to a new column so that each genre of a song was recorded to a separate row. Latitude and longitude was also gathered to abide by the tidy data rules. Dealing with missing values, and outliers followed by steps to transform.

Data

The dataset was originally found on the UCI Machine Learning Repository and came from a project called FMA: A Dataset For Music Analysis Data Set [LINK: <https://archive.ics.uci.edu/ml/datasets/FMA3A4aDatasetForMusicAnalysis>]
(<https://archive.ics.uci.edu/ml/datasets/FMA3A4aDatasetForMusicAnalysis>)

The original data is split in to nine core: the tracks.csv, genres.csv, features.csv, echonest.csv along with the raw_albums.csv, raw_artists.csv, raw_echonest.csv, raw_genres.csv and raw_tracks.csv.

Out of all nine core we had decided to only concentrate on three of them. The echonest.csv, the tracks.csv and the genres.csv. Not all variables were ultimately dealt with, but the data descriptions of the three are below.

TRACKS.CSV is a dataset with 53 variables and 108,574 records. The columns we decided to work with were: track_id [integer] as id used for track - album type [character] a factor categorical data describing the album type - album_title [character] the title of the album - album_tracks [integer] referencing the track on the album - album_isbn [integer] the number of album isbns - artist_name [string] the artist releasing the album/ song - artist_latitude [float] the artist's latitude location - artist_longitude [float] the artist's longitude location - track_instrument [integer] track duration [integer] the length of the song in seconds - track_genre_top [string] the top genre tag describing the song - track_genre [integer] a list of all the other description tags allocated to genre - track_composer [string] the composer working on the track - track_genre_all [integer] all the genre descriptions

GENRES.CSV is a dataset with 5 variables and 163 records - genre_id [integer] as id used to match genres - #tracks [integer] number of tracks recorded against this genre - parent [integer] - title [character] the title of the genre - top_level

ECHONEST.CSV has 250 variables and 13,729 records - audio_features [acousticness, danceability, energy, instrumentalness, liveness, speechiness, tempo, valence columns] [double] all values measured by the echonest project to capture the features of a song - metadata_album_date [date] the album date the album was recorded - metadata_album_name [character] the album title the track belonged to - metadata_artist_latitude [float] the artist's latitude location - metadata_artist_longitude [float] the artist's longitude location - metadata_artist_name [character] the name of the artist releasing the track - metadata_release [character] the title of the release - ranks_artist_discovery_rank, artist_hotness_rank, song_commerce_rank, artist_familiarity, artist_hotness, song_commerce, song_hotness [integer] all values related to the popularity, social feeds, ranks of songs and the artist - temporal_features [000 - 250] temporal features that were recorded using the echonest project to capture the features of a song

```
echonest %>% head() %>% select("track_id", "album_title", "album_tracks", "artist_name", "artist_latitude", "artist_longitude", "track_instrument", "track_duration", "track_genre_top", "track_genre", "track_composer", "track_genre_all")
```

```
X1 audio_features audio_features_1 audio_features_2 audio_features_3 audio_features_4 audio_features_5 audio_features_6 audio_features_7 met
```

```
NA acousticness danceability energy instrumentalness liveness speechiness tempo valence
```

```
tracks %>% head() %>% select("track_id", "album_title", "album_tracks", "artist_name", "artist_latitude", "artist_longitude", "track_instrument", "track_duration", "track_genre_top", "track_genre", "track_composer", "track_genre_all")
```

```
X1 album album_1 album_2 album_3 album_4 album_5 album_6 album_7 album_8 album_9 album_10 album_11 album_12
```

```
NA comments date_created date_released engineer favorites id information isbns producer tags title tracks type
```