# PES University, Bangalore

## UE20CS312 - Data Analytics

## Worksheet 2b : Multiple Linear Regression

Course Anchor : Dr. Gowri Srinivasa

Prepared by : Nishanth M S - nishanthmsathish.23@gmail.com

## Multiple Linear Regression

Multiple Linear Regression (MLR) is a statistical technique that uses several explanatory variables to predict the outcome of response variable.The goal of MLR is to model **a linear relationship** between explanatory(independent) variables and response(dependent) variables.

## Data Dictionary

The data required for this worksheet can be downloaded from this GitHub Link. The data was obtained from this dataset from Kaggle. The dataset contains features of songs on Spotify collected using Spotify API.The features are as follows :

-**acousticness** : A confidence measure from 0.0 to 1.0 of whether the track is acoustic. 1.0 represents high confidence the track is acoustic.

-**danceability** : Danceability describes how suitable a track is for dancing based on a combination of musical elements including tempo, rhythm stability, beat strength, and overall regularity. A value of 0.0 is least danceable and 1.0 is most danceable.

-**duration_ms** : The duration of track in milliseconds.

-**energy** : Energy is a measure from 0.0 to 1.0 and represents a perceptual measure of intensity and activity.Perceptual features contributing to this attribute include dynamic range, perceived loudness, timbre, onset rate, and general entropy.

-**instrumentalness** : Predicts whether a track contains no vocals.The closer the instrumentalness value is to 1.0, the greater likelihood the track contains no vocal content. Values above 0.5 are intended to represent instrumental tracks, but confidence is higher as the value approaches 1.0.

-**key** : The key the track is in. Integers map to pitches using standard Pitch Class notation

-**liveness** : Detects the presence of an audience in the recording. Higher liveness values represent an increased probability that the track was performed live. A value above 0.8 provides strong likelihood that the track is live.

-**loudness** : The overall loudness of a track in decibels (dB). Loudness values are averaged across the entire track and are useful for comparing relative loudness of tracks. Loudness is the quality of a sound that is the primary psychological correlate of physical strength (amplitude). Values typical range between -60 and 0 db.

-**mode** : Mode indicates the modality (major or minor) of a track, the type of scale from which its melodic content is derived. Major is represented by 1 and minor is 0.

-**speechiness** : Speechiness detects the presence of spoken words in a track. The more exclusively speech-like the recording (e.g. talk show, audio book, poetry), the closer to 1.0 the attribute value. Values above 0.66 describe tracks that are probably made entirely of spoken words. Values between 0.33 and 0.66 describe tracks that may contain both music and speech, either in sections or layered, including such cases as rap music. Values below 0.33 most likely represent music and other non-speech-like tracks.

-**tempo** : The overall estimated tempo of a track in beats per minute (BPM). In musical terminology, tempo is the speed or pace of a given piece and derives directly from the average beat duration.

-**time_signature** : An estimated overall time signature of a track. The time signature (meter) is a notational convention to specify how many beats are in each bar (or measure).

-**valence** : A measure from 0.0 to 1.0 describing the musical positiveness conveyed by a track. Tracks with high valence sound more positive (e.g. happy, cheerful, euphoric), while tracks with low valence sound more negative (e.g. sad, depressed, angry).

Throughout the course of this worksheet , our response variable is energy. We shall try and apply the concepts learnt in class to predict the energy of a song using the other features of a song.

## Libraries used

-tidyverse

-corrplot

-olsrr : documentation

## Points

The problems for this worksheet is for a total of 10 points and the weightage is not uniformly distributed.

- *Problem 1* : 0.5 points
- *Problem 2* : 2 points
- *Problem 3* : 2 points
- *Problem 4* : 1 point
- *Problem 5* : 1.5 points
- *Problem 6* : 1 point
- *Problem 7* : 2 points

## Loading the Dataset

After downloading the dataset and ensuring the working directory is right , we read the csv into the dataframe.

```
library(tidyverse)
spotify_df <- read_csv('spotify.csv')
```

## Problem-1 (0.5 Points)

Check for missing values in the dataset and normalize the dataset.

## Problem-2 (2 Points)

Fit a linear model to predict the *energy* rating using *all* other attributes.Get the summary of the model and explain the results in detail.[*Hint* : Use the lm() function. Click here To get the documentation of the same.]

## Problem-3 (2 points)

With the help of a correlogram and scatter plots, choose the features you think are important and model an MLR. Justify your choice and explain the new findings.

## Problem-4 (1 Point)

Conduct a partial F-test to determine if the attributes not chosen by you in *Problem-3* are significant to predict the energy.What are the null and alternate hypotheses? [ *Hint* : Use the anova function between models created in *Problem-2* and *Problem-3*]

## Problem-5 (1.5 Points)

AIC - Akaike Information Criterion is used to compare different models and determine the best fit for the data. The best-fit model according to AIC is the one that explains greatest amount of variation using the fewest number of attributes. Check this resource to learn more about AIC.

Build a model based on AIC using Stepwise AIC regression.Elucidate your observations from the new model. ( *Hint* : Use an appropriate function in olsrr package.)

## Problem-6 (1 Point)

Plot the residuals of the models built till now and comment on it satisfying the assumptions of MLR.

## Problem-7 (2 Points)

For the model built in **Problem-2** , determine the presence of multicollinearity using VIF. Determine if there are outliers in the data using Cook's Distance. If you find any , remove the outliers and fit the model for *Problem-2* and see if the fit improves. [ *Hint* : All the relevant functions can be found in *olsrr* package. An observation can be termed as an outlier if it has a Cook's distance of more than 4/n where n is the number of records.]