



## Part 1, Data Discovery:

The dataset we have chosen is a list of the top 10000 songs on Spotify obtained from Kaggle, from the 1960's till now. Each row which represents a song specifies many details such as the danceability, the artists who participated in creating said song, labels, and which album the song itself belongs to among many other details. The dataset has 10000 entries.

## Normalisation:

### 0NF:

**Song(Track\_URI, Album\_URI, Artist\_URI, Added\_at, Track\_Num, Track\_Duration, Labels, Album\_Name, Album\_release\_date, Album\_Artist\_URI's, Album\_image\_URL, Copyright, Artist\_Name, Artist\_Genres)**

### 1NF:

Song(Track\_URI, Album\_URI, Track\_Num, Track\_Duration, Album\_Name,  
 Album\_image\_URL, Artist\_Name, Copyright\_year)  
 Artists(Track\_URI, Artist\_URI)  
 ArtistGenre(Arist\_URI, Genre)  
 Track\_Copyright(Track\_URI, Company\_ID)  
 Copyright(Company\_ID, Company\_name)  
 Track\_Added\_On(Track\_URI, year, month, day, time)  
 Album\_Artists(Album\_URI, Artist\_URI)  
 Album\_Release\_Date(Album\_URI, year, month, day)  
 Labels(Track\_URI, label)

For getting from 0NF to 1NF, we broke down all compound and multi-valued attributes into their own tables, with the compound and multi-valued attributes specifically being:

- Copyright -> Company\_ID's, Company\_name's, Copyright\_year (Specific to song)
- Labels: A song can have multiple labels
- Artist\_Genres: An artist can have multiple genres associated with them
- Added\_at: Made up of a day, month, and a year
- Album\_release\_date: Also made up of a day, month, and a year
- Company\_ID's: Multiple companies per song, so create a table for determining name from a Company\_ID, and create support table for all Company\_ID's associated with a song
- Album\_Artist\_URI's gets its own table because an album can have multiple artists associated with it

2NF:

Song(Track\_URI, Track\_Num, Track\_Duration, Copyright\_year)  
 Album(Album\_URI, Album\_Name, Album\_image\_URL)  
 Artist(Artist\_URI, Artist\_Name)  
 Artists\_Track(Track\_URI, Artist\_URI)  
 ArtistGenre(Arist\_URI, Genre)  
 Track\_Copyright(Track\_URI, Company\_ID)  
 Copyright(Company\_ID, Company\_name)  
 Album\_Artists(Album\_URI, Artist\_URI)  
 Album\_Release\_Date(Album\_URI, year, month, day)  
 Labels(Track\_URI, label)

For getting from 1NF to 2NF, we removed partial dependency by separating Album\_URI, Album\_Name, and Album\_image\_URL into its own table, and by separating Artist\_URI and Artist\_Name into its own table.

3NF:

Song(Track\_URI, Track\_Num, Track\_Duration, Copyright\_year)  
 Album(Album\_URI, Album\_Name, Album\_image\_URL)  
 Artist(Artist\_URI, Artist\_Name)  
 Artists\_Track(Track\_URI, Artist\_URI)

ArtistGenre(Arist\_URI, Genre)  
 Track\_Copyright(Track\_URI, Company\_name)  
 CopyRight(Company\_ID, Company\_name)  
 Track\_Added\_On(Track\_URI, year, month, day, time)  
 Album\_Artists(Album\_URI, Artist\_URI)  
 Album\_Release\_Date(Album\_URI, year, month, day)  
 Labels(Track\_URI, label)

BCNF: Since all of our determinants are superkeys, we are also in BCNF:

Song(Track\_URI, Track\_Num, Track\_Duration, Copyright\_year)  
 Album(Album\_URI, Album\_Name, Album\_image\_URL)  
 Artist(Artist\_URI, Artist\_Name)  
 Artists\_Track(Track\_URI, Artist\_URI)  
 ArtistGenre(Arist\_URI, Genre)  
 Track\_Copyright(Track\_URI, Company\_name)  
 CopyRight(Company\_ID, Company\_name)  
 Track\_Added\_On(Track\_URI, year, month, day, time)  
 Album\_Artists(Album\_URI, Artist\_URI)  
 Album\_Release\_Date(Album\_URI, year, month, day)  
 Labels(Track\_URI, label)

## Justification for participation/cardinality constraints:

- Song -> Label: A company can have/own many songs and a song could be created by many companies that collaborate to create a song. Therefore there is a many to many relationship. A song has to be owned by a company and not every company needs to own a song, and thus, we have total participation from the song and not from the company.
- Songs -> Album: Every song needs to be in an album and every album needs to have at least one song in it. Thus, there is total participation on both sides between the two. Now, in an album, there could be many songs but one song can only be in one album. Thus, there is a many to one relationship between the two. And the song always determines the album it is in.
- Album -> Artist: In an album there could be many artists performing in them and an Artist can perform in multiple albums. Therefore there is a many to many relationship. Now for the participation, an artist must perform in an album and the album must have at least one artist performing in it. This ensures a total participation.
- Artist -> Song: A song has to be created by an artist and an artist needs to be in every song, so there is a total participation from the song and not the artist. A song can have many artists and so can many artists create many songs, thus there is a many to many relationship between the two.

## References

*Top 10000 Songs on Spotify 1960-Now*. (2023, July). Kaggle. Retrieved October 15, 2023, from <https://www.kaggle.com/datasets/joebeachcapital/top-10000-spotify-songs-1960-now>