

## Descripción de Atributos en cada colección

-Artista:

Atributos: "\_id" "Artist"

Dimensión dataset: 676 2

Haciendo un *unique* se encuentra que este dataset no tiene valores repetidos y haciendo un *is.na* se encuentra que no hay missings.

-Charts

Atributos: "\_id" "Position" "Track\_Name" "Artist" "Streams" "URL" "week\_start" "week\_end"

Dimensión dataset: 63600 8

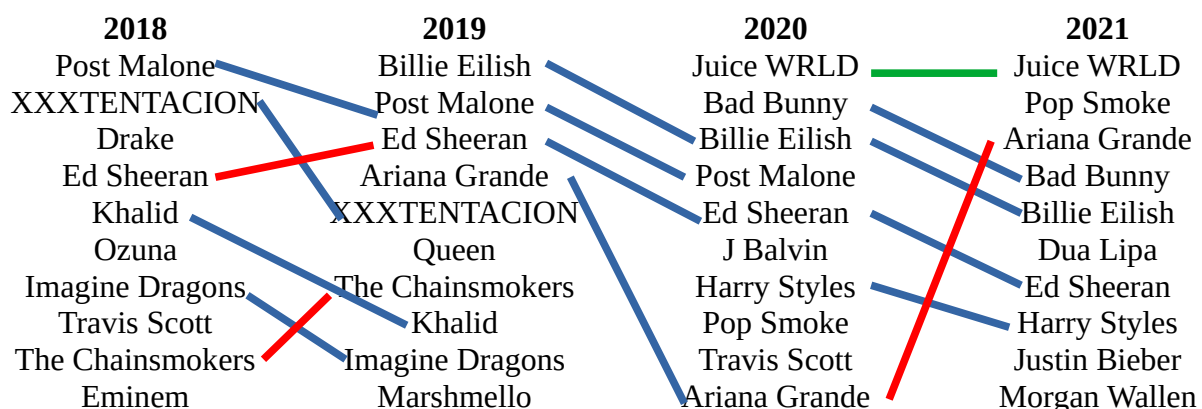
Haciendo un *is.na* se encuentra que no hay missings, pero haciendo una exploración se encuentra que hay filas de datos que si bien tienen distinto identificador, tienen los mismos datos. Unificando estas filas repetidas, el tamaño del nuevo dataset es:

Dimensión dataset: 31404 8

Cantidad de datos por año en dataset charts:

2018 2019 2020 2021  
20400 21200 21200 800

- Artistas con mas tracks en las listas por año (primeros 10 puestos):



El color y pendiente de las líneas da idea de que tan rapido ascienden o descienden en los rankings.

Los artistas que mas meten canciones en el ranking 200 son Post Malone, Billie Eilish, Ed Sheeran y Ariadna Grande (por ahora), que por tres años seguidos han metido varias canciones en este ranking.

## - Streams

Min.	1 <sup>st</sup> Qu.	Median	Mean	3 <sup>rd</sup> Qu.	Max.
3589018	5151455	6536358	8647701	9846342	71467874

El 50% de las reproducciones estan por debajo del orden de las ~6,53 millones de reproducciones y el maximo de reproducciones es del orden de ~71,47 millones de reproducciones. Del histograma se ve que las demas medidas de tendencia central no son representativas. El boxplot de Streams total y por año tiene muchos outliers por tanto no es una tecnica adecuada para analizar outliers.

Falta analisis outliers con otras tecnicas.

## -Artistas con mas reproducciones

Min.	1 <sup>st</sup> Qu.	Median	Mean	3 <sup>rd</sup> Qu.	Max.
3747900	5097888	6162743	6832971	7799853	21176356

El 50% de los artistas tienen en promedio por debajo de ~6,16 millones de reproducciones con un maximo de ~21,18 millones. Del histograma se ve que las demas medidas de tendencia central no son representativas. El boxplot tiene muchos outliers por tanto no es una tecnica adecuada para analizar outliers

Cuando se consideran las reproducciones medias y maximas, el ranking es el siguiente:

Mean	Max
Tones And I	Ariana Grande
SAINT JHN	Drake
Wham!	Shawn Mendes
Internet Money	Olivia Rodrigo
Jawsh 685	Ed Sheeran
Powfu	Mariah Carey
Mariah Carey	Bad Bunny
24kGoldn	The Weeknd
Bobby Helms	Tones And I
Joel Corry	Billie Eilish

Considerando reproducciones medias parecen meterse mas artistas antiguos en las listas, los artistas nuevos son los que meten hits muy reproducidos. Esto se puede considerar al decidir si analizar popularidad por reproducciones medias o maximas para cada artista.

## - Canciones Top

En total durante los años del analisis 39 canciones han ocupado el primer lugar.

Si se discrimina por año se ve que las canciones top cambian año a año, se podria discriminar mas fino para ver como cambian temporada a temporada pero creo que seria mas ruidoso puesto que no se saca tanta informacion debido a que son artistas que pegan un hit y no necesariamente que tengan una popularidad constante o consolidada. Lo que si se puede hacer es usar estas canciones para cruzar con los datos de features y observar porque son exitosas.

2018		2019		2020		2021	
<i>Havana (feat. Young Thug)</i>	Camila Cabello	<i>Sunflower - Spider-Man: Into the Spider-Verse</i>	Post Malone	<i>Dance Monkey</i>	Tones And I	<i>DÁKITI</i>	Bad Bunny
<i>God's Plan</i>	Drake	<i>7 rings</i>	Ariana Grande	<i>The Box</i>	Roddy Ricch	<i>drivers license</i>	Olivia Rodrigo
<i>Call Out My Name</i>	The Weeknd	<i>bad guy</i>	Billie Eilish	<i>Blinding Lights</i>	The Weeknd		
<i>Nice For What</i>	Drake	<i>I Don't Care (with Justin Bieber)</i>	Ed Sheeran	<i>THE SCOTTS</i>	THE SCOTTS		
<i>Better Now</i>	Post Malone	<i>Señorita</i>	Shawn Mendes	<i>Rain On Me (with Ariana Grande)</i>	Lady Gaga		
<i>This Is America</i>	Childish Gambino	<i>Circles</i>	Post Malone	<i>ROCKSTAR (feat. Roddy Ricch)</i>	DaBaby		
<i>SAD!</i>	XXXTENTACION	<i>HIGHEST IN THE ROOM</i>	Travis Scott	<i>cardigan</i>	Taylor Swift		
<i>Nonstop</i>	Drake	<i>Dance Monkey</i>	Tones And I	<i>Savage Love (Laxed - Siren Beat)</i>	Jawsh 685		
<i>In My Feelings</i>	Drake	<i>Lose You To Love Me</i>	Selena Gomez	<i>WAP (feat. Megan Thee Stallion)</i>	Cardi B		
<i>Lucky You (feat. Joyner Lucas)</i>	Eminem			<i>Mood (feat. iann dior)</i>	24kGoldn		

-Artistas con mas tracks top en las listas por año (primeros 5 puestos):

La diferencia con los que meten mas tracks en general,es que estos pueden estar aqui solo porque pegaron un hit. Se pueden cruzar estas listas con las de tracks en general para descartar los que son populares solo por un hit en un solo año y no vuelven a aparecer.

2018		2019		2020		2021	
Drake	4	Post Malone	2	24kGoldn	1	Bad Bunny	1
Ariana Grande	1	Ariana Grande	1	Ariana Grande	1	Olivia Rodrigo	1
Bad Bunny	1	Billie Eilish	1	Bad Bunny	1		
Camila Cabello	1	Ed Sheeran	1	Cardi B	1		
Childish Gambino	1	Selena Gomez	1	DaBaby	1		

- Audio Features

Atributos: "id" "acousticness" "album\_id" "album\_images" "album\_name"  
"album\_release\_date" "album\_release\_date\_precision" "album\_release\_year"  
"album\_type" "analysis\_url" "artist\_id" "artist\_name" "artists" "available\_markets"  
"danceability" "disc\_number" "duration\_ms" "energy" "explicit"  
"external\_urls\_spotify" "instrumentalness" "is\_local" "key"

"key\_mode" "key\_name" "liveness" "loudness" "mode" "mode\_name"  
 "speechiness" "tempo" "time\_signature" "track\_href" "track\_id"  
 "track\_name" "track\_number" "track\_preview\_url" "track\_uri" "type" "valence"

Dimensión dataset: 182049 40

Hay missings:

track\_preview\_url 77280  
 album\_release\_year 415

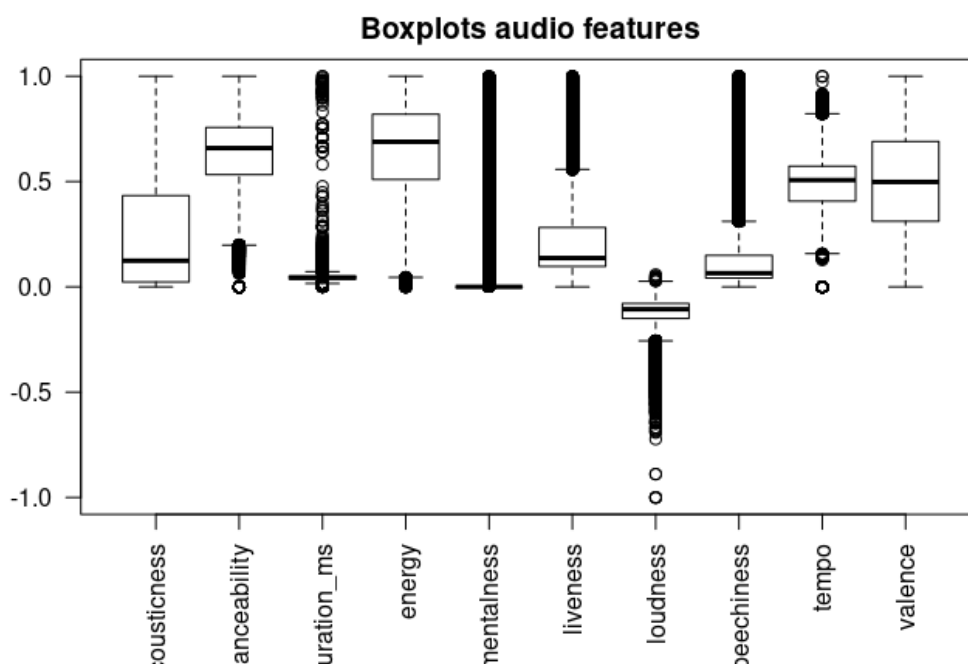
Como no son atributos relevantes podemos no considerarlos.

Haciendo un *unique* el dataset queda de dimensiones:

156621 40

Lo que quiere decir que hay filas 25428 repetidas. Se hace el analisis unificando estas filas.

Los boxplots normalizados de las audio features (para visualizarlos juntos) quedan asi:



para las variables:

"acousticness" "danceability" "duration\_ms" "energy" "instrumentalness"  
 "liveness" "loudness" "speechiness" "tempo" "valence"

Para la mayoría de features no sirve este analisis. Los valores de instrumentalness son muy pequeños la mayoría.

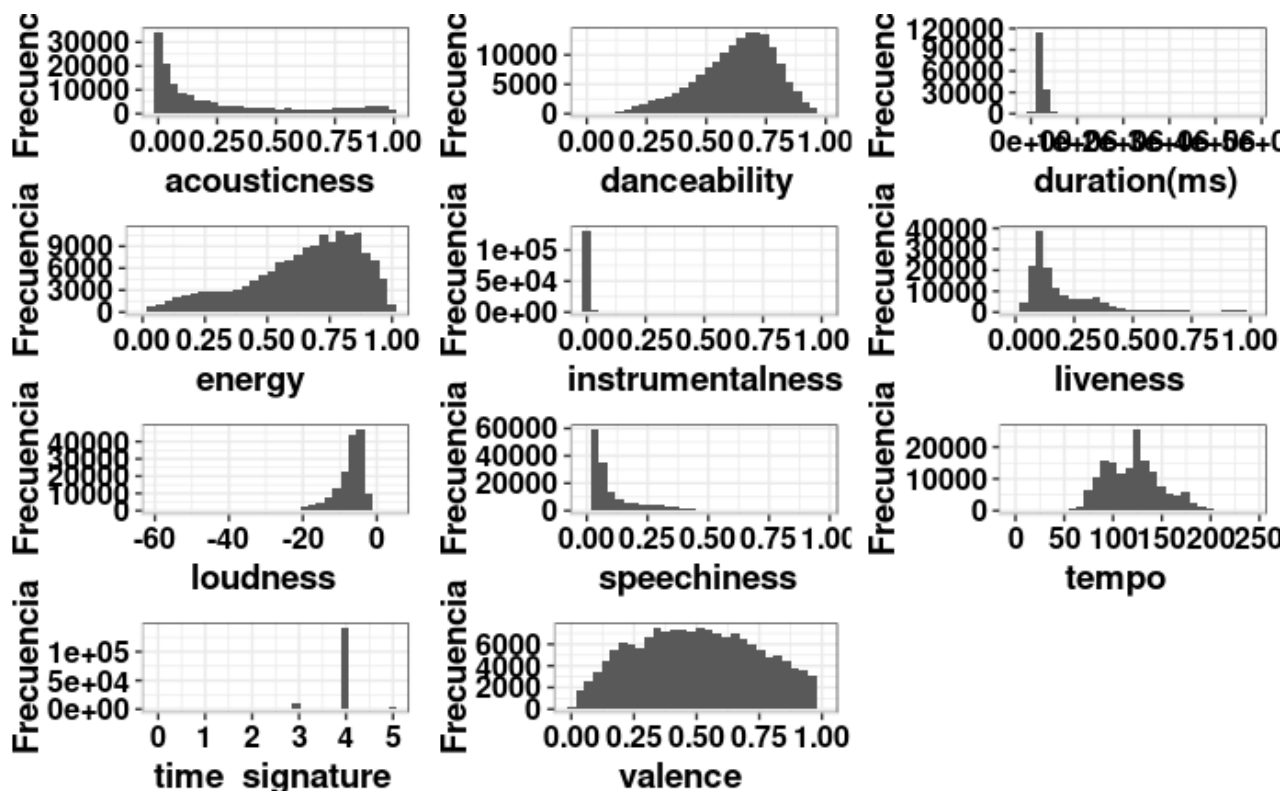
Teniendo la informacion de cuales artistas son los mas populares por numero de tracks en los listados y por tracks numero uno, podemos reducir nuestro analisis de features que hacen que una cancion o artista sea top a estas canciones y artistas.

Antes hay que cruzar la lista de artistas con los del archivo de audio features para que solo esten artistas del archivo de artistas. Al cruzar verificamos que esta todo ok.

Los resultados pueden cambiar al sacar outliers, pero la estructura de analisis que se acaba de describir se puede aplicar igual.

-Distribuciones

Histogramas de las variables audio features:



Observamos de la variable time\_signature que se puede deducir directamente que la gran mayoria de las canciones en los listados estan compuestas en compases simples 4 tiempos. Para las demas analizamos los boxplots y z\_scores para encontrar outliers.

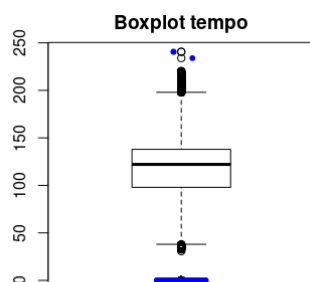
- Outliers

- Z-score

Analisis preliminar con metodos de z-score y z-score modificado para los audio features.

Analizamos con z-score las distribuciones mas simetricas, que vendrian siendo tiempo y valence.

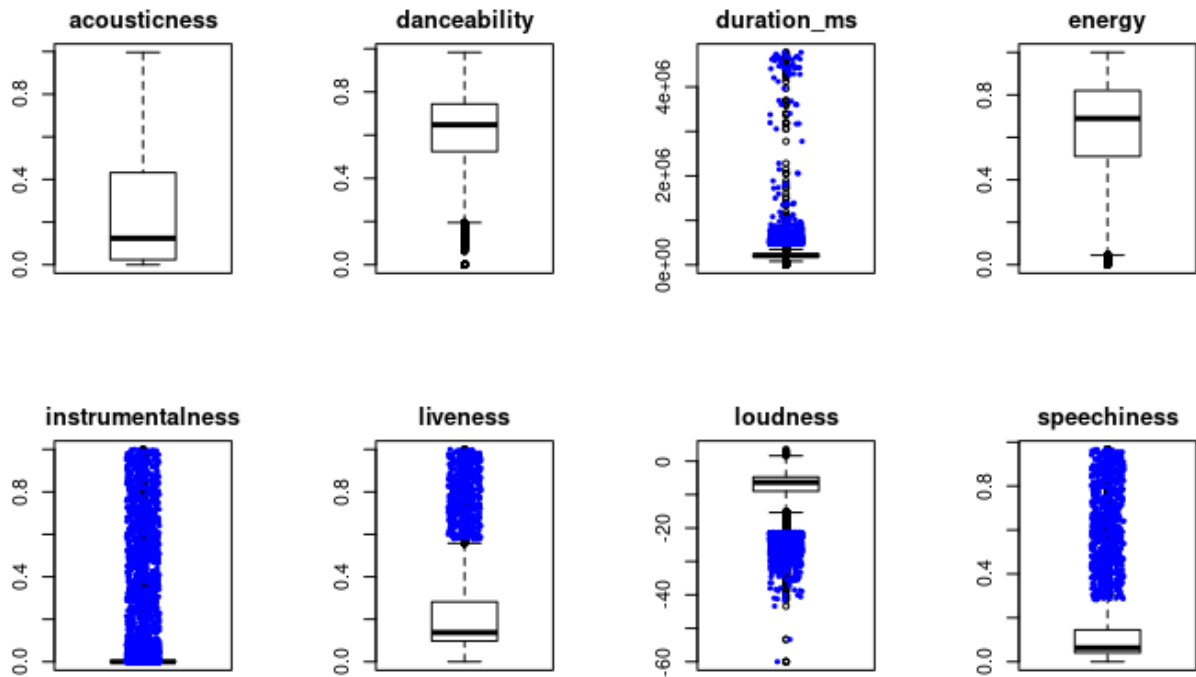
Para valence no se detectan outliers con este metodo, para tiempo:



Los datos que tienen tiempo 0 es porque seguramente los llenaron con ceros por no tener datos. Estos datos se pueden impugnar mas adelante.

-Z-score modificado

Para las otras variables de audio features:



Como se puede observar el metodo marca muchos puntos como outliers, esto es porque las distribuciones son muy asimetricas, por tanto debemos explorar otros metodos.

- Outliers variables dummy

Calculamos el z-score para *trackxartist* (cantidad de tracks que cada artista tiene metidas en las listas, tambien se discrimina por año), *streamsxartist* (reproducciones por artista, estan calculadas, reproducciones medias, donde se meten algunos artistas viejos y reproducciones maximas, donde priman los artistas nuevos), *top\_tracks* (canciones en el top, se discriminan tambien por año, esta es variable categorica), *NumberOnexartist* (artistas con mas canciones en el primer lugar, se discriminan por tambien por año).

Falta crear *streamsxtrack*, para determinar las canciones con mas reproducciones, este tendria que tener resultados parecidos a con *top\_track* que se hace con las posiciones.

Falta agrupar las features por los artistas top en listas (con mas canciones o con mas canciones numero uno) y por las canciones top en las listas. No se como se calcularian outliers aqui, o si mas bien aqui ya se entraria a modelar.

Y ahi voy...

despues de esto metemos lo de Marcos.