# Documentación Billboard.

## Rios Rocio Ayelen.

# Origen_to_raw

**Como primer paso extraemos los archivos de la base de datos y los guardamos en el Datalake.** Este paso lo realicé con una notebook de databricks en donde la conexión se hizo con

• **conexion_to_data-source:**



```python
jdbcHostname = "integrador-srv.database.windows.net"
jdbcPort = 1433
jdbcDatabase = "DB-Bllbd"
jdbcUsername = "Administrador"
jdbcPassword = "Formacion1"
jdbcDriver = "com.microsoft.sqlserver.jdbc.SQLServerDriver"

connectionProperties = {
  "user" : jdbcUsername,
  "password" : jdbcPassword,
  "driver" : jdbcDriver
}
jdbcUrl = "jdbc:sqlserver://{0}:{1};database={2};user={3}@integrador-srv;password={4};encrypt=true;trustServerCertificate=false;hostNameInCertificate=*.database.windows.net;loginTimeout=30;".format(jdbcHostname, jdbcPort, jdbcDatabase, jdbcUsername, jdbcPassword)
```

Luego de esto se generan los dataframes

## Generación de Dataframes

```python
artistDf = spark.read\
.format("jdbc")\
.option("driver", jdbcDriver)\
.option("url", jdbcUrl)\
.option("dbtable", "dbo.artistDf")\
.option("header", True)\
.load()
```

*Se realizó esta acción por cada tabla*
*[ artistDf, billboardHot100, grammyAlbums, grammySongs, riaaAlbumCerts, riaaSingleCerts, songAttribues, spotifyWeeklyTop200Streams]*

Luego se realizó la autenticación en el Datalake.

```python
access_key = "iBjJ9xc/hIfQyzVS+q74rigNxoyegk0qCO0xwLjs5HH6rwCNzczsv610q6QKHtySstLAvlH2cF0v+AStfFQvxg=="
spark.conf.set("fs.azure.account.key.dataintegrador.dfs.core.windows.net", access_key)
```

Se generó la conexión con routes.

```
1   %run "./routes"
```

routes contiene en su notebook:

```
1   raw = "abfss://integrador-rocio-rios@dataintegrador.dfs.core.windows.net/raw"
2   trusted = "abfss://integrador-rocio-rios@dataintegrador.dfs.core.windows.net/trusted"
```

## Luego de esto se hizo el **guardado de tablas en formato csv**

```
artistDf.write.save(f"{raw}/artistDf", format="csv", mode="overwrite",header='True')
```

*Se realizó esta acción con cada dataframe.*

Finalmente con la correcta ejecución de todos los pasos se realiza la **salida de Notebook**

```
1   dbutils.notebook.exit("Todo ok")
```

Notebook exited: Todo ok

# Raw_to_trusted

**En esta etapa comenzamos a limpiar y transformar los datos.**

Importamos librerías a utilizar:

```
 1    import pyspark.sql.functions as F
 2    from pyspark.sql import *
 3    from pyspark.sql.types import *
 4    from pyspark import *
 5    from pyspark.sql.functions import initcap
 6    from pyspark.sql.functions import col
 7    from pyspark.sql.functions import from_unixtime
 8    from pyspark.sql.functions import log10, when, abs
 9    from pyspark.sql.functions import to_date
10    from pyspark.sql.functions import concat, lit
11    from pyspark.sql.functions import coalesce
```

Generamos las conexiones necesarias:

**Conexion a data-source**

Cmd 4

```
 1    %run "./conexion_to_data-source"
```

Command took 1.12 seconds -- by rocioriosayelen@hotmail.com at 31/1/2023, 1:41:21 on Clustercito

Cmd 5

**Routes**

Cmd 6

```
 1    %run "./routes"
```

Command took 1.13 seconds -- by rocioriosayelen@hotmail.com at 31/1/2023, 1:41:21 on Clustercito

Cmd 7

**Configuracion de Storage Account**

Cmd 8

```
 1    access_key = 'iBjJ9xc/hIfQyzVS+q74rigNxoyegk0qCO0xwLjs5HH6rwCNzczsv6l0q6QKHtySstLAv1H2cF0v+AStfFQvxg=='
 2    spark.conf.set("fs.azure.account.key.dataintegrador.dfs.core.windows.net", access_key)
```

# Traer los Path

Cmd 10

```
 1    artistDfPATH = raw + '/artistDf'
 2    billboardhot100PATH = raw + "/billboardHot100"
 3    grammyAlbumsPATH = raw  + "/grammyAlbums"
 4    grammySongsPATH = raw  + "/grammySongs"
 5    riaaAlbumCertsPATH = raw  + "/riaaAlbumCerts"
 6    riaaSingleCertsPATH = raw  + "/riaaSingleCerts"
 7    songAttributesPATH = raw  + "/songAttributes"
 8    spotifyWeeklyTop200PATH = raw  + "/spotifyWeeklyTop200Streams"
```
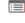
Se cargan los datasets guardados en csv:

**Cargar los Datasets**

Cmd 12

```
1    artistDf = spark.read.load(artistDfPATH, format='csv', sep=',', header='true')
2    billboardHot100 = spark.read.load(billboardhot100PATH, format='csv', sep=',', header='true')
3    grammyAlbums = spark.read.load(grammyAlbumsPATH, format='csv', sep=',', header='true')
4    grammySongs = spark.read.load(grammySongsPATH, format='csv', sep=',', header='true')
5    riaaAlbumCerts = spark.read.load(riaaAlbumCertsPATH, format='csv', sep=',', header='true')
6    riaaSingleCerts = spark.read.load(riaaSingleCertsPATH, format='csv', sep=',', header='true')
7    songAttributes = spark.read.load(songAttributesPATH, format='csv', sep=',', header='true')
8    spotifyWeeklyTop200 = spark.read.load(spotifyWeeklyTop200PATH, format='csv', sep=',', header='true')
```

▸ (8) Spark Jobs

▸ ▤ artistDf: pyspark.sql.dataframe.DataFrame = [X: string, Artist: string ... 6 more fields]

▸ ▤ billboardHot100: pyspark.sql.dataframe.DataFrame = [_c0: string, Unnamed: 0: string ... 8 more fields]

▸ ▤ grammyAlbums: pyspark.sql.dataframe.DataFrame = [_c0: string, Award: string ... 4 more fields]

▸ ▤ grammySongs: pyspark.sql.dataframe.DataFrame = [_c0: string, X: string ... 5 more fields]

▸ ▤ riaaAlbumCerts: pyspark.sql.dataframe.DataFrame = [_c0: string, Album: string ... 3 more fields]

▸ ▤ riaaSingleCerts: pyspark.sql.dataframe.DataFrame = [X: string, Name: string ... 3 more fields]

▸ ▤ songAttributes: pyspark.sql.dataframe.DataFrame = [_c0: string, Acousticness: string ... 16 more fields]

▸ ▤ spotifyWeeklyTop200: pyspark.sql.dataframe.DataFrame = [_c0: string, Name: string ... 4 more fields]

# Transformación y limpieza:

# artistDf

Cmd 15

```
1    artistDf = artistDf.drop("X")
2    artistDf = artistDf.dropDuplicates()
3    artistDf = artistDf.fillna("")
```

- Se eliminó la columna "X"
- Se eliminaron duplicados
- Se rellenaron valores nulls

# billboardHot100

```
1   billboardHot100 = billboardHot100.drop("_c0", "Unnamed: 0")
2   billboardHot100 = billboardHot100.dropDuplicates()
3   billboardHot100= billboardHot100.withColumn('Week',to_date(billboardHot100['Week'],'yyyy-MM-dd'))
4   billboardHot100= billboardHot100.fillna("")
5   billboardHot100 = billboardHot100.withColumnRenamed("Weekly.rank",'WeeklyRank')
6   billboardHot100 = billboardHot100.withColumnRenamed("Weeks.on.chart",'WeeksOnChart')
7   billboardHot100 = billboardHot100.withColumnRenamed("Peak.position",'PeakPosition')
```

▶ 🗐 billboardHot100: pyspark.sql.dataframe.DataFrame = [Artists: string, Name: string … 6 more fields]

- Se eliminaron columnas "_c0" y "Unnamed: 0"
- Se borraron duplicados
- Se cambió el tipo a fecha
- Se rellenaron valores nulls
- Se renombraron columnas

# grammyAlbums

```
1   grammyAlbums = grammyAlbums.drop("_c0")
2   grammyAlbums = grammyAlbums.withColumn("GrammyYear", grammyAlbums["GrammyYear"].cast(DateType()))
3   grammyAlbums = grammyAlbums.dropDuplicates()
```

▶ 🗐 grammyAlbums: pyspark.sql.dataframe.DataFrame = [Award: string, GrammyYear: date … 3 more fields]

- Se eliminó la columna "_c0"
- Se cambió el tipo de dato a tipo fecha
- Se eliminaron duplicados

# grammySongs

```
1   grammySongs = grammySongs.drop("_c0", "X")
2   grammySongs = grammySongs.withColumn("GrammyYear", grammySongs["GrammyYear"].cast(DateType()))
3   grammySongs = grammySongs.dropDuplicates()
4
```

- Se eliminó la columna "_c0"
- Se cambió el tipo de dato a tipo fecha
- Se eliminaron duplicados

# riaaAlbumCerts

Cmd 23

```
1  riaaAlbumCerts = riaaAlbumCerts.drop("_c0", "X")
2  riaaAlbumCerts = riaaAlbumCerts.withColumn("Label", initcap(col('Label')))
3  riaaAlbumCerts = riaaAlbumCerts.withColumn("Artist", initcap(col('Artist')))
4  riaaAlbumCerts = riaaAlbumCerts.withColumn("Album", initcap(col('Album')))
5  riaaAlbumCerts = riaaAlbumCerts.dropDuplicates()
6
```

▸ ▤ riaaAlbumCerts: pyspark.sql.dataframe.DataFrame = [Album: string, Artist: string … 2 more fields]

- Se eliminaron columnas "_c0" y "x"
- Se modificaron mayúsculas y minúsculas
- Se eliminaron duplicados

# riaaSingleCerts

Cmd 25

```
1  riaaSingleCerts = riaaSingleCerts.drop("X")
2  riaaSingleCerts = riaaSingleCerts.dropDuplicates()
3  riaaSingleCerts = riaaSingleCerts.fillna("")
4  riaaSingleCerts = riaaSingleCerts.withColumnRenamed("RiaaStatus",'Status')
```

▸ ▤ riaaSingleCerts: pyspark.sql.dataframe.DataFrame = [Name: string, Artist: string … 2 more fields]

- Se eliminó columna x
- Se eliminaron duplicados
- Se rellenaron nulls
- Se renombró columna "RiaaStatus" a "Status"

## songAttributes

Cmd 27

```
1   songAttributes = songAttributes.drop("_c0")
2   songAttributes = (songAttributes.withColumn("Acousticness", F.col("Acousticness").cast(FloatType())) # Convierte a tipo Float
3                    .withColumn("Danceability", F.col("Danceability").cast(FloatType())) # Convierte a tipo Float
4                    .withColumn("Energy", F.col("Energy").cast(FloatType()))
5                    .withColumn("Instrumentalness", F.col("Instrumentalness").cast(FloatType()))
6                    .withColumn("Liveness", F.col("Liveness").cast(FloatType()))
7                    .withColumn("Loudness", F.col("Loudness").cast(FloatType()))
8                    .withColumn("Speechiness", F.col("Speechiness").cast(FloatType()))
9                    .withColumn("Valence", F.col("Valence").cast(FloatType()))
10                   )
11
```

▸ ▤ songAttributes: pyspark.sql.dataframe.DataFrame = [Acousticness: float, Album: string … 15 more fields]

- Se eliminó columna "_c0"
- Se cambió el tipo de dato según corresponda

## spotifyWeeklyTop200

```
Cmd 29

1
2   spotifyWeeklyTop200 = spotifyWeeklyTop200.drop("_c0")
3   spotifyWeeklyTop200 = spotifyWeeklyTop200.withColumn("Week", to_date(spotifyWeeklyTop200["Week"], 'yyy-MM-dd'))
4   spotifyWeeklyTop200 = spotifyWeeklyTop200.withColumn("artistWithFt", coalesce(col("Artist"), col("Features")))
5   spotifyWeeklyTop200 = spotifyWeeklyTop200.withColumn("artistWithFt", concat(spotifyWeeklyTop200["Artist"], lit(" Ft. "), spotifyWeeklyTop200["Features"]))
6   spotifyWeeklyTop200 = spotifyWeeklyTop200.fillna("")
7
8
```

▶ 🖾 spotifyWeeklyTop200: pyspark.sql.dataframe.DataFrame = [Name: string, Artist: string ... 4 more fields]

- Se eliminó columna "_c0"
- Se cambió el tipo de dato a tipo fecha
- Se crearon columnas nuevas
- Se rellenaron nulls

Luego de la transformación y limpieza se guardan los dataframes en formato parquet:

## Guardar los Dataframes en Trusted

```
Cmd 31

1   artistDf.write.save(f"{trusted}/artistDf", format= "parquet", mode="overwrite")
2   billboardHot100.write.save(f"{trusted}/billboardHot100", format= "parquet", mode="overwrite")
3   grammyAlbums.write.save(f"{trusted}/grammyAlbums", format= "parquet", mode="overwrite")
4   grammySongs.write.save(f"{trusted}/grammySongs", format= "parquet", mode="overwrite")
5   riaaAlbumCerts.write.save(f"{trusted}/riaaAlbumCerts", format= "parquet", mode="overwrite")
6   riaaSingleCerts.write.save(f"{trusted}/riaaSingleCerts", format= "parquet", mode="overwrite")
7   songAttributes.write.save(f"{trusted}/songAttributes", format= "parquet", mode="overwrite")
8   spotifyWeeklyTop200.write.save(f"{trusted}/spotifyWeeklyTop200", format= "parquet", mode="overwrite")
9
```

Finalmente con la correcta ejecución de todos los pasos se realiza la **salida de Notebook**

## Salida de Notebook

```
Cmd 33

1   dbutils.notebook.exit("Proceso OK")
```

Notebook exited: Proceso OK

# Trusted_to_refined

En esta etapa generamos nuevas tablas o tablones (según corresponda el caso) para poder hacer consultas y análisis de los datos. Pueden surgir transformaciones, pero deberían estar en la etapa previa. Luego se exportan a PowerBI para generar los gráficos.

Importamos librerías a utilizar:

## Librerias

Cmd 2

```python
import pyspark.sql.functions as f
from pyspark.sql import SparkSession
from pyspark.sql.types import *
from pyspark.sql.functions import *
```

Creamos las conexiones necesarias

## Configuración de Storage Account

Cmd 4

```python
access_key = 'iBjJ9xc/hIfQyzVS+q74rigNxoyegk0qCO0xwLjs5HH6rwCNzczsv6l0q6QKHtySstLAv1H2cF0v+AStfFQvxg=='
spark.conf.set("fs.azure.account.key.dataintegrador.dfs.core.windows.net", access_key)
```

Command took 0.09 seconds -- by rocioriosayelen@hotmail.com at 31/1/2023, 7:31:21 on Clustercito

Cmd 5

## Traer Path

Cmd 6

```python
%run "./routes"
```

# Crear los Path

```
1  artistDfPATH = trusted + "/artistDf"
2  billboardHot100PATH = trusted + "/billboardHot100"
3  grammyAlbumsPATH = trusted + "/grammyAlbums"
4  grammySongsPATH = trusted + "/grammySongs"
5  riaaAlbumCertsPATH = trusted + "/riaaAlbumCerts"
6  riaaSingleCertsPATH = trusted + "/riaaSingleCerts"
7  songAttributesPATH = trusted + "/songAttributes"
8  spotifyWeeklyTop200PATH = trusted + "/spotifyWeeklyTop200"
9
```

## Cargar los Datasets

```
1  artistDf = spark.read.parquet(artistDfPATH)
2  billboardHot100 = spark.read.parquet(billboardHot100PATH)
3  grammyAlbums = spark.read.parquet(grammyAlbumsPATH)
4  grammySongs  = spark.read.parquet(grammySongsPATH)
5  riaaAlbumCerts= spark.read.parquet(riaaAlbumCertsPATH)
6  riaaSingleCerts= spark.read.parquet(riaaSingleCertsPATH)
7  songAttributes = spark.read.parquet(songAttributesPATH)
8  spotifyWeeklyTop200 = spark.read.parquet(spotifyWeeklyTop200PATH)
```

▶ (8) Spark Jobs

▶ ▤ artistDf:  pyspark.sql.dataframe.DataFrame = [Artist: string, Followers: string ... 5 more fields]

▶ ▤ billboardHot100:  pyspark.sql.dataframe.DataFrame = [Artists: string, Name: string ... 6 more fields]

▶ ▤ grammyAlbums:  pyspark.sql.dataframe.DataFrame = [Award: string, GrammyYear: date ... 3 more fields]

▶ ▤ grammySongs:  pyspark.sql.dataframe.DataFrame = [GrammyAward: string, GrammyYear: date ... 3 more fields]

▶ ▤ riaaAlbumCerts:  pyspark.sql.dataframe.DataFrame = [Album: string, Artist: string ... 2 more fields]

▶ ▤ riaaSingleCerts:  pyspark.sql.dataframe.DataFrame = [Name: string, Artist: string ... 2 more fields]

▶ ▤ songAttributes:  pyspark.sql.dataframe.DataFrame = [Acousticness: float, Album: string ... 15 more fields]

▶ ▤ spotifyWeeklyTop200:  pyspark.sql.dataframe.DataFrame = [Name: string, Artist: string ... 4 more fields]

Se crean las vistas temporales para poder hacer consultas SQL

## Crear Vistas Temporales de los Datasets

Cmd 13

```
1   artistDf.createOrReplaceTempView("artistDf")
2   billboardHot100.createOrReplaceTempView("billboardHot100")
3   grammyAlbums.createOrReplaceTempView("grammyAlbums")
4   grammySongs.createOrReplaceTempView("grammySongs")
5   riaaAlbumCerts.createOrReplaceTempView("riaaAlbumCerts")
6   riaaSingleCerts.createOrReplaceTempView("riaaSingleCerts")
7   songAttributes.createOrReplaceTempView("songAttributes")
8   spotifyWeeklyTop200.createOrReplaceTempView("spotifyWeeklyTop200")
```

En este caso, creamos 3 "mini tablones" para cada caso de uso:

Tablon1:

```
artistGrammy = spark.sql("""
SELECT
a.Artist,
a.Followers,
a.Gender,
ga.Award as AwardAlbum,
ga.GrammyYear as GrammyYearAlbum,
ga.Genre as GenreAlbum,
gs.GrammyAward as AwardSong,
gs.GrammyYear as GrammyYearSong,
gs.Genre as GenreSong
FROM
ArtistDf a
LEFT JOIN grammyAlbums ga ON a.Artist = ga.Artist
LEFT JOIN grammySongs gs ON a.Artist = gs.Artist
""")
```

🗐  artistGrammy:  pyspark.sql.dataframe.DataFrame = [Artist: string, Followers: string ... 7 more fields]

Tablon2:

```
1      billbSpotify = spark.sql("""
2      SELECT
3      b.Artists,
4      b.Name,
5      b.WeeklyRank,
6      b.WeeksOnChart,
7      b.Week as WeekBB,
8      b.Date,
9      b.Genre,
10     s.Name as NameSpoti,
11     s.Features,
12     s.Week as WeekSpoti
13     FROM
14     billboardHot100 b
15     LEFT JOIN spotifyWeeklyTop200 s
16     ON b.Artists = s.Artist
17     """)
18
19
```

▸ ▤ billbSpotify: pyspark.sql.dataframe.DataFrame = [Artists: string, Name: string ... 8 more fields]

Tablon3:

```
1    songs= spark.sql("""
2    SELECT
3    gs.Artist,
4    gs.Name as NameGrammy,
5    gs.GrammyAward,
6    gs.GrammyYear,
7    gs.Genre,
8    sa. Album,
9    sa.Danceability,
10   sa.Energy,
11   sa.Explicit,
12   sa.Name,
13   sa.Popularity
14   FROM
15   grammySongs gs
16   LEFT JOIN songAttributes sa
17   ON gs.Artist = sa.Artist
18   """)
```

▸ ▤ songs: pyspark.sql.dataframe.DataFrame = [Artist: string, NameGrammy: string ... 9 more fields]

Luego de la creación de los mini tablones se guardan en la base de datos esta vez "refinados"

**Guardado en base de datos**

```python
jdbcHostname = "integrador-srv.database.windows.net"
jdbcPort = 1433
jdbcDatabase = "Refined"
jdbcUsername = "Administrador"
jdbcPassword = "Formacion1"
jdbcDriver = "com.microsoft.sqlserver.jdbc.SQLServerDriver"

connectionProperties = {
    "user" : jdbcUsername,
    "password" : jdbcPassword,
    "driver" : jdbcDriver
}
jdbcUrl = "jdbc:sqlserver://{0}:{1};database={2};user={3}@integrador-srv;password={4};encrypt=true;trustServerCertificate=false;hostNameInCertificate=*.database.windows.net;loginTimeout=30;".format(jdbcHostname, jdbcPort, jdbcDatabase, jdbcUsername, jdbcPassword)
```

```python
artistGrammy.write.mode("overwrite")\
  .format("jdbc")\
  .option("driver", "com.microsoft.sqlserver.jdbc.SQLServerDriver")\
  .option("url", f"jdbc:sqlserver://{jdbcHostname}:{jdbcPort};databaseName={jdbcDatabase};user={jdbcUsername};password={jdbcPassword}")\
  .option("dbtable", "Rios_ArtistGrammy")\
  .option("truncate", "true")\
  .save()
```

```python
billbSpotify.write.mode("overwrite")\
  .format("jdbc")\
  .option("driver", "com.microsoft.sqlserver.jdbc.SQLServerDriver")\
  .option("url", f"jdbc:sqlserver://{jdbcHostname}:{jdbcPort};databaseName={jdbcDatabase};user={jdbcUsername};password={jdbcPassword}")\
  .option("dbtable", "Rios_billbSpotify")\
  .option("truncate", "true")\
  .save()
```

```python
songs.write.mode("overwrite")\
  .format("jdbc")\
  .option("driver", "com.microsoft.sqlserver.jdbc.SQLServerDriver")\
  .option("url", f"jdbc:sqlserver://{jdbcHostname}:{jdbcPort};databaseName={jdbcDatabase};user={jdbcUsername};password={jdbcPassword}")\
  .option("dbtable", "Rios_songs")\
  .option("truncate", "true")\
  .save()
```