(https://databricks.com)

## Assignment 5

setup environment (must run first)

%python

# Delete or comment this block out when running the notebook on the full dataset

# users\_df = spark.read.format("json").option("multiline", True).load("/FileStore/tables/zotmusic\_sample/Users.json")

# records\_df = spark.read.format("json").option("multiline", True).load("/FileStore/tables/zotmusic\_sample/Records.json")

# sessions\_df = spark.read.format("json").option("multiline", True).load("/FileStore/tables/zotmusic\_sample/Sessions.json")

# reviews\_df = spark.read.format("json").option("multiline", True).load("/FileStore/tables/zotmusic\_sample/Reviews.json")

# upvotes\_df = spark.read.format("json").option("multiline", True).load("/FileStore/tables/zotmusic\_sample/Upvotes.json")

# users\_df.createOrReplaceTempView("users")

# records\_df.createOrReplaceTempView("records")

# sessions\_df.createOrReplaceTempView("sessions")

# reviews\_df.createOrReplaceTempView("reviews")

# upvotes\_df.createOrReplaceTempView("reviews")

# upvotes\_df.createOrReplaceTempView("upvotes")

```
# %python
# # Uncomment this block when running the notebook on the full dataset (you must do this in the end before submitting the solution!)
users_df = spark.read.format("json").option("multiline", True).load("/FileStore/tables/zotmusic_full/Users.json")
records_df = spark.read.format("json").option("multiline", True).load("/FileStore/tables/zotmusic_full/Records.json")
sessions_df = spark.read.format("json").option("multiline", True).load("/FileStore/tables/zotmusic_full/Reviews.json")
reviews_df = spark.read.format("json").option("multiline", True).load("/FileStore/tables/zotmusic_full/Reviews.json")
upvotes_df = spark.read.format("json").option("multiline", True).load("/FileStore/tables/zotmusic_full/Upvotes.json")

users_df.createOrReplaceTempView("users")
records_df.createOrReplaceTempView("records")
sessions_df.createOrReplaceTempView("reviews")
upvotes_df.createOrReplaceTempView("reviews")
upvotes_df.createOrReplaceTempView("upvotes")
```

```
%sql
   (select "users", count(*) from users)
  union all
   (select "records", count(*) from records)
  union all
   (select "sessions", count(*) from sessions)
   union all
   (select "reviews", count(*) from reviews)
  union all
   (select "upvotes", count(*) from upvotes)
                                                                                                                       QTD
 Table
       ABc users
                     123 count(1)
  1
       users
                               5000
                             100000
  2
       records
  3
       sessions
                             889979
                             962143
  4
       reviews
  5
       upvotes
                            9189735
5 rows
This result is stored as _sqldf and can be used in other Python and SQL cells.
```

## Solution

```
8
  #1.A
  users_df.printSchema()
  #1.B address data type is: struct. It has 4 string fields: city, state, street and zip.
root
|-- address: struct (nullable = true)
     |-- city: string (nullable = true)
     |-- state: string (nullable = true)
     |-- street: string (nullable = true)
      |-- zip: string (nullable = true)
 |-- bio: string (nullable = true)
 |-- email: string (nullable = true)
 |-- genres: array (nullable = true)
     |-- element: string (containsNull = true)
 |-- is_artist: boolean (nullable = true)
 |-- is_listener: boolean (nullable = true)
 |-- joined_date: string (nullable = true)
```

9 %sql --1.C DESCRIBE records; --1.D songs data type is: array<struct<br/>bpm:bigint,length:bigint,mood:string,title:string,track\_number:bigint>> QTD Table ABc col\_name ABc data\_type ABc comment description null string 2 genre string null 3 is\_album boolean null is\_single boolean null 5 record\_id string null 6 released\_by struct<artist\_user\_id:string,release\_date:string> null 7 array<struct<br/>bpm:bigint,length:bigint,mood:string,title:string,track\_number:bigint>> null songs 8 title string null 9 video\_url string null 9 rows This result is stored as \_sqldf and can be used in other Python and SQL cells.



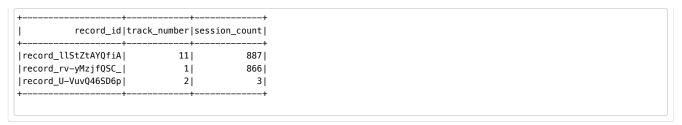
```
12
  #2.B DF
  pop_records_df = records_df.filter(records_df.genre == "Pop")
  pop_record_count_df = (
      pop_records_df.groupBy(pop_records_df.released_by.artist_user_id.alias("user_id"))
      .count()
      .orderBy("count", ascending=False)
  pop_record_count_df.limit(5).show()
          user_id|count|
|user_W5FmJadbR60F|
|user_giDyfRSuSRKk|
                      7|
|user_tylP-mPISHS1|
                      7|
                      7|
|user_h1_M1rkrRs-F|
|user_0lVVujY8RC6Y|
                      7|
```

<del>|-----</del>

```
13
 %sql
 --2.B SQL
 SELECT released_by.artist_user_id AS user_id, COUNT(record_id) AS num_pop_records
 FROM records
 WHERE genre = 'Pop'
 GROUP BY released_by.artist_user_id
 ORDER BY num_pop_records DESC
 LIMIT 5;
                                                                                                                  QVD
Table
     ABc user_id
                        123 num_pop_records
                                            8
1
     user_W5FmJadbR6OF
                                            7
2
     user_giDyfRSuSRKk
                                            7
3
     user_tylP-mPISHS1
                                            7
4
     user_h1_M1rkrRs-F
                                            7
     user_0IVVujY8RC6Y
```

This result is stored as \_sqldf and can be used in other Python and SQL cells.

5 rows



15 %sql --2.C SQL SELECT song.record\_id, song.track\_number, COUNT(\*) AS session\_count FROM sessions WHERE YEAR(session\_duration.initiate\_at) = 2023 AND MONTH(session\_duration.initiate\_at) = 11 GROUP BY song.record\_id, song.track\_number ORDER BY session\_count DESC LIMIT 3; Q70 Table ABc record\_id 123 track\_number 123 session\_count record\_IIStZtAYQfiA 11 887 1 2 record\_rv-yMzjfQSC\_ 1 866 2 record\_U-VuvQ46SD6p 3 3 rows

16

1 This result is stored as \_sqldf and can be used in other Python and SQL cells.



```
#2.E DF
  from pyspark.sql.functions import col, to_date
  filtered_reviews_df = reviews_df.filter(
      (to_date(col("posted_by.posted_time")) >= "2024-08-03") &
      (to_date(col("posted_by.posted_time")) < "2024-09-08")</pre>
  joined_df = filtered_reviews_df.join(users_df, col("posted_by.user_id") == col("user_id"))
  result_df = (
      joined_df.select("review_id", "email", col("posted_by.posted_time").alias("posted_time"))
      .orderBy(col("posted_time").asc())
      .limit(5)
  result_df.show()
          review_id|
                                    email
                                                  posted_time|
|review_aArS-aRvRyG5|robertsaunders@gm...|2024-08-03 00:00:31|
|review_0-V9LubWTBKR|biancadavis@iclou...|2024-08-03 00:01:00|
|review rNdhyxyxQxud|andrewwatkins@fox...|2024-08-03 00:02:01|
|review_xAEdPYgfSGiS|walleraudrey@gmai...|2024-08-03 00:02:07|
|review_2e4fbC-zSW08|riverakelsey@yaho...|2024-08-03 00:02:39|
```

```
19
 %sql
 --2.E SQL
 SELECT reviews.review_id, users.email, reviews.posted_by.posted_time AS posted_time
 FROM reviews
 JOIN users
 ON reviews.posted_by.user_id = users.user_id
 WHERE DATE(reviews.posted_by.posted_time) >= '2024-08-03'
 AND DATE(reviews.posted_by.posted_time) < '2024-09-08'
 ORDER BY reviews.posted_by.posted_time ASC
 LIMIT 5;
                                                                                                                     QTD
Table
      ABc review_id
                          <sup>B</sup>c email
                                                  ABc posted_time
```

5 rows	review_2e4fbC-zSWO8	riverakelsey@yahoo.com	2024-08-03 00:02:39
4	review_xAEdPYgfSGiS	walleraudrey@gmail.com	2024-08-03 00:02:07
3	review_rNdhyxyxQxud	andrewwatkins@foxmail.com	2024-08-03 00:02:01
2	review_O-V9LubWTBKR	biancadavis@icloud.com	2024-08-03 00:01:00
1	review_aArS-aRvRyG5	robertsaunders@gmail.com	2024-08-03 00:00:31

```
20
  #2.F DF
  from pyspark.sql.functions import col, count, desc
  reviews_df = reviews_df.withColumn("posted_by_user_id", col("posted_by.user_id"))
  upvotes_reviews_df = upvotes_df.join(reviews_df, upvotes_df["review_id"] == reviews_df["review_id"])
  upvotes_reviews_users_df = upvotes_reviews_df.join(
      users_df, upvotes_reviews_df["posted_by_user_id"] == users_df["user_id"]
  result_df = (
      upvotes_reviews_users_df.groupBy("email")
      .agg(count("*").alias("total_upvotes"))
      .orderBy(desc("total_upvotes"))
      .limit(3)
  result_df.show()
               email|total_upvotes|
|leejames@protonma...|
                               4682|
|orodriguez@gmail.com|
                               4626|
    qbrown@yahoo.com|
                               4622|
```

%sql --2.F SQL SELECT users.email, COUNT(upvotes.review\_id) AS total\_upvotes JOIN reviews ON upvotes.review\_id = reviews.review\_id JOIN users ON reviews.posted\_by.user\_id = users.user\_id GROUP BY users.email ORDER BY total\_upvotes DESC LIMIT 3;  $Q P \square$ Table <sup>B</sup><sub>C</sub> email 123 total\_upvotes leejames@protonmail.com 4682 1 2 orodriguez@gmail.com 4626 3 qbrown@yahoo.com 4622 3 rows • This result is stored as \_sqldf and can be used in other Python and SQL cells.

```
#2.G DF
  from pyspark.sql.functions import col, avg, count, desc
  active_listeners_df = sessions_df.groupBy("user_id").count().filter(col("count") >= 275)
  certified_reviews_df = (
      upvotes_df.groupBy("review_id").count().filter(col("count") >= 5)
      .join(reviews_df, "review_id")
      .join(active_listeners_df, reviews_df.posted_by.user_id == active_listeners_df.user_id)
  )
  artist_certified_ratings_df = certified_reviews_df.join(
      records_df, "record_id"
  ).select(
      records_df.released_by.artist_user_id.alias("artist_user_id"),
      certified_reviews_df.rating.alias("certified_rating")
  )
  artist_acrr_df = artist_certified_ratings_df.groupBy("artist_user_id").agg(
      avg("certified_rating").alias("acrr"),
      count("certified_rating").alias("total_certified_ratings")
  )
  users_df = users_df.withColumn("first_name", col("real_name.first_name"))
  users_df = users_df.withColumn("last_name", col("real_name.last_name"))
  artist_details_df = artist_acrr_df.join(
      users_df, artist_acrr_df.artist_user_id == users_df.user_id
  ).select(
      "artist_user_id",
      "acrr",
      "total_certified_ratings",
      "first_name",
      "last_name"
  )
  top_artists_df = artist_details_df.orderBy(desc("acrr")).limit(5)
  top_artists_df.show()
  artist_user_id|
                                acrr|total_certified_ratings|first_name| last_name|
|user_BNTubG2nQCa0|3.45833333333333335|
                                                          96| Kristen|
                                                                              Lopez |
                                                          98|
                                                                              NULL
|user_sEETlirSSTC5|3.4285714285714284|
                                                                   NULL
                                                                               NULL|
|user_158lawlPQSSz| 3.382716049382716|
                                                          81|
                                                                   NULL
```

	user_DYywyDQ0TPKC 3.3707865168539324  user_id0U530vTYOJ  3.369047619047619	89  84	Kelly  Baker  Jennifer Richardson			
+-						

```
%sql
--2.G SQL
WITH ActiveListeners AS (
    SELECT user_id
    FROM sessions
    GROUP BY user_id
    HAVING COUNT(*) >= 275
),
CertifiedReviews AS (
    SELECT
        reviews.review_id,
        reviews.rating AS certified_rating,
        reviews.record_id,
        reviews.posted_by.user_id AS listener_id
    FROM reviews
    JOIN (
        SELECT review_id, COUNT(*) AS upvote_count
        FROM upvotes
       GROUP BY review_id
        HAVING COUNT(*) >= 5
    ) AS upvote_counts
    ON reviews.review_id = upvote_counts.review_id
    WHERE reviews.posted_by.user_id IN (SELECT user_id FROM ActiveListeners)
),
ArtistCertifiedRatings AS (
    SELECT
        records.released_by.artist_user_id AS artist_user_id,
        CertifiedReviews.certified_rating
    FROM CertifiedReviews
    JOIN records
    ON CertifiedReviews.record_id = records.record_id
),
ArtistACRR AS (
    SELECT
        artist_user_id,
        AVG(certified_rating) AS acrr,
        COUNT(certified_rating) AS total_certified_ratings
    FROM ArtistCertifiedRatings
    GROUP BY artist_user_id
SELECT
    ArtistACRR.artist_user_id,
    ArtistACRR.acrr,
    ArtistACRR.total_certified_ratings,
    users.real_name.first_name AS first_name,
    users.real_name.last_name AS last_name
FROM ArtistACRR
JOIN users
ON ArtistACRR.artist_user_id = users.user_id
```

ORDER BY acrr DESC LIMIT 5; QTD Table ABc artist\_user\_id 1.2 acrr 123 total\_certified\_ratings ABc first\_name ABc last\_name user\_BNTubG2nQCaO 3.4583333333333335 96 Kristen Lopez user\_sEETlirSSTC5 3.4285714285714284 98 null null user\_158lawlPQSSz 3.382716049382716 81 null null 4 user\_DYywyDQ0TPKC 3.3707865168539324 89 Kelly Baker 5 user\_id0U530vTYOJ 3.369047619047619 84 Jennifer Richardson 5 rows

24

- # 3
- # Go BACK to the beginning and run the whole notebook on the full dataset, both in Databricks Community Version and Databricks Premium Version.
- # Write down your observations about the differences of using single-node community version vs. using multi-node premium version of Databricks.
- # The single-node version of Spark is limited by the hardware resources of a single machine, making it slower for
- # large-scale data processing. In contrast, the multi-node version distributes tasks across a cluster,
- # leveraging parallel computation to significantly improve speed, especially for big data and complex workloads.
- # As a result, the multi-node version is much faster for large-scale processing.

This result is stored as \_sqldf and can be used in other Python and SQL cells.