# TASK 2

In this project, the aim is to explore and analyse a dataset extracted from Steam, named steam-200k.csv, which contains details of games purchased and played by users along with the corresponding behavior and playtime. The dataset consists of four columns:
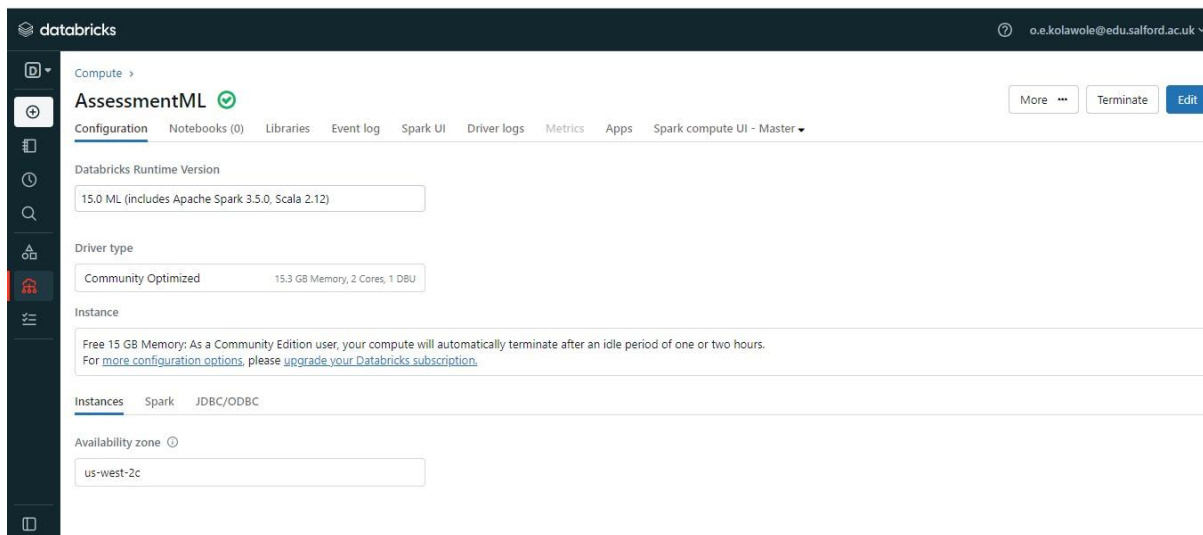
- A unique identifier for each member
- The name of the game they purchased or played
- Details of the member behaviour, either 'purchase' or 'play'. Because a game has to be purchased before it can be played, hence, there will be two entries for the same game/member combination in some instances
- This column is set to 1 for rows where the behavior is 'purchase'. For rows where the behavior is 'play' the value in the fourth column corresponds to the number of hours of play

The objective is to utilize this dataset to build a collaborative filtering recommender system using Apache Spark's MLlib. Collaborative filtering is a popular technique for making personalized recommendations by analyzing user behavior and similarities between users or items. By training a collaborative filtering model on the provided data, we aim to generate accurate recommendations for users based on their purchase and play behavior. By leveraging the power of collaborative filtering and Apache Spark's MLlib, we aim to provide valuable recommendations to Steam users, enhancing their gaming experience and driving engagement on the platform.

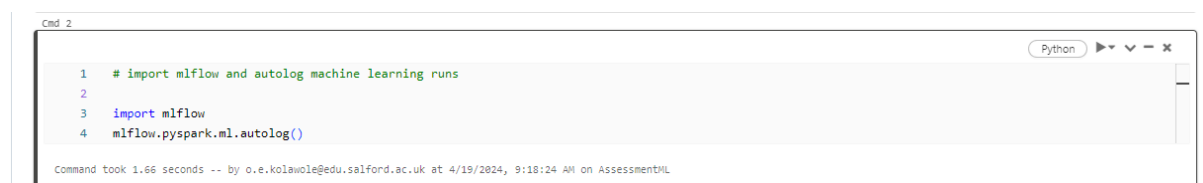To accomplish this task, The following steps were followed:

## SETUP

An ML compute as opposed to a Standard was used for this task. It has every library needed to run machine learning algorithms. The created compute is named "AssesmentML"
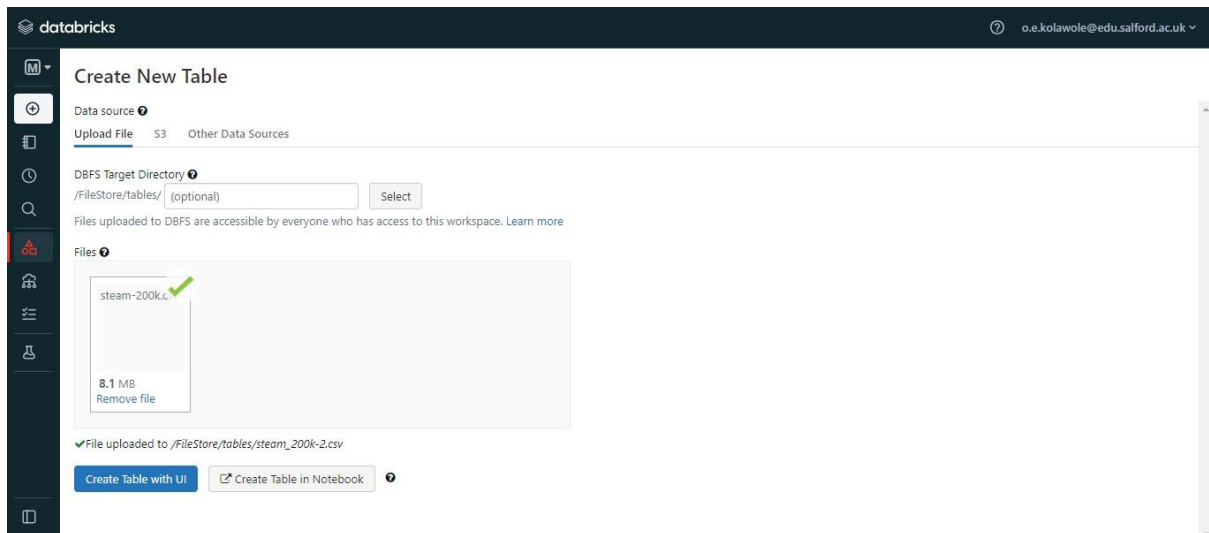
**Fig 3.4 Compute setup.**

The imported mlflow and autolog machine learning will automatically log important parameters, metrics, and models trained during the PySpark machine learning.
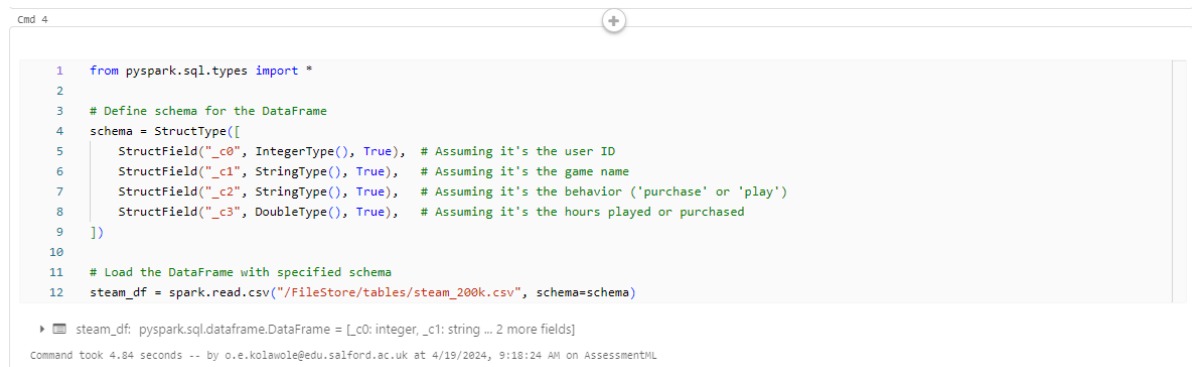


**Fig 3.5 ML flow.**

## LOADING DATA AND EXPLORATORY ANALYSIS

Loading the steam-200k.csv dataset into a Spark DataFrame and conducting initial exploratory analysis of the data to gain insights into the distribution of purchases and play behavior, the most popular games, and other relevant patterns.

**Fig 3.6 Uploading the dataset**



**Fig 3.7 Schema and data loading**

The DataFrame steam_df is created with the specified schema, providing a structured representation of the data. The data is read as a csv file into the databricks environment.

**Fig 3.8 Displaying the data frame**

The show() was used to display the first 20 rows of the steam_df to have a better understanding of the data



**Fig 3.9 Column renaming**

The column name of the data frame is renamed using the renaming mapping method and the new header is UserId, Game_Name, Behavior and Value giving each column better contextual understanding.

```python
1    # Summary statistics
2    steam_200kDF.describe().show()
```

▶ (2) Spark Jobs

```
+-------+--------------------+----------------+--------+------------------+
|summary|              UserId|       Game_Name|Behavior|             Value|
+-------+--------------------+----------------+--------+------------------+
|  count|              200000|          200000|  200000|            200000|
|   mean|   1.0365586594664E8|           140.0|    NULL|17.874383999999914|
| stddev| 7.208073512913968E7|             0.0|    NULL|138.05695165086792|
|    min|                5250|     007 Legends|    play|               0.1|
|    max|           309903146|theHunter Primal|purchase|           11754.0|
+-------+--------------------+----------------+--------+------------------+
```

Command took 14.49 seconds -- by o.e.kolawole@edu.salford.ac.uk at 4/19/2024, 9:18:24 AM on AssessmentML

Cmd 8

**Fig 4.0 Descriptive Statistics**

The descriptive statistics that provide insight into the distribution, central tendency, and variability of the Steam Dataset are a crucial component of the EDA analysis carried out for the task. A summary function was used to produce the mean and the median, which show the fundamental or important value of the findings.

```python
1    # Count the number of unique users and games
2    num_users = steam_200kDF.select('userId').distinct().count()
3    num_games = steam_200kDF.select('game_name').distinct().count()
4    print("Number of unique users:", num_users)
5    print("Number of unique games:", num_games)
```

▶ (6) Spark Jobs

```
Number of unique users: 12393
Number of unique games: 5155
```

Command took 7.55 seconds -- by o.e.kolawole@edu.salford.ac.uk at 4/19/2024, 9:18:24 AM on AssessmentML

**Fig 4.1 Unique features**

Counting the number of distinct users and games in the Data Frame helps in tracking performance metrics.

```python
1    # Number of Behavior records with the details
2    Behavior = steam_200kDF.groupBy('behavior').count()
3    Behavior.show()
```

▶ (2) Spark Jobs

▶ 🔲 Behavior: pyspark.sql.dataframe.DataFrame = [behavior: string, count: long]

```
+--------+------+
|behavior| count|
+--------+------+
|purchase|129511|
|    play| 70489|
+--------+------+
```

Command took 2.79 seconds -- by o.e.kolawole@edu.salford.ac.uk at 4/19/2024, 9:59:22 PM on AssessmentML

**Fig 4.2 Behavior of the members**

By displaying the number of records for every behavior category in your dataset, this code provides insight into the distribution of behaviors within it. It gives a better understanding of the relative prevalence and frequency of various behaviors within

the dataset. It shows the number of games played and purchased. The above code is then visualised using pyplot giving the below bar chart.



```python
import matplotlib.pyplot as plt

# Get member behavior counts
steam_counts = steam_200kDF.groupBy('behavior').count().toPandas()

# Define colors for the bars
colors = ['yellow', 'green']

# Create a bar chart of the counts
plt.bar(steam_counts['behavior'], steam_counts['count'], color=colors)

# Set the chart title and axis labels
plt.title('Behavior Distribution')
plt.xlabel('Behavior')
plt.ylabel('Count')

# Show the chart
plt.show()
```
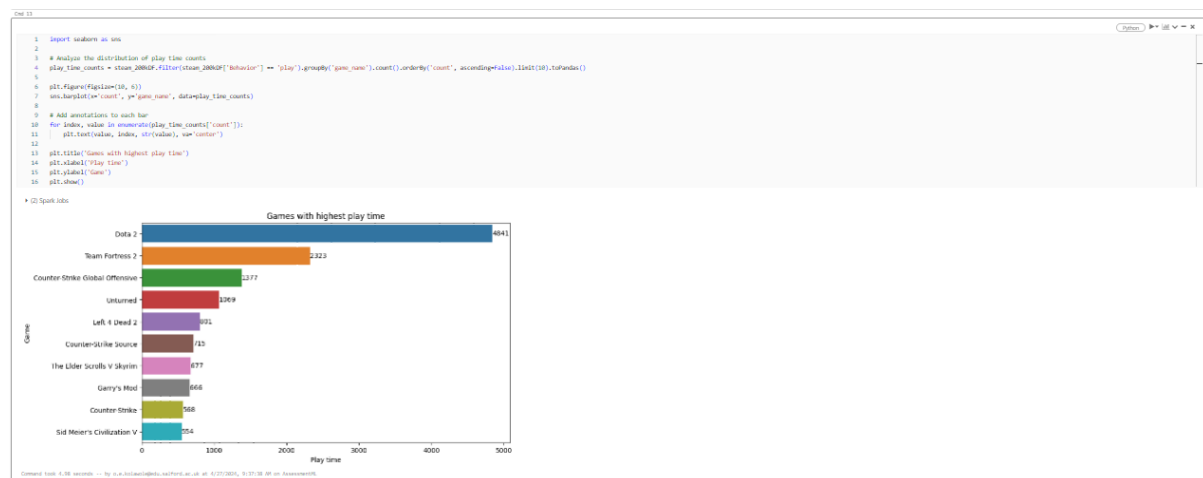
**Fig 4.3 Behavior of the members**

The code creates a bar chart of behavior distribution using Matplotlib, grouping DataFrame steam_200kDF by 'behavior' column and counting occurrences. It converts result to Pandas DataFrame, defines colors, creates bar chart with x-axis and y-axis, sets title, x-label, and y-label, and displays.

```
1   # Visualize distribution of play hours
2   import matplotlib.pyplot as plt
3
4   play_hours = steam_200kDF.filter(steam_200kDF.Behavior == 'play').select('value').rdd.flatMap(lambda x: x).collect()
5   plt.hist(play_hours, bins=20)
6   plt.title("Distribution of play hours")
7   plt.xlabel("Hours")
8   plt.ylabel("Frequency")
9   plt.show()
```

> (1) Spark Jobs



**Fig 4.4 Distribution of play hours in the play behavior of the members**

The code creates a Matplotlib histogram to display the distribution of play hours. It filters the DataFrame `steam_200kDF`, selects the 'value' column, flattens the RDD structure, collects values into a Python list, creates a histogram with 20 bins, sets titles, labels, and displays the histogram. The x-axis represents hours and the y-axis represents hour frequency.



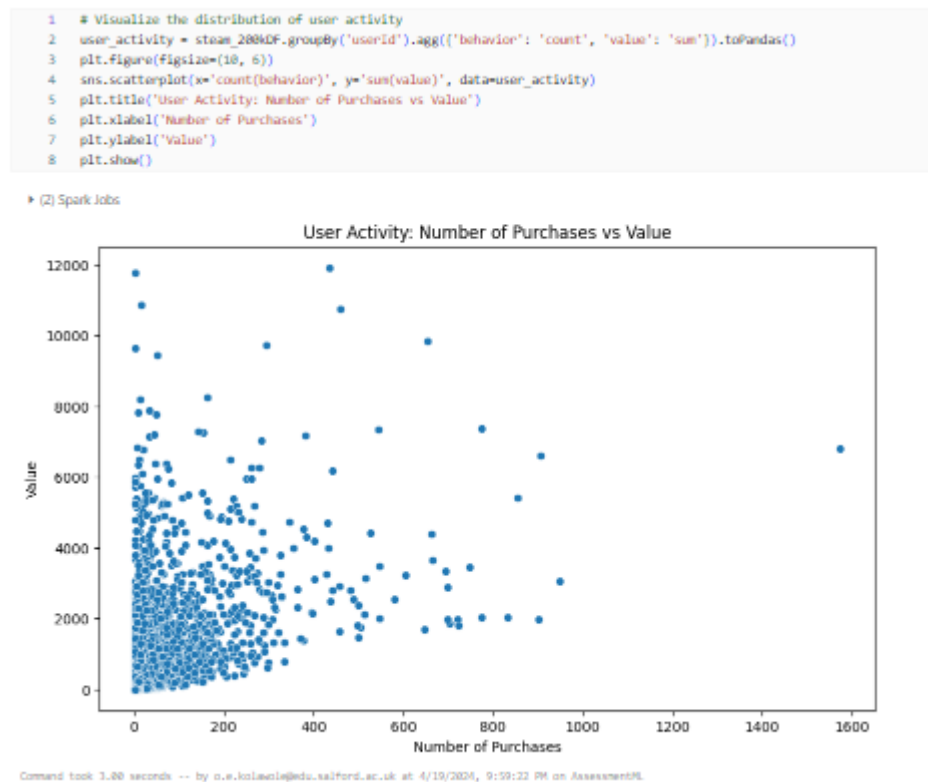**Fig 4.5 Games with the highest play time**

The code creates a bar plot using Seaborn to analyze the distribution of play time counts for the top 10 games. It filters the DataFrame `steam_200kDF`, groups it by 'game_name', and counts the occurrences of each game. The results are converted to a Pandas DataFrame, and a bar plot is created with 'count' as the x-axis and 'game_name' as the y-axis.

This gives additional description of the data, giving an insight of the behavior of the member. The highest play game is Dota 2 of almost 5000 hours.

```
1   import seaborn as sns
2   # Analyze the distribution of purchase counts
3   purchase_counts = steam_200kDF.filter(steam_200kDF['behavior'] == 'purchase').groupBy('game_name').count().orderBy('count', ascending=False).limit(10).toPandas()
4
5   plt.figure(figsize=(10, 6))
6   sns.barplot(x='count', y='game_name', data=purchase_counts)
7
8   # Add annotations to each bar
9   for index, value in enumerate(purchase_counts['count']):
10      plt.text(value, index, str(value), va='center')
11
12  plt.title('Top 10 Most Frequently Purchased Games')
13  plt.xlabel('Purchase Count')
14  plt.ylabel('Game')
15  plt.show()
```

**Fig 4.6 Frequently purchased games**

The code uses Seaborn and Matplotlib to visualize the distribution of purchase counts for the top 10 most frequently purchased games. It filters the DataFrame `steam_200kDF`, groups it by 'game_name', and counts the occurrences of each game. The results are converted to a Pandas DataFrame, and a horizontal bar plot is created with 'count' as the x-axis and 'game_name' as the y-axis.

```
1   # Visualize the distribution of user activity
2   user_activity = steam_200kDF.groupBy('userId').agg({'behavior': 'count', 'value': 'sum'}).toPandas()
3   plt.figure(figsize=(10, 6))
4   sns.scatterplot(x='count(behavior)', y='sum(value)', data=user_activity)
5   plt.title('User Activity: Number of Purchases vs Value')
6   plt.xlabel('Number of Purchases')
7   plt.ylabel('Value')
8   plt.show()
```

**Fig 4.7 User activity**

The code for visualizing user activity distribution using Seaborn and Matplotlib is well-structured. It groups the DataFrame `steam_200kDF` by 'userId', aggregates the count of 'behavior' and the sum of 'value' for each user, converts the result to a Pandas DataFrame, creates a scatter plot, sets title, x-label, and y-label, and displays the plot.



**Fig 4.8 Game Count**

The code correctly groups the DataFrame `steam_200kDF` by the 'Game_Name' column and counts the occurrences of each game. It then sorts the resulting DataFrame `game_counts` by the count column in descending order and displays the distinct records of game names and their respective counts. The sorted DataFrame is then displayed.

The game names along with their respective count was been generated to determine the users interest.

## DATA PREPARATION AND PREPROCESSING

We will preprocess the data to prepare it for training the recommender system. This may involve converting the data into the required format, splitting it into training and test sets, and generating unique integer IDs for the games.



**Fig 4.9 Generating GameId**

To create a distinct numerical identifier (GameId) for every distinct game name in the Data Frame steam_200kDF, StringIndexer from PySpark's MLlib was used. It displays the data frame with the new GameId column.

This process assigns a unique numerical identifier (GameId) to each distinct game name in the DataFrame, which can be useful for machine learning tasks that require numerical inputs.



**Fig 5.0 Generating the distinct behavior data frame**

The data frame is split into two based on the behavior. The filter() method keeps the rows needed in each newly created data frame based on the rows needed (purchase_df and play_df).

We would be using the purchase data frame for the analysis and the value for the purchase data frame is 1.



**Fig 5.1 Displaying the Purchase_df**

The newly created purchase_df is viewed.

## MODEL TRAINING, EVALUATION AND HYPERPARAMETER

The purchase_df DataFrame is split into training and test sets for training a collaborative filtering recommender system. The randomSplit() method is used to randomly split the data into two sets based on the provided weights which is 80% training and 20% testing and the seed parameter ensures reproducibility by specifying a seed for the random number generator.



**Fig 5.2 Splitting the purchase_df**

We will use MLlib's Alternating Least Squares (ALS) algorithm to train a collaborative filtering model on the prepared data which helps to assess its accuracy and effectiveness.

```
Cmd 25

                                                                                    Python  ► ▼ ∨ — ✕
    1    # Convert the DataFrame into the format required by the collaborative filtering algorithm(Training the Model)
    2    from pyspark.ml.recommendation import ALS
    3
    4    als = ALS(maxIter=5, regParam=0.01, userCol="UserId", itemCol="GameId", ratingCol="Value", coldStartStrategy="drop")
    5
    6    # Fit the ALS model on the training data
    7    model = als.fit(trainingDF)
```

▶ (6) Spark Jobs

▼ (1) MLflow run

    Logged 1 run to an experiment in MLflow. Learn more

    2024/04/19 21:01:22 INFO mlflow.utils.autologging_utils: Created MLflow autologging run with ID '2b60a775c1ff40fda2551e953e408f6c', which will track hyperparameters, performance me
    trics, model artifacts, and lineage information for the current pyspark.ml workflow
    2024/04/19 21:01:37 WARNING mlflow.utils.autologging_utils: MLflow autologging encountered a warning: "/databricks/python/lib/python3.11/site-packages/mlflow/types/utils.py:393: Us
    erWarning: Hint: Inferred schema contains integer column(s). Integer columns in Python cannot represent missing values. If your input data contains missing values at inference tim
    e, it will be encoded as floats and will cause a schema enforcement error. The best way to avoid this problem is to infer the model schema based on a realistic data sample (trainin
    g dataset) that includes missing values. Alternatively, you can declare integer columns as doubles (float64) whenever these columns may have missing values. See `Handling Integers
    With Missing Values <https://www.mlflow.org/docs/latest/models.html#handling-integers-with-missing-values>`_ for more details."
    2024/04/19 21:02:02 WARNING mlflow.pyspark.ml: Model ALS_96675b0325df will not be autologged because it is not allowlisted or or because one or more of its nested models are not al
    lowlisted. Call mlflow.spark.log_model() to explicitly log the model, or specify a custom allowlist via the spark.mlflow.pysparkml.autolog.logModelAllowlistFile Spark conf (see mlf
    low.pyspark.ml.autolog docs for more info).

    Command took 41.60 seconds -- by o.e.kolawole@edu.salford.ac.uk at 4/19/2024, 9:59:22 PM on AssessmentML
```

**Fig 5.3 Alternating Least Squares (ALS) algorithm training**

Inspecting the model's predictions on the test data is essential for evaluating the model's performance and understanding how well it generalizes to unseen data. This allows us to examine the model's predictions alongside the original test data. The predicted value can be compared to the purchase value so we can determine the difference between the values.



```
                                                                                    Python  ► ▼ ∨ — ✕
    1    # Make predictions on the test data
    2
    3    predictions = model.transform(testDF)
    4
    5    predictions.show()
```

▶ (4) Spark Jobs

▶ ▦  predictions:  pyspark.sql.dataframe.DataFrame = [UserId: integer, Game_Name: string ... 4 more fields]

```
+------+--------------------+--------+-----+------+----------+
|UserId|           Game_Name|Behavior|Value|GameId|prediction|
+------+--------------------+--------+-----+------+----------+
|  5250|Counter-Strike So...|purchase|  1.0|     5| 0.9743345|
|  5250|Half-Life 2 Death...|purchase|  1.0|    13|0.97621745|
|  5250| Half-Life Blue Shift|purchase|  1.0|    74|0.98232716|
|  5250|Team Fortress Cla...|purchase|  1.0|    79| 0.9821147|
| 76767|Arma 2 Operation ...|purchase|  1.0|   129| 0.9877363|
| 76767|Call of Duty Mode...|purchase|  1.0|    25|0.95887834|
| 76767|       Counter-Strike|purchase|  1.0|     6| 0.9056422|
| 76767|Counter-Strike Gl...|purchase|  1.0|     2| 0.9271606|
| 76767|       Day of Defeat|purchase|  1.0|    28|0.91442585|
| 76767|   Deathmatch Classic|purchase|  1.0|    34| 0.9136086|
| 76767| Half-Life Blue Shift|purchase|  1.0|    74|0.94661915|
| 76767|   Half-Life Opposin...|purchase| 1.0|    75| 0.9456815|
| 76767|             Thief 2|purchase|  1.0|   995| 1.0252678|
|103360|Counter-Strike Co...|purchase|  1.0|    19|0.95088756|
|103360|Counter-Strike Co...|purchase|  1.0|    23|0.95167166|
|103360|             Ricochet|purchase|  1.0|    35| 0.9655145|
|181212|           Half-Life 2|purchase| 1.0|    16| 1.0226392|
|181212|Team Fortress Cla...|purchase|  1.0|    79| 1.0026736|
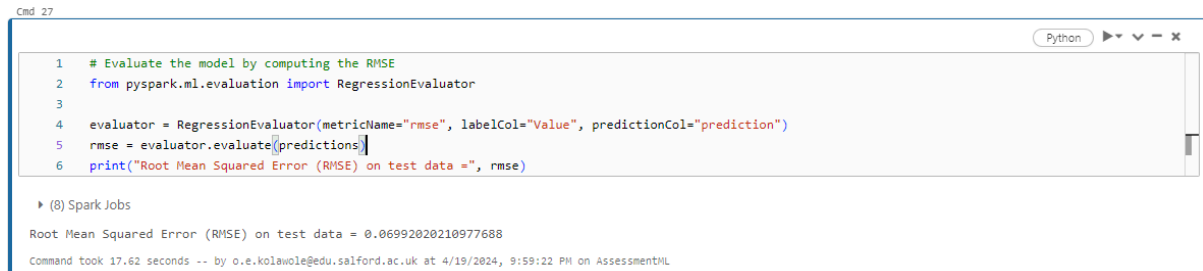+------+--------------------+--------+-----+------+----------+

Command took 8.91 seconds -- by o.e.kolawole@edu.salford.ac.uk at 4/19/2024, 9:59:22 PM on AssessmentML
```

**Fig 5.4 Prediction of test data**

The code uses PySpark's RegressionEvaluator to evaluate a regression model, importing it from pyspark.ml.evaluation. The RegressionEvaluator object is instantiated with metricName set to "rmse", labelCol set to "Value", and predictionCol set to "prediction". The evaluate method computes the RMSE on predictions.

To calculate the RMSE, the actual values (Value column) are compared to the predicted values (prediction column).

The average difference between the actual and predicted values is indicated by the RMSE value that you will receive in the output. Better model performance is indicated by lower RMSE values.

```python
# Evaluate the model by computing the RMSE
from pyspark.ml.evaluation import RegressionEvaluator

evaluator = RegressionEvaluator(metricName="rmse", labelCol="Value", predictionCol="prediction")
rmse = evaluator.evaluate(predictions)
print("Root Mean Squared Error (RMSE) on test data =", rmse)
```

▶ (8) Spark Jobs

Root Mean Squared Error (RMSE) on test data = 0.06992020210977688

Command took 17.62 seconds -- by o.e.kolawole@edu.salford.ac.uk at 4/19/2024, 9:59:22 PM on AssessmentML

**Fig 5.5 Model Evaluation**

A lower RMSE value indicates that the model's predictions are closer to the actual values, suggesting better performance. The RMSE of 0.0699 indicates that, on average, the model's predictions are off by approximately 0.0699 units from the actual values. This means that the model is performing well on the test data.

Generating top recommendations for each user using the ALS model and displaying recommendations for a specific user identified by their UserId. It retrieves the top 5 recommendations for the specified user and joins them with the Data Frame containing game names (purchase_df) to display the recommendations with their corresponding game names.

The result for the UserId 53875128 is displayed below showing the game name and rating given by the ALS model. Higher ratings indicate higher predicted preference or likelihood of the user enjoying the game.

```
1    from pyspark.sql.functions import explode, col
2
3    # Generate top recommendations for each user
4    userRecs = model.recommendForAllUsers(5)
5
6    # Show recommendations for a specific user (replace UserId with any user ID)
7    UserId = 53875128
8    user_recommendations = userRecs.select("UserId", explode("recommendations").alias("recommendation")) \
9                                   .filter(userRecs.UserId == UserId) \
10                                  .select("UserId", "recommendation.GameId", "recommendation.rating")
11
12   # Join with Game names to get the Game names
13   user_recommendations_with_names = user_recommendations.join(purchase_df.select("GameId", "Game_Name").distinct(),
14                                                              user_recommendations["GameId"] == purchase_df["GameId"],
15                                                              "left") \
16                                                     .select(user_recommendations["UserId"],
17                                                             purchase_df["Game_Name"],
18                                                             user_recommendations["rating"])
19
20   # Show recommendations in a tabular form
21   user_recommendations_with_names.show(truncate=False)
```

▶ (6) Spark Jobs

▶ 🔲 userRecs: pyspark.sql.dataframe.DataFrame = [UserId: integer, recommendations: array]

▶ 🔲 user_recommendations: pyspark.sql.dataframe.DataFrame = [UserId: integer, GameId: integer ... 1 more field]

▶ 🔲 user_recommendations_with_names: pyspark.sql.dataframe.DataFrame = [UserId: integer, Game_Name: string ... 1 more field]

```
+--------+------------------------------+---------+
|UserId  |Game_Name                     |rating   |
+--------+------------------------------+---------+
|53875128|IL-2 Sturmovik Cliffs of Dover|1.1497655|
|53875128|The Amazing Spider-Man        |1.1288923|
|53875128|Steel Ocean                   |1.1206661|
|53875128|Terrain Test                  |1.1191045|
|53875128|RACE Caterham Expansion       |1.1186852|
+--------+------------------------------+---------+
```

Command took 34.34 seconds -- by o.e.kolawole@edu.salford.ac.uk at 4/19/2024, 9:59:22 PM on AssessmentML

**Fig 5.6 Model Recommendations**

**Hyperparameter Tuning**

The code uses PySpark to perform hyperparameter tuning for the ALS model using cross-validation. It defines a parameter grid, a regression evaluator, and a cross-validator. The ALS estimator, grid, evaluator, and number of folds are set to 5. Cross-validation and hyperparameter tuning are performed using the `fit()` method. The best model is obtained, predictions are made on test data, and performance is evaluated based on RMSE. The best hyperparameters are printed out.

This ensures that the best-performing model is selected and its performance is evaluated using a reliable evaluation metric (RMSE).

It selects the best-performing model based on cross-validation and evaluates the performance and displays the rmse, and best rank and regularization parameter found during the hyperparameter tuning.

It's a crucial step in building accurate and effective recommendation systems.

```
1    from pyspark.ml.tuning import CrossValidator, ParamGridBuilder
2
3    # training models with different hyperparameters(Hyper parameter tunning)
4
5    # Define the parameter grid for hyperparameter tuning
6    param_grid = ParamGridBuilder() \
7        .addGrid(als.rank, [5, 10, 15]) \
8        .addGrid(als.regParam, [0.01, 0.1, 0.5]) \
9        .build()
10
11   # Define evaluator
12   evaluator = RegressionEvaluator(metricName="rmse", labelCol="Value", predictionCol="prediction")
13
14   # Define cross-validator
15   cross_validator = CrossValidator(estimator=als,
16                                    estimatorParamMaps=param_grid,
17                                    evaluator=evaluator,
18                                    numFolds=5)  # Use 5 folds for cross-validation
19
20   # Perform cross-validation and hyperparameter tuning
21   cv_model = cross_validator.fit(trainingDF)
22
23   # Get the best model from cross-validation
24   best_model = cv_model.bestModel
25
26   # Make predictions on the test data using the best model
27   predictions = best_model.transform(testDF)
28
29   # Evaluate the performance of the best model
30   rmse = evaluator.evaluate(predictions)
31   print("Root Mean Squared Error (RMSE) on test data after hyperparameter tuning:", rmse)
32
33   # Print the best hyperparameters found
34   print("Best Rank:", best_model.rank)
35   print("Best Regularization Parameter:", best_model._java_obj.parent().getRegParam())
```

▸ (9) Spark Jobs

▸ (10) MLflow runs

```
  ▸ (9) Spark Jobs
  ▾ (10) MLflow runs
      Logged 10 runs to an experiment in MLflow. Learn more
  ▸ ☐ predictions: pyspark.sql.dataframe.DataFrame = [UserId: integer, Game_Name: string ... 4 more fields]

  2024/04/19 21:03:05 INFO mlflow.utils.autologging_utils: Created MLflow autologging run with ID '16866ca9c46b4944826fc6cf3082975b', which will track hyperparam
  eters, performance metrics, model artifacts, and lineage information for the current pyspark.ml workflow
  2024/04/19 21:15:31 WARNING mlflow.pyspark.ml: Model CrossValidatorModel_c05f6d2e7a3b will not be autologged because it is not allowlisted or or because one or
  more of its nested models are not allowlisted. Call mlflow.spark.log_model() to explicitly log the model, or specify a custom allowlist via the spark.mlflow.py
  sparkml.autolog.logModelAllowlistFile Spark conf (see mlflow.pyspark.ml.autolog docs for more info).
  Root Mean Squared Error (RMSE) on test data after hyperparameter tuning: 0.07739710217212052
  Best Rank: 5
  Best Regularization Parameter: 0.01
  Command took 12.57 minutes -- by o.e.kolawole@edu.salford.ac.uk at 4/19/2024, 9:59:22 PM on AssessmentML
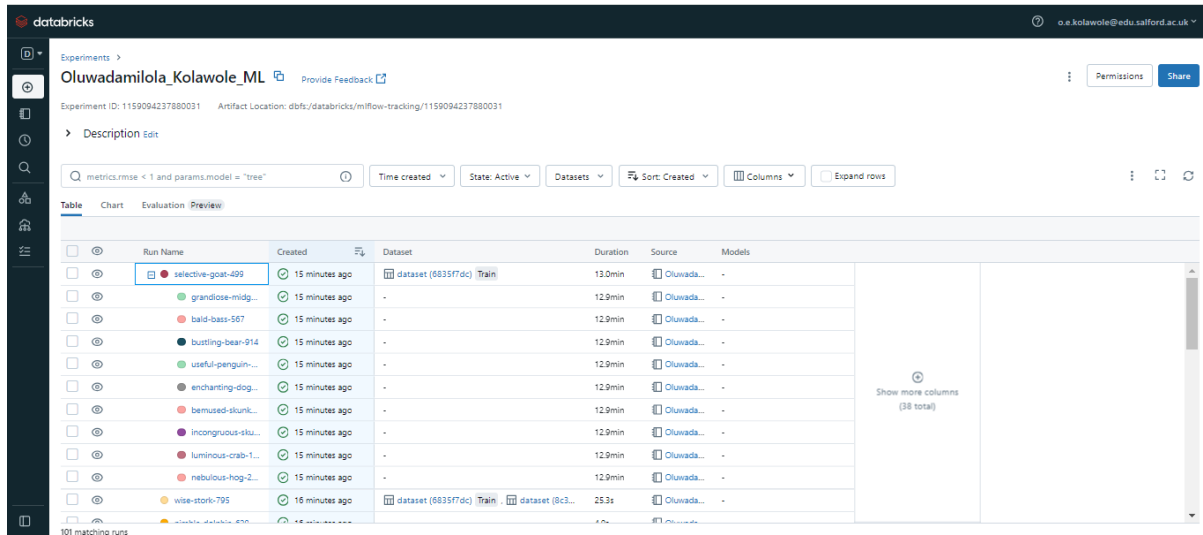```

**Fig 5.7 Hyper parameter tuning**

A lower RMSE indicates better performance, so an RMSE of 0.0774 suggests that the model's predictions are, on average, approximately 0.0774 units away from the actual values. This value could be considered satisfactory depending on the context of your specific problem and the range of the target variable.

Summarily, the tuned ALS model with a rank of 5 and a regularization parameter of 0.01 performs well based on the relatively low RMSE of 0.0774 on the test data. The

high level of accuracy shows that the model can accurately predict user preferences or behaviors because its predictions are typically close to the actual values.



**Fig 5.8: ML Experiment runs**

## DISCUSSION

When the ALS model's hyperparameters are tuned to achieve a low RMSE, it shows that the model is doing a good job of predicting user preferences. By providing specific and relevant recommendations, this precise recommendation system can increase user satisfaction and engagement in real-world applications.

When a platform recognizes its users' preferences and recommends content or products that suit their interests, users are more likely to interact with it. You can provide users with a more customized experience by utilizing the tuned ALS model, which could result in higher conversion rates, longer user retention periods, and happier users.

Furthermore, a precise recommendation system can support other corporate goals like raising revenue, boosting metrics for user engagement, and improving the user experience in general. As a result, implementing the tuned ALS model in practical applications could greatly benefit your platform or company.