## **Machine Learning - IV**

JHANVI PAREKH 60009210033 CSE (DS)

## **Experiment 9**

Link: https://colab.research.google.com/drive/1ci-e-5-bapj2xh8TZ76py8ith1v2iwfA?usp=sharing

```
!pip install pyspark
     Requirement already satisfied: pyspark in /usr/local/lib/python3.10/dist-packages (3.5.3)
     Requirement already satisfied: py4j==0.10.9.7 in /usr/local/lib/python3.10/dist-packages (from pyspark) (0.10.9.7)
import zipfile
import os
from pyspark.sql import SparkSession
# Unzip the dataset
with zipfile.ZipFile('ml-latest-small.zip', 'r') as zip_ref:
   zip_ref.extractall('ml-latest-small')
# Check the contents of the extracted folder
print("Contents of the extracted folder:")
print(os.listdir('ml-latest-small'))
    Contents of the extracted folder:
     ['ml-latest-small']
# Initialize a Spark session
spark = SparkSession.builder \
    .appName("MovieLens Analysis") \
    .getOrCreate()
# Check the contents of the extracted 'ml-latest-small' folder
print("Contents of the 'ml-latest-small' folder:")
print(os.listdir('ml-latest-small/ml-latest-small'))
    Contents of the 'ml-latest-small' folder:
     ['README.txt', 'links.csv', 'tags.csv', 'movies.csv', 'ratings.csv']
import zipfile
import os
# Unzip the dataset
with zipfile.ZipFile('ml-latest-small.zip', 'r') as zip_ref:
   zip_ref.extractall('ml-latest-small')
# Check the contents of the extracted folder
print("Contents of the 'ml-latest-small' folder:")
print(os.listdir('ml-latest-small'))
# Check contents of the inner directory if it exists
if 'ml-latest-small' in os.listdir('ml-latest-small'):
   print("Contents of the inner 'ml-latest-small' folder:")
   print(os.listdir('ml-latest-small/ml-latest-small'))
→ Contents of the 'ml-latest-small' folder:
     ['ml-latest-small']
     Contents of the inner 'ml-latest-small' folder:
     ['README.txt', 'links.csv', 'tags.csv', 'movies.csv', 'ratings.csv']
# Load the ratings dataset with the correct path
ratings_df = spark.read.csv("ml-latest-small/ml-latest-small/ratings.csv", header=True, inferSchema=True)
# Show the dataset structure
ratings_df.show()
     |userId|movieId|rating|timestamp|
                       4.0 964982703
                  3
                       4.0 964981247
          11
                  6
                       4.0 964982224
                  47
                       5.0|964983815
          11
          1
                  50
                       5.0 964982931
                  70
                        3.0 964982400
          1
```

101

5.0 964980868

```
110
                       4.0 964982176
          11
                151
                       5.0 964984041
                157
                       5.0 964984100
          1
          1
                163
                       5.0 964983650
                216
                       5.0 964981208
          1
                223
                       3.0 964980985
          1
                231
                       5.0 964981179
                235
                       4.0 964980908
          1
          1
                260
                       5.0 964981680
                296
                       3.0 964982967
          1
                316
                       3.0|964982310|
          11
          1
                333
                       5.0 964981179
                349
                      4.0 964982563
          1
                ----+
          --+-
     only showing top 20 rows
import os
# Check the contents of the extracted folder
print(os.listdir('ml-latest-small'))
→ ['ml-latest-small']
import zipfile
import os
from pyspark.sql import SparkSession
# Unzip the dataset
with zipfile.ZipFile('ml-latest-small.zip', 'r') as zip_ref:
   zip_ref.extractall('ml-latest-small')
# Check the contents of the extracted folder
print("Contents of the 'ml-latest-small' folder:")
print(os.listdir('ml-latest-small'))
# Check contents of the inner directory
print("Contents of the inner 'ml-latest-small' folder:")
print(os.listdir('ml-latest-small/ml-latest-small'))
# Initialize a Spark session
spark = SparkSession.builder \
    .appName("MovieLens Analysis") \
    .getOrCreate()
# Load the ratings dataset with the correct path
ratings_df = spark.read.csv("ml-latest-small/ml-latest-small/ratings.csv", header=True, inferSchema=True)
# Show the dataset structure
ratings_df.show()
# You can proceed with the analysis from here...
Tontents of the 'ml-latest-small' folder:
     ['ml-latest-small']
     Contents of the inner 'ml-latest-small' folder:
     ['README.txt', 'links.csv', 'tags.csv', 'movies.csv', 'ratings.csv']
     |userId|movieId|rating|timestamp|
        ---+----
                      4.0 964982703
                  3
                      4.0 964981247
          1
          1 |
                  6
                      4.0 964982224
                 47
                      5.0 964983815
                       5.0 964982931
                 50
          11
                 70
                       3.0 964982400
          1
                101
                       5.0 964980868
                110
                       4.0 964982176
                       5.0|964984041|
          11
                151
                       5.0 964984100
          1
                157
          1
                163
                       5.0 964983650
          1
                216
                       5.0 964981208
                223
                       3.0 964980985
          1
          1
                231
                       5.0 964981179
                235
                       4.0 964980908
          1
                260
                       5.0 964981680
          1
                296
                       3.0 964982967
```

```
3.0 964982310
                316
                      5.0 964981179
          11
                333
          1
                349
                     4.0 964982563
             -----+
    only showing top 20 rows
# List all the movies and the number of ratings
movies_count = ratings_df.groupBy("movieId").count()
movies_count.show()
     |movieId|count|
        1580 165
        2366
               25
        3175 l
                75 l
        1088
               42
       32460
                4
       44022
                23
       96488
                4
        1238
                9
        1342
               11
        1591
                26
        1645
                51
        4519
                9
        2142
                10
         471
                40
        3997
                12
        833
                6
        3918
                9
        7982
                4
        1959
                15
       68135
               10
    only showing top 20 rows
# List all the Movie IDs which have been rated
rated_movies = ratings_df.select("movieId").distinct()
rated_movies.show()
    +----+
     movieId
       1580
        2366
        3175
        1088
       32460
       44022
       96488
        1238
        1342
        1591
        1645
        4519
        2142
        471
        3997
         833
        3918
        7982
        1959
       68135
    +----+
    only showing top 20 rows
# List all the Users who have rated the movies
rated_users = ratings_df.select("userId").distinct()
rated_users.show()
    userId
```

```
https://colab.research.google.com/drive/1ci-e-5-bapj2xh8TZ76py8ith1v2iwfA#scrollTo=LV63V0A3E5Mf
```

```
10/27/24, 11:26 PM
```

```
463
        471
        496
        243
        392
        540
         31
        516
        85
        137
        251
        451
        580
        65
        458
         53
        255
        481
        588
     only showing top 20 rows
# List of all the Users with max, min, average ratings
```

```
).withColumnRenamed("max(rating)", "max_rating") \
.withColumnRenamed("min(rating)", "min_rating") \
```

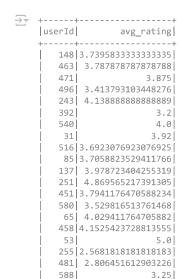
user\_ratings = ratings\_df.groupBy("userId") \

.withColumnRenamed("avg(rating)", "avg\_rating")

{"rating": "max", "rating": "min", "rating": "avg"}

user\_ratings.show()

.agg(



only showing top 20 rows

Start coding or generate with AI.