

## ▼ Import libraries

```
#https://grouplens.org/datasets/movielens/
```

```
pip install pyspark
```

```
Requirement already satisfied: pyspark in /usr/local/lib/python3.10/dist-packages (3.4.1)  
Requirement already satisfied: py4j==0.10.9.7 in /usr/local/lib/python3.10/dist-packages (from pyspark) (0.10.9.7)
```

```
from pyspark import SparkContext  
import pandas as pd  
from pyspark.sql.functions import col, explode  
from pyspark import SparkContext
```

## ▼ Initiate spark session

```
from pyspark.sql import SparkSession  
sc = SparkContext  
# sc.setCheckpointDir('checkpoint')  
spark = SparkSession.builder.appName('Recommendations').getOrCreate()
```

```
from google.colab import drive  
drive.mount('/content/drive')
```

## ▼ 1. Load data

```
movies = spark.read.csv('/content/drive/MyDrive/Colab Notebooks/movies.csv',header=True)  
ratings = spark.read.csv("/content/drive/MyDrive/Colab Notebooks/ratings.csv",header=True)
```

```
movies.show()  
ratings.show()
```

```
+-----+-----+-----+  
|movieId|          title|          genres|  
+-----+-----+-----+  
|      1| Toy Story (1995)|Adventure|Animati...| |
|      2| Jumanji (1995)|Adventure|Childre...|  
|      3|Grumpier Old Men ...|Comedy|Romance|  
|      4|Waiting to Exhale...|Comedy|Drama|Romance|  
|      5|Father of the Bri...|Comedy|  
|      6| Heat (1995)|Action|Crime|Thri...|  
|      7| Sabrina (1995)|Comedy|Romance|  
|      8| Tom and Huck (1995)|Adventure|Children|  
|      9| Sudden Death (1995)|Action|  
|     10| GoldenEye (1995)|Action|Adventure|...|  
|     11|American Presiden...|Comedy|Drama|Romance|  
|     12|Dracula: Dead and...|Comedy|Horror|  
|     13| Balto (1995)|Adventure|Animati...|  
|     14| Nixon (1995)|Drama|  
|     15|Cutthroat Island ...|Action|Adventure|...|  
|     16| Casino (1995)|Crime|Drama|  
|     17|Sense and Sensibi...|Drama|Romance|  
|     18| Four Rooms (1995)|Comedy|  
|     19|Ace Ventura: When...|Comedy|  
|     20| Money Train (1995)|Action|Comedy|Cri...|  
+-----+-----+-----+
```

only showing top 20 rows

```
+-----+-----+-----+-----+  
|userId|movieId|rating|timestamp|  
+-----+-----+-----+-----+  
|      1|      1|   4.0|964982703|  
|      1|      3|   4.0|964981247|  
|      1|      6|   4.0|964982224|  
|      1|     47|   5.0|964983815|  
|      1|     50|   5.0|964982931|  
|      1|     70|   3.0|964982400|  
|      1|    101|   5.0|964980868|  
|      1|    110|   4.0|964982176|  
|      1|    151|   5.0|964984041|  
|      1|    157|   5.0|964984100|  
|      1|    163|   5.0|964983650|  
|      1|    216|   5.0|964981208|
```

1	223	3.0	964980985
1	231	5.0	964981179
1	235	4.0	964980908
1	260	5.0	964981680
1	296	3.0	964982967
1	316	3.0	964982310
1	333	5.0	964981179
1	349	4.0	964982563

+-----+-----+-----+  
only showing top 20 rows

```
ratings.printSchema()
```

```
root
|-- userId: string (nullable = true)
|-- movieId: string (nullable = true)
|-- rating: string (nullable = true)
|-- timestamp: string (nullable = true)
```

```
ratings = ratings.\
  withColumn('userId', col('userId').cast('integer')).\
  withColumn('movieId', col('movieId').cast('integer')).\
  withColumn('rating', col('rating').cast('float')).\
  drop('timestamp')
ratings.show()
```

userId	movieId	rating
1	1	4.0
1	3	4.0
1	6	4.0
1	47	5.0
1	50	5.0
1	70	3.0
1	101	5.0
1	110	4.0
1	151	5.0
1	157	5.0
1	163	5.0
1	216	5.0
1	223	3.0
1	231	5.0
1	235	4.0
1	260	5.0
1	296	3.0
1	316	3.0
1	333	5.0
1	349	4.0

+-----+-----+-----+  
only showing top 20 rows

## ▼ Calculate sparsity

```
# Count the total number of ratings in the dataset
numerator = ratings.select("rating").count()

# Count the number of distinct userIds and distinct movieIds
num_users = ratings.select("userId").distinct().count()
num_movies = ratings.select("movieId").distinct().count()

# Set the denominator equal to the number of users multiplied by the number of movies
denominator = num_users * num_movies

# Divide the numerator by the denominator
sparsity = (1.0 - (numerator * 1.0) / denominator) * 100
print("The ratings dataframe is ", "%.2f" % sparsity + "% empty.")
```

The ratings dataframe is 98.30% empty.

## ▼ Interpret ratings

```
# Group data by userId, count ratings
userId_ratings = ratings.groupBy("userId").count().orderBy('count', ascending=False)
userId_ratings.show()
```

```
+-----+-----+
|userId|count|
+-----+-----+
|  414| 2698|
|  599| 2478|
|  474| 2108|
|  448| 1864|
|  274| 1346|
|  610| 1302|
|   68| 1260|
|  380| 1218|
|  606| 1115|
|  288| 1055|
|  249| 1046|
|  387| 1027|
|  182|  977|
|  307|  975|
|  603|  943|
|  298|  939|
|  177|  904|
|  318|  879|
|  232|  862|
|  480|  836|
+-----+-----+
only showing top 20 rows
```

```
# Group data by userId, count ratings
movieId_ratings = ratings.groupBy("movieId").count().orderBy('count', ascending=False)
movieId_ratings.show()
```

```
+-----+-----+
|movieId|count|
+-----+-----+
|   356|  329|
|   318|  317|
|   296|  307|
|   593|  279|
|  2571|  278|
|   260|  251|
|   480|  238|
|   110|  237|
|   589|  224|
|   527|  220|
|  2959|  218|
|     1|  215|
|  1196|  211|
|    50|  204|
|  2858|  204|
|    47|  203|
|   780|  202|
|   150|  201|
|  1198|  200|
|  4993|  198|
+-----+-----+
only showing top 20 rows
```

## ▼ Build Out An ALS Model

```
# Import the required functions
from pyspark.ml.evaluation import RegressionEvaluator
from pyspark.ml.recommendation import ALS
from pyspark.ml.tuning import ParamGridBuilder, CrossValidator
```

```
# Create test and train set
(train, test) = ratings.randomSplit([0.8, 0.2], seed = 1234)
```

```
# Create ALS model
als = ALS(userCol="userId", itemCol="movieId", ratingCol="rating", nonnegative = True, implicitPrefs = False, coldStartStrategy="drop")
```

```
# Confirm that a model called "als" was created
type(als)
```

```
pyspark.ml.recommendation.ALS
```

## ▼ Tell Spark how to tune your ALS model

```
# Import the requisite items
from pyspark.ml.evaluation import RegressionEvaluator
```

```

from pyspark.ml.tuning import ParamGridBuilder, CrossValidator

# Add hyperparameters and their respective values to param_grid
param_grid = ParamGridBuilder() \
    .addGrid(als.rank, [10, 50, 100, 150]) \
    .addGrid(als.regParam, [.01, .05, .1, .15]) \
    .build()
#         .addGrid(als.maxIter, [5, 50, 100, 200]) \

# Define evaluator as RMSE and print length of evaluator
evaluator = RegressionEvaluator(metricName="rmse", labelCol="rating", predictionCol="prediction")
print ("Num models to be tested: ", len(param_grid))

```

```
Num models to be tested:  16
```

## ▼ Build your cross validation pipeline

```

# Build cross validation using CrossValidator
cv = CrossValidator(estimator=als, estimatorParamMaps=param_grid, evaluator=evaluator, numFolds=5)

# Confirm cv was built
print(cv)

```

```
CrossValidator_49218e193661
```

## ▼ Best Model and Best Model Parameters

```

#Fit cross validator to the 'train' dataset
model = cv.fit(train)

```

```

#Extract best model from the cv model above
best_model = model.bestModel

```

```

# Print best_model
print(type(best_model))

# Complete the code below to extract the ALS model parameters
print("**Best Model**")

# # Print "Rank"
print("  Rank:", best_model._java_obj.parent().getRank())

# Print "MaxIter"
print("  MaxIter:", best_model._java_obj.parent().getMaxIter())

# Print "RegParam"
print("  RegParam:", best_model._java_obj.parent().getRegParam())

```

```

<class 'pyspark.ml.recommendation.ALSModel'>
**Best Model**
  Rank: 50
  MaxIter: 10
  RegParam: 0.15

```

```

# View the predictions
test_predictions = best_model.transform(test)
RMSE = evaluator.evaluate(test_predictions)
print(RMSE)

```

```
0.8692921299296869
```

```
test_predictions.show()
```

```

+-----+-----+-----+-----+
|userId|movieId|rating|prediction|
+-----+-----+-----+-----+
|  580|   1580|   4.0|   3.45169|
|  580|  44022|   3.5|   3.1713147|
|  597|    471|   2.0|   4.163025|
|  108|   1959|   5.0|   3.87913|
|  368|   2122|   2.0|   1.8805976|
|  436|    471|   3.0|   3.657902|
|  587|   1580|   4.0|   3.8653302|
|   27|   1580|   3.0|   3.3503187|
|  606|   1580|   2.5|   3.1741173|

```

```
| 606| 44022| 4.0| 2.7711148|
| 91| 2122| 4.0| 2.2577322|
| 157| 3175| 2.0| 3.4996016|
| 232| 1580| 3.5| 3.3879662|
| 232| 44022| 3.0| 3.107938|
| 246| 1645| 4.0| 3.7548573|
| 599| 2366| 3.0| 2.8609328|
| 111| 1088| 3.0| 3.189333|
| 111| 3175| 3.5| 3.0787396|
| 47| 1580| 1.5| 2.775062|
| 140| 1580| 3.0| 3.3890245|
+-----+-----+-----+
only showing top 20 rows
```

## ▼ Make Recommendations

```
# Generate n Recommendations for all users
nrecommendations = best_model.recommendForAllUsers(10)
nrecommendations.limit(10).show()
```

```
+-----+-----+
|userId| recommendations|
+-----+-----+
| 1| [{3379, 5.8084493...|
| 2| [{131724, 4.79579...|
| 3| [{6835, 4.852526}...|
| 4| [{3851, 4.8515286...|
| 5| [{3379, 4.5817575...|
| 6| [{33649, 4.750657...|
| 7| [{8477, 4.5287256...|
| 8| [{3379, 4.716567}...|
| 9| [{3379, 4.9013042...|
| 10| [{71579, 4.527629...|
+-----+-----+
```

```
nrecommendations = nrecommendations\
.withColumn("rec_exp", explode("recommendations"))\
.select('userId', col("rec_exp.movieId"), col("rec_exp.rating"))

nrecommendations.limit(10).show()
```

```
+-----+-----+-----+
|userId|movieId| rating|
+-----+-----+-----+
| 1| 3379|5.8084493|
| 1| 33649| 5.649271|
| 1| 5490| 5.531026|
| 1| 171495|5.4284987|
| 1| 8477|5.4112787|
| 1| 5915|5.3981776|
| 1| 5416|5.3959126|
| 1| 5328|5.3959126|
| 1| 3951|5.3959126|
| 1| 78836|5.3262763|
+-----+-----+-----+
```

## ▼ Do the recommendations make sense?

Lets merge movie name and genres to teh recommendation matrix for interpretability.

```
nrecommendations.join(movies, on='movieId').filter('userId = 100').show()
```

```
+-----+-----+-----+-----+-----+
|movieId|userId| rating| title| genres|
+-----+-----+-----+-----+-----+
| 67618| 100| 5.148249|Strictly Sexual (...|Comedy|Drama|Romance|
| 33649| 100| 5.053467| Saving Face (2004)|Comedy|Drama|Romance|
| 3379| 100|5.0466537| On the Beach (1959)| Drama|
| 42730| 100|4.9903436| Glory Road (2006)| Drama|
| 93008| 100| 4.919112|Very Potter Seque...| Comedy|Musical|
| 77846| 100| 4.919112| 12 Angry Men (1997)| Crime|Drama|
| 25906| 100| 4.919112|Mr. Skeffington (...| Drama|Romance|
| 184245| 100|4.8997493|De platte jungle ...| Documentary|
| 179135| 100|4.8997493|Blue Planet II (2...| Documentary|
| 138966| 100|4.8997493|Nasu: Summer in A...| Animation|
+-----+-----+-----+-----+-----+
```

```
nrecommendations.join(movies, on='movieId').filter('userId = 100').select('userId', 'movieId', 'rating').show()
```

```
ratings.join(movies, on= 'movieId' ).filter( 'userId' = 100 ).sort( 'rating' , ascending=false).limit(10).show()
```

movieId	userId	rating	title	genres
1101	100	5.0	Top Gun (1986)	Action Romance
1958	100	5.0	Terms of Endearme...	Comedy Drama
2423	100	5.0	Christmas Vacatio...	Comedy
4041	100	5.0	Officer and a Gen...	Drama Romance
5620	100	5.0	Sweet Home Alabam...	Comedy Romance
368	100	4.5	Maverick (1994)	Adventure Comedy ...
934	100	4.5	Father of the Bri...	Comedy
539	100	4.5	Sleepless in Seat...	Comedy Drama Romance
16	100	4.5	Casino (1995)	Crime Drama
553	100	4.5	Tombstone (1993)	Action Drama Western