# Movie Recommendation with MLlib - Collaborative Filtering



CS570 Big Data Processing Project By Feven Araya Instructor: Dr. Chang, Henry

## **Table of contents**

- 1. Introduction
- 2. Design
- 3. Implementation
- 4. Testing
- 5. Enhancement
- 6. Conclusion
- 7. References



# 01 Introduction





#### **Project Context**

- Introduction to the use of machine learning in the entertainment industry, specifically in recommending movies to users based on their preferences.
- Brief overview of the significance of personalized recommendations in enhancing user experience and engagement on movie platforms.

#### **Objective of the Project**

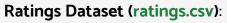
- To develop a machine learning model that accurately predicts user preferences and recommends movies using Collaborative Filtering techniques.
- Aim to leverage Apache Spark's MLlib for efficient processing of large-scale movie rating data from the MovieLens dataset.

### Title: Datasets

#### Content:

- Movies Dataset (movies.csv):
  - Contains metadata about movies.
  - o Columns: movield, title, genres.

movield	title	genres
1	Toy Story (1995)	Adventure Animation Children Comedy Fantasy
2	Jumanji (1995)	Adventure Children Fantasy
3	Grumpier Old Men (1995)	Comedy Romance
4	Waiting to Exhale (1995)	Comedy Drama Romance
5	Father of the Bride Part II (1995)	Comedy
6	Heat (1995)	Action Crime Thriller



- Contains user ratings for movies.
  - Columns: userld, movield, rating, timestamp.

userId	movield	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



# 02 Design

- Google Colab
- GCP





### **Identifying and Understanding the Problems**

- Dataset sparsity
- Cold start problem
- Objective: Provide accurate movie recommendations despite these challenges

#### **Investigating Solutions**

- Traditional collaborative filtering techniques
- Content-based filtering
- Hybrid models
- Our approach: Collaborative Filtering with ALS



#### **Theoretical Comparison and Selection**

- Comparison of different recommendation algorithms
- Benefits of ALS:
  - Handles large datasets efficiently
  - Suitable for implicit feedback
- Selection rationale: ALS for its scalability and effectiveness





# 03 Implementation



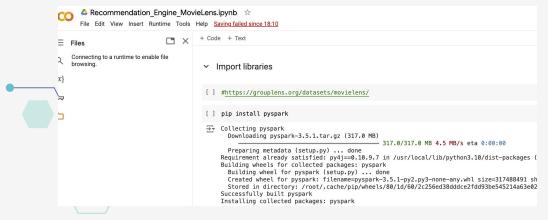


Step 3.1: <u>Download</u> the <u>Pyspark code (ipynb)</u>

Recommendation\_Engine\_MovieLens. ipynb

18.8 KB • Done

Step 3.2: <u>Upload the ipynb file</u> to your <u>Colab</u>





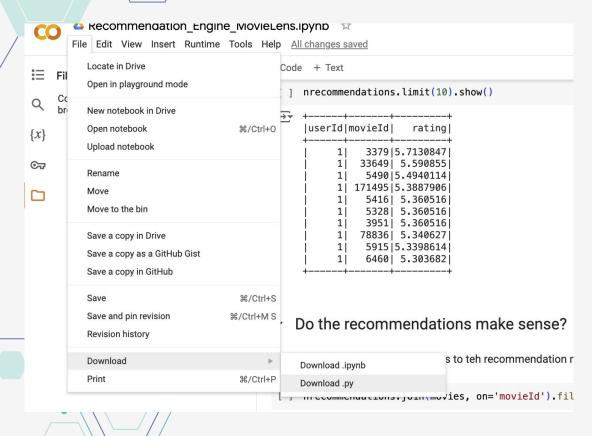
• Step 3.3: Experiment Pyspark code (ipynb) by modifying the ipynb file

```
# Build cross validation using CrossValidator
cv = CrossValidator(estimator=als, estimatorParamMaps=param_grid, evaluator=evaluator, numFolds=3)
# Confirm cv was built
print(cv)
```

Change number of folds to 5



#### Save the modified ipynb file as py and HTML format



• Step 3.5: Run the py file saved at Step 3.4 on GCP

#### **Packages**

```
import pandas as pd
from pyspark.sql.functions import col, explode
from pyspark import SparkContext, SparkConf
from pyspark.sql import SparkSession
from pyspark.ml.evaluation import RegressionEvaluator
from pyspark.ml.recommendation import ALS
from pyspark.ml.tuning import ParamGridBuilder, CrossValidator
```



# Initiate Spark Session

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



```
Load Data
```

```
[ ] movies = spark.read.csv("/content/movies.csv",header=True)
  ratings = spark.read.csv("/content/ratings.csv",header=True)
```



#### Display and Schema of Ratings Data

```
ratings.show()
 |userId|movieId|rating|timestamp|
                      4.0 | 964982703 |
                      4.0 | 964981247 |
                      4.0 | 964982224 |
               47 |
                      5.0 | 964983815 |
               50|
                      5.0 | 964982931 |
               70|
                      3.0|964982400|
       1
              101|
                      5.0 | 964980868 |
              110|
                      4.0 | 964982176 |
       1
              151
                      5.0 | 964984041 |
       1
              157|
                      5.0 | 964984100 |
       1
              163|
                      5.0 | 964983650 |
       1
              216|
                      5.0 | 964981208 |
              223|
                      3.0 | 964980985 |
              231
                      5.0 | 964981179 |
              235|
                      4.0 | 964980908 |
              2601
                      5.0 | 964981680 |
              296
                      3.0 | 964982967 |
       1
              316|
                      3.0 | 964982310 |
       1
              3331
                      5.0 | 964981179 |
                      4.0 | 964982563 |
              349|
only showing top 20 rows
```

# Data Preprocessing

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

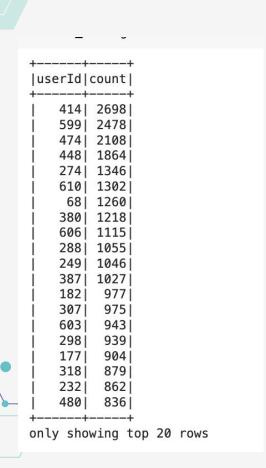


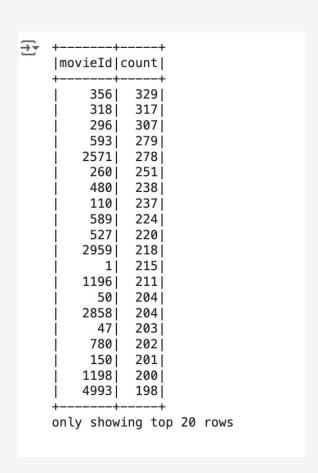
## Calculate Sparsity

```
# 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**





## **Build ALS Model**

```
# 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
 def __init__(*, rank: int=10, maxIter: int=10, regParam: float=0.1, numUserBlocks: int=10,
 numItemBlocks: int=10, implicitPrefs: bool=False, alpha: float=1.0, userCol: str='user',
 itemCol: str='item', seed: Optional[int]=None, ratingCol: str='rating', nonnegative:
 bool=False, checkpointInterval: int=10, intermediateStorageLevel: str='MEMORY AND DISK',
 finalStorageLevel: str='MEMORY_AND_DISK', coldStartStrategy: str='nan', blockSize: int=4096)
 /// Itcm_subsct_fccs - modet.fccommentarofitemsubsct(Itcm_subsct, 5/
 >>> item subset recs.select("recommendations.user", "recommendations.rating").first()
 Row(user=[0, 1, 2], rating=[3.910..., 3.473..., -0.899...])
 >>> als path = temp path + "/als"
```



>>> als.save(als\_path)
>>> als2 = ALS.load(als path)

>>> als.getMaxIter()

## Tune ALS Model

Num models to be tested: 16

```
# 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))
```

## Cross-Validation

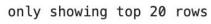
```
# Build cross validation using CrossValidator
cv = CrossValidator(estimator=als, estimatorParamMaps=param_grid, evaluator=evaluator, numFolds=5)
# Confirm cv was built
print(cv)
```

#### Train and Evaluate Model

TrossValidator\_d27636957c4c

### **Make Predictions**

#### |userId|movieId|rating|prediction| 148| 356 4.0 | 3.4951332 | 148 4896| 4.0 | 3.4835334 | 148 4993 | 3.465551 3.0 7153| 3.0 | 3.4216132 | 148 148 8368 4.0 3.591083 148 40629| 5.0 | 3.2217665 | 148 50872| 3.0 | 3.6663907 | 600691 4.5 | 3.695917 | 148 69757| 3.5 | 3.3879697 | 148 148 729981 4.0 | 3.2131975 148 81847 | 4.5 | 3.4920812 | 98491 5.0 | 3.7356784 | 148 3.5 | 3.5717542 | 148 115617 148 122886 3.5 | 3.4257748 | 463 4.149282 296 4.0 463 527| 4.0 | 3.7739785 | 463 4.0 | 3.9446247 | 2019| 471 527 4.5 3.773583 471| 6016| 4.0| 3.9766822| 471 6333| 2.5 | 3.2052839 |



test predictions.show()

#### **Generate Recommendations**

```
# Generate n Recommendations for all users
nrecommendations = best_model.recommendForAllUsers(10)
nrecommendations.limit(10).show()
             recommendations |
|userId|
      1|[{3379, 5.7130847...|
      2|[{131724, 4.79666...
      3|[{6835, 4.8578787...|
      4|[{3851, 4.8525457...|
      5|[{3379, 4.5449133...
      6|[{33649, 4.725941...
      7|[{33649, 4.459244...|
      8|[{3379, 4.635308}...
      9|[{3379, 4.7842216...
     10|[{71579, 4.533425...
```

#### Merge with Movies Data for Interpretability

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

genres	title	rating	userId	movieId
Comedy Drama Romance	Strictly Sexual (	5.0828342	100	67618
Drama	On the Beach (1959)	5.015384	100	3379
Comedy   Drama   Romance	Saving Face (2004)	5.0150394	100	33649
Drama	Glory Road (2006)	4.9038916	100	42730
Children Drama Ro	Anne of Green Gab	4.8903875	100	74282
Documentary	De platte jungle	4.8847737	100	184245
Documentary	Blue Planet II (2	4.8847737	100	179135
Animation	Nasu: Summer in A	4.8847737	100	138966
Documentary	Watermark (2014)	4.8847737	100	117531
Documentary	Connections (1978)	4.8847737	100	86237



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

```
|movieId|userId|rating|
                                        title
                                                              genres
   1101|
            100|
                    5.01
                              Top Gun (1986)|
                                                     Action | Romance |
   1958
            100|
                    5.0|Terms of Endearme...|
                                                        Comedy | Drama |
                   5.0 | Christmas Vacatio...|
   24231
            100|
                                                              Comedy
   40411
            100|
                   5.0|Officer and a Gen...|
                                                       Drama | Romance |
                   5.0|Sweet Home Alabam...|
   5620
            100|
                                                     Comedy | Romance |
    368|
            1001
                             Maverick (1994) | Adventure | Comedy | . . . |
            100|
                   4.5|Father of the Bri...|
    9341
                                                              Comedy
                   4.5|Sleepless in Seat...|Comedy|Drama|Romance|
    539|
            100|
                               Casino (1995)|
                   4.51
                                                         Crime | Drama |
     16
            1001
                            Tombstone (1993) | Action | Drama | Western |
    553|
            100|
                    4.51
```





# 04 Test

Process to test the project



#### Open GCP and upload your the recommendation\_Engine\_MovieLens.py file

```
"""Recommendation Engine MovieLens.ipynb
Automatically generated by Colab.
Original file is located at
    https://colab.research.google.com/drive/1wNSzqsOwDDH6bXQ-I-hC4Zc1nb0SD0WP
### Import libraries
import pandas as pd
from pyspark.sql.functions import col, explode
from pyspark import SparkContext, SparkConf
from pyspark.sql import SparkSession
from pyspark.ml.evaluation import RegressionEvaluator
from pyspark.ml.recommendation import ALS
from pyspark.ml.tuning import ParamGridBuilder, CrossValidator
"""### Initiate spark session"""
from pyspark.sql import SparkSession
sc = SparkContext
spark = SparkSession.builder.appName('Recommendations').getOrCreate()
"""# 1. Load data"""
movies = spark.read.csv("file:///home/faraya85431/movies.csv", header=True)
ratings = spark.read.csv("file:///home/faraya85431/ratings.csv",header=True)
ratings.show()
ratings.printSchema()
```

# Run the py file

```
faraya85431@cloudshell:~ (cs570-big-data-424622)$ spark-submit recommendation engine movielens.py
24/07/17 03:22:36 INFO SparkContext: Running Spark version 3.5.1
24/07/17 03:22:36 INFO SparkContext: 0S info Linux, 6.1.85+, amd64
24/07/17 03:22:36 INFO SparkContext: Java version 17.0.11
24/07/17 03:22:36 WARN NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable
24/07/17 03:22:36 INFO ResourceUtils: ------
24/07/17 03:22:36 INFO ResourceUtils: No custom resources configured for spark.driver.
24/07/17 03:22:36 INFO ResourceUtils: -----
24/07/17 03:22:36 INFO SparkContext: Submitted application: Recommendations
24/07/17 03:22:36 INFO ResourceProfile: Default ResourceProfile created, executor resources: Map(cores -> name: cores, amount: 1, script: , vendor: , memory
 1024, script: , vendor: , offHeap -> name: offHeap, amount: 0, script: , vendor: ), task resources: Map(cpus -> name: cpus, amount: 1.0)
24/07/17 03:22:36 INFO ResourceProfile: Limiting resource is cpu
24/07/17 03:22:36 INFO ResourceProfileManager: Added ResourceProfile id: 0
24/07/17 03:22:36 INFO SecurityManager: Changing view acls to: faraya85431
24/07/17 03:22:36 INFO SecurityManager: Changing modify acls to: faraya85431
24/07/17 03:22:36 INFO SecurityManager: Changing view acls groups to:
24/07/17 03:22:36 INFO SecurityManager: Changing modify acls groups to:
24/07/17 03:22:36 INFO SecurityManager: SecurityManager: authentication disabled; ui acls disabled; users with view permissions: faraya85431; groups with view
ers with modify permissions: faraya85431; groups with modify permissions: EMPTY
24/07/17 03:22:37 INFO Utils: Successfully started service 'sparkDriver' on port 44833.
24/07/17 03:22:37 INFO SparkEnv: Registering MapOutputTracker
24/07/17 03:22:37 INFO SparkEnv: Registering BlockManagerMaster
24/07/17 03:22:37 INFO BlockManagerMasterEndpoint: Using org.apache.spark.storage.DefaultTopologyMapper for getting topology information
```



```
/0//l/ 03:22:52 INFO TaskSchedulerimpi: Killing all running tasks in stage 2: Stage finished
24/07/17 03:22:52 INFO DAGScheduler: Job 2 finished: showString at NativeMethodAccessorImpl.java:0, took 0.255114 s
24/07/17 03:22:52 INFO BlockManagerInfo: Removed broadcast_4 piece0 on cs-2199409848-default:45279 in memory (size: 6.4 KiB, free: 434.3 MiB)
24/07/17 03:22:52 INFO CodeGenerator: Code generated in 25.1205 ms
 |userId|movieId|rating|timestamp|
              3| 4.0|964981247|
             471 5.019649838151
                   5.0[964982931]
             701 3.019649824001
            1011 5.019649808681
            110| 4.0|964982176|
            151| 5.0|964984041|
            1571 5.019649841001
                   5.019649836501
            231| 5.0|964981179|
            235| 4.0|964980908|
            260| 5.0|964981680|
            2961 3.019649829671
                   3.0[964982310]
only showing top 20 rows
 |-- userId: string (nullable = true)
 |-- movieId: string (nullable = true)
 |-- rating: string (nullable = true)
 |-- timestamp: string (nullable = true)
24/07/17 03:22:52 INFO SparkContext: Invoking stop() from shutdown hook
24/07/17 03:22:52 INFO SparkContext: SparkContext is stopping with exitCode 0. 24/07/17 03:22:52 INFO SparkUI: Stopped Spark web UI at http://cs-2199409848-default:4041
```

```
24/0//1/ 03:28:51 INFO CodeGenerator: Code generated in 9.590468 ms
24/07/17 03:28:51 INFO CodeGenerator: Code generated in 8.004966 ms
 +----+
 |userId|count|
 +----+
    414 | 2698 |
    5991 24781
    474| 2108|
    4481 18641
    2741 13461
    610 | 1302 |
     68| 1260|
    380| 1218|
    606| 1115|
    288 | 1055 |
    249 | 1046 |
    387 | 1027 |
    182| 977|
    3071 9751
    603| 943|
         9391
    2981
          9041
    318|
          8791
    2321
          8621
     4801 8361
 +----+
only showing top 20 rows
24/07/17 03:28:51 INFO SparkContext: Invoking stop() from shutdown hook
```

```
24/07/17 03:28:50 INFO DAGScheduler: ResultStage 18 (count at NativeMethodAccessorImpl.java:0) finished in 0.033 s
24/07/17 03:28:50 INFO DAGScheduler: Job 11 is finished. Cancelling potential speculative or zombie tasks for this job
24/07/17 03:28:50 INFO TaskSchedulerImpl: Killing all running tasks in stage 18: Stage finished
24/07/17 03:28:50 INFO DAGScheduler: Job 11 finished: count at NativeMethodAccessorImpl.java:0, took 0.048480 s
The ratings dataframe is 98.30% empty.
24/07/17 03:28:50 INFO FileSourceStrategy: Pushed Filters:
24/07/17 03:28:50 INFO FileSourceStrategy: Post-Scan Filters:
24/07/17 03:28:50 INFO FileSourceStrategy: Post-Scan Filters:
24/07/17 03:28:50 INFO CodeGenerator: Code generated in 62.327674 ms
24/07/17 03:28:50 INFO MemoryStore: Block broadcast_21 stored as values in memory (estimated size 199.3 KiB, free 433.6 MiB)
24/07/17 03:28:50 INFO MemoryStore: Block broadcast_21 piece0 stored as bytes in memory (estimated size 34.3 KiB, free 433.6 MiB)
24/07/17 03:28:50 INFO BlockManagerInfo: Added broadcast_21 piece0 in memory on cs-2199409848-default:43355 (size: 34.3 KiB, free: 4
24/07/17 03:28:50 INFO SparkContext: Created broadcast 21 from showString at NativeMethodAccessorImpl.java:0
24/07/17 03:28:50 INFO DAGScheduler: Registering RDD 58 (showString at NativeMethodAccessorImpl.java:0) as input to shuffle 5
24/07/17 03:28:50 INFO DAGScheduler: Got map stage job 12 (showString at NativeMethodAccessorImpl.java:0) with 1 output partitions
```

```
24/07/17 03:36:25 INFO DAGScheduler: ResultStage 24 (showString at NativeMethodAccessorImpl.java:0) finished in 0.075
24/07/17 03:36:25 INFO DAGScheduler: Job 15 is finished. Cancelling potential speculative or zombie tasks for this job
24/07/17 03:36:25 INFO TaskSchedulerImpl: Killing all running tasks in stage 24: Stage finished
24/07/17 03:36:25 INFO DAGScheduler: Job 15 finished: showString at NativeMethodAccessorImpl.java:0, took 0.092091 s
|movieId|count|
    3561 3291
    318 | 317 |
        3071
   25711
        2781
    2601
         2511
    4801
        2381
    110|
         2371
        224|
    5271
         2201
   29591
        2181
        2151
                                               24/U//1/ U3:34:U5 INFO DAGSCREQUIET: KESULTSTage 9048 (ShowString at NativeMethodAcc
        2111
                                               24/07/17 03:34:05 INFO DAGScheduler: Job 935 is finished. Cancelling potential specu
        204
                                               24/07/17 03:34:05 INFO TaskSchedulerImpl: Killing all running tasks in stage 9048: S
   28581
        2041
     47|
        2031
                                               24/07/17 03:34:05 INFO DAGScheduler: Job 935 finished: showString at NativeMethodAcc
    7801
                                                +----+
         2011
    150|
   11981
         2001
                                                |userId|movieId|
                                                                     rating
   49931 1981
                                                +----+
                                                            3379| 5.763239|
only showing top 20 rows
                                                           33649| 5.598928|
Num models to be tested: 16
                                                            5490|5.5296617|
CrossValidator 25a4a2ee8f51
24/07/17 03:36:25 INFO SparkContext: Invoking stop() from
                                                      1 | 171495 | 5.416649 |
24/07/17 03:36:25 INFO SparkContext: SparkContext is stop
                                                            5416|5.4002886|
24/07/17 03:36:25 INFO SparkUI: Stopped Spark web UI at h
24/07/17 03:36:25 INFO MapOutputTrackerMasterEndpoint: Ma
                                                            5328 | 5.4002886 |
                                                            3951 | 5.4002886 |
                                                      1 | 131724 | 5.363606 |
                                                            5915|5.3629932|
                                                       11 1775931 5.3565161
                                                +----+
                                               24/07/17 03:34:05 INFO FileSourceStrategy: Pushed Filters: IsNotNull (movieId)
                                               24/07/17 03:34:05 INFO FileSourceStrategy: Post-Scan Filters: isnotnull(movieId#17)
                                               24/07/17 03:34:05 INFO DAGScheduler: Registering RDD 19217 (showString at NativeMeth
                                               24/07/17 03:34:05 INFO DAGScheduler: Got map stage job 936 (showString at NativeMeth
                                               24/07/17 03:34:05 INFO DAGScheduler: Final stage: ShuffleMapStage 9072 (showString a
```

```
Z4/U//1/ U3:34:U3 INFO EXECUTOR: FINISHED LASK U.U IN STAGE 3030.U (IID ZZI30). 63/3 DYTES FESUIT SENT TO DILVER
24/07/17 03:34:09 INFO TaskSetManager: Finished task 0.0 in stage 9098.0 (TID 22190) in 371 ms on cs-2199409848-default (executo
24/07/17 03:34:09 INFO TaskSchedulerImpl: Removed TaskSet 9098.0, whose tasks have all completed, from pool
24/07/17 03:34:09 INFO DAGScheduler: ResultStage 9098 (showString at NativeMethodAccessorImpl.java:0) finished in 0.382 s
24/07/17 03:34:09 INFO DAGScheduler: Job 938 is finished. Cancelling potential speculative or zombie tasks for this job
24/07/17 03:34:09 INFO TaskSchedulerImpl: Killing all running tasks in stage 9098: Stage finished
24/07/17 03:34:09 INFO DAGScheduler: Job 938 finished: showString at NativeMethodAccessorImpl.java:0, took 0.390564 s
24/07/17 03:34:09 INFO CodeGenerator: Code generated in 7.276758 ms
  |movieId|userId| rating|
     67618 | 100 | 5.120143 | Strictly Sexual (... | Comedy | Drama | Romance |
       3379| 100| 5.064743| On the Beach (1959)|
     42730| 100| 5.042285| Glory Road (2006)|
     33649| 100|5.0216565| Saving Face (2004)|Comedy|Drama|Romance|
   184245| 100|4.9267745|De platte jungle ...|
                                                                                               Documentary
  179135| 100|4.9267745|Blue Planet II (2...|
                                                                                               Documentary|
   138966| 100|4.9267745|Nasu: Summer in A...|
                                                                                                  Animation
  117531| 100|4.9267745|
                                                  Watermark (2014) |
                                                                                               Documentary|
     86237| 100|4.9267745| Connections (1978)|
                                                                                               Documentary|
     84273| 100|4.9267745|Zeitgeist: Moving...|
                                                                                               Documentary|
  -----
24/07/17 03:34:09 INFO FileSourceStrategy: Pushed Filters: IsNotNull(userId)
24/07/17 03:34:09 INFO FileSourceStrategy: Post-Scan Filters: isnotnull(userId#40), (cast(userId#40 as int) = 100), isnotnull(cast
24/07/17 03:34:09 INFO FileSourceStrategy: Pushed Filters: IsNotNull(movieId)
24/07/17 03:34:09 INFO FileSourceStrategy: Post-Scan Filters: isnotnull(movieId#17)
24/07/17 03:34:09 INFO CodeGenerator: Code generated in 66.457324 ms
24/07/17 03:34:09 INFO MemoryStore: Block broadcast 2993 stored as values in memory (estimated size 199.3 KiB, free 407.4 MiB)
24/07/17 03:34:09 INFO MemoryStore: Block broadcast_2993_piece0 stored as bytes in memory (estimated size 34.3 KiB, free 407.3 24/07/17 03:34:09 INFO BlockManagerInfo: Added broadcast_2993_piece0 in memory24/07/17 03:34:10 INFO CodeGenerator: Code generated in 8.644011 ms
24/07/17 03:34:09 INFO SparkContext: Created broadcast 2993 from $anonfun$with24/07/17 03:34:10 INFO FileScanRDD: Reading File path: file:///home/faraya85431/ratings.cav, range: 0-2483723, partition values: [empty row]
24/07/17 03:34:09 INFO FileSourceScanExec: Planning scan with bin packing, max 24/07/17 03:34:10 INFO BlockManagerInfo: Removed broadcast 2991 place0 on cs-2199409848-default:37827 in memory (size: 55.6 KiB, free: 411.0 MiB) 24/07/17 03:34:10 INFO BlockManagerInfo: Removed broadcast 2983 place0 on cs-2199409848-default:37827 in memory (size: 45.6 KiB, free: 411.1 MiB) 24/07/17 03:34:10 INFO BlockManagerInfo: Removed broadcast 2983 place0 on cs-2199409848-default:37827 in memory (size: 34.3 KiB, free: 411.1 MiB) 34/07/17 03:34:10 INFO BlockManagerInfo: Removed broadcast 2983 place0 on cs-2199409848-default:37827 in memory (size: 34.3 KiB, free: 411.1 MiB) 34/07/17 03:34:10 INFO BlockManagerInfo: Removed broadcast 2983 place0 on cs-2199409848-default:37827 in memory (size: 35.6 KiB, free: 411.0 MiB) 34/07/17 03:34:10 INFO BlockManagerInfo: Removed broadcast 2983 place0 on cs-2199409848-default:37827 in memory (size: 35.6 KiB, free: 411.1 MiB) 34/07/17 03:34:10 INFO BlockManagerInfo: Removed broadcast 2983 place0 on cs-2199409848-default:37827 in memory (size: 35.6 KiB, free: 411.1 MiB) 34/07/17 03:34:10 INFO BlockManagerInfo: Removed broadcast 2983 place0 on cs-2199409848-default:37827 in memory (size: 35.6 KiB, free: 411.1 MiB) 34/07/17 03:34:10 INFO BlockManagerInfo: Removed broadcast 2983 place0 on cs-2199409848-default:37827 in memory (size: 35.6 KiB, free: 411.1 MiB) 34/07/17 03:34:10 INFO BlockManagerInfo: Removed broadcast 2983 place0 on cs-2199409848-default:37827 in memory (size: 35.6 KiB, free: 411.0 MiB) 34/07/17 03:34:10 INFO BlockManagerInfo: Removed broadcast 2983 place0 on cs-2199409848-default:37827 in memory (size: 35.6 KiB, free: 411.0 MiB) 34/07/17 03:34:10 INFO BlockManagerInfo: Removed broadcast 2983 place0 on cs-2199409848-default:37827 in memory (size: 35.6 KiB, free: 411.0 MiB) 34/07/17 03:34:10 INFO BlockManagerInfo: Removed broadcast 2983 place0 on cs-2199409848-default:37827 in memory (size: 35.6 KiB, free: 411.0 MiB) 34/07/17 03:34:10 INFO BlockManagerInfo: Removed broadcast 2983 place
                                                                                                                                     24/07/17 03:34:10 INFO BlockManagerInfo: Removed broadcast 2992 piece0 on cs-2199409848-default:37827 in memory (size: 49.9 KiB, free: 411.1 MiB)
                                                                                                                                     24/07/17 03:34:10 INFO BlockManagerInfo: Removed broadcast 2994 piece0 on cs-2199409848-default:37827 in memory (size: 7.6 KiB, free: 411.2 MiB) 24/07/17 03:34:10 INFO BlockManagerInfo: Removed broadcast 2990 piece0 on cs-2199409848-default:37827 in memory (size: 7.6 KiB, free: 411.2 MiB)
                                                                                                                                     24/07/17 03:34:10 INFO Executor: Finished task 0.0 in stage 9100.0 (TID 22192). 2827 bytes result sent to driver
                                                                                                                                     24/07/17 03:34:10 INFO TaskSetManager: Finished task 0.0 in stage 9100.0 (TID 22192) in 641 ms on cs-2199409848-default (executor driver) (1/1)
                                                                                                                                     24/07/17 03:34:10 INFO TaskSchedulerImpl: Removed TaskSet 9100.0, whose tasks have all completed, from pool 24/07/17 03:34:10 INFO DAGScheduler: ResultStage 9100 (showString at NativeMethodAccessorImpl.java:0) finished in 0.650 s
                                                                                                                                      4/07/17 03:34:10 INFO DAGScheduler: Job 940 is finished. Cancelling potential speculative or zombie tasks for this job
                                                                                                                                      4/07/17 03:34:10 INFO TaskSchedulerImpl: Killing all running tasks in stage 9100: Stage finished
                                                                                                                                      4/07/17 03:34:10 INFO DAGScheduler: Job 940 finished: showString at NativeMethodAccessorImpl.java:0, took 0.653751 s
                                                                                                                                     24/07/17 03:34:10 INFO CodeGenerator: Code generated in 13.564665 ms
                                                                                                                                      |movieId|userId|rating|
                                                                                                                                         1101| 100| 5.0| Top Gun (1986)|
                                                                                                                                                                                               Action | Romance
                                                                                                                                         1958| 100| 5.0|Terms of Endearme...|
                                                                                                                                                                                                 Comedy | Drama |
                                                                                                                                                          5.0|Christmas Vacatio...|
                                                                                                                                         4041| 100| 5.0|Officer and a Gen...|
                                                                                                                                                                                                Drama | Romance |
                                                                                                                                                 100| 5.0|Sweet Home Alabam...|
                                                                                                                                                  100| 4.5| Maverick (1994) | Adventure | Comedy | . . . |
                                                                                                                                                         4.5|Father of the Bri...|
                                                                                                                                                  100| 4.5|Sleepless in Seat...|Comedy|Drama|Romance|
                                                                                                                                                                       Casino (1995)|
                                                                                                                                                  100| 4.5| Tombstone (1993) |Action|Drama|Western|
                                                                                                                                      4/07/17 03:34:11 INFO SparkContext: Invoking stop() from shutdown hook
                                                                                                                                     24/07/17 03:34:11 INFO SparkContext: SparkContext is stopping with exitCode 0
                                                                                                                                     24/07/17 03:34:11 INFO SparkUI: Stopped Spark web UI at http://cs-2199409848-default:4040
```

24/07/17 03:34:11 INFO MapOutputTrackerMasterEndpoint: MapOutputTrackerMasterEndpoint stopped!

4/07/17 03:34:12 INFO OutputCommitCoordinator\$OutputCommitCoordinatorEndpoint: OutputCommitCoordinator stopped!

14/07/17 03:34:12 INFO MemoryStore: MemoryStore cleared 14/07/17 03:34:12 INFO BlockManager: BlockManager stopped 14/07/17 03:34:12 INFO BlockManagerMaster: BlockManagerMaster stopped

4/07/17 03:34:12 INFO SparkContext: Successfully stopped SparkContext





Can we get better result?

05



- Incorporate additional data (e.g., user demographics)
- Use hybrid recommendation models
- Experiment with different machine learning algorithms
- Implement real-time recommendation updates







# 06Conclusion







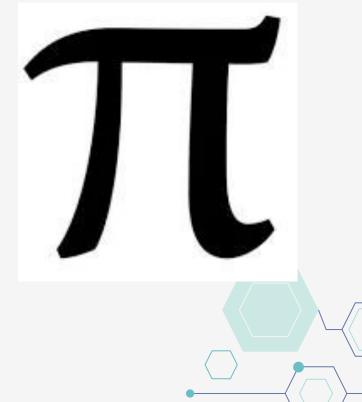
- Importance of data preprocessing
- Benefits of hyperparameter tuning
- Real-world applications of recommendation engines







# 07 References



### How to Build a Movie Recommendation System Based on Collaborative Filtering

## Movie Recommendation with Collaborative Filtering in Pyspark









## Thanks!



Please keep this slide for attribution



