

All result code and output can be seen either in attached jupyter notebook or at the end of this PDF.

Part A:

1. Count of odd and even numbers.

```
1 #reading in text file
2 sourceRDD_A1 = spark.sparkContext.textFile("/FileStore/tables/integer.txt", 4)
```

Command took 0.06 seconds -- by melina.fartaj@mail.utoronto.ca at 10/11/2021, 12:51:59 AM on Cluster_Assign2

Cmd 7

```
1 #counting total number of records in RDD
2 sourceRDD_A1.count()
```

► (1) Spark Jobs

Out[237]: 1010

Command took 0.04 seconds -- by melina.fartaj@mail.utoronto.ca at 10/11/2021, 12:51:59 AM on Cluster_Assign2

Cmd 8

```
1 #Transformation
2 #converting x to floats
3 float_RDD = sourceRDD_A1.map(lambda x: float(x))
```

Command took 0.04 seconds -- by melina.fartaj@mail.utoronto.ca at 10/11/2021, 12:51:59 AM on Cluster_Assign2

Cmd 9

```
1 #Transformation
2 #operations to filter for even and odd numbers
3 even_RDD = float_RDD.filter(lambda x: x%2 == 0)
4 odd_RDD = float_RDD.filter(lambda x: x%2 != 0)
```

Command took 0.04 seconds -- by melina.fartaj@mail.utoronto.ca at 10/11/2021, 12:51:59 AM on Cluster_Assign2

Cmd 10

```
1 #Action
2 #Number of even numbers
3 even_RDD.count()
```

► (1) Spark Jobs

Out[240]: 514

Command took 0.03 seconds -- by melina.fartaj@mail.utoronto.ca at 10/11/2021, 12:51:59 AM on Cluster_Assign2

Cmd 11

```
1 #Action
2 #Number of odd numbers
3 odd_RDD.count()
```

► (1) Spark Jobs

Out[241]: 496

2. Salary sum for each department.

```
1 #reading in text file
2 sourceRDD_A2 = spark.sparkContext.textFile("/FileStore/tables/salary.txt", 4)
3 #viewing source rdd
4 sourceRDD_A2.collect()
```

► (1) Spark Jobs

```
Out[6]: ['Sales 9136',
'Research 13391',
'Developer 22220',
'QA 31888',
'Marketing 22215',
```

Command took 0.89 seconds -- by melina.fartaj@mail.utoronto.ca at 10/11/2021, 2:25:43 PM on Cluster_Assing2_1

Cmd 10

```
1 #creating a list of arrays for each row
2 arrayRDD_A2 = sourceRDD_A2.map(lambda x: x.split(" "))
3 arrayRDD_A2.collect()
```

► (1) Spark Jobs

```
Out[7]: [['Sales', '9136'],
['Research', '13391'],
['Developer', '22220'],
['QA', '31888'],
['Marketing', '22215'],
```

Command took 0.66 seconds -- by melina.fartaj@mail.utoronto.ca at 10/11/2021, 2:25:59 PM on Cluster_Assing2_1

Cmd 11

```
1 #converting the 2nd element (position 1) into float
2 kvRDD_A2 = arrayRDD_A2.map(lambda x: (x[0], float(x[1])))
3 kvRDD_A2.collect()
```

► (1) Spark Jobs

```
Out[8]: [('Sales', 9136.0),
('Research', 13391.0),
('Developer', 22220.0),
('QA', 31888.0),
('Marketing', 22215.0),
```

Command took 0.72 seconds -- by melina.fartaj@mail.utoronto.ca at 10/11/2021, 2:26:49 PM on Cluster_Assing2_1

Cmd 12

```
1 #calculating the salary sum for each department
2 sumRDD_A2 = kvRDD_A2.reduceByKey(lambda x,y: x+y)
3 #Salary sum for each department
4 sumRDD_A2.collect()
```

► (1) Spark Jobs

```
Out[9]: [('Developer', 3221394.0),
('Sales', 3488491.0),
('Research', 3328284.0),
('QA', 3368624.0),
('Marketing', 3158450.0)]
```

Command took 1.76 seconds -- by melina.fartaj@mail.utoronto.ca at 10/11/2021, 2:27:06 PM on Cluster_Assing2_1

3. MapReduce using pyspark.

```
1 #MapReduce -> Mapping step
2 mapper_RDD_A3 = input_RDD_A3.map(lambda word: (word, 1))
3 mapper_RDD_A3.collect()
```

► (1) Spark Jobs

```
Out[14]: [('This', 1),
 ('eBook', 1),
 ('is', 1),
 ('for', 1),
 ('the', 1),
```

Command took 2.75 seconds -- by melina.fartaj@mail.utoronto.ca at 10/11/2021, 2:31:23 PM on Cluster_Assing2_1

Cnd 19

```
1 #MapReduce -> Reducer step
2 reducer_RDD_A3 = mapper_RDD_A3.reduceByKey(lambda x, y: x + y)
3 reducer_RDD_A3.collect()
```

► (1) Spark Jobs

```
Out[15]: [('of', 7957),
 ('no', 1270),
 ('whatsoever', 8),
 ('may', 642),
 ('give', 473),
```

Command took 3.06 seconds -- by melina.fartaj@mail.utoronto.ca at 10/11/2021, 2:32:02 PM on Cluster_Assing2_1

```
1 #convert to dataframe and cache
2 wordcounts = spark.createDataFrame(reducer_RDD_A3).cache()
3 #changing column headers of wordcount dataframe
4 wordcounts = wordcounts.select(col("_1").alias("word"), col("_2").alias("count"))
5 wordcounts.show()
```

► (2) Spark Jobs

► wordcounts: pyspark.sql.dataframe.DataFrame = [word: string, count: long]

```
+-----+-----+
|      word|count|
+-----+-----+
|      of| 7957|
|      no| 1270|
```

Command took 2.92 seconds -- by melina.fartaj@mail.utoronto.ca at 10/11/2021, 2:32:20 PM on Cluster_Assing2_1

Cnd 21

```
1 #filtering the dataframe for the EXACT words in word_list
2 word_list = ['Shakespeare', 'When', 'Lord', 'Library', 'GUTENBERG', 'WILLIAM', 'COLLEGE', 'WORLD']
3 filtered_words = wordcounts.filter(func.col('word').rlike('^(\s|' + '|'.join(word_list) + ')(\s|$)'))
4 #result
5 filtered_words.show()
```

► (2) Spark Jobs

► filtered_words: pyspark.sql.dataframe.DataFrame = [word: string, count: long]

```
+-----+-----+
|      word|count|
+-----+-----+
|Shakespeare|   22|
|  WILLIAM|  127|
|   WORLD|   98|
|   When|  405|
|GUTENBERG|   99|
|  COLLEGE|   98|
|   Lord|  400|
|  Library|    5|
+-----+-----+
```

4. Word Count.

Top 15	Bottom 15																																																																
<pre>1 #Top 15 words with the most count 2 wordcounts.orderBy('count', ascending=False).show(15)</pre> <p>▶ (1) Spark Jobs</p> <table><thead><tr><th>word</th><th>count</th></tr></thead><tbody><tr><td>the</td><td>11402</td></tr><tr><td>I</td><td>8936</td></tr><tr><td>and</td><td>8916</td></tr><tr><td>of</td><td>7957</td></tr><tr><td>to</td><td>7508</td></tr><tr><td>a</td><td>5755</td></tr><tr><td>you</td><td>5306</td></tr><tr><td>my</td><td>4919</td></tr><tr><td>in</td><td>4722</td></tr><tr><td>that</td><td>3754</td></tr><tr><td>And</td><td>3735</td></tr><tr><td>is</td><td>3666</td></tr><tr><td>not</td><td>3588</td></tr><tr><td>me</td><td>3445</td></tr><tr><td>his</td><td>3278</td></tr></tbody></table> <p>only showing top 15 rows</p>	word	count	the	11402	I	8936	and	8916	of	7957	to	7508	a	5755	you	5306	my	4919	in	4722	that	3754	And	3735	is	3666	not	3588	me	3445	his	3278	<pre>1 #Bottom 15 words with the least count 2 wordcounts.orderBy('count', ascending=True).show(15)</pre> <p>▶ (1) Spark Jobs</p> <table><thead><tr><th>word</th><th>count</th></tr></thead><tbody><tr><td>provisions</td><td>1</td></tr><tr><td>INCIDENTAL</td><td>1</td></tr><tr><td>processing</td><td>1</td></tr><tr><td>NOR</td><td>1</td></tr><tr><td>EITHER</td><td>1</td></tr><tr><td>Release</td><td>1</td></tr><tr><td>tilde</td><td>1</td></tr><tr><td>CDROMS</td><td>1</td></tr><tr><td>reader</td><td>1</td></tr><tr><td>Title</td><td>1</td></tr><tr><td>EBCDIC</td><td>1</td></tr><tr><td>EBOOK</td><td>1</td></tr><tr><td>equivalent</td><td>1</td></tr><tr><td>AWAY</td><td>1</td></tr><tr><td>program</td><td>1</td></tr></tbody></table> <p>only showing top 15 rows</p>	word	count	provisions	1	INCIDENTAL	1	processing	1	NOR	1	EITHER	1	Release	1	tilde	1	CDROMS	1	reader	1	Title	1	EBCDIC	1	EBOOK	1	equivalent	1	AWAY	1	program	1
word	count																																																																
the	11402																																																																
I	8936																																																																
and	8916																																																																
of	7957																																																																
to	7508																																																																
a	5755																																																																
you	5306																																																																
my	4919																																																																
in	4722																																																																
that	3754																																																																
And	3735																																																																
is	3666																																																																
not	3588																																																																
me	3445																																																																
his	3278																																																																
word	count																																																																
provisions	1																																																																
INCIDENTAL	1																																																																
processing	1																																																																
NOR	1																																																																
EITHER	1																																																																
Release	1																																																																
tilde	1																																																																
CDROMS	1																																																																
reader	1																																																																
Title	1																																																																
EBCDIC	1																																																																
EBOOK	1																																																																
equivalent	1																																																																
AWAY	1																																																																
program	1																																																																

Part B:

1. Exploratory Data Analysis.

a. Null Values

```
1 #checking for any null values
2 df.select([count(when(col(c).isNull(), c)).alias(c) for c in df.columns]).show()
```

▶ (2) Spark Jobs

movieId	rating	userId
0	0	0

b. Describing data (counts, averages, min/max etc..)

```
1 df.describe().show()
```

▶ (2) Spark Jobs

summary	movieId	rating	userId
count	1501	1501	1501
mean	49.40572951365756	1.7741505662891406	14.383744170552964
stddev	28.937034065088994	1.187276166124803	8.591040424293272
min	0	1	0
max	99	5	29

c. Distinct number of User Ids

```
1 df.select("userId").distinct().count()
```

▶ (3) Spark Jobs

Out[37]: 30

d. Distinct number of Movie Ids

```
1 df.select("movieId").distinct().count()
```

▶ (3) Spark Jobs

Out[38]: 100

e. Describing how many movies users have reviewed (average, min/max, etc..)

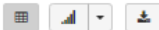
```
1 #calculating the number of times users reviewed movies
2 userrating_counts = df.groupBy("userId").count()
3 display(userrating_counts)
```

▶ (2) Spark Jobs

userrating_counts: pyspark.sql.dataframe.DataFrame = [userId: integer, count: long]

	userId	count
1	12	55
2	13	48
3	14	57
4	18	52
5	25	46
6	6	57
7	9	53

Showing all 30 rows.



```
1 #viewing how many movies ppl have reviewed
2 userrating_counts.describe(['count']).show()
```

▶ (3) Spark Jobs

```
+-----+-----+
|summary|count|
+-----+-----+
|count|30|
|mean|50.03333333333333|
|stddev|3.8460309441017904|
|min|44|
|max|57|
+-----+-----+
```

f. Top 12 users and ratings

Top 12 highest ratings	Top 12 users
<pre>1 #calculating top 12 movies with highest ratings 2 df.orderBy('rating', ascending=False).show(12)</pre> <p>▶ (1) Spark Jobs</p> <pre>+-----+-----+-----+ movieId rating userId +-----+-----+-----+ 55 5 5 30 5 11 25 5 6 8 5 2 25 5 7 39 5 2 29 5 8 51 5 3 52 5 8 7 5 9 27 5 11 49 5 9 +-----+-----+-----+ only showing top 12 rows</pre>	<pre>1 #calculating top 12 highest number of ratings done by users 2 userrating_counts.orderBy('count', ascending=False).show(12)</pre> <p>▶ (2) Spark Jobs</p> <pre>+-----+-----+ userId count +-----+-----+ 6 57 14 57 11 56 22 56 4 55 12 55 7 54 9 53 23 52 24 52 18 52 28 50 +-----+-----+ only showing top 12 rows</pre>

2. Train-test Split

a. Model 1:

```
1 #MODEL 1
2 #first combination split - 80/20
3 (train1, test1) = df.randomSplit([0.80, 0.20])
4 # Train model 1 with Training Data 1
5 als_model_1 = als.fit(train1)
6 #Model 1 predictions on test1
7 pred_1 = als_model_1.transform(test1)
```

► (5) Spark Jobs

- train1: pyspark.sql.dataframe.DataFrame = [movieId: integer, rating: integer ... 1 more field]
- test1: pyspark.sql.dataframe.DataFrame = [movieId: integer, rating: integer ... 1 more field]
- pred_1: pyspark.sql.dataframe.DataFrame = [movieId: integer, rating: integer ... 2 more fields]

b. Model 2

```
1 #MODEL 2
2 #second combination split - 60/40
3 (train2, test2) = df.randomSplit([0.6, 0.4])
4 # Train model 2 with Training Data 2
5 als_model_2 = als.fit(train2)
6 #Model 2 predictions on test1
7 pred_2 = als_model_2.transform(test2)
```

► (5) Spark Jobs

- train2: pyspark.sql.dataframe.DataFrame = [movieId: integer, rating: integer ... 1 more field]
- test2: pyspark.sql.dataframe.DataFrame = [movieId: integer, rating: integer ... 1 more field]
- pred_2: pyspark.sql.dataframe.DataFrame = [movieId: integer, rating: integer ... 2 more fields]

c. Initial ALS Model

```
1 #Create initial ALS model
2 als = ALS(userCol="userId", itemCol="movieId", ratingCol="rating", coldStartStrategy="drop")
```

3. Report Models Performance

a. MSE, RMSE, and MAE.

MSE stands for Mean Squared Error and represents the average of the squared difference between the actual and predicted values. RMSE stands for Root Mean Squared Error and measures the standard deviation of residuals by taking the square root of Mean Squared Error. MAE stands for Mean Absolute Error and represents the average of the absolute difference between the actual and predicted values. The smaller the value of MSE, RMSE and MAE, the better. A smaller value indicates better accuracy.

```
1 #EVALUATORS
2 #Evaluation with RMSE
3 evaluator_RMSE = RegressionEvaluator(metricName="rmse", labelCol="rating", predictionCol="prediction")
4 #Evaluation with MSE
5 evaluator_MSE = RegressionEvaluator(metricName="mse", labelCol="rating", predictionCol="prediction")
6 #Evaluation with MAE
7 evaluator_MAE = RegressionEvaluator(metricName="mae", labelCol="rating", predictionCol="prediction")
```

b. Which one of the two splits worked better?

Model 1 works better as it has a lower RMSE, MSE and MAE than Model 2 indicating higher accuracy. Model 1 works better because of the higher training split. Model 1 trains the model on 80% of the data while Model 2 trains the model on only 60% of the data. However, larger test datasets provide a more accurate representation of the model's performance.

```
1 #MODEL 1 EVALUATION
2 #RMSE Model 1
3 rmse_1 = evaluator_RMSE.evaluate(pred_1)
4 #MSE Model 1
5 mse_1 = evaluator_MSE.evaluate(pred_1)
6 #MAE Model 1
7 mae_1 = evaluator_MAE.evaluate(pred_1)
```

► (15) Spark Jobs

```
1 #MODEL 2 EVALUATION
2 #RMSE Model 2
3 rmse_2 = evaluator_RMSE.evaluate(pred_2)
4 #MSE Model 2
5 mse_2 = evaluator_MSE.evaluate(pred_2)
6 #MAE Model 2
7 mae_2 = evaluator_MAE.evaluate(pred_2)
```

► (15) Spark Jobs

```
1 print ("RMSE of model 1: ", rmse_1, "\nRMSE of model 2: ", rmse_2)
2 print ("MSE of model 1: ", mse_1, "\nMSE of model 2: ", mse_2)
3 print ("MAE of model 1: ", mae_1, "\nMAE of model 2: ", mae_2)
```

```
RMSE of model 1:  1.0223004355696874
RMSE of model 2:  1.1692273433760312
MSE of model 1:  1.0450981805659725
MSE of model 2:  1.3670925804981717
MAE of model 1:  0.6846297040877752
MAE of model 2:  0.8261748726854657
```

4. Model Tuning

a. Combination selected:

Model 1, 80/20 split (however model 2 was also tuned and run as per assignment instructions given during lecture)

b. What are the hyperparameters of ALS?

The hyperparameters of ALS are rank, maxIter, regParam, numBlocks, and alpha.

c. Which hyperparameters did you tune?

I chose to tune the hyperparameters rank and regParam. Rank is the number of latent factors in the model and defaults to 10. This hyperparameter helps to group/categorize how similar and different movies are. I decided to test out the following different rank values: 10, 50, 100 and 200. regParam specifies the regularization parameter and defaults to 1.0. I decided to test out the following different values: 0.1, 0.15, 0.2, 0.25.

```

1 # Create ParamGrid for Cross Validation
2 parameters = (ParamGridBuilder()
3               .addGrid(als.rank, [10, 50, 100, 200])
4               .addGrid(als.regParam, [.1, .15, 0.2, 0.25])
5               .build())
6
7 #Build train validation split
8 tranvs = TrainValidationSplit(estimator=als, estimatorParamMaps=parameters, evaluator=evaluator_RMSE)

```

```

1 #BEST MODEL 1
2 #Run train validation split
3 tuned_model_1 = tranvs.fit(train1)
4 #Get best Model
5 best_model = tuned_model_1.bestModel

```

► (9) Spark Jobs

d. Best Model and train-test error.

The Best Model uses a rank of 200 and a regParam of 0.15 to get a RMSE of 0.994.

```

1 #Using best model to predict
2 tuned_pred_1 = best_model.transform(test1)
3 tuned_pred_2 = best_model2.transform(test2)

```

►  tuned_pred_1: pyspark.sql.dataframe.DataFrame = [movieId: integer, rating: integer ... 2 more fields]
 ►  tuned_pred_2: pyspark.sql.dataframe.DataFrame = [movieId: integer, rating: integer ... 2 more fields]

Command took 0.11 seconds -- by melina.fartaj@mail.utoronto.ca at 10/11/2021, 12:51:59 AM on Cluster_Assign2

Cmd 66

```

1 #RMSE of Best Model 1
2 tuned_rmse_1 = evaluator_RMSE.evaluate(tuned_pred_1)
3 #RMSE of Best Model 2
4 tuned_rmse_2 = evaluator_RMSE.evaluate(tuned_pred_2)

```

► (10) Spark Jobs

Command took 7.37 seconds -- by melina.fartaj@mail.utoronto.ca at 10/11/2021, 12:51:59 AM on Cluster_Assign2

Cmd 67

```

1 print ("RMSE of new model 1: ", tuned_rmse_1)
2 print ("RMSE of new model 2: ", tuned_rmse_2)

```

RMSE of new model 1: 0.9944893950996873
 RMSE of new model 2: 1.1479176736055905

Command took 0.03 seconds -- by melina.fartaj@mail.utoronto.ca at 10/11/2021, 12:51:59 AM on Cluster_Assign2

5. Recommended Movies

The Best Model (best_model) was used for this question.

```

1 # Generate Recommendations for all users (from best to worst)
2 movie_recommendations = best_model.recommendForAllUsers(100)
3 movie_recommendations.show()

```

► (2) Spark Jobs

►  movie_recommendations: pyspark.sql.dataframe.DataFrame = [userId: integer, recommendations: array]

```

+-----+-----+
|userId|recommendations|
+-----+-----+
| 20|[22, 3.9048154],...|
| 10|[92, 3.6125093],...|
| 0|[92, 3.11709], [...]|
| 1|[62, 3.160814],...|

```


a. Top 12 movies recommendations for User Id 1

```
1 #filtering for only User Id 1
2 movie_recommendations_ID1 = movie_recommendations.filter(movie_recommendations.userId == 1)
3 movie_recommendations_ID1.show()
```

▶ (2) Spark Jobs

▶ movie_recommendations_ID1: pyspark.sql.dataframe.DataFrame = [userid: integer, recommendations: array]

```
-----+-----
|userId| recommendations|
-----+-----
| 1| [[62, 3.1668814], ...]|
-----+-----
```

Command took 17.41 seconds -- by melina.fartaj@mail.utoronto.ca at 10/11/2021, 1:42:46 AM on Cluster_Assign2

Cmd 57

```
1 #User Id 1
2 #Creating a list of all the movies that user Id 1 reviewed
3 userID1movies = df.filter(col("userId") == 1).select('movieId').rdd.flatMap(lambda x: x).collect()
4
5 #Using explode to format the array into columns
6 #filtering out the movies that the user has already reviewed and showing top 12
7 recommendations_ID1 = movie_recommendations_ID1.withColumn("recommendations", explode("recommendations")).select('userId', col("recommendations.movieId"), col("recommendations.rating"))
8 recommendations_ID1.filter(~recommendations_ID1.movieId.isin(userID1movies)).show(12)
```

▶ (2) Spark Jobs

▶ recommendations_ID1: pyspark.sql.dataframe.DataFrame = [userid: integer, movieId: integer ... 1 more field]

```
-----+-----+-----
|userId|movieId| rating|
-----+-----+-----
| 1| 22| 2.38624|
| 1| 7| 2.273389|
| 1| 51| 2.232317|
| 1| 95| 2.0818558|
| 1| 31| 2.0718462|
| 1| 80| 2.0705826|
| 1| 15| 2.0507987|
| 1| 19| 1.9632709|
| 1| 98| 1.9516469|
| 1| 23| 1.9421765|
| 1| 24| 1.8715812|
| 1| 32| 1.8000154|
-----+-----+-----
```

only showing top 12 rows

b. Top 12 movie recommendations for User Id 12

```
1 #filtering for only User Id 12
2 movie_recommendations_ID12 = movie_recommendations.filter(movie_recommendations.userId == 12)
3 movie_recommendations_ID12.show()
```

▶ (2) Spark Jobs

▶ movie_recommendations_ID12: pyspark.sql.dataframe.DataFrame = [userid: integer, recommendations: array]

```
-----+-----
|userId| recommendations|
-----+-----
| 12| [[55, 4.326746], ...]|
-----+-----
```

Command took 17.83 seconds -- by melina.fartaj@mail.utoronto.ca at 10/11/2021, 1:42:46 AM on Cluster_Assign2

Cmd 59

```
1 #User Id 12
2 #Creating a list of all the movies that user Id 12 reviewed
3 userID12movies = df.filter(col("userId") == 12).select('movieId').rdd.flatMap(lambda x: x).collect()
4
5 #Using explode to format the array into columns
6 #filtering out the movies that the user has already reviewed and showing top 12
7 recommendations_ID12 = movie_recommendations_ID12.withColumn("recommendations", explode("recommendations")).select('userId', col("recommendations.movieId"), col("recommendations.rating"))
8 recommendations_ID12.filter(~recommendations_ID12.movieId.isin(userID12movies)).show(12)
```

▶ (2) Spark Jobs

▶ recommendations_ID12: pyspark.sql.dataframe.DataFrame = [userid: integer, movieId: integer ... 1 more field]

```
-----+-----+-----
|userId|movieId| rating|
-----+-----+-----
| 12| 55| 4.326746|
| 12| 46| 3.583644|
| 12| 49| 3.520106|
| 12| 65| 3.4627342|
| 12| 48| 3.369004|
| 12| 20| 3.341531|
| 12| 32| 3.2379308|
| 12| 90| 3.1749177|
| 12| 10| 2.9882762|
| 12| 68| 2.9059153|
| 12| 1| 2.8913836|
| 12| 28| 2.647625|
-----+-----+-----
```

only showing top 12 rows