# Experiment-6

Q Delving deeper into Spark API.

1. Load the rating.txt file into Databricks cloud.

2. Load the movies.txt file into Databricks cloud.

3. create an RDD from rating.txt file using textfile().
→ rating = sc.textfile("/FileStore/tables/rating.txt")

4. Split the RDD based on ':::' as a delimeter.
→ rating_split = rating.map(lambda line: line.split(":::"))

rating_split.take(2).

5. create a pair RDD with key being movieID and value being its rating.
→ pair_rating = rating_split.map(lambda x: (int(x[0]), int(x[2])))

pair_rating.collect().

6. Group the rating by key and display the output using take().
→ rating_group = pair_rating.groupByKey()
rating_group.take(5).

7 calculate the average rating for all movies
using mapValues().

→ rating_average = rating_group. mapValues ( lambda
x: sum(x)/ len(x))

8 Display 10 elements using take()

→ rating_average. take(10).

9 create a Pair RDD with rating being the key and
movie ID being its value and count the number
of movies for a certain rating.

→ rating_id = rating_split. map( lambda a: (int(x[2]),
int(x[0]))), groupByKey()

rating_rdd = rating_id. mapValues ( lambda v: len(v))

10 Review the output.

→ rating_rdd. collect()

11 create a RDD from movies.txt file using textFile().

→ movies = sc. textFile ("FileStore/tables/movie.txt")

movies. take (5).

12 split the RDD based upon ":" as delimeter.

→ movies_split = movies. map( lambda line: line.split(":")

movies. split. take (5).

13. create a PairRDD with movie IP being the key and movie name being its value and display the output using take()

→ movie_pair = movie_split.map( lambda x: (str(x[0]), str(x[1])).

movie_pair. take (5)

14. Now join the movies.txt in step 13 and ratings-rdd in steps to find the average rating of movies along with their names.

→ rdd_join = rating_average. join (movie_pair)

15. Display 10 elements and using take()

→ rdd_join.take(5).

16. Save the RDD as text file using save as TextFile()

→ rdd_join. saveAsTextFile ("result. txt")

rdd = sc. textFile ("result. txt")

rdd. take(5).

Outputs:

out[1]

4 [['2', '1193', '5', '978390760'],
    ['2', '662', '3', '678302109']]

out[2]

5 [(2,5),
    (1,8),
    (1,3),
    (2,4),
    - - -
    - - -
    (10,5),
    - - . ]

6 Out[3]:
[('1', <pyspark, resultiterable. ResultIterable at 0x72676203)
('4', <pyspark, resultiterable. ResultIterable at 0x7f2c765204),
('8', <pyspark, resultiterable. ResultIterable at 0x7f2c762053),
('3', <pyspark, resultiterable. ResultIterable at 0x7f2c7206>),
('10', <pyspark, resultiterable. ResultIterable at 0x7f2c7202>)]

```
8 Out[7]: [('2', 4.18),
          ('4', 4.19),
          ('8', 3.88),
          ('9', 3.77),
          ('10', 4.11),
          ('12', 3.92),
          ('14', 3.13),
          ('16', 3.02),
          ('17', 4.03),
          ('19', 3.57)]
```

```
10 Out[0]: [(4, 227696), (2, 720174), (5, 147187), (3, 17512),
           (1, 37523)]
```

```
11 Out[0]: [['1', 'Toy Story (1995)', 'Adventure|children|comedy'],
           ['2', 'Jumanji (1995)', 'Adventure|children|comedy'],
           ['3', 'Greetenpais (1995)', 'comedy|Romance'],
           ['4', 'Waiting to Extale (1995)', 'comedy|Drama'],
           ['5', 'Father of the Bride (1995)', 'comedy']]
```

13. out[13]: [[('1', 'Toy Story (1995)'),
('2', 'Jumanji (1995)',
'('3', 'Greenptan (1995)'
('4', 'waiting to Exhale (1995)'
('5', 'father of Bride (1995)')]

15. out[14]: [('4', (4.19, 'waiting to exhale (1995)')),
('10', (4.11, 'GoldenEye (1995)')),
('12', (3.82, 'Dracula: Dead and loving it (1995)')
('16', (3.02, 'Casino (1995)')),
('20', (4.08, 'Money Train (1995))',
('24', (3.94, 'Birder (1995)')),
('26', (2.93, 'Othello (1995)')),
('40', (3.44, 'Guy, the Beyond country (1995)')),
('44', (3.63, 'Mortal Kombat (1995)'))
('50', (3.06, 'usual suspects, the (1995)))]