

Experiment-4

Q. Transformations

1 Remove data file from Databricks file system.

→ `dbutils.fs.rm("FileStore/tables/data.txt")`

2 FlatMap

→ `data = sc.textFile("FileStore/tables/mydata.txt")`
`rdd = data.flatMap(lambda x: x.split(" "))`
`rdd.collect()`

3 Filter

→ `ls = [1, 2, 3, 4, 5]`

`data = sc.parallelize(ls)`

`f = data.filter(lambda x: x != 5)`

`f.collect()`

4 Distinct

→ `data = sc.textFile("FileStore/tables/mydata.txt")`

`rdd = data.flatMap(lambda x: x.split(" ")).distinct()`

`rdd.collect()`

5 MapValues

→ `rdd = sc.parallelize([(1, 2), (2, 4), (3, 5), (4, 7), (5, 8),
(6, 23), (7, 45)])`

`rdd.mapValues(lambda x: x + 1).collect()`

6 flatMapValues

```
→ data = sc.textFile("FileStore/tables/user_address.txt")  
a = data.map(lambda x: x.split("\t"))  
b = a.map(lambda x: (x[0], x[1]))  
c = b.flatMapValues(lambda x: x.split(";"))  
c.collect()
```

7 Key

```
→ data = sc.textFile("FileStore/tables/spark.txt")  
a = data.flatMap(lambda x: x.split(" "))  
b = a.map(lambda x: (x, 1)).keyBy()  
b.collect()
```

8 Values

```
→ data = sc.textFile("FileStore/tables/spark.txt")  
a = data.flatMap(lambda x: x.split(" "))  
b = a.map(lambda x: (x, 1)).valuesBy()  
b.collect()
```

9 sortByKey

```
→ data = sc.textFile("FileStore/tables/spark.txt")  
a = data.flatMap(lambda x: x.split(" "))  
b = a.map(lambda x: (x, 1)).sortByKey()  
b.collect()
```


10 groupByKey

```

→ data = sc.textFile("FileStore/tables/spark.txt")
a = data.flatMap(lambda x: x.split(" "))
b = a.map(lambda x: (x, 1)).groupByKey()
b.collect()

```

11 reduceByKey

```

→ data = sc.textFile("FileStore/tables/spark.txt")
a = data.flatMap(lambda x: x.split(" "))
b = a.map(lambda x: (x, 1)).reduceByKey(lambda v1, v2:
                                         v1 + v2)
b.collect()

```

12 Joins on RDDs (Inner joins between two RDDs)

```

→ rdd = sc.textFile("FileStore/tables/spark.txt").flatMap(
    lambda x: x.split(" ")).map(lambda
    x: (x, 1))
rdd1 = sc.textFile("FileStore/tables/hadoop.txt").flatMap(
    lambda x: x.split(" ")).map(lambda x: (x, 1))
rdd2 = rdd.join(rdd1)
rdd2.collect()

```

13 Joins on files

```

→ rdd = sc.textFile("FileStore/tables/log.txt").map(lambda
    l: l.split(" ")).map(lambda k: (k[5], k[0]))
rdd1 = sc.textFile("FileStore/tables/userdetails.csv").map(
    lambda l: l.split(",")).map(lambda k: (k[0], k[1]))
rdd2 = rdd.join(rdd1)
rdd2.collect()

```

Date - 01/02/2021

Outputs:

1. Remove data file

Out[1]: true

2. FlatMap.

Out[2]: ['1', '2', '3', ..., 'a', 'b', 'c', 'd', 'd', 'b', 's', 'd', 't']

3. Filter.

Out[3]: [1, 2, 3, 4]

4. Distinct.

Out[4]: ['1', '2', '3', ..., 'a', 'b', 's', 'd', 't', 'w', 't']

5. MapValues

Out[5]: [(4, 3), (2, 5), (3, 6), (4, 8), (5, 9), (6, 24), (7, 46)]

6. FlatMapValues

Out[6]: [('usr001', 'PMB101'),
('usr001', '228, Park ave'),
('usr001', 'New York'),
...,
('usr089', 'New York'),
('usr089', 'NY')]

Teacher's Signature

11 reduceByKey

11 spark jobs

```
Out[11]: [('spark', 1),  
          ('run', 1),  
          ('programs', 2),  
          ('in', 1),  
          ('easy', 1),  
          ...  
          ('write', 1),  
          ('applications', 1)]
```

12 inner join b/w two RDDs

```
Out[12]: [('to', (1, 1)), ('and', (1, 1)), ('data', (1, 1))]
```

13 join on file

```
Out[13]: [('7352'), ('64.242.88.11', 'Steven'),  
          ('6291'), ('64.242.88.10', 'Mark')]
```