# Ritika Kumari(A20414073)

## CSP554—Big Data Technologies

## Assignment #7

Exercise 1)

<u>Step A</u>

Use the TestDataGen program from previous assignments to generate new data files

Copy the files to HDFS.

**Command Executed:**

**java TestDataGen**

**hdfs dfs -copyFromLocal foodratings53475.txt /user/maria_dev/foodratings53475.csv**

**hdfs dfs -copyFromLocal foodplaces53475.txt hdfs:///user/maria_dev/foodplaces53475.csv**

**Magic Number = 53475**

## Step B

Load the 'foodratings' file as a 'csv' file into a DataFrame called ex1_foodratings. When doing so specify a schema having fields of the following names and types:

| Field Name | Field Type |
| --- | --- |
| name | String |
| food1 | Integer |
| food1 | Integer |
| food1 | Integer |
| food1 | Integer |
| placeid | Integer |

As the results of this exercise provide the magic number, the code you execute and screen shots of the following commands:

foodratings.printSchema()

foodratings.head(5)

**vi assign_7_q1.py**

from pyspark.sql.types import *

struct1 = StructType(

    [

        StructField("name", StringType(), True),

        StructField("food1",IntegerType(), True),

        StructField("food2",IntegerType(), True),

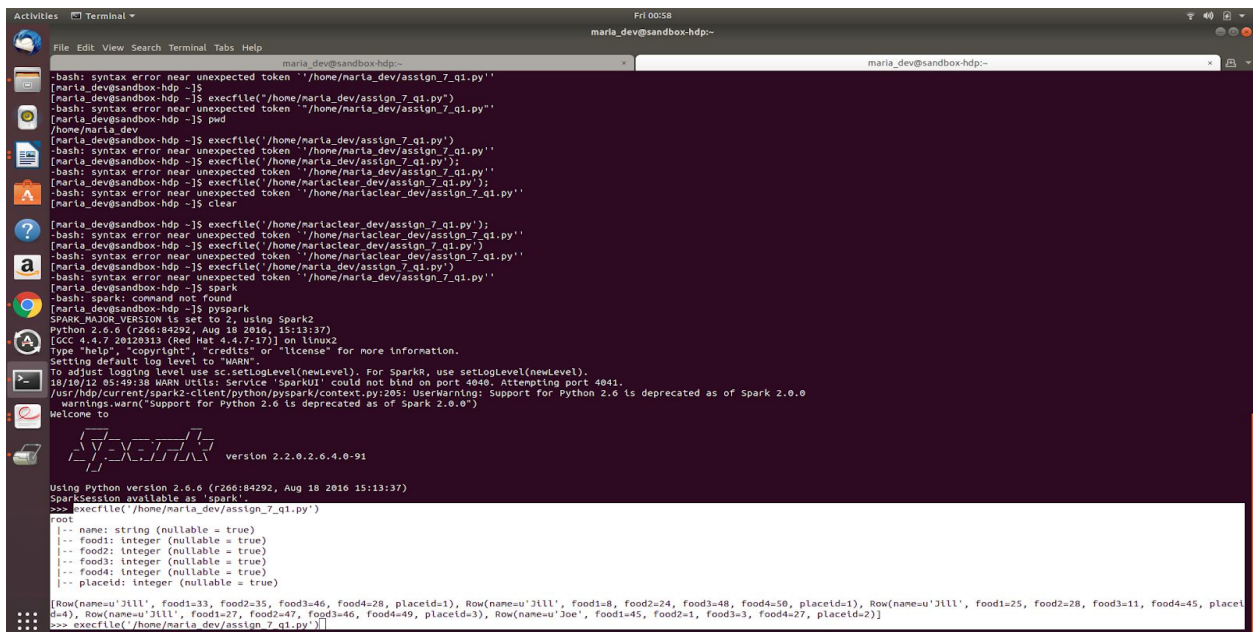        StructField("food3",IntegerType(), True),

        StructField("food4",IntegerType(), True),

        StructField("placeid",IntegerType(), True)

    ]

)

**execfile('/home/maria_dev/assign_7_q1.py')**

Load the 'foodplaces' file as a 'csv' file into a DataFrame called foodplaces. When doing so specify a schema having fields of the following names and types:

| Field Nampee | Field Ty |
|---|---|
| placeid | integer |
| placename | string |

As the results of this exercise provide the code you execute and screen shots of the following commands:

foodratings.printSchema()

foodratings.head(5)

**Command Executed:**
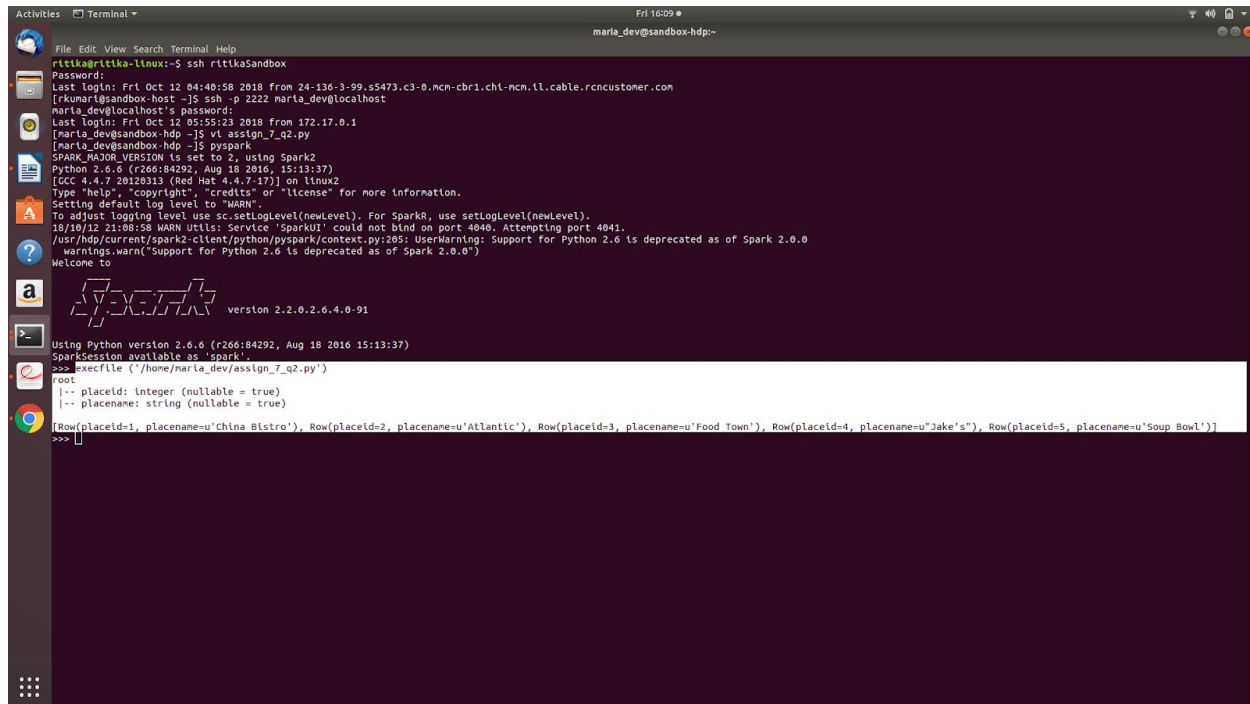
**vi assign_7_q2.py**

**from pyspark.sql.types import \***

**struct1 = StructType().add("placeid", IntegerType(), True).add("placename",StringType(), True)**

**foodplaces = spark.read.schema(struct1).csv('/user/maria_dev/foodplaces53475.csv')**

**foodplaces.printSchema()**

**print foodplaces.head(5)**

**execfile ('/home/maria_dev/assign_7_q2.py')**

Exercise 3)

Step A

Register the DataFrames created in exercise 1 and 2 as tables called "foodratingsT" and "foodplacesT"

Step B

Use a SQL query on the table "foodratingsT" to create a new DataFrame called foodratings_ex3 holding records which meet the following condition: food2 < 25 and food4 > 40

As the results of this step provide the code you execute and screen shots of the following commands:

foodratings.printSchema()

foodratings.head(5)

Step C

Use a SQL query on the table "foodplacesT" to create a new DataFrame called foodplaces_ex3 holding records which meet the following condition: placeid > 3

As the results of this step provide the code you execute and screen shots of the following commands:

```
        foodratings.printSchema()

        foodratings.head(5)
```

**vi assign_7_q3.py**

```python
from pyspark.sql.types import *

structfr = StructType(

    [

        StructField("name", StringType(), True),

        StructField("food1",IntegerType(), True),

        StructField("food2",IntegerType(), True),

        StructField("food3",IntegerType(), True),

        StructField("food4",IntegerType(), True),

        StructField("placeid",IntegerType(), True)

    ]

)


structfp = StructType().add("placeid", IntegerType(), True).add("placename",StringType(), True)

foodratings = spark.read.schema(structfr).csv('/user/maria_dev/foodratings53475.csv')

foodplaces = spark.read.schema(structfp).csv('/user/maria_dev/foodplaces53475.csv')


foodratings.createOrReplaceTempView("foodratingsT")

foodplaces.createOrReplaceTempView("foodplacesT")


foodratings_ex3 = spark.sql("SELECT * FROM foodratingsT WHERE food2 < 25 AND food4 > 40")

foodratings_ex3.printSchema()
```

**print foodratings_ex3.head(5)**


**foodplaces_ex3 = spark.sql("SELECT * FROM foodplacesT WHERE placeid > 3")**

**foodplaces_ex3.printSchema()**

**print foodplaces_ex3.head(5)**

**execfile ('/home/maria_dev/assign_7_q3.py')**



Exercise 4)

Use an operation (not a SQL query) on the DataFrame 'foodratings' create in exercise 1 to create a new DataFrame called foodratings_ex4 that includes only those records (rows) where the 'name' field is "Mel" and food3 < 25.

As the results of this step provide the code you execute and screen shots of the following commands:

foodratings.printSchema()

foodratings.head(5)

```
vi assign_7_q4.py

from pyspark.sql.types import *

struct1 = StructType(

    [

        StructField("name", StringType(), True),

        StructField("food1",IntegerType(), True),

        StructField("food2",IntegerType(), True),

        StructField("food3",IntegerType(), True),

        StructField("food4",IntegerType(), True),

        StructField("placeid",IntegerType(), True)

    ]

)


foodratings = spark.read.schema(struct1).csv('/user/maria_dev/foodratings53475.csv')

foodratings_ex4 = foodratings.filter((foodratings['name'] == "Mel") & (foodratings['food3'] < 25))

foodratings_ex4.printSchema()

print foodratings_ex4.head(5)


execfile ('/home/maria_dev/assign_7_q4.py')
```

## Exercise 5)

Use an operation (not a SQL query) on the DataFrame 'foodratings' create in exercise 1 to create a new DataFrame called foodratings_ex5 that includes only the columns (fields) 'name' and 'placeid'

As the results of this step provide the code you execute and screen shots of the following commands:

foodratings.printSchema()

foodratings.head(5)

**Command Executed:**

**vi assign_7_q5.py**

**from pyspark.sql.types import ***

**struct1 = StructType(**

**[**

**StructField("name", StringType(), True),**

**StructField("food1",IntegerType(), True),**

```
        StructField("food2",IntegerType(), True),

        StructField("food3",IntegerType(), True),

        StructField("food4",IntegerType(), True),

        StructField("placeid",IntegerType(), True)

    ]

)


foodratings = spark.read.schema(struct1).csv('/user/maria_dev/foodratings53475.csv')

foodratings_ex5 = foodratings.select(foodratings['name'],foodratings['placeid'])

foodratings_ex5.printSchema()

print foodratings_ex5.head(5)

 execfile ('/home/maria_dev/assign_7_q5.py')
```

Exercise 6)

Use an operation on the DataFrame 'to create a new DataFrame called ex6 which is the inner join, on placeid, of the DataFrames 'foodratings; and 'foodplaces' created in exercises 1 and 2

As the results of this step provide the code you execute and screen shots of the following commands:

> ex6.printSchema()

> ex6.head(5)

**vi assign_7_q6.py**

**from pyspark.sql.types import ***

**structfr = StructType(**

**[**

**StructField("name", StringType(), True),**

**StructField("food1",IntegerType(), True),**

**StructField("food2",IntegerType(), True),**

**StructField("food3",IntegerType(), True),**

**StructField("food4",IntegerType(), True),**

**StructField("placeid",IntegerType(), True)**

**]**

**)**

**structfp = StructType().add("placeid", IntegerType(), True).add("placename",StringType(), True)**

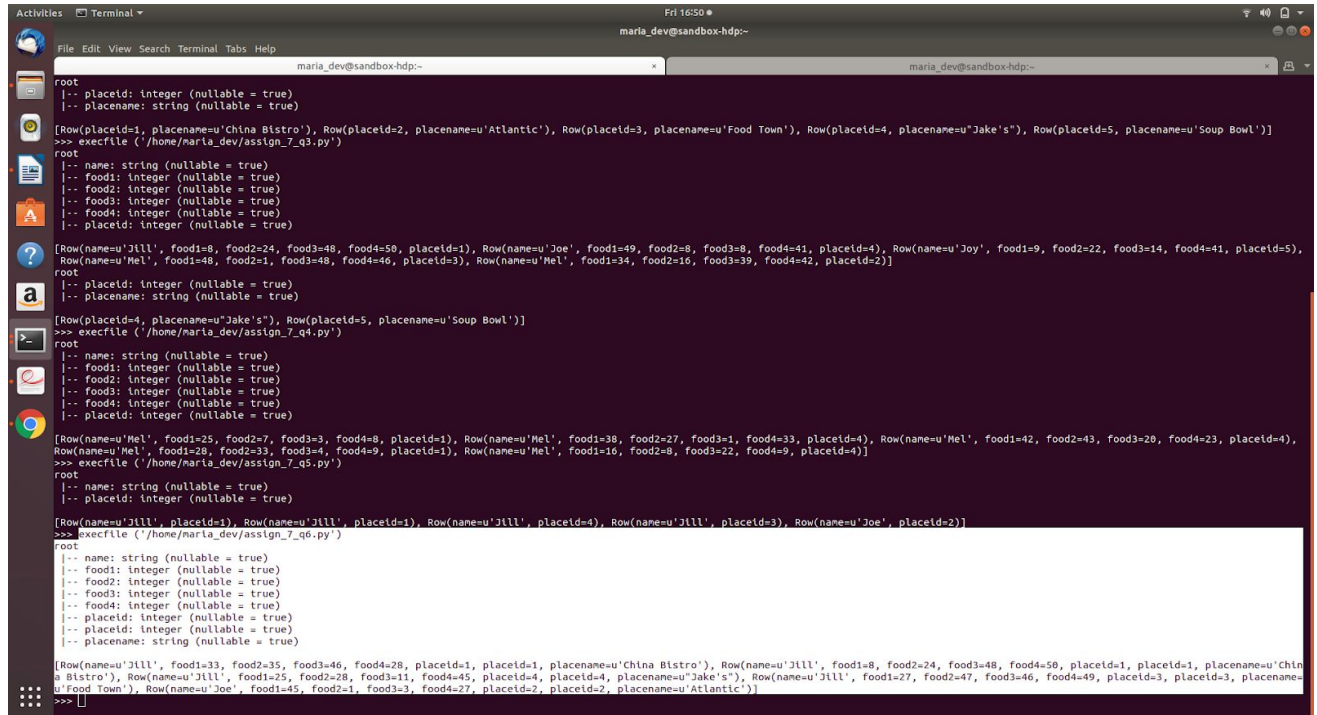**foodratings = spark.read.schema(structfr).csv('/user/maria_dev/foodratings53475.csv')**

**foodplaces = spark.read.schema(structfp).csv('/user/maria_dev/foodplaces53475.csv')**

**ex6 = foodratings.join(foodplaces, foodratings.placeid == foodplaces.placeid, 'inner')**

**ex6.printSchema()**

**print ex6.head(5)**

**execfile ('/home/maria_dev/assign_7_q6.py')**



```
root
 |-- placeid: integer (nullable = true)
 |-- placename: string (nullable = true)

[Row(placeid=1, placename=u'China Bistro'), Row(placeid=2, placename=u'Atlantic'), Row(placeid=3, placename=u'Food Town'), Row(placeid=4, placename=u"Jake's"), Row(placeid=5, placename=u'Soup Bowl')]
>>> execfile ('/home/maria_dev/assign_7_q3.py')
root
 |-- name: string (nullable = true)
 |-- food1: integer (nullable = true)
 |-- food2: integer (nullable = true)
 |-- food3: integer (nullable = true)
 |-- food4: integer (nullable = true)
 |-- placeid: integer (nullable = true)

[Row(name=u'Jill', food1=8, food2=24, food3=48, food4=50, placeid=1), Row(name=u'Joe', food1=49, food2=8, food3=8, food4=41, placeid=4), Row(name=u'Joy', food1=9, food2=22, food3=14, food4=41, placeid=5), Row(name=u'Mel', food1=48, food2=1, food3=48, food4=46, placeid=3), Row(name=u'Mel', food1=34, food2=16, food3=39, food4=42, placeid=2)]
root
 |-- placeid: integer (nullable = true)
 |-- placename: string (nullable = true)

[Row(placeid=4, placename=u"Jake's"), Row(placeid=5, placename=u'Soup Bowl')]
>>> execfile ('/home/maria_dev/assign_7_q4.py')
root
 |-- name: string (nullable = true)
 |-- food1: integer (nullable = true)
 |-- food2: integer (nullable = true)
 |-- food3: integer (nullable = true)
 |-- food4: integer (nullable = true)
 |-- placeid: integer (nullable = true)

[Row(name=u'Mel', food1=25, food2=7, food3=3, food4=8, placeid=1), Row(name=u'Mel', food1=38, food2=27, food3=1, food4=33, placeid=4), Row(name=u'Mel', food1=42, food2=43, food3=20, food4=23, placeid=4), Row(name=u'Mel', food1=28, food2=33, food3=4, food4=9, placeid=1), Row(name=u'Mel', food1=16, food2=8, food3=22, food4=9, placeid=4)]
>>> execfile ('/home/maria_dev/assign_7_q5.py')
root
 |-- name: string (nullable = true)
 |-- placeid: integer (nullable = true)

[Row(name=u'Jill', placeid=1), Row(name=u'Jill', placeid=1), Row(name=u'Jill', placeid=4), Row(name=u'Jill', placeid=3), Row(name=u'Joe', placeid=2)]
>>> execfile ('/home/maria_dev/assign_7_q6.py')
root
 |-- name: string (nullable = true)
 |-- food1: integer (nullable = true)
 |-- food2: integer (nullable = true)
 |-- food3: integer (nullable = true)
 |-- food4: integer (nullable = true)
 |-- placeid: integer (nullable = true)
 |-- placeid: integer (nullable = true)
 |-- placename: string (nullable = true)

[Row(name=u'Jill', food1=33, food2=35, food3=46, food4=28, placeid=1, placeid=1, placename=u'China Bistro'), Row(name=u'Jill', food1=8, food2=24, food3=48, food4=50, placeid=1, placeid=1, placename=u'China Bistro'), Row(name=u'Jill', food1=25, food2=28, food3=11, food4=45, placeid=4, placeid=4, placename=u"Jake's"), Row(name=u'Jill', food1=27, food2=47, food3=46, food4=49, placeid=3, placeid=3, placename=u'Food Town'), Row(name=u'Joe', food1=45, food2=1, food3=3, food4=27, placeid=2, placeid=2, placename=u'Atlantic')]
>>>
```