# CSP-554 Big Data Technologies

## Bairi Rohith Reddy - Assignment #7

**Exercise 1)**

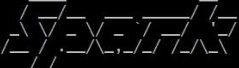Magic Number : 216304



foodratingsClass = StructType().add("name", StringType(), True).add ("food1", IntegerType(), True).add ("food2", IntegerType(), True).add("food3", StringType(), True). add ("food4", StringType(), True).add("placeid", StringType(), True)

foodratings = spark.read.schema(foodratingsClass).csv('/user/hadoop/foodratings216304.txt')

foodratings.printSchema()

foodratings.show(5)

```
>>> foodratings = spark.read.schema(foodratingsClass).csv('/user/hadoop/foodratings216304.txt')
>>>
>>> foodratings.printSchema()
root
 |-- name: string (nullable = true)
 |-- food1: integer (nullable = true)
 |-- food2: integer (nullable = true)
 |-- food3: string (nullable = true)
 |-- food4: string (nullable = true)
 |-- placeid: string (nullable = true)

>>> foodratings.show(5)
+----+-----+-----+-----+-----+-------+
|name|food1|food2|food3|food4|placeid|
+----+-----+-----+-----+-----+-------+
| Joy|   39|   33|   46|   50|      3|
| Sam|   15|   13|   20|   28|      1|
| Joe|   39|   46|   15|   21|      3|
| Sam|   36|   32|   40|   21|      4|
| Mel|   49|   37|   33|   23|      3|
+----+-----+-----+-----+-----+-------+
only showing top 5 rows
```

## Exercise 2)

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

foodplacesClass

foodplaces = spark.read.schema(foodplacesClass).csv('/user/hadoop/foodplaces216304.txt')

foodplaces.printSchema()

```
hadoop@ip-172-31-58-16:~
>>>
>>> foodplacesClass = StructType().add("placeid", StringType(), True).add ("placename", StringType(), True)
>>>
>>> foodplacesClass
StructType(List(StructField(placeid,StringType,true),StructField(placename,StringType,true)))
>>>
>>> foodplaces = spark.read.schema(foodplacesClass).csv('/user/hadoop/foodplaces216304.txt')
>>>
>>> foodplaces.printSchema()
root
 |-- placeid: string (nullable = true)
 |-- placename: string (nullable = true)
```

foodplaces.show(5)

```
hadoop@ip-172-31-58-16:~
>>>
>>> foodplacesClass = StructType().add("placeid", StringType(), True).add ("placename", StringType(), True)
>>>
>>> foodplacesClass
StructType(List(StructField(placeid,StringType,true),StructField(placename,StringType,true)))
>>>
>>> foodplaces = spark.read.schema(foodplacesClass).csv('/user/hadoop/foodplaces216304.txt')
>>>
>>> foodplaces.printSchema()
root
 |-- placeid: string (nullable = true)
 |-- placename: string (nullable = true)

>>> foodplaces.show(5)
+-------+------------+
|placeid|   placename|
+-------+------------+
|      1|China Bistro|
|      2|    Atlantic|
|      3|   Food Town|
|      4|      Jake's|
|      5|   Soup Bowl|
+-------+------------+
```

**Exercise 3)**

Register the DataFrames

foodratings.registerTempTable('foodratingsT');

foodplaces.registerTempTable('foodplacesT');



foodratings_ex3a = spark.sql("select * from foodratingsT where food2 < 25 and food4 > 40")

foodratings_ex3a.printSchema()



foodratings_ex3a.show(5)



foodplaces_ex3b = spark.sql("select * from foodplacesT where placeid > 3")

foodplaces_ex3b.printSchema()

```
hadoop@ip-172-31-58-16:~
>>>
>>> foodplaces_ex3b = spark.sql("select * from foodplacesT where placeid > 3")
>>>
>>> foodplaces_ex3b.printSchema()
root
 |-- placeid: string (nullable = true)
 |-- placename: string (nullable = true)

>>>
```

foodplaces_ex3b.show(5)

```
hadoop@ip-172-31-58-16:~
>>>
>>> foodplaces_ex3b.show(5)
+-------+---------+
|placeid|placename|
+-------+---------+
|      4|   Jake's|
|      5|Soup Bowl|
+-------+---------+
```

## Exercise 4)

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

foodratings_ex4.printSchema()

```
hadoop@ip-172-31-58-16:~
>>>
>>> foodratings_ex4 = foodratings.filter(foodratings['name'] == "Mel").filter(foodratings['food3'] < 25)
>>>
>>> foodratings_ex4.printSchema()
root
 |-- name: string (nullable = true)
 |-- food1: integer (nullable = true)
 |-- food2: integer (nullable = true)
 |-- food3: string (nullable = true)
 |-- food4: string (nullable = true)
 |-- placeid: string (nullable = true)

>>>
```

foodratings_ex4.show(5)

```
hadoop@ip-172-31-58-16:~
>>>
>>> foodratings_ex4.show(5)
+----+-----+-----+-----+-----+-------+
|name|food1|food2|food3|food4|placeid|
+----+-----+-----+-----+-----+-------+
| Mel|    8|    3|    2|    8|      4|
| Mel|   10|    6|    6|   50|      3|
| Mel|   43|    2|   10|   28|      2|
| Mel|   11|   43|   10|   44|      5|
| Mel|    8|   31|   11|   21|      2|
+----+-----+-----+-----+-----+-------+
only showing top 5 rows

>>>
```

**Exercise 5)**

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

foodratings_ex5.printSchema()

```
hadoop@ip-172-31-58-16:~
>>>
>>> foodratings_ex5 = foodratings.select(foodratings['name'], foodratings['placeid'])
>>>
>>> foodratings_ex5.printSchema()
root
 |-- name: string (nullable = true)
 |-- placeid: string (nullable = true)

>>>
```

foodratings_ex4.show(5)

```
hadoop@ip-172-31-58-16:~
>>>
>>> foodratings_ex5.show(5)
+----+-------+
|name|placeid|
+----+-------+
| Joy|      3|
| Sam|      1|
| Joe|      3|
| Sam|      4|
| Mel|      3|
+----+-------+
only showing top 5 rows

>>>
```

**Exercise 6)**

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

ex6.printSchema()

```
hadoop@ip-172-31-58-16:~
>>>
>>> ex6 = foodratings.join(foodplaces, foodratings.placeid == foodplaces.placeid, 'inner')
>>>
>>> ex6.printSchema()
root
 |-- name: string (nullable = true)
 |-- food1: integer (nullable = true)
 |-- food2: integer (nullable = true)
 |-- food3: string (nullable = true)
 |-- food4: string (nullable = true)
 |-- placeid: string (nullable = true)
 |-- placeid: string (nullable = true)
 |-- placename: string (nullable = true)

>>>
```

ex6.show(5)

```
hadoop@ip-172-31-58-16:~
>>>
>>> ex6.show(5)
+----+-----+-----+-----+-----+-------+-------+-----------+
|name|food1|food2|food3|food4|placeid|placeid|  placename|
+----+-----+-----+-----+-----+-------+-------+-----------+
| Joy|   39|   33|   46|   50|      3|      3|  Food Town|
| Sam|   15|   13|   20|   28|      1|      1|China Bistro|
| Joe|   39|   46|   15|   21|      3|      3|  Food Town|
| Sam|   36|   32|   40|   21|      4|      4|     Jake's|
| Mel|   49|   37|   33|   23|      3|      3|  Food Town|
+----+-----+-----+-----+-----+-------+-------+-----------+
only showing top 5 rows

>>> |
```