**Week 5 Assignment**

Guruprasad Velikadu Krishnamoorthy

College of Science and Technology, Bellevue University

**DSC650**-T301: Big Data

Professor. Nasheb Ismaily
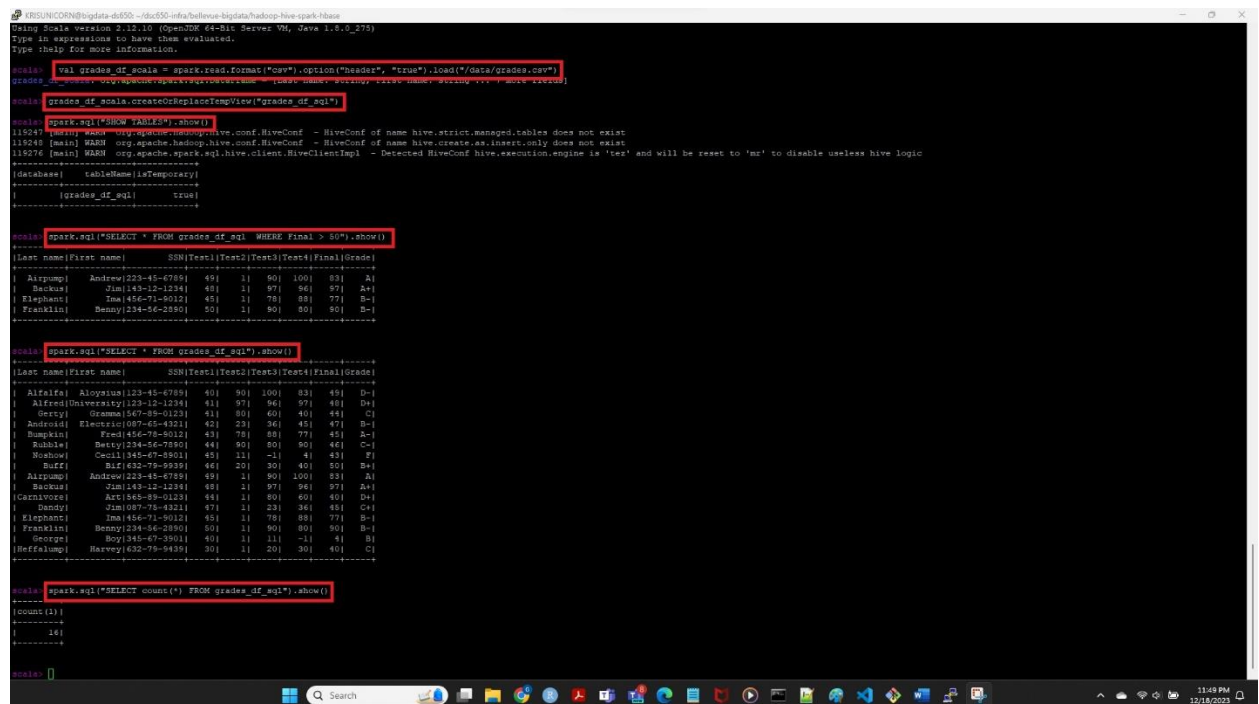
December 19, 2023

1. **SparkSQL with Scala:**

   val grades_df_scala = spark.read.format("csv").option("header", "true").load("/data/grades.csv")

   grades_df_scala.createOrReplaceTempView("grades_df_sql")

   spark.sql("SHOW TABLES").show()

   spark.sql("SELECT * FROM grades_df_sql WHERE Final > 50").show()

   spark.sql("SELECT * FROM grades_df_sql").show()



**Additional SparkSQL commands used:**

   1.1. spark.sql("SELECT * from grades_df_sql order by Final desc;").show()
   **Significance:** This command lists the in the Descending order of the values of column "Final"

   1.2. spark.sql("SELECT * from grades_df_sql where Test1>=0 and Test2>=0 and Test3>=0 and Test4 >=0 and Final>=0").show()
   **Significance:** This returns the records where the Test1,Test2, Test3 ,Test4 and Final are only positive values.

   1.3. spark.sql("SELECT AVG(Test1),min(Test1),max(Test1),std(Test1) from grades_df_sql ").show()
   **Significance:** Returns the statistical quantities of Test1 values.

2. **SparkSQL with Python:**

   grades_df_python = spark.read.format('csv').option('header', 'true').load('/data/grades.csv')

   grades_df_python.show()

   grades_df_python.createOrReplaceTempView('grades_df_py_sql')

   spark.sql('SHOW TABLES').show()

   spark.sql('SELECT * FROM grades_df_py_sql WHERE Final > 50').show()

**Additional SparkSQL commands used:**

2.1. spark.sql("SELECT * from grades_df_py_sql order by Final desc;").show()
   **Significance:** This command lists the in the Descending order of the values of column "Final"

2.2. spark.sql("SELECT * from grades_df_py_sql where Test1>=0 and Test2>=0 and Test3>=0 and Test4 >=0 and Final>=0").show()
   **Significance:** This returns the records where the Test1,Test2, Test3 ,Test4 and Final are only positive values.

2.3. spark.sql("SELECT AVG(Test1),min(Test1),max(Test1),std(Test1) from grades_df_py_sql ").show()
   **Significance:** Returns the statistical quantities of Test1 values.

2.4. spark.sql("SELECT Grade,count(*) from grades_df_py_sql group by Grade order by Grade").show()
   **Significance:** Returns the count of the number of students grouped by the Grade.

2.5. spark.sql("SELECT * from grades_df_py_sql where (Test1+Test2+Test3+Test4)/4  < Final ").show()
   **Significance:** Returns the records where the average of Test1-4 scores is lesser than the Final scores.

2.6. spark.sql("SELECT *, DENSE_RANK() OVER (ORDER BY Grade) AS ROW_RANK FROM grades_df_py_sql").show()
   **Significance:** This assigns Rank to the students ordered based on the Grades.

```
>>> spark.sql("SELECT * from grades_df_py_sql order by Final desc;").show()
```

| Last name | First name | SSN | Test1 | Test2 | Test3 | Test4 | Final | Grade |
|---|---|---|---|---|---|---|---|---|
| Backus | Jim | 143-12-1234 | 48 | 1 | 97 | 96 | 97 | A+ |
| Franklin | Benny | 234-56-2890 | 50 | 1 | 90 | 80 | 90 | B- |
| Airpump | Andrew | 223-45-6789 | 49 | 1 | 90 | 100 | 83 | A |
| Elephant | Ima | 456-71-9012 | 45 | 1 | 78 | 88 | 77 | B- |
| Buff | Bif | 632-79-9939 | 46 | 20 | 30 | 40 | 50 | B+ |
| Alfalfa | Aloysius | 123-45-6789 | 40 | 90 | 100 | 83 | 49 | D- |
| Alfred | University | 123-12-1234 | 41 | 97 | 96 | 97 | 48 | D+ |
| Android | Electric | 087-65-4321 | 42 | 23 | 36 | 45 | 47 | B- |
| Rubble | Betty | 234-56-7890 | 44 | 90 | 80 | 90 | 46 | C- |
| Dandy | Jim | 087-75-4321 | 47 | 1 | 23 | 36 | 45 | C+ |
| Bumpkin | Fred | 456-78-9012 | 43 | 78 | 88 | 77 | 45 | A- |
| Gerty | Gramma | 567-89-0123 | 41 | 80 | 60 | 40 | 44 | C |
| Noshow | Cecil | 345-67-8901 | 45 | 11 | -1 | 4 | 43 | F |
| Carnivore | Art | 565-89-0123 | 44 | 1 | 80 | 60 | 40 | D+ |
| Heffalump | Harvey | 632-79-9439 | 30 | 1 | 20 | 30 | 40 | C |
| George | Boy | 345-67-3901 | 40 | 1 | 111 | -1 | 4 | B |

```
>>> spark.sql("SELECT * from grades_df_py_sql where Test1>=0 and Test2>=0 and Test3>=0 and Test4 >=0 and Final>=0").show()
```

| Last name | First name | SSN | Test1 | Test2 | Test3 | Test4 | Final | Grade |
|---|---|---|---|---|---|---|---|---|
| Alfalfa | Aloysius | 123-45-6789 | 40 | 90 | 100 | 83 | 49 | D- |
| Alfred | University | 123-12-1234 | 41 | 97 | 96 | 97 | 48 | D+ |
| Gerty | Gramma | 567-89-0123 | 41 | 80 | 60 | 40 | 44 | C |
| Android | Electric | 087-65-4321 | 42 | 23 | 36 | 45 | 47 | B- |
| Bumpkin | Fred | 456-78-9012 | 43 | 78 | 88 | 77 | 45 | A- |
| Rubble | Betty | 234-56-7890 | 44 | 90 | 80 | 90 | 46 | C- |
| Buff | Bif | 632-79-9939 | 46 | 20 | 30 | 40 | 50 | B+ |
| Airpump | Andrew | 223-45-6789 | 49 | 1 | 90 | 100 | 83 | A |
| Backus | Jim | 143-12-1234 | 48 | 1 | 97 | 96 | 97 | A+ |
| Carnivore | Art | 565-89-0123 | 44 | 1 | 80 | 60 | 40 | D+ |
| Dandy | Jim | 087-75-4321 | 47 | 1 | 23 | 36 | 45 | C+ |
| Elephant | Ima | 456-71-9012 | 45 | 1 | 78 | 88 | 77 | B- |
| Franklin | Benny | 234-56-2890 | 50 | 1 | 90 | 80 | 90 | B- |
| Heffalump | Harvey | 632-79-9439 | 30 | 1 | 20 | 30 | 40 | C |

```
>>> spark.sql("SELECT AVG(Test1),min(Test1),max(Test1),std(Test1) from grades_df_py_sql ").show()
```

| avg(CAST(Test1 AS DOUBLE)) | min(Test1) | max(Test1) | std(CAST(Test1 AS DOUBLE)) |
|---|---|---|---|
| 43.4375 | 30 | 50 | 4.766899931166679 |



```
>>> spark.sql("SELECT Grade,count(*) from grades_df_py_sql group by Grade order by Grade").show()
```

| Grade | count(1) |
|---|---|
| A | 1 |
| A+ | 1 |
| A- | 1 |
| B | 1 |
| B+ | 1 |
| B- | 3 |
| C | 2 |
| C+ | 1 |
| C- | 1 |
| D+ | 2 |
| D- | 1 |
| F | 1 |

```
>>>
>>>
>>> spark.sql("SELECT * from grades_df_py_sql where (Test1+Test2+Test3+Test4)/4  < Final ").show()
```

| Last name | First name | SSN | Test1 | Test2 | Test3 | Test4 | Final | Grade |
|---|---|---|---|---|---|---|---|---|
| Android | Electric | 087-65-4321 | 42 | 23 | 36 | 45 | 47 | B- |
| Noshow | Cecil | 345-67-8901 | 45 | 11 | -1 | 4 | 43 | F |
| Buff | Bif | 632-79-9939 | 46 | 20 | 30 | 40 | 50 | B+ |
| Airpump | Andrew | 223-45-6789 | 49 | 1 | 90 | 100 | 83 | A |
| Backus | Jim | 143-12-1234 | 48 | 1 | 97 | 96 | 97 | A+ |
| Dandy | Jim | 087-75-4321 | 47 | 1 | 23 | 36 | 45 | C+ |
| Elephant | Ima | 456-71-9012 | 45 | 1 | 78 | 88 | 77 | B- |
| Franklin | Benny | 234-56-2890 | 50 | 1 | 90 | 80 | 90 | B- |
| Heffalump | Harvey | 632-79-9439 | 30 | 1 | 20 | 30 | 40 | C |

```
>>> spark.sql("SELECT *, DENSE_RANK() OVER (ORDER BY Grade) AS ROW_RANK FROM grades_df_py_sql").show()
```
23/... WARN org.apache.spark.sql.execution.window.WindowExec: No Partition Defined for window operation! Moving all data to a single partition, this can cause serious performance degradation.

| Last name | First name | SSN | Test1 | Test2 | Test3 | Test4 | Final | Grade | ROW_RANK |
|---|---|---|---|---|---|---|---|---|---|
| Airpump | Andrew | 223-45-6789 | 49 | 1 | 90 | 100 | 83 | A | 1 |
| Backus | Jim | 143-12-1234 | 48 | 1 | 97 | 96 | 97 | A+ | 2 |
| Bumpkin | Fred | 456-78-9012 | 43 | 78 | 88 | 77 | 45 | A- | 3 |
| George | Boy | 345-67-3901 | 40 | 1 | 111 | -1 | 4 | B | 4 |
| Buff | Bif | 632-79-9939 | 46 | 20 | 30 | 40 | 50 | B+ | 5 |
| Android | Electric | 087-65-4321 | 42 | 23 | 36 | 45 | 47 | B- | 6 |
| Elephant | Ima | 456-71-9012 | 45 | 1 | 78 | 88 | 77 | B- | 6 |
| Franklin | Benny | 234-56-2890 | 50 | 1 | 90 | 80 | 90 | B- | 6 |
| Gerty | Gramma | 567-89-0123 | 41 | 80 | 60 | 40 | 44 | C | 7 |
| Heffalump | Harvey | 632-79-9439 | 30 | 1 | 20 | 30 | 40 | C | 7 |
| Dandy | Jim | 087-75-4321 | 47 | 1 | 23 | 36 | 45 | C+ | 8 |
| Rubble | Betty | 234-56-7890 | 44 | 90 | 80 | 90 | 46 | C- | 9 |
| Alfred | University | 123-12-1234 | 41 | 97 | 96 | 97 | 48 | D+ | 10 |
| Carnivore | Art | 565-89-0123 | 44 | 1 | 80 | 60 | 40 | D+ | 10 |
| Alfalfa | Aloysius | 123-45-6789 | 40 | 90 | 100 | 83 | 49 | D- | 11 |
| Noshow | Cecil | 345-67-8901 | 45 | 11 | -1 | 4 | 43 | F | 12 |

```
>>> []
```

3. **SparkSQL** with Custom Data Set:

For this exercise, the Hollywood movies dataset used in the prior Hive assignment is reused.

**3.1. The description of some of the commonly used columns is explained below:**

- *Movie:* Name of the movie
- *LeadStudio:* The Studio that produced the movie.
- *RottenTomatoes:* Represents the score rated by Rotten Tomatoes critics.
- *AudienceScore*: Represents the score rated by Audience critics.
- *Genre*: Represents the Genre of the movie
- *TheatersOpenWeek*: Number of theatres where the movie was released in the opening week in the US and Canada. This is an integer value.

- **OpeningWeekend:** Box Office collection in the opening week in the US and Canada in USD. This is a floating number value.
- **DomesticGross:** *Total domestic collection in the US and Canada represented in USD.*
- **ForeignGross:** *Total collection from countries other than the US and Canada represented in USD.*
- **WorldGross:** Sum of Domestic and Foreign Gross collection represented in USD.
- **Budget:** Total movie budget represented in USD.
- **Profitability:** Total profits made from the movie represented in Percentage.
- **Year**: The year in which the movie was released.

3.2. Commands to copy the dataset into HDFS, start Spark shell and on Scala, and load the data into Data frame.

*# Copying file to HDFS*
```
hdfs dfs -put /data/HollywoodMovies.csv /data/HollywoodMovies.csv
hdfs dfs -ls  /data/HollywoodMovies.csv
```

*# Starting Spark shell on Scala*
```
spark-shell
```

*# Creating a Dataframe using Spark Scala and printing sample rows:*
```
val movies_df_scala=spark.read.format("csv").option("header","true").load("/data/HollywoodMovies.csv")
movies_df_scala.show(4)
```

*# Creating a View on the Dataframe to use Spark Sql:*
```
movies_df_scala.createOrReplaceTempView("movies_df_sql")
spark.sql("show tables").show()
```

### 3.3. Additional Scala Sql Commands used:

3.3.1.   spark.sql("select count(*) from movies_df_sql").show()
         **Significance:**  Returns the count of number of records  in the Dataframe.

3.3.2.   spark.sql("select Year,count(Movie) as movie_count from movies_df_sql where Year >=2007
         group by Year Order by Year").show()
         **Significance:** The query returns the count of the number of movies from each year in the
         dataframe.

3.3.3.   spark.sql("select max(Profitability),Year from movies_df_sql where Year >=2007 group by
         Year Order by Year").show()
         **Significance:** Returns the maximum profitability achieved by a movie over the period of
         2007-2013.

3.3.4.   spark.sql("select Genre,Year,avg(OpenProfit) as profits from movies_df_sql where Year in
         (2007,2009) and Year is not Null and Genre is not Null group by Genre,Year order by
         Year,profits").show()
         **Significance:** Returns  the average Profits made in the years 2007 and 2009 by Genre, the
         results are ordered in ascending order of Year and average profit.

3.3.5.   spark.sql("select round(avg(DomesticGross),2) as
         Avg_DomesticGross,round(var_pop(DomesticGross),2) as
         Var_DomesticGross,round(stddev_pop(DomesticGross),2) as Std_DomesticGross,Year from
         movies_df_sql where Year >=2007 group by Year order by Year").show()

**Significance:** Returns the statistical quantities such as Mean, Variance, and standard deviation of Domestic Gross , arranged by Year.

3.3.6. spark.sql("SELECT Movie,Year,Genre,Profitability as `Profit%` from movies_df_sql where Profitability in (select max(Profitability) from movies_df_sql where Year >=2007 group by Year) order by Year").show(false)
**Significance:** Using Subqueries, the SQL returns the movie name and its Genre with the highest profit % for each year.



3.3.7. spark.sql("SELECT Movie,RottenTomatoes,Year, RANK() OVER (PARTITION BY Year ORDER BY RottenTomatoes DESC) AS ROWNUM FROM movies_df_sql").show(false)
**Significance:** Ranks the Movies based on the Rotten Tomatoes rating for each year.

3.3.8. spark.sql("SELECT *, round(PERCENT_RANK() OVER (ORDER BY OpenProfit),2) AS Percentile from (SELECT Movie,OpenProfit,Year from movies_df_sql where Year=2010 and Genre='Action' and OpenProfit is not null)").show(50, false)
**Significance:** Returns the percentile value of Action movies of 2010, based on the Profits made.



```
scala> spark.sql("SELECT *, round(PERCENT_RANK() OVER (ORDER BY OpenProfit),2) AS Percentile from (SELECT Movie,OpenProfit,Year from movies_df_sql where Year=2010 and Genre='Action' and OpenProfit is not null)").show(50,false)
21:05:58 [main] WARN  org.apache.spark.sql.execution.window.WindowExec  - No Partition Defined for window operation! Moving all data to a single partition, this can cause serious performance degradation.
+------------------------------------------------------+----------+----+----------+
|Movie                                                 |OpenProfit|Year|Percentile|
+------------------------------------------------------+----------+----+----------+
|Jonah Hex                                             |11.49     |2010|0.0       |
|The Sorcerer's Apprentice                             |11.73     |2010|0.03      |
|Skyline                                               |116.9     |2010|0.05      |
|The Karate Kid                                        |139.25    |2010|0.08      |
|Green Zone                                            |14.3      |2010|0.11      |
|Cats & Dogs: The Revenge of Kitty Galore              |14.47     |2010|0.13      |
|Prince of Persia: The Sands of Time                   |15        |2010|0.16      |
|From Paris with Love                                  |15.77     |2010|0.18      |
|Knight and Day                                        |17.26     |2010|0.21      |
|Scott Pilgrim vs. the World                           |17.67     |2010|0.24      |
|Robin Hood                                            |18        |2010|0.26      |
|Repo Men                                              |19.06     |2010|0.29      |
|Killers                                               |21.07     |2010|0.32      |
|Unstoppable                                           |23.88     |2010|0.34      |
|The A-Team                                            |25.7      |2010|0.37      |
|Tron: Legacy                                          |25.9      |2010|0.39      |
|The Last Airbender                                    |26.87     |2010|0.42      |
|Salt                                                  |32.73     |2010|0.45      |
|Percy Jackson & the Olympians: The Lightning Thief    |32.84     |2010|0.47      |
|Faster                                                |35.5      |2010|0.5       |
|The Other Guys                                        |35.5      |2010|0.5       |
|Red                                                   |37.59     |2010|0.55      |
|The Losers                                            |37.6      |2010|0.58      |
|Inception                                             |39.25     |2010|0.61      |
|MacGruber                                             |40        |2010|0.63      |
|The Book of Eli                                       |40.88     |2010|0.66      |
|The Expendables                                       |42.44     |2010|0.68      |
|Resident Evil: Afterlife                              |44.5      |2010|0.71      |
|Clash of the Titans                                   |48.96     |2010|0.74      |
|The Bounty Hunter                                     |51.75     |2010|0.76      |
|Machete                                               |55        |2010|0.79      |
|Cop Out                                               |60.7      |2010|0.82      |
|Predators                                             |62        |2010|0.84      |
|Iron Man 2                                            |64.05     |2010|0.87      |
|Takers                                                |64.06     |2010|0.89      |
|Kick-Ass                                              |66        |2010|0.92      |
|Legion                                                |67.31     |2010|0.95      |
|Daybreakers                                           |75.5      |2010|0.97      |
|Brooklyn's Finest                                     |78.82     |2010|1.0       |
+------------------------------------------------------+----------+----+----------+

scala>
```