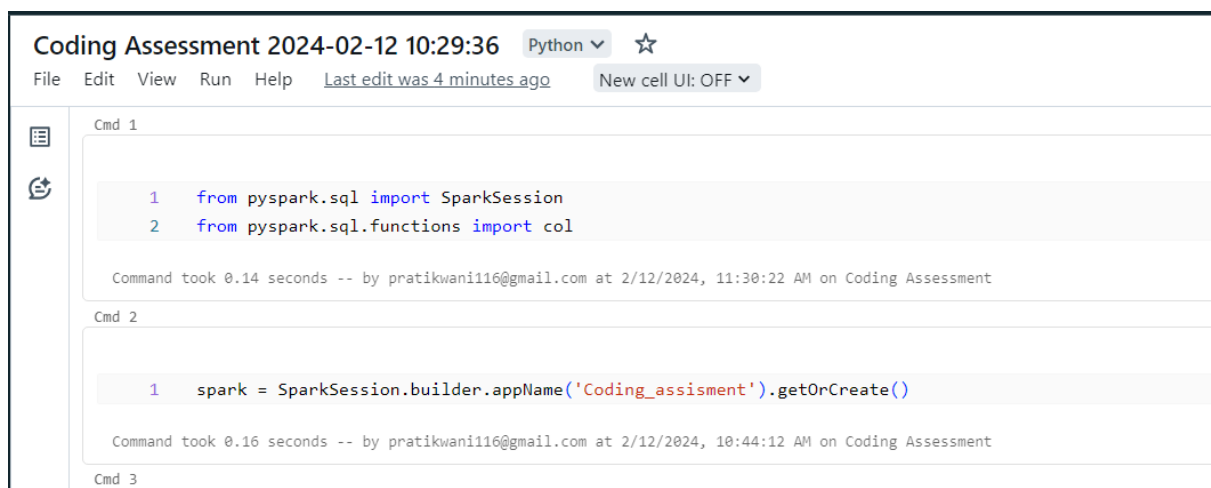# PySpark Coding Challenge

Pratik Wani

## Question No 1

- Created two dataframes with the named as df1 and df2
- All the operations (Manipulating, Droping, Sorting, Aggregations, Joining, GroupeBy) performed on these dataframes

- Creating Session:
  - After importing necessary Libraries Spark Session is created using getOrCreate()

```
Coding Assessment 2024-02-12 10:29:36   Python ∨   ☆
File  Edit  View  Run  Help   Last edit was 4 minutes ago       New cell UI: OFF ∨

Cmd 1

    1   from pyspark.sql import SparkSession
    2   from pyspark.sql.functions import col

  Command took 0.14 seconds -- by pratikwani116@gmail.com at 2/12/2024, 11:30:22 AM on Coding Assessment

Cmd 2

    1   spark = SparkSession.builder.appName('Coding_assisment').getOrCreate()

  Command took 0.16 seconds -- by pratikwani116@gmail.com at 2/12/2024, 10:44:12 AM on Coding Assessment

Cmd 3
```

- Creating DataFrames:

   o Two DataFrames are created as df1 and df2 using the method createDataFrame by passing the data and the header

Cmd 3

```
1   # Dataframe 1
2   Data1=[(1, "Pratik", 22, "Nashik"),
3          (2, "Vikas", 24, "Pune"),
4          (3, "Rushi", 23, "Chicago")]
5   Header1=["ID", "Name", "Age", "City"]
6   df1=spark.createDataFrame(Data1, Header1)
7   df1.show()
```

▶ (3) Spark Jobs

▶ 🖳 df1: pyspark.sql.dataframe.DataFrame = [ID: long, Name: string ... 2 more fields]

```
+---+------+---+-------+
| ID|  Name|Age|   City|
+---+------+---+-------+
|  1|Pratik| 22| Nashik|
|  2| Vikas| 24|   Pune|
|  3| Rushi| 23|Chicago|
+---+------+---+-------+
```

Command took 8.06 seconds -- by pratikwani116@gmail.com at 2/12/2024, 11:17:11 AM on Coding Assessment

```
1   #DataFrame 2
2   Data2 = [(1, "Data Engineer", 4000000),
3           (2, "Software Developer", 3500000),
4           (3, "Tester", 2000000)]
5   Header2 = ["ID", "Role", "Salary"]
6   df2=spark.createDataFrame(Data2, Header2)
7   df2.show()
```

▶ (3) Spark Jobs

▶ 🖳 df2: pyspark.sql.dataframe.DataFrame = [ID: long, Role: string ... 1 more field]

```
+---+------------------+-------+
| ID|              Role| Salary|
+---+------------------+-------+
|  1|     Data Engineer|4000000|
|  2|Software Developer|3500000|
|  3|            Tester|2000000|
+---+------------------+-------+
```

Command took 2.38 seconds -- by pratikwani116@gmail.com at 2/12/2024, 11:21:50 AM on Coding Assessment

- Manupulating Data:
    - After creating the df1 added the new data with ID 4 and Name Yogita in df1
    - To add the data I created new dataframe with newdata and then combined it with df1 using union method

```
1    # Manipulating data
2
3    newData=(4,"Yogita",30, "Nashik")
4    df1 = df1.union(spark.createDataFrame([newData]))
5    df1.show()
```

▶ (3) Spark Jobs

▶ 🖾 df1: pyspark.sql.dataframe.DataFrame = [ID: long, Name: string … 2 more fields]

```
+---+------+---+-------+
| ID|  Name|Age|   City|
+---+------+---+-------+
|  1|Pratik| 22| Nashik|
|  2| Vikas| 24|   Pune|
|  3| Rushi| 23|Chicago|
|  4|Yogita| 30| Nashik|
+---+------+---+-------+
```

Command took 1.93 seconds -- by pratikwani116@gmail.com at 2/12/2024, 11:39:35 AM on Coding Assessment

- Sorting:
  - The df1 is sorted by Age in Ascending Order while the df2 is sorted by Salary in Descending Order
  - To sort the data I used orderBy() method

```
1    # Sorting the Dataframe 1 by age and Dataframe 2 by salary
2
3    sorted_df1 = df1.orderBy("Age")
4    sorted_df2 = df2.orderBy(col("Salary").desc())
5    print("Sorted DataFrame1 by Age in Ascending Order:")
6    sorted_df1.show()
7    print("Sorted Dataframe3 by Salary in Descending Order:")
8    sorted_df2.show()
```

▶ (2) Spark Jobs

▶ 🗎 sorted_df1: pyspark.sql.dataframe.DataFrame = [ID: long, Name: string … 2 more fields]
▶ 🗎 sorted_df2: pyspark.sql.dataframe.DataFrame = [ID: long, Role: string … 1 more field]

```
Sorted DataFrame1 by Age in Ascending Order:
+---+------+---+-------+
| ID|  Name|Age|   City|
+---+------+---+-------+
|  1|Pratik| 22| Nashik|
|  3| Rushi| 23|Chicago|
|  2| Vikas| 24|   Pune|
|  4|Yogita| 30| Nashik|
+---+------+---+-------+

Sorted Dataframe3 by Salary in Descending Order:
+---+-----------------+-------+
| ID|             Role| Salary|
+---+-----------------+-------+
|  1|    Data Engineer|4000000|
|  2|Software Developer|3500000|
|  3|           Tester|2000000|
+---+-----------------+-------+
```

- Joining:
    - The df1 and df2 both having one column in common which is ID based on column ID both DataFrames are joined
    - To join the DataFrames we have to use the join() method and pass the name of second DataFrame with the Column name and type of join
    - In this example Inner join is performed

```
1    # Joining two Dataframes bases on the ID column ( INNER JOIN )
2
3    joined_df = df1.join(df2, "ID", "inner")
4    print("Joined DataFrame1 and Dataframe2:")
5    joined_df.show()
```

▸ (3) Spark Jobs

▸ ▤ joined_df: pyspark.sql.dataframe.DataFrame = [ID: long, Name: string ... 4 more fields]

```
Joined DataFrame1 and Dataframe2:
+---+------+---+-------+------------------+-------+
| ID|  Name|Age|   City|              Role| Salary|
+---+------+---+-------+------------------+-------+
|  1|Pratik| 22| Nashik|     Data Engineer|4000000|
|  2| Vikas| 24|   Pune|Software Developer|3500000|
|  3| Rushi| 23|Chicago|            Tester|2000000|
+---+------+---+-------+------------------+-------+
```

Command took 2.57 seconds -- by pratikwani116@gmail.com at 2/12/2024, 11:39:58 AM on Coding Assessment

- Group By:
  - To group the data bases on some field we have to use groupBy() method by passing the column name
  - We can perform different aggregate functions on such grouped data
  - In this example DataFrame df1 is groupby citywise and citywise total number of employees are counted.

```
1    # Group by Dataframe to count the total employee from each City
2
3    grouped_df = df1.groupBy("City").count()
4    print("Grouped DataFrame:")
5    grouped_df.show()
```

▸ (2) Spark Jobs

▸ ▤ grouped_df: pyspark.sql.dataframe.DataFrame = [City: string, count: long]

```
Grouped DataFrame:
+-------+-----+
|   City|count|
+-------+-----+
| Nashik|    2|
|   Pune|    1|
|Chicago|    1|
+-------+-----+
```

Command took 2.83 seconds -- by pratikwani116@gmail.com at 2/12/2024, 11:41:45 AM on Coding Assessment

- Aggregation:

  ○ After grouping the data we can perform different aggregations on the data

  ○ Such as Sum(), Max(), Min(), Avg()

  ○ Here 2 examples are performed one is finding Grand Total and another one is finding Citywise Average age of employees

```python
1    # Aggregations
2    # Grand Total Salary
3    Grand_total=df2.groupBy().sum("Salary")
4    print("Grand Total: ")
5    Grand_total.show()
6
7    # Calulate City wise Avarage age of emplyees
8    Avg_age=df1.groupBy().avg("Age")
9    print("City wise Avg age: ")
10   Avg_age.show()
11
```

▶ (4) Spark Jobs

▶ ▦  Grand_total: pyspark.sql.dataframe.DataFrame = [sum(Salary): long]

▶ ▦  Avg_age: pyspark.sql.dataframe.DataFrame = [avg(Age): double]

```
Grand Total:
+-----------+
|sum(Salary)|
+-----------+
|    9500000|
+-----------+


City wise Avg age:
+--------+
|avg(Age)|
+--------+
|   24.75|
+--------+
```

Command took 2.46 seconds -- by pratikwani116@gmail.com at 2/12/2024, 11:51:16 AM on Coding Assessment

- Droping:

  o To drop any column from the DataFrame we have to use drop method by passing the name of the columns we want to delete.

  o In this example column name City is deleted from the DataFrame df1

```
1    # Droping a column city from DataFrame1
2    Droped_df1=df1.drop("City")
3    print("New DataFrame1: ")
4    Droped_df1.show()
```

▸ (3) Spark Jobs

▸ ▤ Droped_df1: pyspark.sql.dataframe.DataFrame = [ID: long, Name: string … 1 more field]

```
New DataFrame1:
+---+------+---+
| ID|  Name|Age|
+---+------+---+
|  1|Pratik| 22|
|  2| Vikas| 24|
|  3| Rushi| 23|
|  4|Yogita| 30|
+---+------+---+
```

Command took 1.91 seconds -- by pratikwani116@gmail.com at 2/12/2024, 11:53:14 AM on Coding Assessment

# Question No 2

- The Database Student_info is created in that database two tables named as Students and Courses are created

- Join operations performed on these tables

- Create Database:



```
Coding Assessment_Question 2   2024-02-12 12:24:19   Python ∨   ☆
File   Edit   View   Run   Help    Last edit was 14 minutes ago      New cell UI: OFF ∨

Cmd 1

    1    from pyspark.sql import SparkSession

  Command took 0.07 seconds -- by pratikwani116@gmail.com at 2/12/2024, 12:24:50 PM on Coding Assessment

Cmd 2

    1    spark = SparkSession.builder.appName('Coding_assisment').getOrCreate()

  Command took 0.10 seconds -- by pratikwani116@gmail.com at 2/12/2024, 12:24:53 PM on Coding Assessment

Cmd 3

    1    # Creating DataBase
    2
    3    spark.sql("CREATE DATABASE IF NOT EXISTS Student_info")

  Out[3]: DataFrame[]

  Command took 3.35 seconds -- by pratikwani116@gmail.com at 2/12/2024, 12:25:57 PM on Coding Assessment

Cmd 4
```

- Create Tables

```
1   # Creating Tables
2
3   spark.sql("CREATE TABLE if not exists Student_info.Students (Roll_no Int, Name String, Course_id Int)")
4   spark.sql("CREATE TABLE if not exists Student_info.Courses (Course_name String, Course_id Int)")
5   spark.sql("Select * from Student_info.Students").show()
6   spark.sql("Select * from Student_info.Courses").show()
7
```

▶ (2) Spark Jobs

```
+-------+----+---------+
|Roll_no|Name|Course_id|
+-------+----+---------+
+-------+----+---------+


+-----------+---------+
|Course_name|Course_id|
+-----------+---------+
+-----------+---------+
```

Command took 4.33 seconds -- by pratikwani116@gmail.com at 2/12/2024, 12:29:17 PM on Coding Assessment

- Insert Records

```
1    # Inserting Dummy Data
2
3    spark.sql("INSERT INTO Student_info.Students (Roll_no, Name, Course_id) VALUES\
4    (1, 'Amit Kumar', 101),\
5    (2, 'Priya Patel', 102),\
6    (3, 'Rahul Sharma', 103),\
7    (4, 'Neha Gupta', 104),\
8    (5, 'Sandeep Singh', 105);")
9
10   spark.sql("INSERT INTO Student_info.Courses (Course_name, Course_id) VALUES\
11   ('Computer Science', 101),\
12   ('Electrical Engineering', 102),\
13   ('Mechanical Engineering', 103),\
14   ('Civil Engineering', 104),\
15   ('Mathematics', 105),\
16   ('Physics', 106),\
17   ('Chemistry', 107);")
18
```

```
1    # After Inserting Records
2
3    spark.sql("Select * from Student_info.Students").show()
4    spark.sql("Select * from Student_info.Courses").show()
```

▶ (2) Spark Jobs

```
|Roll_no|         Name|Course_id|
+-------+-------------+---------+
|      1|   Amit Kumar|      101|
|      2|  Priya Patel|      102|
|      3| Rahul Sharma|      103|
|      4|   Neha Gupta|      104|
|      5|Sandeep Singh|      105|
+-------+-------------+---------+


+-------------------+---------+
|        Course_name|Course_id|
+-------------------+---------+
|   Computer Science|      101|
|Electrical Engine...|      102|
|Mechanical Engine...|      103|
|   Civil Engineering|      104|
|        Mathematics|      105|
|            Physics|      106|
|          Chemistry|      107|
+-------------------+---------+
```

Command took 2.13 seconds -- by pratikwani116@gmail.com at 2/12/2024, 12:35:16 PM on Coding Assessment

- Inner Join:
  - Inner Join takes all the matching rows from the both tables
  - student is enrolled in courses with course_id 106, 107 so in the result set of inner join these two coursed will not include

```
1   # Performing Joins
2   #inner Join
3   spark.sql("SELECT * FROM Student_info.Students INNER JOIN Student_info.Courses\
4            ON Student_info.Students.Course_id = Student_info.Courses.Course_id" ).show()
```

▶ (2) Spark Jobs

```
+-------+-------------+---------+--------------------+---------+
|Roll_no|         Name|Course_id|         Course_name|Course_id|
+-------+-------------+---------+--------------------+---------+
|      1|   Amit Kumar|      101|    Computer Science|      101|
|      2|  Priya Patel|      102|Electrical Engine...|      102|
|      3| Rahul Sharma|      103|Mechanical Engine...|      103|
|      4|   Neha Gupta|      104|    Civil Engineering|     104|
|      5|Sandeep Singh|      105|         Mathematics|      105|
+-------+-------------+---------+--------------------+---------+
```

Command took 3.06 seconds -- by pratikwani116@gmail.com at 2/12/2024, 12:58:28 PM on Coding Assessment

- Left Join:
  - Left Join takes all the rows from Left Table and Matching Rows from the Right Table
  - The Unmatched fields are filled with Null values

```
1   # Left Join
2
3   spark.sql("SELECT * FROM Student_info.Students LEFT JOIN Student_info.Courses\
4           ON Student_info.Students.Course_id = Student_info.Courses.Course_id" ).show()
```

▶ (2) Spark Jobs

```
+-------+-------------+---------+--------------------+---------+
|Roll_no|         Name|Course_id|         Course_name|Course_id|
+-------+-------------+---------+--------------------+---------+
|      1|   Amit Kumar|      101|    Computer Science|      101|
|      2|  Priya Patel|      102|Electrical Engine...|      102|
|      3| Rahul Sharma|      103|Mechanical Engine...|      103|
|      4|   Neha Gupta|      104|    Civil Engineering|     104|
|      5|Sandeep Singh|      105|         Mathematics|      105|
+-------+-------------+---------+--------------------+---------+
```

Command took 2.42 seconds -- by pratikwani116@gmail.com at 2/12/2024, 1:00:13 PM on Coding Assessment

Cmd 0

- Right Join:
  - Right Join takes all the rows from Right Table and Matching Rows from the Left Table
  - The Unmatched fields are filled with Null values

```
1    # Right Join
2
3    spark.sql("SELECT * FROM Student_info.Students RIGHT JOIN Student_info.Courses\
4         |     |     ON Student_info.Students.Course_id = Student_info.Courses.Course_id" ).show()
```

▶ (2) Spark Jobs

```
+-------+-------------+---------+--------------------+---------+
|Roll_no|         Name|Course_id|         Course_name|Course_id|
+-------+-------------+---------+--------------------+---------+
|      1|  Amit Kumar|      101|    Computer Science|      101|
|      2|  Priya Patel|      102|Electrical Engine...|      102|
|      3| Rahul Sharma|      103|Mechanical Engine...|      103|
|      4|   Neha Gupta|      104|   Civil Engineering|      104|
|      5|Sandeep Singh|      105|         Mathematics|      105|
|   null|         null|     null|             Physics|      106|
|   null|         null|     null|           Chemistry|      107|
+-------+-------------+---------+--------------------+---------+
```

Command took 1.92 seconds -- by pratikwani116@gmail.com at 2/12/2024, 1:00:31 PM on Coding Assessment

- Left Anti Join:
  - Left Anti Join takes all the unmatched rows from left Table
  - The Unmatched fields are filled with Null values
  - In this example there is no such unmatched row hence result set is empty

```
1    # Left Anti join
2
3    spark.sql("SELECT * FROM Student_info.Students Left ANTI JOIN Student_info.Courses\
4            ON Student_info.Students.Course_id = Student_info.Courses.Course_id" ).show()
5
```

▶ (2) Spark Jobs

```
+-------+----+---------+
|Roll_no|Name|Course_id|
+-------+----+---------+
+-------+----+---------+
```

Command took 1.96 seconds -- by pratikwani116@gmail.com at 2/12/2024, 1:02:32 PM on Coding Assessment

- Left Semi Join:
  - Left Semi Join takes all the matched rows from left Table

```
1    # Left Semi Join
2
3    spark.sql("SELECT * FROM Student_info.Students LEFT SEMI JOIN Student_info.Courses\
4           ON Student_info.Students.Course_id = Student_info.Courses.Course_id" ).show()
```

▸ (2) Spark Jobs

```
+-------+-------------+---------+
|Roll_no|         Name|Course_id|
+-------+-------------+---------+
|      1|  Amit Kumar|      101|
|      2| Priya Patel|      102|
|      3|Rahul Sharma|      103|
|      4|  Neha Gupta|      104|
|      5|Sandeep Singh|     105|
+-------+-------------+---------+
```

Command took 2.11 seconds -- by pratikwani116@gmail.com at 2/12/2024, 1:03:10 PM on Coding Assessment

- Applying Functions in a Pandas DataFrame

  o DataFrame is created in which we have weekly incomes

  o The function is created to get the total income of each month

**Coding Assessment_Question 2**  2024-02-12 12:24:19   Python ∨   ☆

File   Edit   View   Run   Help      Last edit was 8 minutes ago      New cell UI: OFF ∨

Command took 2.11 seconds -- by pratikwani116@gmail.com at 2/12/2024, 1:03:10 PM on Coding Assessment

Cmd 12

```
1    import pyspark.pandas as pd
2    import numpy as np
```

Command took 0.08 seconds -- by pratikwani116@gmail.com at 2/12/2024, 1:23:12 PM on Coding Assessment

Cmd 13

```
1
2    # DataFrame Created
3    Data = ({
4
5        'Four Week Income of March':[1000,2500,1500,1200],
6        'Four Week Income of may':[1000,5008,1004,2006]
7               })
8
9    pandasDataFrame= pd.DataFrame(Data)
10   print(type(pandasDataFrame))
11   print(pandasDataFrame)
12
```

▶ (1) Spark Jobs

```
<class 'pyspark.pandas.frame.DataFrame'>
  Four Week Income of March  Four Week Income of may
0                      1000                     1000
1                      2500                     5008
2                      1500                     1004
3                      1200                     2006
```

Command took 0.53 seconds -- by pratikwani116@gmail.com at 2/12/2024, 1:23:15 PM on Coding Assessment

```
1    # Function for total income
2
3    def Total_income_Monthwise(DF):
4        return DF[0]+DF[1]
5
6    Final_DF = pandasDataFrame.apply(Total_income_Monthwise)
7    print(Final_DF)
```

▶ (2) Spark Jobs

```
Four Week Income of March    3500
Four Week Income of may      6008
dtype: int64
```

💡 1

Command took 0.82 seconds -- by pratikwani116@gmail.com at 2/12/2024, 1:23:19 PM on Coding Assessment