# PySpark Assignment 4

Pratik Wani

- Notes

---

| Day - 4 |

* reading files into dataframe :-

.read.format().load()
.read.csv()
.read.text()

* Reading ways :-

spark.read.format ("text").load (Path, format = N,
schema = None)

* Adding new column

* we have to use withcolumn() function with lit() tas parameter to give constant values to the column.

dataframe. withcolumn ("column-name, lit(vare)

Lit funtion is present in pyspark.sql.funchion library.

* concrate two columns

concat_us ()

Syntax → concat us ( "seprator", 'onecol', 'twocol'

concat-us present in pyspark. sql.functione)

* Group By and Aggregations:-

° group by is used to group to resut of dataframe based on some columns.
° after grouping the data we can perform different aggregate Functions on the resut set.

. such as sum(), min(), count(), max(), Avg()

Syntax:-

   dataframe. groupBy ('Department'). sum ('salary')
                              ↗
              grouping the data based on
              department

Note
   → ① inferschema = True,
                         , df infers the data
                         according to different
                         datatypes.

+ Pivot()
        → It can convert some row values
          into the columns.

+ Missing values:-

   1) To drop the NULL values we can use
      drop() function.

      dataframe.NA. drop()

   3 parameters:-

① How = 'all'
→ only drop if all the values are NULL

② How = 'any'
→ drop if any value is NULL

③ thresh = 2
→ means if at least two values must be Not NULL then only data is not droped.

④ subset = ['col name']
→ only check null values for these column.

\* <u>Fill NULL values :-</u>

• fill ('string') ~ for string value
• fill (o) ~ for int values

\* <u>Sort :-</u>

use to sort the data based on perticular column.

3 ways :-

ⓐ
df. sort ⑤ → "col1"
→ ascending order

ⓑ
a df. sort df["col1"]. desc().

↳ descending order

© multiple column

    df. Sort (" col1", "col2")

* Order by

    we can also use orderby to Sort values

    df. orderby (" col")

* joins

    similar as saL

    Syntax 1 :-

- df1 . join (df2, df1.col_name = df2.col_name, 'type')

- New types

    ⓐ Left semi :- Add a Match from right

    ⓑ Left Anti :- unmatch from right

    ⓒ Left outer :- semi + Anti

- Reading Text File and Adding New Column:

+ Code  + Markdown  ▷ Run All  ↻ Restart  ☰ Clear All Outputs  Variables  ☰ Outline  ⋯

```python
from pyspark.sql import SparkSession
```
[1]  ✓ 0.7s

```python
spark=SparkSession.builder.appName("Practice").getOrCreate()
```
[2]  ✓ 15.1s

```python
data=spark.read.csv(r"DEPT_DATA.CSV",header=True,inferSchema=True)
```
[3]  ✓ 18.5s

```python
data.show()
```
[4]  ✓ 0.7s

```
+----------+-------+
|dept_name"|dept_id|
+----------+-------+
|   Finance|     10|
| Marketing|     20|
|     Sales|     30|
|        IT|     40|
+----------+-------+
```

+ Code  + Markdown  ▷ Run All  ↻ Restart  ☰ Clear All Outputs  Variables  ☰ Outline  ⋯

```python
from pyspark.sql.functions import lit
data.withColumn('Total_Staff',lit(10)).show()
```
[5]  ✓ 0.4s

```
+----------+-------+-----------+
|dept_name"|dept_id|Total_Staff|
+----------+-------+-----------+
|   Finance|     10|         10|
| Marketing|     20|         10|
|     Sales|     30|         10|
|        IT|     40|         10|
+----------+-------+-----------+
```

```python
data.withColumn('Total_Staff',data.dept_id*5).show()
```
[10]  ✓ 0.2s

```
+----------+-------+-----------+
|dept_name"|dept_id|Total_Staff|
+----------+-------+-----------+
|   Finance|     10|         50|
| Marketing|     20|        100|
|     Sales|     30|        150|
|        IT|     40|        200|
+----------+-------+-----------+
```

```python
from pyspark.sql.functions import concat_ws
```

[14]  ✓  0.0s

```python
data.withColumn('Details',concat_ws('-','dept_id','dept_name"')).show()
```

[17]  ✓  0.4s

```
+----------+-------+-----------+
|dept_name"|dept_id|    Details|
+----------+-------+-----------+
|   Finance|     10| 10-Finance|
| Marketing|     20|20-Marketing|
|     Sales|     30|   30-Sales|
|        IT|     40|      40-IT|
+----------+-------+-----------+
```

```python
# Reading Text file
df1 = spark.read.text("sample.txt")

df1.selectExpr("split(value, ' ') as Text_Data").show(4,False)
```

[21]  ✓  0.2s

```
+----------------------------------------------------------+
|Text_Data                                                 |
+----------------------------------------------------------+
|[Hii!!, My, name, is, pratik, arun, wani., I, am, IT, Engineer.]|
+----------------------------------------------------------+
```

```python
# Reading CSV file
df2 = spark.read.text("DEPT_DATA.CSV")

df2.selectExpr("split(value, ' ') as Text_Data").show(4,False)
```

[22]  ✓  0.3s

```
+----------------------+
|Text_Data             |
+----------------------+
|[dept_name","dept_id"]|
|["Finance",10]        |
|["Marketing",20]      |
|["Sales",30]          |
+----------------------+
only showing top 4 rows
```

- Group By, Handling Missing Values



```python
from pyspark.sql import SparkSession
```

```python
spark=SparkSession.builder.appName("Practice").getOrCreate()
```

```python
data=spark.read.csv(r"C:\Users\dell\OneDrive\Documents\Desktop\Hexaware_Data_Engineering_Batch1\PySpark\Codes\salary.csv",
```

```python
data.show()
```

```
+------+----+------+
|  Name|Dept|Salary|
+------+----+------+
|Pratik|   A|500000|
| Vikki|   A| 60000|
| Rushi|   B| 90000|
|Shravi|   B|500000|
|Sweeha|   C|100000|
| Sanvi|   C|200000|
|Spruha|NULL|100000|
|  Arun|   B|  NULL|
|  NULL|   A|500000|
+------+----+------+
```



```python
data.groupBy('Dept').avg('Salary').show()
```

```
+----+-----------------+
|Dept|      avg(Salary)|
+----+-----------------+
|NULL|         100000.0|
|   B|         295000.0|
|   C|         150000.0|
|   A|353333.3333333333|
+----+-----------------+
```

```python
data.groupBy('Dept').agg({'Dept':'count','Salary':'sum'}).show()
```

```
+----+-----------+-----------+
|Dept|sum(Salary)|count(Dept)|
+----+-----------+-----------+
|NULL|     100000|          0|
|   B|     590000|          3|
|   C|     300000|          2|
|   A|    1060000|          3|
+----+-----------+-----------+
```

```python
from pyspark.sql import functions as F
data.groupBy('Dept').sum("Salary").select(F.col("Dept").alias("Department_Name"),
                                           F.col("sum(Salary)").alias("Departmentwise_sum")).show()
```

[37]

```
+---------------+------------------+
|Department_Name|Departmentwise_sum|
+---------------+------------------+
|           NULL|            100000|
|              B|            590000|
|              C|            300000|
|              A|           1060000|
+---------------+------------------+
```

```python
data.groupBy("Dept").pivot("Name").sum("salary").show()
```

[38]

```
+----+------+----+------+------+------+------+------+------+-----+
|Dept|  null|Arun|Pratik| Rushi| Sanvi|Shravi|Spruha|Sweeha|Vikki|
+----+------+----+------+------+------+------+------+------+-----+
|NULL|  NULL|NULL|  NULL|  NULL|  NULL|  NULL|100000|  NULL| NULL|
|   B|  NULL|NULL|  NULL| 90000|  NULL|500000|  NULL|  NULL| NULL|
|   C|  NULL|NULL|  NULL|  NULL|200000|  NULL|  NULL|100000| NULL|
|   A|500000|NULL|500000|  NULL|  NULL|  NULL|  NULL|  NULL|60000|
+----+------+----+------+------+------+------+------+------+-----+
```

#entered null values here

+ Code  + Markdown  | ▷ Run All  ⟳ Restart  ≡ Clear All Outputs  | 🔢 Variables  ≣ Outline  ⋯

```python
#entered null values here
data.show()
```

[39]

```
+------+----+------+
|  Name|Dept|Salary|
+------+----+------+
|Pratik|   A|500000|
| Vikki|   A| 60000|
| Rushi|   B| 90000|
|Shravi|   B|500000|
|Sweeha|   C|100000|
| Sanvi|   C|200000|
|Spruha|NULL|100000|
|  Arun|   B|  NULL|
|  NULL|   A|500000|
+------+----+------+
```

```python
data.na.drop().show()
```

[40]

```
+------+----+------+
|  Name|Dept|Salary|
+------+----+------+
|Pratik|   A|500000|
| Vikki|   A| 60000|
| Rushi|   B| 90000|
|Shravi|   B|500000|
|Sweeha|   C|100000|
| Sanvi|   C|200000|
+------+----+------+
```

+ Code   + Markdown   ▷ Run All   ↺ Restart   ☰ Clear All Outputs   (x)

```python
data.na.drop(how="all").show()
```

[53]

```
+------+----+------+
|  Name|Dept|Salary|
+------+----+------+
|Pratik|   A|500000|
| Vikki|   A| 60000|
| Rushi|   B| 90000|
|Shravi|   B|500000|
|Sweeha|   C|100000|
| Sanvi|   C|200000|
|Spruha|NULL|100000|
|  Arun|   B|  NULL|
|  NULL|   A|500000|
+------+----+------+
```

```python
data.na.drop(how="any",thresh=2).show()
data.na.drop(how="any",thresh=3).show()
```

[54]

```
+------+----+------+
|  Name|Dept|Salary|
+------+----+------+
|Pratik|   A|500000|
| Vikki|   A| 60000|
| Rushi|   B| 90000|
|Shravi|   B|500000|
|Sweeha|   C|100000|
| Sanvi|   C|200000|
|Spruha|NULL|100000|
|  Arun|   B|  NULL|
|  NULL|   A|500000|
+------+----+------+


+------+----+------+
|  Name|Dept|Salary|
+------+----+------+
|Pratik|   A|500000|
| Vikki|   A| 60000|
| Rushi|   B| 90000|
|Shravi|   B|500000|
|Sweeha|   C|100000|
| Sanvi|   C|200000|
+------+----+------+
```

```python
data.na.drop(how='any',subset=['Salary']).show()
```

[55]

```
+------+----+------+
|  Name|Dept|Salary|
+------+----+------+
|Pratik|   A|500000|
| Vikki|   A| 60000|
| Rushi|   B| 90000|
|Shravi|   B|500000|
|Sweeha|   C|100000|
| Sanvi|   C|200000|
|Spruha|NULL|100000|
|  NULL|   A|500000|
+------+----+------+
```

```python
(data.na.fill('Not known')).na.fill(0).show()
```

[60]

```
+---------+---------+------+
|     Name|     Dept|Salary|
+---------+---------+------+
|   Pratik|        A|500000|
|    Vikki|        A| 60000|
|    Rushi|        B| 90000|
|   Shravi|        B|500000|
|   Sweeha|        C|100000|
|    Sanvi|        C|200000|
|   Spruha|Not known|100000|
|     Arun|        B|     0|
|Not known|        A|500000|
```

- Sorting



```
data.sort('Salary','Dept').show()
```

```
+------+----+------+
|  Name|Dept|Salary|
+------+----+------+
|  Arun|   B|  NULL|
| Vikki|   A| 60000|
| Rushi|   B| 90000|
|Spruha|NULL|100000|
|Sweeha|   C|100000|
| Sanvi|   C|200000|
|Pratik|   A|500000|
|  NULL|   A|500000|
|Shravi|   B|500000|
+------+----+------+
```

```
data.sort(data['Salary'].desc()).show()
```

```
+------+----+------+
|  Name|Dept|Salary|
+------+----+------+
|Pratik|   A|500000|
|Shravi|   B|500000|
|  NULL|   A|500000|
| Sanvi|   C|200000|
|Sweeha|   C|100000|
|Spruha|NULL|100000|
| Rushi|   B| 90000|
| Vikki|   A| 60000|
```

- Joins

```python
#Inner Join
data1.join(data2,data1.emp_dept_id==data2.dept_id,"inner").na.fill('UnKnown').show()
```

[9]  ✓  0.6s

```
+------+--------+-------------+-----------+-----------+-------+------+----------+-------+
|emp_id|    name|superior_emp_id|year_joined|emp_dept_id|gender|salary|dept_name"|dept_id|
+------+--------+-------------+-----------+-----------+-------+------+----------+-------+
|     1|   Smith|           -1|       2018|         10|      M|  3000|   Finance|     10|
|     2|    Rose|            1|       2010|         20|      M|  4000| Marketing|     20|
|     3|Williams|            1|       2010|         10|      M|  1000|   Finance|     10|
|     4|   Jones|            2|       2005|         10|      F|  2000|   Finance|     10|
|     5|   Brown|            2|       2010|         40|UnKnown|    -1|        IT|     40|
+------+--------+-------------+-----------+-----------+-------+------+----------+-------+
```

```python
#Outer Join
data1.join(data2,data1.emp_dept_id==data2.dept_id,"outer").show()
```

[10]  ✓  1.5s

```
+------+--------+-------------+-----------+-----------+------+------+----------+-------+
|emp_id|    name|superior_emp_id|year_joined|emp_dept_id|gender|salary|dept_name"|dept_id|
+------+--------+-------------+-----------+-----------+------+------+----------+-------+
|     1|   Smith|           -1|       2018|         10|     M|  3000|   Finance|     10|
|     3|Williams|            1|       2010|         10|     M|  1000|   Finance|     10|
|     4|   Jones|            2|       2005|         10|     F|  2000|   Finance|     10|
|     2|    Rose|            1|       2010|         20|     M|  4000| Marketing|     20|
|  NULL|    NULL|         NULL|       NULL|       NULL|  NULL|  NULL|     Sales|     30|
|     5|   Brown|            2|       2010|         40|  NULL|    -1|        IT|     40|
|     6|   Brown|            2|       2010|         50|  NULL|    -1|      NULL|   NULL|
+------+--------+-------------+-----------+-----------+------+------+----------+-------+
```

```python
# Full join
data1.join(data2,data1.emp_dept_id==data2.dept_id,"full").show()
```

[11]  ✓  0.8s

```
+------+--------+-------------+-----------+-----------+------+------+----------+-------+
|emp_id|    name|superior_emp_id|year_joined|emp_dept_id|gender|salary|dept_name"|dept_id|
+------+--------+-------------+-----------+-----------+------+------+----------+-------+
|     1|   Smith|           -1|       2018|         10|     M|  3000|   Finance|     10|
|     3|Williams|            1|       2010|         10|     M|  1000|   Finance|     10|
|     4|   Jones|            2|       2005|         10|     F|  2000|   Finance|     10|
|     2|    Rose|            1|       2010|         20|     M|  4000| Marketing|     20|
|  NULL|    NULL|         NULL|       NULL|       NULL|  NULL|  NULL|     Sales|     30|
|     5|   Brown|            2|       2010|         40|  NULL|    -1|        IT|     40|
|     6|   Brown|            2|       2010|         50|  NULL|    -1|      NULL|   NULL|
+------+--------+-------------+-----------+-----------+------+------+----------+-------+
```

```python
# left join
data1.join(data2,data1.emp_dept_id==data2.dept_id,"left").show()
```

[12]  ✓  0.5s

```
+------+--------+-------------+-----------+-----------+------+------+----------+-------+
|emp_id|    name|superior_emp_id|year_joined|emp_dept_id|gender|salary|dept_name"|dept_id|
+------+--------+-------------+-----------+-----------+------+------+----------+-------+
|     1|   Smith|           -1|       2018|         10|     M|  3000|   Finance|     10|
|     2|    Rose|            1|       2010|         20|     M|  4000| Marketing|     20|
|     3|Williams|            1|       2010|         10|     M|  1000|   Finance|     10|
|     4|   Jones|            2|       2005|         10|     F|  2000|   Finance|     10|
|     5|   Brown|            2|       2010|         40|  NULL|    -1|        IT|     40|
|     6|   Brown|            2|       2010|         50|  NULL|    -1|      NULL|   NULL|
+------+--------+-------------+-----------+-----------+------+------+----------+-------+
```

```python
# Right join
data1.join(data2,data1.emp_dept_id==data2.dept_id,"right").show()
```
[13] ✓ 0.6s

```
+------+--------+-------------+-----------+-----------+------+------+---------+-------+
|emp_id|    name|superior_emp_id|year_joined|emp_dept_id|gender|salary|dept_name"|dept_id|
+------+--------+-------------+-----------+-----------+------+------+---------+-------+
|     4|   Jones|            2|       2005|         10|     F|  2000|  Finance|     10|
|     3|Williams|            1|       2010|         10|     M|  1000|  Finance|     10|
|     1|   Smith|           -1|       2018|         10|     M|  3000|  Finance|     10|
|     2|    Rose|            1|       2010|         20|     M|  4000|Marketing|     20|
|  NULL|    NULL|         NULL|       NULL|       NULL|  NULL|  NULL|    Sales|     30|
|     5|   Brown|            2|       2010|         40|  NULL|    -1|       IT|     40|
+------+--------+-------------+-----------+-----------+------+------+---------+-------+
```

```python
# LeftSemi join
data1.join(data2,data1.emp_dept_id==data2.dept_id,"leftsemi").show()
```
[14] ✓ 0.6s

```
+------+--------+-------------+-----------+-----------+------+------+
|emp_id|    name|superior_emp_id|year_joined|emp_dept_id|gender|salary|
+------+--------+-------------+-----------+-----------+------+------+
|     1|   Smith|           -1|       2018|         10|     M|  3000|
|     2|    Rose|            1|       2010|         20|     M|  4000|
|     3|Williams|            1|       2010|         10|     M|  1000|
|     4|   Jones|            2|       2005|         10|     F|  2000|
|     5|   Brown|            2|       2010|         40|  NULL|    -1|
+------+--------+-------------+-----------+-----------+------+------+
```

```python
data1.join(data2,data1.emp_dept_id==data2.dept_id,"leftanti").show()
```
[15] ✓ 0.4s

```
+------+-----+-------------+-----------+-----------+------+------+
|emp_id| name|superior_emp_id|year_joined|emp_dept_id|gender|salary|
+------+-----+-------------+-----------+-----------+------+------+
|     6|Brown|            2|       2010|         50|  NULL|    -1|
+------+-----+-------------+-----------+-----------+------+------+
```