

Read Text file into PySpark Dataframe.

1) Some file formats other than csv and textfile are:-
parquet, orc, avro etc.

2) `spark.read.format("text").load("output.txt")` OR `spark.read.csv`
or `spark.read.text`.

Process of execution:

→ It converts textfile into dataframe.
When running on the pyspark notebook, it stores as RDD.

Eg:-

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

```
spark = SparkSession.builder.getOrCreate()
```

```
df = spark.read.format("text").load("output.txt")
```

1) Using selectExpr and split method.

```
df = selectExpr("split(Value, '')
```

or.

```
df = selectExpr as Text - Data In Rows - Using format load";  
show(4, false)
```

2) For eg:- calculating avg of age:-

```
df.selectExpr("AVG(age) as avg-age").show()
```


Method to add new Column with Constant Value.

→ Using LIT

• `dataframe.withColumn("column name", lit(Value)).show`

LIT → function to download column.

→ Using Concat

`dataframe.withColumn("column name", concat_ws("-", "Name", "Company"))`
• shows

→ Add Column when not exists on DF.

Syntax: -

if 'column name' not in `dataframe.columns`:

`dataframe.withColumn("column name", lit(Value))`

Manipulating Dataframes.

→ GroupBy and Aggregate Functions.

Group by is used to collect identical data

① Create a dataframe

② Syntax for `groupBy()`:-

`df = pyspark.groupBy("column name").sum("value").show()`

`df = pyspark.groupBy("Department").sum("Salary").show()`

Using PIVOT / UNPIVOT

Used to rotate data from one column into multiple columns.
(transpose row to columns).

Example :- `df = spark.groupBy("Department").pivot("Name").sum("Salary").show()`

Handling missing values

Dropping rows based on null values

`df = spark.na.drop().show()`

Parameters of drop :-

`df = spark.na.drop(how = "all").show()`

`df = spark.na.drop(how = "any", thresh = 2).show()`

`df = spark.na.drop(how = "any", subset = ["Salary"]).show()`

Subset = optional list of column names to consider. Columns specified in subset that do not have matching data type are ignored.

For eg:- If value is a string, non-string value will be ignored.

Order By and Sort.

① `sort()` → default is asc order

```
salary = df_pyspark.sort("salary").show()
          = df_pyspark.sort(df_pyspark["salary"].desc()).show()
          = df_pyspark.sort("salary", "name").show()
```

② `OrderBy()`

system → `df_pyspark.orderBy("salary").show()`

Joins

- ① Inner Join
- ② Outer Join
- ③ left Join / left Outer Join
- ④ Right Join / Right Outer Join
- ⑤ left Semi Join
- ⑥ left Anti Join -

Step 1:- Create two diff. dataframes.