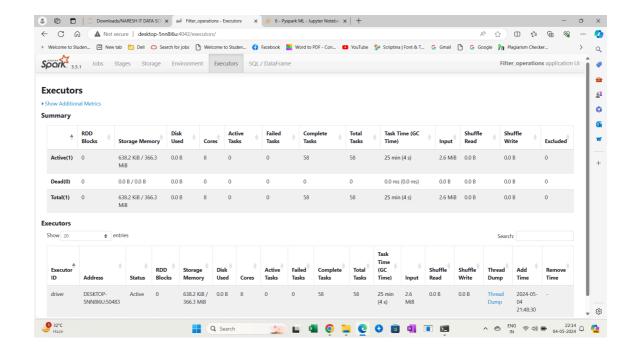
```
In [1]:
          1 from pyspark.sql import SparkSession
           2 spark = SparkSession.builder.appName('Filter_operations').getOrCreate()
           3 spark
 Out[1]: SparkSession - in-memory
         SparkContext
         Spark UI (http://DESKTOP-5NN8I6U:4042)
         Version
         v3.5.1
         Master
         local[*]
         AppName
          Filter_operations
In [10]:
            df=spark.read.csv(r'C:\Users\Neha\Downloads\NARESH IT DATA SCIENCE\MAY\
           2
           3
In [11]:
           1 df
Out[11]: DataFrame[Name: string, age: int, Experience: int, Salary: int]
In [12]:
           1 df.show()
           -----+
              Name | age | Experience | Salary |
            . - - - - - + - - - + - - - - - - - + - - - - - +
                            10| 30000|
              jack| 31|
              alex| 30|
                              8 25000
         |caroline| 29|
                              4 | 20000 |
                              3 20000
              paul| 24|
                              1 | 15000 |
            sandra 21
                              2 | 18000 |
         |casandra| 23|
         +----+
In [13]:
           1 df.printSchema()
         root
          |-- Name: string (nullable = true)
          |-- age: integer (nullable = true)
          |-- Experience: integer (nullable = true)
          |-- Salary: integer (nullable = true)
In [15]:
          1 df.columns
Out[15]: ['Name', 'age', 'Experience', 'Salary']
```

## **Vector Assembler:**

```
In [20]:
           from pyspark.ml.feature import VectorAssembler
           feature_assembler = VectorAssembler(inputCols=['age', 'Experience'],
                                         outputCol = 'Independent Features')
         3
In [21]:
         1 | feature_assembler
Out[21]: VectorAssembler_ed2086a260f0
In [22]:
         1 output= feature_assembler.transform(df)
In [23]:
           output.show()
         -----+---+---+
            Name | age | Experience | Salary | Independent Features |
         -----
            jack 31
                         10 | 30000 |
                                         [31.0,10.0]
                          8| 25000|
4| 20000|
                                          [30.0,8.0]
            alex 30
        |caroline| 29|
                                          [29.0,4.0]
            paul | 24|
                          3 20000
                                          [24.0,3.0]
         sandra| 21|
                          1 | 15000 |
                                          [21.0,1.0]
                                          [23.0,2.0]
        |casandra| 23|
                          2 18000
        +----+
In [24]:
         1 output.columns
Out[24]: ['Name', 'age', 'Experience', 'Salary', 'Independent Features']
In [25]:
           finalised_data= output.select ('Independent Features', 'Salary')
         2
In [26]:
         1 finalised data.show()
        +----+
        |Independent Features|Salary|
        +----+----
                [31.0,10.0]| 30000|
                 [30.0,8.0] | 25000|
                 [29.0,4.0] | 20000|
                 [24.0,3.0]| 20000|
                 [21.0,1.0]| 15000|
                 [23.0,2.0] | 18000 |
        +----+
In [ ]:
```

## **Linear Regression Model:**

```
In [72]:
             from pyspark.ml.regression import LinearRegression
             train_data, test_data = finalised_data.randomSplit([0.75,0.25])
           2
             regressor = LinearRegression(featuresCol='Independent Features', labelC
           3
            regressor = regressor.fit(train_data)
In [73]:
             regressor
Out[73]: LinearRegressionModel: uid=LinearRegression_2d0ec61160d5, numFeatures=2
In [74]:
             pred_results = regressor.evaluate(test_data)
In [75]:
             pred_results.predictions.show()
         |Independent Features|Salary|
                                             prediction|
                    [21.0,1.0] | 15000 | 19285.71428571376 |
                    [23.0,2.0] | 18000 | 18571.428571428507 |
                    [29.0,4.0] | 20000 | 12857.142857144696 |
           -----+
```



```
In [ ]: 1
```