Name : Shailesh A. Tagadghar

Roll No : 31031423034

# Big Data Analytics

# Journal

Submitted By

**Tagadghar Shailesh Tagadghar**

**Roll No: 31031523034**

**MSc CS – Part II**

**Department Of Computer Science**

**Somaiya Vidyavihar University**

**SK Somaiya college**

**S K Somaiya College**

# Index :

**Name : Shailesh A. Tagadghar**

**Roll No : 31031423034**

## Practical 1

<u>**Aim :**</u> Installation of Hadoop in Windows

<u>**Steps :**</u>

1. Download Binary File for Windows. https://hadoop.apache.org/releases.html

Download

Hadoop is released as source code tarballs with corresponding binary tarballs for convenience. The downloads are distributed via mirror sites and should be checked for tampering using GPG or SHA-512.

| Version | Release date | Source download | Binary download | Release notes |
|---------|-------------|-----------------|-----------------|---------------|
| 3.4.0 | 2024 Mar 17 | source (checksum signature) | binary (checksum signature)<br>binary-aarch64 (checksum signature) | Announcement |
| 3.3.6 | 2023 Jun 23 | source (checksum signature) | binary (checksum signature)<br>binary-aarch64 (checksum signature) | Announcement |
| 2.10.2 | 2022 May 31 | source (checksum signature) | binary (checksum signature) | Announcement |

We suggest the following location for your download:

https://dlcdn.apache.org/hadoop/common/hadoop-3.4.0/hadoop-3.4.0.tar.gz

Alternate download locations are suggested below.

It is essential that you verify the integrity of the downloaded file using the PGP signature ( .asc file) or a hash ( .md5 or .sha* file).

## HTTP

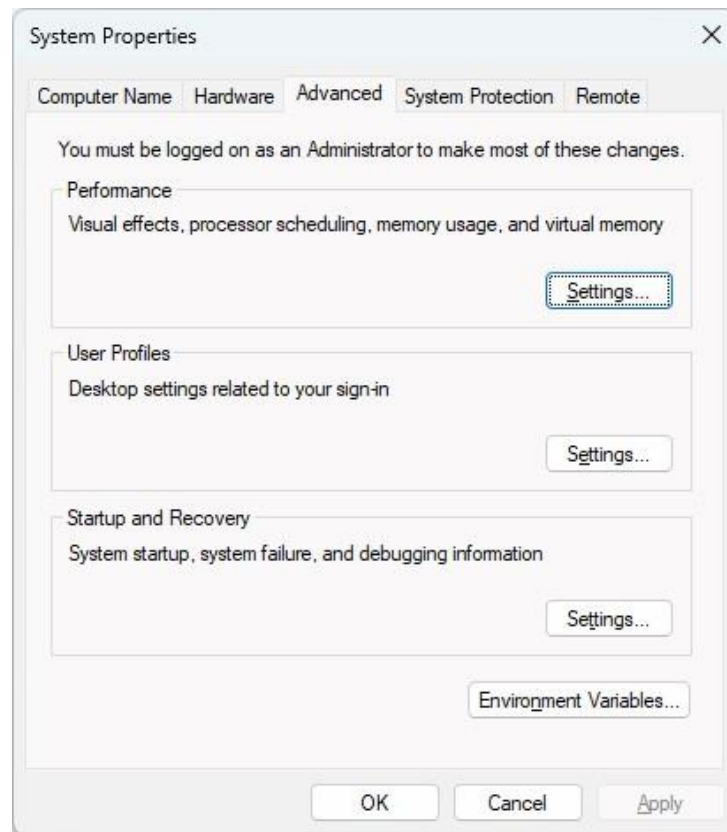https://dlcdn.apache.org/hadoop/common/hadoop-3.4.0/hadoop-3.4.0.tar.gz

https://www.oracle.com/java/technologies/javase-downloads.html

2. Extract the file using Winrar.
3. Go to "Edit Environment Variables" and Click Environment Variables.

**S K Somaiya College**

4. Under System Variables click "New" and set "Variable name" as JAVA_HOME and "Variable value" as the path of your JAVA JDK.

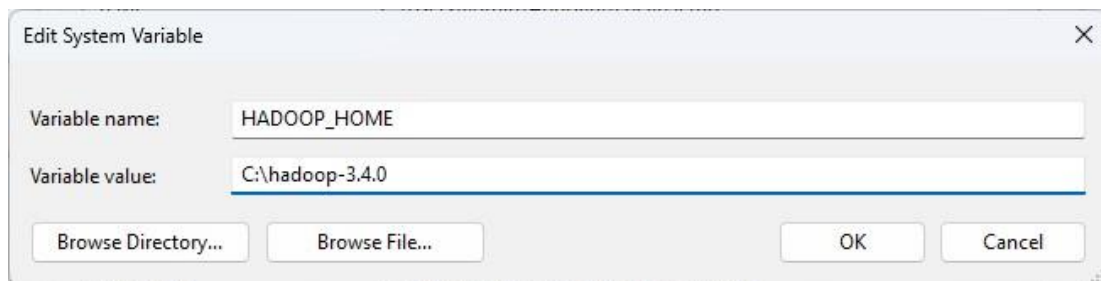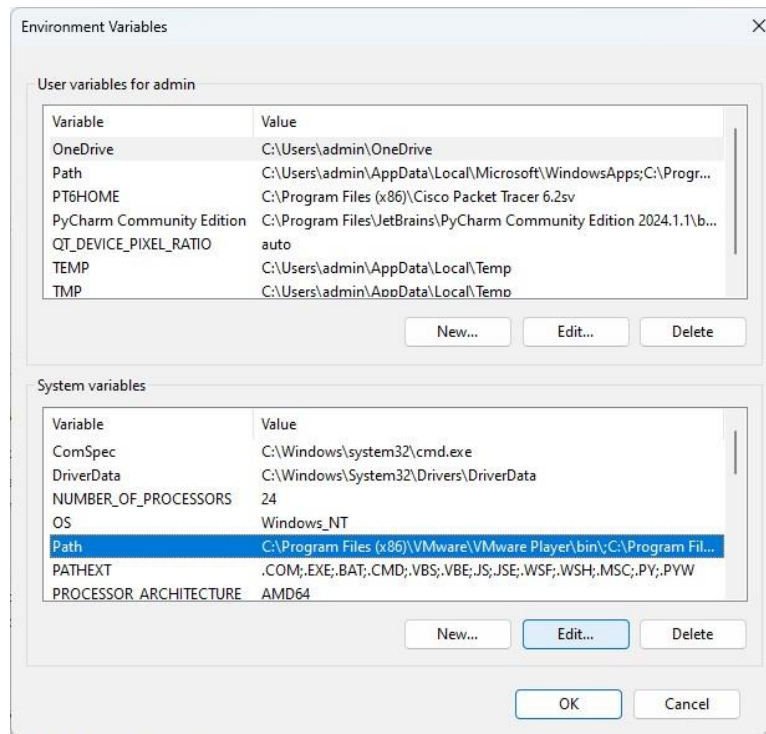C:\Program Files\Java\jdk-21

5. Similarly add "HADOOP_HOME" variable.
6. Download the bin folder from the below link
   https://drive.google.com/drive/folders/1iURNbow2IglhAhSy3sfY5xxVfAg33NBW
7. Extract the bin archive and replace the bin folder in Hadoop folder with the bin folder in this archive.
8. Check if "winutils" is working. If you get any dll error then download that dll and paste in the Windows -> System32 folder.
9. Move the Hadoop folder to C drive.
10. Create a data folder in the hadoop home directory and add the folders datanode and namenode to it.
11. Add the following path to "Path" under "System Variables" in "Edit Environment Variables"

C:\hadoop-3.4.0\sbin



```
%JAVA_HOME%\bin
%HADOOP_HOME%bin
C:\hadoop-3.4.0\bin
C:\hadoop-3.4.0\sbin
```

12. If your PC username has spaces in it then go to hadoop-env.cmd and find this line

   set HADOOP_IDENT_STRING=%USERNAME%

13. Change the above line to your PC username instead of %USERNAME% but WITHOUT SPACES set HADOOP_IDENT_STRING=DeepShah

**S K Somaiya College**

14. Make the changes to the following files as given, in "etc/hadoop" folder of hadoop home.

| core-site.xml | 31-07-2024 08:44 | xmlfile | 1 KB |
|---|---|---|---|
| hadoop-env | 04-03-2024 12:06 | Windows Comma... | 4 KB |
| hadoop-env | 04-03-2024 13:35 | SH Source File | 17 KB |
| hadoop-metrics2 | 04-03-2024 12:06 | Properties Source ... | 4 KB |
| hadoop-policy.xml | 04-03-2024 12:06 | xmlfile | 14 KB |
| hadoop-user-functions.sh.example | 04-03-2024 12:06 | EXAMPLE File | 4 KB |
| hdfs-rbf-site.xml | 04-03-2024 12:37 | xmlfile | 1 KB |
| hdfs-site.xml | 04-03-2024 12:13 | xmlfile | 1 KB |
| httpfs-env | 04-03-2024 12:22 | SH Source File | 2 KB |
| httpfs-log4j | 04-03-2024 12:22 | Properties Source ... | 2 KB |
| httpfs-site.xml | 04-03-2024 12:22 | xmlfile | 1 KB |
| kms-acls.xml | 04-03-2024 12:08 | xmlfile | 4 KB |
| kms-env | 04-03-2024 12:08 | SH Source File | 2 KB |
| kms-log4j | 04-03-2024 12:08 | Properties Source ... | 2 KB |
| kms-site.xml | 04-03-2024 12:08 | xmlfile | 1 KB |
| log4j | 04-03-2024 12:06 | Properties Source ... | 15 KB |
| mapred-env | 04-03-2024 13:00 | Windows Comma... | 1 KB |
| mapred-env | 04-03-2024 13:00 | SH Source File | 2 KB |
| mapred-queues.xml.template | 04-03-2024 13:00 | TEMPLATE File | 5 KB |
| mapred-site.xml | 04-03-2024 13:00 | xmlfile | 1 KB |

core-site.xml

```
<configuration>

        <property>

                <name>fs.default.name</name>

                <value>hdfs://localhost:9000</value>

        </property>
</configuration>
```

mapred-site.xml

```
<configuration>

        <property>

                <name>mapred.framework.name</name>
                        <value>yarn</value>
```

**S K Somaiya College**

```
        </property>

</configuration> hdfs-

site.xml <configuration>

        <property>

                <name>dfs.replication</name>

                <value>1</value>

        </property>

        <property>

                <name>dfs.namenode.name.dir</name>

                <value>C:\hadoop-3.4.0\data\namenode</value>

        </property>

        <property>

                <name>dfs.datanode.data.dir</name>

                <value>C:\hadoop-3.4.0\data\datanode</value>

        </property>
</configuration>


yarn-site.xml
<configuration>

        <property>

                <name>yarn.nodemanager.aux-services</name>

                <value>mapreduce_shuffle</value>

        </property>
        <property>

                <name>yarn.nodemanager.auxservice.mapreduce.shuffle.class</name>
```

<value>org.apache.hadoop.mapred.shuffleHandler</v
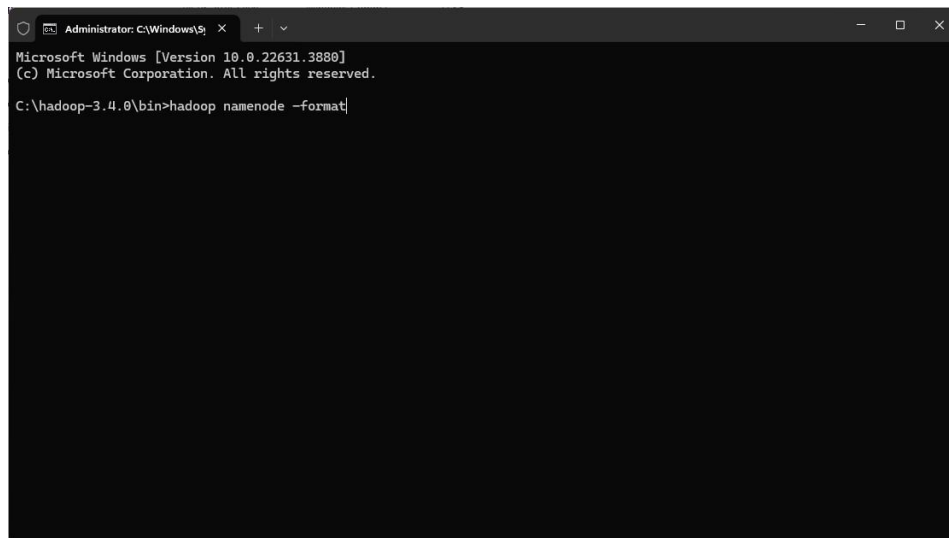alue> </property> </configuration>

15. Go to hadoop-env.cmd file in /etc/hadoop folder and replace the set
    JAVA_HOME=%JAVA_HOME% line with the following:

    set JAVA_HOME=C:\Progra~1\Java\jdk-21
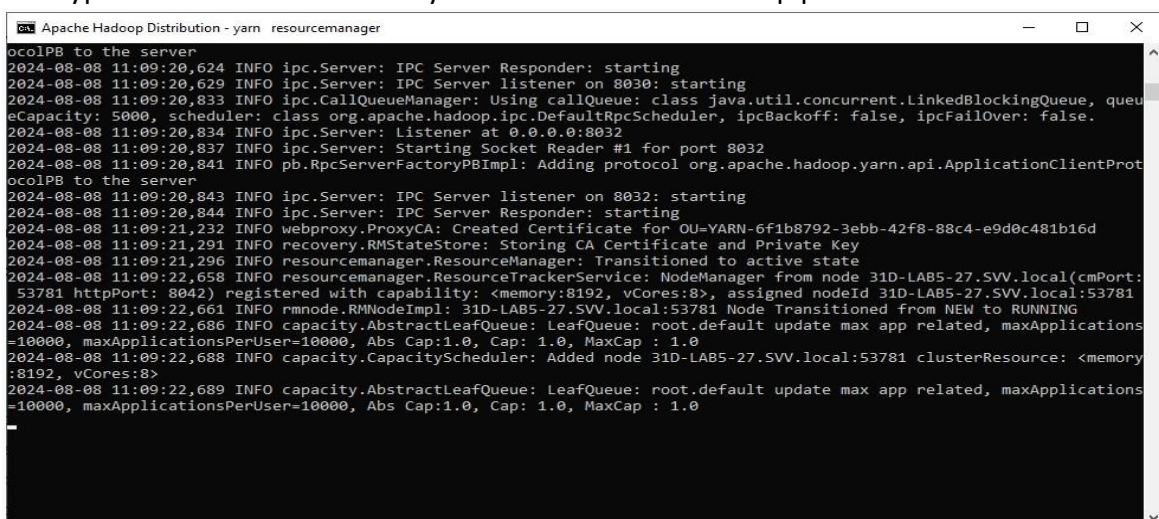
16. Restart your PC for the changes to take effect.
17. Go to Admin Command prompt and type "hadoop" to see if the server is recognized.
18. Type "hdfs namenode -format" to format the namenode.



19. Type start-dfs.cmd and start-yarn.cmd to start all hadoop processes

**S K Somaiya College**

20. Go to your browser and type localhost:9870 to view Hadoop Page.



21. If your yarn process doesnt start then use java version 11.

22. After the yarn process has started visit localhost:8088

# Practical 2

## Aim : Getting started with Scala

## Steps :

1. Get Scala for Windows from the below link Getting Started | Scala Documentation (scala-lang.org)

### Install Scala on your computer

Installing Scala means installing various command-line tools such as the Scala compiler and build tools. We recommend using the Scala installer tool "Coursier" that automatically installs all the requirements, but you can still manually install each tool.

### Using the Scala Installer (recommended way)

The Scala installer is a tool named Coursier, whose main command is named `cs`. It ensures that a JVM and standard Scala tools are installed on your system. Install it on your system with the following instructions.

| macOS | Linux | Windows | Other |
|---|---|---|---|

Download and execute the Scala installer for Windows based on Coursier, and follow the on-screen instructions.

> ℹ You may need to restart your terminal, log out, or reboot in order for the changes to take effect.

Testing your setup ❯

2. Extract the given archive and run the executable file.

```
Checking if a JVM is installed
Found a JVM installed under C:\Program Files\Java\jdk1.8.0_191.

Checking if ~\AppData\Local\Coursier\data\bin is in PATH
  Should we add ~\AppData\Local\Coursier\data\bin to your PATH? [Y/n] |
```

3. Go to the following path and then copy it.

**S K Somaiya College**

4. Set the SCALA_HOME environment variable.



**S K Somaiya College**

5. Open a command prompt on this path and type "scala". It will download the Scala CLI.

C:\Users\admin\AppData\Local\Coursier\data\bin

6. Test the working with a command or prompt 'println("Hello")'.



'Var' before a variable name is used to make the variable mutable.

```
scala> var c: Int = 10
var c: Int = 10

scala> c = c - 5
c: Int = 5
```

```
scala> var a: Int = 10
var a: Int = 10

scala> var b: Int = 20
var b: Int = 20

scala> var c: Int = a + b
var c: Int = 30
```

'Val' is used to make the variable immutable

```
scala> val a: Int = 10
val a: Int = 10

scala> val b: Int = 90
val b: Int = 90

scala> val c: Int = a + b
val c: Int = 100

scala> c = c - 10
-- [E052] Type Error: ----------------------------------------------------------
1 |c = c - 10
  |^^^^^^^^^^
  |Reassignment to val c
  |
  | longer explanation available when compiling with '-explain'
1 error found
```

Installing Spark

Steps

1.  Download Spark from the following link Downloads | Apache Spark

## Download Apache Spark™

1. Choose a Spark release: `3.4.3 (Apr 18 2024)`

2. Choose a package type: `Pre-built for Apache Hadoop 3.3 and later`

3. Download Spark: spark-3.4.3-bin-hadoop3.tgz

4. Verify this release using the 3.4.3 signatures, checksums and project release KEYS by following these procedures.

Note that Spark 3 is pre-built with Scala 2.12 in general and Spark 3.2+ provides additional pre-built distribution with Scala 2.13.

## Link with Spark

Spark artifacts are hosted in Maven Central. You can add a Maven dependency with the following coordinates:

```
groupId: org.apache.spark
artifactId: spark-core_2.12
version: 3.5.2
```

We suggest the following location for your download:

https://dlcdn.apache.org/spark/spark-3.4.3/spark-3.4.3-bin-hadoop3.tgz

Alternate download locations are suggested below.

It is essential that you verify the integrity of the downloaded file using the PGP signature ( `.asc` file) or a hash ( `.md5` or `.sha*` file).

## HTTP

https://dlcdn.apache.org/spark/spark-3.4.3/spark-3.4.3-bin-hadoop3.tgz

## BACKUP SITES

https://dlcdn.apache.org/spark/spark-3.4.3/spark-3.4.3-bin-hadoop3.tgz

2. Extract the files and set the path for its home directory



**S K Somaiya College**

New System Variable

| | |
|---|---|
| Variable name: | SPARK_HOME |
| Variable value: | C:\spark\spark-3.4.3-bin-hadoop3 |

Browse Directory...    Browse File...    OK    Cancel

%SCALA_HOME%
C:\spark\spark-3.4.3-bin-hadoop3\bin

3. Check the installation with 'spark-shell'

**S K Somaiya College**

Lets try to show a sample data from this path

C:\spark\spark-3.4.3-bin-hadoop3\examples\src\main\resources\people.json

```
scala> x.select($"name",$"age").show()
+-------+----+
|   name| age|
+-------+----+
|Michael|null|
|   Andy|  30|
| Justin|  19|
+-------+----+


scala> x.filter($"age">20).show()
+---+----+
|age|name|
+---+----+
| 30|Andy|
+---+----+
```

Reading CSV/Excel File

```
scala> val y = spark.read.csv("C:/spark/spark-3.4.3-bin-hadoop3/examples/src/main/resources/people.csv").show()
+------------------+
|               _c0|
+------------------+
|      name;age;job|
|Jorge;30;Developer|
|  Bob;32;Developer|
+------------------+

y: Unit = ()
```

Creating an SQL Tempory View

```
scala> x.createOrReplaceTempView("people")

scala> val sqlDF = spark.sql("Select * from people")
sqlDF: org.apache.spark.sql.DataFrame = [age: bigint, name: string]

scala> sqlDF.show()
+----+-------+
| age|   name|
+----+-------+
|null|Michael|
|  30|   Andy|
|  19| Justin|
+----+-------+
```

**S K Somaiya College**

Temporary views in Spark SQL are session-scoped and will disappear if the session that creates it terminates. If you want to have a temporary view that is shared among all sessions and keep alive until the Spark application terminates, you can create a global temporary view.

Global temporary view is tied to a system preserved database global_temp, and we must use the qualified name to refer it, e.g. SELECT * FROM global_temp.view.

```
scala> x.createGlobalTempView("people")

scala> spark.sql("SELECT * FROM global_temp.people").show()
+----+-------+
| age|   name|
+----+-------+
|null|Michael|
|  30|   Andy|
|  19| Justin|
+----+-------+


scala> spark.newSession().sql("SELECT * FROM global_temp.people").show()
+----+-------+
| age|   name|
+----+-------+
|null|Michael|
|  30|   Andy|
|  19| Justin|
+----+-------+
```

Creating Datasets

Datasets are similar to RDDs, however, instead of using Java serialization or Kryo they use a specialized Encoder to serialize the objects for processing or transmitting over the network. While both encoders and standard serialization are responsible for turning an object into bytes, encoders are code generated dynamically and use a format that allows Spark to perform many operations like filtering, sorting and hashing without deserializing the bytes back into an object.

case class Person(name: String, age: Long)

// Encoders are created for case classes val
caseClassDS = Seq(Person("Andy", 32)).toDS()
caseClassDS.show()

// +----+---+

// |name|age|

// +----+---+

// |Andy| 32|

**S K Somaiya College**

```
// +----+---+
```

// Encoders for most common types are automatically provided by importing spark.implicits._ val primitiveDS = Seq(1, 2, 3).toDS() primitiveDS.map(_ + 1).collect() // Returns: Array(2, 3, 4)

// DataFrames can be converted to a Dataset by providing a class. Mapping will be done by name

```
val path = "examples/src/main/resources/people.json" val
peopleDS = spark.read.json(path).as[Person]
peopleDS.show()
```

```
// +----+-------+
```

```
// | age|   name|
```

```
// +----+-------+
```

```
// |null|Michael|
```

```
// | 30|   Andy|
```

```
// | 19| Justin|
```

```
// +----+-------+
```

```
scala> case class Person(name: String, age: Long)
defined class Person

scala> val caseClassDS = Seq(Person("Andy", 32)).toDS()
caseClassDS: org.apache.spark.sql.Dataset[Person] = [name: string, age: bigint]

scala> caseClassDS.show()
+----+---+
|name|age|
+----+---+
|Andy| 32|
+----+---+
```

```
scala> val primitiveDS = Seq(1, 2, 3).toDS()
primitiveDS: org.apache.spark.sql.Dataset[Int] = [value: int]

scala> primitiveDS.map(_ + 1).collect()
res6: Array[Int] = Array(2, 3, 4)
```

```
scala> val peopleDS = spark.read.json("C:/spark/spark-3.4.3-bin-hadoop3/examples/src/main/resources/people.json").as[Person]
peopleDS: org.apache.spark.sql.Dataset[Person] = [age: bigint, name: string]

scala> peopleDS.show()
+----+-------+
| age|   name|
+----+-------+
|null|Michael|
|  30|   Andy|
|  19| Justin|
+----+-------+
```

Inferring the Schema Using Reflection

The Scala interface for Spark SQL supports automatically converting an RDD containing case classes to a DataFrame. The case class defines the schema of the table. The names of the arguments to the case class are read using reflection and become the names of the columns.

Case classes can also be nested or contain complex types such as Seqs or Arrays. This RDD can be implicitly converted to a DataFrame and then be registered as a table. Tables can be used in subsequent SQL statements.

// For implicit conversions from RDDs to DataFrames import
spark.implicits._

// Create an RDD of Person objects from a text file, convert it to a Dataframe val
peopleDF = spark.sparkContext

  .textFile("examples/src/main/resources/people.txt")
  .map(_.split(","))

```
  .map(attributes => Person(attributes(0), attributes(1).trim.toInt)) .toDF()
```

// Register the DataFrame as a temporary view peopleDF.createOrReplaceTempView("people")

// SQL statements can be run by using the sql methods provided by Spark

val teenagersDF = spark.sql("SELECT name, age FROM people WHERE age BETWEEN 13 AND 19")

// The columns of a row in the result can be accessed by field index teenagersDF.map(teenager => "Name: " + teenager(0)).show()

// +------------+

// |      value|

// +------------+

// |Name: Justin|

// +------------+

// or by field name teenagersDF.map(teenager => "Name: " + teenager.getAs[String]("name")).show()

// +------------+

// |      value|

// +------------+

// |Name: Justin|

// +------------+

// No pre-defined encoders for Dataset[Map[K,V]], define explicitly

implicit val mapEncoder = org.apache.spark.sql.Encoders.kryo[Map[String, Any]]

// Primitive types and case classes can be also defined as

// implicit val stringIntMapEncoder: Encoder[Map[String, Any]] = ExpressionEncoder()

// row.getValuesMap[T] retrieves multiple columns at once into a Map[String, T]

teenagersDF.map(teenager => teenager.getValuesMap[Any](List("name", "age"))).collect()

// Array(Map("name" -> "Justin", "age" -> 19))

**S K Somaiya College**

```
scala> case class Person(name: String, age: Long)
defined class Person

scala> val peopleDF = spark.sparkContext.textFile("C:/spark/spark-3.4.3-bin-hadoop3/examples/src/main/resources/people.t
xt").map(_.split(",")).map(attributes => Person(attributes(0), attributes(1).trim.toInt)).toDF()
peopleDF: org.apache.spark.sql.DataFrame = [name: string, age: bigint]

scala> peopleDF.createOrReplaceTempView("people")

scala> val teenagersDF = spark.sql("SELECT name, age FROM people WHERE age BETWEEN 13 AND 19")
teenagersDF: org.apache.spark.sql.DataFrame = [name: string, age: bigint]

scala> teenagersDF.map(teenager => "Name: " + teenager(0)).show()
+------------+
|       value|
+------------+
|Name: Justin|
+------------+
```

```
scala> teenagersDF.map(teenager => "Name: " + teenager.getAs[String]("name")).show()
+------------+
|       value|
+------------+
|Name: Justin|
+------------+


scala> implicit val mapEncoder = org.apache.spark.sql.Encoders.kryo[Map[String, Any]]
     |
     | ]
mapEncoder: org.apache.spark.sql.Encoder[Map[String,Any]] = class[value[0]: binary]

scala> teenagersDF.map(teenager => teenager.getValuesMap[Any](List("name", "age"))).collect()
res3: Array[Map[String,Any]] = Array(Map(name -> Justin, age -> 19))
```

Programmatically Specifying the Schema

When case classes cannot be defined ahead of time (for example, the structure of records is encoded
in a string, or a text dataset will be parsed and fields will be projected differently for different users),
a DataFrame can be created programmatically with three steps.

Create an RDD of Rows from the original RDD;

Create the schema represented by a StructType matching the structure of Rows in the RDD created in
Step 1.

Apply the schema to the RDD of Rows via createDataFrame method provided by SparkSession.

import org.apache.spark.sql.Row import org.apache.spark.sql.types._

// Create an RDD

val peopleRDD = spark.sparkContext.textFile("examples/src/main/resources/people.txt")

// The schema is encoded in a string val
schemaString = "name age"

**S K Somaiya College**

```
// Generate the schema based on the string of schema val
fields = schemaString.split(" ")

  .map(fieldName => StructField(fieldName, StringType, nullable = true))



schema = StructType(fields)

// Convert records of the RDD (people) to Rows val
rowRDD = peopleRDD

  .map(_.split(","))

  .map(attributes => Row(attributes(0), attributes(1).trim))

// Apply the schema to the RDD val peopleDF =
spark.createDataFrame(rowRDD, schema)

// Creates a temporary view using the DataFrame peopleDF.createOrReplaceTempView("people")

// SQL can be run over a temporary view created using DataFrames val
results = spark.sql("SELECT name FROM people")

// The results of SQL queries are DataFrames and support all the normal RDD operations // The
columns of a row in the result can be accessed by field index or by field name
results.map(attributes => "Name: " + attributes(0)).show()

// +-------------+

// |       value|

// +-------------+

// |Name: Michael|

// |   Name: Andy|

// | Name: Justin|

// +-------------+

import org.apache.spark.sql.Row
Import org.apache.spark.sql.types._
```

**S K Somaiya College**

```
scala> import org.apache.spark.sql.Row
import org.apache.spark.sql.Row

scala>

scala> import org.apache.spark.sql.types._
import org.apache.spark.sql.types._
```

val peopleRDD =

spark.sparkContext.textFile("C:/spark/spark-3.4.3-bin/hadoop/examples/src/main/resources/peo ple.txt")

```
scala> val peopleRDD = spark.sparkContext.textFile("eC:/spark/spark-3.4.3-bin-hadoop3/examples/src/main/resources/people
.txt")
peopleRDD: org.apache.spark.rdd.RDD[String] = eC:/spark/spark-3.4.3-bin-hadoop3/examples/src/main/resources/people.txt M
apPartitionsRDD[14] at textFile at <console>:27
```

fields = schemaString.split(" ").map(fieldName => StructField(fieldName, StringType, nullable = true))

```
scala> val fields = schemaString.split(" ").map(fieldName => StructField(fieldName, StringType, nullab
le = true))
fields: Array[org.apache.spark.sql.types.StructField] = Array(StructField(name,StringType,true), Struc
tField(age,StringType,true))
```

val schema = StructType(fields)

```
scala> val schema = StructType(fields)
schema: org.apache.spark.sql.types.StructType = StructType(StructField(name,StringType,true),StructFie
ld(age,StringType,true))
```

val rowRDD = peopleRDD.map(_.split(",")).map(attributed => Row(attributes(0), attributes(1).trim))

```
scala> val rowRDD = peopleRDD.map(_.split(",")).map(attributes => Row(attributes(0), attributes(1).tri
m))
rowRDD: org.apache.spark.rdd.RDD[org.apache.spark.sql.Row] = MapPartitionsRDD[3] at map at <console>:2
7
```

val peopleDF = spark.createDataFrame(rowRDD, schema)

```
scala> val peopleDF = spark.createDataFrame(rowRDD, schema)
peopleDF: org.apache.spark.sql.DataFrame = [name: string, age: string]
```

peopleDF.createOrReplaceTempView("people")

```
scala> peopleDF.createOrReplaceTempView("people")
```

val results = spark.sql("SELECT name FROM people")

**S K Somaiya College**

```
scala> val results = spark.sql("SELECT name FROM people")
results: org.apache.spark.sql.DataFrame = [name: string]
```

results.map(attributes => "Name: " + attributes(0)).show()

```
scala> results.map(attributes => "Name: " + attributes(0)).show()
+--------------+
|         value|
+--------------+
|Name: Michael|
|   Name: Andy|
| Name: Justin|
+--------------+
```

Basic Operations with csv file

    myData =
spark.read.format("csv").option("inferSchema","true").option("header","true").option("delimeter",":

").load("C:/spark/spark-3.4.3-hadoop3/examples/src/main/resources/people.csv")

```
scala> val myData = spark.read.format("csv").option("inferSchema", "true").option("header", "true").option("delimiter", ";").load("C:/spark/spark-3.4.3-bin-hadoop3/examples/src/main/resources/people.csv")
myData: org.apache.spark.sql.DataFrame = [name: string, age: int ... 1 more field]
```

myData.show()

```
scala> myData.show()
+--------------------+
|        name;age;job|
+--------------------+
|Jorge;30;Developer|
|   Bob;32;Developer|
+--------------------+
```

myData.select($"name","$age").show()

```
scala> myData.select($"name",$"age").show()
+-----+---+
| name|age|
+-----+---+
|Jorge| 30|
|  Bob| 32|
+-----+---+
```

myData.count()

```
scala> myData.count()
res8: Long = 2
```

**S K Somaiya College**

myData.count().toDouble

```scala
scala> myData.count().toDouble
res10: Double = 2.0
```

# Practical 03

## Aim : GraphX

```
import
org.apache.spark._
import
org.apache.spark.rdd.RDD
import
org.apache.spark.graphx.

_
```

```
scala> import org.apache.spark._
import org.apache.spark._

scala> import org.apache.spark.rdd.RDD
import org.apache.spark.rdd.RDD

scala> import org.apache.spark.graphx._
import org.apache.spark.graphx._
```

val vertices = Array((1L,("A")),(2L,("B")),(3L,("C")))

```
scala> val vertices = Array((1L,("A")),(2L,("B")),(3L,("C")))
vertices: Array[(Long, String)] = Array((1,A), (2,B), (3,C))
```

val vRDD = sc.parallelize(vertices)

```
scala> val vRDD = sc.parallelize(vertices)
vRDD: org.apache.spark.rdd.RDD[(Long, String)] = ParallelCollectionRDD[0] at parallelize at <console>:31
```

```
vRDD.take(1)
vRDD.take(2)
vRDD.take(3)
```

```
scala> vRDD.take(1)
res0: Array[(Long, String)] = Array((1,A))

scala> vRDD.take(2)
res1: Array[(Long, String)] = Array((1,A), (2,B))

scala> vRDD.take(3)
res2: Array[(Long, String)] = Array((1,A), (2,B), (3,C))
```

val edges = Array(Edge(1L,2L,1800),Edge(2L,3L,800),Edge(3L,1L,1400))

```
scala> val edges = Array(Edge(1L,2L,1800),Edge(2L,3L,800),Edge(3L,1L,1400))
edges: Array[org.apache.spark.graphx.Edge[Int]] = Array(Edge(1,2,1800), Edge(2,3,800), Edge(3,1,1400))
```

val eRDD = sc.parallelize(edges)

```
scala> val eRDD = sc.parallelize(edges)
eRDD: org.apache.spark.rdd.RDD[org.apache.spark.graphx.Edge[Int]] = ParallelCollectionRDD[1] at parallelize at <console>
:31
```

eRDD.take(2)

```
scala> eRDD.take(2)
res3: Array[org.apache.spark.graphx.Edge[Int]] = Array(Edge(1,2,1800), Edge(2,3,800))
```

val nowhere = "nowhere"

```
scala> val nowhere = "nowhere"
nowhere: String = nowhere
```

val graph = Graph(vRDD,eRDD,nowhere)

```
scala> val graph = Graph(vRDD,eRDD,nowhere)
graph: org.apache.spark.graphx.Graph[String,Int] = org.apache.spark.graphx.impl.GraphImpl@3e1e7aa2
```

#To check number of Airports
val numairports =
graph.numVertices

```
scala> val numairports = graph.numVertices
numairports: Long = 3
```

#To check routes val
numairports =
graph.numEdges

**S K Somaiya College**

```
scala> val numairports = graph.numEdges
numairports: Long = 3
```

#Route having distance > 1000

(graph.edges.filter{case
Edge(src,dst,prop)=>prop>1000}.collect.foreach(println))

```
scala> (graph.edges.filter{case Edge(src,dst,prop)=>prop>1000}.collect.foreach(println))
Edge(1,2,1800)
Edge(3,1,1400)
```

#Triplet Information
graph.triplets.take(3).foreach(println)

```
scala> graph.triplets.take(3).foreach(println)
((1,A),(2,B),1800)
((2,B),(3,C),800)
((3,C),(1,A),1400)
```

#Indegree val i =
graph.inDegrees

```
scala> val i = graph.inDegrees
i: org.apache.spark.graphx.VertexRDD[Int] = VertexRDDImpl[25] at RDD at VertexRDD.scala:57
```

i.collect()

```
scala> i.collect()
res6: Array[(org.apache.spark.graphx.VertexId, Int)] = Array((1,1), (2,1), (3,1))
```

#Outdegrees val o =
graph.outDegrees

```
scala> val o = graph.outDegrees
o: org.apache.spark.graphx.VertexRDD[Int] = VertexRDDImpl[29] at RDD at VertexRDD.scala:57
```

o.collect()

```
scala> o.collect()
res7: Array[(org.apache.spark.graphx.VertexId, Int)] = Array((1,1), (2,1), (3,1))
```

#Total Degree val
t = graph.degrees

```
scala> val t = graph.degrees
t: org.apache.spark.graphx.VertexRDD[Int] = VertexRDDImpl[33] at RDD at VertexRDD.scala:57
```

t.collect()

```
scala> t.collect()
res8: Array[(org.apache.spark.graphx.VertexId, Int)] = Array((1,2), (2,2), (3,2))
```

# Practical 04

**Aim :** PySpark

**Steps :**

1. Make a CSV file with data related to Name Age Salary and Experience.

| | A | B | C | D |
|---|---|---|---|---|
| 1 | Name | Age | Salary (INR) | Experience |
| 2 | Aditi Sharma | 28 | 5,50,000 | 1 |
| 3 | Raj Patel | 34 | 7,50,000 | 3 |
| 4 | Neha Gupta | 26 | 4,80,000 | 2 |
| 5 | Vikram Singh | 40 | 12,00,000 | 4 |
| 6 | Priya Rao | 30 | 6,20,000 | 1 |
| 7 | Anil Kumar | 45 | 15,00,000 | 5 |
| 8 | Kavita Joshi | 29 | 5,80,000 | 2 |
| 9 | Rohan Mehta | 32 | 7,00,000 | 3 |
| 10 | Sneha Desai | 27 | 5,20,000 | 1 |
| 11 | Amit Verma | 38 | 9,00,000 | 2 |

2. Code

```
from pyspark.sql import SparkSession

# Create a Spark session spark =
SparkSession.builder.appName("Read CSV").getOrCreate()

# Path to your CSV file
csv_file_path = "Student.csv"
```

**S K Somaiya College**

# Read the CSV file into a DataFrame. inferSchema tries to determine the datatype of values in the fields.

```
df = spark.read.csv(csv_file_path, header=True,
                                    inferSchema=True)
```

# Display the DataFrame
df.show()

Output

```
+-------------+---+-------------+----------+
|         Name|Age|Salary (INR)|Experience|
+-------------+---+-------------+----------+
| Aditi Sharma| 28|     5,50,000|         1|
|    Raj Patel| 34|     7,50,000|         3|
|   Neha Gupta| 26|     4,80,000|         2|
| Vikram Singh| 40|    12,00,000|         4|
|    Priya Rao| 30|     6,20,000|         1|
|   Anil Kumar| 45|    15,00,000|         5|
| Kavita Joshi| 29|     5,80,000|         2|
|  Rohan Mehta| 32|     7,00,000|         3|
|  Sneha Desai| 27|     5,20,000|         1|
|   Amit Verma| 38|     9,00,000|         2|
+-------------+---+-------------+----------+
```

type(df)



```
type(df)

pyspark.sql.dataframe.DataFrame
def __init__(jdf: JavaObject, sql_ctx: Union['SQLContext', 'SparkSession'])

/usr/local/lib/python3.10/dist-packages/pyspark/sql/dataframe.py
A distributed collection of data grouped into named columns.

.. versionadded:: 1.3.0

.. versionchanged:: 3.4.0
```

df.printSchema()

```
df.printSchema()
```

```
root
 |-- Name: string (nullable = true)
 |-- Age: integer (nullable = true)
 |-- Salary (INR): string (nullable = true)
 |-- Experience: integer (nullable = true)
```

df.head(5) or df.show(5)

```
df.head(5)
```

```
[Row(Name='Aditi Sharma', Age=28, Salary (INR)='5,50,000', Experience=1),
 Row(Name='Raj Patel', Age=34, Salary (INR)='7,50,000', Experience=3),
 Row(Name='Neha Gupta', Age=26, Salary (INR)='4,80,000', Experience=2),
 Row(Name='Vikram Singh', Age=40, Salary (INR)='12,00,000', Experience=4),
 Row(Name='Priya Rao', Age=30, Salary (INR)='6,20,000', Experience=1)]
```

df.columns

```
df.columns
```

```
['Name', 'Age', 'Salary (INR)', 'Experience']
```

df.select('Name').show()

```
df.select('Name').show()
```

```
+------------+
|        Name|
+------------+
|Aditi Sharma|
|   Raj Patel|
|  Neha Gupta|
|Vikram Singh|
|   Priya Rao|
|  Anil Kumar|
|Kavita Joshi|
| Rohan Mehta|
| Sneha Desai|
|  Amit Verma|
+------------+
```

df.select(['Name','Experience']).show()

**S K Somaiya College**

```
df.select(['Name','Experience']).show()
```

```
+------------+----------+
|        Name|Experience|
+------------+----------+
|Aditi Sharma|         1|
|   Raj Patel|         3|
|  Neha Gupta|         2|
|Vikram Singh|         4|
|   Priya Rao|         1|
|  Anil Kumar|         5|
|Kavita Joshi|         2|
| Rohan Mehta|         3|
| Sneha Desai|         1|
|  Amit Verma|         2|
+------------+----------+
```

df.dtypes

```
df.dtypes
```

```
[('Name', 'string'),
 ('Age', 'int'),
 ('Salary (INR)', 'string'),
 ('Experience', 'int')]
```

df.describe().show()

```
+-------+------------+-----------------+------------+------------------+
|summary|        Name|              Age|Salary (INR)|        Experience|
+-------+------------+-----------------+------------+------------------+
|  count|          10|               10|          10|                10|
|   mean|        NULL|             32.9|        NULL|               2.4|
| stddev|        NULL|6.279596590015423|        NULL|1.3498971154211057|
|    min|Aditi Sharma|               26|  12,00,000|                 1|
|    max|Vikram Singh|               45|   9,00,000|                 5|
+-------+------------+-----------------+------------+------------------+
```

Adding columns to the dataframe df = df.withColumn('Experience after 2 years',df['Experience']+2)

```
df.withColumn('Experience after 2 years',df['Experience']+2)
```

```
DataFrame[Name: string, Age: int, Salary (INR): string, Experience: int, Experience after 2 years: int]
```

**S K Somaiya College**

df.show()

```
+------------+---+-----------+----------+-----------------------+
|        Name|Age|Salary (INR)|Experience|Experience after 2 years|
+------------+---+-----------+----------+-----------------------+
|Aditi Sharma| 28|   5,50,000|         1|                      3|
|   Raj Patel| 34|   7,50,000|         3|                      5|
|  Neha Gupta| 26|   4,80,000|         2|                      4|
|Vikram Singh| 40|  12,00,000|         4|                      6|
|   Priya Rao| 30|   6,20,000|         1|                      3|
|  Anil Kumar| 45|  15,00,000|         5|                      7|
|Kavita Joshi| 29|   5,80,000|         2|                      4|
| Rohan Mehta| 32|   7,00,000|         3|                      5|
| Sneha Desai| 27|   5,20,000|         1|                      3|
|  Amit Verma| 38|   9,00,000|         2|                      4|
+------------+---+-----------+----------+-----------------------+
```

Dropping columns from dataframe
df.drop('Experience after 2 years').show()

```
df.drop('Experience after 2 years').show()

+------------+---+-----------+----------+
|        Name|Age|Salary (INR)|Experience|
+------------+---+-----------+----------+
|Aditi Sharma| 28|   5,50,000|         1|
|   Raj Patel| 34|   7,50,000|         3|
|  Neha Gupta| 26|   4,80,000|         2|
|Vikram Singh| 40|  12,00,000|         4|
|   Priya Rao| 30|   6,20,000|         1|
|  Anil Kumar| 45|  15,00,000|         5|
|Kavita Joshi| 29|   5,80,000|         2|
| Rohan Mehta| 32|   7,00,000|         3|
| Sneha Desai| 27|   5,20,000|         1|
|  Amit Verma| 38|   9,00,000|         2|
+------------+---+-----------+----------+
```

df.withColumnRenamed('Name','New Name').show()

```
df.withColumnRenamed('Name','New Name').show()
```

```
+------------+---+-----------+----------+----------------------+
|    New Name|Age|Salary (INR)|Experience|Experience after 2 years|
+------------+---+-----------+----------+----------------------+
|Aditi Sharma| 28|   5,50,000|         1|                     3|
|   Raj Patel| 34|   7,50,000|         3|                     5|
|  Neha Gupta| 26|   4,80,000|         2|                     4|
|Vikram Singh| 40|  12,00,000|         4|                     6|
|   Priya Rao| 30|   6,20,000|         1|                     3|
|  Anil Kumar| 45|  15,00,000|         5|                     7|
|Kavita Joshi| 29|   5,80,000|         2|                     4|
| Rohan Mehta| 32|   7,00,000|         3|                     5|
| Sneha Desai| 27|   5,20,000|         1|                     3|
|  Amit Verma| 38|   9,00,000|         2|                     4|
+------------+---+-----------+----------+----------------------+
```

With new data that has null values

| Name | Age | Salary (INR) | Experience |
|---|---|---|---|
| Aditi Sharma | 28 | 5,50,000 | 1 |
| Raj Patel | 34 | 7,50,000 | 3 |
| Neha Gupta | 26 | 4,80,000 | 2 |
| Vikram Singh | 40 | 12,00,000 | 4 |
| Priya Rao | 30 | 6,20,000 | 1 |
| Anil Kumar | 45 | 15,00,000 | 5 |
| Kavita Joshi | 29 | 5,80,000 | 2 |
| Rohan Mehta | 32 | 7,00,000 | 3 |
| Sneha Desai | 27 | 5,20,000 | 1 |
| Amit Verma | 38 | 9,00,000 | 2 |
| Sumit Milind | | 1,00,000 | |
| Sunita Shinde | | | 3 |

```
df.show()
```

```
+-------------+----+------------+----------+
|     New Name| Age|Salary (INR)|Experience|
+-------------+----+------------+----------+
| Aditi Sharma|  28|    5,50,000|         1|
|    Raj Patel|  34|    7,50,000|         3|
|   Neha Gupta|  26|    4,80,000|         2|
| Vikram Singh|  40|   12,00,000|         4|
|    Priya Rao|  30|    6,20,000|         1|
|   Anil Kumar|  45|   15,00,000|         5|
| Kavita Joshi|  29|    5,80,000|         2|
|  Rohan Mehta|  32|    7,00,000|         3|
|  Sneha Desai|  27|    5,20,000|         1|
|   Amit Verma|  38|    9,00,000|         2|
| Sumit Milind|NULL|    1,00,000|      NULL|
|Sunita Shinde|NULL|        NULL|         3|
+-------------+----+------------+----------+
```

**S K Somaiya College**

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

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

```
+------------+---+------------+----------+
|   New Name|Age|Salary (INR)|Experience|
+------------+---+------------+----------+
|Aditi Sharma| 28|    5,50,000|         1|
|   Raj Patel| 34|    7,50,000|         3|
|  Neha Gupta| 26|    4,80,000|         2|
|Vikram Singh| 40|   12,00,000|         4|
|   Priya Rao| 30|    6,20,000|         1|
|  Anil Kumar| 45|   15,00,000|         5|
|Kavita Joshi| 29|    5,80,000|         2|
| Rohan Mehta| 32|    7,00,000|         3|
| Sneha Desai| 27|    5,20,000|         1|
|  Amit Verma| 38|    9,00,000|         2|
+------------+---+------------+----------+
```

#Drops entries that have all columns as null df =
df.na.drop(how="all") df.show()

Here none will be dropped

```
+-------------+----+------------+----------+
|         Name| Age|Salary (INR)|Experience|
+-------------+----+------------+----------+
| Aditi Sharma|  28|    5,50,000|         1|
|    Raj Patel|  34|    7,50,000|         3|
|   Neha Gupta|  26|    4,80,000|         2|
| Vikram Singh|  40|   12,00,000|         4|
|    Priya Rao|  30|    6,20,000|         1|
|   Anil Kumar|  45|   15,00,000|         5|
| Kavita Joshi|  29|    5,80,000|         2|
|  Rohan Mehta|  32|    7,00,000|         3|
|  Sneha Desai|  27|    5,20,000|         1|
|   Amit Verma|  38|    9,00,000|         2|
| Sumit Milind|NULL|    1,00,000|      NULL|
|Sunita Shinde|NULL|        NULL|         3|
+-------------+----+------------+----------+
```

#Only drop the rows that have given threshold number of NULL COLUMNS df
= df.na.drop(how="any", thresh = 2) df.show()

| Name | Age | Salary (INR) | Experience |
|---|---|---|---|
| Aditi Sharma | 28 | 5,50,000 | 1 |
| Raj Patel | 34 | 7,50,000 | 3 |
| Neha Gupta | 26 | 4,80,000 | 2 |
| Vikram Singh | 40 | 12,00,000 | 4 |
| Priya Rao | 30 | 6,20,000 | 1 |
| Anil Kumar | 45 | 15,00,000 | 5 |
| Kavita Joshi | 29 | 5,80,000 | 2 |
| Rohan Mehta | 32 | 7,00,000 | 3 |
| Sneha Desai | 27 | 5,20,000 | 1 |
| Amit Verma | 38 | 9,00,000 | 2 |
| Sunita Shinde | NULL | NULL | 3 |

| Name | Age | Salary (INR) | Experience |
|---|---|---|---|
| Aditi Sharma | 28 | 5,50,000 | 1 |
| Raj Patel | 34 | 7,50,000 | 3 |
| Neha Gupta | 26 | 4,80,000 | 2 |
| Vikram Singh | 40 | 12,00,000 | 4 |
| Priya Rao | 30 | 6,20,000 | 1 |
| Anil Kumar | 45 | 15,00,000 | 5 |
| Kavita Joshi | 29 | 5,80,000 | 2 |
| Rohan Mehta | 32 | 7,00,000 | 3 |
| Sneha Desai | 27 | 5,20,000 | 1 |
| Amit Verma | 38 | 9,00,000 | 2 |
| Sumit Milind | NULL | 1,00,000 | NULL |
| Sunita Shinde | NULL | NULL | 3 |

```
# Drops rows that has Experience as NULL df =
df.na.drop(how = "any",subset = ['Experience'])
df.show()
```

```
df = df.na.fill('Missing',['Age','Experience']).show()
```

**S K Somaiya College**

```
+-------------+-------+-----------+----------+
|         Name|    Age|Salary (INR)|Experience|
+-------------+-------+-----------+----------+
| Aditi Sharma|     28|   5,50,000|         1|
|    Raj Patel|     34|   7,50,000|         3|
|   Neha Gupta|     26|   4,80,000|         2|
| Vikram Singh|     40|  12,00,000|         4|
|    Priya Rao|     30|   6,20,000|         1|
|   Anil Kumar|     45|  15,00,000|         5|
| Kavita Joshi|     29|   5,80,000|         2|
|  Rohan Mehta|     32|   7,00,000|         3|
|  Sneha Desai|     27|   5,20,000|         1|
|   Amit Verma|     38|   9,00,000|         2|
| Sumit Milind|Missing|   1,00,000|   Missing|
|Sunita Shinde|Missing|       NULL|         3|
|   Pankaj Rao|Missing|       NULL|   Missing|
|         NULL|     30|   1,00,000|         1|
+-------------+-------+-----------+----------+
```

#Fills in the said text in the string columns that are NULL df = df.na.fill("Missing") df.show()

```
+-------------+-------+-----------+----------+
|         Name|    Age|Salary (INR)|Experience|
+-------------+-------+-----------+----------+
| Aditi Sharma|     28|   5,50,000|         1|
|    Raj Patel|     34|   7,50,000|         3|
|   Neha Gupta|     26|   4,80,000|         2|
| Vikram Singh|     40|  12,00,000|         4|
|    Priya Rao|     30|   6,20,000|         1|
|   Anil Kumar|     45|  15,00,000|         5|
| Kavita Joshi|     29|   5,80,000|         2|
|  Rohan Mehta|     32|   7,00,000|         3|
|  Sneha Desai|     27|   5,20,000|         1|
|   Amit Verma|     38|   9,00,000|         2|
| Sumit Milind|Missing|   1,00,000|   Missing|
|Sunita Shinde|Missing|    Missing|         3|
|   Pankaj Rao|Missing|    Missing|   Missing|
+-------------+-------+-----------+----------+
```

from pyspark.ml.feature import Imputer imputer = Imputer(

inputCols = ['Age','Experience'], outputCols = ["{}_imputed".format(c) for c in

['Age','Experience']]).setStrategy("mean")

imputer.fit(df).transform(df).show()

```
+-------------+----+-------------+----------+-----------+------------------+
|         Name| Age|Salary (INR)|Experience|Age_imputed|Experience_imputed|
+-------------+----+-------------+----------+-----------+------------------+
| Aditi Sharma|  28|     5,50,000|         1|         28|                 1|
|    Raj Patel|  34|     7,50,000|         3|         34|                 3|
|   Neha Gupta|  26|     4,80,000|         2|         26|                 2|
| Vikram Singh|  40|    12,00,000|         4|         40|                 4|
|    Priya Rao|  30|     6,20,000|         1|         30|                 1|
|   Anil Kumar|  45|    15,00,000|         5|         45|                 5|
|  Kavita Joshi|  29|     5,80,000|         2|         29|                 2|
|   Rohan Mehta|  32|     7,00,000|         3|         32|                 3|
|   Sneha Desai|  27|     5,20,000|         1|         27|                 1|
|    Amit Verma|  38|     9,00,000|         2|         38|                 2|
| Sumit Milind|NULL|     1,00,000|      NULL|         32|                 2|
|Sunita Shinde|NULL|         NULL|         3|         32|                 3|
|   Pankaj Rao|NULL|         NULL|      NULL|         32|                 2|
|         NULL|  30|     1,00,000|         1|         30|                 1|
+-------------+----+-------------+----------+-----------+------------------+
```