

## Hadoop Distributed File System

- 1) YARN - Yet Another Resource Negotiator
  - Operating System for HDFS

### Commands :-

- 1) jps - To check which are the process running.
- 2) start-dfs.sh && start-yarn.sh
  - 2 sessions to be started .dfs and yarn
- 3) Create directory in HDFS
 

```
hdfs dfs -mkdir /user/test_dir
```
- 4) Check if directory created
 

```
hdfs dfs -lsR /user/
```
- 5) Put files to HDFS
 

```
hdfs dfs -put /home/sois/lab2.txt /user/test_dir
```
- 6) hdfs dfs -ls /user/
- 7) Put multiple files
 

```
hdfs dfs -put t1.txt t2.txt /user/test_dir
```
- 8) Copying files from HDFS
 

```
hdfs dfs -get /user/test_dir/cab2.txt /home/sois/Desktop
```
- 9) hdfs dfs -cat /user/test\_dir/lab2.txt
  - To view content
- 10) To view in web
  - localhost:9780/
  - Go to utilities
  - Browse File System
- 11) To copy files within HDFS
 

```
hdfs dfs -cp /user/test_dir/lab2.txt /user/test_dir/cab3.txt
```
- 12) Move files within HDFS
 

```
hdfs dfs -mv /user/test_dir/lab2.txt /user/test_dir/cab4.txt
```
- 13) Delete in HDFS
 

```
hdfs dfs -rm /user/test_dir/lab3.txt
```

23/08/2023

Sqoop :

1) Transferring structured data from RDBMS to HDFS

1) sudo mysql -u root -P

2) use retail\_db;

3) show tables;

Step 1:-

\* Grant permission:-

GRANT ALL PRIVILEGES ON \*.\* TO 'msis'@'localhost'  
IDENTIFIED BY 'bda23'

Step 2:-

Start Hadoop.

Step 3:-

# List the DBs

sqoop list-databases --connect jdbc:mysql://localhost/retail\_db  
?useSSL=false --username msis --password bda23

# List the tables

sqoop list-tables --connect jdbc:mysql://localhost/retail\_db  
?useSSL=false --username msis --password bda23

Step 4:-

# Data transfer from local to HDFS

target-dir:- It will not create a folder. Directly dumps all data together

sqoop import --connect jdbc:mysql://localhost/retail\_db  
?useSSL=false --username msis --password bda23 --table orders  
--target-dir '/user/sqoop1' → folder should not be present

30/08/2023

## Warehouse

```
sqoop import --connect jdbc:mysql://localhost/retail_db?  
useSSL=false --username msis --password bda23  
--table orders --warehouse-dir '/user/warehouse'
```

Warehouse copies and stores the data creating folder of name of the table in HDFS.

### # Importing partial data

```
sqoop import --connect jdbc:mysql://localhost/retail_db?  
useSSL=false --username msis --password bda23 --  
table orders --target-dir '/user/partial'
```

⇒ Primary key should be there to copy, or else it will fail

### # copying selected columns

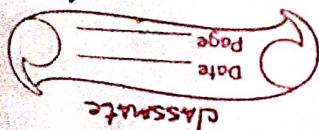
```
sqoop import --connect jdbc:mysql://localhost/retail_db?  
useSSL=false --username msis --password bda23  
--columns product_id,product_name,product_price --table  
products --target-dir '/user/rel/col'
```

### # copying all tables excluding few:- Works for Warehouse

```
sqoop import-all-tables --connect jdbc:mysql://localhost/retail_db?  
useSSL=false --username msis --password bda23  
--warehouse-dir '/test/tables' --exclude-tables products,  
customers
```

### # export data/table from HDFS to mysql:-

```
sqoop export --connect jdbc:mysql://localhost/new_db?useSSL  
=false --username msis --password bda23 --table orders  
--export-dir '/user/warehouse/orders'
```



HDFS folder

↓  
newly  
created  
table (mysql)

## HIVE

Hive is a tool used to perform queries on the data stored in HDFS. In HDFS data is stored as simple texts on which queries can not be operated directly. So, Hive can be used for this.

Step 1: Start HDFS node

start-dfs.sh && start-yarn.sh

Step 2: Use sqoop command to import mysql data from Local to HDFS

Step 3: In mysql check for the schema of copied table.

DESC orders;

Step 4: Open new terminal and use the command hive.

Step 5: It will prompt the hive prompts.

Step 6: In hive create new database and tables

a) create database Test;

b) create table Test.orders (order\_id int, department\_id int, order\_status varchar(45)) row format delimited fields terminated by ',' stored as textfile;

Step 7: In HDFS create a backup of order table.

a) hdfs dfs -mkdir /user/TempTables

b) hdfs dfs -cp /user/orders/ /user/TempTables

Step 8: Go to hive prompt load data from HDFS

a) load data inpath '/user/orders' into table Test.orders;

This will move the data from HDFS to Hive (This is like cut-copy-paste) This data will not be present in HDFS unless we explicitly create a copy of it as given in step 7.

We can not move data back into HDFS from Hive

Step 9: All mysql commands except for update works fine in Hive.

a) select \* from orders; → To see the data

## Map-Reduce :-

- 1) In Eclipse create a new Java project
- 2) In project create a Java class.
- 3) Copy wordcount code to java class
- 4) Right click on project and select Build Path → Configure Build Path
- 5) Select Libraries → Add external libraries (JARs)
- 6) In <sup>hadoop</sup> opt → common → hadoop-common-3.2.1.jar (4.1MB)
- 7) In hadoop → mapreduce → hadoop-mapreduce-client-core-3.2.1.jar
- 8) Select those 2 jars apply
- 9) Select the project → Right click → Export → JAR file → next → Select the jar file name & save it in eclipse-workspace → Next → Next → Finish
- 10) Open terminal → gedit test.txt # creating an input text file
- 11) Copy this file to HDFS, create directory & save  

```
hdfs dfs -mkdir /user/map_reduce/input
```
- 12) hdfs dfs -put /home/sois/test.txt /user/map\_reduce/input  

```
ls
```
- 13) Open new terminal
- 14) # check if the created jar is present  

```
ls
```
- 15) Executing the map-reduce task  

```
hadoop jar jar-name.jar package-name.java-classname
/user/map_reduce/input/test.txt /user/map_reduce/output1
```

Where the output file is generated with count  
 input text file.

Spark :-

```

import pyspark
from pyspark.sql import SparkSession
spark = SparkSession.builder.appName('RDD example')
    ↗ # spark session created
    ↗ getOrCreate

```

```

# type of spark
type(spark)

```

# first way of creating RDD is parallelize method.

# second way of creating is reading a datasource and putting it in an RDD

# third way of creating is using transformations.

03/10/2023

RDD using parallelizing method:-

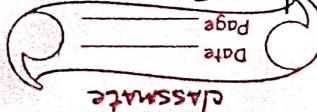
- parRDD = spark.sparkContext.parallelize(['Hope all is well? I'm Nidhi'])  
 → It can have any number of rows, but only one column.  
 → Type of data inside RDD will always be string
- parRDD.collect() → To view the data
- parRDD.count() → Returns the total no. of elements in RDD  
 Count and collect are examples of Actions.

- flatMap() is a transformation  
 → Use it whenever we need 1:n mapping

- map() → use it when we need 1:1 mapping  
 $\text{parRDD.flatMap(lambda data: data.split(' ')) .collect()}$   
 $['Hope', 'all', 'is', 'well', 'I', 'm', 'Nidhi']$

- parRDD.flatMap(lambda data: data.split(' ')).map(lambda d: (d, 1)).collect()

[('Hope', 1), ('all', 1), ('is', 1), ('well', 1),  
 ('I', 1), ('m', 1), ('Nidhi', 1)]



~~parRDD.flatMap(lambda data: data.split(' ')).map(lambda d:(d,1)).reduceByKey(lambda a,b:a+b).collect()~~

~~parRDD.flatMap(lambda d:d.split(' ')).map(lambda d:(len(d),1)).reduceByKey(lambda a,b:a+b).collect()~~

~~parRDD.flatMap(lambda d:d.split(' ')).map(lambda d:(len(d),1)).max()~~

→ It will give the tuple with largest length

5/10/2023

RDD using datasource method :-

→ numRDD = spark.sparkContext.textFile("nums.txt")

→ numRDD.collect()

→ numRDD.flatMap(lambda d:d.split(' ')).split()

→ q1 = numRDD.flatMap(lambda d:d.split(' ')).filter(lambda d: len(d)>0).collect()

→ def clean(data):

    pat = ['0', '1', '2', '3', '4', '5', '6', '7', '8', '9']

    out = ''

    for ele in data:

        if ele in pat:

            out = out + ele

    return out

→ q1.map(lambda d:len(clean(d))).filter(lambda d:len(d)>0).collect()

→ q1.map(lambda d:clean(d)).filter(lambda d:len(d)>0).map(lambda d:int(d)).collect()

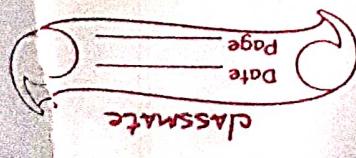
2/10/2023

# Converting an RDD to dataframe

df = qdd.toDF

df.show()

# Showing the structured data as rows & columns



# Renaming the column name  
~~df = df.withColumnRenamed("1", "currency")~~

# Creating a dataframe with column names all together  
~~column = ["currency", "value"]~~

~~df1 = spark.createDataFrame(ROD).toDF(\*column)~~

~~df1.columns() # Gives all column names~~

~~df1.printSchema() # Gives the schema~~

# Directly creating a dataframe if we have a structured dataset  
~~df2 = spark.createDataFrame(data=inputdata, schema=column)~~

CSVdf = spark.read.csv(' ', header=True)

Textdf = spark.read.text(' ')

JSONdf = spark.read.json(' ')

Exceldf = spark.read.excel(' ')

03

# Create dataframe using schema:-

from pyspark.sql.types import \*

sch = "Name STRING, RegNo STRING, marks INT"

CSV2 = spark.read.csv("path", header=True, schema=sch)

# Infer schema

# Need not go write DDL statements

# Use it only for small size dataset (rows are less)

# Can be used if more number of columns are present

CSV3 = spark.read.csv("path", inferSchema=True)

# Built-in Date-Time Operations:-

from pyspark.sql.functions import \*

df = spark.range(5).withColumn('Today', current\_date()).  
 withColumn('Now', current\_timestamp())

df = df.withColumn('New Date', date\_add('Today', 5))

```
df = df.withColumn('Previous date', date_sub('Today', 5))
```

df = ?e

```
df.select(abs(datediff('prev_date', 'New Date'))).show()
```

# Getting the days difference

# lit is the literal function

```
df = df.withColumn('String Date', lit('07-12-2022'))
```

```
df.withColumn('Proper Date', to_date('String Date', 'dd-mm-yyyy'))
```

# converts string type to date data type.

# Changing the default date format from 'yyyy-mm-dd' to  
'dd-mm-yyyy'.

```
df = df.withColumn('Proper Date', date_format(to_date('String Date', 'dd-mm-yyyy'), 'dd-mm-yyyy'))
```

```
df.select(quarter(to_date(lit('2022-12-07')))).show()
```

```
df.select(year(to_date(lit('2022-12-07')))).show()
```

```
df.select(month(to_date(lit('2022-12-07')))).show()
```

```
df.select(weeks(to_date(lit('2022-12-07')))).show()
```

17/10/2023

Fire dataset :-

```
import pyspark
```

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

```
spark = SparkSession.builder.appName('Fire DF').getOrCreate()
```

```
from pyspark.sql.functions import *
```

```
from pyspark.sql.types import *
```

```
df = spark.read.csv('dataset/Fire.csv', header=True, inferSchema=True)
```

```
df.show()
```

```
df.printSchema()
```

```
df1 = df.select('callType', 'Call Date', 'City', 'Zipcode', 'Neighborhood', 'Delay')
```

```
df1.show()
```

17/10/2023

$df_2 = df_1.withColumn('Date', to_date)$

$df_2 = df_1.withColumn('Date', to_date(col('callDate'), 'dd/MM/yyyy'))$   
drop C('callType')

$df_3 = df_2.withColumn('Year', year(col('Date'))).show()$

$df_3 = df_2.withColumn('Year', year(col('Date')))$

•  $withColumn('Month', month(col('Date')))$

•  $withColumn('Week', weekofyear(col('Date')))$

Any python method should be stored as UDF - User Defined Function

def method(data):

newUDF = udf(method, StringType())

df\_3.withColumn('Season', newUDF(col('month'))).show()

#converting pyspark's dataframe to pandas dataframe

import pandas as pd

import matplotlib.pyplot as plt

panda\_df = clean\_df.select('Year')

• groupBy('Year').count()

• orderBy('Year', ascending=True).toPandas()

panda\_df = plot().line(x='Year', y='Count')

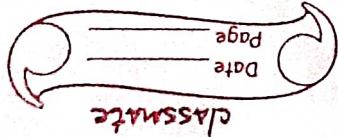
plt.show()

20/10/2023

from pyspark.sql.window import Window

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

window = Window.partitionBy('Program').orderBy(col('Marks').desc())



20/10/2023

```
programRanks = df.withColumn('Ranks', rank().over(windows).denseRank())
```

```
programRanks.filter(col('Rank') < 1).show()
```

```
season_count = clean_df.select('callType', 'Season')\n    .filter(Col('Year') == 2014)\n    .groupBy('Season', 'callType')\n    .count()\n    .withColumn('Ranks', dense_rank().over(Window.partitionBy('Season').orderBy(Col('count').desc))).\n    orderBy('Season', 'Ranks', ascending=[1, 1])
```

Which weeks in 2018 had most fire calls?

```
clean_df.select('week')\n    .where(Col('Year') == 2018)\n    .groupBy('Weeks')\n    .count()\n    .orderBy('count', ascending=False).show()
```

```
max_month.select('weeks', 'count').filter(Col('count') == max_month.\n    agg(CC('count': 'max')).collect()[0][0])
```

28/10/2023

```
clean_df.createOrReplaceTempView('TestView')
```

```
spark.sql("SELECT CallType, city FROM TestView WHERE Year == 2014\n    and Month == 10").show()
```

Streaming data :-

```
import pyspark\nfrom time import sleep\nfrom pyspark.sql import SparkSession\nfrom pyspark.sql.functions import *
```

```
spark = SparkSession.builder.appName('Basic streaming').getOrCreate()
```

28/10/2023

# Get schema of dataset

state = spark.read.json("D:/Dataset/activity-data")

Location of the file. (Don't give filename)

static.printSchema()

dataschema = static.schema.

# Create streaming dataset

streaming = spark.readStream.schema(dataschema).option("maxFilesPerTrigger", 1)\n.json("D:/Dataset/activity-data")

# Get required info from dataset

activityCounts = streaming.groupBy("gt").count()

# Specify the number of shuffle partitions

spark.conf.set("spark.sql.shuffle.partitions", 1)

Since we have only one system, partition = 1. In clusters we can give any value

# Generate query by specifying output sink using writeSchema  
activityQuery = activityCounts.writeStream.queryName("activity-count")\n-

.format("memory").outputMode("complete")

.start()

# Any number of queries (queries) can be created

userCounts = streaming.groupBy('User').count()

userQuery = userCounts.writeStream.queryName('user\_count')\n-

.format("memory").outputMode("complete")

.start()

# Display result

```

for x in range(5):
    spark.sql("SELECT * FROM activity_count").show()
    spark.sql("SELECT * FROM user_count").show()
    sleep(10)

```

# Repeat the process non-stop  
activity.awaitTermination()

31/10/2023 Bank Data generator :-

Count = 30!

```
for x in range(0, 100):
```

data = {}

readings = []

```
for iter in range(0, 10):
```

```
for i in range(0, 20):
```

data['Acc No'] = accNumbers[random.randint(0, 5)]

data['Transaction'] = transType[random.randint(0, 17)]

data['Time stamp'] = str(datetime.now())

data['Amount'] = random.randint(1, 1000000)

readings.append(data.copy())

time.sleep(1)

fname = 'trans' + str(count) + '.json'

files = open('.\data\'+fname, 'w')

files.write(json.dumps(readings))

files.close()

count += 1

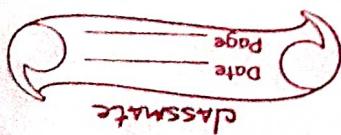
Streaming application :-

```

userSchema = StructType().add('Acc No', 'integer').add('Transaction', 'string')
.add('Time stamp', 'string').add('Amount', 'integer')
streaming = spark.readStream.schema(userSchema)\
```

```
.option("maxFilesPerTrigger", 1)\
```

```
.json('c:\users\DATA\data')
```



31/10/2023

accountDF = streaming.withColumn('Time', to\_timestamp(col('Time stamp'))).drop('Time stamp')  
spark.conf.set("spark.sql.shufflePartitions", 1)

accCount = accountDF.groupBy('Acc No').count()

accountQuery = account.writeStream.queryName('count')  
.format('console')  
.outputMode('complete')  
.start()

accountQuery.awaitTermination()

If the dataframe is result of aggregated data, then only we can select 'complete' output mode.

For filter and other cases, select 'update' or 'append' mode  
non-aggregated operation

accountQuery = account.writeStream.queryName('min amount')  
.trigger(processingTime = "20 seconds")

It waits for time specified before processing the new data.

windowedCount = accountDF  
.withWatermark("Time", "10 minutes")  
.groupBy(col('Acc No'), "Time")  
.count()

watermarks is used when we want the data in out of order

query = windowedCount.writeStream()  
.format("csv")  
.trigger(-----)  
.option("checkpointLocation", ".checkpoint1")  
.option("path", ".outputDir1")  
.outputMode("append")  
.start()