

Let's start our Practical Session 2 :

## Task 1

### Our first streaming using rate source

#### Initialization SparkSession

```
In [ ]: 1 import findspark  
        2 findspark.init()
```

```
In [ ]: 1 import pyspark  
        2 from pyspark.sql import SparkSession  
        3 spark=SparkSession.builder.getOrCreate()
```

#### Let's first import the required libraries.

- Pyspark functions
- Pyspark Types

```
In [ ]: 1 from pyspark.sql.functions import *  
        2 from pyspark.sql.types import *
```

```
In [ ]: 1
```

#### Create streaming DataFrame

Let's create our first Spark Streaming DataFrame using rate source. Here we have specified the format as rate and specified

rowsPerSecond = 1 to generate 1 row for each micro-batch and load the data into initDF streaming DataFrame.

```
In [ ]: 1 first_df = (spark
        2   .readStream
        3   .format("rate")
        4   .option("rowsPerSecond", 1)
        5   .load()
        6   )
```

## Check if DataFrame is streaming or Not.

```
In [ ]: 1 print("Streaming Status : " ,first_df.isStreaming)
```

Streaming DataFrame : True

## Transformation

Perform transformation on initDF to generate another column result by just adding 1 to column value :

```
In [ ]: 1 first_df = first_df.withColumn("result", col("value") + lit(1))
```

## Output

Use append output mode to output only newly generated data and format as console to print the result on the console.

```
In [ ]: 1 first_df.writeStream.outputMode("append").option("truncate", False).format("console").start().awaitTerminat
```

## Our Second streaming using rate source

## Create Streaming DataFrame

## Create Streaming DataFrame using socket source. Also, check if DataFrame isStreaming.

```
In [ ]: 1 host = "128.0.0.2"
        2 port = "8080"
```

```
In [ ]: 1 wordCount.writeStream.outputMode("complete").option("truncate", False).format("console").start().awaitTermination()
```

```
In [ ]: 1 second_df = (spark.readStream.format("socket").option("host", host).option("port", port).load())
```

```
In [ ]: 1 print("Streaming Status : ", second_df.isStreaming)
```

Streaming DataFrame : True

## Transformation: Word Count

```
In [ ]: 1 wordCount = second_df.select(explode(split(col("value"), " ")).alias("words")).groupBy("words").count()
```

Here we count words from a stream of data coming from this socket. Also, we check the schema of our streaming DataFrame.

```
In [ ]: 1 wordCount.writeStream.outputMode("complete").option("truncate", False).format("console").start().awaitTermination()
```

## Input Sources — File

With file input source, our application will wait for available data in the specified directory. We will use some of the stock data available here. For example, Apple stock data present in this file: AAPL\_2006-01-01\_to\_2018-01-01.csv. We will take the data for a few years like 2015, 2016, and 2017 and manually save it to a different file like AAPL\_2015.csv, AAPL\_2016.csv and AAPL\_2017.csv respectively. Similarly, we will create the sample data for Google, Amazon, and Microsoft as well. We will keep all the CSV files locally under data/stocks folder. Also, create another folder data/stream which we will use to simulate the streaming data.

Schema Our data contains the fields Date, Open, High, Low, Close, Adj Close, Volume and we will extract Name from the filename using a custom function. Here we define the schema and write a custom function to extract the stock ticker symbol.



```
In [ ]: 1 third_df.createOrReplaceTempView("stockView")
        2 fourth_df = spark.sql("""select year(Date) as Year, Name, max(High) as Max from stockView group by Name, Ye
```

## Output

Print the contents of streaming DataFrame to console using update mode.

Work with any dataset from data file

```
In [ ]: 1 fourth_df.writeStream.outputMode("update").option("truncate", False).option("numRows", 3).format("console")
```

## Task 2

### Some aggregations transformation

Select "Name", "Date", "Open", "High", "Low" and group by Name and Date and get the avg and use WithColumn to get current\_timestamp

```
In [ ]: 1 fifth_df=third_df.select("Name", "Date", "Open", "High", "Low").groupBy('Name', 'Date').avg().withColumn('cur
```

## Checkpoints

```
In [ ]: 1 fifth_df.writeStream.outputMode("complete").option("truncate", False).format("console").start().awaitTermin
```

## Task 3

### One-time micro-batch

With a once trigger, our query will execute a single micro-batch. It will process all available data and then stop the application. This trigger is useful when you would like to spin-up a cluster periodically, Execute our streaming application with the once trigger.

```
In [ ]: 1 output_df = third_df.groupBy("Name", year("Date").as("Year"))
        2           .agg(max("High").as("Max"))
```

### Users can create GraphFrames from vertex and edge DataFrames.

- Vertex DataFrame: A vertex DataFrame should contain a special column named "id" which specifies unique IDs for each vertex in the graph.
- Edge DataFrame: An edge DataFrame should contain two special columns: "src" (source vertex ID of edge) and "dst" (destination vertex ID of edge).

id	name	age
a	Alice	34
b	Bob	36
c	Charlie	30
d	David	29
e	Esther	32
f	Fanny	36

src	dst	relationship
a	e	friend
f	b	follow
c	e	friend
a	b	friend
b	c	follow
c	b	follow
f	c	follow
e	f	follow
e	d	friend
d	a	friend

```
In [ ]: 1 vertex = spark.createDataFrame([("a", "Alice", 34),("b", "Bob", 36),("c", "Charlie", 30),("d", "David", 29)
```

```
In [ ]: 1 Edge = spark.createDataFrame([ ("a", "e", "friend"),
2   ("f", "b", "follow"),
3   ("c", "e", "friend"),
4   ("a", "b", "friend"),
5   ("b", "c", "follow"),
6   ("c", "b", "follow"),
7   ("f", "c", "follow"),
8   ("e", "f", "follow"),
9   ("e", "d", "friend"),
10  ("d", "a", "friend")],
11      ["src", "dst", "relationship"])
```

```
In [ ]: 1 from graphframes import *
```

## Create a GraphFrame from vertex and edge DataFrames

```
In [ ]: 1 graph = GraphFrame(vertex, Edge)
```

## Take a look at the DataFrames

Get vertices, edges, and check the number of edges of each vertex



In [ ]: 1 vertex.show()

```
+---+-----+---+
| id|   name|age|
+---+-----+---+
|  a|  Alice| 34|
|  b|   Bob| 36|
|  c|Charlie| 30|
|  d|  David| 29|
|  e| Esther| 32|
|  f|  Fanny| 36|
+---+-----+---+
```

In [ ]: 1 Edge.show()

```
+---+-----+-----+
|src|dst|relationship|
+---+-----+-----+
|  a|  e|      friend|
|  f|  b|      follow|
|  c|  e|      friend|
|  a|  b|      friend|
|  b|  c|      follow|
|  c|  b|      follow|
|  f|  c|      follow|
|  e|  f|      follow|
|  e|  d|      friend|
|  d|  a|      friend|
+---+-----+-----+
```

In [ ]: 1 graph.vertices.count()

Out[20]: 6

In [ ]: 1 graph.edges.count()

Out[22]: 10

## Create UDF Functions

convert your graph by mapping a function over the edges DataFrame that deletes the row if  $src \geq dst$  return "Delete" else "Keep"

```
In [ ]: 1 Edge.filter('src>=dst').show()
```

```
+---+---+-----+
|src|dst|relationship|
+---+---+-----+
| f | b |      follow|
| c | b |      follow|
| f | c |      follow|
| e | d |      friend|
| d | a |      friend|
+---+---+-----+
```

## Filtering and connected components

Check vertices when "age" greater than 30 and check edges "relationship" equal "friend"

```
In [ ]: 1 graph2 = graph.filterVertices("age > 30").filterEdges("relationship = 'friend']").dropIsolatedVertices()
```

GraphFrames requires you to set a directory where it can save checkpoints. Create such a folder in your working directory

```
In [ ]: 1 sc=spark.sparkContext
```

```
In [ ]: 1 sc.setCheckpointDir(dirName="graphframes")
        2
        3 output = graph.connectedComponents()
        4 output.select("id", "name", "age").orderBy("id").show()
```

```
+---+-----+---+
| id|   name|age|
+---+-----+---+
| a|  Alice| 34|
| b|   Bob| 36|
| c|Charlie| 30|
| d|  David| 29|
| e| Esther| 32|
| f|  Fanny| 36|
+---+-----+---+
```

Then, the connected components can easily be computed with the connectedComponents-function.

```
In [ ]: 1 output.select("id", "name", "age").orderBy("id").show()
```

```
+---+-----+---+
| id|   name|age|
+---+-----+---+
| a|  Alice| 34|
| b|   Bob| 36|
| c|Charlie| 30|
| d|  David| 29|
| e| Esther| 32|
| f|  Fanny| 36|
+---+-----+---+
```

## Motif finding

Search for pairs of vertices a,b connected by edge e and pairs of vertices b,c connected by edge e2. It will return a DataFrame of all such structures in the graph

In [ ]: 1 graph.find("(a)-[e]->(b); (b)-[e2]->(c)").show()

```

+-----+-----+-----+-----+-----+
|          a|          e|          b|          e2|          c|
+-----+-----+-----+-----+-----+
|[c, Charlie, 30]| [c, e, friend]| [e, Esther, 32]| [e, f, follow]| [f, Fanny, 36]|
|[a, Alice, 34]| [a, e, friend]| [e, Esther, 32]| [e, f, follow]| [f, Fanny, 36]|
|[f, Fanny, 36]| [f, c, follow]| [c, Charlie, 30]| [c, e, friend]| [e, Esther, 32]|
|[b, Bob, 36]| [b, c, follow]| [c, Charlie, 30]| [c, e, friend]| [e, Esther, 32]|
|[d, David, 29]| [d, a, friend]| [a, Alice, 34]| [a, e, friend]| [e, Esther, 32]|
|[c, Charlie, 30]| [c, e, friend]| [e, Esther, 32]| [e, d, friend]| [d, David, 29]|
|[a, Alice, 34]| [a, e, friend]| [e, Esther, 32]| [e, d, friend]| [d, David, 29]|
|[e, Esther, 32]| [e, f, follow]| [f, Fanny, 36]| [f, c, follow]| [c, Charlie, 30]|
|[f, Fanny, 36]| [f, b, follow]| [b, Bob, 36]| [b, c, follow]| [c, Charlie, 30]|
|[c, Charlie, 30]| [c, b, follow]| [b, Bob, 36]| [b, c, follow]| [c, Charlie, 30]|
|[a, Alice, 34]| [a, b, friend]| [b, Bob, 36]| [b, c, follow]| [c, Charlie, 30]|
|[e, Esther, 32]| [e, f, follow]| [f, Fanny, 36]| [f, b, follow]| [b, Bob, 36]|
|[f, Fanny, 36]| [f, c, follow]| [c, Charlie, 30]| [c, b, follow]| [b, Bob, 36]|
|[b, Bob, 36]| [b, c, follow]| [c, Charlie, 30]| [c, b, follow]| [b, Bob, 36]|
|[d, David, 29]| [d, a, friend]| [a, Alice, 34]| [a, b, friend]| [b, Bob, 36]|
|[e, Esther, 32]| [e, d, friend]| [d, David, 29]| [d, a, friend]| [a, Alice, 34]|
+-----+-----+-----+-----+-----+

```

In [ ]:

1