

E10-2 (2)

October 31, 2023

Name : Kattirsitti Jeet Govindrao

Roll no. : 22B0010

0.1 E10-2

This Notebook is about using SPARK Dataframe functions to process nsedata.csv.

0.1.1 Problem

- Write SPARK code to solve the problem stated at the end this Notebook (do not use the createTempView function!)

0.1.2 Submission

Create and upload a PDF of this Notebook after completing your assignment. BEFORE CONVERTING TO PDF and UPLOADING ENSURE THAT YOU REMOVE / TRIM LENGTHY DEBUG OUTPUTS . Short debug outputs of up to 5 lines are acceptable.

```
[19]: import findspark
      findspark.init()
```

```
[20]: import pyspark
      from pyspark.sql.types import *
```

```
[21]: sc = pyspark.SparkContext(appName="E10-2")
```

```
[22]: ss = pyspark.sql.Session(sc)
```

```
[23]: dfr = ss.read
```

```
[24]: schemaStruct = StructType()
      schemaStruct.add("SYMBOL", StringType(), True)
      schemaStruct.add("SERIES", StringType(), True)
      schemaStruct.add("OPEN", DoubleType(), True)
      schemaStruct.add("HIGH", DoubleType(), True)
      schemaStruct.add("LOW", DoubleType(), True)
      schemaStruct.add("CLOSE", DoubleType(), True)
      schemaStruct.add("LAST", DoubleType(), True)
```

```

schemaStruct.add("PREVCLOSE", DoubleType(), True)
schemaStruct.add("TOTTRDQTY", LongType(), True)
schemaStruct.add("TOTTRDVAL", DoubleType(), True)
schemaStruct.add("TIMESTAMP", StringType(), True)
schemaStruct.add("ADDNL", StringType(), True)

```

```

[24]: StructType([StructField('SYMBOL', StringType(), True), StructField('SERIES',
StringType(), True), StructField('OPEN', DoubleType(), True),
StructField('HIGH', DoubleType(), True), StructField('LOW', DoubleType(), True),
StructField('CLOSE', DoubleType(), True), StructField('LAST', DoubleType(),
True), StructField('PREVCLOSE', DoubleType(), True), StructField('TOTTRDQTY',
LongType(), True), StructField('TOTTRDVAL', DoubleType(), True),
StructField('TIMESTAMP', StringType(), True), StructField('ADDNL', StringType(),
True)])

```

```

[25]: df = dfr.csv("/home/hduser/spark/nsedata.csv", schema=schemaStruct, header=True)

```

```

[26]: df.printSchema()

```

```

root
 |-- SYMBOL: string (nullable = true)
 |-- SERIES: string (nullable = true)
 |-- OPEN: double (nullable = true)
 |-- HIGH: double (nullable = true)
 |-- LOW: double (nullable = true)
 |-- CLOSE: double (nullable = true)
 |-- LAST: double (nullable = true)
 |-- PREVCLOSE: double (nullable = true)
 |-- TOTTRDQTY: long (nullable = true)
 |-- TOTTRDVAL: double (nullable = true)
 |-- TIMESTAMP: string (nullable = true)
 |-- ADDNL: string (nullable = true)

```

```

[27]: from pyspark.sql.functions import col, date_format, to_date

df1 = df.withColumn("TIMESTAMP2", date_format(to_date(col("TIMESTAMP")),
↪ "dd-MMM-yyyy"), "yyyy-MM"))

```

```

[28]: df1.printSchema()

```

```

root
 |-- SYMBOL: string (nullable = true)
 |-- SERIES: string (nullable = true)
 |-- OPEN: double (nullable = true)
 |-- HIGH: double (nullable = true)
 |-- LOW: double (nullable = true)
 |-- CLOSE: double (nullable = true)

```

```

|-- LAST: double (nullable = true)
|-- PREVCLOSE: double (nullable = true)
|-- TOTTRDQTY: long (nullable = true)
|-- TOTTRDVAL: double (nullable = true)
|-- TIMESTAMP: string (nullable = true)
|-- ADDNL: string (nullable = true)
|-- TIMESTAMP2: string (nullable = true)

```

0.2 Problem Statement

Using SPARK Dataframe functions write code to create the data shown below for all the traded companies. Save this data in an output file in ascending order of the company names, year and month.

SYMBOL | Month-Year | min(CLOSE) | max(CLOSE) | avg(CLOSE) | stddev(CLOSE) | traded-Count |

The output should appear as follows

SYMBOL	TIMESTAMP2	min(OPEN)	max(OPEN)	avg(OPEN)	stddev(OPEN)	count(OPEN)
20MICRONS	2010-08	51.6	54.0	52.81666666666667	0.9266876496425305	9
20MICRONS	2010-09	54.9	64.3	59.11428571428571	2.514614426564382	21
20MICRONS	2010-10	55.05	60.0	57.166666666666664	1.3035848009751156	21
20MICRONS	2010-11	53.6	61.75	55.98809523809524	2.2001650370997603	21
20MICRONS	2010-12	38.8	61.0	45.66590909090909	5.796599708606606	22
20MICRONS	2011-01	38.3	48.2	44.042500000000004	2.357310856396376	20
20MICRONS	2011-02	35.15	45.9	41.635	2.3022929074248895	20
20MICRONS	2011-03	35.2	40.9	37.83636363636364	1.735770846886316	22
20MICRONS	2011-04	37.75	42.9	40.66388888888889	1.4290891335511524	18
20MICRONS	2011-05	40.1	47.3	42.304545454545455	2.2407433445021625	22

tradedCount = number of times the company shares have been traded in that month

Notes and Hints:

- use the functions `groupBy` (based on `SYMBOL` and `TIMESTAMP2`) and `agg` to create the individual statistics like `min`, `max`, `avg`, etc.
- use `join` (based on `SYMBOL` and `TIMESTAMP2`) to combine the individual dataframes into a single table

This is just one method of solving the problem! You can discover of any other method, using any other combination of Dataframe functions-

```

[29]: df1.show()
# Current df1

```

```

+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
| SYMBOL | SERIES | OPEN | HIGH | LOW | CLOSE | LAST | PREVCLOSE | TOTTRDQTY |

```

TOTTRDVAL	TIMESTAMP	ADDNL	TIMESTAMP2						
20MICRONS	EQ 37.75	37.75	36.35	37.45	37.3	37.15	38638		
1420968.1 01-APR-2011	0	2011-04							
3IINFOTECH	EQ 43.75	45.3	43.75	44.9	44.8	43.85	1239690		
5.531120435E7 01-APR-2011	0	2011-04							
3MINDIA	EQ 3374.0 3439.95	3338.0	3397.5	3400.0	3364.7	871			
2941547.35 01-APR-2011	0	2011-04							
A2ZMES	EQ 281.8	294.45	279.8	289.2	287.2	281.3	140643		
4.02640755E7 01-APR-2011	0	2011-04							
AARTIDRUGS	EQ 127.0	132.0	126.55	131.3	130.6	127.6	2972		
384468.2 01-APR-2011	0	2011-04							
AARTIIND	EQ 50.0	50.0	49.0	49.25	49.35	49.05	24056		
1188195.85 01-APR-2011	0	2011-04							
AARVEEDEN	EQ 58.45	58.45	56.6	56.65	56.6	56.55	123		
7000.1 01-APR-2011	0	2011-04							
ABAN	EQ 620.0	645.95	617.0	643.3	644.0	616.25			
1192421 7.5745251715E8 01-APR-2011	0	2011-04							
ABB	EQ 796.8	796.8	777.35	785.2	780.2	796.8	58038		
4.562089595E7 01-APR-2011	0	2011-04							
ABBOTINDIA	EQ 1379.0	1379.0 1335.05	1353.2	1355.0	1343.05	587			
793494.8 01-APR-2011	0	2011-04							
ABCIL	EQ 129.55	130.8	128.35	130.0	130.0	129.7	1941		
251299.4 01-APR-2011	0	2011-04							
ABGSHIP	EQ 367.0	374.0	335.6	370.0	370.0	363.75	307293		
1.134908749E8 01-APR-2011	0	2011-04							
ABHISHEK	EQ 15.0	16.0	15.0	15.95	16.0	15.2	6360		
100264.3 01-APR-2011	0	2011-04							
ABIRLANUVU	EQ 816.45	844.7	812.4	824.85	824.9	814.35	70865		
5.86104648E7 01-APR-2011	0	2011-04							
ABSHEKINDS	EQ 14.4	15.25	14.2	15.05	15.2	14.2	159188		
2365626.1 01-APR-2011	0	2011-04							
ACC	EQ 1070.0	1098.0 1069.95	1091.85	1091.15	1074.55	240346			
2.598602339E8 01-APR-2011	0	2011-04							
ACE	EQ 43.2	44.9	42.0	44.5	44.5	43.1	142292		
6273701.1 01-APR-2011	0	2011-04							
ACKRUTI	EQ 228.0	228.0	223.1	224.5	225.5	226.4	23053		
5187103.8 01-APR-2011	0	2011-04							
ACROPETAL	EQ 58.75	62.7	58.35	58.35	58.35	61.4	2507573		
1.513034327E8 01-APR-2011	0	2011-04							
ADANIENT	EQ 666.0	668.2	652.3	661.0	665.0	666.55			
210396 1.3885739135E8 01-APR-2011	0	2011-04							

only showing top 20 rows

23/10/31 17:17:35 WARN CSVHeaderChecker: Number of column in CSV header is not equal to number of fields in the schema:

Header length: 14, schema size: 12

CSV file: file:///home/hduser/spark/nsedata.csv

```
[30]: df1_min = df1.groupBy('SYMBOL', 'TIMESTAMP2').agg({'CLOSE': 'min'})
df1_max = df1.groupBy('SYMBOL', 'TIMESTAMP2').agg({'CLOSE': 'max'})
df1_avg = df1.groupBy('SYMBOL', 'TIMESTAMP2').agg({'CLOSE': 'avg'})
df1_stddev = df1.groupBy('SYMBOL', 'TIMESTAMP2').agg({'CLOSE': 'stddev'})
df1_count = df1.groupBy('SYMBOL', 'TIMESTAMP2').agg({'CLOSE': 'count'})
```

```
[31]: df1 = df1_min.join(df1_max, on=["SYMBOL", "TIMESTAMP2"], how="inner") \
      .join(df1_avg, on=["SYMBOL", "TIMESTAMP2"], how="inner") \
      .join(df1_stddev, on=["SYMBOL", "TIMESTAMP2"], how="inner") \
      .join(df1_count, on=["SYMBOL", "TIMESTAMP2"], how="inner")
# Backslashes are only for readability.
```

```
[32]: df1.show()
```

[Stage 13:> (0 + 1) / 1]

```
+-----+-----+-----+-----+-----+-----+
--+-----+
|  SYMBOL|TIMESTAMP2|min(CLOSE)|max(CLOSE)|          avg(CLOSE)|
stddev(CLOSE)|count(CLOSE)|
+-----+-----+-----+-----+-----+-----+
--+-----+
|20MICRONS|  2010-08|      51.55|      54.3|          52.75|
1.0647769719523452|          9|
|20MICRONS|  2010-09|      54.9|      60.9|  58.4547619047619|
1.7269123285436907|         21|
|20MICRONS|  2010-11|      53.35|      60.3|  55.69047619047619|
1.8280193549043067|         21|
|20MICRONS|  2011-01|      41.3|      47.75|43.917500000000004|
1.9892656010646148|         20|
|20MICRONS|  2011-03|      35.85|      40.1|  37.70227272727272|
1.3524648813966484|         22|
|20MICRONS|  2011-04|      37.45|      41.65|          40.425|
1.0581963011486593|         18|
|20MICRONS|  2011-08|      46.6|      54.7|  51.05952380952381|
2.2866798555776837|         21|
|20MICRONS|  2011-09|      51.85|      59.35|56.145238095238106|
1.7756621836588724|         21|
|20MICRONS|  2012-01|      61.7|      66.0|  63.04523809523808|
1.2234076141974053|         21|
|20MICRONS|  2012-03|      79.05|      85.1|  81.84761904761905|
1.578011105036673|         21|
|20MICRONS|  2012-05|      84.35|      93.1|  88.00227272727273|
2.7357630289645605|         22|
```

```
|20MICRONS| 2012-06| 84.7| 89.95| 87.38095238095237|
1.7464590025589037| 21|
|20MICRONS| 2012-07| 84.45| 97.7| 92.38181818181819|
4.4267731302067475| 22|
|20MICRONS| 2012-08| 96.85| 120.85|114.29999999999998|
7.240700932920781| 21|
|20MICRONS| 2012-10| 106.55| 123.05|115.22857142857144|
4.723784802162661| 21|
|20MICRONS| 2012-11| 125.6| 163.9|145.45499999999998|
11.469064799207526| 20|
|20MICRONS| 2012-12| 135.7| 156.95| 143.685|
6.708284352005988| 20|
|20MICRONS| 2013-01| 68.1| 160.8|136.45434782608694|
31.77105569912895| 23|
|20MICRONS| 2013-02| 30.15| 68.65| 42.2325|
14.365723780472818| 20|
|20MICRONS| 2013-05| 30.4| 31.5|
30.72954545454546|0.26306756823019073| 22|
+-----+-----+-----+-----+-----+-----+
--+-----+
only showing top 20 rows
```

```
[33]: ss.stop()
      sc.stop()
```