Big Data 230 : EMERGING TECHNOLOGIES IN BIG DATA PROJECT

STOCK ANALYSIS QUERIES

//Moving Files to VM

scp *.csv BigData230@52.183.35.71:/tmp

scp -P 2222 *.csv root@localhost:/tmp

//Creating Data Frames from files

val

s1=spark.read.format("csv").option("header","true").option("inferSchema","true").load("file:///tmp/daily_AAPL.csv").withColumn("Stock",lit("AAPL"))

val

s2=spark.read.format("csv").option("header","true").option("inferSchema","true").load("file:///tmp/daily_ACN.csv").withColumn("Stock",lit("ACN"))

val

s3=spark.read.format("csv").option("header","true").option("inferSchema","true").load("file:///tmp/daily_ADBE.csv").withColumn("Stock",lit("ADBE"))

val

s4=spark.read.format("csv").option("header","true").option("inferSchema","true").load("file:///tmp/daily_AMZN.csv").withColumn("Stock",lit("AMZN"))

val

s5=spark.read.format("csv").option("header","true").option("inferSchema","true").load("file:///tmp/daily_AVGO.csv").withColumn("Stock",lit("AVGO"))

val

s6=spark.read.format("csv").option("header","true").option("inferSchema","true").load("file:///tmp/daily_BABA.csv").withColumn("Stock",lit("BABA"))

val

s7=spark.read.format("csv").option("header","true").option("inferSchema","true").load("file:///tmp/daily_CRM.csv").withColumn("Stock",lit("CRM"))

val

s8=spark.read.format("csv").option("header","true").option("inferSchema","true").load("file:///tmp/daily_CSCO.csv").withColumn("Stock",lit("CSCO"))

```
val
s9=spark.read.format("csv").option("header","true").option("inferSchema","true").load("file:///tmp/
daily EXPE.csv").withColumn("Stock",lit("EXPE"))
val
s10=spark.read.format("csv").option("header","true").option("inferSchema","true").load("file:///tm
p/daily_FB.csv").withColumn("Stock",lit("FB"))
val
s11=spark.read.format("csv").option("header","true").option("inferSchema","true").load("file:///tm
p/daily_GOOGL.csv").withColumn("Stock",lit("GOOGL"))
val
s12=spark.read.format("csv").option("header","true").option("inferSchema","true").load("file:///tm
p/daily_IBM.csv").withColumn("Stock",lit("IBM"))
val
s13=spark.read.format("csv").option("header","true").option("inferSchema","true").load("file:///tm
p/daily INTC.csv").withColumn("Stock",lit("INTC"))
val
s14=spark.read.format("csv").option("header", "true").option("inferSchema", "true").load("file:///tm
p/daily MSFT.csv").withColumn("Stock",lit("MSFT"))
val
s15=spark.read.format("csv").option("header", "true").option("inferSchema", "true").load("file:///tm
p/daily_MU.csv").withColumn("Stock",lit("MU"))
val
s16=spark.read.format("csv").option("header","true").option("inferSchema","true").load("file:///tm
p/daily_NFLX.csv").withColumn("Stock",lit("NFLX"))
val
s17=spark.read.format("csv").option("header","true").option("inferSchema","true").load("file:///tm
p/daily_NVDA.csv").withColumn("Stock",lit("NVDA"))
val
s18=spark.read.format("csv").option("header","true").option("inferSchema","true").load("file:///tm
p/daily_ORCL.csv").withColumn("Stock",lit("ORCL"))
```

val

s19=spark.read.format("csv").option("header","true").option("inferSchema","true").load("file:///tmp/daily_sbux.csv").withColumn("Stock",lit("SBUX"))

val

s20=spark.read.format("csv").option("header","true").option("inferSchema","true").load("file:///tmp/daily_TXN.csv").withColumn("Stock",lit("TXN"))

val StockData =

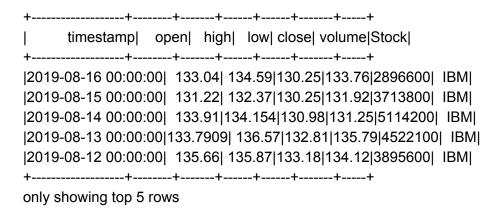
s1.unionAll(s2).unionAll(s3).unionAll(s4).unionAll(s5).unionAll(s6).unionAll(s7).unionAll(s8).unionAll(s9).unionAll(s10).unionAll(s11).unionAll(s12).unionAll(s13).unionAll(s14).unionAll(s15).unionAll(s16).unionAll(s17).unionAll(s18).unionAll(s19).unionAll(s20)

//Writing to parquet file

StockData.write.parguet("hdfs:///tmp/AllStockData.parguet")

val StockDataDF= spark.read.parquet("hdfs:///tmp/AllStockData.parquet")

StockDataDF.show(5)



Q1: Tech stock with max movement in 20 years

StockDataDF.createOrReplaceTempView("StockData")

spark.sql("select * from StockData where Stock = 'SBUX'").show(3)

val StockReturn =

StockDataDF.groupBy(\$"Stock").agg(max(\$"high").as("max_stock"),min(\$"low").as("min_stock")).withColumn("sold",coalesce(\$"max_stock"-\$"min_stock",lit(0.0))).withColumn("Return", coalesce(\$"sold"*100/\$"min_stock",lit(0.0)))

StockReturn.orderBy(\$"Return"desc).show

Return	sold	min_stock	max_stock	Stock
37114.15607985481	2044.99	5.51	2050.5	AMZN
14666.18556701031	711.31	4.85	716.16	NFLX
9039.182389937107	143.72299999999998	1.59	145.313	MU
5443.003144654088	692.35	12.72	705.07	AAPL
4991.478260869565	287.01	5.75	292.76	NVDA
2583.333333333333	155.0	6.0	161.0	EXPE
2155.4082344731332	308.87	14.33	323.2	AVGO
1988.222222222222	178.94	9.0	187.94	CRM
1894.331210191083	297.41	15.7	313.11	ADBE
1707.266009852217	138.63	8.12	146.75	CSCO
1647.5221238938052	186.17	11.3	197.47	ACN
1312.4645892351275	92.66	7.06	99.72	SBUX
1251.5787828261775	1201.0149999999999	95.96	1296.975	GOOGL
1145.6980056980055	201.07	17.55	218.62	FB
895.1908396946567	117.270000000000001	13.1	130.37	TXN
852.7908540685945	126.81	14.87	141.68	MSFT
734.4827586206897	53.25	7.25	60.5	ORCL
529.2780082987551	63.7780000000000006	12.05	75.828	INTC
299.7267268879636	161.8881	54.0119	215.9	IBM
270.1048951048951	154.5	57.2	211.7	BABA

Q2: Which year did each of the individual stocks have the most change (+ or -) in their price

scala> val StockDataDF= spark.read.parquet("hdfs:///tmp/AllStockData.parquet")

StockDataDF: org.apache.spark.sql.DataFrame = [timestamp: timestamp, open: double ... 5 more fields]

scala> StockDataDF.createOrReplaceTempView("StockData")

scala> val YearlyDataDF = spark.sql("select max(high) as maxhigh, min(low) as minlow, Stock, year(timestamp) as yearstock from StockData group by Stock, yearstock order by yearstock")

YearlyDataDF: org.apache.spark.sql.DataFrame = [maxhigh: double, minlow: double ... 2 more fields]

scala> YearlyDataDF.createOrReplaceTempView("YearData")

scala> val maxYearlyFluctuationDF = spark.sql("select max((maxhigh - minlow)/minlow) as range, Stock from YearData group by Stock")

maxYearlyFluctuationDF: org.apache.spark.sql.DataFrame = [range: double, Stock: string]

scala> maxYearlyFluctuationDF.createOrReplaceTempView("maxYearlyFluctuation")

scala> val maxYearChanges = spark.sql("select y.*, range from YearData as y join maxYearlyFluctuation as m on y.Stock = m.Stock where (y.maxhigh - y.minlow)/y.minlow = m.range").show(50)

range	yearstock	Stock	minlow	maxhigh
5.151260504201681	2000	AMZN	14.875	91.5
2.1847748792993746	2000	ADBE	53.438	170.188
1.8508857142857145	2000	TXN	35.0	99.781
3.17425190579133	2000	CSC0	35.156	146.75
10.036697247706423	2000	AAPL	13.625	150.375
1.9472049689440993	2000	MSFT	40.25	118.625
2.4448818897637796	2001	ORCL	10.16	35.0
2.8122585438335808	2001	SBUX	13.46	51.313
9.09166666666665	2002	NVDA	7.2	72.66
1.3400398801004965	2002	IBM	54.0119	126.39
1.8401544401544403	2002	INTC	12.95	36.78
7.59891891891892	2004	NFLX	9.25	79.54
4.3133333333333333	2008	EXPE	6.0	31.88
1.7087806867619637	2008	GOOGL	257.44	697.3485
4.6415094339622645	2008	MU	1.59	8.97
1.899098083427283	2010	ACN	17.74	51.43
1.5840317600352887	2013	FB	22.67	58.58
4.207536713771127	2013	CRM	36.09	187.94
1.0235112738485257		AVGO	51.89	105.0
1.1768770810888132	2017	BABA	88.0849	191.75