

Data: NYSE.CSV

- Write a series of Pig Latin commands to identify:

1. The total volume traded for each stock for each month.

```
grunt> data = LOAD 'Sudi/NYSE.csv' USING PigStorage(',') AS (
>> exchange: chararray,
>> stock_symbol: chararray,
>> date: chararray,
>> stock_price_open: float,
>> stock_price_high: float,
>> stock_price_low: float,
>> stock_price_close: float,
>> stock_volume: int,
>> stock_price_adj_close: float
>> );
2024-03-14 00:35:07,187 INFO Configuration.deprecation: yarn.resourcemanager.system-metrics-publisher.enabled is deprecated. Instead, use yarn.system-metrics-publisher.enabled
grunt> []
```

```
grunt> data_with_month = FOREACH data GENERATE
>> stock_symbol,
>> GetMonth(ToDate(date, 'MM/dd/yyyy')) AS month,
>> stock_volume;
```

```
grunt> NYSE_Data = LOAD 'Sudi/NYSE.csv' USING PigStorage(',') AS (
>> exchange: chararray,
>> stock_symbol: chararray,
>> date: chararray,
>> stock_price_open: float,
>> stock_price_high: float,
>> stock_price_low: float,
>> stock_price_close: float,
>> stock_volume: int,
>> stock_price_adj_close: float
>> );
2024-03-13 23:56:42,383 INFO Configuration.deprecation: yarn.resourcemanager.system-metrics-publisher.enabled is deprecated. Instead, use yarn.system-metrics-publisher.enabled
grunt>
grunt> NYSE_with_date = FOREACH NYSE_Data GENERATE
>> exchange,
>> stock_symbol,
>> ToDate(date, 'MM/dd/yyyy') AS date,
>> stock_price_open,
>> stock_price_high,
>> stock_price_low,
>> stock_price_close,
>> stock_volume,
>> stock_price_adj_close;
grunt>
grunt>
```

Extract the month from the datetime value.

```
grunt> data_with_month = FOREACH NYSE_with_date GENERATE
>> exchange,
>> stock_symbol,
>> GetMonth(date) AS month,
>> stock_volume;
grunt>
```

```
grunt> volume_per_month = FOREACH (GROUP data_with_month BY (stock_symbol, month)) GENERATE
>> FLATTEN(group) AS (stock_symbol, month),
>> SUM(data_with_month.stock_volume) AS total_volume;
grunt>
```

2. The total volume traded for each stock.

```
grunt> volume_per_stock = FOREACH (GROUP NYSE_with_date BY stock_symbol) GENERATE
>> group AS stock_symbol,
>> SUM(NYSE_with_date.stock_volume) AS total_volume;
```

```
grunt> -- Calculate total volume traded for each stock
grunt> volume_per_stock = FOREACH (GROUP data BY stock_symbol) GENERATE
>> group AS stock_symbol,
>> SUM(data.stock_volume) AS total_volume;
grunt>
```

3. Store the results in HDFS and display them.

```
STORE volume_per_month INTO 'output_per_month' USING PigStorage(',');
```

```
DUMP volume_per_month;
```

```
STORE volume_per_stock INTO 'output_per_stock' USING PigStorage(',');
```

```
DUMP volume_per_stock;
```

```
FileBytesRead: 1633
FileBytesWritten: 1569
HdfsBytesRead: 268137
HdfsBytesWritten: 1793

SpillableMemoryManager spill count: 0
Bags proactively spilled: 0
Records proactively spilled: 0

DAG Plan:
Tez vertex scope-141 -> Tez vertex scope-142,
Tez vertex scope-142

Vertex Stats:
VertexId Parallelism TotalTasks InputRecords ReduceInputRecords OutputRecords FileBytesRead FileBytesWritten
scope-141 1 1 5001 0 5001 64
scope-142 1 1 0 155 155 1569
c2.internal:8020/user/hadoop/output_per_stock, GROUP_BY hdfs://ip-172-31-25-235.ec2.internal:8020/user/hadoop/output_per_stock"

Input(s):
Successfully read 5001 records (268137 bytes) from: "hdfs://ip-172-31-25-235.ec2.internal:8020/user/hadoop/p/Sudi/NYSE.csv"

Output(s):
Successfully stored 155 records (1793 bytes) in: "hdfs://ip-172-31-25-235.ec2.internal:8020/user/hadoop/output_per_stock"

grunt> []
```

```
2024-03-14 01:12:21,438 INFO util.MapRedUtil: Total input paths to process : 1
(GA,9531100)
(GB,7681400)
(GD,146828800)
(GE,2141242700)
(GF,6057900)
(GG,198083800)
(GH,52768300)
(GI,276912700)
(GJ,123810700)
(GK,52211500)
(GL,14105700)
(GM,15470700)
(GN,425600)
(GO,675700)
(GP,275000)
(GQ,2837400)
(GR,14445200)
(GS,7100)
(GT,16713800)
(GU,116700)
(GV,13700)
(GW,2833700)
(GX,75400)
(GY,1701400)
(GZ,4668200)
(HA,3401700)
(HB,191300)
(HC,5151000)
(HD,62263600)
(HE,20883600)
(HF,1587000)
(HG,748100)
(HH,1747500)
(HI,14871600)
(HJ,602100)
(HK,8000)
(HL,8523900)
```

need to use ToDate operator to change the time format from String to DateTime