

Use Case 4

Gabriel D Hofer

July 18, 2021

Apply data frames on Walmart dataset and solve the below problems:

1. Create a new dataframe with a column called HV Ratio that is the ratio of the High Price versus volume of stock traded for a day.

```
1 scala> val df = spark.read.option("header",true).csv("/home/gabriel/Data/walmart_stock.txt")
2
3
4 scala> df.withColumn("HV Ratio", col("High") / col("Volume"))
```

2. What day had the Peak High in Price?

```
1 scala> df.sort(col("High").desc).show(1)
2 +-----+
3 |      Date|      Open|      High|      Low|      Close| Volume|Adj Close|
4 +-----+
5 |2015-01-13|90.800003|90.970001|88.93|89.309998|8215400|83.825448|
6 +-----+
7 only showing top 1 row
```

3. What is the mean of the Close column?

```
1 scala> df.select(mean(col("Close"))).show()
2 +-----+
3 |      avg(Close)|
4 +-----+
5 |72.38844998012726|
6 +-----+
```

4. What is the max and min of the Volume column?

```
1 scala> df.agg(max(col("Volume")), min(col("Volume"))).show()
2 +-----+
3 |max(Volume)|min(Volume)|
4 +-----+
5 |    9994400|    10010500|
6 +-----+
```

5. How many days was the Close lower than 60 dollars?

```
1 scala> df.filter(col("Close") < lit(60.0)).count()
2 res10: Long = 81
```

6. What percentage of the time was the High greater than 80 dollars ?

```
1 scala> df.agg((sum(when(col("High") > lit(80.0),1).otherwise(0)) / count("*")
   ).as("High > 80")).show()
2 +-----+
3 |           High > 80|
4 +-----+
5 |0.09141494435612083|
6 +-----+
```

7. What is the Pearson correlation between High and Volume?

$$r = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum (x_i - \bar{x})^2 \sum (y_i - \bar{y})^2}}$$

Where,

r = Pearson Correlation Coefficient

x_i = x variable samples y_i = y variable sample

\bar{x} = mean of values in x variable \bar{y} = mean of values in y variable

```
1 scala> val highAvg = df.select(mean(col("High"))).first().getDouble(0)
2
3 scala> val volumeAvg = df.select(mean(col("Volume"))).first().getDouble(0)
4
5 scala> val numerator = df.agg(sum((col("High") - highAvg) * (col("Volume") -
   volumeAvg))).first().getDouble(0)
6 numerator: Double = -1.3013592106342264E10
7
8 scala> val denominator = math.sqrt(df.agg( sum((col("High") - highAvg) * (col
   ("High") - highAvg)) * sum((col("Volume") - volumeAvg) * (col("Volume") -
   volumeAvg))).first().getDouble(0))
9 denominator: Double = 3.845253639556948E10
10
11 scala> val r = numerator / denominator
12 r: Double = -0.33843260617371645
```

8. What is the max High per year?

```
1 scala> df.groupBy(substring(col("Date"),0,4).as("Year")).agg(max(col("High"))
   ).orderBy(col("Year")).show()
2 +-----+
3 |Year|max(High)|
4 +-----+
5 |2012|77.599998|
6 |2013|81.370003|
7 |2014|88.089996|
8 |2015|90.970001|
9 |2016|75.190002|
10 +-----+
```

9. What is the average Close for each Calendar Month?

```

1
2 scala> df.groupBy(substring(col("Date"),0,7).as("Year-Month")).agg(avg(col("
3   Close"))).show(100)
4 +-----+
5 |Year-Month|      avg(Close)|
6 +-----+
7 | 2013-05| 77.81636368181817|
8 | 2013-09|      74.4395005|
9 | 2013-12| 78.7752382857143|
10 | 2013-06| 74.97800020000001|
11 | 2016-02| 66.24800044999999|
12 | 2015-05| 77.33599970000002|
13 | 2012-08| 73.04478265217392|
14 | 2015-12| 59.98681827272728|
15 | 2012-02|      60.898|
16 | 2012-04| 60.149000150000006|
17 | 2016-12| 70.51904728571428|
18 | 2012-05| 61.456363409090905|
19 | 2016-09| 72.00857180952381|
20 | 2016-03| 67.55499963636365|
21 | 2016-10| 69.23952366666666|
22 | 2012-12| 69.71100009999999|
23 | 2012-07| 72.40666661904763|
24 | 2014-01| 76.53142833333334|
25 | 2015-02| 85.52315805263159|
26 | 2015-08| 69.2866677142857|
27 | 2014-03| 75.30238076190474|
28 | 2014-08| 74.67666623809525|
29 | 2013-11| 78.97300075000001|
30 | 2014-02| 74.05578978947368|
31 | 2012-01|      60.2354999|
32 | 2015-10| 61.564545636363626|
33 | 2012-11| 71.10952333333333|
34 | 2013-02| 70.62315857894738|
35 | 2015-11| 58.911999949999995|
36 | 2014-09| 76.33619004761903|
37 | 2016-05| 68.05285676190476|
38 | 2012-06| 67.50380961904762|
39 | 2013-07| 77.11545418181818|
40 | 2014-05| 77.38095276190477|
41 | 2012-09| 74.18157921052631|
42 | 2015-01|      87.60949975|
43 | 2016-04| 68.82523861904761|
44 | 2016-11| 70.30476261904762|
45 | 2013-08| 75.22409204545455|
46 | 2014-11| 81.88526321052632|
47 | 2016-01| 63.22105263157895|
48 | 2016-07| 73.54149939999999|
49 | 2015-07| 72.75000036363635|
50 | 2016-06| 71.34636304545454|
51 | 2014-04| 77.80857085714285|
52 | 2014-10| 76.48869486956522|
53 | 2015-09| 64.25238128571429|
54 | 2014-12| 85.1259102727273|
55 | 2013-03| 73.43649940000002|
56 | 2013-10| 74.97913104347826|
57 | 2014-06| 76.01000033333332|
58 | 2012-10| 75.30619061904761|
59 | 2014-07| 76.21090877272728|

```

59		2015-04		79.56047561904764	
60		2015-06		72.79727304545456	
61		2013-01		69.09476142857143	
62		2016-08		72.8300000869565	
63		2013-04		77.68954572727273	
64		2012-03		60.433636818181796	
65		2015-03		82.47318172727273	
66		-----		-----	