

Spark DataFrames Project Exercise

Let's get some quick practice with your new Spark DataFrame skills, you will be asked some basic questions about some stock market data, in this case Walmart Stock from the years 2012-2017. This exercise will just ask a bunch of questions, unlike the future machine learning exercises, which will be a little looser and be in the form of "Consulting Projects", but more on that later!

For now, just answer the questions and complete the tasks below.

Use the walmart_stock.csv file to Answer and complete the tasks below!

Start a simple Spark Session

```
In [1]: from pyspark.sql import SparkSession
spark = SparkSession.builder.appName('assignment').getOrCreate()
```

Load the Walmart Stock CSV File, have Spark infer the data types.

```
In [2]: df = spark.read.csv('D:\Python\PySpark\Python-and-Spark-for-Big-Data-master\Spark_DataFrame_Project_Exercise\walmart_stock.csv',inferSchema= True, header= True)
```

```
In [87]: df.show()
```

Date	Open	High	Low	Close	Volume	Adj Close
2012-01-03	59.970001	61.060001	59.869999	60.330002	12668800	52.619234999999996
2012-01-04	60.209998999999996	60.349998	59.470001	59.709998999999996	9593300	52.078475
2012-01-05	59.349998	59.619999	58.369999	59.419998	12768200	51.825539
2012-01-06	59.419998	59.450001	58.869999	59.0	8069400	51.45922
2012-01-09	59.029999	59.549999	58.919998	59.18	6679300	51.616215000000004
2012-01-10	59.43	59.709998999999996	58.98	59.040001000000004	6907300	51.494109
2012-01-11	59.060001	59.529999	59.040001000000004	59.400002	6365600	51.808098
2012-01-12	59.790001000000004	60.0	59.400002	59.5	7236400	51.895315999999994
2012-01-13	59.18	59.610001000000004	59.009997999999996	59.540001000000004	7729300	51.930203999999996
2012-01-17	59.869999	60.110001000000004	59.52	59.849998	8500000	52.200581
2012-01-18	59.790001000000004	60.029999	59.650002	60.009997999999996	5911400	52.340131
2012-01-19	59.93	60.73	59.75	60.610001000000004	9234600	52.863447
2012-01-20	60.75	61.25	60.669998	61.009997999999996	10378800	53.212320999999996
2012-01-23	60.810001	60.98	60.509997999999996	60.91	7134100	53.125104
2012-01-24	60.75	62.0	60.75	61.389998999999996	7362800	53.543754000000001
2012-01-25	61.18	61.610001000000004	61.040001000000004	61.470001	5915800	53.613531000000001
2012-01-26	61.799999	61.84	60.77	60.970001	7436200	53.177436
2012-01-27	60.860001000000004	61.119999	60.540001000000004	60.709998999999996	6287300	52.950665
2012-01-30	60.470001	61.32	60.349998	61.299999	7636900	53.465256999999994
2012-01-31	61.529999	61.57	60.580002	61.360001000000004	9761500	53.517590000000006

only showing top 20 rows

What are the column names?

```
In [10]: df.columns

Out[10]: ['Date', 'Open', 'High', 'Low', 'Close', 'Volume', 'Adj Close']
```

What does the Schema look like?

```
In [13]: df.printSchema()

root
|-- Date: string (nullable = true)
|-- Open: double (nullable = true)
|-- High: double (nullable = true)
|-- Low: double (nullable = true)
|-- Close: double (nullable = true)
|-- Volume: integer (nullable = true)
|-- Adj Close: double (nullable = true)
```

Print out the first 5 columns.

```
In [25]: for row in df.head(5):
          print(row)
          print('\n')
```

```
Row(Date='2012-01-03', Open=59.970001, High=61.060001, Low=59.869999, Close=60.330002, Volume=12668800, Adj Close=52.619234999999996)
```

```
Row(Date='2012-01-04', Open=60.209998999999996, High=60.349998, Low=59.470001, Close=59.709998999999996, Volume=9593300, Adj Close=52.078475)
```

```
Row(Date='2012-01-05', Open=59.349998, High=59.619999, Low=58.369999, Close=59.419998, Volume=12768200, Adj Close=51.825539)
```

```
Row(Date='2012-01-06', Open=59.419998, High=59.450001, Low=58.869999, Close=59.0, Volume=8069400, Adj Close=51.45922)
```

```
Row(Date='2012-01-09', Open=59.029999, High=59.549999, Low=58.919998, Close=59.18, Volume=6679300, Adj Close=51.616215000000004)
```

Use describe() to learn about the DataFrame.

```
In [22]: df.describe().show()
```

-----+							
summary	Date	Open	High	Low	Close	Volume	
Adj Close							
-----+							
count	1258	1258	1258	1258	1258	1258	
1258							
mean	null	72.35785375357709	72.83938807631165	71.9186009594594	72.38844998012726	8222093.481717011	67.23883848728146
stddev	null	6.76809024470826	6.768186808159218	6.744075756255496	6.756859163732991	4519780.8431556	6.722609449996857
min	2012-01-03	56.389998999999996	57.060001	56.299999	56.419998	2094900	50.363689
max	2016-12-30	90.800003	90.970001	89.25	90.470001	80898100	84.91421600000001
-----+							
-----+							

Bonus Question!

There are too many decimal places for mean and stddev in the describe() dataframe. Format the numbers to just show up to two decimal places. Pay careful attention to the datatypes that .describe() returns, we didn't cover how to do this exact formatting, but we covered something very similar. [Check this link for a hint \(http://spark.apache.org/docs/latest/api/python/pyspark.sql.html#pyspark.sql.Column.cast\)](http://spark.apache.org/docs/latest/api/python/pyspark.sql.html#pyspark.sql.Column.cast)

If you get stuck on this, don't worry, just view the solutions.

```
In [27]: df.describe().show()
```

+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+									
-----+									
summary	Date	Open	High	Low	Close	Volume			
Adj Close									
+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+									
-----+									
count	1258	1258	1258	1258	1258	1258	1258		
1258									
mean	null	72.35785375357709	72.83938807631165	71.9186009594594	72.38844998012726	8222093.481717011	67.238	83848728146	
stddev	null	6.76809024470826	6.768186808159218	6.744075756255496	6.756859163732991	4519780.8431556	6.7226	09449996857	
min	2012-01-03	56.389998999999996	57.060001	56.299999	56.419998	2094900	50.363689		
max	2016-12-30	90.800003	90.970001	89.25	90.470001	80898100	84.914	21600000001	
+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+									
-----+									

```
In [28]: result = df.describe()
```

```
In [33]: from pyspark.sql.functions import format_number
```

```
In [44]: result.select(result['summary'],
                        format_number(result['Open'].cast('float'),2).alias('Open'),
                        format_number(result['High'].cast('float'),2).alias('High'),
                        format_number(result['Low'].cast('float'),2).alias('Low'),
                        format_number(result['Close'].cast('float'),2).alias('Close'),
                        format_number(result['Volume'].cast('float'),2).alias('Volume')).show()
```

summary	Open	High	Low	Close	Volume
count	1,258.00	1,258.00	1,258.00	1,258.00	1,258.00
mean	72.36	72.84	71.92	72.39	8,222,093.50
stddev	6.77	6.77	6.74	6.76	4,519,781.00
min	56.39	57.06	56.30	56.42	2,094,900.00
max	90.80	90.97	89.25	90.47	80,898,096.00

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.

```
In [48]: df_new = df.withColumn("HV Ratio",df['High']/df['Volume'])
```

```
In [55]: df_new.select('HV Ratio').show()
```

HV Ratio
4.819714653321546E-6
6.290848613094555E-6
4.669412994783916E-6
7.367338463826307E-6
8.915604778943901E-6
8.644477436914568E-6
9.351828421515645E-6
8.29141562102703E-6
7.712212102001476E-6
7.071764823529412E-6
1.015495466386981E-5
6.576354146362592...
5.90145296180676E-6
8.547679455011844E-6
8.420709512685392E-6
1.041448341728929...
8.316075414862431E-6
9.721183814992126E-6
8.029436027707578E-6
6.307432259386365E-6

only showing top 20 rows

What day had the Peak High in Price?

```
In [69]: df.orderBy(df['High'].desc()).head(1)[0][0]
Out[69]: '2015-01-13'
```

What is the mean of the Close column?

```
In [75]: df.agg({'Close': 'avg'}).show()
```

avg(Close)
72.38844998012726

What is the max and min of the Volume column?

```
In [85]: from pyspark.sql.functions import max, min
```

```
In [86]: df.select(max('Volume'),min('Volume')).show()
```

```
+-----+-----+
|max(Volume)|min(Volume)|
+-----+-----+
| 80898100| 2094900|
+-----+-----+
```

How many days was the Close lower than 60 dollars?

```
In [92]: df.filter(df['Close']<60).count()
```

```
Out[92]: 81
```

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

In other words, (Number of Days High>80)/(Total Days in the dataset)

```
In [117]: (df.filter(df['high'] > 80).count() / df.select(df['Date']).count() *100)
```

```
Out[117]: 9.141494435612083
```

What is the Pearson correlation between High and Volume?

Hint (<http://spark.apache.org/docs/latest/api/python/pyspark.sql.html#pyspark.sql.DataFrameStatFunctions.corr>).

```
In [120]: from pyspark.sql.functions import corr
df.select(corr(df['High'],df['Volume'])).show()
```

```
+-----+
| corr(High, Volume)|
+-----+
|-0.3384326061737161|
+-----+
```

What is the max High per year?

```
In [122]: from pyspark.sql.functions import year
yeardf = df.withColumn("Year",year(df["Date"]))
max_df = yeardf.groupBy('Year').max()
max_df.select('Year','max(High)').show()
```

```
+----+-----+
|Year|max(High)|
+----+-----+
|2015|90.970001|
|2013|81.370003|
|2014|88.089996|
|2012|77.599998|
|2016|75.190002|
+----+-----+
```

What is the average Close for each Calendar Month?

In other words, across all the years, what is the average Close price for Jan,Feb, Mar, etc... Your result will have a value for each of these months.

```
In [123]: from pyspark.sql.functions import month
monthdf = df.withColumn("Month",month("Date"))
monthavgs = monthdf.select("Month","Close").groupBy("Month").mean()
monthavgs.select("Month","avg(Close)").orderBy('Month').show()
```

+-----+-----+	
Month	avg(Close)
+-----+-----+	
1	71.44801958415842
2	71.306804443299
3	71.77794377570092
4	72.97361900952382
5	72.30971688679247
6	72.4953774245283
7	74.43971943925233
8	73.02981855454546
9	72.18411785294116
10	71.57854545454543
11	72.1110893069307
12	72.84792478301885
+-----+-----+	

Great Job!