```python
[1]: import findspark
     findspark.init()
     #inporting module to locate spark in VM

[2]: import pyspark
     from pyspark.sql.types import *
     #importing functions to run spark using python

[3]: sc = pyspark.SparkContext(appName="E10")
     #creates a spark context, using which we assigns the cluster of computer on which our code will run
```

Setting default log level to "WARN".
To adjust logging level use sc.setLogLevel(newLevel). For SparkR, use setLogLevel(newLevel).
23/10/31 18:10:08 WARN NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable

```python
[4]: rdd1 = sc.textFile("/home/hduser/spark/nsedata.csv")
     #importing csv file as rdd

[5]: rdd1 = rdd1.filter(lambda x: "SYMBOL" not in x)
     #this is used to remove header row

[6]: rdd2 = rdd1.map(lambda x : x.split(","))
     #this is used to create a new rdd from the rdd1 by splitting each record

[7]: # Helper comment!: The goal is to find out the mean of the OPEN prices and the mean of the CLOSE price in one batch of tasks ...

[8]: rdd_open = rdd2.map(lambda x : (x[0]+"_open",float(x[2])))
     rdd_close = rdd2.map(lambda x : (x[0]+"_close",float(x[5])))
     #this creates two new rdds which contain the opening and closing prices
     #x[2] and x[5] are basically the positions containing the opening and closing prices
```

```python
[3]: sc = pyspark.SparkContext(appName="E10")
     #creates a spark context, using which we assigns the cluster of computer on which our code will run
```

Setting default log level to "WARN".
To adjust logging level use sc.setLogLevel(newLevel). For SparkR, use setLogLevel(newLevel).
23/10/31 18:10:08 WARN NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable

```python
[4]: rdd1 = sc.textFile("/home/hduser/spark/nsedata.csv")
     #importing csv file as rdd

[5]: rdd1 = rdd1.filter(lambda x: "SYMBOL" not in x)
     #this is used to remove header row

[6]: rdd2 = rdd1.map(lambda x : x.split(","))
     #this is used to create a new rdd from the rdd1 by splitting each record

[7]: # Helper comment!: The goal is to find out the mean of the OPEN prices and the mean of the CLOSE price in one batch of tasks ...

[8]: rdd_open = rdd2.map(lambda x : (x[0]+"_open",float(x[2])))
     rdd_close = rdd2.map(lambda x : (x[0]+"_close",float(x[5])))
     #this creates two new rdds which contain the opening and closing prices
     #x[2] and x[5] are basically the positions containing the opening and closing prices

[9]: rdd_united = rdd_open.union(rdd_close)
     #here we are combining both the rdds into a single rdd to do their analysis

[10]: reducedByKey = rdd_united.reduceByKey(lambda x,y: x+y)
      #to find the average we first need to find the sum of all
      #hence, here we are basically summing up the opening and closing prices for each entry

[11]: temp1 = rdd_united.map(lambda x: (x[0],1)).countByKey()
      countOfEachSymbol = sc.parallelize(temp1.items())
      #finding out the number of times a given entry occured
      #we creates a new key value pair rdd
      #also we are counting the number of occurence of each stock using the "countbykey" function
      #Parallelize is a method to create an RDD from an existing collection (For e.g Array) present in the driver.
```

```python
[12]: symbol_sum_count = reducedByKey.join(countOfEachSymbol)
      #here we are combining the two rdds into a single one

[13]: averages = symbol_sum_count.map(lambda x : (x[0], x[1][0]/x[1][1]))
      # here we are finally finding the average of each of the symbols
      # so x[1]x[0]/x[1]x[1] is basically calculating sum of opening prices and closing prices of that particular symbol

[14]: averagesSorted = averages.sortByKey()
      #we are sorting the records present in the rdd based on symbol name
```

```python
[15]: averagesSorted.saveAsTextFile("/home/hduser/spark/averages")
      # saving the results in folder named averages
```

```python
[16]: sc.stop()
```

```
1 ('20MICRONS_close', 53.004122877930484)
2 ('20MICRONS_open', 53.32489894907032)
3 ('3IINFOTECH_close', 18.038803556992725)
4 ('3IINFOTECH_open', 18.17417138237672)
5 ('3MINDIA_close', 4520.343977364591)
6 ('3MINDIA_open', 4531.084518997574)
7 ('3RDROCK_close', 173.2137755102041)
8 ('3RDROCK_open', 173.18316326530612)
9 ('8KMILES_close', 480.73622047244095)
10 ('8KMILES_open', 481.63858267716535)
11 ('A2ZINFRA_close', 18.609433962264156)
12 ('A2ZINFRA_open', 18.73553459119497)
13 ('A2ZMES_close', 89.69389505549951)
14 ('A2ZMES_open', 90.46271442986883)
15 ('AANJANEYA_close', 441.84030249110316)
16 ('AANJANEYA_open', 440.93959074733095)
17 ('AARTIDRUGS_close', 312.94446240905427)
18 ('AARTIDRUGS_open', 312.832012934519)
19 ('AARTIIND_close', 127.70270816491507)
20 ('AARTIIND_open', 127.76463217461601)
```

```python
In [2]: import pyspark
        from pyspark.sql.types import *
```

```python
In [3]: sc = pyspark.SparkContext(appName="E10-01")
```

```
Setting default log level to "WARN".
To adjust logging level use sc.setLogLevel(newLevel). For SparkR, use setLogL
evel(newLevel).
23/10/28 17:12:24 WARN NativeCodeLoader: Unable to load native-hadoop library
for your platform... using builtin-java classes where applicable
```

```python
In [4]: rdd7 = sc.textFile("/home/hduser/spark/nsedata.csv")
```

```python
In [5]: rdd8 = rdd7.filter(lambda x: "SYMBOL" not in x)
```

```python
In [6]: rdd9 = rdd8.map(lambda x : x.split(","))
```

```python
In [17]: rdd_10 = rdd9.filter(lambda x: x[10] in "Apr" or  "May" or  "Jul" )
```

```python
In [18]: rdd_high_1 = rdd_10.map(lambda x : (x[0]+"_high",float(x[3])))
```

```python
In [19]: reducedByKey_1 = rdd_high_1.reduceByKey(lambda x,y: x+y)
```

```python
In [20]: temp_1 = rdd_high_1.map(lambda x: (x[0],1)).countByKey()
         countOfEachSymbol_1 = sc.parallelize(temp_1.items())
```

```python
In [21]: symbol_sum_count_1 = reducedByKey_1.join(countOfEachSymbol_1)
```

```python
In [22]: averages_high = symbol_sum_count_1.map(lambda x : (x[0], x[1][0]/x[1][1]))
```

```python
In [23]: averagesSorted = averages_high.sortByKey(ascending=False)
```

```python
In [24]: averagesSorted.saveAsTextFile("/home/hduser/spark/averages_high_4")
```