# ASBD Problem Set 8

5. Generate FIMs using FPgrowth algorithm on pyspark setup using a benchmark dataset available on the FIMI website.

```
] !pip install pyspark py4j

   Collecting pyspark
     Downloading pyspark-3.2.1.tar.gz (281.4 MB)
     |████████████████████████████████| 281.4 MB 33 kB/s
   Collecting py4j
     Downloading py4j-0.10.9.5-py2.py3-none-any.whl (199 kB)
     |████████████████████████████████| 199 kB 49.1 MB/s
     Downloading py4j-0.10.9.3-py2.py3-none-any.whl (198 kB)
     |████████████████████████████████| 198 kB 57.6 MB/s
   Building wheels for collected packages: pyspark
     Building wheel for pyspark (setup.py) ... done
     Created wheel for pyspark: filename=pyspark-3.2.1-py2.py3-none-any.whl size=281853642 sha256=43666033e484fcdd1c583847678ca34947c97d
     Stored in directory: /root/.cache/pip/wheels/9f/f5/07/7cd8017084dce4e93e84e92efd1e1d5334db05f2e83bcef74f
   Successfully built pyspark
   Installing collected packages: py4j, pyspark
   Successfully installed py4j-0.10.9.3 pyspark-3.2.1

] from pyspark.sql import SparkSession

   import numpy as np
   import pandas as pd
   data = pd.read_excel(DDDD.xlsx)
   data.head()
```

```python
import os
os.environ["JAVA_HOME"] = "/usr/lib/jvm/java-8-openjdk-amd64"
os.environ["SPARK_HOME"] = "/content/spark-3.0.0-bin-hadoop3.2"
import findspark
findspark.init()
from pyspark.sql import SparkSession
spark = SparkSession.builder.master("local[*]").getOrCreate()
spark = SparkSession
\.builder
\.appName('fpgrowth')
\.getOrCreate()spark
```

```python
from pyspark.sql import functions as F
from pyspark.ml.fpm import FPGrowth
import pandas
sparkdata = spark.createDataFrame(data)
```

```python
basketdata = sparkdata.dropDuplicates(['SalesTransactionID',
'SalesItem']).sort('SalesTransactionID')
```

```
basketdata =
basketdata.groupBy("SalesTransactionID").agg(F.collect_list("SalesItem")).
sort('SalesTransactionID')
```

```
+-------------------+----+
|              items|freq|
+-------------------+----+
|              [257]| 432|
|               [20]|2837|
|              [104]|2417|
|          [104, 20]| 981|
|             [1491]| 432|
|              [110]|2172|
|         [110, 104]| 745|
|     [110, 104, 20]| 476|
|          [110, 20]| 765|
|             [1495]| 431|
|              [103]|2123|
|         [103, 110]| 671|
|    [103, 110, 104]| 445|
|[103, 110, 104, 20]| 348|
|     [103, 110, 20]| 444|
|         [103, 104]| 885|
|     [103, 104, 20]| 572|
|          [103, 20]| 861|
|              [179]| 431|
|               [67]|1975|
+-------------------+----+
only showing top 20 rows
```

```
+------------+----------+------------------+------------------+
|  antecedent|consequent|        confidence|              lift|
+------------+----------+------------------+------------------+
|       [128]|      [67]| 0.3379978471474704|  8.28753607390552|
|       [128]|      [91]|0.34230355220667386|10.666918802548512|
|       [128]|     [104]|  0.38751345532831| 7.764057338737584|
|       [128]|      [92]| 0.3153928955866523| 9.648273128034885|
|       [128]|     [103]|0.35522066738428415| 8.102645331489093|
|[91, 92, 67]|     [104]| 0.7577197149643705| 15.18135495112313|
|[91, 92, 67]|     [103]| 0.7197149643705463|16.416823770423022|
|   [83, 103]|      [67]| 0.519327731092437|12.733653015636635|
|   [83, 103]|     [110]| 0.5210084033613446|11.616184595385116|
|   [83, 103]|     [104]| 0.5495798319327732|11.011151403051912|
|   [83, 103]|      [20]| 0.6201680672268908|10.585921333637438|
|   [83, 103]|     [120]| 0.5277310924369748| 14.42206878236622|
|   [83, 103]|     [108]| 0.5478991596638656|16.031761151590544|
|       [316]|     [514]| 0.3029087261785356|10.440325960086666|
|       [316]|      [83]|0.35907723169508526| 8.804391909906936|
|       [316]|      [67]|0.30090270812437314| 7.377982047408047|
|       [316]|     [104]| 0.3350050150451354| 6.7120202145534655|
|       [316]|     [103]|0.32898696088264795| 7.504249914132412|
|       [316]|      [63]| 0.5536609829488466|20.250443172417555|
|       [316]|      [64]| 0.506519558676028|21.109050041691336|
+------------+----------+------------------+------------------+
only showing top 20 rows
```

```
+------------------+---------------------+--------------------+
|SalesTransactionID|collect_list(SalesItem)|          prediction|
+------------------+---------------------+--------------------+
|                 0|                  [0]|                  []|
|                 1|            [0, 1, 2]|                  []|
|                 2|                  [1]|                  []|
|                 3|                  [0]|                  []|
|                 4|                  [0]|                  []|
|                 5|                  [0]|                  []|
|                 6|                  [2]|                  []|
|                 7|                  [2]|                  []|
|                 8|                  [0]|                  []|
|                10|               [1, 0]|                  []|
|                11|                  [0]|                  []|
|                12|               [4, 3]|                  []|
|                13|                  [5]|                  []|
|                15|            [7, 8, 9]|[83, 20, 108, 514...|
|                16|     [17, 15, 16, 20, ...|[101, 104, 110, 1...|
|                17|                 [21]|                  []|
|                18|                 [22]|                  []|
|                19|     [33, 41, 27, 34, ...|                  []|
|                20|     [59, 50, 57, 60, ...|                  []|
|                21|                 [62]|                  []|
+------------------+---------------------+--------------------+
only showing top 20 rows
```

3. Compute correlation between the given two series using Pearson's and Spearman's Method.
a. (Use the Spark MLlib libraries and helper functions available)
i. Series A: 35, 23, 47, 17, 10, 43, 9, 6, 28
ii. Series B: 30, 33, 45, 23, 8, 49, 12, 4, 31

```
[2]  from pyspark.sql import SparkSession
```

```
[8]  import numpy as np
     import pandas as pd
```

```
     from pyspark.ml.stat import Correlation
     from pyspark.ml.linalg import DenseMatrix, Vectors
     from pyspark.ml.feature import VectorAssembler
     from pyspark.sql.functions import *
```

```
[11] spark = SparkSession.builder.appName('demo').master('local').enableHiveSupport().getOrCreate()
```

```
[12] spark
```

```
5]
    features = df.rdd.map(lambda row: row[0:])

    from pyspark.mllib.stat import Statistics

    corr_mat=Statistics.corr(features, method="pearson")
```

```
5] corr_mat

   array([[ 1., -1.],
          [-1.,  1.]])
```

```
7] corr_mat1=Statistics.corr(features, method="spearman")
```

```
8] corr_mat1

   array([[ 1., -1.],
          [-1.,  1.]])
```

2. Randomly populate 1000 numbers and calculate mean, variance, standard deviation for the generated data.

```
[29] data1 = pd.read_csv('rand.csv')

[30] df1 = spark.createDataFrame(data1)

[31] from pyspark.sql.functions import mean as _mean, stddev as _stddev, col

[32]  df1.select(mean('list1')).collect()

     [Row(avg(list1)=984.2171816977666)]

▶   df1.select(variance('list1')).collect()

     [Row(var_samp(list1)=332657.1194870868)]

[34]  df1.select(stddev('list1')).collect()

     [Row(stddev_samp(list1)=576.7643535163098)]

[ ]
```

*(Question 1 and 4 are not working as pyspark master throws an error which i was unable to resolve )*

```
paleti@paleti-Lenovo-ideapad-330-15ICH:~$ start-master.sh
starting org.apache.spark.deploy.master.Master, logging to /opt/spark/logs/spark-paleti-org.apache.spark.de
ploy.master.Master-1-paleti-Lenovo-ideapad-330-15ICH.out
failed to launch: nice -n 0 /opt/spark/bin/spark-class org.apache.spark.deploy.master.Master --host paleti-
Lenovo-ideapad-330-15ICH --port 7077 --webui-port 8080
  /opt/spark/bin/spark-class: line 71: /usr/local/java/jdk1.8.0_121/bin/java: No such file or directory
  /opt/spark/bin/spark-class: line 96: CMD: bad array subscript
full log in /opt/spark/logs/spark-paleti-org.apache.spark.deploy.master.Master-1-paleti-Lenovo-ideapad-330-
15ICH.out
```