**Week 4 Assignment**

Guruprasad Velikadu Krishnamoorthy

College of Science and Technology, Bellevue University

**DSC650**-T301: Big Data

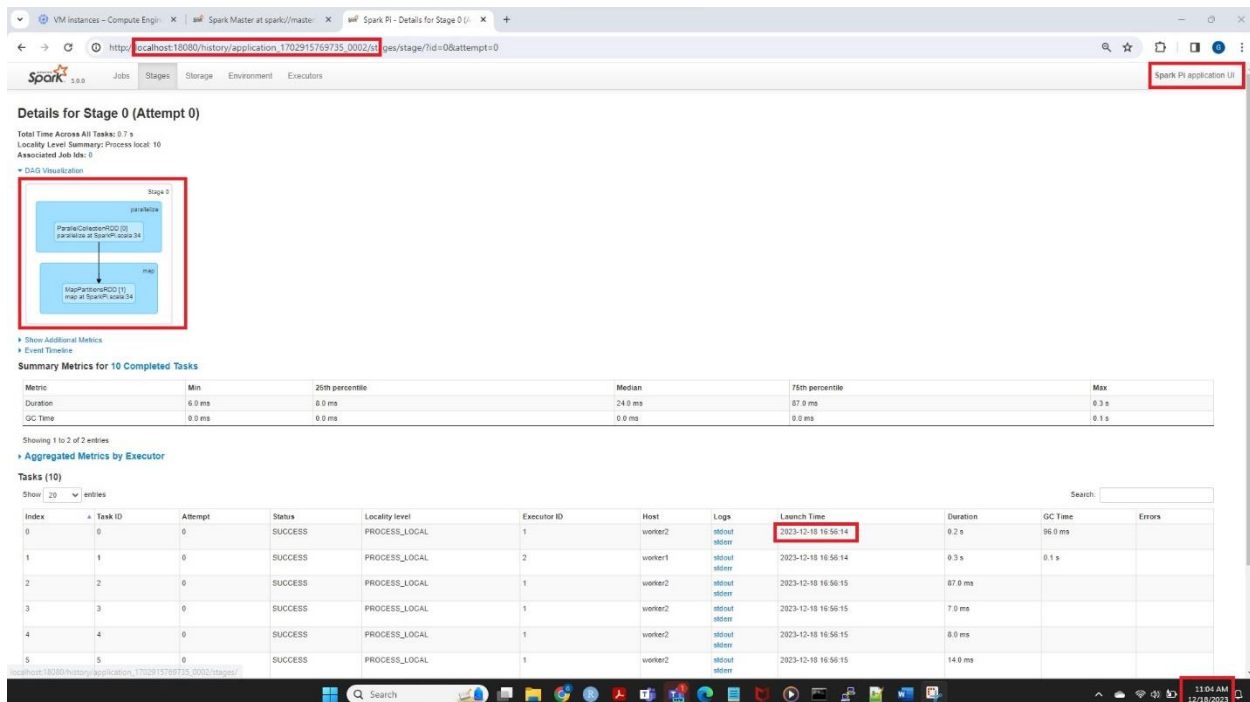Professor. Nasheb Ismaily

December 18, 2023

## 1. SparkPi Program:

### 1.1 Screenshot of SparkPi program:



### 1.2 Output from Spark History Server:

### 1.3 Significance of the SparkPi program:

In the Monte Carlo method for estimating pi, random points are generated within a square, and the number of points falling inside a quarter-circle inside that square is counted. The ratio of the points inside the circle to the total points is then used to estimate the pi, based on the formula of the area of the circle and square.

In SparkPi Program, x and y coordinates of the "n" random points are calculated, and the number of points where $x*2+y*2 < 1$ is calculated. Spark parallelize creates a RDD and distributes the calculation between the executors and returns the result. The Pi is then calculated by using the formula $4 * count/ n$.

Note, in our example the value of n passed is 10 and the program is executed in "client" mode with 2 GB of driver memory, 1 GB of executor memory, and 1 core. Yarn is specified as the Resource manager.

```
spark-submit --class org.apache.spark.examples.SparkPi \
    --master yarn \
    --deploy-mode client \
    --driver-memory 2g \
    --executor-memory 1g \
    --executor-cores 1 \
    $SPARK_HOME/examples/jars/spark-examples*.jar \
    10
```

### 2. Screenshot of the first 100 generated random numbers.

3. **Screenshots of Custom Transformation and their significance:**

Due to space constraints while taking the screenshot of the output, only 10 rows are displayed for each transformation. 7 transformations were used, and their significance is listed below:

**3.1. toUpperCase**:

val transformedSentences = sentencesRDD.map(sentence => sentence.**toUpperCase**)

**Significance**:

This command will convert all the characters to Upper case and return the Uppercase converted values of the RDD. Similarly, **toLowerCase** can also be used.

**3.2. startsWith with Map method:**

val transformedSentences2 = sentencesRDD.**map**(sentence => sentence.startsWith("appl"))

**Significance:**

This command returns a Boolean value if the Array starts with the specified value. In the example, all arrays that start with "appl" returns True and everything else returns False.

**3.3. startsWith used with filter method:**

val transformedSentences2 = sentencesRDD.**filter**(sentence => sentence.**startsWith**("appl"))

**Significance:**

This is similar to the command above; the difference is it applies the filters and returns the array elements that start with the specified value.

### 3.4. Contains:
val transformedSentences3 = sentencesRDD.filter(sentence => sentence.contains("erry"))
**Significance:**
This command returns the list of elements that contains the string **"erry"**

### 3.5. substring:
val transformedSentences4 = sentencesRDD.map(sentence => sentence.substring(0,5))
**Significance:**
This transformation returns the substring of the element from position 0 to 4.

### 3.6. Replace:
val transformedSentences5 = sentencesRDD.map(sentence => sentence.replace(" ","-"))
**Significance:**
This transformation replaces the spaces with "-" in the Array of elements.

### 3.7. indexOf:
val transformedSentences6 = sentencesRDD.map(sentence => sentence.indexOf("fig"))
**Significance:**
This transformation returns the position of the string "fig" in the array. If the string does not exist, it returns -1.

```
scala> val transformedSentences3 = sentencesRDD.filter(sentence => sentence.contains("erry"))
transformedSentences3: org.apache.spark.rdd.RDD[String] = MapPartitionsRDD[12] at filter at <console>:25

scala> transformedSentences3.take(10).foreach(println)
banana elderberry honeydew.
elderberry cherry date apple fig grape.
elderberry fig grape cherry date.
fig cherry banana grape date honeydew.
banana apple elderberry cherry.
fig elderberry date grape banana.
honeydew date apple elderberry banana cherry.
elderberry fig cherry apple.
elderberry apple date cherry.
elderberry honeydew apple fig banana date.

scala>

scala> val transformedSentences4 = sentencesRDD.map(sentence => sentence.substring(0,5))
transformedSentences4: org.apache.spark.rdd.RDD[String] = MapPartitionsRDD[13] at map at <console>:25

scala> transformedSentences4.take(10).foreach(println)
banan
elder
elder
fig c
banan
fig e
honey
date
elder
elder

scala> val transformedSentences5 = sentencesRDD.map(sentence => sentence.replace(" ","-"))
transformedSentences5: org.apache.spark.rdd.RDD[String] = MapPartitionsRDD[14] at map at <console>:25

scala> transformedSentences5.take(10).foreach(println)
banana-elderberry-honeydew.
elderberry-cherry-date-apple-fig-grape.
elderberry-fig-grape-cherry-date.
fig-cherry-banana-grape-date-honeydew.
banana-apple-elderberry-cherry.
fig-elderberry-date-grape-banana.
honeydew-date-apple-elderberry-banana-cherry.
date-fig.
elderberry-fig-cherry-apple.
elderberry-apple-date-cherry.

scala> val transformedSentences6 = sentencesRDD.map(sentence => sentence.indexOf("fig"))
transformedSentences6: org.apache.spark.rdd.RDD[Int] = MapPartitionsRDD[15] at map at <console>:25

scala> transformedSentences6.take(10).foreach(println)
-1
29
11
0
-1
0
-1
5
11
-1

scala>
```