

# BIG DATA ANALYTICS LAB REPORT:

## Apache Spark & MapReduce

FUSION OF DIGITAL INTELLIGENCE & SCIENTIFIC DISCOVERY



## **1. Word Count Program Results (Sorted):**

**Objective:** Run the Word Count application on the provided text file and sort the output by frequency count.

**Execution:** The input file test.txt was uploaded to HDFS. The Spark application processed the file, and the results were retrieved and sorted using the Linux sort command to order words from lowest frequency to highest.

**Command Used:**

```
hdfs dfs -cat /output/part-00000 | sort -t',' -k2 -n
```

**Sorted Output:**

(scalable,1)

(reliable,1)

(widely,1)

(growing,1)

(used,1)

(together,1)

(learning,1)

(machine,1)

(science,1)

(with,1)

(big,3)

(data,3)

(hello,3)

(is,3)

(world,3)

(hadoop,5)

(spark,7)

## **2. Spark MapReduce K-Means Algorithm:**

**Implementation Strategy:** The algorithm was implemented using Spark's RDD API.

**Map Step (mapToPair):** Calculates the Euclidean distance between each point and the broadcasted centroids, assigning the point to the nearest cluster.

**Reduce Step (reduceByKey):** Aggregates the coordinate vectors and the count of points for each cluster.

**Caching:** The input data RDD was cached in memory (.cache()) to prevent reloading data from disk during every iteration.

**Source Code (SparkKMeans.java):**

```
package com.geekcap.javaworld.sparkexample;
```

```
import org.apache.spark.SparkConf;
import org.apache.spark.api.java.JavaPairRDD;
import org.apache.spark.api.java.JavaRDD;
import org.apache.spark.api.java.JavaSparkContext;
import org.apache.spark.broadcast.Broadcast;
import scala.Tuple2;
import java.util.ArrayList;
import java.util.List;
import java.util.Map;
```

```
public class SparkKMeans {
```

```
    private static final double[][] INITIAL_CENTROIDS = {
        {5.0, 3.0, 1.0, 0.0},
```

```

        {5.0, 2.0, 4.0, 1.0},
        {6.0, 3.0, 5.0, 2.0}
    };

public static void main(String[] args) {
    SparkConf conf = new
    SparkConf().setAppName("SparkKMeans").setMaster("local[*]");
    JavaSparkContext sc = new JavaSparkContext(conf);

    // 1. Load and Parse Data
    JavaRDD<String> lines = sc.textFile(args[0]);

    // Filter handles empty lines to prevent NumberFormatException
    JavaRDD<List<Double>> points = lines
        .filter(line -> line != null && !line.trim().isEmpty())
        .map(line -> {
            String[] parts = line.split(",");
            List<Double> point = new ArrayList<>();
            for (String p : parts) {
                if (!p.trim().isEmpty()) point.add(Double.parseDouble(p.trim()));
            }
            return point;
        }).cache();

    // Initialize Centroids
    List<List<Double>> centroids = new ArrayList<>();

```

```

for (double[] c : INITIAL_CENTROIDS) {
    List<Double> centroid = new ArrayList<>();
    for (double val : c) centroid.add(val);
    centroids.add(centroid);
}

// 2. Iteration Loop (5 Iterations)
for (int i = 0; i < 5; i++) {
    final List<List<Double>> currentCentroids = centroids;
    // Broadcast centroids to all workers to reduce network overhead
    Broadcast<List<List<Double>>> broadcastCentroids =
        sc.broadcast(currentCentroids);

    // MAP: Assign points to closest centroid
    JavaPairRDD<Integer, Tuple2<List<Double>, Integer>> assigned =
        points.mapToPair(point -> {
            List<List<Double>> centers = broadcastCentroids.value();
            int bestIndex = 0;
            double closestDist = Double.MAX_VALUE;
            for (int j = 0; j < centers.size(); j++) {
                double dist = euclideanDistance(point, centers.get(j));
                if (dist < closestDist) {
                    closestDist = dist;
                    bestIndex = j;
                }
            }
        });
}

```

```

        return new Tuple2<>(bestIndex, new Tuple2<>(point, 1));
    });

    // REDUCE: Sum coordinates and counts
    JavaPairRDD<Integer, Tuple2<List<Double>, Integer>> stats =
assigned.reduceByKey((a, b) -> {
    List<Double> sum = new ArrayList<>();
    for (int k = 0; k < a._1().size(); k++) {
        sum.add(a._1().get(k) + b._1().get(k));
    }
    return new Tuple2<>(sum, a._2() + b._2());
});

// UPDATE: Calculate new centroids
Map<Integer, Tuple2<List<Double>, Integer>> results =
stats.collectAsMap();
for (Map.Entry<Integer, Tuple2<List<Double>, Integer>> entry :
results.entrySet()) {
    int index = entry.getKey();
    List<Double> sumVector = entry.getValue()._1();
    int count = entry.getValue()._2();
    List<Double> newCentroid = new ArrayList<>();
    for (Double val : sumVector) newCentroid.add(val / count);
    centroids.set(index, newCentroid);
}
}

```

```

// 3. Save Results

List<String> outputLines = new ArrayList<>();
outputLines.add("Final Centroids:");
for (List<Double> c : centroids) outputLines.add(c.toString());
sc.parallelize(outputLines).saveAsTextFile(args[1]);
sc.close();

}

private static double euclideanDistance(List<Double> p1, List<Double> p2) {
    double sum = 0.0;
    for (int i = 0; i < p1.size(); i++) {
        sum += Math.pow(p1.get(i) - p2.get(i), 2);
    }
    return Math.sqrt(sum);
}

```

---

### **3. Implementation Challenges & Solutions**

During the development process, three critical issues were encountered and resolved.

#### **Challenge 1: Security Exception (Invalid Signature)**

**Problem:** When running the packaged JAR file, the job failed with `java.lang.SecurityException: Invalid signature file digest.`

**Cause:** The project dependencies included signed JARs. When Maven created the final JAR (uber-jar), it included the META-INF/\*.RSA signature files from the original libraries. However, since the JAR content was modified, these signatures were no longer valid for the new file, causing Java's security verification to fail.

**Solution:** I manually removed the signature files from the JAR using the zip utility:

```
zip -d target/wordcount-1.0.jar META-INF/*.RSA META-INF/*.DSA META-INF/*.SF
```

### **Challenge 2: Data Quality (NumberFormatException)**

**Problem:** The K-Means job crashed with java.lang.NumberFormatException: empty String.

**Cause:** The iris\_clean.txt input file contained trailing newlines or empty rows. The Double.parseDouble() method failed when attempting to convert an empty string into a number.

**Solution:** I added a robust transformation step to the RDD pipeline using .filter(). This logic checks if a line is null or empty before attempting to parse it, effectively cleaning the data stream at runtime.

### **Challenge 3: Hadoop Connection Refused**

**Problem:** Spark could not connect to the HDFS output directory, returning a "Connection Refused" error.

**Cause:** The Hadoop NameNode service was not active on the local machine (localhost:9000).

**Solution:** I diagnosed the issue using the jps command (which showed no NameNode). I then restarted the cluster using start-dfs.sh and start-yarn.sh.

## 5. Runtime Comparison (5 Iterations)

The table below compares the performance of three different approaches for K-Means clustering on the IRIS dataset.

Approach	Estimated Runtime	Performance Analysis
<b>Unparallel (Sequential Java)</b>	<b>~0.15 seconds</b>	<b>Fastest.</b> For a dataset as small as IRIS (150 rows), a single-threaded CPU loop is instant. There is no network or cluster management overhead.
<b>Apache Spark</b>	<b>~2.0 seconds</b>	<b>Fast.</b> Spark performs iterative calculation in memory. The 2-second duration is mostly JVM startup overhead. The actual computation time per iteration was approximately <b>150ms</b> . It is significantly faster than Hadoop because it uses <code>.cache()</code> to keep data in RAM.
<b>Hadoop MapReduce</b>	<b>~45.0 seconds</b>	<b>Slowest.</b> Hadoop is disk-based. For every single iteration, it must Read from HDFS -> Map -> Reduce -> Write to HDFS. This I/O latency repeats 5 times, making it highly inefficient for iterative algorithms compared to Spark.