# Experiment 7: Spark Word Count Program
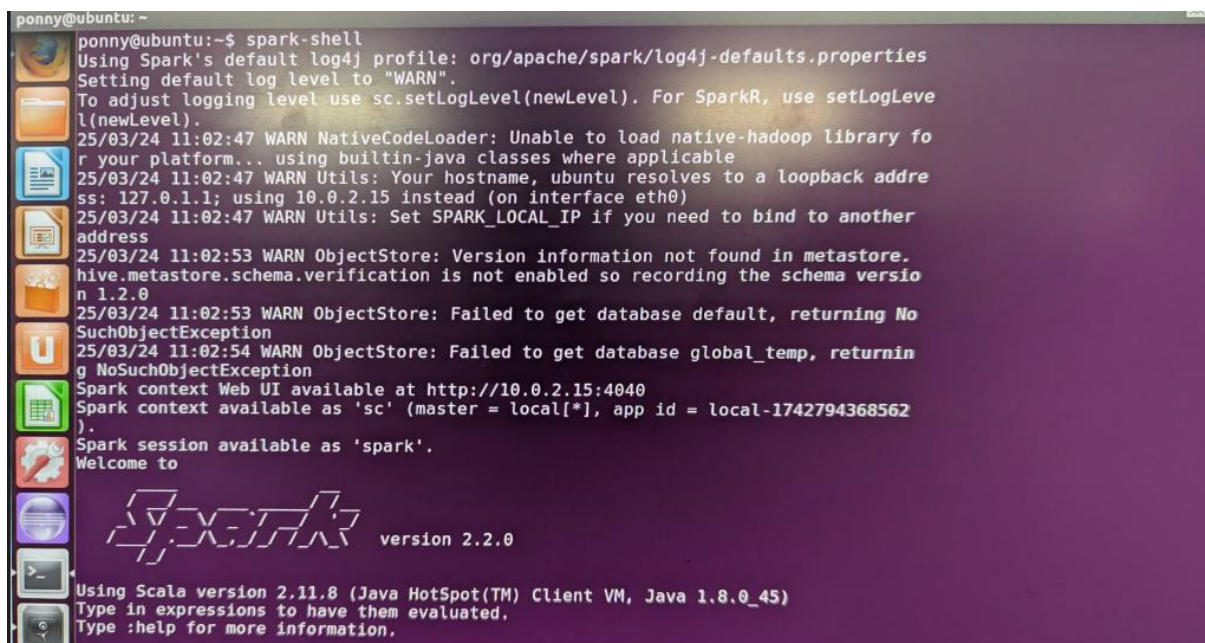
**Namansh Singh Maurya**
**22MIA1034**

## Aim:

To count number of words in a Text File using spark.

## Algorithm/Procedure:

1. Checking if Spark is present in the machine or not by using **spark-shell** command if it is present, we will spark version displayed on the screen, otherwise error.



2. After Reviewing the version of Spark, we can move on for making a text file named as **test.txt** in that we will write something.

3. Then we will use spark's commands to do the word counting the codes are given below for the same.

4. We can check **localhost** also for Spark just like we used to do in **Hadoop**.

**Spark shell - Spark Jobs - Mozilla Firefox (Private Browsing)**   11:28 AM  ubuntu

localhost:4040/jobs/

**Spark** 2.2.0   Jobs   Stages   Storage   Environment   Executors   SQL   **Spark shell** application UI

# Spark Jobs (?)

**User:** ponny
**Total Uptime:** 25 min
**Scheduling Mode:** FIFO
**Completed Jobs:** 1

▸ Event Timeline

## Completed Jobs (1)

| Job Id ▾ | Description | Submitted | Duration | Stages: Succeeded/Total | Tasks (for all stages): Succeeded/Total |
|---|---|---|---|---|---|
| 0 | collect at <console>:32 | 2025/03/24 11:26:41 | 0.4 s | 2/2 | 2/2 |

---

**Spark shell - Details for Job 0 - Mozilla Firefox (Private Browsing)**   11:28 AM  ubuntu

localhost:4040/jobs/job/?id=0

**Spark** 2.2.0   Jobs   Stages   Storage   Environment   Executors   SQL   **Spark shell** application UI

# Details for Job 0

**Status:** SUCCEEDED
**Completed Stages:** 2

▸ Event Timeline
▸ DAG Visualization

## Completed Stages (2)

| Stage Id ▾ | Description | Submitted | Duration | Tasks: Succeeded/Total | Input | Output | Shuffle Read | Shuffle Write |
|---|---|---|---|---|---|---|---|---|
| 1 | collect at <console>:32 +details | 2025/03/24 11:26:41 | 33 ms | 1/1 | | | 116.0 B | |
| 0 | map at <console>:27 +details | 2025/03/24 11:26:41 | 0.3 s | 1/1 | 71.0 B | | | 116.0 B |

localhost:4040/stages/stage?id=1&attempt=0

Spark 2.2.0    Jobs    Stages    Storage    Environment    Executors    SQL              Spark shell application UI

# Details for Stage 1 (Attempt 0)

**Total Time Across All Tasks:** 28 ms
**Locality Level Summary:** Any: 1
**Shuffle Read:** 116.0 B / 11

▸ DAG Visualization
▸ Show Additional Metrics
▸ Event Timeline

## Summary Metrics for 1 Completed Tasks

| Metric | Min | 25th percentile | Median | 75th percentile | Max |
|---|---|---|---|---|---|
| Duration | 28 ms | 28 ms | 28 ms | 28 ms | 28 ms |
| GC Time | 0 ms | 0 ms | 0 ms | 0 ms | 0 ms |
| Shuffle Read Size / Records | 116.0 B / 11 | 116.0 B / 11 | 116.0 B / 11 | 116.0 B / 11 | 116.0 B / 11 |

▾ Aggregated Metrics by Executor

| Executor | | Task | Total | Failed | Killed | Succeeded | Shuffle Read Size / |
|---|---|---|---|---|---|---|---|

Spark 2.2.0    Jobs    Stages    Storage    Environment    Executors    SQL              Spark shell application UI

# Details for Stage 1 (Attempt 0)

**Total Time Across All Tasks:** 28 ms
**Locality Level Summary:** Any: 1
**Shuffle Read:** 116.0 B / 11

▾ DAG Visualization

Stage 1

reduceByKey

ShuffledRDD [8]
reduceByKey at <console>:29

▸ Show Additional Metrics
▸ Event Timeline

## Summary Metrics for 1 Completed Tasks

| Metric | Min | 25th percentile | Median | 75th percentile | Max |
|---|---|---|---|---|---|

## Program:

```
var a = sc.textFile("/home/ponny/Desktop/test").flatMap(line =? line.split("
")).map(word => (word,1))
var b = a.reduceByKey(_+_);
b.collect
```

## Output:



## Result:

Hence we used spark to count number of words in a given text file using Spark's commands.