

**Project : PageRank using Scala and PySpark**

**Student : ID**

**SFBU- 19599, Manickam Ravisekar - MSCS  
San Francisco Bay University , Fremont,CA, USA**

Professor : Dr Henry Chung

TA: LIANG GU

## **Contents**

**Abstract**

**PageRank GRAPH and Matrix**

**PageRank Formula**

**Google Cloud Setup**

**Scala program to find the Iterations values**

**PySpark program to find the Iteration values**

**Conclusion**

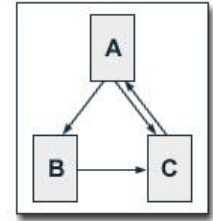
## **Abstract**

In this project to learn basic graph in Pyspark and Scala used in Big Data to find the PageRank of a given graph. The primary learning goal of the project is to gain familiarity with the syntax, data structures to learn scala , pyspark. Also learning the computation involved in finding the pagerank of a given graph.

Thank full to professor Dr Henry who encouraged us to work on this assignment on google cloud platform.

## Adjacency Matrix of the Graph

Row-Column	A	B	C	No of Links
A	-	1	1	2
B	-	-	1	1
C	1	-	-	1



Page Rank Iterations Values

Iteration	PR(A)	PR(B)	PR(C)
0	1	1	1
1	1	0.575	1.425
2	1.36125	0.575	1.06375
3	1.0541875	0.72853125	1.21728125

**Process of Calculating PageRank** : Initialize each page's rank to 1.0

On each iteration, have page p send a contribution of  $\text{rank}(p) / \text{numNeighbors}(p)$  to its neighbors (the pages it has links to). Set each page's rank to  $0.15 + 0.85 * \text{contributionsReceived}$ .

Note: 0.85 is the damping factor

### **PageRank overview**

If The initial PageRank value for each webpage is 1.

$\text{PR}(A) = 1$   $\text{PR}(B) = 1$   $\text{PR}(C) = 1$

Page B has a link to pages C and A ,Page C has a link to page A ,Page D has links to all three pages

And A's PageRank is

$\text{PR}(A) = (1-d) + d * (\text{PR}(B) / 2 + \text{PR}(C) / 1 + \text{PR}(D) / 3)$

B's PageRank is

$\text{PR}(B) = (1-d) + d * (\text{PR}(D) / 3)$

C's PageRank is

$\text{PR}(C) = (1-d) + d * (\text{PR}(B) / 2 + \text{PR}(D) / 3)$

D's PageRank is

$\text{PR}(D) = 1-d$

Damping factor is 0.85

Then after 1st iteration

Output

Page B would transfer half of its existing value, or 0.5, to page A and the other half, or 0.5, to page C.

Page C would transfer all of its existing value, 1, to the only page it links to, A.

Since D had three outbound links, it would transfer one third of its existing value, or approximately 0.33, to A.

Input

$\text{PR}(A)$

$= (1-d) + d * (\text{PR}(B) / 2 + \text{PR}(C) / 1 + \text{PR}(D) / 3)$

$= (1-0.85) + 0.85 * (0.5 + 1 + 0.33)$

$= 1.71$

$\text{PR}(B)$

$= (1-d) + d * (\text{PR}(D) / 3)$

$= (1-0.85) + 0.85 * 0.33$

$= 0.43$

$\text{PR}(C)$

$= (1-d) + d * (\text{PR}(B) / 2 + \text{PR}(D) / 3)$

$= (1-0.85) + 0.85 * (0.5 + 0.33)$

$= 0.86$

$\text{PR}(D)$

$= 1-d$

$= 0.15$

# Cluster creation of Google Cloud Platform

The screenshot shows the Google Cloud Platform dashboard for a project named 'scala-3'. The browser address bar displays the URL <https://console.cloud.google.com/home/dashboard?project=scala-3>. The dashboard header includes a navigation menu with 'DASHBOARD', 'ACTIVITY', and 'RECOMMENDATIONS', and a search bar. The main content area is divided into several sections:

- Project info:** Displays project details for 'scala-3', including the Project name, Project number (179637453175), and Project ID (scala-3). It includes a link to 'ADD PEOPLE TO THIS PROJECT' and a button to 'Go to project settings'.
- Resources:** Lists various Google Cloud services available in the project, such as BigQuery, SQL, Compute Engine, Storage, Cloud Functions, and App Engine.
- API APIs:** A section for API usage, showing a graph of 'Requests (requests/sec)' over time. A message indicates 'No data is available for the selected time frame.' and a link to 'Go to APIs overview' is provided.
- Google Cloud Platform status:** Provides information about the Google Compute Engine status, noting that VMs using Local SSD are experiencing intermittent terminations. It includes a link to 'Go to Cloud status dashboard'.
- Billing:** Shows the estimated charges for the billing period Oct 1 - 31, 2022, with a total of USD \$0.00. It includes a link to 'View detailed charges'.
- Monitoring:** Offers options to 'Create my dashboard', 'Set up alerting policies', and 'Create uptime checks'.

The dashboard also features a top navigation bar with a 'Sign in' button and a user profile icon.

Snoozed - manickam@student.s...Student IndexRPI APIs & Services - APIs & Service...

https://console.cloud.google.com/apis/dashboard?project=scala-3&show=all

Start your Free Trial with \$300 in credit. Don't worry—you won't be charged if you run out of credits. [Learn more](#)

DISMISSACTIVATE

Google Cloudscala-3Search for resources, docs, products, and more

API APIs & Services

Enabled APIs & services

Library

Credentials

OAuth consent screen

Domain verification

Page usage agreements

APIs & Services+ ENABLE APIS AND SERVICES

1 hour6 hours12 hours1 day2 days4 days7 days14 days30 days

Traffic

No data is available for the selected time frame.

Errors

No data is available for the selected time frame.

Median latency

No data is available for the selected time frame.

Filter

Filter

Name	Requests	Errors (%)	Latency, median (ms)	Latency, 95% (ms)
<a href="#">BigQuery API</a>				
<a href="#">BigQuery Migration API</a>				
<a href="#">BigQuery Storage API</a>				
<a href="#">Cloud Datastore API</a>				
<a href="#">Cloud Debugger API</a>				
<a href="#">Cloud Logging API</a>				
<a href="#">Cloud Monitoring API</a>				
<a href="#">Cloud SQL</a>				

Snoozed - manickam@student.s...Student IndexRPI APIs & Services - APIs & Service...+https://console.cloud.google.com/apis/dashboard?project=scala-3&show=allA<sup>®</sup>🔖🔄🌟🔒Sign in...DISMISSACTIVATE

Google Cloudscala-3Search for resources, docs, products, and moreQ Search🗨🔔🔗⋮M

RPI APIs & ServicesAPIs & Services+ ENABLE APIS AND SERVICES

Enabled APIs & servicesLibraryCredentialsOAuth consent screenDomain verificationPage usage agreements

1 hour6 hours12 hours✓ 1 day2 days4 days7 days14 days30 days

Traffic

No data is available for the selected time frame.

Errors

No data is available for the selected time frame.

Median latency

No data is available for the selected time frame.

FilterFilter?

Name	Requests	Errors (%)	Latency, median (ms)	Latency, 95% (ms)
<a href="#">BigQuery API</a>				
<a href="#">BigQuery Migration API</a>				
<a href="#">BigQuery Storage API</a>				
<a href="#">Cloud Datastore API</a>				
<a href="#">Cloud Debugger API</a>				
<a href="#">Cloud Logging API</a>				
<a href="#">Cloud Monitoring API</a>				
<a href="#">Cloud SQL</a>				



📁 Start your Free Trial with \$300 in credit. Don't worry—you won't be charged if you run out of credits. [Learn more](#)

DISMISS

ACTIVATE

☰ Google Cloud scala-3 ▾

🔍 📄 🔔 ? ⋮ M



## Cloud Dataproc API

[Google Enterprise API](#)

Manages Hadoop-based clusters and jobs on Google Cloud Platform.



TRY THIS API ↗

OVERVIEW

DOCUMENTATION

### Overview

Manages Hadoop-based clusters and jobs on Google Cloud Platform.

### Additional details

Type: [SaaS & APIs](#)

Last updated: 7/21/22

Category: [Google Enterprise APIs](#)

Service name: dataproc.googleapis.com

### Tutorials and documentation

[Learn more](#) ↗

Terms of Service

Snoozed - manickam@student.s x Student Index Clusters - Dataproc - scala-3 x

https://console.cloud.google.com/dataproc/clusters?project=scala-3

Start your Free Trial with \$300 in credit. Don't worry—you won't be charged if you run out of credits. [Learn more](#)

DISMISS ACTIVATE

Google Cloud scala-3 Search for resources, docs, products, and more Search

Dataproc

Jobs on Clusters

Clusters

Jobs

Workflows

Autoscaling policies

Serverless

Batches

Metastore Services

Metastore

Federation

Utilities

Component exchange

Workbench

Release Notes

Clusters

CREATE CLUSTER REFRESH START STOP DELETE REGIONS + 5 RECOMMENDED ALERTS

Cluster

Cloud Dataproc

Google Cloud Dataproc lets you provision Apache Hadoop clusters and connect to underlying analytic data stores.

There are no clusters in the currently selected Cloud Dataproc region(s). Create a cluster to get started.

CREATE CLUSTER

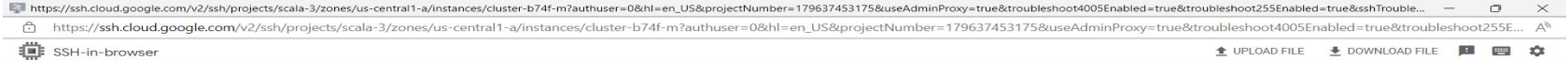


<b>Dataproc</b>	<b>Clusters</b> <a href="#">CREATE CLUSTER</a> <a href="#">REFRESH</a> <a href="#">START</a> <a href="#">STOP</a> <a href="#">DELETE</a> <a href="#">REGIONS</a> <a href="#">+ 5 RECOMMENDED ALERTS</a> <a href="#">SHOW INFO PANEL</a>
Jobs on Clusters	<a href="#">Filter</a> Search clusters, press Enter
<a href="#">Clusters</a>	
<a href="#">Jobs</a>	
<a href="#">Workflows</a>	
<a href="#">Autoscaling policies</a>	
Serverless	
<a href="#">Batches</a>	
Metastore Services	
<a href="#">Metastore</a>	
<a href="#">Federation</a>	
Utilities	
<a href="#">Component exchange</a>	
<a href="#">Workbench</a>	
<a href="#">Release Notes</a>	

<input type="checkbox"/>	Name <span>↑</span>	Status	Region	Zone	Total worker nodes	Scheduled deletion	Cloud Storage staging bucket	Created
<input type="checkbox"/>	<a href="#">cluster-b74f</a>	Running	us-central1	us-central1-a	2	Off	<a href="#">dataproc-staging-us-central1-179637453175-deik4x2r</a>	Oct 31, 2022, 7:55:15 AM



## Ssh the virtual session



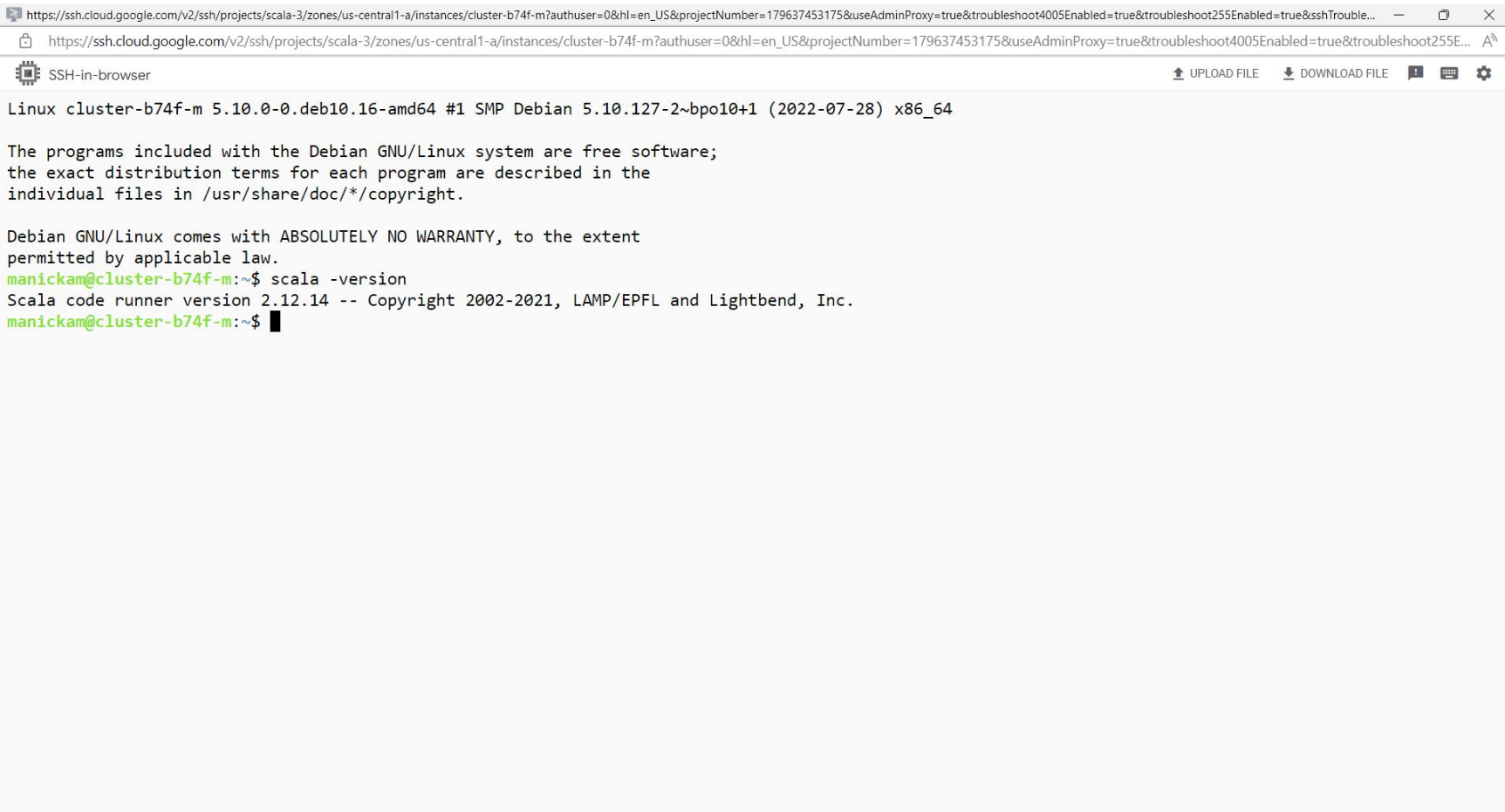
```
Linux cluster-b74f-m 5.10.0-0.deb10.16-amd64 #1 SMP Debian 5.10.127-2~bpo10+1 (2022-07-28) x86_64
```

```
The programs included with the Debian GNU/Linux system are free software;  
the exact distribution terms for each program are described in the  
individual files in /usr/share/doc/*/copyright.
```

```
Debian GNU/Linux comes with ABSOLUTELY NO WARRANTY, to the extent  
permitted by applicable law.
```


```
manickam@cluster-b74f-m:~$ █
```

# Verify Scala Version



```
https://ssh.cloud.google.com/v2/ssh/projects/scala-3/zones/us-central1-a/instances/cluster-b74f-m?authuser=0&hl=en_US&projectNumber=179637453175&useAdminProxy=true&troubleshoot4005Enabled=true&troubleshoot255Enabled=true&sshTrouble...  
https://ssh.cloud.google.com/v2/ssh/projects/scala-3/zones/us-central1-a/instances/cluster-b74f-m?authuser=0&hl=en_US&projectNumber=179637453175&useAdminProxy=true&troubleshoot4005Enabled=true&troubleshoot255E...  
SSH-in-browser  
Linux cluster-b74f-m 5.10.0-0.deb10.16-amd64 #1 SMP Debian 5.10.127-2~bpo10+1 (2022-07-28) x86_64  
  
The programs included with the Debian GNU/Linux system are free software;  
the exact distribution terms for each program are described in the  
individual files in /usr/share/doc/*/copyright.  
  
Debian GNU/Linux comes with ABSOLUTELY NO WARRANTY, to the extent  
permitted by applicable law.  
manickam@cluster-b74f-m:~$ scala -version  
Scala code runner version 2.12.14 -- Copyright 2002-2021, LAMP/EPFL and Lightbend, Inc.  
manickam@cluster-b74f-m:~$
```

Input file for PageRank : pagerank.txt



The screenshot shows a web browser window with two tabs. The active tab displays the URL `https://ssh.cloud.google.com/v2/ssh/projects/scala-3/zones/us-central1-a/instances/cluster-b74f-m?authuser=0&hl=en_US&projectNumber=179637453175&useAdminProxy=true&troubleshoot4005Enabled=true&troubleshoot255Enabled=true&sshTrouble...`. Below the address bar, there is a toolbar with an SSH icon, the text "SSH-in-browser", and buttons for "UPLOAD FILE", "DOWNLOAD FILE", and a settings icon. The main content area is a terminal window with a green prompt `manickam@cluster-b74f-m:~/PageRank$`. The command `cat pagerank.txt` has been executed, resulting in the following output:

```
A B
A C
B C
C A
manickam@cluster-b74f-m:~/PageRank$
```



## Iteration - 1

```
https://ssh.cloud.google.com/v2/ssh/projects/scala-3/zones/us-central1-a/instances/cluster-b74f-m?authuser=0&hl=en_US&projectNumber=179637453175&useAdminProxy=true&troubleshoot4005Enabled=true&troubleshoot255Enabled=true&sshTrouble...
https://ssh.cloud.google.com/v2/ssh/projects/scala-3/zones/us-central1-a/instances/cluster-b74f-m?authuser=0&hl=en_US&projectNumber=179637453175&useAdminProxy=true&troubleshoot4005Enabled=true&troubleshoot255E...
SSH-in-browser
UPLOAD FILE
DOWNLOAD FILE
/

Using Scala version 2.12.14 (OpenJDK 64-Bit Server VM, Java 1.8.0_345)
Type in expressions to have them evaluated.
Type :help for more information.

scala> import org.apache.spark.sql.SparkSession
import org.apache.spark.sql.SparkSession

scala> import org.apache.spark.HashPartitioner
import org.apache.spark.HashPartitioner

scala> val links = sc.parallelize(List(("A",List("B","C")),("B", List("C")),("C",List("A")))).partitionBy(new HashPartitioner(3)).persist()
links: org.apache.spark.rdd.RDD[(String, List[String])] = ShuffledRDD[1] at partitionBy at <console>:25

scala> var ranks = links.mapValues(v => 1.0) // Initialized
ranks: org.apache.spark.rdd.RDD[(String, Double)] = MapPartitionsRDD[2] at mapValues at <console>:25

scala>
scala> for (i <- 0 to 0) {
  | val contributions = links.join(ranks).flatMap { case (url, (links, rank)) => links.map(dest => (dest, rank / links.size)) }
  | ranks = contributions.reduceByKey((x, y) => x + y).mapValues(v => 0.15 + 0.85*v)
  | ranks.collect
  | }

scala> ranks.collect
res1: Array[(String, Double)] = Array((B,0.575), (C,1.4249999999999998), (A,1.0))

scala>
scala> :quit
manickam@cluster-b74f-m:~/PageRank$
```

## Iteration 2

```
https://ssh.cloud.google.com/v2/ssh/projects/scala-3/zones/us-central1-a/instances/cluster-b74f-m?authuser=0&hl=en_US&projectNumber=179637453175&useAdminProxy=true&troubleshoot4005Enabled=true&troubleshoot255Enabled=true&sshTrouble...  
https://ssh.cloud.google.com/v2/ssh/projects/scala-3/zones/us-central1-a/instances/cluster-b74f-m?authuser=0&hl=en_US&projectNumber=179637453175&useAdminProxy=true&troubleshoot4005Enabled=true&troubleshoot255E...  
SSH-in-browser  
UPLOAD FILE  
DOWNLOAD FILE  
/_/  
  
Using Scala version 2.12.14 (OpenJDK 64-Bit Server VM, Java 1.8.0_345)  
Type in expressions to have them evaluated.  
Type :help for more information.  
  
scala> import org.apache.spark.sql.SparkSession  
import org.apache.spark.sql.SparkSession  
  
scala> import org.apache.spark.HashPartitioner  
import org.apache.spark.HashPartitioner  
  
scala> val links = sc.parallelize(List(("A",List("B","C")),("B", List("C")),("C",List("A")))).partitionBy(new HashPartitioner(3)).persist()  
links: org.apache.spark.rdd.RDD[(String, List[String])] = ShuffledRDD[1] at partitionBy at <console>:25  
  
scala> var ranks = links.mapValues(v => 1.0) // Initialized  
ranks: org.apache.spark.rdd.RDD[(String, Double)] = MapPartitionsRDD[2] at mapValues at <console>:25  
  
scala>  
  
scala> for (i <- 0 to 1) {  
  | val contributions = links.join(ranks).flatMap { case (url, (links, rank)) => links.map(dest => (dest, rank / links.size)) }  
  | ranks = contributions.reduceByKey((x, y) => x + y).mapValues(v => 0.15 + 0.85*v)  
  | ranks.collect  
  | }  
  
scala> ranks.collect  
res1: Array[(String, Double)] = Array((B,0.575), (C,1.06375), (A,1.3612499999999996))  
  
scala>  
  
scala> :quit  
manickam@cluster-b74f-m:~/PageRank$
```

## Iteration 3

```
https://ssh.cloud.google.com/v2/ssh/projects/scala-3/zones/us-central1-a/instances/cluster-b74f-m?authuser=0&hl=en_US&projectNumber=179637453175&useAdminProxy=true&troubleshoot4005Enabled=true&troubleshoot255Enabled=true&sshTrouble...  
https://ssh.cloud.google.com/v2/ssh/projects/scala-3/zones/us-central1-a/instances/cluster-b74f-m?authuser=0&hl=en_US&projectNumber=179637453175&useAdminProxy=true&troubleshoot4005Enabled=true&troubleshoot255E...  
SSH-in-browser  
/_/  
Using Scala version 2.12.14 (OpenJDK 64-Bit Server VM, Java 1.8.0_345)  
Type in expressions to have them evaluated.  
Type :help for more information.  
  
scala> import org.apache.spark.sql.SparkSession  
import org.apache.spark.sql.SparkSession  
  
scala> import org.apache.spark.HashPartitioner  
import org.apache.spark.HashPartitioner  
  
scala> val links = sc.parallelize(List(("A",List("B","C")),("B", List("C")),("C",List("A")))).partitionBy(new HashPartitioner(3)).persist()  
links: org.apache.spark.rdd.RDD[(String, List[String])] = ShuffledRDD[1] at partitionBy at <console>:25  
  
scala> var ranks = links.mapValues(v => 1.0) // Initialized  
ranks: org.apache.spark.rdd.RDD[(String, Double)] = MapPartitionsRDD[2] at mapValues at <console>:25  
  
scala>  
  
scala> for (i <- 0 to 2) {  
  | val contributions = links.join(ranks).flatMap { case (url, (links, rank)) => links.map(dest => (dest, rank / links.size)) }  
  | ranks = contributions.reduceByKey((x, y) => x + y).mapValues(v => 0.15 + 0.85*v)  
  | ranks.collect  
  | }  
  
scala> ranks.collect  
res1: Array[(String, Double)] = Array((B,0.7285312499999999), (C,1.2172812499999999), (A,1.0541874999999998))  
  
scala>  
  
scala> :quit  
manickam@cluster-b74f-m:~/PageRank$
```

## Iterations using PySpark program

### Iteration 0

```
hduser@cs570bigdata: ~  
22/11/02 00:04:11 INFO TaskSetManager: Finished task 0.0 in stage 2.0 (TID 2) in 581 ms on cs570bigdata (executor driver) (1/1)  
22/11/02 00:04:11 INFO TaskSchedulerImpl: Removed TaskSet 2.0, whose tasks have all completed, from pool  
22/11/02 00:04:11 INFO DAGScheduler: ResultStage 2 (collect at /home/hduser/PythonPageRank.py:172) finished in 0.705 s  
22/11/02 00:04:11 INFO DAGScheduler: Job 0 is finished. Cancelling potential speculative or zombie tasks for this job  
22/11/02 00:04:11 INFO TaskSchedulerImpl: Killing all running tasks in stage 2: Stage finished  
22/11/02 00:04:11 INFO DAGScheduler: Job 0 finished: collect at /home/hduser/PythonPageRank.py:172, took 4.118938 s  
22/11/02 00:04:11 INFO SparkUI: Stopped Spark web UI at http://cs570bigdata:4040  
A has rank: 1.0.  
B has rank: 1.0.  
C has rank: 1.0.  
22/11/02 00:04:12 INFO MapOutputTrackerMasterEndpoint: MapOutputTrackerMasterEndpoint stopped!
```

## Iteration 1

```
hduser@cs570bigdata: ~  
22/11/02 00:23:22 INFO TaskSetManager: Finished task 0.0 in stage 4.0 (TID 6) in 773 ms on cs570bigdata (executor driver) (1/2)  
22/11/02 00:23:22 INFO Executor: Finished task 1.0 in stage 4.0 (TID 7). 1673 bytes result sent to driver  
22/11/02 00:23:22 INFO TaskSetManager: Finished task 1.0 in stage 4.0 (TID 7) in 822 ms on cs570bigdata (executor driver) (2/2)  
22/11/02 00:23:22 INFO TaskSchedulerImpl: Removed TaskSet 4.0, whose tasks have all completed, from pool  
22/11/02 00:23:22 INFO DAGScheduler: ResultStage 4 (collect at /home/hduser/PythonPageRank.py:172) finished in 0.853 s  
22/11/02 00:23:22 INFO DAGScheduler: Job 0 is finished. Cancelling potential speculative or zombie tasks for this job  
22/11/02 00:23:22 INFO TaskSchedulerImpl: Killing all running tasks in stage 4: Stage finished  
22/11/02 00:23:22 INFO DAGScheduler: Job 0 finished: collect at /home/hduser/PythonPageRank.py:172, took 8.724759 s  
C has rank: 1.4249999999999998.  
A has rank: 1.0.  
B has rank: 0.575.
```

## Iteration 2

```
hduser@cs570bigdata: ~  
22/11/02 00:28:33 INFO TaskSetManager: Finished task 0.0 in stage 6.0 (TID 12) in 1085 ms on cs570bigdata (executor driver) (3/3)  
22/11/02 00:28:33 INFO TaskSchedulerImpl: Removed TaskSet 6.0, whose tasks have all completed, from pool  
22/11/02 00:28:33 INFO DAGScheduler: ResultStage 6 (collect at /home/hduser/PythonPageRank.py:172) finished in 1.229 s  
22/11/02 00:28:33 INFO DAGScheduler: Job 0 is finished. Cancelling potential speculative or zombie tasks for this job  
22/11/02 00:28:33 INFO TaskSchedulerImpl: Killing all running tasks in stage 6: Stage finished  
22/11/02 00:28:33 INFO DAGScheduler: Job 0 finished: collect at /home/hduser/PythonPageRank.py:172, took 10.690277 s  
C has rank: 1.06375.  
B has rank: 0.575.  
A has rank: 1.3612499999999996.  
22/11/02 00:28:34 INFO SparkUI: Stopped Spark web UI at http://cs570bigdata:4040  
22/11/02 00:28:36 INFO MapOutputTrackerMasterEndpoint: MapOutputTrackerMasterEndpoint stopped!
```

## **Conclusion**

PageRank is a system developed in 1997 by Google founders Larry Page and Sergey Brin. It was designed to evaluate the quality and quantity of links to a page. Along with other factors, the score determined pages' positions in search engine rankings.

It helps Google to decide the importance of a page and it is the main reason behind which the PageRank for a website is determined in the search results.