

# Assignment 1- Data Analysis

Srijith Unni

DCU Student ID: 20211114

MSc in Computing – Data Analytics

Cloud Technologies CA675

## Complete chronological screenshots:

The screenshot shows the StackExchange Data Explorer interface. At the top, there's a navigation bar with 'Home', 'Queries', and 'Users'. A 'Compose Query' button is on the right. Below the navigation bar, the 'Viewing Query' section shows a query editor with a title 'Enter a title for your query' and a description 'Enter a description'. The query is: `1 to top 50000 from Posts where Posts.ViewCount > 115000 order by Posts.ViewCount desc`. The database schema is visible on the right, showing tables like 'Posts', 'Answers', 'Comments', etc. Below the query editor, there are buttons for 'Run Query', 'Cancel', and 'Options'. The 'Results' tab is selected, showing a table of query results. The table has columns: Id, PostTypeId, AcceptedAnswerId, ParentId, CreationDate, DeletionDate, Score, ViewCount, Body, OwnerUserId, OwnerDisplayName, LastEditorUserId, LastEditorDisplayName, LastEditDate, and LastActivityDate. The results are sorted by ViewCount in descending order.

Id	PostTypeId	AcceptedAnswerId	ParentId	CreationDate	DeletionDate	Score	ViewCount	Body	OwnerUserId	OwnerDisplayName	LastEditorUserId	LastEditorDisplayName	LastEditDate	LastActivityDate
927358	1	927358		2009-05-29 18:09:14		21842	9043120	<p>I accidentally <strong>committed the vno...	89904		8674599		2020-11-03 04:39:59	2020-11-03 05:20:08
2003505	1	2003515		2010-01-05 01:12:15		17427	8490909	<p>I want to delete a branch both locally and ...	95592		13114227		2020-08-06 12:08:14	2020-10-01 05:08:44
5767325	1	5767357		2011-04-23 22:17:18		8789	7328939	<p>I have an array of numbers and I'm using ...	364969		594916		2020-05-21 04:20:03	2020-11-08 13:50:18
169568	1	16957078		2013-06-06 06:06:45		5567	7079498	<p>I'm trying to find a way to scan my entire ...	954806		63550		2017-05-21 11:40:16	2020-08-09 07:33:09
2906582	1	2906586		2010-05-25 16:39:47		2031	8891424	<p>I would like to create an HTML button that ...	48523		792066		2020-07-24 26:57:44	2020-07-24 26:57:44
503093	1	506004		2009-02-02 12:54:16		7718	8488955	<p>How can I redirect the user from one pag...	44984	venkatasubram	63550	shager	2019-02-06 12:40:14	2020-09-27 04:49:44

44177 rows returned in 897 ms (cached)

Run Query Cancel Options: ☐ Text-only results ☐ Include execution plan

Switch sites:  search by name or url

Hold tight while we fetch your results

Results Messages

Download CSV

Id	PostTypeId	AcceptedAnswerId	ParentId	CreationDate	DeletionDate	Score	ViewCount	Body	OwnerUserId	OwnerDisplayName
21376645	1	21376885		2014-01-27 09:01:51		7	69999	<p>As i know, i can create an array with item ...	3122881	
20829348	1	20829349		2013-12-29 19:43:17		162	69999	<p>How do i join two strings in a list with a sp...	1139353	
5323349	1	5323350		2011-03-16 05:59:12		42	69999	<p>What is the <code>#error</code> directiv...	390452	Shine
29574732	1			2015-04-11 06:09:43		44	69997	<p>As of 2015, I see that Android studio is an ...	1800583	
22842691	1	22844164		2014-04-03 15:46:59		159	69996	<p>Since I first saw a <code>dist</code> dir...	2391795	
14870596	1	14871317		2013-02-14 08:27:44		36	69996	<p>I'm trying to run a sample android code in...	893411	
26307920	1	26308081		2014-10-10 20:56:46		22	69995	<p>I'm learning node Js and i'm trying to do a...	3595813	
26033239	1	26033300		2014-09-25 07:46:20		31	69995	<p>i have a problem converting <code>Obje...	4046638	
5193973	1	5193991		2011-03-04 12:55:02		17	69995	<p>How would i loop through all the td eleme...	517406	
8707753	1			2012-01-03 04:00:10		9	69995	<p>i was thinking of developing an iPhone an...	1080731	
6993132	1	6993177		2011-08-09 08:05:02		56	69994	<p>If i have an integer variable i can use <co...	184759	
6825834	1	6826076		2011-07-26 05:52:19		32	69993	<p>i have a string as shown below.</p> <p>pre...	683898	
3192095	1			2010-07-07 05:17:09		50	69993	<p>i was just curious where exactly the singl...	311498	
11260659	1	11262192		2012-06-29 11:05:59		14	69992	<p>i have a question about if statement in tcl ...	user707549	
12336234	1	12340302		2012-09-09 03:11:19		28	69990	<p>i have data stored in a CSV where the firs...	907714	
12999899	1	13000013		2012-10-21 16:44:03		7	69990	<p>i am trying to have a scanner take input i...	605328	

46925 rows returned in 84053 ms

Google Cloud PlatformMy Project 46318Search products and resources

Dataproc

ClustersJobsWorkflowsAutoscaling policiesComponent exchangeNotebooks

Create a cluster

- Set up cluster
  - Create Cluster
- Configure nodes (optional)
  - Nodes
- Customise cluster (optional)
  - Customise
- Manage security (optional)
  - Manage network security

CREATECANCEL

Equivalent REST or command line

Name

- Cluster Name
  - cluster-677c

Location

Region

- europa-west1

Zone

- europa-west1-b

Cluster type

☒ Standard (1 master, N workers)

☐ Single Node (1 master, 0 workers)

Provides one node that acts as both master and worker. Good for proof-of-concept or small-scale processing.

☐ High availability (3 masters, N workers)

Hadoop High Availability mode provides uninterrupted YARN and HDFS operations despite any single-node failures/reboots.

Autoscaling

Automates cluster resource management based on an autoscaling policy.

Creating cluster cluster-677c

```
srijith_unni2@cluster-677c-m: ~ - Google Chrome
ssh.cloud.google.com/projects/focused-mote-292715/zones/europe-west1-b/instances/cluster-677c-m?authuser=1&hl=en_GB&projectNumber=611159050419&useAdminProxy=true

System load: 0.37      Processes:      126
Usage of /: 46.3% of 14.37GB   Users logged in: 0
Memory usage: 65%      IP address for ens4: 10.132.0.11
Swap usage: 0%

* Introducing self-healing high availability clustering for MicroK8s!
  Super simple, hardened and opinionated Kubernetes for production.

  https://microk8s.io/high-availability

16 packages can be updated.
16 updates are security updates.

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

Ubuntu comes with ABSOLUTELY NO WARRANTY, to the extent permitted by
applicable law.

srijith_unni2@cluster-677c-m:~$ hadoop fs -ls /
Found 3 items
drwx-----   - mapred  hadoop          0 2020-11-20 07:49 /hadoop
drwxrwxrwt   - hddfs  hadoop          0 2020-11-20 07:50 /tmp
drwxrwxrwt   - hddfs  hadoop          0 2020-11-20 07:49 /user
srijith_unni2@cluster-677c-m:~$ hadoop fs -mkdir /stackex
srijith_unni2@cluster-677c-m:~$ hadoop fs -ls /
Found 4 items
drwx-----   - mapred  hadoop          0 2020-11-20 07:49 /hadoop
drwxr-xr-x   - srijith_unni2 hadoop      0 2020-11-20 07:54 /stackex
drwxrwxrwt   - hddfs  hadoop          0 2020-11-20 07:50 /tmp
drwxrwxrwt   - hddfs  hadoop          0 2020-11-20 07:49 /user
srijith_unni2@cluster-677c-m:~$ hadoop fs -mkdir /stackex/results
srijith_unni2@cluster-677c-m:~$ hadoop fs -ls /stackex
Found 1 items
drwxr-xr-x   - srijith_unni2 hadoop      0 2020-11-20 07:55 /stackex/results
srijith_unni2@cluster-677c-m:~$
```

```
srijith_unni2@cluster-677c-m: ~ - Google Chrome
ssh.cloud.google.com/projects/focused-mote-292715/zones/europe-west1-b/instances/cluster-677c-m?authuser=1&hl=en_GB&projectNumber=611159050419&useAdminProxy=true

System information as of Fri Nov 20 13:10:43 UTC 2020

System load: 0.0      Processes:      130
Usage of /: 46.5% of 14.37GB   Users logged in: 1
Memory usage: 66%      IP address for ens4: 10.132.0.16
Swap usage: 0%

* Introducing self-healing high availability clustering for MicroK8s!
  Super simple, hardened and opinionated Kubernetes for production.

  https://microk8s.io/high-availability

16 packages can be updated.
16 updates are security updates.

New release '20.04.1 LTS' available.
Run 'do-release-upgrade' to upgrade to it.

Last login: Fri Nov 20 12:49:07 2020 from 35.235.240.98
srijith_unni2@cluster-677c-m:~$ ls
ViewCount1.csv ViewCount2.csv ViewCount3.csv csupload ViewCount3_1.csv csupload
srijith_unni2@cluster-677c-m:~$ rm ViewCount3
rm: cannot remove 'ViewCount3': No such file or directory
srijith_unni2@cluster-677c-m:~$ rm ViewCount3.csv csupload
srijith_unni2@cluster-677c-m:~$ ls
ViewCount1.csv ViewCount2.csv ViewCount3_1.csv csupload
srijith_unni2@cluster-677c-m:~$ hadoop fs -ls /stackex
Found 3 items
-rw-r--r--  2 srijith_unni2 hadoop  9652198 2020-11-20 12:59 /stackex/ViewCount1.csv
-rw-r--r--  2 srijith_unni2 hadoop  8747424 2020-11-20 13:03 /stackex/ViewCount2.csv
drwxr-xr-x  - srijith_unni2 hadoop          0 2020-11-20 12:50 /stackex/results
srijith_unni2@cluster-677c-m:~$ hadoop fs -put /home/srijith_unni2/ViewCount3_1.csv /stackex
put: '/home/srijith_unni2/ViewCount3_1.csv': No such file or directory
srijith_unni2@cluster-677c-m:~$ hadoop fs -put /home/srijith_unni2/ViewCount3_1.csv /stackex
srijith_unni2@cluster-677c-m:~$ hadoop fs -ls /stackex
Found 4 items
-rw-r--r--  2 srijith_unni2 hadoop  9652198 2020-11-20 12:59 /stackex/ViewCount1.csv
-rw-r--r--  2 srijith_unni2 hadoop  8747424 2020-11-20 13:03 /stackex/ViewCount2.csv
-rw-r--r--  2 srijith_unni2 hadoop  9194085 2020-11-20 13:26 /stackex/ViewCount3_1.csv
drwxr-xr-x  - srijith_unni2 hadoop          0 2020-11-20 12:50 /stackex/results
srijith_unni2@cluster-677c-m:~$ hadoop fs -put /home/srijith_unni2/ViewCount4.csv /stackex
srijith_unni2@cluster-677c-m:~$ hadoop fs -ls /stackex
Found 5 items
-rw-r--r--  2 srijith_unni2 hadoop  9652198 2020-11-20 12:59 /stackex/ViewCount1.csv
-rw-r--r--  2 srijith_unni2 hadoop  8747424 2020-11-20 13:03 /stackex/ViewCount2.csv
-rw-r--r--  2 srijith_unni2 hadoop  9194085 2020-11-20 13:26 /stackex/ViewCount3_1.csv
-rw-r--r--  2 srijith_unni2 hadoop  8622108 2020-11-20 13:32 /stackex/ViewCount4.csv
drwxr-xr-x  - srijith_unni2 hadoop          0 2020-11-20 12:50 /stackex/results
srijith_unni2@cluster-677c-m:~$
```

File Transfer	Cancel
ViewCount5.csv	14%

File upload destination: /home/srijith\_unni2

## Browse Directory

Show 25 entries

Search:

<input type="checkbox"/>	Permission	Owner	Group	Size	Last Modified	Replication	Block Size	Name	<input type="checkbox"/>
<input type="checkbox"/>	-rw-r--r--	srjith_unni2	hadoop	9.21 MB	Nov 20 18:29	2	128 MB	ViewCount1.csv	<input type="checkbox"/>
<input type="checkbox"/>	-rw-r--r--	srjith_unni2	hadoop	9.39 MB	Nov 20 22:53	2	128 MB	ViewCount1.txt	<input type="checkbox"/>
<input type="checkbox"/>	-rw-r--r--	srjith_unni2	hadoop	9.39 MB	Nov 20 20:15	2	128 MB	ViewCount1_1.csv	<input type="checkbox"/>
<input type="checkbox"/>	-rw-r--r--	srjith_unni2	hadoop	9.39 MB	Nov 20 20:38	2	128 MB	ViewCount1_1.txt	<input type="checkbox"/>
<input type="checkbox"/>	-rw-r--r--	srjith_unni2	hadoop	8.34 MB	Nov 20 18:33	2	128 MB	ViewCount2.csv	<input type="checkbox"/>
<input type="checkbox"/>	-rw-r--r--	srjith_unni2	hadoop	8.34 MB	Nov 20 23:02	2	128 MB	ViewCount2.txt	<input type="checkbox"/>
<input type="checkbox"/>	-rw-r--r--	srjith_unni2	hadoop	8.76 MB	Nov 20 23:20	2	128 MB	ViewCount3.txt	<input type="checkbox"/>
<input type="checkbox"/>	-rw-r--r--	srjith_unni2	hadoop	8.77 MB	Nov 20 18:56	2	128 MB	ViewCount3_1.csv	<input type="checkbox"/>
<input type="checkbox"/>	-rw-r--r--	srjith_unni2	hadoop	8.22 MB	Nov 20 19:02	2	128 MB	ViewCount4.csv	<input type="checkbox"/>
<input type="checkbox"/>	-rw-r--r--	srjith_unni2	hadoop	8.22 MB	Nov 20 23:14	2	128 MB	ViewCount4.txt	<input type="checkbox"/>
<input type="checkbox"/>	-rw-r--r--	srjith_unni2	hadoop	3.03 MB	Nov 20 19:06	2	128 MB	ViewCount5.csv	<input type="checkbox"/>
<input type="checkbox"/>	-rw-r--r--	srjith_unni2	hadoop	3.03 MB	Nov 20 23:21	2	128 MB	ViewCount5.txt	<input type="checkbox"/>

```
sshcloud.google.com/projects/focused-mote-292715/zones/europe-west1-b/instances/cluster-677c-m?authuser=1&hl=en_GB&projectNumber=611159050419&useAdminProxy=true
$ pig
20/11/20 17:53:10 INFO pig.ExecTypeProvider: Trying ExecType : LOCAL
20/11/20 17:53:10 INFO pig.ExecTypeProvider: Trying ExecType : MAPREDUCE
20/11/20 17:53:10 INFO pig.ExecTypeProvider: Picked MAPREDUCE as the ExecType
2020-11-20 17:53:10,503 [main] INFO org.apache.pig.Main - Apache Pig version 0.17.0 (r: unknown) compiled Oct 18 2020, 10:13:54
2020-11-20 17:53:10,503 [main] INFO org.apache.pig.Main - Logging error messages to: /home/srjith_unni2/pig_160584790495.log
2020-11-20 17:53:12,404 [main] INFO org.apache.pig.PigServer - Pig Script ID for the session: PIG-default-d926933d-8a37-4d46-9d31-d3deb51e52e
2020-11-20 17:53:12,457 [main] INFO org.apache.hadoop.yarn.client.api.impl.TimelineClientImpl - Timeline service address: null
2020-11-20 17:53:13,030 [main] INFO org.apache.pig.backend.hadoop.PigATSCClient - Created ATS Hook
2020-11-20 17:53:13,083 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - yarn.resourcemanager.system-metrics-publisher.enabled is deprecated. Instead, use yarn.system-metrics-publis
her.enabled
grunt> A = LOAD '/stackex/ViewCount1.txt' USING PigStorage('(t') AS (Id:chararray, PostTypeId:chararray, AcceptedAnswerId:chararray, ParentId:chararray, CreationDate:chararray, DeletionDate:chararr
ay, Score:chararray, ViewCount:chararray, OwnerUserId:chararray, OwnerDisplayName:chararray, LastEditorUserId:chararray, LastEditorDisplayName:chararray, LastEditDate:chararray, LastActivityTime:
chararray, Title:chararray, Tags:chararray, AnswerCount:chararray, CommentCount:chararray, FavouriteCount:chararray, ClosedDate:chararray, CommunityOwnedDate:chararray, ContentLicense:chararray);
2020-11-20 17:53:34,437 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - yarn.resourcemanager.system-metrics-publisher.enabled is deprecated. Instead, use yarn.system-metrics-publis
her.enabled
grunt> B = LOAD '/stackex/ViewCount2.txt' USING PigStorage('(t') AS (Id:chararray, PostTypeId:chararray, AcceptedAnswerId:chararray, ParentId:chararray, CreationDate:chararray, DeletionDate:chararr
ay, Score:chararray, ViewCount:chararray, OwnerUserId:chararray, OwnerDisplayName:chararray, LastEditorUserId:chararray, LastEditorDisplayName:chararray, LastEditDate:chararray, LastActivityTime:
chararray, Title:chararray, Tags:chararray, AnswerCount:chararray, CommentCount:chararray, FavouriteCount:chararray, ClosedDate:chararray, CommunityOwnedDate:chararray, ContentLicense:chararray);
2020-11-20 17:53:41,531 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - yarn.resourcemanager.system-metrics-publisher.enabled is deprecated. Instead, use yarn.system-metrics-publis
her.enabled
grunt> C = LOAD '/stackex/ViewCount3.txt' USING PigStorage('(t') AS (Id:chararray, PostTypeId:chararray, AcceptedAnswerId:chararray, ParentId:chararray, CreationDate:chararray, DeletionDate:chararr
ay, Score:chararray, ViewCount:chararray, OwnerUserId:chararray, OwnerDisplayName:chararray, LastEditorUserId:chararray, LastEditorDisplayName:chararray, LastEditDate:chararray, LastActivityTime:
chararray, Title:chararray, Tags:chararray, AnswerCount:chararray, CommentCount:chararray, FavouriteCount:chararray, ClosedDate:chararray, CommunityOwnedDate:chararray, ContentLicense:chararray);
2020-11-20 17:53:49,143 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - yarn.resourcemanager.system-metrics-publisher.enabled is deprecated. Instead, use yarn.system-metrics-publis
her.enabled
grunt> D = LOAD '/stackex/ViewCount4.txt' USING PigStorage('(t') AS (Id:chararray, PostTypeId:chararray, AcceptedAnswerId:chararray, ParentId:chararray, CreationDate:chararray, DeletionDate:chararr
ay, Score:chararray, ViewCount:chararray, OwnerUserId:chararray, OwnerDisplayName:chararray, LastEditorUserId:chararray, LastEditorDisplayName:chararray, LastEditDate:chararray, LastActivityTime:
chararray, Title:chararray, Tags:chararray, AnswerCount:chararray, CommentCount:chararray, FavouriteCount:chararray, ClosedDate:chararray, CommunityOwnedDate:chararray, ContentLicense:chararray);
2020-11-20 17:53:56,091 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - yarn.resourcemanager.system-metrics-publisher.enabled is deprecated. Instead, use yarn.system-metrics-publis
her.enabled
grunt> E = LOAD '/stackex/ViewCount5.txt' USING PigStorage('(t') AS (Id:chararray, PostTypeId:chararray, AcceptedAnswerId:chararray, ParentId:chararray, CreationDate:chararray, DeletionDate:chararr
ay, Score:chararray, ViewCount:chararray, OwnerUserId:chararray, OwnerDisplayName:chararray, LastEditorUserId:chararray, LastEditorDisplayName:chararray, LastEditDate:chararray, LastActivityTime:
chararray, Title:chararray, Tags:chararray, AnswerCount:chararray, CommentCount:chararray, FavouriteCount:chararray, ClosedDate:chararray, CommunityOwnedDate:chararray, ContentLicense:chararray);
2020-11-20 17:54:03,063 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - yarn.resourcemanager.system-metrics-publisher.enabled is deprecated. Instead, use yarn.system-metrics-publis
her.enabled
grunt>
```

```
sshcloud.google.com/projects/focused-mote-292715/zones/europe-west1-b/instances/cluster-677c-m?authuser=1&hl=en_GB&projectNumber=611159050419&useAdminProxy=true
grunt> I = FOREACH A GENERATE $0,$2,$4,$6,$7,$8,$10,$12,$13,$14,$15,$16,$17,$18;
grunt> J = FOREACH B GENERATE $0,$2,$4,$6,$7,$8,$10,$12,$13,$14,$15,$16,$17,$18;
grunt> K = FOREACH C GENERATE $0,$2,$4,$6,$7,$8,$10,$12,$13,$14,$15,$16,$17,$18;
grunt> L = FOREACH D GENERATE $0,$2,$4,$6,$7,$8,$10,$12,$13,$14,$15,$16,$17,$18;
grunt> M = FOREACH E GENERATE $0,$2,$4,$6,$7,$8,$10,$12,$13,$14,$15,$16,$17,$18;
grunt> N = UNION I,J,K,L,M;
grunt>
```

```
ssh.cloud.google.com/projects/foocused-mote-292715/zones/europe-west1-b/instances/cluster-677c-m?authuser=1&hl=en_GB&projectNumber=611159050419&useAdminProxy=true

grunt> I = FOREACH A GENERATE $0,$2,$4,$6,$7,$8,$10,$12,$13,$14,$15,$16,$17,$18;
grunt> J = FOREACH B GENERATE $0,$2,$4,$6,$7,$8,$10,$12,$13,$14,$15,$16,$17,$18;
grunt> K = FOREACH C GENERATE $0,$2,$4,$6,$7,$8,$10,$12,$13,$14,$15,$16,$17,$18;
grunt> L = FOREACH D GENERATE $0,$2,$4,$6,$7,$8,$10,$12,$13,$14,$15,$16,$17,$18;
grunt> M = FOREACH E GENERATE $0,$2,$4,$6,$7,$8,$10,$12,$13,$14,$15,$16,$17,$18;
grunt> R = UNION I,J,K,L,M;
grunt> STORE R INTO '/stackex/output_final.txt' USING PigStorage('\t');
2020-11-20 17:55:34,425 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - yarn.resourcemanager.system-metrics-publisher.enabled is deprecated. Instead, use yarn.system-metrics-publisher.enabled
2020-11-20 17:55:34,542 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - mapred.textoutputformat.separator is deprecated. Instead, use mapreduce.output.textoutputformat.separator
2020-11-20 17:55:34,652 [main] INFO org.apache.pig.tools.pigstats.ScriptState - Pig features used in the script: UNION
2020-11-20 17:55:34,696 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - yarn.resourcemanager.system-metrics-publisher.enabled is deprecated. Instead, use yarn.system-metrics-publisher.enabled
2020-11-20 17:55:34,707 [main] INFO org.apache.pig.data.SchemaTupleBackend - Key [pig.schematuple] was not set... will not generate code.
2020-11-20 17:55:34,762 [main] INFO org.apache.pig.newplan.logical.optimizer.LogicalPlanOptimizer - (RULES PRUNED=[AddForEach, ColumnMapKeyPrune, ConstantCalculator, GroupByConstParallelSetter, JoinOptimizer, LoadTypeCastInserter, MergeFilter, MergeForEach, NestedLimitOptimizer, PartitionFilterOptimizer, PredicatePushdownOptimizer, PushDownForEachFlatten, PushUpFilter, SplitFilter, StreamTypeCastInserter])
2020-11-20 17:55:34,864 [main] INFO org.apache.pig.newplan.logical.rules.ColumnPruneVisitor - Columns pruned for A: $1, $3, $5, $9, $11, $19, $20, $21
2020-11-20 17:55:34,868 [main] INFO org.apache.pig.newplan.logical.rules.ColumnPruneVisitor - Columns pruned for B: $1, $3, $5, $9, $11, $19, $20, $21
2020-11-20 17:55:34,868 [main] INFO org.apache.pig.newplan.logical.rules.ColumnPruneVisitor - Columns pruned for C: $1, $3, $5, $9, $11, $19, $20, $21
2020-11-20 17:55:34,869 [main] INFO org.apache.pig.newplan.logical.rules.ColumnPruneVisitor - Columns pruned for D: $1, $3, $5, $9, $11, $19, $20, $21
2020-11-20 17:55:34,869 [main] INFO org.apache.pig.newplan.logical.rules.ColumnPruneVisitor - Columns pruned for E: $1, $3, $5, $9, $11, $19, $20, $21
2020-11-20 17:55:34,955 [main] INFO org.apache.pig.impl.util.SpillableMemoryManager - Selected heap (Tenured Gen) of size 699072912 to monitor. collectionUsageThreshold = 489350752, usageThreshold = 489350752
2020-11-20 17:55:35,107 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MRCCompiler - File concatenation threshold: 100 optimistic? false
2020-11-20 17:55:35,172 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MultiQueryOptimizer - MR plan size before optimization: 1
2020-11-20 17:55:35,173 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MultiQueryOptimizer - MR plan size after optimization: 1
2020-11-20 17:55:35,212 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - yarn.resourcemanager.system-metrics-publisher.enabled is deprecated. Instead, use yarn.system-metrics-publisher.enabled
2020-11-20 17:55:35,301 [main] INFO org.apache.hadoop.yarn.client.RMProxy - Connecting to ResourceManager at cluster-677c-m/10.132.0.16:8032
2020-11-20 17:55:35,571 [main] INFO org.apache.hadoop.yarn.client.AHSProxy - Connecting to Application History server at cluster-677c-m/10.132.0.16:10200
2020-11-20 17:55:35,653 [main] INFO org.apache.pig.tools.pigstats.mapreduce.MRScriptState - Pig script settings are added to the job
2020-11-20 17:55:35,665 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - mapred.job.reduce.markreset.buffer.percent is deprecated. Instead, use mapreduce.reduce.markreset.buffer.percent
2020-11-20 17:55:35,669 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.JobControlCompiler - mapred.job.reduce.markreset.buffer.percent is not set, set to default 0.3
2020-11-20 17:55:35,673 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - mapred.output.compress is deprecated. Instead, use mapreduce.output.fileoutputformat.compress
2020-11-20 17:55:35,679 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.JobControlCompiler - This job cannot be converted run in-process
2020-11-20 17:55:35,705 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - mapred.submit.replication is deprecated. Instead, use mapreduce.client.submit.file.replication
2020-11-20 17:55:36,043 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.JobControlCompiler - Added jar file:/usr/lib/pig/pig-0.17.0-core-h2.jar to DistributedCache through /tmp/temp-381137324/tmp-693856074/pig-0.17.0-core-h2.jar
2020-11-20 17:55:36,085 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.JobControlCompiler - Added jar file:/usr/lib/pig/lib/automaton-1.11-8.jar to DistributedCache through /tmp/temp-381137324/tmp-30294847/automaton-1.11-8.jar
2020-11-20 17:55:36,123 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.JobControlCompiler - Added jar file:/usr/lib/pig/lib/antlr-runtime-3.4.jar to DistributedCache through /tmp/temp-381137324/tmp-472282150/antlr-runtime-3.4.jar
2020-11-20 17:55:36,491 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.JobControlCompiler - Added jar file:/usr/lib/hive/lib/hive-exec-2.3.7.jar to DistributedCache through /tmp/temp-381137324/tmp-125294893/hive-exec-2.3.7.jar
2020-11-20 17:55:36,520 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.JobControlCompiler - Setting up single store job
2020-11-20 17:55:36,543 [main] INFO org.apache.pig.data.SchemaTupleFrontend - Key [pig.schematuple] is false, will not generate code.
2020-11-20 17:55:36,549 [main] INFO org.apache.pig.data.SchemaTupleFrontend - Starting process to move generated code to distributed cache.
2020-11-20 17:55:36,549 [main] INFO org.apache.pig.data.SchemaTupleFrontend - Setting key [pig.schematuple.classes] with classes to deserialize {}
```

```
ssh.cloud.google.com/projects/foocused-mote-292715/zones/europe-west1-b/instances/cluster-677c-m?authuser=1&hl=en_GB&projectNumber=611159050419&useAdminProxy=true

2020-11-20 17:56:19,053 [main] INFO org.apache.hadoop.yarn.client.RMProxy - Connecting to ResourceManager at cluster-677c-m/10.132.0.16:8032
2020-11-20 17:56:19,054 [main] INFO org.apache.hadoop.yarn.client.AHSProxy - Connecting to Application History server at cluster-677c-m/10.132.0.16:10200
2020-11-20 17:56:19,059 [main] INFO org.apache.hadoop.mapred.ClientServiceDelegate - Application state is completed. FinalApplicationStatus=SUCCEEDED. Redirecting to job history server
2020-11-20 17:56:19,064 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - mapred.reduce.tasks is deprecated. Instead, use mapreduce.job.reduces
2020-11-20 17:56:19,075 [main] INFO org.apache.hadoop.yarn.client.AHSProxy - Connecting to Application History server at cluster-677c-m/10.132.0.16:10200
2020-11-20 17:56:19,080 [main] INFO org.apache.hadoop.yarn.client.AHSProxy - Connecting to Application History server at cluster-677c-m/10.132.0.16:10200
2020-11-20 17:56:19,084 [main] INFO org.apache.hadoop.mapred.ClientServiceDelegate - Application state is completed. FinalApplicationStatus=SUCCEEDED. Redirecting to job history server
2020-11-20 17:56:19,157 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - 100% complete
2020-11-20 17:56:19,162 [main] INFO org.apache.pig.tools.pigstats.mapreduce.SimplePigStats - Script Statistics:

HadoopVersion PigVersion UserId StartedAt FinishedAt Features
2.9.2 0.17.0 srjith_umn2 2020-11-20 17:55:35 2020-11-20 17:56:19 UNION

Success!

Job Stats (time in seconds):
JobID Maps Reduces MaxMapTime MinMapTime AvgMapTime MedianMapTime MaxReduceTime MinReduceTime AvgReduceTime MedianReduceTime Alias Feature Outputs
job_1605876410595_0004 5 0 23 14 20 23 0 0 0 A,B,C,D,E,I,J,K,L,M,R MAP_ONLY /stackex/output_final.txt,

Input(s):
Successfully read 16295 records from: "/stackex/ViewCount5.txt"
Successfully read 48135 records from: "/stackex/ViewCount1.txt"
Successfully read 44473 records from: "/stackex/ViewCount2.txt"
Successfully read 46924 records from: "/stackex/ViewCount3.txt"
Successfully read 44177 records from: "/stackex/ViewCount4.txt"

Output(s):
Successfully stored 200004 records (34673280 bytes) in: "/stackex/output_final.txt"

Counters:
Total records written : 200004
Total bytes written : 34673280
Spillable Memory Manager spill count : 0
Total bags proactively spilled: 0
Total records proactively spilled: 0

Job DAG:
job_1605876410595_0004

2020-11-20 17:56:19,167 [main] INFO org.apache.hadoop.yarn.client.RMProxy - Connecting to ResourceManager at cluster-677c-m/10.132.0.16:8032
2020-11-20 17:56:19,176 [main] INFO org.apache.hadoop.yarn.client.AHSProxy - Connecting to Application History server at cluster-677c-m/10.132.0.16:10200
2020-11-20 17:56:19,185 [main] INFO org.apache.hadoop.mapred.ClientServiceDelegate - Application state is completed. FinalApplicationStatus=SUCCEEDED. Redirecting to job history server
2020-11-20 17:56:19,242 [main] INFO org.apache.hadoop.yarn.client.RMProxy - Connecting to ResourceManager at cluster-677c-m/10.132.0.16:8032
2020-11-20 17:56:19,248 [main] INFO org.apache.hadoop.yarn.client.AHSProxy - Connecting to Application History server at cluster-677c-m/10.132.0.16:10200
2020-11-20 17:56:19,261 [main] INFO org.apache.hadoop.mapred.ClientServiceDelegate - Application state is completed. FinalApplicationStatus=SUCCEEDED. Redirecting to job history server
2020-11-20 17:56:19,291 [main] INFO org.apache.hadoop.yarn.client.RMProxy - Connecting to ResourceManager at cluster-677c-m/10.132.0.16:8032
2020-11-20 17:56:19,294 [main] INFO org.apache.hadoop.yarn.client.AHSProxy - Connecting to Application History server at cluster-677c-m/10.132.0.16:10200
2020-11-20 17:56:19,300 [main] INFO org.apache.hadoop.mapred.ClientServiceDelegate - Application state is completed. FinalApplicationStatus=SUCCEEDED. Redirecting to job history server
2020-11-20 17:56:19,352 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Success!

grunt>
```

## Browse Directory

Show 25 entries

Search:

<input type="checkbox"/>	Permission	Owner	Group	Size	Last Modified	Replication	Block Size	Name	
<input type="checkbox"/>	-rw-r--r--	<a href="#">srijith_unni2</a>	<a href="#">hadoop</a>	0 B	Nov 20 23:26	2	128 MB	<a href="#">_SUCCESS</a>	
<input type="checkbox"/>	-rw-r--r--	<a href="#">srijith_unni2</a>	<a href="#">hadoop</a>	8.23 MB	Nov 20 23:26	2	128 MB	<a href="#">part-m-00000</a>	
<input type="checkbox"/>	-rw-r--r--	<a href="#">srijith_unni2</a>	<a href="#">hadoop</a>	7.67 MB	Nov 20 23:26	2	128 MB	<a href="#">part-m-00001</a>	
<input type="checkbox"/>	-rw-r--r--	<a href="#">srijith_unni2</a>	<a href="#">hadoop</a>	7.3 MB	Nov 20 23:26	2	128 MB	<a href="#">part-m-00002</a>	
<input type="checkbox"/>	-rw-r--r--	<a href="#">srijith_unni2</a>	<a href="#">hadoop</a>	7.21 MB	Nov 20 23:26	2	128 MB	<a href="#">part-m-00003</a>	
<input type="checkbox"/>	-rw-r--r--	<a href="#">srijith_unni2</a>	<a href="#">hadoop</a>	2.65 MB	Nov 20 23:26	2	128 MB	<a href="#">part-m-00004</a>	

Showing 1 to 6 of 6 entries

Hadoop, 2018.

## Browse Directory

Show 25 entries

Search:

<input type="checkbox"/>	Permission	Owner	Group	Size	Last Modified	Replication	Block Size	Name	
<input type="checkbox"/>	-rw-r--r--	<a href="#">srijith_unni2</a>	<a href="#">hadoop</a>	0 B	Nov 20 23:26	2	128 MB	<a href="#">_SUCCESS</a>	
<input checked="" type="checkbox"/>	-rw-r--r--	<a href="#">srijith_unni2</a>	<a href="#">hadoop</a>	8.23 MB	Nov 20 23:26	2	128 MB	<a href="#">part-m-00000</a>	
<input checked="" type="checkbox"/>	-rw-r--r--	<a href="#">srijith_unni2</a>	<a href="#">hadoop</a>	7.67 MB	Nov 20 23:26	2	128 MB	<a href="#">part-m-00001</a>	
<input checked="" type="checkbox"/>	-rw-r--r--	<a href="#">srijith_unni2</a>	<a href="#">hadoop</a>	7.3 MB	Nov 20 23:26	2	128 MB	<a href="#">part-m-00002</a>	
<input checked="" type="checkbox"/>	-rw-r--r--	<a href="#">srijith_unni2</a>	<a href="#">hadoop</a>	7.21 MB	Nov 20 23:26	2	128 MB	<a href="#">part-m-00003</a>	
<input checked="" type="checkbox"/>	-rw-r--r--	<a href="#">srijith_unni2</a>	<a href="#">hadoop</a>	2.65 MB	Nov 20 23:26	2	128 MB	<a href="#">part-m-00004</a>	

Showing 1 to 6 of 6 entries

Hadoop, 2018.

## Browse Directory

Show 25 entries

Search:

<input type="checkbox"/>	Permission	Owner	Group	Size	Last Modified	Replication	Block Size	Name	
--------------------------	------------	-------	-------	------	---------------	-------------	------------	------	--

No data available in table

Showing 0 to 0 of 0 entries

Hadoop, 2018.



## Browse Directory

Show 25 entries

Search:

	Permission	Owner	Group	Size	Last Modified	Replication	Block Size	Name
<input type="checkbox"/>	-rw-r--r--	srijith_unni2	hadoop	8.23 MB	Nov 20 23:26	2	128 MB	part-m-00000
<input type="checkbox"/>	-rw-r--r--	srijith_unni2	hadoop	7.67 MB	Nov 20 23:26	2	128 MB	part-m-00001
<input type="checkbox"/>	-rw-r--r--	srijith_unni2	hadoop	7.3 MB	Nov 20 23:26	2	128 MB	part-m-00002
<input type="checkbox"/>	-rw-r--r--	srijith_unni2	hadoop	7.21 MB	Nov 20 23:26	2	128 MB	part-m-00003
<input type="checkbox"/>	-rw-r--r--	srijith_unni2	hadoop	2.65 MB	Nov 20 23:26	2	128 MB	part-m-00004

Showing 1 to 5 of 5 entries

Previous

1

Next

Hadoop, 2018.

```
srijith_unni2@cluster-677c-mc - Google Chrome
ssh.cloud.google.com/projects/focused-mote-292715/zones/europe-west1-b/instances/cluster-677c-m?authuser=1&hl=en_GB&projectNumber=611159050419&useAdminProxy=true
srijith_unni2@cluster-677c-m:~$ hadoop fs -ls /stackex/final_result
Found 5 items
-rw-r--r-- 2 srijith_unni2 hadoop 8630105 2020-11-20 17:56 /stackex/final_result/part-m-00000
-rw-r--r-- 2 srijith_unni2 hadoop 8044423 2020-11-20 17:56 /stackex/final_result/part-m-00001
-rw-r--r-- 2 srijith_unni2 hadoop 7658192 2020-11-20 17:56 /stackex/final_result/part-m-00002
-rw-r--r-- 2 srijith_unni2 hadoop 7556999 2020-11-20 17:56 /stackex/final_result/part-m-00003
-rw-r--r-- 2 srijith_unni2 hadoop 2783561 2020-11-20 17:56 /stackex/final_result/part-m-00004
srijith_unni2@cluster-677c-m:~$ hadoop fs -getmerge /stackex/final_result/part-m-00000 /stackex/final_result/part-m-00001 /stackex/final_result/part-m-00002 /stackex/final_result/part-m-00003 /stackex/final_result/part-m-00004 /home/srijith_unni2/result_final.txt
srijith_unni2@cluster-677c-m:~$ hadoop fs -put /home/srijith_unni2/result_final.txt /stackex/final_result
srijith_unni2@cluster-677c-m:~$
```

## Browse Directory

Show 25 entries

Search:

	Permission	Owner	Group	Size	Last Modified	Replication	Block Size	Name
<input type="checkbox"/>	-rw-r--r--	srijith_unni2	hadoop	8.23 MB	Nov 20 23:26	2	128 MB	part-m-00000
<input type="checkbox"/>	-rw-r--r--	srijith_unni2	hadoop	7.67 MB	Nov 20 23:26	2	128 MB	part-m-00001
<input type="checkbox"/>	-rw-r--r--	srijith_unni2	hadoop	7.3 MB	Nov 20 23:26	2	128 MB	part-m-00002
<input type="checkbox"/>	-rw-r--r--	srijith_unni2	hadoop	7.21 MB	Nov 20 23:26	2	128 MB	part-m-00003
<input type="checkbox"/>	-rw-r--r--	srijith_unni2	hadoop	2.65 MB	Nov 20 23:26	2	128 MB	part-m-00004
<input type="checkbox"/>	-rw-r--r--	srijith_unni2	hadoop	33.07 MB	Nov 20 23:37	2	128 MB	result_final.txt

Showing 1 to 6 of 6 entries

Previous

1

Next

Hadoop, 2018.

focused-mote-292715 > cluster-677c Sign out

Hadoop Overview Datanodes Pathways Volume Failures Snapshot Status Progress Metrics

## Browse Directory

/stackex/final\_result

Show 25 entries

Permission	Owner
-rW-r--r--	srijith_uni2
-rW-r--r--	srijith_uni2
-rW-r--r--	srijith_uni2
-rW-r--r--	srijith_uni2
-rW-r--r--	srijith_uni2
-rW-r--r--	srijith_uni2
-rW-r--r--	srijith_uni2

Showing 1 to 6 of 6 entries

Hadoop, 2018.

File information - result\_final.txt

Download Head the file (first 32K) Tail the file (last 32K)

Block information - Block 0

Block ID: 1073741902  
Block Pool ID: BP-842087565-10.132.0.16-1605876406184  
Generation Stamp: 1078  
Size: 34673280  
Availability:

- cluster-677c-w-1.c.focused-mote-292715.internal
- cluster-677c-w-0.c.focused-mote-292715.internal

Close

Block Size Name

MB	part-m-00000
MB	part-m-00001
MB	part-m-00002
MB	part-m-00003
MB	part-m-00004
MB	result_final.txt

Search:

Go

Previous 1 Next

```
srijith_uni2@cluster-677c-m: ~ - Google Chrome
ssh.cloud.google.com/projects/focused-mote-292715/zones/europe-west1-b/instances/cluster-677c-m?authuser=1&hl=en_GB&projectNumber=611159050419&useAdminProxy=true
srijith_uni2@cluster-677c-m:~$ hive

Logging initialized using configuration in file:/etc/hive/conf.dist/hive-log4j2.properties Async: true
hive> CREATE DATABASE stackexdb;
OK
Time taken: 0.649 seconds
hive> CREATE EXTERNAL TABLE stackexdb.posts ( id int comment 'Id', acceptedanswerid int comment "AcceptedAnswerId", creationdate string comment "CreationDate", score int comment "Score", viewcount int comment "ViewCount", owneruserid int comment "OwnerUserId", lasteditoruserid int comment "LastEditorUserId", lasteditdate string comment "LastEditDate", lastactivitytime string comment "LastActivityTime", title string comment "Title", tags string comment "Tags", answercount int comment "AnswerCount", commentcount int comment "CommentCount", favouritecount int comment "FavouriteCount") ROW FORMAT DELIMITED FIELDS TERMINATED BY '\t' LOCATION '/stackex';
OK
Time taken: 0.429 seconds
hive> LOAD DATA LOCAL INPATH 'result_final.txt' INTO TABLE stackexdb.posts;
Loading data to table stackexdb.posts
OK
Time taken: 1.818 seconds
hive>
```

```
srijith_uni2@cluster-e7d2-m: ~ - Google Chrome
ssh.cloud.google.com/projects/focused-mote-292715/zones/europe-west1-b/instances/cluster-e7d2-m?authuser=1&hl=en_GB&projectNumber=611159050419&useAdminProxy=true
hive> set hive.cli.print.header=true;
hive>
```

```
srijith_uni2@cluster-e7d2-m: ~ - Google Chrome
ssh.cloud.google.com/projects/focused-mote-292715/zones/europe-west1-b/instances/cluster-e7d2-m?authuser=1&hl=en_GB&projectNumber=611159050419&useAdminProxy=true
hive> CREATE TABLE stackexdb.top10_post ROW FORMAT DELIMITED FIELDS TERMINATED BY '\t' AS SELECT DISTINCT id, score, viewcount, answercount, commentcount, favouritecount, title FROM stackexdb.posts ORDER BY score DESC LIMIT 10;
Query ID = srijith_uni2_20201121084719_0e5c1451-0629-4490-a21b-cfde34ce200f
Total jobs = 1
Launching Job 1 out of 1
Status: Running (Executing on YARN cluster with App id application_1605944478645_0003)

-----
VERTICES      MODE        STATUS      TOTAL  COMPLETED  RUNNING  PENDING  FAILED  KILLED
-----
Map 1 ..... container SUCCEEDED      1         1         0         0         0         0
Reducer 2 ..... container SUCCEEDED      1         1         0         0         0         0
Reducer 3 ..... container SUCCEEDED      1         1         0         0         0         0
-----
VERTICES: 03/03 [=====] 100% ELAPSED TIME: 11.30 s
-----
Moving data to directory hdfs://cluster-e7d2-m/user/hive/warehouse/stackexdb.db/top10_post
OK
id      score  viewcount  answercount  commentcount  favouritecount  title
Time taken: 13.27 seconds
hive> SELECT * FROM stackexdb.top10_post;
OK
top10_post.id  top10_post.score  top10_post.viewcount  top10_post.answercount  top10_post.commentcount  top10_post.favouritecount  top10_post.title
11227809      24969  1541165  26      3      11149  "Why is processing a sorted array faster than processing an unsorted array?"
927358  21777  9000893  86      14      6970  "How do I undo the most recent local commits in Git?"
2003505  17395  8435222  40      7      5477  "How do I delete a Git branch locally and remotely?"
293357  12200  2468645  36      9      2382  "What is the difference between 'git pull' and 'git fetch'?"
231767  10627  2333118  42      0      5902  "What does the 'yield' keyword do?"
477816  10467  2946855  36      0      1463  "What is the correct JSON content type?"
348170  9309   3316501  37      11     1590  "How do I undo 'git add' before commit?"
1642028  9174   820505  24      26     2108  "What is the '***' operator in C++?"
6591213  8919   3122723  34      0      1293  "How do I rename a local Git branch?"
5767325  8762   7275538  97      1     1349  "How can I remove a specific item from an array?"
Time taken: 0.289 seconds, Fetched: 10 row(s)
hive>
```



```
srjith_umn2@cluster-e7d2-nc - Google Chrome
ssh.cloud.google.com/projects/focused-mote-292715/zones/europe-west1-b/instances/cluster-e7d2-m?authuser=1&hl=en_GB&projectNumber=611159050419&useAdminProxy=true
hive> CREATE TABLE stackexdb.top10_user ROW FORMAT DELIMITED FIELDS TERMINATED BY '\t' AS SELECT DISTINCT id, score, viewcount, owneruserid, title FROM stackexdb.posts ORDER BY score DESC LIMIT 10;
Query ID = srjith_umn2_20201121084834_7e765ae0-afe9-4294-b089-ace04b632f05
Total Jobs = 1
Launching Job 1 out of 1
Status: Running (Executing on YARN cluster with App id application_1605944478645_0003)

-----
VERTICES      MODE      STATUS  TOTAL  COMPLETED  RUNNING  PENDING  FAILED  KILLED
-----
Map 1 ..... container  SUCCEEDED  1      1      0      0      0      0
Reducer 2 ..... container  SUCCEEDED  1      1      0      0      0      0
Reducer 3 ..... container  SUCCEEDED  1      1      0      0      0      0
-----
VERTICES: 03/03 [=====] 100% ELAPSED TIME: 11.67 s
-----
Moving data to directory hdfs://cluster-e7d2-m/user/hive/warehouse/stackexdb.db/top10_user
OK
id      score  viewcount  owneruserid  title
Time taken: 13.469 seconds
hive> SELECT * FROM stackexdb.top10_user;
OK
top10_user.id  top10_user.score  top10_user.viewcount  top10_user.owneruserid  top10_user.title
11227809      24965   1541145  87234   "Why is processing a sorted array faster than processing an unsorted array?"
927358      21777   9000893  89904   "How do I undo the most recent local commits in Git?"
2003505      17395   8435222  95592   "How do I delete a Git branch locally and remotely?"
292357      12200   2868445  6068    "What is the difference between 'git pull' and 'git fetch'?"
231767      10627   2333118  18300   "What does the 'yield' keyword do?"
477816      10467   2946855  12870   "What is the correct JSON content type?"
348170      9309    3316501  14069   "How do I undo 'git add' before commit?"
1642028      9174    820505   87234   "What is the '--->' operator in C++?"
4591213      8919    3122723  338204   "How do I rename a local Git branch?"
5767325      8762    7275538  364969   "How can I remove a specific item from an array?"
Time taken: 0.301 seconds, Fetched: 10 row(s)
hive>
```

```
srjith_umn2@cluster-e7d2-nc - Google Chrome
ssh.cloud.google.com/projects/focused-mote-292715/zones/europe-west1-b/instances/cluster-e7d2-m?authuser=1&hl=en_GB&projectNumber=611159050419&useAdminProxy=true
hive> CREATE TABLE stackexdb.users_hadoop ROW FORMAT DELIMITED FIELDS TERMINATED BY '\t' AS SELECT DISTINCT id, owneruserid, title, tags FROM stackexdb.posts where tags rlike '.*(Hadoop|hadoop)';
Query ID = srjith_umn2_20201121084156_001eb096-2b82-4a80-a879-2f64173f5ed8
Total Jobs = 1
Launching Job 1 out of 1
Tex session was closed. Reopening...
Session re-established.
Status: Running (Executing on YARN cluster with App id application_1605944478645_0003)

-----
VERTICES      MODE      STATUS  TOTAL  COMPLETED  RUNNING  PENDING  FAILED  KILLED
-----
Map 1 ..... container  SUCCEEDED  1      1      0      0      0      0
Reducer 2 ..... container  SUCCEEDED  1      1      0      0      0      0
-----
VERTICES: 02/02 [=====] 100% ELAPSED TIME: 15.88 s
-----
Moving data to directory hdfs://cluster-e7d2-m/user/hive/warehouse/stackexdb.db/users_hadoop
OK
id      owneruserid  title  tags
Time taken: 30.578 seconds
hive> SELECT * FROM stackexdb.users_hadoop;
OK
users_hadoop.id  users_hadoop.owneruserid  users_hadoop.title  users_hadoop.tags
24179   2588   How does Hive compare to HBase? <hadoop><hbase><hive>
339344  30861   Is there a .NET equivalent to Apache Hadoop? <c#><.net><hadoop><mapreduce>
1182732 114196   How does the MapReduce sort algorithm work? <algorithm><sorting><parallel-processing><hadoop><mapreduce>
1482282 123067   Java vs Python on Hadoop <java><python><hadoop>
1533330 41717   Writing data to Hadoop <hadoop><hdfs>
2354525 77308   What should be hadoop.tmp.dir ? <hadoop><hdfs><config>
2358402 246677   Where HDFS stores files locally by default? <hadoop><hdfs>
2489585 114196   Chaining multiple MapReduce jobs in Hadoop <hadoop><mapreduce>
2669800 305105   Change block size of dfs file <hadoop>
2674421 152253   Free Large datasets to experiment with Hadoop <resources><hadoop><opendata>
2821507 154886   How does Hadoop perform input splits? <hadoop><mapreduce><hdfs>
3207218 218900   Where does hadoop mapreduce framework send its System.out.print() statements ? (stdout) <hadoop><mapreduce>
3356259 68920   "Difference between Pig and Hive? Why have both?" <hadoop><hive><apache-pig>
3515481 68920   Pig Latin: Load multiple files from a date range (part of the directory structure) <hadoop><apache-pig>
3548259 428495   Merging multiple files into one within Hadoop <hadoop><apache-pig>
4065999 211528   "Does Hive have a String split function?" <hadoop><hive>
4716961 215971   Hadoop copy a directory? <hadoop><hdfs>
5058400 2819    "Where does Hive store files in HDFS?" <hadoop><hive><hdfs>
5293446 71834   "HDFS error: could only be replicated to 0 nodes, instead of 1" <amazon-ec2><hadoop>
5377118 447952   How to convert .txt file to Hadoop's sequence file format <java><file><hadoop><type-conversion><hive>
5385163 319013   Create temporary table in Hive? <hadoop><hive>
5571156 463286   "Hadoop, how to compress mapper output but not the reducer output" <compression><hadoop><hdfs>
5700048 458560   merge output files after reduce phase <hadoop><mapreduce>
5864589 572138   "How to fix 'task attempt 201104251139_0295_r_000006.0 failed to report status for 600 seconds,'" <hadoop><mapreduce>
6153560 145360   Hbase client ConnectionLoss for /hbase error <java><ruby><hadoop><hbase><thrift>
6297533 324968   Search/Find a file and file content in Hadoop <file><filesystems><hadoop><distributed><distributed-computing>
6445339 578961   "COLLECT_SET() in Hive, keep duplicates?" <java><hadoop><user-defined-function><hive>
```

```
srjith_umn2@cluster-e7d2-nc - Google Chrome
ssh.cloud.google.com/projects/focused-mote-292715/zones/europe-west1-b/instances/cluster-e7d2-m?authuser=1&hl=en_GB&projectNumber=611159050419&useAdminProxy=true
srjith_umn2@cluster-e7d2-nc:~$ hadoop fs -ls /user/hive/warehouse/stackexdb.db:
Found 3 items
drwxrwxrwt - srjith_umn2 hadoop 0 2020-11-21 08:47 /user/hive/warehouse/stackexdb.db/top10_post
drwxrwxrwt - srjith_umn2 hadoop 0 2020-11-21 08:48 /user/hive/warehouse/stackexdb.db/top10_user
drwxrwxrwt - srjith_umn2 hadoop 0 2020-11-21 08:42 /user/hive/warehouse/stackexdb.db/users_hadoop
srjith_umn2@cluster-e7d2-nc:~$
```

Browse Directory

Go!

Show 

25

 entries 

Search:

<input type="checkbox"/>	<div><div></div><div></div></div> Permission	<div><div></div><div></div></div> Owner	<div><div></div><div></div></div> Group	<div><div></div><div></div></div> Size	<div><div></div><div></div></div> Last Modified	<div><div></div><div></div></div> Replication	<div><div></div><div></div></div> Block Size	<div><div></div><div></div></div> Name	<div><div></div><div></div></div>
<input type="checkbox"/>	drwxrwxrwt	<a href="#">srijith_unni2</a>	<a href="#">hadoop</a>	0 B	Nov 21 14:17	0	0 B	<a href="#">top10_post</a>	<div></div>
<input type="checkbox"/>	drwxrwxrwt	<a href="#">srijith_unni2</a>	<a href="#">hadoop</a>	0 B	Nov 21 14:18	0	0 B	<a href="#">top10_user</a>	<div></div>
<input type="checkbox"/>	drwxrwxrwt	<a href="#">srijith_unni2</a>	<a href="#">hadoop</a>	0 B	Nov 21 14:12	0	0 B	<a href="#">users_hadoop</a>	<div></div>

Showing 1 to 3 of 3 entries

Previous

1

Next