

INFSCI 2750 - Mini Project 02

Yichi Zhang (yiz141@pitt.edu)

Quan Zhou (quz3@pitt.edu)

Part 1: Setting up Spark

In the previous project, we've already finished setting up the Hadoop cluster. In order to set up a Spark environment on top the Hadoop cluster, the following steps were taken:

Add new environment variables to the `~/.bashrc` file.

```
# Add Spark installation path to $PATH
export PATH="/usr/lib/jvm/java-8-oracle:/usr/local/spark/bin":$PATH
# Add three new environment variable
export HADOOP_CONF_DIR=/usr/local/hadoop/etc/hadoop
export SPARK_HOME=/usr/local/spark
export LD_LIBRARY_PATH=/usr/local/hadoop/lib/native:$LD_LIBRARY_PATH
```

Download and install the pre-built version of Spark and install it on the cluster.

```
cd /usr/local/
wget http://mirrors.gigenet.com/apache/spark/spark-2.3.0/spark-2.3.0-bin-hadoop2.7.tgz
tar -xzf spark-2.3.0-bin-hadoop2.7.tgz
ln -s spark-2.3.0-bin-hadoop2.7/ spark
```

Edit the configuration file `$SPARK_HOME/conf/spark-defaults.conf` for Spark

<code>spark.master</code>	<code>yarn</code>
<code>spark.driver.memory</code>	<code>1g</code>
<code>spark.yarn.am.memory</code>	<code>1g</code>
<code>spark.executor.memory</code>	<code>1g</code>
<code>spark.eventLog.enabled</code>	<code>true</code>
<code>spark.eventLog.dir</code>	<code>hdfs://master:9000/spark-logs</code>
<code>spark.history.provider</code>	<code>org.apache.spark.deploy.history.FsHistoryProvider</code>
<code>spark.history.fs.logDirectory</code>	<code>hdfs://master:9000/spark-logs</code>
<code>spark.history.fs.update.interval</code>	<code>10s</code>
<code>spark.history.ui.port</code>	<code>18080</code>
<code>spark.executor.cores</code>	<code>2</code>

With all the preparation steps done, we can start the Hadoop cluster and Spark History Server by commands:

```
$HADOOP_HOME/sbin/start-all.sh
$HADOOP_HOME/bin/mapred --daemon start historyserver
```

```
$SPARK_HOME/sbin/start-history-server.sh
```

We can track all the Hadoop related processes with the `jps` command.

```
Isolachine — root@master: /usr/local/spark — ssh root@159.89.43.89 — 80x11
Last login: Sun Mar 11 15:18:58 2018 from 24.3.23.164
[root@master:~# jps
10690 NameNode
12003 HistoryServer
25411 NodeManager
11908 JobHistoryServer
10839 DataNode
25257 ResourceManager
11067 SecondaryNameNode
8255 Jps
[root@master:~# cd /usr/local/spark
```

Screenshot of running processes related to Hadoop

All the services on our cluster are running properly. We can monitor the cluster's status and performance through the following webpages.

Monitor	Link
HDFS NameNode	http://159.89.43.89:9870
ResourceManager of Yarn	http://159.89.43.89:8088
Spark History Server	http://159.89.43.89:18080
Hadoop Job History Server	http://159.89.43.89:19888

To start a Spark Shell on the cluster, simply run:

```
$SPARK_HOME/bin/spark-shell --master yarn --deploy-mode client
```

```
root@master:/usr/local/spark# $SPARK_HOME/bin/spark-shell --master yarn --deploy-mode client --driver-memory 2g --executor-memory 2g --executor-cores 1
Setting default log level to "WARN".
To adjust logging level use sc.setLogLevel(newLevel). For SparkR, use setLogLevel(newLevel).
2018-03-18 15:59:44 WARN Client:66 - Neither spark.yarn.jars nor spark.yarn.archive is set, falling back to uploading libraries under SPARK_HOME.
Spark context Web UI available at http://master:4040
Spark context available as 'sc' (master = yarn, app id = application_1521384446685_0008).
Spark session available as 'spark'.
Welcome to

  ____      __
 / ___/____/ /_  __
/ /  / __/ __/ / / /
/ /__/ /_/ /_/ /_/ /
/____/____/_/_/_/_/

version 2.3.0

Using Scala version 2.11.8 (Java HotSpot(TM) 64-Bit Server VM, Java 1.8.0_161)
Type in expressions to have them evaluated.
Type :help for more information.

scala>
```

Screenshot of Spark Shell running on VM

```
Isolachine — root@master: /usr/local/spark — ssh root@159.89.43.89 — 80x10
scala> val textFile = spark.read.textFile("input/hadoop-env.sh")
textFile: org.apache.spark.sql.Dataset[String] = [value: string]

scala> textFile.count()
res3: Long = 414

scala> textFile.first()
res4: String = export JAVA_HOME=/usr/lib/jvm/java-8-oracle

scala> █
```

Simple test run in Spark Shell

The whole project's source code was written in JAVA and used Maven as dependency management and build tool.

The JAVA code and `pom.xml` file for the project is in the `source_code` folder. The `mini-project-02-1.0.jar` was built locally with maven to include all the source code provided and was uploaded to the VM server for running.

```
mvn package
scp target/mini-project-02-1.0.jar root@159.89.43.89:/usr/local/spark
```

To track the log of Spark program run on the cluster, we added edited `yarn-site.xml` make logs all output to the HDFS in a aggregated manner.

```
<property>
  <name>yarn.log-aggregation-enable</name>
  <value>true</value>
</property>
```

So one can simply print the aggregated `stdout` logs with the following command:

```
$HADOOP_HOME/bin/yarn logs -log_files stdout -applicationId <AppID>
```

Part 2 : Last FM Data Analysis

The `ListeningCount.java` file is the source code for counting total listening counts of each artist.

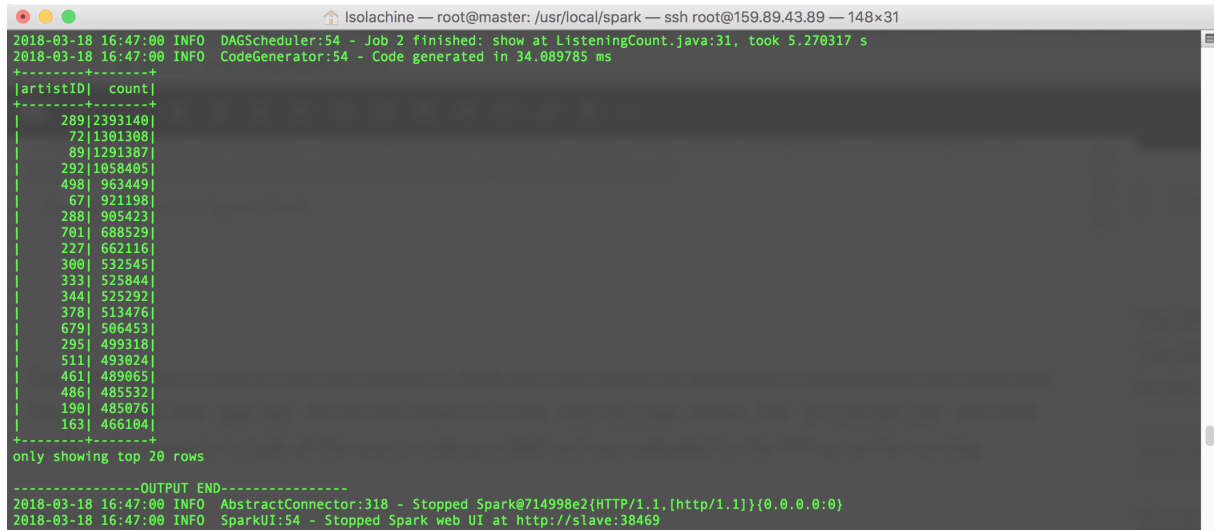
The program can be launched by:

```
$SPARK_HOME/bin/spark-submit \
```

```
--deploy-mode cluster \  
--class ListeningCount \  
mini-project-02-1.0.jar
```

The logs can be viewed by:

```
$SHADOOP_HOME/bin/yarn logs -log_files stdout -applicationId application_1521391514647_0001
```



The screenshot shows a terminal window with the following content:

```
2018-03-18 16:47:00 INFO DAGScheduler:54 - Job 2 finished: show at ListeningCount.java:31, took 5.270317 s  
2018-03-18 16:47:00 INFO CodeGenerator:54 - Code generated in 34.089785 ms  
+-----+  
|artistID| count|  
+-----+  
| 289|2393140|  
| 72|1301308|  
| 89|1291387|  
|292|1058405|  
|498| 963449|  
| 67| 921198|  
|288| 905423|  
|701| 688529|  
|227| 662116|  
|300| 532545|  
|333| 525844|  
|344| 525292|  
|378| 513476|  
|679| 506453|  
|295| 499318|  
|511| 493024|  
|461| 489065|  
|486| 485532|  
|190| 485076|  
|163| 466104|  
+-----+  
only showing top 20 rows  
-----OUTPUT END-----  
2018-03-18 16:47:00 INFO AbstractConnector:318 - Stopped Spark@714998e2{HTTP/1.1,[http/1.1]}{0.0.0.0:0}  
2018-03-18 16:47:00 INFO SparkUI:54 - Stopped Spark web UI at http://slave:38469
```

Screenshot of Top 20 Listened Artists and Count

The job summary can be viewed at:

```
http://159.89.43.89:18080/history/application_1521391514647_0001/1/jobs/
```

Part 3: Access Log Analysis

Problem 1

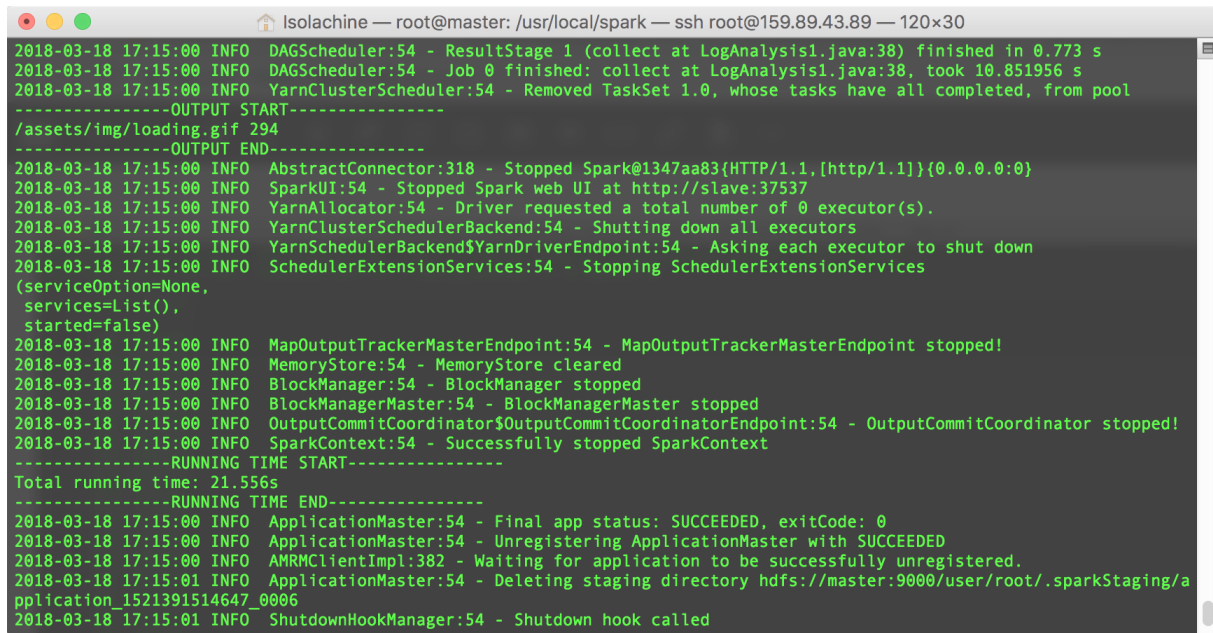
The `LogAnalysis1.java` file is the source code for problem 1.

The program can be launched by:

```
$SPARK_HOME/bin/spark-submit \  
--deploy-mode cluster \  
--class LogAnalysis1 \  
mini-project-02-1.0.jar
```

The logs can be viewed by:

```
$SHADOOP_HOME/bin/yarn logs -log_files stdout -applicationId application_1521391514647_0006
```



```
2018-03-18 17:15:00 INFO DAGScheduler:54 - ResultStage 1 (collect at LogAnalysis1.java:38) finished in 0.773 s
2018-03-18 17:15:00 INFO DAGScheduler:54 - Job 0 finished: collect at LogAnalysis1.java:38, took 10.851956 s
2018-03-18 17:15:00 INFO YarnClusterScheduler:54 - Removed TaskSet 1.0, whose tasks have all completed, from pool
-----OUTPUT START-----
/assets/img/loading.gif 294
-----OUTPUT END-----
2018-03-18 17:15:00 INFO AbstractConnector:318 - Stopped Spark@1347aa83(HTTP/1.1,[http/1.1]){0.0.0.0:0}
2018-03-18 17:15:00 INFO SparkUI:54 - Stopped Spark web UI at http://slave:37537
2018-03-18 17:15:00 INFO YarnAllocator:54 - Driver requested a total number of 0 executor(s).
2018-03-18 17:15:00 INFO YarnClusterSchedulerBackend:54 - Shutting down all executors
2018-03-18 17:15:00 INFO YarnSchedulerBackend$YarnDriverEndpoint:54 - Asking each executor to shut down
2018-03-18 17:15:00 INFO SchedulerExtensionServices:54 - Stopping SchedulerExtensionServices
(serviceOption=None,
 services=List(),
 started=false)
2018-03-18 17:15:00 INFO MapOutputTrackerMasterEndpoint:54 - MapOutputTrackerMasterEndpoint stopped!
2018-03-18 17:15:00 INFO MemoryStore:54 - MemoryStore cleared
2018-03-18 17:15:00 INFO BlockManager:54 - BlockManager stopped
2018-03-18 17:15:00 INFO BlockManagerMaster:54 - BlockManagerMaster stopped
2018-03-18 17:15:00 INFO OutputCommitCoordinator$OutputCommitCoordinatorEndpoint:54 - OutputCommitCoordinator stopped!
2018-03-18 17:15:00 INFO SparkContext:54 - Successfully stopped SparkContext
-----RUNNING TIME START-----
Total running time: 21.556s
-----RUNNING TIME END-----
2018-03-18 17:15:00 INFO ApplicationMaster:54 - Final app status: SUCCEEDED, exitCode: 0
2018-03-18 17:15:00 INFO ApplicationMaster:54 - Unregistering ApplicationMaster with SUCCEEDED
2018-03-18 17:15:00 INFO AMRMClientImpl:382 - Waiting for application to be successfully unregistered.
2018-03-18 17:15:01 INFO ApplicationMaster:54 - Deleting staging directory hdfs://master:9000/user/root/.sparkStaging/a
pplication_1521391514647_0006
2018-03-18 17:15:01 INFO ShutdownHookManager:54 - Shutdown hook called
```

Screenshot of Results and Running Time

As the screen shot shows, `/assets/img/loading.gif` was accessed 294 times.

The job summary can be viewed at:

```
http://159.89.43.89:18080/history/application_1521391514647_0006/1/jobs/
```

Problem 2

The `LogAnalysis2.java` file is the source code for problem 2.

The program can be launched by:

```
$SPARK_HOME/bin/spark-submit \
--deploy-mode cluster \
--class LogAnalysis2 \
mini-project-02-1.0.jar
```

The logs can be viewed by:

```
$SHADOOP_HOME/bin/yarn logs -log_files stdout -applicationId application_1521391514647_0009
```

```
Isolachine — root@master: /usr/local/spark — ssh root@159.89.43.89 — 120x30
2018-03-18 17:24:37 INFO TaskSetManager:54 - Finished task 2.0 in stage 1.0 (TID 6) in 772 ms on master (executor 2) (4 /4)
2018-03-18 17:24:37 INFO YarnClusterScheduler:54 - Removed TaskSet 1.0, whose tasks have all completed, from pool
2018-03-18 17:24:37 INFO DAGScheduler:54 - ResultStage 1 (collect at LogAnalysis2.java:38) finished in 0.791 s
2018-03-18 17:24:37 INFO DAGScheduler:54 - Job 0 finished: collect at LogAnalysis2.java:38, took 9.387901 s
-----OUTPUT START-----
/assets/js/lightbox.js 297
-----OUTPUT END-----
2018-03-18 17:24:37 INFO AbstractConnector:318 - Stopped Spark@dca7cdc{HTTP/1.1,[http/1.1]}{0.0.0.0:0}
2018-03-18 17:24:37 INFO SparkUI:54 - Stopped Spark web UI at http://slave:35217
2018-03-18 17:24:37 INFO YarnAllocator:54 - Driver requested a total number of 0 executor(s).
2018-03-18 17:24:37 INFO YarnClusterSchedulerBackend:54 - Shutting down all executors
2018-03-18 17:24:37 INFO YarnSchedulerBackend$YarnDriverEndpoint:54 - Asking each executor to shut down
2018-03-18 17:24:37 INFO SchedulerExtensionServices:54 - Stopping SchedulerExtensionServices
(serviceOption=None,
 services=List(),
 started=false)
2018-03-18 17:24:37 INFO MapOutputTrackerMasterEndpoint:54 - MapOutputTrackerMasterEndpoint stopped!
2018-03-18 17:24:37 INFO MemoryStore:54 - MemoryStore cleared
2018-03-18 17:24:37 INFO BlockManager:54 - BlockManager stopped
2018-03-18 17:24:37 INFO BlockManagerMaster:54 - BlockManagerMaster stopped
2018-03-18 17:24:37 INFO OutputCommitCoordinator$OutputCommitCoordinatorEndpoint:54 - OutputCommitCoordinator stopped!
2018-03-18 17:24:37 INFO SparkContext:54 - Successfully stopped SparkContext
-----RUNNING TIME START-----
Total running time: 19.993s
-----RUNNING TIME END-----
2018-03-18 17:24:37 INFO ApplicationMaster:54 - Final app status: SUCCEEDED, exitCode: 0
2018-03-18 17:24:37 INFO ApplicationMaster:54 - Unregistering ApplicationMaster with SUCCEEDED
2018-03-18 17:24:37 INFO AMRMClientImpl:382 - Waiting for application to be successfully unregistered.
2018-03-18 17:24:37 INFO ApplicationMaster:54 - Deleting staging directory hdfs://master:9000/user/root/.sparkStaging/a
```

Screenshot of Results and Running Time

As the screen shot shows, `/assets/js/lightbox.js` was accessed 297 times.

The job summary can be viewed at:

http://159.89.43.89:18080/history/application_1521391514647_0009/1/jobs/

Problem 3

The `LogAnalysis3.java` file is the source code for problem 3.

The program can be launched by:

```
$SPARK_HOME/bin/spark-submit \
--deploy-mode cluster \
--class LogAnalysis3 \
mini-project-02-1.0.jar
```

The logs can be viewed by:

```
$HADOOP_HOME/bin/yarn logs -log_files stdout -applicationId application_1521391514647_0012
```

```
Isolachine — root@master: /usr/local/spark — ssh root@159.89.43.89 — 120x30
2018-03-18 18:02:17 INFO TaskSetManager:54 - Finished task 2.0 in stage 1.0 (TID 6) in 841 ms on master (executor 2) (4 /4)
2018-03-18 18:02:17 INFO YarnClusterScheduler:54 - Removed TaskSet 1.0, whose tasks have all completed, from pool
2018-03-18 18:02:17 INFO DAGScheduler:54 - ResultStage 1 (collect at LogAnalysis3.java:38) finished in 0.859 s
2018-03-18 18:02:17 INFO DAGScheduler:54 - Job 0 finished: collect at LogAnalysis3.java:38, took 10.062426 s
-----OUTPUT START-----
/assets/css/combined.css 117348
-----OUTPUT END-----
2018-03-18 18:02:17 INFO AbstractConnector:318 - Stopped Spark@76cid1e6{HTTP/1.1,[http/1.1]}{0.0.0.0:0}
2018-03-18 18:02:17 INFO SparkUI:54 - Stopped Spark web UI at http://slave:34324
2018-03-18 18:02:17 INFO YarnAllocator:54 - Driver requested a total number of 0 executor(s).
2018-03-18 18:02:17 INFO YarnClusterSchedulerBackend:54 - Shutting down all executors
2018-03-18 18:02:17 INFO YarnSchedulerBackend$YarnDriverEndpoint:54 - Asking each executor to shut down
2018-03-18 18:02:17 INFO SchedulerExtensionServices:54 - Stopping SchedulerExtensionServices
(serviceOption=None,
 services=List(),
 started=false)
2018-03-18 18:02:17 INFO MapOutputTrackerMasterEndpoint:54 - MapOutputTrackerMasterEndpoint stopped!
2018-03-18 18:02:18 INFO MemoryStore:54 - MemoryStore cleared
2018-03-18 18:02:18 INFO BlockManager:54 - BlockManager stopped
2018-03-18 18:02:18 INFO BlockManagerMaster:54 - BlockManagerMaster stopped
2018-03-18 18:02:18 INFO OutputCommitCoordinator$OutputCommitCoordinatorEndpoint:54 - OutputCommitCoordinator stopped!
2018-03-18 18:02:18 INFO SparkContext:54 - Successfully stopped SparkContext
-----RUNNING TIME START-----
Total running time: 21.309s
-----RUNNING TIME END-----
2018-03-18 18:02:18 INFO ApplicationMaster:54 - Final app status: SUCCEEDED, exitCode: 0
2018-03-18 18:02:18 INFO ApplicationMaster:54 - Unregistering ApplicationMaster with SUCCEEDED
2018-03-18 18:02:18 INFO AMRMClientImpl:382 - Waiting for application to be successfully unregistered.
2018-03-18 18:02:18 INFO ApplicationMaster:54 - Deleting staging directory hdfs://master:9000/user/root/.sparkStaging/a
```

Screenshot of Results and Running Time

As the screen shot shows, `/assets/css/combined.css` was the most accessed resource, with 117348 times.

The job summary can be viewed at:

http://159.89.43.89:18080/history/application_1521391514647_0012/1/jobs/

Problem 4

The `LogAnalysis4.java` file is the source code for problem 4.

The program can be launched by:

```
$SPARK_HOME/bin/spark-submit \
--deploy-mode cluster \
--class LogAnalysis4 \
mini-project-02-1.0.jar
```

The logs can be viewed by:

```
$HADOOP_HOME/bin/yarn logs -log_files stdout -applicationId application_1521391514647_0015
```

```
Isolachine — root@master: /usr/local/spark — ssh root@159.89.43.89 — 120x30
2018-03-18 18:09:54 INFO BlockManagerInfo:54 - Removed taskresult_4 on master:38119 in memory (size: 2.2 MB, free: 366.3 MB)
2018-03-18 18:09:54 INFO DAGScheduler:54 - ResultStage 1 (collect at LogAnalysis4.java:38) finished in 2.294 s
2018-03-18 18:09:54 INFO DAGScheduler:54 - Job 0 finished: collect at LogAnalysis4.java:38, took 11.927303 s
2018-03-18 18:09:54 INFO YarnClusterScheduler:54 - Removed TaskSet 1.0, whose tasks have all completed, from pool
-----OUTPUT START-----
10.216.113.172 158614
-----OUTPUT END-----
2018-03-18 18:09:54 INFO AbstractConnector:318 - Stopped Spark@6265dbb4{HTTP/1.1,[http/1.1]}{0.0.0.0:0}
2018-03-18 18:09:54 INFO SparkUI:54 - Stopped Spark web UI at http://slave:36963
2018-03-18 18:09:54 INFO YarnAllocator:54 - Driver requested a total number of 0 executor(s).
2018-03-18 18:09:54 INFO YarnClusterSchedulerBackend:54 - Shutting down all executors
2018-03-18 18:09:54 INFO YarnSchedulerBackend$YarnDriverEndpoint:54 - Asking each executor to shut down
2018-03-18 18:09:54 INFO SchedulerExtensionServices:54 - Stopping SchedulerExtensionServices
(serviceOption=None,
 services=List(),
 started=false)
2018-03-18 18:09:54 INFO MapOutputTrackerMasterEndpoint:54 - MapOutputTrackerMasterEndpoint stopped!
2018-03-18 18:09:54 INFO MemoryStore:54 - MemoryStore cleared
2018-03-18 18:09:54 INFO BlockManager:54 - BlockManager stopped
2018-03-18 18:09:54 INFO BlockManagerMaster:54 - BlockManagerMaster stopped
2018-03-18 18:09:54 INFO OutputCommitCoordinator$OutputCommitCoordinatorEndpoint:54 - OutputCommitCoordinator stopped!
2018-03-18 18:09:54 INFO SparkContext:54 - Successfully stopped SparkContext
-----RUNNING TIME START-----
Total running time: 22.58s
-----RUNNING TIME END-----
2018-03-18 18:09:54 INFO ApplicationMaster:54 - Final app status: SUCCEEDED, exitCode: 0
2018-03-18 18:09:54 INFO ApplicationMaster:54 - Unregistering ApplicationMaster with SUCCEEDED
2018-03-18 18:09:54 INFO AMRMClientImpl:382 - Waiting for application to be successfully unregistered.
2018-03-18 18:09:55 INFO ApplicationMaster:54 - Deleting staging directory hdfs://master:9000/user/root/.sparkStaging/a
```

Screenshot of Results and Running Time

As the screen shot shows, 10.216.113.172 was the IP which accessed the website most frequently, with 158614 times.

The job summary can be viewed at:

```
http://159.89.43.89:18080/history/application_1521391514647_0015/1/jobs/
```

RDD Cache Used

The LogAnalysisCache.java file is the source code for problem with RDD cache.

The program can be launched by:

```
$SPARK_HOME/bin/spark-submit \
--deploy-mode cluster \
--class LogAnalysisCache \
mini-project-02-1.0.jar
```

The logs can be viewed by:

```
$HADOOP_HOME/bin/yarn logs -log_files stdout -applicationId application_1521391514647_0017
```



```
Isolachine — root@master: /usr/local/spark — ssh root@159.89.43.89 — 120x30
2018-03-18 18:19:16 INFO TaskSetManager:54 - Finished task 0.0 in stage 1.0 (TID 4) in 733 ms on master (executor 2) (4 /4)
2018-03-18 18:19:16 INFO YarnClusterScheduler:54 - Removed TaskSet 1.0, whose tasks have all completed, from pool
2018-03-18 18:19:16 INFO DAGScheduler:54 - ResultStage 1 (collectAsMap at LogAnalysisCache.java:38) finished in 0.752 s
2018-03-18 18:19:16 INFO DAGScheduler:54 - Job 0 finished: collectAsMap at LogAnalysisCache.java:38, took 10.372129 s
-----OUTPUT START-----
/assets/img/loading.gif 294
/assets/js/lightbox.js 297
-----OUTPUT END-----
2018-03-18 18:19:16 INFO AbstractConnector:318 - Stopped Spark@76c1d1e6(HTTP/1.1,[http/1.1]){0.0.0.0:0}
2018-03-18 18:19:16 INFO SparkUI:54 - Stopped Spark web UI at http://slave:34657
2018-03-18 18:19:16 INFO YarnAllocator:54 - Driver requested a total number of 0 executor(s).
2018-03-18 18:19:16 INFO YarnClusterSchedulerBackend:54 - Shutting down all executors
2018-03-18 18:19:16 INFO YarnSchedulerBackend$YarnDriverEndpoint:54 - Asking each executor to shut down
2018-03-18 18:19:16 INFO SchedulerExtensionServices:54 - Stopping SchedulerExtensionServices
(serviceOption=None,
services=List(),
started=false)
2018-03-18 18:19:16 INFO MapOutputTrackerMasterEndpoint:54 - MapOutputTrackerMasterEndpoint stopped!
2018-03-18 18:19:16 INFO MemoryStore:54 - MemoryStore cleared
2018-03-18 18:19:16 INFO BlockManager:54 - BlockManager stopped
2018-03-18 18:19:16 INFO BlockManagerMaster:54 - BlockManagerMaster stopped
2018-03-18 18:19:16 INFO OutputCommitCoordinator$OutputCommitCoordinatorEndpoint:54 - OutputCommitCoordinator stopped!
2018-03-18 18:19:17 INFO SparkContext:54 - Successfully stopped SparkContext
-----RUNNING TIME START-----
Total running time: 21.055s
-----RUNNING TIME END-----
2018-03-18 18:19:17 INFO ApplicationMaster:54 - Final app status: SUCCEEDED, exitCode: 0
2018-03-18 18:19:17 INFO ApplicationMaster:54 - Unregistering ApplicationMaster with SUCCEEDED
2018-03-18 18:19:17 INFO AMRMClientImpl:382 - Waiting for application to be successfully unregistered.
```

Screenshot of Results and Running Time

As the screen shot shows, when using RDD cache as intermediate results for further computation, the running time was 21.055 seconds.

The job summary can be viewed at:

```
http://159.89.43.89:18080/history/application_1521391514647_0017/1/jobs/
```

No Cache Used

The `LogAnalysisNoCache.java` file is the source code for problem with no cache used.

The program can be launched by:

```
$SPARK_HOME/bin/spark-submit \
--deploy-mode cluster \
--class LogAnalysisNoCache \
mini-project-02-1.0.jar
```

The logs can be viewed by:

```
$SHADOOP_HOME/bin/yarn logs -log_files stdout -applicationId application_1521391514647_0028
```

```
Isolachine — root@master: /usr/local/spark — ssh root@159.89.43.89 — 120x30
2018-03-18 18:47:20 INFO DAGScheduler:54 - ResultStage 3 (collectAsMap at LogAnalysisNoCache.java:52) finished in 0.327 s
2018-03-18 18:47:20 INFO DAGScheduler:54 - Job 1 finished: collectAsMap at LogAnalysisNoCache.java:52, took 5.921248 s
2018-03-18 18:47:20 INFO YarnClusterScheduler:54 - Removed TaskSet 3.0, whose tasks have all completed, from pool
-----OUTPUT START-----
/assets/img/loading.gif 294
/assets/js/lightbox.js 297
-----OUTPUT END-----
2018-03-18 18:47:20 INFO AbstractConnector:318 - Stopped Spark@3dffa792{HTTP/1.1,[http/1.1]}{0.0.0.0:0}
2018-03-18 18:47:20 INFO SparkUI:54 - Stopped Spark web UI at http://slave:41834
2018-03-18 18:47:20 INFO YarnAllocator:54 - Driver requested a total number of 0 executor(s).
2018-03-18 18:47:20 INFO YarnClusterSchedulerBackend:54 - Shutting down all executors
2018-03-18 18:47:20 INFO YarnSchedulerBackend$YarnDriverEndpoint:54 - Asking each executor to shut down
2018-03-18 18:47:20 INFO SchedulerExtensionServices:54 - Stopping SchedulerExtensionServices
(serviceOption=None,
 services=List(),
 started=false)
2018-03-18 18:47:20 INFO MapOutputTrackerMasterEndpoint:54 - MapOutputTrackerMasterEndpoint stopped!
2018-03-18 18:47:20 INFO MemoryStore:54 - MemoryStore cleared
2018-03-18 18:47:20 INFO BlockManager:54 - BlockManager stopped
2018-03-18 18:47:20 INFO BlockManagerMaster:54 - BlockManagerMaster stopped
2018-03-18 18:47:20 INFO OutputCommitCoordinator$OutputCommitCoordinatorEndpoint:54 - OutputCommitCoordinator stopped!
2018-03-18 18:47:20 INFO SparkContext:54 - Successfully stopped SparkContext
-----RUNNING TIME START-----
Total running time: 28.567s
-----RUNNING TIME END-----
2018-03-18 18:47:20 INFO ApplicationMaster:54 - Final app status: SUCCEEDED, exitCode: 0
2018-03-18 18:47:20 INFO ApplicationMaster:54 - Unregistering ApplicationMaster with SUCCEEDED
2018-03-18 18:47:20 INFO AMRMClientImpl:382 - Waiting for application to be successfully unregistered.
2018-03-18 18:47:20 INFO ApplicationMaster:54 - Deleting staging directory hdfs://master:9000/user/root/.sparkStaging/a
```

Screenshot of Results and Running Time

As the screen shot shows, when using RDD cache as intermediate results for further computation, the running time was 28.567 seconds.

The job summary can be viewed at:

http://159.89.43.89:18080/history/application_1521391514647_0028/1/jobs/

RDD Cache Efficiency Analysis

From the previous sections, we were able to get the total running time of the program generating access counts on two different resource links. Both the mechanism of running with RDD cache and running without RDD cache was used.

The program with RDD cache was able to finish with 21.055 seconds, while the one without cache took 28.567 seconds. Caching made the whole program 25% faster than on demand loading.

We were able to see how time was consumed on each sub-task from the job summary web interface.

← → ↻

159.89.43.89:18080/history/application_1521391514647_0017/1/jobs/

☆ ⋮

APACHE

Spark

2.3.0

Log Analysis Cache application UI

Jobs

Stages

Storage

Environment

Executors

Spark Jobs (?)

User: root
Total Uptime: 21 s
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	collectAsMap at LogAnalysisCache.java:38 collectAsMap at LogAnalysisCache.java:38	2018/03/18 18:19:06	10 s	2/2	8/8

Job Summary of Program with RDD Cache

Despite the total running time of 21 seconds, the actual text processing time was about 10 seconds for caching the file and output the desired results.

← → ↻

159.89.43.89:18080/history/application_1521391514647_0028/1/jobs/

☆ ⋮

APACHE

Spark

2.3.0

Log Analysis No Cache application UI

Jobs

Stages

Storage

Environment

Executors

Spark Jobs (?)

User: root
Total Uptime: 28 s
Scheduling Mode: FIFO
Completed Jobs: 2

▶ Event Timeline

Completed Jobs (2)

Job Id ▼	Description	Submitted	Duration	Stages: Succeeded/Total	Tasks (for all stages): Succeeded/Total
1	collectAsMap at LogAnalysisNoCache.java:52 collectAsMap at LogAnalysisNoCache.java:52	2018/03/18 18:47:14	6 s	2/2	8/8
0	collectAsMap at LogAnalysisNoCache.java:37 collectAsMap at LogAnalysisNoCache.java:37	2018/03/18 18:47:02	11 s	2/2	8/8

Job Summary of Program with No Cache Used

Despite the total running time of 28 seconds, the actual text processing time was about 11 seconds and 6 seconds for the two separate process.

This comparison shows caching has great advantage when dealing with data that can be used repeatedly. Although the second loading time is shorter than the first one (maybe due to loading mechanism of HDFS?), it could still save a lot of time with cache on a repetition of two jobs using the same intermediate result.

If the task was to do ad hoc search like database queries rather than ones we did above, caching would further save much more time with a load once, use many mechanism. This can be easily achieved by using the Spark Shell and Scala language to stream tasks.