

## Mini Project 2

### Group:

Sweta Rawal (swr22@pitt.edu)

Tazeen Khan (tfk12@pitt.edu)

Neha Shah (nes95@pitt.edu)

The zip contains two folders

### Program 1:

Artist.java (jar is /home/ubuntu/target folder in cc-project-13)

### Program 2:

p2.py and dependencies3.zip for running program inside the vm cluster

Folder jupyternotebook contains spark program for running it on jupyter along with save model and html of the output

### Part 1: Setting up spark

Installed version 3.0.2 in instance cc-project-13

To check enter command: ./hadoop/spark/bin/spark-submit --version

```
*** System restart required ***
Last login: Tue Mar 30 18:17:18 2021 from 158.212.127.39
ubuntu@cc-project-13:~$ ./hadoop/spark/bin/spark-submit --version
Welcome to

    _/ \
   / \ \_ \_ \_ \_ \_ \_ \_ \
  / \ \_ \_ \_ \_ \_ \_ \_ \_ \
 / \ \_ \_ \_ \_ \_ \_ \_ \_ \
/ \ \_ \_ \_ \_ \_ \_ \_ \_ \
 \_ \_ \_ \_ \_ \_ \_ \_ \
                           version 3.0.2

Using Scala version 2.12.10, OpenJDK 64-Bit Server VM, 1.8.0_282
Branch HEAD
Compiled by user centos on 2021-02-16T06:09:22Z
Revision 648457905c4ea7d00e3d88048c63f360045f0714
Url https://gitbox.apache.org/repos/asf/spark.git
Type --help for more information.
ubuntu@cc-project-13:~$
```

### Part 2: Developing Spark programs

1. To run the program use below command (Since the output is very big I have the command to store the output in consoleoutfile.txt file.)-

```
./hadoop/spark/bin/spark-submit --class artists.Artists --master yarn target/artists-1.0-SNAPSHOT.jar > /home/ubuntu/consoleoutfile.txt 2>&1
```

2. To see the output on the command line run below-

```
./hadoop/spark/bin/spark-submit --class artists.Artists --master yarn target/artists-1.0-SNAPSHOT.jar
```

```
To see these additional updates run: apt list --upgradable
```

```
*** System restart required ***
Last login: Tue Mar 30 18:17:22 2021 from 150.212.127.39
ubuntu@cc-project-13:~$ ./hadoop/spark/bin/spark-submit --class artists.Artists --master yarn target/artists-1.0-SNAPSHOT.jar
```

```
- Reduce system reboots and improve kernel security. Activate at:
  https://ubuntu.com/livepatch

24 packages can be updated.
0 of these updates are security updates.
To see these additional updates run: apt list --upgradable

*** System restart required ***
Last login: Tue Mar 30 18:17:22 2021 from 150.212.127.39
ubuntu@cc-project-13:~$ ./hadoop/spark/bin/spark-submit --class artists.Artists --master yarn target/artists-1.0-SNAPSHOT.jar
2021-03-30 19:47:02,130 WARN util.NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable
2021-03-30 19:47:02,357 INFO spark.SparkContext: Running Spark version 3.0.2
2021-03-30 19:47:02,446 INFO resource.ResourceUtils: ======
2021-03-30 19:47:02,448 INFO resource.ResourceUtils: ======
2021-03-30 19:47:02,449 INFO spark.SparkContext: Submitted application: total_count
2021-03-30 19:47:02,512 INFO spark.SecurityManager: Changing view acls to: ubuntu
2021-03-30 19:47:02,512 INFO spark.SecurityManager: Changing modify acls to: ubuntu
2021-03-30 19:47:02,512 INFO spark.SecurityManager: Changing view acls groups to:
2021-03-30 19:47:02,512 INFO spark.SecurityManager: Changing modify acls groups to:
2021-03-30 19:47:02,512 INFO spark.SecurityManager: SecurityManager: authentication disabled; ui acls disabled; users with view permissions: Set(ubuntu); groups with view permissions: Set()
2021-03-30 19:47:02,752 INFO util.Utils: Successfully started service 'sparkDriver' on port 34379.
2021-03-30 19:47:02,784 INFO spark.SparkEnv: Registering MapOutputTracker
2021-03-30 19:47:02,817 INFO spark.SparkEnv: Registering BlockManagerMaster
2021-03-30 19:47:02,835 INFO storage.BlockManagerMasterEndpoint: Using org.apache.spark.storage.DefaultTopologyMapper for getting topology information
2021-03-30 19:47:02,836 INFO storage.BlockManagerMasterEndpoint: BlockManagerMasterEndpoint up
2021-03-30 19:47:02,872 INFO spark.SparkEnv: Registering BlockManagerMasterHeartbeat
2021-03-30 19:47:02,885 INFO storage.DiskBlockManager: Created local directory at /tmp/blockmgr-f932f2d0-5fb9-47db-bb52-77acbb8be5c1
2021-03-30 19:47:02,916 INFO memory.MemoryStore: MemoryStore started with capacity 366.3 MiB
2021-03-30 19:47:02,963 INFO spark.SparkEnv: Registering OutputCommitCoordinator
```

```
2021-03-30 19:47:40,929 INFO scheduler.TaskSetManager: Finished task 70.0 in stage 18.0 (TID 599) in 10 ms on cc-project-14.cc-s21-project.is2750-pg0.utah.cloudlab.us (executed by executor 2, part 0, 7336 bytes)
2021-03-30 19:47:40,938 INFO scheduler.TaskSetManager: Starting task 72.0 in stage 18.0 (TID 601, cc-project-14.cc-s21-project.is2750-pg0.utah.cloudlab.us, executor 2, part 1, 7336 bytes)
2021-03-30 19:47:40,946 INFO scheduler.TaskSetManager: Finished task 71.0 in stage 18.0 (TID 600) in 8 ms on cc-project-14.cc-s21-project.is2750-pg0.utah.cloudlab.us (executed by executor 2, part 0, 7336 bytes)
2021-03-30 19:47:40,946 INFO cluster.YarnScheduler: Removed TaskSet 18.0, whose tasks have all completed, from pool
2021-03-30 19:47:40,947 INFO scheduler.DAGScheduler: ResultStage 18 (show at Artists.java:21) finished in 0.660 s
2021-03-30 19:47:40,947 INFO scheduler.DAGScheduler: Job 7 is finished. Cancelling potential speculative or zombie tasks for this job
2021-03-30 19:47:40,947 INFO cluster.YarnScheduler: Killing all running tasks in stage 18: Stage finished
2021-03-30 19:47:40,948 INFO scheduler.DAGScheduler: Job 7 finished: show at Artists.java:21, took 0.673690 s
2021-03-30 19:47:40,980 INFO codegen.CodeGenerator: Code generated in 14.877918 ms
+-----+
| artistID | sum(weight) |
+-----+-----+
| 289 | 2393140 |
| 72 | 1301308 |
| 89 | 1291387 |
| 292 | 1058485 |
| 498 | 963449 |
| 67 | 921198 |
| 288 | 985423 |
| 701 | 688529 |
| 227 | 662116 |
| 300 | 532545 |
| 333 | 525844 |
| 344 | 525292 |
| 378 | 513476 |
| 679 | 506453 |
| 295 | 499318 |
| 511 | 493024 |
| 461 | 489665 |
| 486 | 485532 |
| 190 | 485076 |
| 163 | 466104 |
| 55 | 449292 |
| 154 | 385306 |
| 466 | 384405 |
| 257 | 384207 |
| 707 | 371916 |
| 917 | 368710 |
| 792 | 350635 |
| 51 | 348919 |
| 65 | 330757 |
| 475 | 321011 |
| 203 | 318221 |
| 157 | 296882 |
| 207 | 288520 |
| 198 | 277397 |
| 377 | 265362 |
| 291 | 253027 |
| 614 | 251440 |
| 173 | 245878 |
| 503 | 237148 |
| 687 | 215777 |
| 903 | 213103 |
| 302 | 207761 |
| 187 | 205195 |
| 1412 | 203665 |
| 1008 | 1902178 |
```

```

[13822 |1
|7630 |1
|13762 |1
|14374 |1
|15749 |1
|1173 |1
|9581 |1
|11751 |1
|14369 |1
|4551 |1
|9584 |1
|17168 |1
|7631 |1
|17272 |1
|9586 |1
|14372 |1
|17169 |1
|17468 |1
|8052 |1
|9579 |1
|13818 |1
|17269 |1
|16804 |1
|11746 |1
|13823 |1
|2897 |1
|9490 |1
|15943 |1
|7752 |1
|13760 |1
|13820 |1
|13828 |1
|18432 |1
|15937 |1
|12402 |1
|17273 |1
|17266 |1
|7627 |1
|13817 |1
+-----+
2021-03-30 19:47:41.192 INFO spark.SparkContext: Invoking stop() from shutdown hook
2021-03-30 19:47:41.200 INFO server.AbstractConnector: Stopped Spark@28c4205{HTTP/1.1, {<http://1.1>}}{0.0.0.0:4040}
2021-03-30 19:47:41.202 INFO ui.SparkUI: Stopped Spark web UI at http://cc-project-13.cc-s21-project.1s2750-pg0.utah.cloudlab.us:4040
2021-03-30 19:47:41.207 INFO cluster.YarnClientSchedulerBackend: Interrupting monitor thread
2021-03-30 19:47:41.229 INFO cluster.YarnClientSchedulerBackend: Shutting down all executors
2021-03-30 19:47:41.229 INFO cluster.YarnSchedulerBackend$YarnDriverEndpoint: Asking each executor to shut down
2021-03-30 19:47:41.231 INFO cluster.YarnClientSchedulerBackend: YARN client scheduler backend Stopped
2021-03-30 19:47:41.252 INFO spark.MapOutputTrackerMasterEndpoint: MapOutputTrackerMasterEndpoint stopped!
2021-03-30 19:47:41.261 INFO memory.MemoryStore: MemoryStore cleared
2021-03-30 19:47:41.265 INFO storage.BlockManager: BlockManager stopped
2021-03-30 19:47:41.273 INFO storage.BlockManagerMaster: BlockManagerMaster stopped
2021-03-30 19:47:41.276 INFO scheduler.OutputCommitCoordinator$OutputCommitCoordinatorEndpoint: OutputCommitCoordinator stopped!
2021-03-30 19:47:41.287 INFO util.ShutdownHookManager: Successfully stopped SparkContext
2021-03-30 19:47:41.287 INFO util.ShutdownHookManager: Shutdown hook called
2021-03-30 19:47:41.288 INFO util.ShutdownHookManager: Deleting directory /tmp/spark-288e2798-8785-4251-adc9-12a2641974c6
2021-03-30 19:47:41.292 INFO util.ShutdownHookManager: Deleting directory /tmp/spark-163e428d-fa87-4ed5-99d4-d6be4fe0fdf10
ubuntu@cc-project-13:~$ 

```

To see full output check the pre generated output in consoleoutfile.txt or run again.

### Part 3: Developing Spark programs 2

- Access\_Log file is placed at hdfs://user/root/accesslog/access\_log

(Path is hardcoded into the p2.py file, if path is changed please update in line 17)

```

ubuntu@cc-project-13:~$ hdfs dfs -ls /user/root/accesslog
Found 3 items
-rw-r--r--   1      supergroup  504941532 2021-03-28 20:39 /user/root/accesslog/access_log
drwxr-xr-x  1000  Termius      supergroup          0 2021-03-28 20:39 /user/root/accesslog/home
-rw-r--r--   1      supergroup  1296455 2021-03-24 16:21 /user/root/accesslog/user_artists.dat
ubuntu@cc-project-13:~$ 

```

- To run the program external dependencies3.zip is required which is placed at /home/ubuntu
- Execute the following command to run the program:
  - Cd Hadoop
  - Cd spark
  - ./bin/spark-submit --py-files /home/ubuntu/dependencies3.zip /home/ubuntu/p2.py

```

ubuntu@cc-project-13:~/hadoop/spark$ ./bin/spark-submit --py-files /home/ubuntu/dependencies3.zip /home/ubuntu/p2.py
2021-03-28 20:42:50,037 WARN util.NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable
2021-03-28 20:42:50,037 INFO util.NativeCodeLoader: Using builtin-java classes for Hadoop library
2021-03-28 20:42:50,084 INFO spark.SparkContext: Running Spark version 3.0.2
2021-03-28 20:42:50,084 INFO resource.ResourceUtils: ****
2021-03-28 20:42:50,084 INFO resource.ResourceUtils: Resources for spark.driver:
2021-03-28 20:42:50,084 INFO resource.ResourceUtils: ****
2021-03-28 20:42:50,084 INFO spark.SparkContext: Submitted application: p2.py
2021-03-28 20:42:50,090 INFO spark.SecurityManager: Changing view acls to: ubuntu
2021-03-28 20:42:50,090 INFO spark.SecurityManager: Changing modify acls to: ubuntu
2021-03-28 20:42:50,090 INFO spark.SecurityManager: Changing view acls groups to:
2021-03-28 20:42:50,090 INFO spark.SecurityManager: Changing modify acls groups to:
2021-03-28 20:42:50,090 INFO spark.SecurityManager: SecurityManagers authentication disabled; ui acls disabled; users with view permissions: Set(ubuntu);
groups with view permissions: Set(); users with modify permissions: Set(ubuntu); groups with modify permissions: Set()
2021-03-28 20:42:51,141 INFO util.Utils: Successfully started service 'sparkDriver' on port 38441.
2021-03-28 20:42:51,176 INFO spark.SparkContext: Registering RDD 0
2021-03-28 20:42:51,200 INFO storage.BlockManagerMasterEndpoint: Registered BlockManagerMaster
2021-03-28 20:42:51,229 INFO storage.BlockManagerMasterEndpoint: Using org.apache.spark.storage.DefaultTopologyMapper for getting topology information
2021-03-28 20:42:51,230 INFO storage.BlockManagerMasterEndpoint: BlockManagerMasterEndpoint up
2021-03-28 20:42:51,233 INFO spark.SparkEnv: Registering BlockManagerMasterHeartbeat
2021-03-28 20:42:51,243 INFO storage.DiskBlockManager: Created local directory at /tmp/blockmgr-508f310f-5d41-4520-8208-8022585c331e
2021-03-28 20:42:51,274 INFO memory.MemoryStore: MemoryStore started with capacity 368.3 M
2021-03-28 20:42:51,300 INFO spark.SparkEnv: Registering OutputCommitCoordinator
2021-03-28 20:42:51,376 INFO util.ILogging: Logging initialized @2752ms to org.apache.spark.storage.DefaultTopologyMapper
2021-03-28 20:42:51,448 INFO server.Server: Server started on jetty-9.4.24.v20201020; built: 2020-11-02T14:15:29.302Z; git: e46af88704a893fc12cb0e3bf46e2c7b48a099e7; jvm: 1.8
._.282-8u28-b08-ubuntu18.04-04008
2021-03-28 20:42:51,471 INFO server.Server: Started @2844ms
2021-03-28 20:42:51,512 INFO server.AbstractConnector: Started ServerConnector@25052dbf[HTTP/1.1, ((http://1.0.0.0:4040)]
2021-03-28 20:42:51,515 INFO util.Utils: Successfully started service 'SparkUI' on port 4040
2021-03-28 20:42:51,516 INFO util.Utils: Registered application sparkDriver
2021-03-28 20:42:51,549 INFO handler.ContextHandler: Started o.s.j.s.ServletContextHandler@133f96c/[jobs/json,null,AVAILABLE,@Spark]
2021-03-28 20:42:51,550 INFO handler.ContextHandler: Started o.s.j.s.ServletContextHandler@6ac8c3f86ff/[jobs/null,AVAILABLE,@Spark]
2021-03-28 20:42:51,553 INFO handler.ContextHandler: Started o.s.j.s.ServletContextHandler@3934903bf/[jobs/job/json,null,AVAILABLE,@Spark]
2021-03-28 20:42:51,555 INFO handler.ContextHandler: Started o.s.j.s.ServletContextHandler@b5ee2ff8/[stages/null,AVAILABLE,@Spark]
2021-03-28 20:42:51,558 INFO handler.ContextHandler: Started o.s.j.s.ServletContextHandler@42229377/[stages/json,null,AVAILABLE,@Spark]
2021-03-28 20:42:51,559 INFO handler.ContextHandler: Started o.s.j.s.ServletContextHandler@9764d7fb/[stage/age/null,AVAILABLE,@Spark]
2021-03-28 20:42:51,562 INFO handler.ContextHandler: Started o.s.j.s.ServletContextHandler@2933d61a/[stages/null,AVAILABLE,@Spark]
2021-03-28 20:42:51,564 INFO handler.ContextHandler: Started o.s.j.s.ServletContextHandler@2933d61a/[stages/json,null,AVAILABLE,@Spark]
2021-03-28 20:42:51,566 INFO handler.ContextHandler: Started o.s.j.s.ServletContextHandler@487d422ef/[stages/pool/json,null,AVAILABLE,@Spark]
2021-03-28 20:42:51,567 INFO handler.ContextHandler: Started o.s.j.s.ServletContextHandler@8d66e010f/[storage/null,AVAILABLE,@Spark]
2021-03-28 20:42:51,568 INFO handler.ContextHandler: Started o.s.j.s.ServletContextHandler@73ee629f/[storage/json,null,AVAILABLE,@Spark]
2021-03-28 20:42:51,570 INFO handler.ContextHandler: Started o.s.j.s.ServletContextHandler@43c434a8f/[storage/rd,null,AVAILABLE,@Spark]
2021-03-28 20:42:51,571 INFO handler.ContextHandler: Started o.s.j.s.ServletContextHandler@242104f9/[storage/threadDump/json,null,AVAILABLE,@Spark]
2021-03-28 20:42:51,572 INFO handler.ContextHandler: Started o.s.j.s.ServletContextHandler@2c2f17f/[storage/threadDump/json,null,AVAILABLE,@Spark]
2021-03-28 20:42:51,574 INFO handler.ContextHandler: Started o.s.j.s.ServletContextHandler@9541061bf/[environments/json,null,AVAILABLE,@Spark]
2021-03-28 20:42:51,574 INFO handler.ContextHandler: Started o.s.j.s.ServletContextHandler@2d3e38ccf/[executors/null,AVAILABLE,@Spark]
2021-03-28 20:42:51,575 INFO handler.ContextHandler: Started o.s.j.s.ServletContextHandler@2967799e/[executors/threadDump/json,null,AVAILABLE,@Spark]
2021-03-28 20:42:51,576 INFO handler.ContextHandler: Started o.s.j.s.ServletContextHandler@2180bf8/[executors/threadDump/json,null,AVAILABLE,@Spark]
2021-03-28 20:42:51,577 INFO handler.ContextHandler: Started o.s.j.s.ServletContextHandler@5e8dd195f/[executors/threadDump/json,null,AVAILABLE,@Spark]
2021-03-28 20:42:51,578 INFO handler.ContextHandler: Started o.s.j.s.ServletContextHandler@8570ffdb9/[api/json,null,AVAILABLE,@Spark]
2021-03-28 20:42:51,580 INFO handler.ContextHandler: Started o.s.j.s.ServletContextHandler@2acc8cd7f/[api/json,null,AVAILABLE,@Spark]
2021-03-28 20:42:51,587 INFO handler.ContextHandler: Started o.s.j.s.ServletContextHandler@2719728e/[jobs/job/kill,null,AVAILABLE,@Spark]
2021-03-28 20:42:51,590 INFO handler.ContextHandler: Started o.s.j.s.ServletContextHandler@7f7313e4f/[stages/stage/kill/null,AVAILABLE,@Spark]
2021-03-28 20:42:51,593 INFO handler.ContextHandler: Started o.s.j.s.ServletContextHandler@77333553f/[stages/stage/kill/null,AVAILABLE,@Spark]

```

```

2021-03-28 20:43:01,130 INFO codegen.CodeGenerator: Code generated in 16.357893 ms
2021-03-28 20:43:01,153 INFO executor.Executor: Finished task 0.0 in stage 2.0 (TID 5). 2403 bytes result sent to driver
2021-03-28 20:43:01,154 INFO scheduler.TaskSetManager: Finished task 0.0 in stage 2.0 (TID 5) in 60 ms on cc-project-13.cc-s21-project.is2750-pg0.utah.cloudlab.us (executor driver) (1/1)
2021-03-28 20:43:01,154 INFO scheduler.TaskSchedulerImpl: Removed TaskSet 2.0, whose tasks have all completed, from pool
2021-03-28 20:43:01,155 INFO scheduler.DAGScheduler: ResultStage 2 (showString at NativeMethodAccessorImpl.java:0) finished in 0.078 s
2021-03-28 20:43:01,155 INFO scheduler.DAGScheduler: Job 2 is finished. Cancelling potential speculative or zombie tasks for this job
2021-03-28 20:43:01,155 INFO scheduler.TaskSchedulerImpl: Killing all running tasks in stage 2: Stage finished
2021-03-28 20:43:01,155 INFO scheduler.DAGScheduler: Job 2 finished: showString at NativeMethodAccessorImpl.java:0, took 0.082848 s
2021-03-28 20:43:01,184 INFO codegen.CodeGenerator: Code generated in 20.345887 ms
+-----+
| _c0|_c1|_c2|      _c3| _c4|      _c5|_c6| _c7|
+-----+
| 10.223.157.186| -|[15/Jul/2009:14:5...| -|0700|| GET / HTTP/1.1|403| 202|
| 10.223.157.186| -|[15/Jul/2009:14:5...| -|0700|| GET /favicon.ico ...|404| 209|
| 10.223.157.186| -|[15/Jul/2009:15:5...| -|0700|| GET / HTTP/1.1|200| 9157|
| 10.223.157.186| -|[15/Jul/2009:15:5...| -|0700|| GET /assets/js/lo...|200| 10469|
| 10.223.157.186| -|[15/Jul/2009:15:5...| -|0700|| GET /assets/css/r...|200| 1014|
| 10.223.157.186| -|[15/Jul/2009:15:5...| -|0700|| GET /assets/css/9...|200| 6206|
| 10.223.157.186| -|[15/Jul/2009:15:5...| -|0700|| GET /assets/css/t...|200| 15779|
| 10.223.157.186| -|[15/Jul/2009:15:5...| -|0700|| GET /assets/js/th...|200| 4492|
| 10.223.157.186| -|[15/Jul/2009:15:5...| -|0700|| GET /assets/js/l...|200| 25960|
| 10.223.157.186| -|[15/Jul/2009:15:5...| -|0700|| GET /assets/img/s...|200| 168|
| 10.223.157.186| -|[15/Jul/2009:15:5...| -|0700|| GET /assets/img/d...|200| 5604|
| 10.223.157.186| -|[15/Jul/2009:15:5...| -|0700|| GET /assets/img/d...|200| 10556|
| 10.223.157.186| -|[15/Jul/2009:15:5...| -|0700|| GET /assets/img/d...|200| 9925|
| 10.223.157.186| -|[15/Jul/2009:15:5...| -|0700|| GET /assets/img/c...|200| 979|
| 10.223.157.186| -|[15/Jul/2009:15:5...| -|0700|| GET /assets/img/h...|200| 3892|
| 10.223.157.186| -|[15/Jul/2009:15:5...| -|0700|| GET /assets/img/d...|200| 5397|
| 10.223.157.186| -|[15/Jul/2009:15:5...| -|0700|| GET /assets/img/l...|200| 2767|
| 10.223.157.186| -|[15/Jul/2009:15:5...| -|0700|| GET /assets/img/d...|200| 5766|
| 10.223.157.186| -|[15/Jul/2009:15:5...| -|0700|| GET /assets/img/h...|200| 68831|
| 10.223.157.186| -|[15/Jul/2009:15:5...| -|0700|| GET /assets/img/d...|200| 5766|
+-----+
only showing top 20 rows

2021-03-28 20:43:01,249 INFO datasources.FileSourceStrategy: Pushed Filters:
2021-03-28 20:43:01,249 INFO datasources.FileSourceStrategy: Post-Scan Filters:
2021-03-28 20:43:01,249 INFO datasources.FileSourceStrategy: Output Data Schema: struct<_c0: string, _c3: string>
2021-03-28 20:43:01,285 INFO codegen.CodeGenerator: Code generated in 17.556135 ms
2021-03-28 20:43:01,290 INFO memory.MemoryStore: Block broadcast_0 stored as values in memory (estimated size 307.7 Kib, free 365.0 MiB)
2021-03-28 20:43:01,299 INFO memory.MemoryStore: Block broadcast_6_piece0 stored as bytes in memory (estimated size 27.9 Kib, free 364.9 MiB)
2021-03-28 20:43:01,300 INFO storage.BlockManagerInfo: Added broadcast_6_piece0 in memory on cc-project-13.cc-s21-project.is2750-pg0.utah.cloudlab.us:38469
(size: 27.9 Kib, free: 366.2 MiB)
2021-03-28 20:43:01,301 INFO spark.SparkContext: Created broadcast 6 from showString at NativeMethodAccessorImpl.java:0
2021-03-28 20:43:01,302 INFO execution.FileSourceScanExec: Planning scan with bin packing, max size: 127283959 bytes, open cost is considered as scanning 4 194304 bytes.
2021-03-28 20:43:01,310 INFO spark.SparkContext: Starting job: showString at NativeMethodAccessorImpl.java:0
2021-03-28 20:43:01,311 INFO scheduler.DAGScheduler: Got job 3 (showString at NativeMethodAccessorImpl.java:0) with 1 output partitions
2021-03-28 20:43:01,311 INFO scheduler.DAGScheduler: Final stage: ResultStage 3 (showString at NativeMethodAccessorImpl.java:0)
2021-03-28 20:43:01,311 INFO scheduler.DAGScheduler: Parents of final stage: List()
2021-03-28 20:43:01,312 INFO scheduler.DAGScheduler: Submitting ResultStage 3 (MapPartitionsRDD[16] at showString at NativeMethodAccessorImpl.java:0), which has no missing parents
2021-03-28 20:43:01,315 INFO memory.MemoryStore: Block broadcast_7 stored as values in memory (estimated size 13.6 Kib, free 364.9 MiB)
2021-03-28 20:43:01,317 INFO memory.MemoryStore: Block broadcast_7_piece0 stored as bytes in memory (estimated size 6.7 Kib, free 364.9 MiB)
2021-03-28 20:43:01,318 INFO storage.BlockManagerInfo: Added broadcast_7_piece0 in memory on cc-project-13.cc-s21-project.is2750-pg0.utah.cloudlab.us:38469
(size: 6.7 Kib, free: 366.2 MiB)
2021-03-28 20:43:01,318 INFO spark.SparkContext: Created broadcast 7 from broadcast at DAGScheduler.scala:1223

```

```
udlab.us (executor driver) (1/4)
2021-03-28 20:43:09,254 INFO executor.Executor: Finished task 3.0 in stage 8.0 (TID 16). 3764 bytes result sent to driver
2021-03-28 20:43:09,255 INFO executor.Executor: Finished task 1.0 in stage 8.0 (TID 14). 3754 bytes result sent to driver
2021-03-28 20:43:09,256 INFO scheduler.TaskSetManager: Finished task 3.0 in stage 8.0 (TID 16) in 30 ms on cc-project-13.cc-s21-project.is2750-pg0.utah.cloudlab.us (executor driver) (2/4)
2021-03-28 20:43:09,257 INFO scheduler.TaskSetManager: Finished task 1.0 in stage 8.0 (TID 14) in 32 ms on cc-project-13.cc-s21-project.is2750-pg0.utah.cloudlab.us (executor driver) (3/4)
2021-03-28 20:43:09,262 INFO executor.Executor: Finished task 0.0 in stage 8.0 (TID 13). 3785 bytes result sent to driver
2021-03-28 20:43:09,263 INFO scheduler.TaskSetManager: Finished task 0.0 in stage 8.0 (TID 13) in 39 ms on cc-project-13.cc-s21-project.is2750-pg0.utah.cloudlab.us (executor driver) (4/4)
2021-03-28 20:43:09,263 INFO scheduler.TaskSchedulerImpl: Removed TaskSet 8.0, whose tasks have all completed, from pool
2021-03-28 20:43:09,263 INFO scheduler.DAGScheduler: ResultStage 8 (showString at NativeMethodAccessorImpl.java:0) finished in 0.845 s
2021-03-28 20:43:09,264 INFO scheduler.DAGScheduler: Job 6 is finished. Cancelling potential speculative or zombie tasks for this job
2021-03-28 20:43:09,264 INFO scheduler.TaskSchedulerImpl: Killing all running tasks in stage 8: Stage finished
2021-03-28 20:43:09,265 INFO scheduler.DAGScheduler: Job 6 finished: showString at NativeMethodAccessorImpl.java:0, took 0.053355 s
2021-03-28 20:43:09,275 INFO codegen.CodeGenerator: Code generated in 7.431512 ms
+-----+
| year|count|
+-----+
| 0.1258848| 2851|
| 0.12814848| 7885|
| 0.1249776| 115|
| 0.12559104| 32|
| 0.12852864| 4203|
| 0.12962592| 3238|
| 0.12963456| 4868|
| 0.130464| 6501|
| 0.13179456| 9857|
| 0.12712896| 4564|
| 0.12819168| 4848|
| 0.12957408| 3536|
| 0.13099968| 8722|
| 0.13109472| 8993|
| 0.13138848| 5632|
| 0.13215744| 7051|
| 0.12794976| 7376|
| 0.12923712| 4842|
| 0.1266624| 2222|
| 0.12712032| 6862|
+-----+
only showing top 20 rows

2021-03-28 20:43:09,565 INFO datasources.FileSourceStrategy: Pushed Filters:
2021-03-28 20:43:09,565 INFO datasources.FileSourceStrategy: Post-Scan Filters:
2021-03-28 20:43:09,566 INFO datasources.FileSourceStrategy: Output Data Schema: struct<`c3: string>
2021-03-28 20:43:09,642 INFO codegen.CodeGenerator: Code generated in 38.393172 ms
2021-03-28 20:43:09,652 INFO memory.MemoryStore: Block broadcast_14 stored as values in memory (estimated size 307.7 KiB, free 364.8 MiB)
2021-03-28 20:43:09,660 INFO memory.MemoryStore: Block broadcast_14_piece0 stored as bytes in memory (estimated size 27.9 KiB, free 364.8 MiB)
2021-03-28 20:43:09,661 INFO storage.BlockManagerInfo: Added broadcast_14_piece0 in memory on cc-project-13.cc-s21-project.is2750-pg0.utah.cloudlab.us:38469 (size: 27.9 KiB, free: 366.1 MiB)
2021-03-28 20:43:09,662 INFO spark.SparkContext: Created broadcast 14 from showString at NativeMethodAccessorImpl.java:0
2021-03-28 20:43:09,663 INFO execution.FileSourceScanExec: Planning scan with bin packing, max size: 127283959 bytes, open cost is considered as scanning 194304 bytes.
2021-03-28 20:43:09,715 INFO spark.SparkContext: Starting job: showString at NativeMethodAccessorImpl.java:0
2021-03-28 20:43:09,720 INFO scheduler.DAGScheduler: Registering RDD 31 (showString at NativeMethodAccessorImpl.java:0) as input to shuffle 1
2021-03-28 20:43:09,721 INFO scheduler.DAGScheduler: Got job 7 (showString at NativeMethodAccessorImpl.java:0) with 1 output partitions
2021-03-28 20:43:09,721 INFO storage.BlockManagerInfo: Removed broadcast_12_piece0 on cc-project-13.cc-s21-project.is2750-pg0.utah.cloudlab.us:38469 in memory (size: 16.5 KiB, free: 366.1 MiB)
2021-03-28 20:43:09,721 INFO scheduler.DAGScheduler: Final stage: ResultStage 10 (showString at NativeMethodAccessorImpl.java:0)
```

```

2021-03-28 20:43:15,361 INFO executor.Executor: Finished task 1.0 in stage 12.0 (TID 23). 3824 bytes result sent to driver
2021-03-28 20:43:15,361 INFO executor.Executor: Finished task 3.0 in stage 12.0 (TID 25). 3840 bytes result sent to driver
2021-03-28 20:43:15,362 INFO executor.Executor: Finished task 2.0 in stage 12.0 (TID 24). 3692 bytes result sent to driver
2021-03-28 20:43:15,363 INFO scheduler.TaskSetManager: Finished task 1.0 in stage 12.0 (TID 23) in 23 ms on cc-project-13.cc-s21-project.is2750-pg0.utah.cl
oudlab.us (executor driver) (1/4)
2021-03-28 20:43:15,363 INFO scheduler.TaskSetManager: Finished task 3.0 in stage 12.0 (TID 25) in 23 ms on cc-project-13.cc-s21-project.is2750-pg0.utah.cl
oudlab.us (executor driver) (2/4)
2021-03-28 20:43:15,364 INFO scheduler.TaskSetManager: Finished task 2.0 in stage 12.0 (TID 24) in 23 ms on cc-project-13.cc-s21-project.is2750-pg0.utah.cl
oudlab.us (executor driver) (3/4)
2021-03-28 20:43:15,364 INFO executor.Executor: Finished task 0.0 in stage 12.0 (TID 22). 3856 bytes result sent to driver
2021-03-28 20:43:15,365 INFO scheduler.TaskSetManager: Finished task 0.0 in stage 12.0 (TID 22) in 26 ms on cc-project-13.cc-s21-project.is2750-pg0.utah.cl
oudlab.us (executor driver) (4/4)
2021-03-28 20:43:15,365 INFO scheduler.TaskSchedulerImpl: Removed TaskSet 12.0, whose tasks have all completed, from pool
2021-03-28 20:43:15,366 INFO scheduler: ResultStage 12 (showString at NativeMethodAccessorImpl.java:0) finished in 0.033 s
2021-03-28 20:43:15,366 INFO scheduler.DAGScheduler: Job 8 is finished. Cancelling potential speculative or zombie tasks for this job
2021-03-28 20:43:15,366 INFO scheduler.TaskSchedulerImpl: Killing all running tasks in stage 12: Stage finished
2021-03-28 20:43:15,366 INFO scheduler.DAGScheduler: Job 8 finished: showString at NativeMethodAccessorImpl.java:0, took 0.036463 s
2021-03-28 20:43:15,380 INFO codegen.CodeGenerator: Code generated in 11.157409 ms
+-----+
| year|count|      input|
+-----+
| 0.1258848| 2051|[0.1258848]|
| 0.12814848| 7885|[0.12814848]|
| 0.1249776| 115|[0.1249776]|
| 0.12559104| 32|[0.12559104]|
| 0.12852864| 4203|[0.12852864]|
| 0.12962592| 3238|[0.12962592]|
| 0.12963456| 4868|[0.12963456]|
| 0.130464| 6501|[0.130464]|
| 0.13179456| 9857|[0.13179456]|
| 0.12712896| 4564|[0.12712896]|
| 0.12819168| 4848|[0.12819168]|
| 0.12957408| 3536|[0.12957408]|
| 0.13099968| 8722|[0.13099968]|
| 0.13109472| 8993|[0.13109472]|
| 0.13138848| 5632|[0.13138848]|
| 0.13215744| 7051|[0.13215744]|
| 0.12794976| 7376|[0.12794976]|
| 0.12923712| 4842|[0.12923712]|
| 0.1266624| 2222|[0.1266624]|
| 0.12712032| 6862|[0.12712032]|
+-----+
only showing top 20 rows

2021-03-28 20:43:15,522 INFO datasources.FileSourceStrategy: Pushed Filters:
2021-03-28 20:43:15,523 INFO datasources.FileSourceStrategy: Post-Scan Filters:
2021-03-28 20:43:15,523 INFO datasources.FileSourceStrategy: Output Data Schema: struct<_c3: string>
2021-03-28 20:43:15,587 INFO codegen.CodeGenerator: Code generated in 40.303887 ms
2021-03-28 20:43:15,597 INFO memory.MemoryStore: Block broadcast_18 stored as values in memory (estimated size 307.7 KiB, free 364.4 MiB)
2021-03-28 20:43:15,605 INFO memory.MemoryStore: Block broadcast_18_piece0 stored as bytes in memory (estimated size 27.9 KiB, free 364.4 MiB)
2021-03-28 20:43:15,605 INFO storage.BlockManagerInfo: Added broadcast_18_piece0 in memory on cc-project-13.cc-s21-project.is2750-pg0.utah.cloudlab.us:3846
9 (size: 27.9 KiB, free: 366.1 MiB)
2021-03-28 20:43:15,607 INFO spark.SparkContext: Created broadcast 18 from rdd at Predictor.scala:80
2021-03-28 20:43:15,607 INFO execution.FileSourceScanExec: Planning scan with bin packing, max size: 127283959 bytes, open cost is considered as scanning 4
194304 bytes.
2021-03-28 20:43:15,626 INFO util.Instrumentation: [Slad6109] Stage class: LinearRegression
2021-03-28 20:43:15,627 INFO util.Instrumentation: [Slad6109] Stage uid: LinearRegression_8cc66242dc17
2021-03-28 20:43:15,651 INFO datasources.FileSourceStrategy: Pushed Filters:
2021-03-28 20:43:15,651 INFO datasources.FileSourceStrategy: Post-Scan Filters:

```

```

2021-03-28 20:43:52,441 INFO executor.Executor: Finished task 0.0 in stage 35.0 (TID 1116). 3855 bytes result sent to driver
2021-03-28 20:43:52,441 INFO executor.Executor: Finished task 8.0 in stage 35.0 (TID 1124). 3855 bytes result sent to driver
2021-03-28 20:43:52,442 INFO scheduler.TaskSetManager: Finished task 0.0 in stage 35.0 (TID 1116) in 24 ms on cc-project-13.cc-s21-project.is2750-pg0.utah.cl
oudlab.us (executor driver) (6/9)
2021-03-28 20:43:52,442 INFO scheduler.TaskSetManager: Finished task 8.0 in stage 35.0 (TID 1124) in 6 ms on cc-project-13.cc-s21-project.is2750-pg0.utah.cl
oudlab.us (executor driver) (7/9)
2021-03-28 20:43:52,442 INFO executor.Executor: Finished task 7.0 in stage 35.0 (TID 1123). 3855 bytes result sent to driver
2021-03-28 20:43:52,443 INFO executor.Executor: Finished task 6.0 in stage 35.0 (TID 1122). 3953 bytes result sent to driver
2021-03-28 20:43:52,443 INFO scheduler.TaskSetManager: Finished task 7.0 in stage 35.0 (TID 1123) in 10 ms on cc-project-13.cc-s21-project.is2750-pg0.utah.cl
oudlab.us (executor driver) (8/9)
2021-03-28 20:43:52,443 INFO scheduler.TaskSetManager: Finished task 6.0 in stage 35.0 (TID 1122) in 13 ms on cc-project-13.cc-s21-project.is2750-pg0.utah.cl
oudlab.us (executor driver) (9/9)
2021-03-28 20:43:52,443 INFO scheduler.TaskSchedulerImpl: Removed TaskSet 35.0, whose tasks have all completed, from pool
2021-03-28 20:43:52,443 INFO scheduler.DAGScheduler: ResultStage 35 (showString at NativeMethodAccessorImpl.java:0) finished in 0.029 s
2021-03-28 20:43:52,443 INFO scheduler.DAGScheduler: Job 17 is finished. Cancelling potential speculative or zombie tasks for this job
2021-03-28 20:43:52,443 INFO scheduler.TaskSchedulerImpl: Killing all running tasks in stage 35: Stage finished
2021-03-28 20:43:52,444 INFO scheduler.DAGScheduler: Job 17 finished: showString at NativeMethodAccessorImpl.java:0, took 0.031195 s
2021-03-28 20:43:52,452 INFO codegen.CodeGenerator: Code generated in 7.517669 ms
+-----+
| year|count|      prediction|
+-----+
| 0.1258848| 2051|[0.1258848]2206.8133996202274|
| 0.1266624| 2222|[0.1266624]3119.6718135910924|
| 0.1302048| 10496|[0.1302048] 7278.249073791638|
| 0.12797568| 8927|[0.12797568] 4661.388261275191|
| 0.13038624| 5853|[0.13038624] 7491.249372484803|
| 0.13101696| 5085|[0.13101696] 8231.678982227837|
| 0.12664512| 6546|[0.12664512] 3099.386070858396|
| 0.1271376| 3256|[0.1271376] 3677.5297387399187|
| 0.1301184| 5622|[0.1301184] 7176.828360128186|
| 0.13042944| 4668|[0.13042944] 7541.963729316543|
| 0.1248912| 34|[0.1248912] 1040.383183490805|
| 0.13021344| 7708|[0.13021344] 7288.391945157986|
| 0.13004064| 8814|[0.13004064] 7085.534517831111|
| 0.12794112| 4848|[0.12794112] 4620.816775809799|
| 0.1258416| 662|[0.1258416] 2156.099033788516|
| 0.1281312| 7004|[0.1281312] 4843.959945869341|
| 0.1288288| 4756|[0.1288288] 5726.38975471191|
| 0.12807936| 7000|[0.12807936] 4783.1027176713105|
| 0.12877056| 3512|[0.12877056] 5594.532426078723|
| 0.13005792| 5772|[0.13005792] 7105.820260563778|
+-----+
only showing top 20 rows

2021-03-28 20:43:52,508 INFO Configuration.deprecation: mapred.output.dir is deprecated. Instead, use mapreduce.output.fileoutputformat.outputdir
2021-03-28 20:43:52,511 INFO io.HadoopMapRedCommitProtocol: Using output committer class org.apache.hadoop.mapred.FileOutputCommitter
2021-03-28 20:43:52,513 INFO output.FileOutputCommitter: File Output Committer Algorithm version is 1
2021-03-28 20:43:52,513 INFO output.FileOutputCommitter: FileOutputCommitter skip cleanup _temporary folders under output directory:false, ignore cleanup f
ailures: false

```

```

8 MiB)
2021-03-28 20:43:53,013 INFO parquet.ParquetWriteSupport: Initialized Parquet WriteSupport with Catalyst schema:
{
  "type" : "struct",
  "fields" : [ {
    "name" : "intercept",
    "type" : "double",
    "nullable" : false,
    "metadata" : { }
  }, {
    "name" : "coefficients",
    "type" : {
      "type" : "udt",
      "class" : "org.apache.spark.ml.linalg.VectorUDT",
      "pyClass" : "pyspark.ml.linalg.VectorUDT",
      "sqlType" : {
        "type" : "struct",
        "fields" : [ {
          "name" : "type",
          "type" : "byte",
          "nullable" : false,
          "metadata" : { }
        }, {
          "name" : "size",
          "type" : "integer",
          "nullable" : true,
          "metadata" : { }
        }, {
          "name" : "indices",
          "type" : {
            "type" : "array",
            "elementType" : "integer",
            "containsNull" : false
          },
          "nullable" : true,
          "metadata" : { }
        }, {
          "name" : "values",
          "type" : {
            "type" : "array",
            "elementType" : "double",
            "containsNull" : false
          },
          "nullable" : true,
          "metadata" : { }
        }, {
          "name" : "scale",
          "type" : "double",
          "nullable" : false,
          "metadata" : { }
        } ]
      }
    },
    "nullable" : true,
    "metadata" : { }
  }, {
    "name" : "scale",
    "type" : "double",
    "nullable" : false,
    "metadata" : { }
  } ]
}
and corresponding Parquet message type:
  "type" : "integer",
  "containsNull" : false
},
"nullable" : true,
"metadata" : { }
}, {
  "name" : "values",
  "type" : {
    "type" : "array",
    "elementType" : "double",
    "containsNull" : false
  },
  "nullable" : true,
  "metadata" : { }
} ]
},
"nullable" : true,
"metadata" : { }
}, {
  "name" : "scale",
  "type" : "double",
  "nullable" : false,
  "metadata" : { }
} ]
}
and corresponding Parquet message type:
message spark_schema {
  required double intercept;
  optional group coefficients {
    required int32 type (INT_8);
    optional int32 size;
    optional group indices (LIST) {
      repeated group list {
        required int32 element;
      }
    }
    optional group values (LIST) {
      repeated group list {
        required double element;
      }
    }
  }
  required double scale;
}

2021-03-28 20:43:53,013 INFO storage.BlockManagerInfo: Removed broadcast_45_piece0 on cc-project-13.cc-s21-project.is2750-pg0.utah.cloudlab.us:38469 in memory (size: 30.1 KiB, free: 365.9 MiB)
2021-03-28 20:43:53,019 INFO storage.BlockManagerInfo: Removed broadcast_43_piece0 on cc-project-13.cc-s21-project.is2750-pg0.utah.cloudlab.us:38469 in memory (size: 25.6 KiB, free: 365.9 MiB)
2021-03-28 20:43:53,026 INFO storage.BlockManagerInfo: Removed broadcast_42_piece0 on cc-project-13.cc-s21-project.is2750-pg0.utah.cloudlab.us:38469 in memory (size: 25.5 KiB, free: 365.9 MiB)
2021-03-28 20:43:53,051 INFO compress.CodecPool: Got brand-new compressor [, snappy]
2021-03-28 20:43:53,284 INFO hadoop.InternalParquetRecordWriter: Flushing mem columnstore to file, allocated memory: 44
2021-03-28 20:43:53,505 INFO output.FileOutputCommitter: Saved output of task attempt_202103282043526025176865859917874_0038_m_000000_1127' to hdfs://cc-project-13:9000/user/ubuntu/linear_regression/part3_model/data/_temporary/_0/task_202103282043526025176865859917874_0038_m_000000_1127
2021-03-28 20:43:53,505 INFO mapred.SparkHadoopMapReduceUtil: attempt_202103282043526025176865859917874_0038_m_000000_1127: Committed
2021-03-28 20:43:53,597 INFO executor.Executor: Finished task 0.0 in stage 38.0 (TID 1127). 3281 bytes result sent to driver

```

The model is saved with name linear\_reg\_spark\_part3\_model

To check the saved model execute: hdfs dfs -ls linear\_reg\_spark\_part3\_model

```
"rw-r--r-- 3 ubuntu supergroup 1298455 2021-03-24 16:21 /user/root/accesslog/user_artists.dat
ubuntu@cc-project-13:~$ hdfs dfs -ls linear_reg_spark_part3_model
Found 2 items
drwxr-xr-x  - ubuntu supergroup          0 2021-03-28 20:46 linear_reg_spark_part3_model/data
drwxr-xr-x  - ubuntu supergroup          0 2021-03-28 20:46 linear_reg_spark_part3_model/metadata
ubuntu@cc-project-13:~$
```

execute: hdfs dfs -ls linear\_reg\_spark\_part3\_model/data

```
ubuntu@cc-project-13:~$ hdfs dfs -ls linear_reg_spark_part3_model/data/
Found 2 items
-rw-r--r-- 3 ubuntu supergroup          0 2021-03-28 20:46 linear_reg_spark_part3_model/data/_SUCCESS
-rw-r--r-- 3 ubuntu supergroup      1943 2021-03-28 20:46 linear_reg_spark_part3_model/data/part-00000-3bcc953f-7374-4e09-95fa-c8c2084df42a-c000.snappy.parquet
ubuntu@cc-project-13:~$
```