

Cloud Computing Capstone - Data Extraction, Batch Processing with Hadoop

Task 1 Report


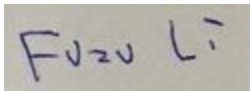
Instructor:

Indranil Gupta, Roy H. Campbell, P. Brighten Godfrey, Ankit Singla, Reza Farivar



Submitted as a partial fulfillment of course completion.

Submitted by

Avatar	Name	Signature
	Fuzu Li	

Video Demonstration: <https://www.youtube.com/watch?v=j7KBRNtLO4M>

A brief overview of how extracted and cleaned the data.

According to the problem, as analyzed, I found the on-time performance related data is most useful for the problem domain, so I picked up the on-time performance data from the specific folder, then uncompress them to my local folder, then copy them into hadoop file system.

The principle for cleaning the raw data is to prevent errors, reduce duplication, and improve efficiency. For this task, I keep trying best to maintain the original format of the data, remove the unstructured data avoid breaking the data process.

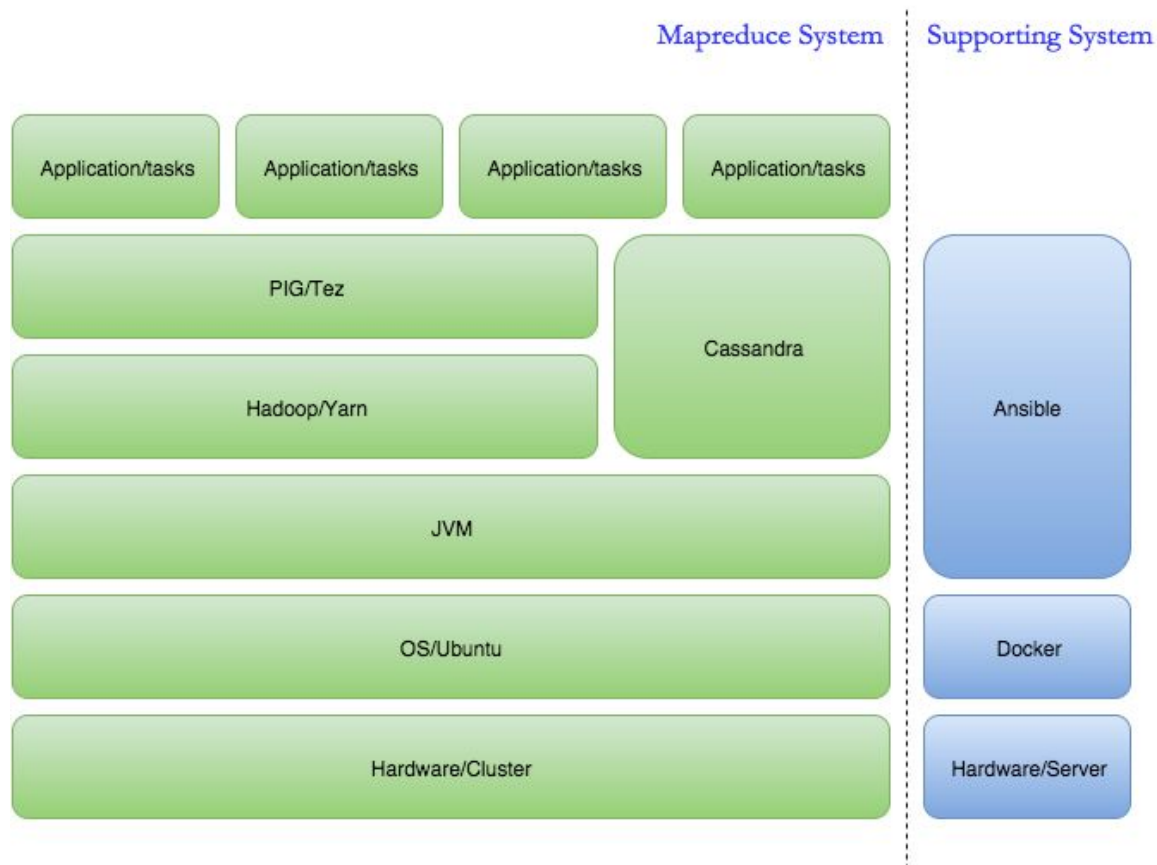
With the principle in mind, I keep the clean structured csv file for On_Time_On_Time_Performance_{1988-2008}_{1-12}.csv, remove the readme.html, in order to be imported by hadoop/pig scripts.

I created an open source repository on github to save some scripts to cleanup the raw data and import to hadoop file system. Please refer to [Reference #1](#) for details.

A brief overview of how integrated each system.

Here is a brief overview for my local hadoop cluster infrastructure:

Hadoop Infrastructure for Cloud Computing Capstone



I have 6 Nodes for Hadoop/PIG, includes 1 for Namenode/ResourceManager, 1 for secondary Namenode, 4 for Datanode/TaskManager.



Nodes of the cluster

Logged in as: dr.who

Cluster

About

Nodes

Node Labels

Applications

NEW

NEW SAVING

SUBMITTED

ACCEPTED

RUNNING

FINISHED

FAILED

KILLED

Scheduler

Tools

Cluster Metrics

Apps Submitted	Apps Pending	Apps Running	Apps Completed	Containers Running	Memory Used	Memory Total	Memory Reserved	VCores Used	VCores Total	VCores Reserved	Active Nodes	Decommissioned Nodes	Lost Nodes	Unhealthy Nodes	Rebooted Nodes
17	0	0	17	0	0 B	32 GB	0 B	0	32	0	4	0	0	0	0

Scheduler Metrics

Scheduler Type		Scheduling Resource Type		Minimum Allocation		Maximum Allocation						
Capacity Scheduler		[MEMORY]		<memory:1024, vCores:1>		<memory:8192, vCores:8>						
Show 20 entries						Search:						
Node Labels	Rack	Node State	Node Address	Node HTTP Address	Last health-update	Health-report	Containers	Mem Used	Mem Avail	VCores Used	VCores Avail	Version
/default-rack		RUNNING	dnjplatbuild10.barnesandnoble.com:37831	dnjplatbuild10.barnesandnoble.com:8042	Thu Feb 04 00:17:30 -0500 2016		0	0 B	8 GB	0	8	2.7.1
/default-rack		RUNNING	dnjplatbuild06.barnesandnoble.com:42747	dnjplatbuild06.barnesandnoble.com:8042	Thu Feb 04 00:17:30 -0500 2016		0	0 B	8 GB	0	8	2.7.1
/default-rack		RUNNING	dnjplatbuild04.barnesandnoble.com:50061	dnjplatbuild04.barnesandnoble.com:8042	Thu Feb 04 00:17:30 -0500 2016		0	0 B	8 GB	0	8	2.7.1
/default-rack		RUNNING	dnjplatbuild08.barnesandnoble.com:45613	dnjplatbuild08.barnesandnoble.com:8042	Thu Feb 04 00:17:30 -0500 2016		0	0 B	8 GB	0	8	2.7.1

Showing 1 to 4 of 4 entries

First Previous 1 Next Last

I have 4 nodes for Cassandra Cluster, some of them share same servers with Hadoop cluster, as following:

```
build@dnjplatbuild02:~$ nodetool status
Datacenter: datacenter1
=====
Status=Up/Down
|/ State=Normal/Leaving/Joining/Moving
-- Address            Load             Tokens           Owns            Host ID                               Rack
UN  10.5.178.84         367.75 MB        256              ?               dc866979-3427-4bd0-8344-49beea35f5a7 rack1
UN  10.5.178.82         345.01 MB        256              ?               2bfb70fe-99ed-4be8-a5f4-ae8e8629caf8 rack1
UN  10.5.178.79         442.47 MB        256              ?               f0abdf90-a967-4798-b2e6-9d89c3f7206b rack1
UN  10.5.178.91         345.88 MB        256              ?               35b4b2b8-8f54-433d-a1ca-0a8b37075575 rack1
```

For creating and maintaining the cluster, I created another open source repository on github to automatically manage the Hadoop/Cassandra cluster to install Java, Hadoop, PIG, Cassandra, Storm and Spark (for task 2), then config them in a cluster.

Please take a deep look at [Reference #2](#).

The approaches and algorithms to answer each question?

Group 1.1



Group 1.2



The results of each question is as following:

Group 1.1's result:

<https://gist.github.com/lifuzu/352989a7a40db1e53b7f>

Group 1.2's result:

<https://gist.github.com/lifuzu/511b649777b876ba9723>

Group 2.1's full result:

<https://gist.github.com/lifuzu/c844c01e9824c6d32bc1>

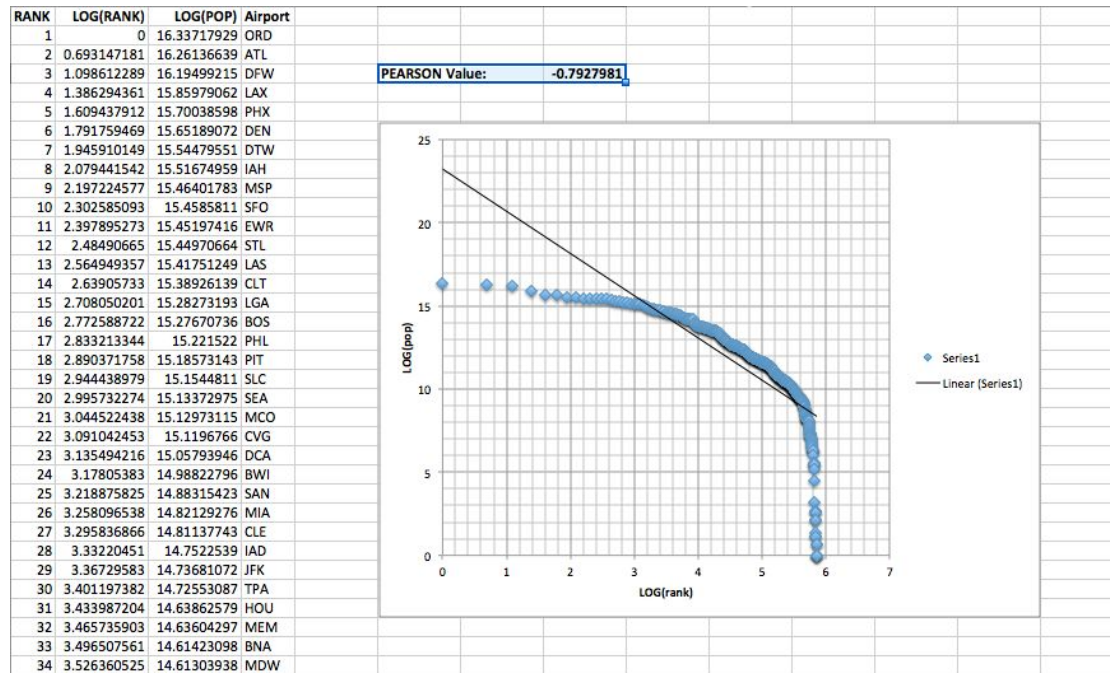
Group 2.2's full result:

<https://gist.github.com/lifuzu/439933ee2718b92c36c1>

Group 2.3's full result:

<https://gist.github.com/lifuzu/629354f6cd916f264997>

Group 3.1's result:



Pearson Product-Moment Correlation Coefficient: **-0.7927981**

Group 3.2's full result:

<https://gist.github.com/lifuzu/7080eadf72d38778f6c7>

Some optimizations employed on the task:

I employed the following optimizations on the Hadoop cluster and application level:

1. Project Early and Often
2. Filter Early and Often
3. Drop Nulls Before a Join
4. Use the PARALLEL Clause - SET DEFAULT_PARALLEL 10;

Opinion about whether the results make sense and are useful in any way:

Definitely, the results came from the real data should be very useful to analyze the aviation history and forecast the future to provide convenience to potential passengers.

References:

1. <https://github.com/lifuzu/cccapstone-tasks>
2. <https://github.com/lifuzu/cccapstone-cluster-manager>