# Lab 2

## Divya K Raman, EE15B085

**Answer 1**

CAP Theorem is comprised of three technical terms. C stands for Consistency where all nodes see the data in homogeneous form i.e. every node has the same knowledge of data at any instant of time. A stands for availability i.e. A guarantee that every request receives a response which may be processed or failed. P stands for Partition Tolerance i.e. the system continues to operate even if a message is lost or part of the system fails Any system which satisfies two out of these three properties is a distributed system. Hadoop supports the Availability and Partition Tolerance property. The Consistency property is not supported because only namenode has the information of where the replicas are placed. This information is not available with each and every node of the cluster.

**Answer 2**

Apache Hive provides a database query interface which resembles SQL to data stored in databases and file systems that integrate with Hadoop. Pig enables data workers to write complex data transformations and manipulations without knowing Java. Both Pig and Hive are feature complete. Either of them can be used to do a particular task. Depending on their primary use case, data scientists generally prefer one over the other. Hive is easier for users proficient with SQL to query data whereas Pig has data flow strengths. Pig can bring data into Hadoop and get it into the form suitable for querying.

**Answer 3**

Before uploading any of the 4 files:

After uploading trucks.csv

After uploading seasons.csv

After products.tsv

As files are added to the server, number of blocks increases which is understandable as new files are being added to the server. Under replicated remains the same. Under replicated blocks are blocks that do not meet their target replication for the file they belong to. This number indicates the number of copies of blocks missing from the environment. While uploading files, there were probably no copies of blocks which went missing and hence this number remains unchanged. Heap size first increases and then decreases. In Java, heap memory is an area of memory reserved for data that is created at runtime that is, when the program actually executes. NameNode heap size depends on many factors, such as the number of files, the number of blocks, and the load on the system.

**Answer 4**

In map, each worker node applies the map function to the local data, and writes the output to a temporary storage. A master node ensures that only one copy of redundant input data is processed. In Reduce, worker nodes now process each group of output data, per key, in parallel. In the given problem, each row describes a match officiated by a referee. In the map stage, we map the matches to the corresponding referees. In the reduce stage, we combine the matches officiated by each referee and take a total count for each referee.

**Answer 5**

**Zeppelin**   Notebook ▾   Job

lab 2 ▷ ⤢ ▤ ✎ ⟲ ⬇   🔍  🗑

```
val hiveContext = new org.apache.spark.sql.SparkSession.Builder().getOrCreate();
val riskFactorDataFrame = spark.read.format("csv").option("header", "true").load("hdfs:///tmp/data/trucks.csv");
riskFactorDataFrame.createOrReplaceTempView("trucks");
```

```
hiveContext: org.apache.spark.sql.SparkSession = org.apache.spark.sql.SparkSession@dfc138e
riskFactorDataFrame: org.apache.spark.sql.DataFrame = [driverid: string, truckid: string ... 109 more fields]
```

Took 2 sec. Last updated by anonymous at February 04 2019, 3:57:34 PM.

SPARK JOB FINISHED ▷ ⤢ ▤ ⚙

```
val ans = hiveContext.sql("SELECT driverid, jun13_miles FROM trucks LIMIT 15");
ans.show();
```

```
+--------+-----------+
|driverid|jun13_miles|
+--------+-----------+
|      A1|       9217|
|      A2|      12058|
|      A3|      13652|
|      A4|      12687|
|      A5|      10233|
|      A6|      14488|
|      A7|      10938|
|      A8|      11392|
|      A9|      12601|
|     A10|      13699|
|     A11|      12447|
|     A12|      10006|
|     A13|       9740|
```

SPARK JOB FINISHED ▷ ⤢ ▤ ⚙

📄 lab 2.json        ∧                                          Show all   ✕

---

```
|     A10|      13699|
|     A11|      12447|
|     A12|      10006|
|     A13|       9740|
|     A14|      10608|
|     A15|       9659|
```

Took 3 sec. Last updated by anonymous at February 04 2019, 3:57:37 PM.

SPARK JOB FINISHED ▷ ⤢ ▤ ⚙

```
%pyspark
from pyspark.sql import SQLContext
from pyspark.sql.types import DoubleType

sqlContext = SQLContext(sc)
dfObj = spark.read.csv("hdfs:///tmp/data/trucks.csv", header=True, mode="DROPMALFORMED")
dfObj.columns
```

```
['driverid', 'truckid', 'model', 'jun13_miles', 'jun13_gas', 'may13_miles', 'may13_gas', 'apr13_miles', 'apr13_gas', 'mar13_miles', 'mar13_gas', 'feb13_miles', 'feb13_gas', 'jan13_miles', 'jan13_
gas', 'dec12_miles', 'dec12_gas', 'nov12_miles', 'nov12_gas', 'oct12_miles', 'oct12_gas', 'sep12_miles', 'sep12_gas', 'aug12_miles', 'aug12_gas', 'jul12_miles', 'jul12_gas', 'jun12_miles', 'jun12
_gas', 'may12_miles', 'may12_gas', 'apr12_miles', 'apr12_gas', 'mar12_miles', 'mar12_gas', 'feb12_miles', 'feb12_gas', 'jan12_miles', 'jan12_gas', 'dec11_miles', 'dec11_gas', 'nov11_miles', 'nov1
1_gas', 'oct11_miles', 'oct11_gas', 'sep11_miles', 'sep11_gas', 'aug11_miles', 'aug11_gas', 'jul11_miles', 'jul11_gas', 'jun11_miles', 'jun11_gas', 'may11_miles', 'may11_gas', 'apr11_miles', 'apr
11_gas', 'mar11_miles', 'mar11_gas', 'feb11_miles', 'feb11_gas', 'jan11_miles', 'jan11_gas', 'dec10_miles', 'dec10_gas', 'nov10_miles', 'nov10_gas', 'oct10_miles', 'oct10_gas', 'sep10_miles', 'se
p10_gas', 'aug10_miles', 'aug10_gas', 'jul10_miles', 'jul10_gas', 'jun10_miles', 'jun10_gas', 'may10_miles', 'may10_gas', 'apr10_miles', 'apr10_gas', 'mar10_miles', 'mar10_gas', 'feb10_miles', 'f
eb10_gas', 'jan10_miles', 'jan10_gas', 'dec09_miles', 'dec09_gas', 'nov09_miles', 'nov09_gas', 'oct09_miles', 'oct09_gas', 'sep09_miles', 'sep09_gas', 'aug09_miles', 'aug09_gas', 'jul09_miles', '
jul09_gas', 'jun09_miles', 'jun09_gas', 'may09_miles', 'may09_gas', 'apr09_miles', 'apr09_gas', 'mar09_miles', 'mar09_gas', 'feb09_miles', 'feb09_gas', 'jan09_miles', 'jan09_gas']
```

Took 2 sec. Last updated by anonymous at February 04 2019, 3:58:03 PM.

FINISHED ▷ ⤢ ▤ ⚙

```
%pyspark
from pyspark.sql.types import DoubleType
changedTypedfTmp = dfObj.withColumn("jun13_miles", dfObj["jun13_miles"].cast(DoubleType()))
changedTypedf = changedTypedfTmp.withColumn("jun13_gas", dfObj["jun13_gas"].cast(DoubleType()))
changedTypedf.registerTempTable("trucks")
```

Took 0 sec. Last updated by anonymous at February 04 2019, 3:58:46 PM.

📄 lab 2.json        ∧                                          Show all   ✕

```
%pyspark
from pyspark.sql.types import DoubleType
changedTypedfTmp = dfObj.withColumn("jun13_miles", dfObj["jun13_miles"].cast(DoubleType()))
changedTypedf = changedTypedfTmp.withColumn("jun13_gas", dfObj["jun13_gas"].cast(DoubleType()))
changedTypedf.registerTempTable("trucks")
```

FINISHED

Took 0 sec. Last updated by anonymous at February 04 2019, 3:58:46 PM.

```
%sql
SELECT jun13_miles, jun13_gas FROM trucks
```

SPARK JOB FINISHED

settings ▾

| jun13_miles | jun13_gas |
|---|---|
| 9217 | 1914 |
| 12058 | 2335 |
| 13652 | 2899 |
| 12687 | 2439 |
| 10233 | 1825 |
| 14488 | 2883 |
| 10938 | 2231 |
| 11392 | 2280 |

Took 1 sec. Last updated by anonymous at February 04 2019, 3:58:59 PM. (outdated)

---

Took 1 sec. Last updated by anonymous at February 04 2019, 3:58:59 PM. (outdated)

```
val hiveContext = new org.apache.spark.sql.SparkSession.Builder().getOrCreate();
val riskFactorDataFrame = spark.read.format("csv").option("header", "true").load("hdfs:///tmp/data/season-1213_csv.csv");
riskFactorDataFrame.createOrReplaceTempView("seasons");
```

SPARK JOB FINISHED

```
hiveContext: org.apache.spark.sql.SparkSession = org.apache.spark.sql.SparkSession@dfc138e
riskFactorDataFrame: org.apache.spark.sql.DataFrame = [Date: string, HomeTeam: string ... 20 more fields]
hiveContext: org.apache.spark.sql.SparkSession = org.apache.spark.sql.SparkSession@dfc138e
riskFactorDataFrame: org.apache.spark.sql.DataFrame = [Date: string, HomeTeam: string ... 20 more fields]
```

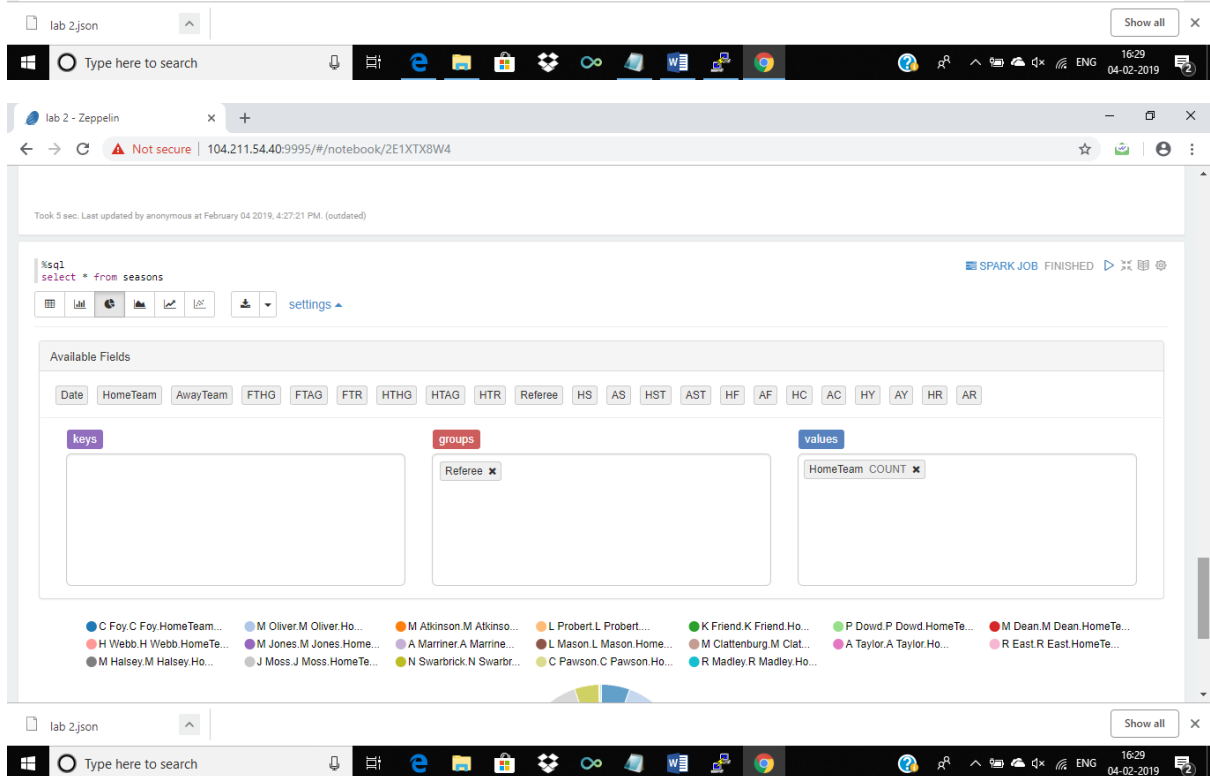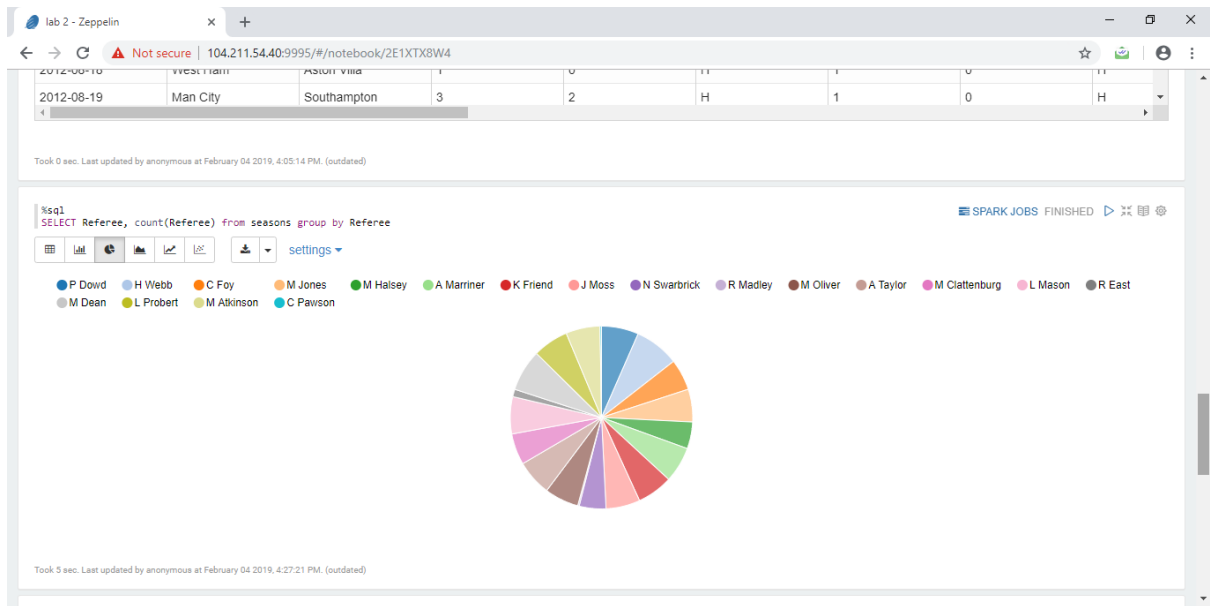Took 1 sec. Last updated by anonymous at February 04 2019, 4:01:42 PM.
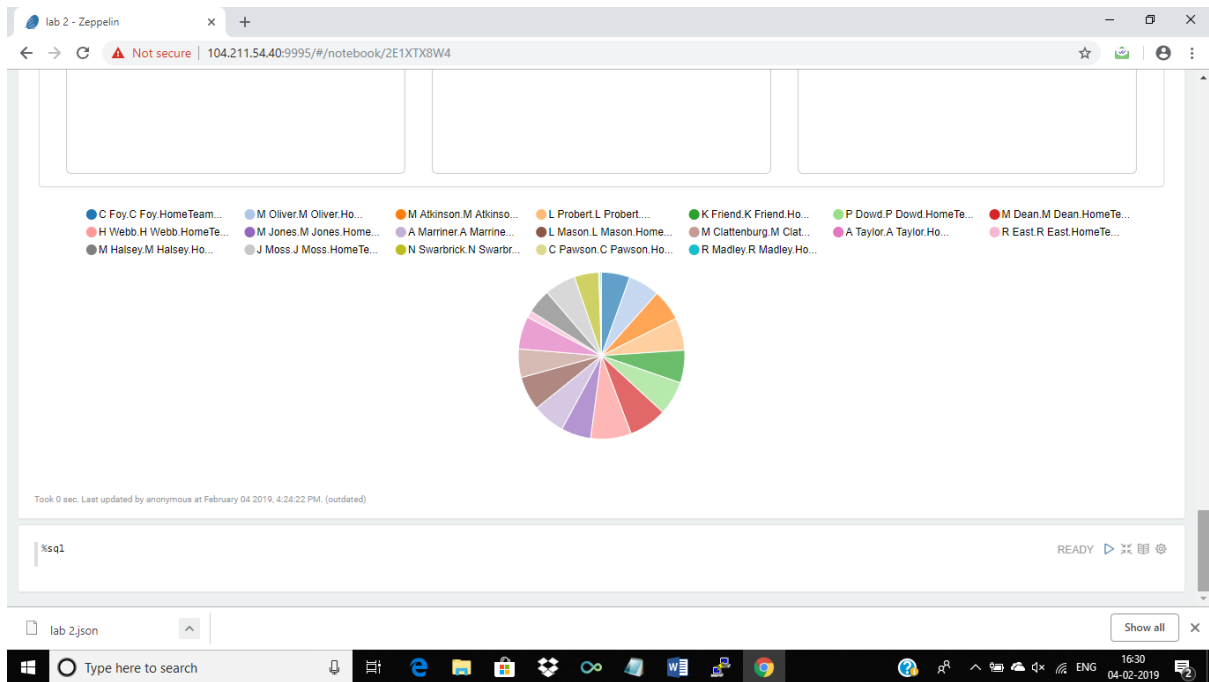
```
%sql
SELECT * from seasons
```

SPARK JOB FINISHED

settings ▾

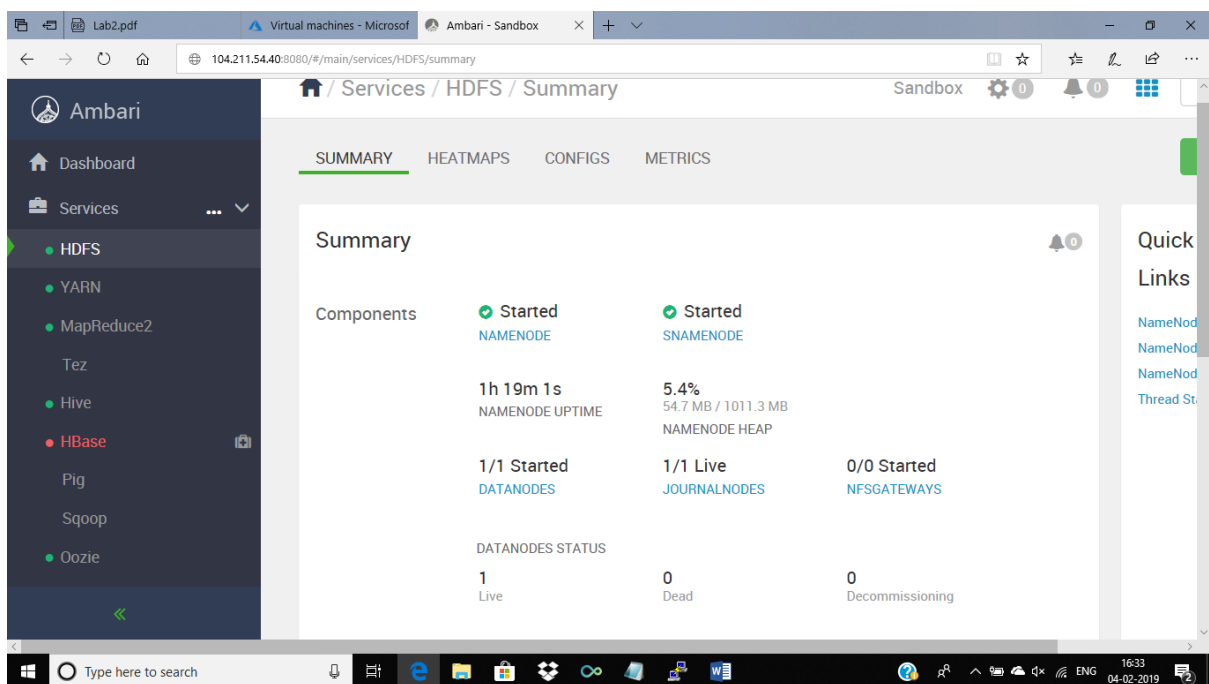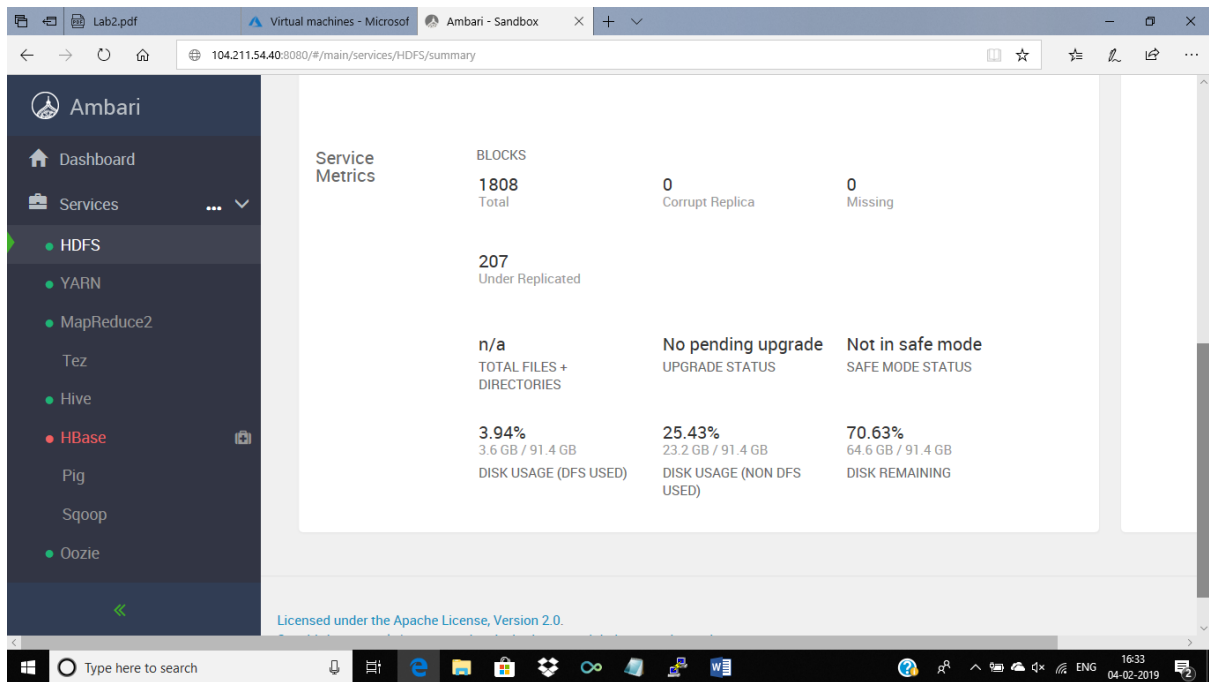| Date | Home Team | Away Team | FTHG | FTAG | FTR | HTHG | HTAG | HTR |
|---|---|---|---|---|---|---|---|---|
| 2012-08-18 | Arsenal | Sunderland | 0 | 0 | D | 0 | 0 | D |
| 2012-08-18 | Fulham | Norwich | 5 | 0 | H | 2 | 0 | H |
| 2012-08-18 | Newcastle | Tottenham | 2 | 1 | H | 0 | 0 | D |
| 2012-08-18 | QPR | Swansea | 0 | 5 | A | 0 | 1 | A |
| 2012-08-18 | Reading | Stoke | 1 | 1 | D | 0 | 1 | A |
| 2012-08-18 | West Brom | Liverpool | 3 | 0 | H | 1 | 0 | H |
| 2012-08-18 | West Ham | Aston Villa | 1 | 0 | H | 1 | 0 | H |
| 2012-08-19 | Man City | Southampton | 3 | 2 | H | 1 | 0 | |

Number of blocks after running the notebooks:

References:

https://data-flair.training/forums/topic/what-is-cap-theorem-what-aspects-hadoop-supports-from-this-theorem/

https://stackoverflow.com/questions/19923196/cap-with-distributed-system

https://hortonworks.com/tutorial/how-to-process-data-with-apache-hive/

https://hortonworks.com/tutorial/beginners-guide-to-apache-pig/

https://stackoverflow.com/questions/36977746/in-hadoop-whats-under-replication-and-over-replication-mean-and-how-does-it-work

https://docs.hortonworks.com/HDPDocuments/HDP2/HDP-2.6.5/bk_command-line-installation/content/configuring-namenode-heap-size.html