

Name: Alvyn Abranches

Class: FY MSc Data Science and Big Data Analytics

Roll No 1

Exercise 1

```
[cloudera@quickstart ~]$ hive
```

Logging initialized using configuration in file:/etc/hive/conf.dist/hive-log4j.properties

WARNING: Hive CLI is deprecated and migration to Beeline is recommended.

```
hive> show databases;
```

OK

default

Time taken: 5.675 seconds, Fetched: 1 row(s)

```
hive> create database alvyn;
```

OK

Time taken: 1.16 seconds

```
hive> use alvyn;
```

OK

Time taken: 0.079 seconds

```
hive> create table hospitals(DRGDefinition string, ProviderId int, ProviderName string, ProviderStreetAddress string, ProviderCity string, ProviderState string, ProviderZipCode int, HospitalReferralRegionDescription string, TotalDischarges int, AverageCoveredCharges double, AverageTotalPayments double, AverageMedicarePayments double) row format delimited fields terminated by ',' stored as textfile;
```

OK

Time taken: 0.782 seconds

```
hive> load data local inpath 'hospital_data.csv' into table hospitals;
```

Loading data to table alvyn.hospitals

Table alvyn.hospitals stats: [numFiles=1, totalSize=26841576]

OK

Time taken: 1.205 seconds

```
hive> select * from hospitals limit 5;
```

OK

DRGDefinition	NULL	ProviderName	ProviderStreetAddress	ProviderCity	ProviderState	NULL	HospitalReferralRegionDescription	NULL	NULL	NULL	NULL		
839 - EXTRACRANIAL PROCEDURES W/O CC/MCC	10001	SOUTHEAST ALABAMA MEDICAL CENTER	1188 ROSS CLARK CIRCLE	DOTHAN	AL	36301	AL - Dothan	91	32963.07	5777.24	4763.73		
839 - EXTRACRANIAL PROCEDURES W/O CC/MCC	10005	MARSHALL MEDICAL CENTER SOUTH	2505 U S HIGHWAY 431 NORTH	BOAZ	AL	35937	AL - Birmingham	14	15131.85	3787.57	4976.71		
839 - EXTRACRANIAL PROCEDURES W/O CC/MCC	10006	ELIZA COFFEE MEMORIAL HOSPITAL	205 MARENGO STREET	FLORENCE	AL	35631	AL - Birmingham	24	37560.37	5434.95	4453.79		
839 - EXTRACRANIAL PROCEDURES W/O CC/MCC	10011	ST VINCENT'S EAST	50 MEDICAL PARK EAST DRIVE	BIRMINGHAM	AL	35235	AL - Birmingham	25	13998.28	5417.56	4129.16		

Time taken: 0.063 seconds, Fetched: 5 row(s)

Exercise 2

```

hive> select AverageTotalPayments, ProviderState from hospitals order by AverageTotalPayments desc, ProviderState asc limit 5;
FAILED: SemanticException [Error 10001]: Line 1:48 Table not found 'hospitals'
hive> select AverageTotalPayments, ProviderState from hospitals order by AverageTotalPayments desc, ProviderState asc limit 5;
Query ID = cloudera_20191118195454_e3c1a9c0-f3c9-449a-a23b-9383194f10b9
Total jobs = 1
Launching Job 1 out of 1
Number of reduce tasks determined at compile time: 1
In order to change the average load for a reducer (in bytes):
  set hive.exec.reducers.bytes.per.reducer=<number>
In order to limit the maximum number of reducers:
  set hive.exec.reducers.max=<number>
In order to set a constant number of reducers:
  set mapreduce.job.reduces=<number>
Starting Job = job_1574131451414_0001, Tracking URL = http://quickstart.cloudera:8088/proxy/application_1574131451414_0001/
Kill Command = /usr/lib/hadoop/bin/hadoop job -kill job_1574131451414_0001
Hadoop job information for Stage-1: number of mappers: 1; number of reducers: 1
2019-11-18 19:54:30,615 Stage-1 map = 0%, reduce = 0%
2019-11-18 19:54:37,972 Stage-1 map = 100%, reduce = 0%, Cumulative CPU 3.94 sec
2019-11-18 19:54:44,204 Stage-1 map = 100%, reduce = 100%, Cumulative CPU 5.93 sec
MapReduce Total cumulative CPU time: 5 seconds 930 msec
Ended Job = job_1574131451414_0001
MapReduce Jobs Launched:
Stage-Stage-1: Map: 1 Reduce: 1 Cumulative CPU: 5.93 sec HDFS Read: 26850270 HDFS Write: 111 SUCCESS
Total MapReduce CPU Time Spent: 5 seconds 930 msec
OK
447714.88      SAN FRANCISCO
414555.91      SAN FRANCISCO
400675.86      GREENBRAE
391446.0       GREENBRAE
381799.81      SAN FRANCISCO
Time taken: 25.371 seconds, Fetched: 5 row(s)

```

Exercise 3

```

hive> select ProviderState, avg(AverageMedicarePayments) from hospitals group by ProviderState limit 5;
Query ID = cloudera_20191118195656_e377defc-c730-427e-a320-64f3577a56c4
Total jobs = 1
Launching Job 1 out of 1
Number of reduce tasks not specified. Estimated from input data size: 1
In order to change the average load for a reducer (in bytes):
  set hive.exec.reducers.bytes.per.reducer=<number>
In order to limit the maximum number of reducers:
  set hive.exec.reducers.max=<number>
In order to set a constant number of reducers:
  set mapreduce.job.reduces=<number>
Starting Job = job_1574131451414_0003, Tracking URL = http://quickstart.cloudera:8088/proxy/application_1574131451414_0003/
Kill Command = /usr/lib/hadoop/bin/hadoop job -kill job_1574131451414_0003
Hadoop job information for Stage-1: number of mappers: 1; number of reducers: 1
2019-11-18 19:57:01,749 Stage-1 map = 0%, reduce = 0%
2019-11-18 19:57:06,957 Stage-1 map = 100%, reduce = 0%, Cumulative CPU 2.64 sec
2019-11-18 19:57:12,143 Stage-1 map = 100%, reduce = 100%, Cumulative CPU 4.3 sec
MapReduce Total cumulative CPU time: 4 seconds 300 msec
Ended Job = job_1574131451414_0003
MapReduce Jobs Launched:
Stage-Stage-1: Map: 1 Reduce: 1 Cumulative CPU: 4.3 sec HDFS Read: 26851870 HDFS Write: 117 SUCCESS
Total MapReduce CPU Time Spent: 4 seconds 300 msec
OK
2ND FLOOR"    41004.6
INC" 18.0
P O BOX 280"  8306.91
PO BOX 788"   12882.036666666667
10 EAST 31ST ST" 24371.192
Time taken: 17.135 seconds, Fetched: 5 row(s)

```

Exercise 4

```

hive> CREATE TABLE airline (YEAR int, MONTH int, DAY int, DAY OF WEEK int, AIRLINE string, FLIGHT NUMBER int, TAIL NUMBER string, ORIGIN AIRPORT string, DESTINATION AIRPORT string, SCHEDULED DEPARTURE int, DEPARTURE TIME int, DEPARTURE D
ELAY int, TAXI OUT int, WHEELS OFF int, SCHEDULED TIME int, ELAPSED TIME int, AIR TIME int, DISTANCE int, WHEELS ON int, TAXI N int, SCHEDULED ARRIVAL int, ARRIVAL TIME int, ARRIVAL DELAY int, DIVERTED int, CANCELLED int) row format deli
mit fields terminated by ',';
OK
Time taken: 0.129 seconds

```

```
hive> load data local inpath 'airline_data.csv' into table airline;
Loading data to table alvyn.airline
Table alvyn.airline stats: [numFiles=1, totalSize=97368314]
OK
Time taken: 0.724 seconds
```

```
cloudera-quickstart-vm-5.13.0-0-virtualbox (Running) - Oracle VM VirtualBox
File Machine View Input Devices Help
Applications Places System
cloudera@quickstart:~
at org.apache.hadoop.util.RunJar.main(RunJar.java:136)
FAILED: ParseException line 1:83 cannot recognize input near 'United' 'Air' 'Lines' in constant
hive> load data local inpath 'airlines.csv' overwrite into table air partition (AIRLINES='United Air Lines Inc.');
```

IATA_CODE	United Air Lines Inc.
UA	United Air Lines Inc.
AA	United Air Lines Inc.
US	United Air Lines Inc.
F9	United Air Lines Inc.
B6	United Air Lines Inc.
OO	United Air Lines Inc.
AS	United Air Lines Inc.
NK	United Air Lines Inc.
WN	United Air Lines Inc.
DL	United Air Lines Inc.
EV	United Air Lines Inc.
HA	United Air Lines Inc.
MQ	United Air Lines Inc.
VX	United Air Lines Inc.

```
Time taken: 0.582 seconds
hive> select * from air;
OK
IATA_CODE      United Air Lines Inc.
UA             United Air Lines Inc.
AA             United Air Lines Inc.
US             United Air Lines Inc.
F9             United Air Lines Inc.
B6             United Air Lines Inc.
OO             United Air Lines Inc.
AS             United Air Lines Inc.
NK             United Air Lines Inc.
WN             United Air Lines Inc.
DL             United Air Lines Inc.
EV             United Air Lines Inc.
HA             United Air Lines Inc.
MQ             United Air Lines Inc.
VX             United Air Lines Inc.
Time taken: 0.125 seconds, Fetched: 15 row(s)
hive> describe air;
OK
iata_code      string
airlines       string

# Partition Information
# col_name     data_type      comment
airlines       string
Time taken: 0.138 seconds, Fetched: 7 row(s)
hive>
```

Exercise 5

```
hive> create table realestate(street varchar(40),city string,zip int,state string,beds int,baths int,sq_ft int,type string,price int) c
3 buckets row format delimited fields terminated by ',';
OK
Time taken: 0.094 seconds
```

Exercise 6

Map Side Join