#1:

Firm size: small to mid-cap size based on the amount of data collected.

Sector: Ecommerce business

Competitors: Amazon, eBay, Alibaba, JD.com, AliExpress, Wish, Shopify, and Coupang

Technology: Hadoop Ecosystem, Spark, GCP, Excel, and OpenRefine Human

Capital:

- Requires 1-2 personnel
- Several skills are required like data cleaner, data miner, and data visualizer.
- Training will take 5 weeks.

Technologies Deployed:

- Google Cloud Platform
- Hive2 and Hive Meta store
- Spark
- HQL
- HDFS system
- YARN
- OpenRefine
- Excel

**#2.1:**

Action: Removed the in-between white space with excel by manually going to each column cells and deleting them.

**#2.2:**

Action: Removed the commas on each cell of "Product Names" by using OpenRefine's "replace" function which lets you replace a character for another, and I replaced "," with a blank argument.
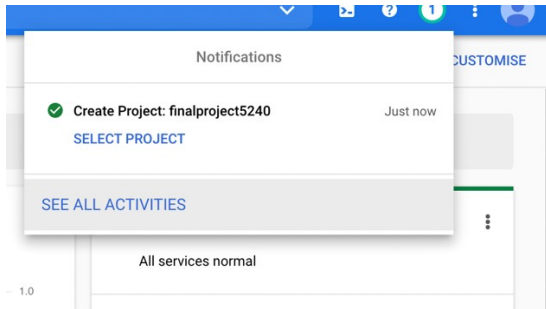
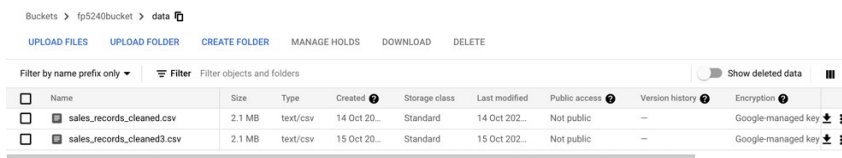**#3:**

*Figure 1: CREATE NEW PROJECT*



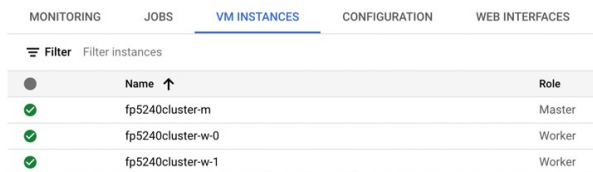*Figure 2: MAKE NEW BUCKET AND UPLOAD THE CLEANED DATASET*



*Figure 3: CREATE THE CLUSTER WITH 1 MASTER & 2 WORKER NODES & SSH IN*



*Figure 4: COPY THE DATASET USING THE "gsutil cp -n gs://fp5240bucket/data/filename" COMMAND TO THE LINUX DIRECTORY*



*Figure 5: PUT THE DATASET INTO THE HDFS STORAGE FROM THE LINUX DIRECTORY I CREATED*



*Figure 6: CONNECT TO HIVE TO START THE SCHEME CREATION*

```
0: jdbc:hive2://localhost:10000/> CREATE EXTERNAL TABLE IF NOT EXISTs sales_records
. . . . . . . . . . . . . . . . . . > (salesrecords string)
. . . . . . . . . . . . . . . . . . > ROW FORMAT DELIMITED
. . . . . . . . . . . . . . . . . . > STORED AS TEXTFILE
. . . . . . . . . . . . . . . . . . > LOCATION '/user/feysele14fy/data/';
No rows affected (2.771 seconds)
0: jdbc:hive2://localhost:10000/> show tables;
+---------------+
|    tab_name   |
+---------------+
| sales_records |
+---------------+
1 row selected (0.241 seconds)
```

*Figure 7: CREATE A SIMPLE SCHEME TO SEE IF EVERYTHIN IS WORKING IN THE METASTORE*

```
0: jdbc:hive2://localhost:10000> CREATE EXTERNAL TABLE IF NOT EXISTS sales_records_1
. . . . . . . . . . . . . . . . . > ( `rowid` string, `orderid` string, `orderdate` string,`shipdate` string, `shipmode` string, `customerid` string, `customername` st
ring, `segment` string, `country` string, `city` string, `state` string, `postalcode` string, `region` string, `productid` string, `category` string, `sub_category`
 string, `product` string, `sales` string, `quantity` string, `discount` string, `profit` string)
. . . . . . . . . . . . . . . . .> ROW FORMAT DELIMITED
. . . . . . . . . . . . . . . . .> FIELDS TERMINATED BY ','
. . . . . . . . . . . . . . . . .> STORED AS TEXTFILE
. . . . . . . . . . . . . . . . .> LOCATION '/user/feysele14fy/data/sales/';
No rows affected (0.222 seconds)
0: jdbc:hive2://localhost:10000> show tables;
+-----------------+
|    tab_name     |
+-----------------+
| sales_records   |
| sales_records_1 |
+-----------------+
2 rows selected (0.106 seconds)
0: jdbc:hive2://localhost:10000>
```

*Figure 8: CREATE THE COMPLEX SCHEME & TERMINATE BY ','*

```
0: jdbc:hive2://localhost:10000> SELECT state, COUNT(customerid) AS NumCustomers FROM sales_records_1
. . . . . . . . . . . . . . . . .> GROUP BY state
. . . . . . . . . . . . . . . . .> ORDER BY NumCustomers DESC LIMIT 5;
+--------------+--------------+
|     state    | numcustomers |
+--------------+--------------+
| California   | 2001         |
| New York     | 1128         |
| Texas        | 985          |
| Pennsylvania | 587          |
| Washington   | 506          |
+--------------+--------------+
5 rows selected (17.438 seconds)
0: jdbc:hive2://localhost:10000>
```

*Figure 9: QUERY FOR TOP 5 STATE BASED ON # OF CUSTOMER RESIDING*

```
0: jdbc:hive2://localhost:10000> SELECT postalcode, SUM(sales) AS TotalSales FROM sales_records_1
. . . . . . . . . . . . . . . . .> GROUP BY postalcode
. . . . . . . . . . . . . . . . .> ORDER BY TotalSales DESC LIMIT 10;
+------------+--------------------+
| postalcode |     totalsales     |
+------------+--------------------+
| 10024      | 78697.182          |
| 10035      | 77357.88500000001  |
| 10009      | 54761.49599999996  |
| 94122      | 52667.46700000001  |
| 10011      | 45551.59800000001  |
| 98105      | 41838.00799999998  |
| 98115      | 41160.90800000001  |
| 19134      | 39390.292999999976 |
| 32216      | 39133.327999999994 |
| 90049      | 37961.012          |
+------------+--------------------+
10 rows selected (8.066 seconds)
0: jdbc:hive2://localhost:10000>
```

*Figure 10: QUERY FOR TOP TOTAL SALES BASED ON TOP 10 ZIP CODES*

#4:

*Figure 11: STARTING SPARK-SQL*



*Figure 12: QUERY FOR TOP 5 STATE BASED ON # OF CUSTOMER RESIDING (SPARK)*



*Figure 13: QUERY FOR TOTAL SALES BASED ON TOP 10 ZIP CODES (SPARK)*

#5:

```
0: jdbc:hive2://localhost:10000>
. . . . . . . . . . . . . . . .>
. . . . . . . . . . . . . . . .>
. . . . . . . . . . . . . . . .>
+---------------+---------------+
|     state     |  numcustomer  |
+---------------+---------------+
| California    | 2001          |
| New York      | 1128          |
| Texas         | 985           |
| Pennsylvania  | 587           |
| Washington    | 506           |
+---------------+---------------+
5 rows selected (28.662 seconds)
0: jdbc:hive2://localhost:10000>
. . . . . . . . . . . . . . . .>
. . . . . . . . . . . . . . . .>
. . . . . . . . . . . . . . . .>
+---------------+---------------+
|     state     |  numcustomer  |
+---------------+---------------+
| California    | 2001          |
| New York      | 1128          |
| Texas         | 985           |
| Pennsylvania  | 587           |
| Washington    | 506           |
+---------------+---------------+
5 rows selected (6.626 seconds)
0: jdbc:hive2://localhost:10000>
. . . . . . . . . . . . . . . .>
. . . . . . . . . . . . . . . .>
. . . . . . . . . . . . . . . .>
+---------------+---------------+
|     state     |  numcustomer  |
+---------------+---------------+
| California    | 2001          |
| New York      | 1128          |
| Texas         | 985           |
| Pennsylvania  | 587           |
| Washington    | 506           |
+---------------+---------------+
5 rows selected (1.808 seconds)
0: jdbc:hive2://localhost:10000>
```

*Figure 14: QUERY SPEEDTEST ON TOP 5 STATES (HIVE)*

```
| 94122        | 52667.46700000001   |
| 10011        | 45551.59800000001   |
| 98105        | 41838.00799999998   |
| 98115        | 41160.90800000001   |
| 19134        | 39390.292999999976  |
| 32216        | 39133.327999999994  |
| 90049        | 37961.012           |
+------------+--------------------+
10 rows selected (17.603 seconds)
0: jdbc:hive2://localhost:10000> SELECT
. . . . . . . . . . . . . . . .> FROM s
. . . . . . . . . . . . . . . .> GROUP
. . . . . . . . . . . . . . . .> ORDER
+------------+--------------------+
| postalcode |     totalsales     |
+------------+--------------------+
| 10024        | 78697.182           |
| 10035        | 77357.88500000001   |
| 10009        | 54761.49599999996   |
| 94122        | 52667.46700000001   |
| 10011        | 45551.59800000001   |
| 98105        | 41838.00799999998   |
| 98115        | 41160.90800000001   |
| 19134        | 39390.292999999976  |
| 32216        | 39133.327999999994  |
| 90049        | 37961.012           |
+------------+--------------------+
10 rows selected (17.548 seconds)
0: jdbc:hive2://localhost:10000> SELECT
. . . . . . . . . . . . . . . .> FROM s
. . . . . . . . . . . . . . . .> GROUP
. . . . . . . . . . . . . . . .> ORDER
+------------+--------------------+
| postalcode |     totalsales     |
+------------+--------------------+
| 10024        | 78697.182           |
| 10035        | 77357.88500000001   |
| 10009        | 54761.49599999996   |
| 94122        | 52667.46700000001   |
| 10011        | 45551.59800000001   |
| 98105        | 41838.00799999998   |
| 98115        | 41160.90800000001   |
| 19134        | 39390.292999999976  |
| 32216        | 39133.327999999994  |
| 90049        | 37961.012           |
+------------+--------------------+
10 rows selected (7.017 seconds)
```

*Figure 15: QUERY SPEEDTEST ON TOP 10 ZIPCODES (HIVE)*

*Figure 16: QUERY SPEEDTEST ON TOP 5 STATE (SPARK)*

```
21/10/16 02:18:49 INFO org.a
10024    78697.182
10035    77357.885
10009    54761.49600000001
94122    52667.467
10011    45551.59799999998
98105    41838.007999999994
98115    41160.90799999999
19134    39390.293000000005
32216    39133.328
90049    37961.012
Time taken: 17.718 seconds,
spark-sql> SELECT postalcode
         > FROM sales_recor
         > GROUP BY postalco
         > ORDER BY TotalSa
21/10/16 02:20:36 INFO org.a
10024    78697.182
10035    77357.885
10009    54761.496
94122    52667.467000000004
10011    45551.59799999998
98105    41838.008
98115    41160.907999999996
19134    39390.293
32216    39133.32800000001
90049    37961.012
Time taken: 3.451 seconds,
spark-sql> SELECT postalcode
         > FROM sales_recor
         > GROUP BY postalco
         > ORDER BY TotalSa
21/10/16 02:20:53 INFO org.a
10024    78697.182
10035    77357.88499999998
10009    54761.49559999999
94122    52667.467000000004
10011    45551.59799999998
98105    41838.007999999994
98115    41160.907999999996
19134    39390.293000000005
32216    39133.328
90049    37961.011999999995
Time taken: 6.759 seconds,
```

*Figure 17: QUERY SPEEDTEST ON TOP 10 ZIPCODES (SPARK)* 1$^{ST}$ QUERY SPEED (sec):
HIVE

1. 28.7
2. 6.6
3. 1.8

SPARK

1. 20.6
2. 10.9
3. 7.4

**Report**:
Spark initially had a 28.2% decrease in time taken on the 1st round, but hive was able to the cut the time by 93.7% at the last round compared to SPARK's 64.1%.

2ND QUERY SPEED (sec):
HIVE
1. 17.6
2. 17.5
3. 7.0
SPARK
1. 17.7
2. 3.5
3. 6.8

**Report**:
There was only 0.1 second difference at the 1st round for the second query test, but at the second and last round Spark made noticeable strides in the time. However, Hive was able to close the gap to only 0.2 seconds which is not significant.

Although Spark is technically faster because it loads the dataset in ram rather than on a solid drive, there is no big gap between the two systems because the dataset is not big enough. In addition, the queries are not to intensive which results in less variances. As a result, I recommend the company uses Hive for the time being until we are presented with enormous data.