# RETAIL GROCERY ANALYSIS

Documented By

Aruna Balamurugan

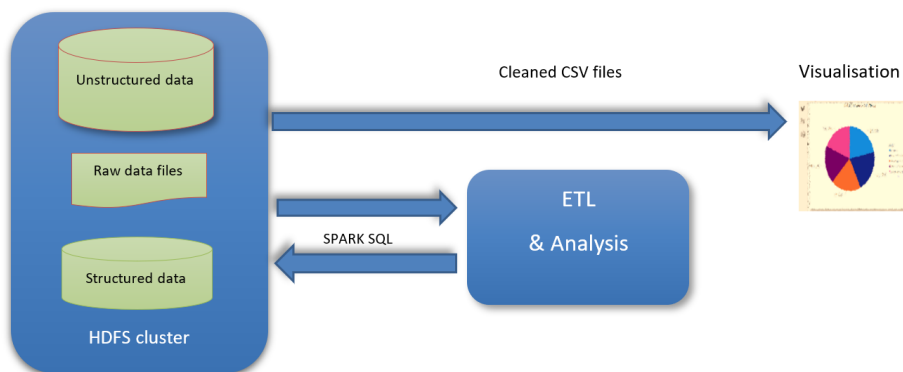F07064

# RETAIL GROCERY ANALYSIS

## Objective:

- To perform analysis on a retail grocery company.
- To extract and build ETL pipeline and produce aggregated data into data lake.
- To analyse and provide business insights on improving sales volume and customer retention.
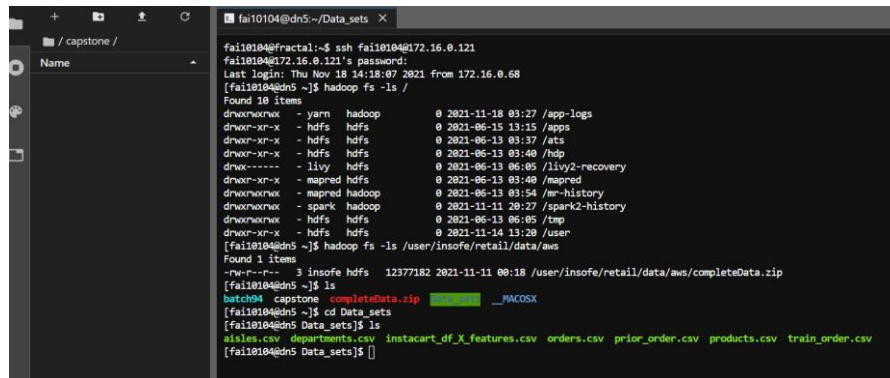
## Solution:

- Extract unstructured raw data from HDFS located in: /user/insofe/retail/data/aws
- Use SparkSQL to read CSV file by specifying path location.
- Create schema and dataframes for each schema.
- Check for null values and redundancies.
- Merge dataframes and perform aggregations. One final dataframe for cleaned data.
- Remove duplicate columns
- Write dataframes into CSV file and load it back on HDFS. ETL pipelining is now achieved.
- Load the output CSV file into local machine and import to Power BI for visualisation.
- Remove duplicate rows and columns, null values, perform necessary transformation for analysis.
- Analyse for product-order insights by creating visuals.
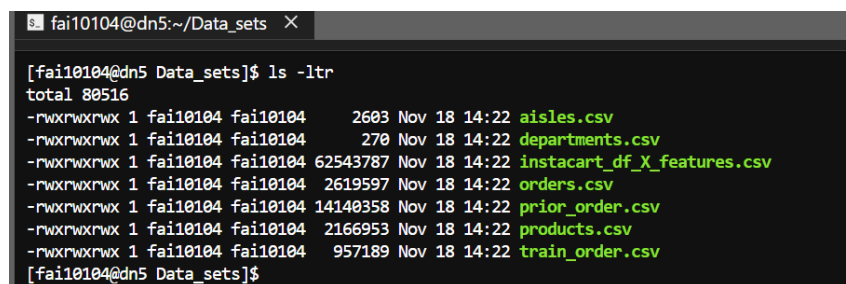
## Solution Architecture:

**Steps:**

1. Open Terminal in JupyterLab

2. Connect to Hadoop by using ssh connection:
   a. Provide IP address and password.

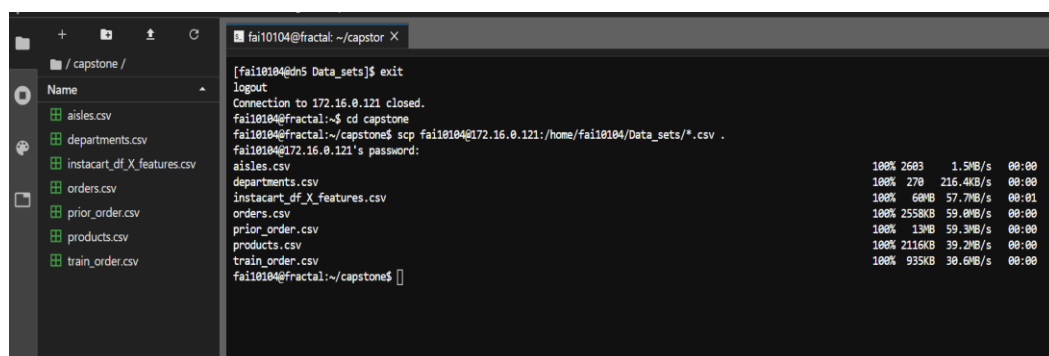   b. Connect to HDFS by using hadoop fs -ls /path location



   c. Use ls -ltr for listing all files and directories in the cluser



   d. Exit the connection.

   e. Use scp cluster_location local_path . for loading all CSV files into local machine (Optional)

3. Open another terminal to start connection and use Pyspark.
   (Note: PySpark is used inside the connection throughout the ETL pipeline).



4. Start a spark session by importing necessary packages.
5. Use spark.read.csv for reading the loaded CSV file:



Fig.1: df.show() is used to display data in aisle.csv

6. Schema Creation:

   - StructType() is used to create schema by defining fields to allow nullable(True) or not (False)
   - To convert schema to dataframe, use
     df=spark.read.csv("path",header=True,schema=schema)
     This dataframe df reads file located in "path", defines schema as mentioned in "schema" variable above and keeps the first row as the header.

a. aisle.csv



Fig.2: schema_ais = schema for aisle; df_ais = dataframe for aisle schema.

b. department.csv



Fig.3: schema_dept = schema for dept data; df_dept = dataframe for dept schema.

c. order.csv



Fig.4: schema_orders = schema for orders data; df_order = dataframe for orders schema.

d. prior_order.csv



Fig.5: schema_prior_order = schema for prior_order data; df_po = dataframe for prior_order schema.

e. product.csv



Fig.6: schema_pdt = schema for product data; df_pdt = dataframe for product schema.

f. train_order.csv



Fig.7: schema_train_order = schema for train_orders data; df_to = dataframe for train_order schema

7. Displaying column names:

Use print(df.columns) to print column names as a list:



8. Display the datatypes of the columns:

Use df.dypes for diplaying data types of all columns in a dataframe.
For loop is used to display column and its datatypes one after the other, separated by a comma.



9. Check for Null values in the columns:

Use df.column_name.isNull() in each column of a dataframe to display all the null values present.
Use df.filter() to display according to the specified condition (null values, in this case)
Use count() to display the count of null values in each column.

a. Orders.csv (df_order)

b.  Department.csv (df_dept) and aisle.csv (df_ais)



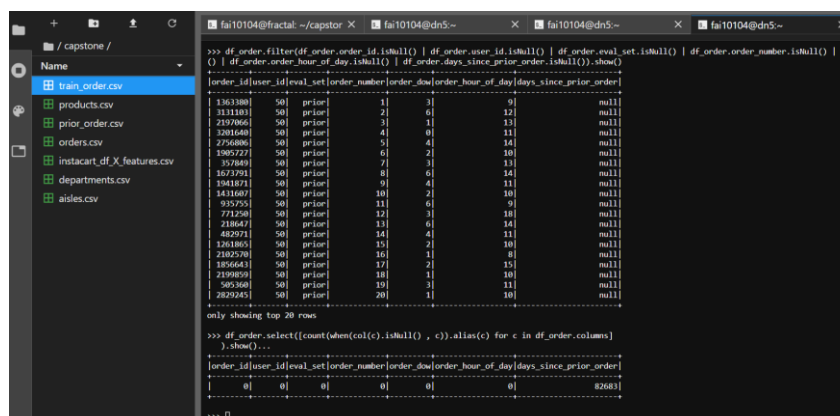c.  Prior_order.csv (df_po) and train_orders.csv (df_to)



d.  Products.csv (df_pdt)



10. Merging of datasets:
    Use df.join() to merge two or more datasets.
    Join conditions used:
    (i)     Left Join Product with aisle using aisle_id.
    (ii)    Left join above joined product with department using department_id.
    (iii)   Inner join above new joined product with prior_orders and train_orders with product_id.
    (iv)    Left join Orders with prior_order and train_orders using order_id.

a. Product – Left Join - aisle:



b. Product_aisle – Left Join - department:



c. Prior_order – Inner Join – joined_product :

d. Train_order – Inner Join – joined_product:



e. Order – Left Join – prior_orders:





f. Order – Left Join – Train_orders:

11. <u>Aggregation of datasets:</u>

Use df1.union(df2) to combine two or more datasets.

    a.   Combining prior_orders and train_orders.



    b.   Combine all the datasets into one dataframe.



12. Loading CSV file into cluster and local machine:





13. Visualisation:

Transformation made before loading:

    a.   Removed duplicated rows and columns
    b.   Removed empty/null values
    c.   Created new table and relationship with the other table.

| dow_number | day |
|---|---|
| 0 | Sunday |
| 1 | Monday |
| 2 | Tuesday |
| 3 | Wednesday |
| 4 | Thursday |
| 5 | Friday |
| 6 | Saturday |

**table1**

Σ add_to_cart_order
  aisle
Σ aisle_id
Σ days_since_prior_order
  department
Σ department_id
  eval_set
  order_dow
Σ order_hour_of_day
Collapse ∧

**dow_to_days**

  day
  dow_number

Collapse ∧

d. Created visuals:

## Top 5 orders by Department

Department ● produce ● dairy eggs ● snacks ● beverages ● frozen

54.3K (9.8%)
67.1K (12.1%)
232.9K (41.9%)
71.6K (12.9%)
130.2K (23.4%)

Insight: 'Produce' department accounted 41.9% of top 5 orders.

## Department wise orders

| Department | Orders |
|---|---|
| fresh fruits | 87.07K |
| fresh vegetables | 85.53K |
| packaged vegetables fruits | 43.86K |
| yogurt | 33.28K |
| packaged cheese | 23.59K |
| water seltzer sparkling water | 21.07K |
| milk | 20.84K |
| chips pretzels | 18.05K |
| soy lactosefree | 16.52K |
| refrigerated | 14.72K |
| bread | 14.70K |
| frozen produce | 12.54K |
| ice cream ice | 11.95K |
| eggs | 11.40K |
| crackers | 11.27K |
| energy granola bars | 11.14K |
| lunch meat | 9.77K |
| fresh herbs | 9.53K |
| cereal | 9.23K |
| frozen meals | 9.13K |

Insight: Fresh Fruits accounted around 87 thousand orders and frozen meals having lower sales. Improving cereals, frozen meals sales can drive business orders.

## Total orders by hour of day

Select all | Sunday | Monday | Tuesday | Wednesday | Thursday | Friday | Saturday

order_hour_of_day  11
Count of order_id  68971

## Total average by week

| Weekday | Total orders |
|---|---|
| Sunday | 156K |
| Monday | 137K |
| Tuesday | 103K |
| Wednesday | 91K |
| Thursday | 91K |
| Friday | 105K |
| Saturday | 114K |

Insight: The highest orders are made around 11AM in an average week. 8AM-6PM is peak hours on average. Sunday has the highest sales while Wednesday has the lowest. Improving sales on middle of weekdays can increase business growth.

**Customers with more than 1000 orders**

user_id 164055
Count of order_number 3060

Insight: Customer_number = 164055 contributes the most accounting 3060 orders. Offering sales discounts, freebies for other customers can increase sales orders.



**Top 10 Reordered products by Department**

Reorders by product_name

| Department |
|---|
| Select all |
| alcohol |
| babies |
| bakery |
| beverages |
| breakfast |
| bulk |
| canned goods |
| dairy eggs |
| deli |
| dry goods pasta |
| frozen |
| household |
| international |
| meat seafood |
| missing |
| other |
| pantry |
| personal care |
| pets |
| produce |
| snacks |

| Product | Reorders |
|---|---|
| Banana | 9.8K |
| Bag of Organic Bananas | 8.4K |
| Organic Strawberries | 5.0K |
| Organic Baby Spinach | 5.0K |
| Organic Hass Avocado | 4.3K |
| Organic Avocado | 3.5K |
| Organic Raspberries | 2.8K |
| Large Lemon | 2.8K |
| Organic Whole Milk | 2.7K |
| Strawberries | 2.6K |

Reorders
Total product orders: 3.11K — 7.05K — 10.99K

Insight: Banana has the highest reorders accounting 9800 orders.

Top 10 Reordered products by Department

Insight: Selecting by department slicer, pantry for example provides corresponding sales orders.

## Conclusion:

The sales orders can be increased by focusing on beverages and frozen department. Irregular customers can avail more discounts and o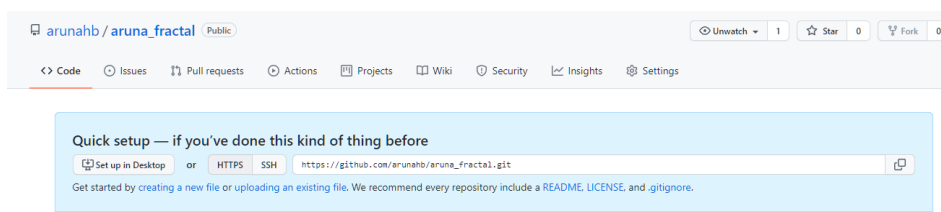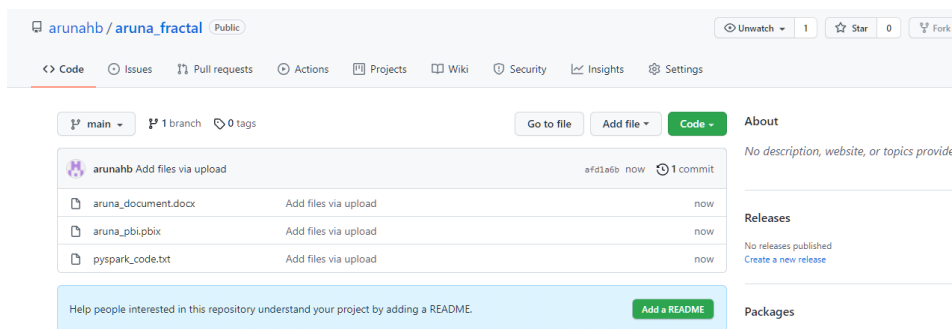ffers than regular customers, on less peak hours on least sales day of the week (Wednesday). This can maintain business growth throughout.

## Github repository:

a. Repository creation:



b. Upload supporting files:



Github_link: https://github.com/arunahb/aruna_fractal

**Files submitted:**

aruna_document, aruna_pbi, pyspark_code, aruna_ouput.csv