# Process followed for the analysis of fetch AE assignment dataset
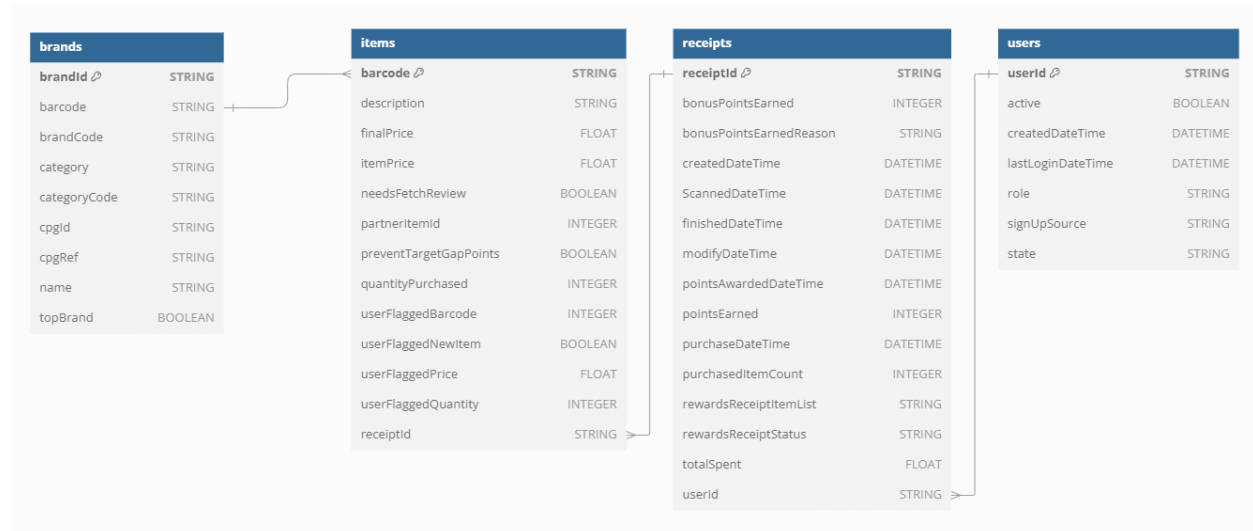
At the beginning , I have looked at the data set and realized these are mongodb json files with three datasets , after close examination , I found that receipts has nested data inside , so I thought this could be divided into 2 making the over all count to 4 datasets, which are

- Users
- Brands
- Receipts
- Items

I decided to use python to convert the json files to csv files since GCP BIG QUERY can import these files with ease, additionally I have use python to do a basic data profiling task using the library ProfileReport from ydata-profiling library, This profilereport creates html files about the profile of the dataset , this report includes dataset statistics, Distribution of values , missing values etc., it is used to determines the structure of the date , detect patterns.

once I have the overall idea of how data looked like , I have exported them converted csv files from json to BIGQUERY

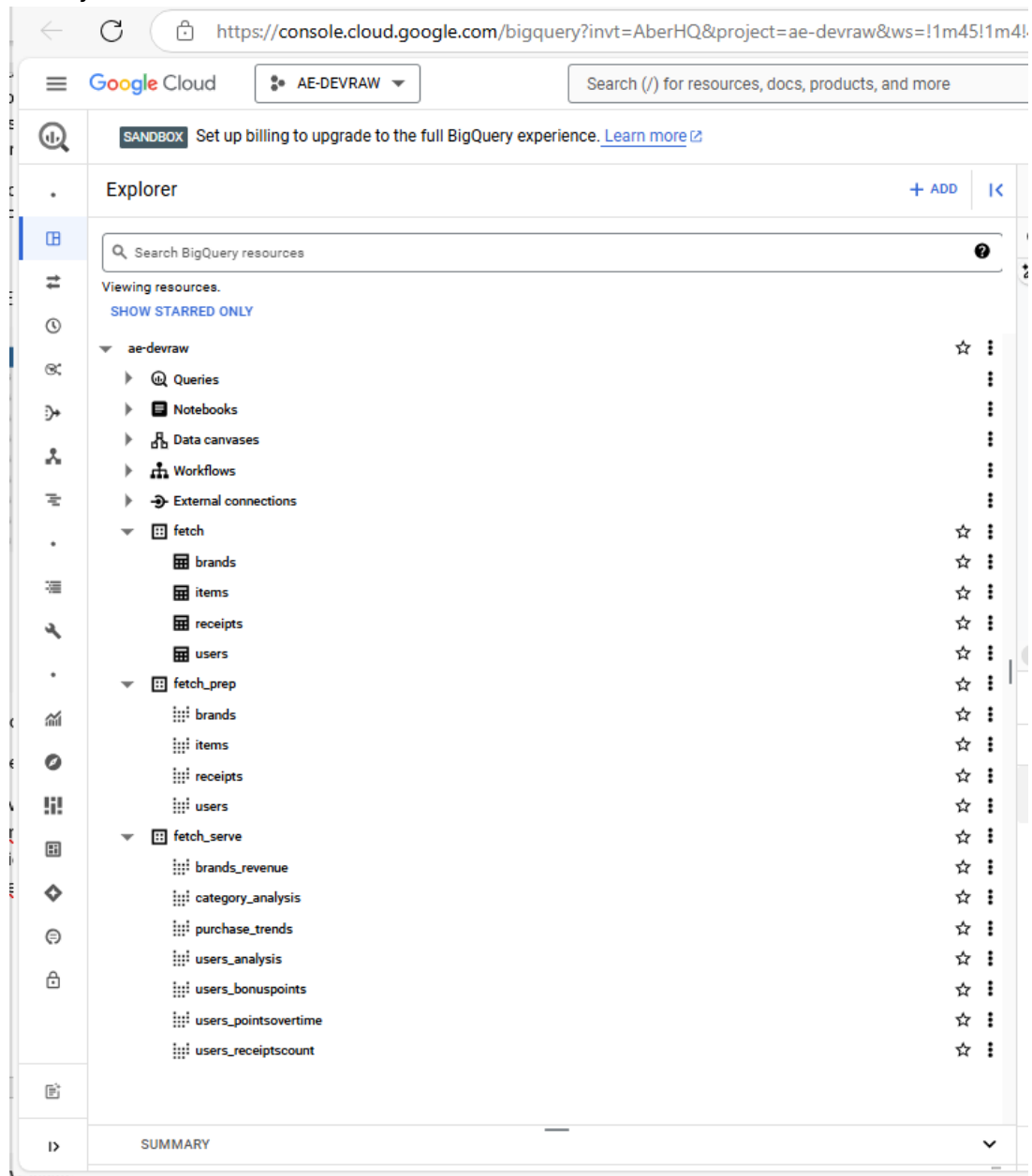Here is the ERD that can be used to represent the data



Here , I thought of following medallion design in the Big query Datalake

So , I have created three datasets

- Fetch ( where data is brought in as is )
- Fetch_prep ( here, the transformation takes place including de duplication , datatypes conversion and column names changes etc )
- Fetch_serve ( once the cleaned data is available in prep , we do the aggregation and analysis in this layer )

The layers can be seen below



Following link should take you to the dataset in Big query
fetch :
https://console.cloud.google.com/bigquery?ws=!1m4!1m3!3m2!1sae-devraw!2sfetch
fetch_prep:
https://console.cloud.google.com/bigquery?ws=!1m4!1m3!3m2!1sae-devraw!2sfetch_prep
fetch_serve:
https://console.cloud.google.com/bigquery?ws=!1m4!1m3!3m2!1sae-devraw!2sfetch_prep

I have attached the data mapping excel file ( named Mapping_Datatypes.xlsx ) to the repository wherever datatype and other changes were done.

I have used the queries in analysis to draw insights in Serve dataset.

User Analysis : SELECT
    u.userId,
    u.active,
    u.createdDateTime,
    u.lastLoginDateTime,
    u.role,
    u.signUpSource,
    u.state,
    COUNT(r.receiptId) AS receiptCount
FROM
    `fetch_prep.users` AS u
LEFT JOIN
    `fetch_prep.receipts` AS r ON u.userId = r.userId
GROUP BY
    u.userId, u.active, u.createdDateTime, u.lastLoginDateTime, u.role,
u.signUpSource, u.state
ORDER BY
    receiptCount DESC;



2 )Brand Revenue

```
1   SELECT
2       b.name AS brandName,
3       SUM(r.totalSpent) AS totalRevenue
4   FROM
5       `fetch_prep.brands` b
6   JOIN
7       `fetch_prep.items` i ON b.barcode = i.barcode
8   JOIN
9       `fetch_prep.receipts` r ON i.receiptId = r.receiptId
10  GROUP BY
11      brandName
12  ORDER BY
13      totalRevenue DESC;
14
```

Query results

| Row | brandName | totalRevenue |
|-----|-----------|--------------|
| 1 | Pepperidge Farm | 23298.18999999... |
| 2 | Diet Chris Cola | 20511.16 |
| 3 | Prego | 20511.16 |
| 4 | V8 | 18576.23999999... |
| 5 | Tostitos | 15799.37 |
| 6 | Cracker Barrel Cheese | 15509.03999999... |
| 7 | Cheetos | 13854.28999999... |
| 8 | Jell-O | 9320.539999999... |
| 9 | Swanson | 7187.139999999... |
| 10 | Quaker | 5781.690000000... |
| 11 | Mountain Dew | 4566.17 |
| 12 | Rice A Roni | 3071.4 |
| 13 | Kraft | 2484.46 |
| 14 | Kettle Brand | 2400.91 |

## 3) Category Analysis



```
1   SELECT
2       b.category,
3       SUM(r.totalSpent) AS totalSpent
4   FROM
5       `fetch_prep.receipts` r
6   JOIN
7       `fetch_prep.items` i ON r.receiptId = i.receiptId
8   JOIN
9       `fetch_prep.brands` b ON i.barcode = b.barcode
10  GROUP BY
11      b.category
12  ORDER BY
13      totalSpent DESC
14
15  ;
```

Query results

| Row | category | totalSpent |
|-----|----------|------------|
| 1 | Snacks | 32054.57000000... |
| 2 | Grocery | 27031.0 |
| 3 | Beverages | 23142.41 |
| 4 | Condiments & Sauces | 21254.95 |
| 5 | null | 20511.16 |
| 6 | Dairy | 17993.5 |
| 7 | Baking | 9320.539999999... |
| 8 | Canned Goods & Soups | 7187.139999999... |
| 9 | Breakfast & Cereal | 5781.690000000... |
| 10 | Dairy & Refrigerated | 944.15 |

Job history

## 4) Purchase Trends over time

⊕ **purchase_trends**   ▶ RUN   💾 SAVE VIEW ▾   OPEN IN ▾   ⚙ MORE ▾                                ✔ Query completed.

```
 1  SELECT
 2      FORMAT_TIMESTAMP('%Y-%m', r.purchaseDateTime) AS purchaseMonth,
 3      COUNT(r.receiptId) AS totalPurchases,
 4      SUM(r.totalSpent) AS totalRevenue,
 5      SUM(r.pointsEarned) as totalpointsEarned
 6  FROM
 7      `fetch_prep.receipts` r
 8  GROUP BY
 9      purchaseMonth
10  ORDER BY
11      purchaseMonth;
12
```
Press Alt+F1 for Accessibility Options.

**Query results**                                        ⬇ SAVE RESULTS ▾   📊 EXPLORE DATA ▾   ⟳

JOB INFORMATION | **RESULTS** | CHART | JSON | EXECUTION DETAILS | EXECUTION GRAPH

| Row | purchaseMonth ▾ | totalPurchases ▾ | totalRevenue ▾ | totalpointsEarned ▾ |
|---|---|---|---|---|
| 1 | *null* | 448 | 0.16 | *null* |
| 2 | 2017-10 | 9 | 27.0 | *null* |
| 3 | 2019-01 | 1 | 0.0 | *null* |
| 4 | 2020-08 | 40 | 1398.400000000... | 1000 |
| 5 | 2020-09 | 1 | 34.96 | 750 |
| 6 | 2020-10 | 2 | 15.0 | 1139 |
| 7 | 2020-11 | 6 | 322.0 | 11984 |
| 8 | 2020-12 | 25 | 602.48 | 5230 |
| 9 | 2021-01 | 498 | 47578.60999999... | 269022 |
| 10 | 2021-02 | 87 | 3232.440000000... | 67207 |
| 11 | 2021-03 | 2 | 2.0 | 525 |

Results per page:  50 ▾   1 – 11 of 11   |<  <  >  >|

5)

## User bonus points

⊕ **users_bonuspoints**   ▶ RUN   💾 SAVE VIEW ▾   OPEN IN ▾   ⚙ MORE ▾                            ✔ Query completed.

```
 1  SELECT
 2      u.userId,
 3      SUM(r.bonusPointsEarned) AS totalBonusPoints
 4  FROM
 5      `fetch_prep.users` AS u
 6  JOIN
 7      `fetch_prep.receipts` AS r ON u.userId = r.userId
 8  GROUP BY
 9      u.userId
10  ORDER BY
11      totalBonusPoints DESC;
12
```
Press Alt+F1 for Accessibility Options.

**Query results**                                        ⬇ SAVE RESULTS ▾   📊 EXPLORE DATA ▾   ⟳

JOB INFORMATION | **RESULTS** | CHART | JSON | EXECUTION DETAILS | EXECUTION GRAPH

| Row | userId ▾ | totalBonusPoints ▾ |
|---|---|---|
| 1 | 5fbc35711d967d1222cbfefc | 2500 |
| 2 | 600fb1ac73c60b12049027bb | 2147 |
| 3 | 5ff1e194b6a9d73a3a9f1052 | 2040 |
| 4 | 6000b75bbe5fc96dfee1d4d3 | 1800 |
| 5 | 5fb0a078be5fc9775c1f3945 | 1500 |
| 6 | 6000d46cfb296c121a81b20c | 1500 |
| 7 | 5ff5d15aeb7c7d12096d91a2 | 1487 |
| 8 | 6010bddaa4b74c120bd19dfb | 1480 |
| 9 | 6009e60450b3311194385009 | 1450 |
| 10 | 600987d77d983a11f63cfa92 | 1420 |
| 11 | 5ffc8f9704929111f6e922bf | 1305 |
| 12 | 6008f02fb6310511daa4f314 | 1250 |
| 13 | 60099c1450b33111fd61f702 | 1250 |
| 14 | 5ff1e1eacfcf6c399c274ae6 | 1220 |

Results per page:  50 ▾   1 – 50 of 141   |<  <  >  >|

Further , these can be used in a Visualization platform like Power BI , Tableau and looker to show more insightful data , this attempt is just showing the overall process outline we usually follow.

Thank you very for your consideration !!!