



Second: Write a query that directly answers a predetermined question from a business stakeholder

We would be using 'sqlite' throughout as an engine where we make use of sql_alchemy library to perform sql queries along with python/on dataframes.

```
In [481]: engine = create_engine('sqlite://', echo=False)
receipts_final.columns = receipts_final.columns.str.replace(".", "_")

<ipython-input-481-0-dc88936c079d>: FutureWarning: The default value of regex will change from True to False in a future version. In addition, single character regular expressions will "not" be treated as literal strings when regex=True.
receipts_final.columns = receipts_final.columns.str.replace(".", "_")

In [482]: receipts_final.to_sql('receipts', con=engine)
brands_final.to_sql('brands', con=engine)
```

- a. When considering average spend from receipts with 'rewardsReceiptStatus' of 'Accepted' or 'Rejected', which is greater?
- b. When considering total number of items purchased from receipts with 'rewardsReceiptStatus' of 'Accepted' or 'Rejected', which is greater?
- To keep the code efficient, I have performed both of the queries in one single to reduce the redundancy and keep it clean.
- The below code shows counts of unique values present in the rewardsReceiptStatus column and we can see that theres no status as "ACCEPTED" in the original dataset itself. Based on the data we acquired, we can calculate for ACCEPTED, assuming theres no more status apart from the below mentioned.
- ACCEPTED = FINISHED - REJECTED assuming that receipts are flagged before being processed and final results are only accepted/rejected.

```
In [483]: receipts_data['rewardsReceiptStatus'].value_counts()

Out[483]: FINISHED    5920
          FLAGGED     810
          SUBMITTED    434
          REJECTED     167
          PENDING       50
          Name: rewardsReceiptStatus, dtype: int64

In [484]: x = engine.execute("SELECT rewardsReceiptStatus as STATUS, avg(totalSpent) as AVG_SPENT, count(purchasedItemCo
x

Out[484]: [(('FINISHED', 1244.372934121628, 5920),
('FLAGGED', 2635.570247360689, 810),
('SUBMITTED', 28.0324497959184, 0),
('REJECTED', 19.544970059880253, 167),
('PENDING', None, 0))

Here we have calculated the avg_spend and total_count for rejected and accepted status and can notice that status 'ACCEPTED' is greater than 'REJECTED' in both the conditions.
```

```
In [485]: status = pd.DataFrame(x, columns=['STATUS', 'AVG_SPENT', 'TOTAL_COUNT']).set_index("STATUS")
status.loc['ACCEPTED'] = status.apply(lambda x: x['FINISHED']=x['REJECTED'])
status

Out[485]:
```

	AVG_SPENT	TOTAL_COUNT
STATUS		
FINISHED	1244.372934	5920.0
FLAGGED	2635.570247	810.0
PENDING	28.032449	0.0
REJECTED	19.544970	167.0
SUBMITTED	NaN	0.0
ACCEPTED	1224.827964	5753.0

Third: Evaluate Data Quality Issues in the Data Provided

Finding the unique values of brandcodes in the receipts final dataset to compare with brands dataset

```
In [486]: receipts_final = receipts_final.dropna(subset=['rewardsReceiptItemList_brandCode'])
receipts_final = receipts_final.drop_duplicates(subset=['rewardsReceiptItemList_brandCode'])
receipts_final = receipts_final.reset_index(drop=True)
len(receipts_final.index)

Out[486]: 227

In [487]: brands_final = brands_final.drop_duplicates(subset=['brandCode'])
brands_final = brands_final.reset_index(drop=True)
len(brands_final.index)

Out[487]: 540

Copying it into a new dataframe so I dont mess with the original dataset

In [488]: receipts_dropped = receipts_final.copy()
receipts_dropped['rewardsReceiptItemList_brandCode'] = receipts_dropped['rewardsReceiptItemList_brandCode'].app

In [489]: brands_final['barcode'] = brands_final['brandCode'].apply(lambda x: str(x))
brands_final['brandCode'] = brands_final['brandCode'].apply(lambda x: str(x))

Here, we are extracting the brandcodes that are present in the receipts dataset and comparing it with brands dataset. We can notice that there are almost 186 brandcodes present in receipts dataset but not in brands dataset which is data quality issue and will create issues when querying through these two tables/joining them. Also for the company if their brands record doesn't have these items but are still recognised in the receipts for items. Note- theres a high possibility that these brandcodes could be a part of those 'TEST BRANDCODE @number' present in brand codes. This can only be clear when the brands table is updated with the correct brand names.
```

```
In [490]: brand_code_not_included = receipts_dropped[(receipts_dropped['rewardsReceiptItemList_brandCode'].isin(brands_f
#DISPLAYING ONLY TOP 10 TO NOT CLUTTER THE OUTPUT
print(brand_code_not_included[0:10])
print(len(brand_code_not_included.index))
receipts_dropped['rewardsReceiptItemList_brandCode'].is_unique
brand_code_not_included.is_unique

0      MISSION
1      BRAND
2      KRAFT EASY CHEESE
6      WINGSTOP
7      GIBSON
8      BEN AND JERRYS
9      BROWN
11     KLABBERONN
12     HY-VEE
13     LIGHT & FIT GREEK
186
Name: rewardsReceiptItemList_brandCode, dtype: object

Out[490]: True

Out[491]: True
```

Fourth: Communicate with Stakeholders

Hi,

This is Rishabh from the Analytics team who is responsible for processing all the incoming data, convert them into a readable format for the machine and the user to find meaningful insights. In this process, there are times we tumble across many issues related to the data/data source which I will be discussing below for receipts, users and brands dataset.

1. I tried to merge these three datasets together using the unique ids(identifier for every row) present in every table but soon noticed that there are multiple records per id["id"] after cleaning them which negates the use of having "id" column. Receipts column had userId which can be used to relate to the users data however, there isn't any brandId which can be used to relate with brands table. So I would like to know whether the connection between these two are related to each other by any column or not, and if not, then how are we identifying the brandcode for verification and if yes, whats the field(parameter) thats being used to create the connection.

2. There were multiple data quality issues which I found and how I can fix them if I get any additional insights/data -

a. Id, date columns have been recorded with additional key-value pair or in a dictionary format which I did not understand the reason for it when the column/field were good enough to provide enough insight on what that column is used for. Directly storing the values as a string can reduce the processing steps, which indeed will keep the overall process optimised.

b. Missing of brandcodes in brands dataset that are present in receipts dataset questions the relationship between these two table in case theres a need for verification/restocking/keeping count of sales per brand. So, there are additional brandcodes named "TEST BRANDCODE @number" which can be verified to a valid brandcode in the brands table that can solve the issue of the missing/mismatch of brandcodes between these two datasets.

c.While going through the "rewardsReceiptStatus" column, which is basically the status of the receipt through receipt validation and processing, I tried to count the number of occurrences of each status and found the following - FINISHED 5920 FLAGGED 810 SUBMITTED 434 REJECTED 167 PENDING 50

The status "ACCEPTED" is no where to be found in the original dataset however, we have FINISHED status which means that the processing of these receipts for rewards has finished status. I had to make assumptions that there are no other status and only accepted/rejected status the only ones after receipt is processed completely. Accepted = Finished - Rejected. This thing could be made straightforward by adding another status as ACCEPTED to avoid us creating assumptions or invalid calculations that can create issues.

3.The consistency in data types is a very important aspect of data analysis especially while collecting and storing the data. There can be multiple production issues if there are not proper try/catch method - this is way to handle errors, foreign data type, data structure without hindering the overall process. There are times when these are not properly handled can create server issues/server stoppages. One other thing that can be added to the tables is a way to connect different tables using primary and foreign keys - which can be done during the processing and storing the data, however, proper data dictionary of what the values mean inside specific columns, such as rewards status could make things more clear.

Looking forward to discussing the above in detail.

Regards, Rishabh