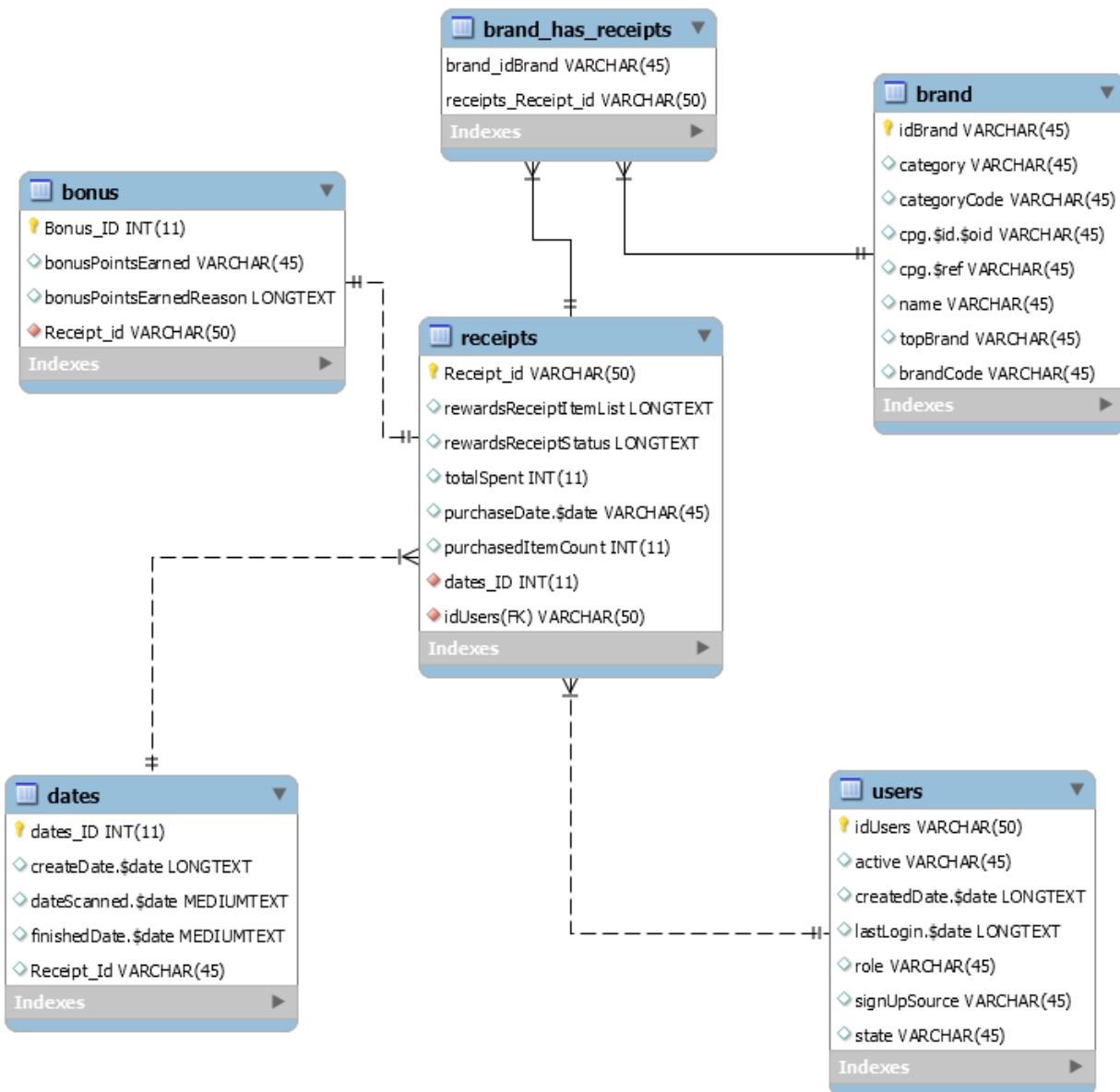


Fetch Rewards Analytics Engineer Assessment

Before moving ahead, I have converted all JSON files to CSV for convenience and to import in MySQL workbench

First: Review Existing Unstructured Data and Diagram a New Structured Relational Data Model



I have created ER model with the help of MySQL workbench using ETL queries and then reverse engineering it in the platform.

To make it simpler, I have normalized the data and divided receipts CSV into three tables i.e receipt, date and bonus.

Relationships:

Each user will have multiple receipts

Each date will have multiple receipts

Each receipt is unique

Brand will have multiple receipts, or many receipts will have multiple brands

Each receipt will have one bonus(considering you get only one bonus point after uploading the receipt)

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

Write a SQL query against your new structured relational data model that answers one of the following bullet points below of your choosing.

Even though one question was expected I tried to answer all of the questions mentioned

- 1) When considering average spend from receipts with 'rewardsReceiptStatus' of 'Accepted' or 'Rejected', which is greater?

In this query I got Finished as greater but if you can run it in the systems you will get answer according to the data that loads.

```
5  -- When considering average spend from receipts with 'rewardsReceiptStatus' of 'Accepted' or 'Rejected'
6  select
7  case
8      when (avg(a.totalSpent)) > (avg(b.totalSpent)) Then 'FINISHED'
9      when (avg(b.totalSpent)) > (avg(a.totalSpent)) Then 'REJECTED'
10     else 'Equal'
11  end AS Compare
12  from receipts a cross join receipts b
13  where a.rewardsReceiptStatus = 'FINISHED' and b.rewardsReceiptStatus = 'REJECTED'
14  group by a.rewardsReceiptStatus,b.rewardsReceiptStatus;
15
```



Compare
FINISHED

- 2) What are the top 5 brands by receipts scanned for most recent month?

I have included auto increment date id in both table... you can use the same query without join if you do not decide to separate Receipts in three parts

```

16
17 -- What are the top 5 brands by receipts scanned for most recent month?
18 select brand.name as Brand, Count(receipts.Receipt_id) as Total Receipt from brand
19 left join receipts on brand.Receipt_id = receipts.Receipt_id
20 inner join dates on receipts.dates_ID = dates.dates_ID -- assuming dates and receipts have common column
21 group by brand.name
22 having datepart(mm,date.dateScanned.$date) = month(getdate())
23 order by Total Receipt desc limit 5;
24
25 -- How does the ranking of the top 5 brands by receipts scanned for the recent month compare to the ranking for the
26
27 select *, LAG(Count(receipts.Receipt_id)) OVER (PARTITION BY brand.name ORDER BY month(getdate()))
28 as previousmonth
29 FROM brand

```

idBrand	category	categoryCode	cpq.\$id.\$oid	cpq.\$ref	name	topBrand	brandCode	idUsers
54943462e4b07e684157a532	true	1418996882381	1614963143204	fetch-staff	----	----	NIUS	NIUS
55308179e4b0eabd8f99caa2	true	1429242233186	1525713820003	consumer	----	WI	NIUS	NIUS
60544ab713a2b713a64b7c7f767b	try sa	1400785001771	1614084860739	fetch-staff	----	TI	NIUS	NIUS

- 3) Which brand has the most spend among users who were created within the past 6 months?

I have included auto increment date id in both table... you can use the same query without join if you do not decide to separate Receipts in three parts

I have used curdate() function to get the current date

```

-- Which brand has the most spend among users who were created within the past 6 months?
select brand.idBrand as BrandID, brand.name as Brand, sum(receipts.totalSpent) as TotalSpent, date.createDate.$date as Date from brand
inner join receipts on brand.Receipt_id = receipts.Receipt_id
inner join dates on receipts.dates_ID = dates.dates_ID -- assuming dates and receipts have common column
where Date > curdate() - interval (dayofmonth(curdate()) - 1) day - interval 6 month
group by BrandID, Brand
order by TotalSpent desc limit 1;

```

- 4) Which brand has the most transactions among users who were created within the past 6 months?

```

39
40 -- Which brand has the most transactions among users who were created within the past 6 months?
41 select brand.idBrand as BrandID, brand.name as Brand, count(date.dateScanned.$date) as No_of_Transaction from brand
42 inner join receipts on brand.Receipt_id = receipts.Receipt_id
43 inner join dates on receipts.dates_ID = dates.dates_ID -- assuming dates and receipts have common column
44 where date.dateScanned.$date > curdate() - interval (dayofmonth(curdate()) - 1) day - interval 6 month
45 group by BrandID, Brand
46 order by No_of_Transaction desc limit 1;
47

```

- 5) When considering total number of items purchased from receipts with 'rewardsReceiptStatus' of 'Accepted' or 'Rejected', which is greater?

I have included count because Receipt item list does not contain numbers they have list of items. If you have number, you may use sum.

```

15
16 -- When considering total number of items purchased from receipts with 'rewardsReceiptStatus' of 'Accepted' or 'Rejected', which is greater?
17
18 select
19 *
20 case
21 when (count(a.rewardsReceiptItemList)) > (count(b.rewardsReceiptItemList)) Then 'FINISHED'
22 when (count(b.rewardsReceiptItemList)) > (count(a.rewardsReceiptItemList)) Then 'REJECTED'
23 else 'Equal'
24 end AS Compare
25 from receipts a cross join receipts b
26 where a.rewardsReceiptStatus = 'FINISHED' and b.rewardsReceiptStatus = 'REJECTED'
27 group by a.rewardsReceiptStatus, b.rewardsReceiptStatus;
28
29 -- What are the top 5 brands by receipts scanned for most recent month?
30 select brand.name as Brand, Count(receipts.Receipt_id) as Total_Receipt from brand
31 left join receipts on brand.Receipt_id = receipts.Receipt_id

```

Third: Evaluate Data Quality Issues in the Data Provided:

I have decided to quick quality check using SQL. One can do Profiling with the help of Python to get more in-depth idea.

- 1) All ids such as Brand id, Receipt id is not in one format. Because of this reason I have decided to include to include VARCHAR() while loading the data.
- 2) All Dates are not in one format. We cannot convert it into DDMMYYYY or any other format

use fetchrewards

```
select * from receipts;
```

Receipt_id	rewardsReceiptItemList	rewardsReceiptStatus	totalSpent	purchaseDate.\$date	purchasedItemCount	dates_ID
5F9c74f70a7214ad07000037	[{"barcode": "075925306254", "competitiveProd...	REJECTED	1	1604002679000	11	2021
5F9c74f90a7214ad07000038	[{"barcode": "075925306254", "competitiveProd...	FINISHED	14	1604002681000	6	2021
5Fa5ad370a720f05ef000089	[{"barcode": "075925306254", "competitiveProd...	FINISHED	291	1604006903000	11	2021
5Fa5b0ca0a720f05ef0000bf	[{"barcode": "075925306254", "competitiveProd...	FINISHED	14	1604521418000	6	2021
5Fffe1960a720f0523000567	[{"barcode": "4011", "description": "ITEM NOT FO...	FINISHED	1	1609601046000	1	2021
5Fffe1a10a720f0523000568	[{"barcode": "013962300631", "description": "An...	FINISHED	50	1609545600000	5	2021
5Fffe1a40a720f0523000569	[{"barcode": "046000832517", "brandCode": "BR...	FINISHED	10	1609027200000	1	2021
5Fffe1b20a7214ada1000055a	[{"barcode": "4011", "description": "ITEM NOT FO...	FINISHED	1	1612365875000	1	2021
5Fffe1b40a7214ada1000055b	[{"barcode": "075925306254", "competitiveProd...	FLAGGED	1	1609601076000	1	2021
5Fffe1b60a7214ada1000055c	[{"barcode": "034300573065", "description": "MIL...	FLAGGED	290	1612365878000	10	2021

- 3) There are many "" values instead of NULL values this indicate it does not contain anything. We can replace these values with "NULL" in python.

```

1 use fetchrewards
2
3 select * from brand;

```

	idBrand	category	categoryCode	cpg.\$id.\$oid	cpg.\$ref	name	topBrand	brandCode
▶	54943462e4b07e684157a532	true	1418998882381	1614963143204	fetch-staff	****	****	NULL
	55308179e4b0eabd8f99caa2	true	1429242233186	1525713820003	consumer	****	WI	NULL
	5964eb07e4b03efd0c0f267b	true	1499785991771	1614884869770	fetch-staff	****	IL	NULL
	59c124bae4b0299e55b0f330	true	1505830074302	1612802578117	fetch-staff	****	WI	NULL
	5a43c08fe4b014fd6b6a0612	true	1514389647059	1613146957155	consumer	****	****	NULL
	5e27526d0bdb6a138c32b556	true	1579635309795	****	consumer	Google	WI	NULL
	5f2068904928021530f8fc34	true	1595959440905	1612452605375	fetch-staff	Email	WI	NULL
	5fa32b4d898c7a11a6bcebce	true	1604528973309	1614842518047	fetch-staff	Google	AL	NULL
	5fa41775898c7a11a6bcef3e	true	1604589429396	1614873722026	fetch-staff	Email	****	NULL
	5fb0a078be5fc9775c1f3945	true	1605410936818	****	consumer	Google	AL	NULL

Fourth: Communicate with Stakeholders

Dear XYZ

Hope you are doing well!

I am writing this email regarding the recent data task. I have created data model and used ETL queries to answer business questions mentioned in the meeting. During the task, I observed a few irregularities in the data and would like to put light on them. They are as follows:

- 1) While importing the data in SQL, I observed that there are many irregularities in data. I wanted to make sure the source of the data. We have used JSON files to load but are there any other sources apart from that? We would like to check authenticity of date format and null values with different sources.
- 2) I would like to know if there are more CSVs apart from these three. It would be great to get more information to solve the business problems.
- 3) How would you like to view the data? Do they want to see it in a tabular format, interactive graphs and charts, and trending tables?
- 4) Are there any new KPIs dimension to view more detailed information that we can focus?

We can jump on the call to discuss this further. I would like to schedule 1 hr with the team. Please let me know your availability so we can connect!

Thanking you in advance!

Regards,

Chinmay Arolkar

Please refer to SQL files for the code.