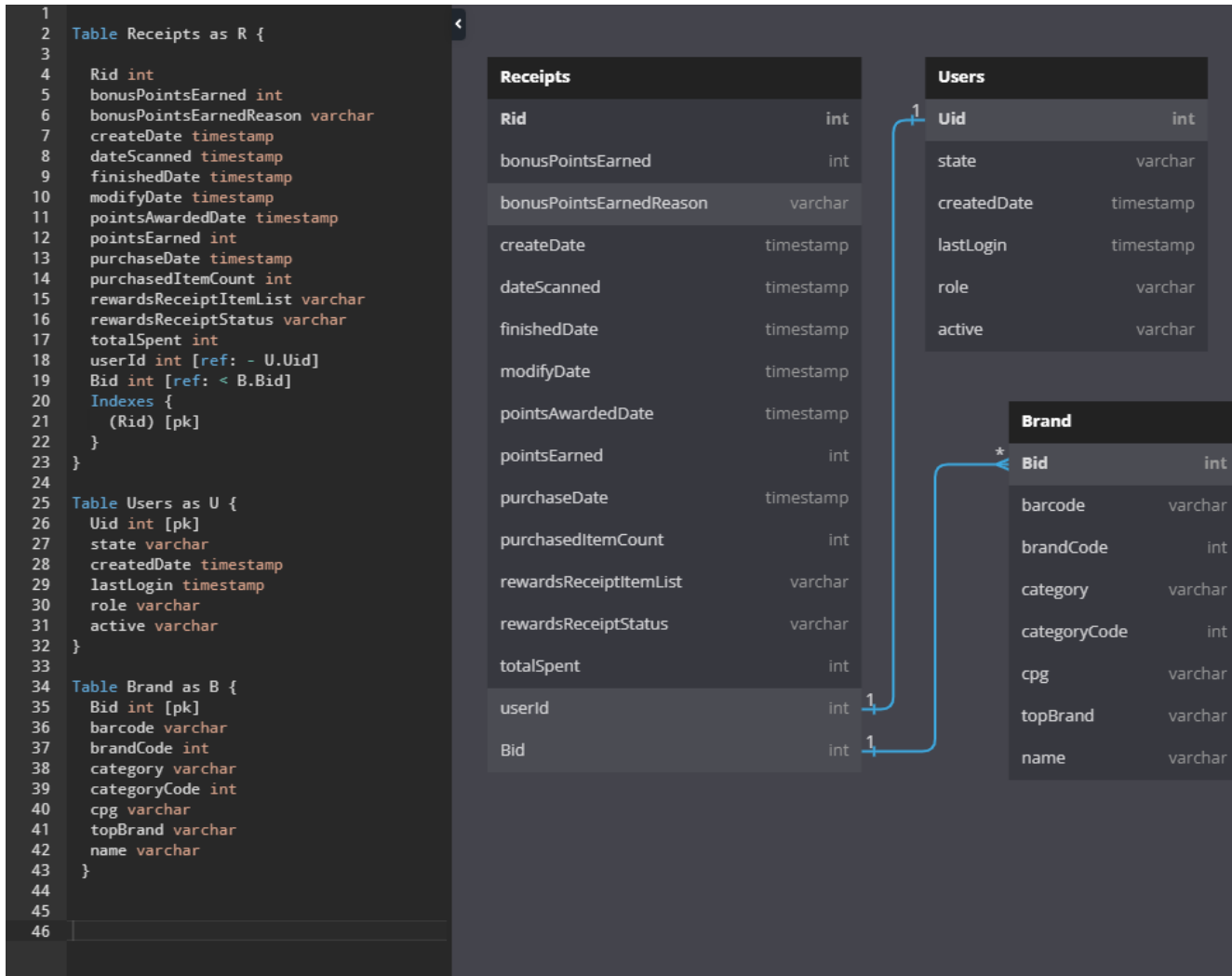


Yuhan (Caroline) Gao

11/18/2022

## Fetch Reward DA Challenge

1. Relational Schema between 3 tables: 'Receipts', 'Users', and 'Brand', by using 'dbdiagram.io'.



## 2. SQL Queries for question 1, 5 and 6.

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

```
SELECT TOP 5
R.BID
,B.[NAME]
,SUM(R. TOTALSPENT) AS [TOP SPEND AMOUNT]
FROM RECEIPTS R
INNER JOIN BRAND B
ON B.BID =R.BID
WHERE MONTH(R.DATESCANNED) = MONTH(GETDATE())
GROUP BY
R.BID
,B.[NAME]
ORDER BY [TOP SPEND AMOUNT] DESC
```

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

```
SELECT
MAX(R.TOTALSPENT) AS [SPEND AMOUNT]
,R.BID
,B.[NAME]
FROM RECEIPTS R
INNER JOIN USERS U
ON U.[UID] = R.USERID
INNER JOIN BRAND B
ON B.BID=R.BID
WHERE DATEDIFF(MONTH,U.CREATEDDATE, GETDATE())<=6
GROUP BY
R.BID
```

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

```
SELECT
MAX(C.[TRANSACTION AMOUNT])
,C.BRANDNAME
FROM (
SELECT
COUNT(DISTINCT R.RID) AS [TRANSACTION AOMUNT]
,R.BID AS [BRNADID]
,B.[NAME] AS [BRANDNAME]
FROM RECEIPTS R
INNER JOIN USERS U
ON U.[UID] = R.USERID
INNER JOIN BRAND B
ON B.BID=R.BID
WHERE DATEDIFF(MONTH,U.CREATEDDATE, GETDATE())<=6
)C
GROUP BY
C.BRANDNAME
```

---

## 3. Evaluate Data Quality Issues

By using R and Python.

import pandas as pd

from pathlib import Path

import json

```
p = Path(r'C:\Users\Carol\OneDrive\Desktop\receipts.json')
```

```
# read json file
```

```
with p.open('r', encoding='utf-8') as f:
```

```
    data = json.loads(f.read())
```

```
#dataframe
```

```
df = pd.json_normalize(data)
```

```
1  import pandas as pd
2
3  from pathlib import Path
4  import json
5
6  p = Path(r'C:\Users\Carol\OneDrive\Desktop\receipts.json')
7  # read json file
8  with p.open('r', encoding='utf-8') as f:
9      data = json.loads(f.read())
10
11 #dataframe
12 df = pd.json_normalize(data)
13 #data quality issue: exist 'extra data' error: means the input JSON file has more than one object
14 #per line. In general, there would be only one object per line.
15
16
17 import pandas as pd
18 from pathlib import Path
19 import json
20
21 # set path to file
22 p = Path(r'C:\Users\Carol\OneDrive\Desktop\brands.json')
23
24 # read json
25 with p.open('r', encoding='utf-8') as f:
26     data = json.loads(f.read())
27
28 # create dataframe
29 df = pd.json_normalize(data)
30
31
32 #####
33 import json
34
35 with open('brands.json', 'r') as openfile:
36     json_object = json.load(openfile)
37
38 print(json_object)
39 print(type(json_object))
40
```

Data quality issue (1): exist 'extra data' error, means the input JSON file has more than one object per line. In general, there would be only one object per line.

Issue (2): huge amount of unstructured & malformed data and exist duplicate data.

Issue (3): need to replace newlines and wrap the whole file in a pair of { }, since the json file has a large amount of statements separated by newlines instead of a single unified statement.

Issue (4): For the 'receipts' data and the 'brand' data, there is no connection between them, which are unable to join together. To solve this problem, we can add the brand uuid into receipts data, and therefore, we can use inner join and other join functions.

```
1 install.packages("rjson")
2 install.packages("jsonlite")
3 install.packages("RJSONIO")
4 library(jsonlite)
5 library("rjson")
6 json_file <- "C:/Users/Carol/OneDrive/Desktop/receipts.json"
7
8 result <- fromJSON(file=json_file)
9
10
11 print(result)
12
13
14 #####
15 library(RJSONIO)
16
17 D2 <- RJSONIO::fromJSON("C:/Users/Carol/OneDrive/Desktop/receipts.json")
18
19 # convert the numeric vector helpful to one string
20 D2$helpful <- paste(D2$helpful, collapse = " ")
21
22 D2
23 reviewerID      asin      reviewerName      helpful
24 [1,] "A3TS466QBAWB9D" "0014072149" "Silver Pencil" "0 0"
25
26 D3 <- do.call(cbind, D2)
27
28 write.csv(D3, "D3.csv")
29 #####2222
30 json_file2 <- "C:/Users/Carol/OneDrive/Desktop/brands.json"
31
32
33 out <- lapply(readLines("C:/Users/Carol/OneDrive/Desktop/brands.json"), fromJSON)
34 l <- replicate(
35   132,
36   as.list(sample(letters, 20)),
37   simplify = FALSE
38 )
39 df <- data.frame(matrix(unlist(l), nrow=length(l), byrow=TRUE))|
40
```

#### 4. Communication-Email Format

Dear XXX manager,

Hope you are doing well!

There were some data quality issues existed in all the three json files after I put them in detail analysis. All of them needed to be converted to csv format by using R, Python, or other programming languages, in order to use SQL or other data analytical tools to do better & further analysis such as find which customers purchased the most in a specific period, etc. The three datasets needed to be modified by many ways, such as replace newlines and wrap the whole file in a pair of {}, since the json file has a large amount of statements separated by newlines instead of a single unified statement; or add the brand uuid into receipts data, and therefore, we can use inner join and other join functions to better determine which specific product need to be promoted more and attract more sales, etc.

Best,  
Caroline