

EBUS3030 Assignment 1

Stavros Karmaniolos c3160280@uon.edu.au

Jay Rovacsek c3146220@uon.edu.au

Jacob Litherland c3263482@uon.edu.au

Edward Lonsdale c3252144@uon.edu.au

September 5, 2018

Contents

| | | |
|----------|---|-----------|
| 1 | Assignment Overview & Requirements | 2 |
| 1.1 | Datamart Business Rules | 4 |
| 2 | Data Model | 5 |
| 3 | Data Load Process (ETL/ELT) | 7 |
| 3.1 | Quality Assurance Processes | 8 |
| 3.2 | Assumptions and Reasoning | 12 |
| 3.2.1 | Item Table | 12 |
| 3.2.2 | Price Table | 12 |
| 3.2.3 | ReceiptItem | 12 |
| 3.2.4 | Receipt | 12 |
| 3.2.5 | Staff | 12 |
| 3.2.6 | Customer | 12 |
| 4 | Base Analysis | 13 |
| 4.1 | Raw Results | 13 |
| 4.1.1 | Total Number of Sales | 13 |
| 4.1.2 | Total Items Sold | 13 |
| 4.1.3 | Discounted Sales Ratio | 14 |
| 4.1.4 | Total Sales Value per Staff Member | 16 |
| 4.1.5 | Average Value Per Sale | 17 |
| 4.2 | Notes on Analysis | 18 |
| 5 | Executive Summary | 19 |
| | References | 20 |
| 6 | Appendix | 21 |
| 6.1 | CTE Raw Results | 21 |
| 6.2 | Python Script | 22 |

1 Assignment Overview & Requirements

EBUS3030 – Assignment 1

Business Intelligence - EBUS3030 Assignment 1

Due: Assignment One TurnItIn drop folder by 12 noon on Thursday 6th September
Paper copy at the beginning of week 6 workshop.

Assignment Outcomes

This assignment requires multiple outputs to be created to exhibit your understanding of business intelligence/data analysis through an example ‘real world’ question that is comparable to what you may be asked of you as you become an IT professional.

Key outcomes to be delivered are: Data Modelling of the provided dataset, Extract Transform Load (ETL) processing undertaken to make the data usable, the Output of your analysis, a Report summarising your findings and a presentation to the class of your work. The presentation is expected to concentrate more on your findings/recommendations as if it were a situation where you are presenting the response to the head sales executive’s question.

Assignment Question

The head Sales Executive of ‘BIA Inc’ comes to you as the lead Business/Data Analyst and asks you to help with a problem they have.

“I’ve heard that people aren’t motivated at the moment and sales aren’t as good as we had hoped. To try and provide incentives for staff, I want to provide an award (and probably associated cash prize) to my best performer for sales from this Office, I need you to tell me who that is?”

“As part of your response I want you to provide the justification as to why the particular sales officer was selected because we need governance over things like this.

.... By the way, we don’t currently have any of this information stored centrally in a database thingy, but I have gotten the Office Business Manager to collate a summary of the recent sales into a rough excel file that can be used as a starting basis. As part of the processes of getting me an answer on my best salesperson, can you also create a database as part of the preparation of the answer. We will then use that as the base of further reporting into the future. We haven’t ever had people with your skills working with us before so I expect there will be lots of questions that will come up as we utilise your expertise.”

Assignment Deliverables

Using the data file provided in Excel and associated notes about the data, (*AssOneData.xlsx* and *Datamart Business Notes*) you are required to complete the following elements as part of the assignment.

- Data Model
 - Using the information made available to you and your understanding of concepts around data mart design in the labs, design a “Sales” DataMart to store the information in a format that will allow the information to be expanded and one that would enable analysis to occur.
- Data Load Process undertaken
 - Provide an overview of the ETL/ELT process completed and what (if any) Quality Assurance processes you undertook as part of this.
 - Ensure you record any assumptions you have made as part of this component and your reasoning behind the assumption.
- Output of Analysis (including SQL used)
 - Once the data loaded and is available and ready for use, you need to create a set of sql scripts to be used to generate the results to the business question provided to you from the Head Sales Executive
 - Provide a snapshot of the raw results of your analysis that provides the basis of your recommendations
 - Ensure you record any assumptions you have made as part of this analysis component and your reasoning behind the assumption.

EBUS3030 – Assignment 1

- Executive Summary in response to business question.
 - Provide a short Executive brief/summary that presents a clear concise response back to the Sales Executive's question about possible incentives to the best salesperson. This should clearly detail the recommendation and any key assumptions/restrictions the executive need to be aware of.
- Team Presentation
 - All members of the team need to participate in a (10-15 minute) presentation to be delivered as part of the lab in Week 6. This needs to be presented in a format as if you were summoned to the board room with the Head Sales Executive to provide a formal response to their question.
 - Please be aware that the Head Sales Executive may ask any of the team members questions as you present your analysis.

NB: As part of your responses, you should also specifically include any assumptions you have made throughout the process.

Breakup of assignment Marks (total course mark for assignment = Assignment Part A submission (20% + Presentation One (5%) = 25%).

| Assignment Component | Percentage Allocation |
|----------------------|-----------------------|
| Data Model | 30% |
| ETL | 10% |
| Base Analysis | 30% |
| Executive Summary | 10% |
| Team Presentation | 20% |
| Assumptions | 100% |

Key Documents Required & Format

You are required to upload all files in a single zip file (including any presentation items for the team delivery within the lab) via blackboard to the Assignment One TurnItIn drop folder by 12 noon on Thursday 6th September. You will also be required to submit a paper copy of your deliverables at the workshop (make sure this is printed well before the workshop).

NB: Only 1 load per team only but it should contain all of the deliverable items in a .zip file.

1.1 Datamart Business Rules

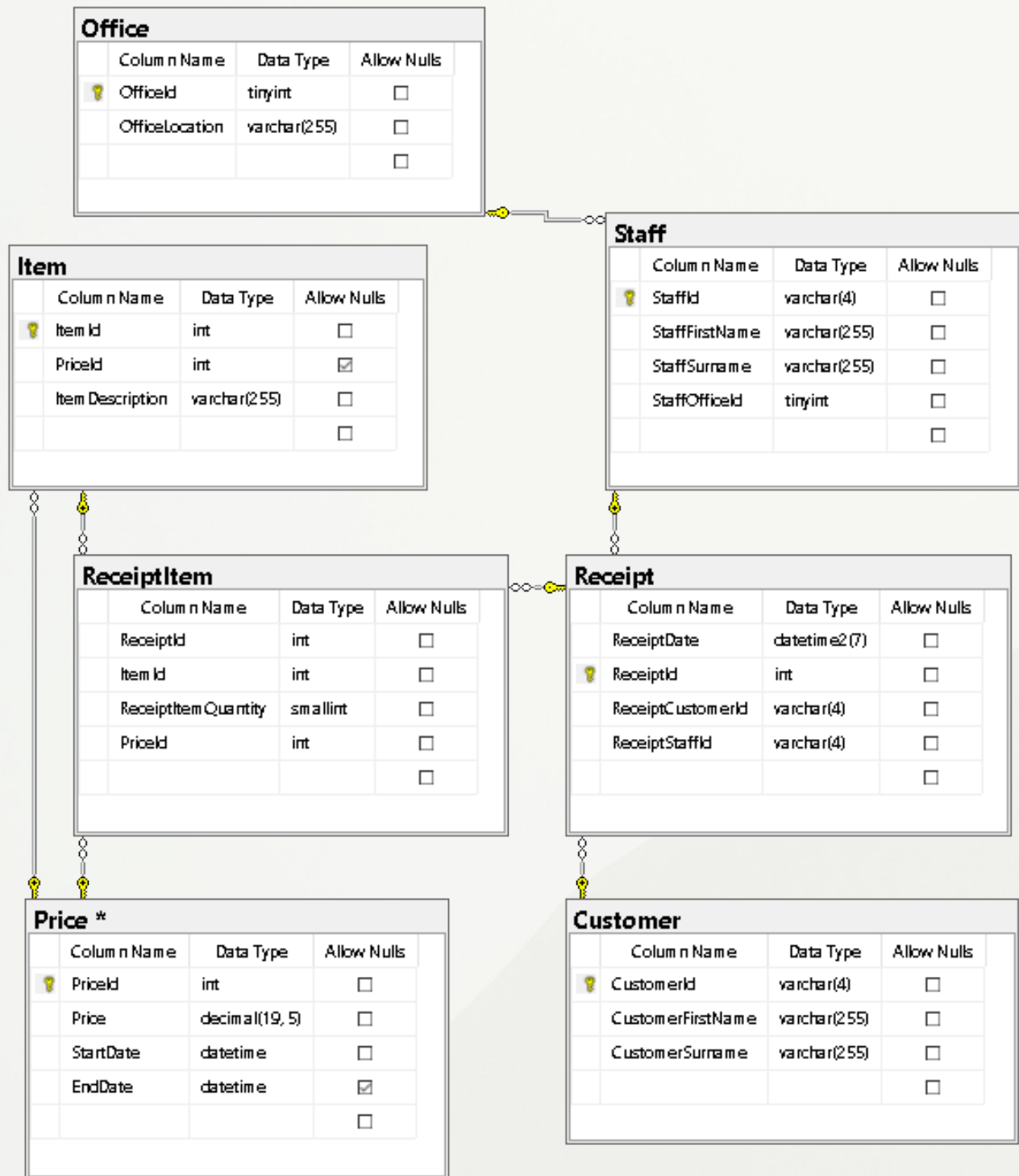
The following business rules were provided to be used in the context of this assignment:

- * At BIA all customers interacts are in an online environment, there are no orders outside of electronic.
- * Returning customers can provide POI information via the web interface and look up their record and that will flow with the sale.
- * The sales associate can complete the order form/sale for the client.
- * Each sale will have a receipt number/id.
- * A receipt can have many line items.
- * Each line item can only be for a single item, but the customer can purchase multiples of the same item.
- * Where a customer has multiple line items, any sale with more than 5 row items (containing at least 5 different items) is provided a 15% discount.
- * The system automatically handles the total for the sale by looking up the item, then multiplying the costs per item by number purchased, and then should store this final field total as a record in the system (but should also be able to see clearly sales that were provided a discount.
- * Item prices can change at any point, and the price the customer pays is the amount listed for the item on the sale date. We need to keep a record of all item prices historically.
- * Only 1 BIA sales assistant can be attributed to any receipt.

With these considerations in mind, the following report was created to outline the discovery, creation and polish to satisfy the assignment requirements.

2 Data Model

The below data model is only a suggestion and is still subject to change into the future. A full create script can be found in the [appendix](#)



It must be noted that the structure of this data model is less than efficient, and it would be expected in a datamart situation that only at lower levels of data would this schema remain responsive in the manner it is now, as the outline suggests the datamart is not necessarily the most suitable design for future use, however suits very well currently.

It would be expected that only at extremely large data sets would this model prove a bad design. In such cases a model more representative of the snowflake or star schema would be heavily advised.

3 Data Load Process (ETL/ELT)

Initial import of the data supplied in the xlsx file generated a very basic table that allowed us to analyze the data for potential outliers, confirm the business requirements of the data and then create tables from which the data model was derived.

The Imported table structure was as follows:

| Assignment1Data | | | |
|-----------------|----------------------------|--------------|--------------------------|
| | Column Name | Data Type | Allow Nulls |
| | Sale_Date | datetime2(7) | <input type="checkbox"/> |
| | Reciept_Id | int | <input type="checkbox"/> |
| | Customer_ID | nvarchar(50) | <input type="checkbox"/> |
| | Customer_First_Name | nvarchar(50) | <input type="checkbox"/> |
| | Customer_Surname | nvarchar(50) | <input type="checkbox"/> |
| | Staff_ID | nvarchar(50) | <input type="checkbox"/> |
| | Staff_First_Name | nvarchar(50) | <input type="checkbox"/> |
| | Staff_Surname | nvarchar(50) | <input type="checkbox"/> |
| | Staff_office | int | <input type="checkbox"/> |
| | Office_Location | nvarchar(50) | <input type="checkbox"/> |
| | Reciept_Transaction_Row_ID | int | <input type="checkbox"/> |
| | Item_ID | int | <input type="checkbox"/> |
| | Item_Description | nvarchar(50) | <input type="checkbox"/> |
| | Item_Quantity | int | <input type="checkbox"/> |
| | Item_Price | float | <input type="checkbox"/> |
| | Row_Total | float | <input type="checkbox"/> |
| | | | <input type="checkbox"/> |

A decision to leave this initial import table as default was made to allow easy reference to the initially supplied excel data file.

In the following sections of [Quality Assurance Processes](#), [Assumptions and Reasoning](#) and [Base Analysis](#) we intend to clarify the reasoning behind leaving the imported data in the default table suggested by SSMS.

3.1 Quality Assurance Processes

A number of queries were written to look for data which did not adhere to the spec outlined in business requirements and to ensure data was "clean" before entry. The first instance of potential issues were encountered with a basic python script which checked validity of column data, it was found that cells starting at B13777 to the end of file in the originally supplied excel file were formula values and not static values, this would not have caused an issue with importing into SSMS however certainly broke the script temporarily.

After clarifying the issues with the aforementioned cells with Peter, a data file without the offending formula was supplied and used for the remainder of the assignment.

Discrepancies with some cell formatting were noted in Reciept_Id column in the raw data presented to us in the excel file. After testing both an unmodified and a modified version of the excel file it was recognised that these cells did not impact the import of data into SSMS. The offending cells in question were: B13776 - B13865.

The next potential issue encountered was not until a suggested schema structure was complete and data was being scripted to be added to the new schema for analysis. The issue encountered was that receipt number 52136 seemed to be an incorrect entry, this was discovered when running the import query for the new schema:

```
1 INSERT INTO Receipt(ReceiptId , ReceiptCustomerId , ReceiptStaffId )
2 SELECT DISTINCT(Reciept_Id) , Customer_ID , Staff_ID
3 FROM Assignment1Data
4 ORDER BY Reciept_Id
```

Which resulted in the error:

```
Violation of PRIMARY KEY constraint 'PK_Receipt'. Cannot insert duplicate key in object
'dbo.Receipt'. The duplicate key value is (52136).
```

Leading us to recognise that either one of the entries could be incorrect, therefore best to investigate both records of the customer Id against the rest of the database:

```
1 SELECT * FROM Assignment1Data
2 WHERE Customer_ID='C32'
3 AND Staff_ID='S15'
4 AND Sale_Date='2017-11-12 00:00:00.0000000' ;
5
6 SELECT * FROM Assignment1Data
7 WHERE Customer_ID='C13'
8 AND Staff_ID='S4'
9 AND Sale_Date='2017-12-30 00:00:00.0000000' ;
```

When both queries were performed it was apparent that the data associated with C32 was the likely broken record and modification of the data occurred:

```
1 UPDATE Assignment1Data
2 SET Reciept_Id=51585,
3 Reciept_Transaction_Row_ID=(
4     SELECT MAX(Reciept_Transaction_Row_ID)+1
5     FROM Assignment1Data
6     WHERE Reciept_Id=51585)
7 WHERE Customer_ID='C32'
8 AND Staff_ID='S15'
9 AND Sale_Date='2017-11-12 00:00:00.0000000'
10 AND Item_ID='14' ;
```

The next issue arose when again, attempting to run the aforementioned query to import into the new Receipt table, this time not one stray record was found, but a complete collision on the ReceiptId of 52137, this time as neither record seemed to have records that were correct, it was decided to move one to the maximum ReceiptId + 1:

```
1 UPDATE Assignment1Data
2 SET Reciept_Id=(
3     SELECT MAX(Reciept_Id)+1
```

```

4      FROM Assignment1Data)
5 WHERE Customer_ID='C27'
6 AND Staff_ID='S4'
7 AND Sale_Date='2017-12-30 00:00:00.0000000';

```

The same issue was replicated on ReceiptId 52138, resolved via:

```

1 UPDATE Assignment1Data
2 SET Reciept_Id=(
3     SELECT MAX(Reciept_Id)+1
4     FROM Assignment1Data)
5 WHERE Customer_ID='C30'
6 AND Staff_ID='S19'
7 AND Sale_Date='2017-05-16 00:00:00.0000000';

```

At this point we recognised the broken data likely continued for a while, and evaluated our hypothesis by looking at the original excel file. It turned out that data with ReceiptId from 52137-52145 was all broken in the same manner. The following query shows this well:

```

1 SELECT Reciept_Id , Customer_ID , Staff_ID
2 FROM Assignment1Data
3 WHERE Reciept_Id BETWEEN 52137 AND 52150
4 GROUP BY Reciept_Id , Customer_ID , Staff_ID
5 ORDER BY Reciept_Id;

```

In order to clean this data we looked at a number of potential methods, with an emphasis on avoiding effort in the task if possible but not breaking the data further, which to this point just appeared to be a collision of a number of receipts.

We knew a structure such as a CTE [3] would allow us to easily split distinct records which shared a receiptId and filter by a value such as row number.

```

1 WITH CTE AS
2 (
3     SELECT ROWNUMBER() OVER (ORDER BY Reciept_Id) AS RowNumber ,
4           Reciept_Id ,
5           Customer_ID ,
6           Staff_ID
7     FROM Assignment1Data
8     WHERE Reciept_Id BETWEEN 52137 AND 52150
9     GROUP BY Reciept_Id , Customer_ID , Staff_ID
10 )
11 SELECT Reciept_Id , Customer_ID , Staff_ID FROM CTE WHERE (RowNumber % 2 = 0)

```

Results of the above query yielded:

| Reciept_Id | Customer_Id | Staff_Id |
|------------|-------------|----------|
| 52137 | C59 | S2 |
| 52138 | C30 | S19 |
| 52139 | C31 | S20 |
| 52140 | C52 | S10 |
| 52141 | C42 | S7 |
| 52142 | C47 | S6 |
| 52143 | C8 | S13 |
| 52144 | C50 | S4 |
| 52145 | C40 | S15 |
| 52146 | C38 | S5 |
| 52147 | C9 | S19 |
| 52148 | C43 | S16 |
| 52149 | C45 | S11 |
| 52150 | C57 | S7 |

Whereas the original result without a modulo comparison on the row would have yielded a much different result, the raw table supplied in the [appendix](#)

With this known, an additional section was added to the [python](#) script to generate update statements that would be easy to add to the current migrations.sql script we were prototyping.

The generated update statements appeared as:

```
1  — Auto-generated query to fix error of type: Staff.Id Mismatch
2  — Resolved error identified by UUID: dcf16fba08c63ecc85556c385204d9524ec359cf
3  UPDATE Assignment1Data
4  SET Reciept_Id=(
5  SELECT MAX(Reciept_Id)+1
6  FROM Assignment1Data)
7  WHERE Reciept_Id=52136
8  AND Customer.Id = 'C13' AND Staff.Id = 'S4'
9  GO
```

Determining now potential entries that broke further rules was our next objective. We pursued the idea that entries of receipts could potentially have duplicate items recorded against the ReceiptItem table. A simple script was generated to check our assumptions of this:

```
1  — Verify that no receipt has duplicate ItemIds and all are unique per order
2  SELECT *
3  FROM
4  (
5      SELECT [ReceiptItem].[ReceiptId],
6      COUNT([ReceiptItem].[ReceiptId]) AS 'ItemCount',
7      COUNT(DISTINCT [ReceiptItem].[ItemId]) AS 'ItemIdCount',
8      FROM [ReceiptItem]
9      GROUP BY [ReceiptItem].[ReceiptId]) AS SubQuery
10 WHERE [SubQuery].[ItemIdCount] != [SubQuery].[ItemCount]
11 ORDER BY [SubQuery].[ReceiptId]
12 GO
```

This query returned a result of 912 rows out of the total 2514, which we believed was a large amount given the issues identified earlier numbered in only the teens, however on manual inspection of a number of the reported issue records, it was apparent this figure was actually correct.

Given the large task associated with the entries, an additional module was written for generation of SQL in [python](#) which resulted in two queries for each duplicate item entry per receipt, the first query updating the total of one of the records to reflect the real item quantity, the later dropping the non-altered entry after the first had been completed.

The script was as follows:

```
1  — Auto-generated query to fix error of type: Item.Id Duplicate
2  — Resolved error identified by UUID: 8b34383524a00eb2097c1c22f870ef2ad104b6b8
3  UPDATE Assignment1Data
4  SET [Item_Quantity]=(
5  SELECT SUM([Item_Quantity])
6  FROM Assignment1Data
7  WHERE Reciept_Id=52316
8  AND Item.ID = 8)
9  WHERE Reciept_Id=52316
10 AND Item.ID = 8
11 AND Item_Quantity = 10
12 GO
13
14 — Auto-generated query to fix error of type: Item.Id Duplicate
15 — Resolved error identified by UUID: 8b34383524a00eb2097c1c22f870ef2ad104b6b8
16 DELETE FROM Assignment1Data
17 WHERE Reciept_Id=52316
18 AND Item.ID = 8
19 AND Item_Quantity < (
20     SELECT MAX([Item_Quantity])
21     FROM Assignment1Data
22     WHERE Reciept_Id=52316
23     AND Item.ID = 8
24 )
25 GO
```

Having now cleaned what we believed to be all discrepancies, we could finally start to look at evaluating data, our analysis outlined in [base analysis](#)

3.2 Assumptions and Reasoning

3.2.1 Item Table

An assumption of the ItemId never needing to be larger than a smallint was followed, as a basic query into the maximum range within the test data suggested that the maximum Id that currently existed was 30:

```
1 — Some basic queries for us to determine potential outlier data:
2 — What is the max of each column where datatype is int?
3 SELECT MAX(Item_ID) AS 'Max Item_ID'
4 FROM Assignment1Data;
```

With the results:

ItemDescription underwent some size optimisation, as the max data length that currently existed within the supplied data was 52, and we are to assume that into the future more items may be added, a value of 255 should allow for a varied range of descriptions.

SQL queried to determine to above assumption:

```
1 — Determine current max varchar used in Item_Description
2 SELECT MAX(DATALength(Item_Description))
3 FROM Assignment1Data;
```

We do recognise the requirements for optimisation may not require such measures, and acknowledge that a varchar(max)/text datatype would also be reasonable.

3.2.2 Price Table

The price table was designed to hold historical data as required by the business rules, an effective range can be used here to determine item pricing for time frames, current items having no end date or an end date as some point in time into the future.

Accuracy on the pricing was important, we decided to use a decimal(19,5) structure to ensure no problems should arise at any point with calculation of totals. [1]

Another notable feature of the price table is the relationships with both item and receiptItem, which allows the receiptItem table to point at a price value that can either be current or historical in nature.

3.2.3 ReceiptItem

The receipt item table acts as a line-item style associative entity, the quantity and historical priceId used at time of transaction can allow an item's price to be updated and still maintain historical pricing associated with the receipt.

As noted above, to facilitate historical value lookup, this table also holds relationships with the price table.

3.2.4 Receipt

The receipt table acts as a meta-table in this instance, other tables associate with this table with the receiptId field. Due to this it made it extremely easy to use a number of joins/inner joins to determine some of the metrics outlined in the base analysis.

3.2.5 Staff

Staff was left in a non-normalised state to ensure efficiency of queries into the future, normalising the table further would yield little value to the business based on the requirements. The office table is referenced by the staff table. This is merely to satisfy the assumption that, while the only office to exist was Newcastle in this setting, the requirement of more offices into the future is a possibility and the required join would be little impact on speed of queries in a datamart.

3.2.6 Customer

Customer, just like staff could be normalised further requiring more joins and potentially causing a performance issue into the future, for simplicity we kept only the supplied data in mind, and assumed no more data would be required by the datamart into the future.

4 Base Analysis

4.1 Raw Results

A number of metrics were considered to satisfy the request related to the best salesperson, as we are not certain if this is determined by a specific metric or a set of metrics we included a number of analyzed points for the project:

- Total receipts attributed to a staff member
- Total items sold by a staff member
- Ratio of discounted sales to normal sales for each staff member
- Total sale value per staff member
- Average sale value per staff member
- Average item value per staff member

4.1.1 Total Number of Sales

The total number of sales per staff member were considered with the following sql query:

```
1  -- Sales count per staff member (Receipt Count)
2  SELECT COUNT(*) AS 'Sales Count', s.StaffId ,s.StaffFirstName ,s.StaffSurname
3  FROM Receipt r
4  INNER JOIN ReceiptItem ri ON r.ReceiptId = ri.ReceiptId
5  INNER JOIN Item i ON i.ItemId = ri.ItemId
6  INNER JOIN Price p ON p.PriceId = ri.PriceId
7  INNER JOIN Staff s ON s.StaffId = r.ReceiptStaffId
8  GROUP BY s.StaffId ,s.StaffFirstName ,s.StaffSurname
9  ORDER BY 'Sales Count' DESC;
```

Leading to a range of 700 to 478, the top five staff were:

| Sales Count | StaffId | StaffFirstName | StaffSurname |
|-------------|---------|----------------|--------------|
| 700 | S17 | Daniel | Baker |
| 682 | S19 | Kaitlyn | Ortiz |
| 676 | S8 | Michelle | Miller |
| 664 | S5 | Stephanie | Watson |
| 664 | S6 | Evan | Hill |

4.1.2 Total Items Sold

The total items attributed to each staff member were considered also, determined by the query:

```
1  -- Item count per staff member
2  SELECT SUM(ri.ReceiptItemQuantity) AS 'Item Count', s.StaffId ,s.StaffFirstName ,s.
3  StaffSurname
4  FROM Receipt r
5  INNER JOIN ReceiptItem ri ON r.ReceiptId = ri.ReceiptId
6  INNER JOIN Staff s ON s.StaffId = r.ReceiptStaffId
7  GROUP BY s.StaffId ,s.StaffFirstName ,s.StaffSurname
8  ORDER BY 'Item Count' DESC;
```

Yielding a range of 4217 to 2813, with the top five staff members in this analysis:

| Item Count | StaffId | StaffFirstName | StaffSurname |
|------------|---------|----------------|--------------|
| 4217 | S19 | Kaitlyn | Ortiz |
| 4212 | S17 | Daniel | Baker |
| 4144 | S8 | Michelle | Miller |
| 4052 | S1 | Lauren | Martin |
| 4036 | S5 | Stephanie | Watson |

4.1.3 Discounted Sales Ratio

Consideration of the number of sales made by each staff member was also made, the following query yielding the results we required:

```
1  — Sales metrics for discounted and standard sales per staff member
2  SELECT s.StaffId ,s.StaffFirstName ,s.StaffSurname ,
3  SUM(SubQuery.[Discounted Sales]) AS 'Discounted Sales',
4  SUM(SubQuery.[Standard Sales]) AS 'Standard Sales'
5  FROM (
6      SELECT CAST(
7          CASE
8              WHEN COUNT(ri.[ReceiptItemQuantity]) >= 5
9              THEN 1
10             ELSE 0
11             END AS int) AS 'Discounted Sales',
12      CAST(
13          CASE
14              WHEN COUNT(ri.[ReceiptItemQuantity]) >= 5
15              THEN 0
16              ELSE 1
17             END AS int) AS 'Standard Sales',
18      r.ReceiptId
19  FROM Receipt r
20  INNER JOIN ReceiptItem ri ON r.ReceiptId = ri.ReceiptId
21  INNER JOIN Item i ON i.ItemId = ri.ItemId
22  INNER JOIN Price p ON p.PriceId = ri.PriceId
23  GROUP BY r.ReceiptId
24 ) AS SubQuery
25 INNER JOIN Receipt r ON SubQuery.ReceiptId = r.ReceiptId
26 INNER JOIN ReceiptItem ri ON r.ReceiptId = ri.ReceiptId
27 INNER JOIN Staff s ON s.StaffId = r.ReceiptStaffId
28 GROUP BY s.StaffId ,s.StaffFirstName ,s.StaffSurname
```

Results from the query yielded:

| StaffId | StaffFirstName | StaffSurname | Discounted Sales | Standard Sales | Discount Rate (%) |
|---------|----------------|--------------|------------------|----------------|-------------------|
| S4 | Robert | Wood | 518 | 115 | 81.83% |
| S14 | Noah | Brooks | 533 | 119 | 81.75% |
| S1 | Lauren | Martin | 533 | 126 | 80.88% |
| S16 | Jordan | Turner | 520 | 123 | 80.87% |
| S15 | Bailey | Green | 500 | 124 | 80.13% |
| S13 | Molly | Carter | 527 | 131 | 80.09% |
| S17 | Daniel | Baker | 556 | 144 | 79.43% |
| S20 | Dylan | Hall | 505 | 132 | 79.28% |
| S6 | Evan | Hill | 524 | 140 | 78.92% |
| S10 | Jonathan | Jenkins | 454 | 123 | 78.68% |
| S5 | Stephanie | Watson | 520 | 144 | 78.31% |
| S18 | Megan | James | 508 | 142 | 78.15% |
| S19 | Kaitlyn | Ortiz | 531 | 151 | 77.86% |
| S7 | Molly | Jackson | 474 | 141 | 77.07% |
| S9 | Mélissa | Garcia | 489 | 147 | 76.89% |
| S8 | Michelle | Miller | 509 | 167 | 75.30% |
| S12 | Leah | Harris | 356 | 122 | 74.48% |
| S11 | Gavin | Thompson | 395 | 137 | 74.25% |
| S2 | Joseph | Reed | 447 | 160 | 73.64% |
| S3 | Amber | Hill | 396 | 168 | 70.21% |

4.1.4 Total Sales Value per Staff Member

Consideration of the total sales per staff member was considered a highly important metric to consider also, we did consider comparing the results of this to the results of a query that did not include discount to see whom would be considered the best performer if discounts were not relevant, however we also recognise this to be too speculative in nature. The required query was as follows:

```
1  — Sales total per staff with discounts applied ($)
2  SELECT CAST(
3      CASE
4          WHEN COUNT(ri.[ReceiptItemQuantity]) >= 5
5              THEN SUM(p.[Price] * ri.[ReceiptItemQuantity]) * 0.85
6          ELSE SUM(p.[Price] * ri.[ReceiptItemQuantity])
7          END AS decimal(19,5)) AS 'Sales Totals',
8      s.StaffId ,s.StaffFirstName ,s.StaffSurname
9  FROM Receipt r
10 INNER JOIN ReceiptItem ri ON r.ReceiptId = ri.ReceiptId
11 INNER JOIN Item i ON i.ItemId = ri.ItemId
12 INNER JOIN Price p ON p.PriceId = ri.PriceId
13 INNER JOIN Staff s ON s.StaffId = r.ReceiptStaffId
14 INNER JOIN Customer c ON c.CustomerId = r.ReceiptCustomerId
15 GROUP BY s.StaffId ,s.StaffFirstName ,s.StaffSurname
16 ORDER BY 'Sales Totals' DESC;
```

Resulting in the following results:

| Sales Totals | StaffId | StaffFirstName | StaffSurname |
|--------------|---------|----------------|--------------|
| 78572.21500 | S8 | Michelle | Miller |
| 73847.10750 | S19 | Kaitlyn | Ortiz |
| 72764.88750 | S17 | Daniel | Baker |
| 71699.66750 | S14 | Noah | Brooks |
| 70514.68250 | S3 | Amber | Hill |
| 69182.56250 | S2 | Joseph | Reed |
| 69051.19500 | S13 | Molly | Carter |
| 68831.81000 | S5 | Stephanie | Watson |
| 68133.83250 | S4 | Robert | Wood |
| 66267.19000 | S6 | Evan | Hill |
| 65018.37000 | S10 | Jonathan | Jenkins |
| 64002.06750 | S1 | Lauren | Martin |
| 63760.66750 | S16 | Jordan | Turner |
| 62304.70250 | S18 | Megan | James |
| 61862.44750 | S9 | Mélissa | Garcia |
| 61832.27250 | S15 | Bailey | Green |
| 61536.85500 | S20 | Dylan | Hall |
| 58996.16250 | S7 | Molly | Jackson |
| 52102.66250 | S11 | Gavin | Thompson |
| 50259.35250 | S12 | Leah | Harris |

4.1.5 Average Value Per Sale

The average receipt value per staff member was another metric we considered would add value to the decision to be suggested in the [executive summary](#). The required query to determine this metric was as follows:

```
1  — Sales average per staff with discounts applied
2  SELECT (CAST(
3      CASE
4          WHEN COUNT(ri.[ReceiptItemQuantity]) >= 5
5              THEN SUM(p.[Price] * ri.[ReceiptItemQuantity]) * 0.85
6          ELSE SUM(p.[Price] * ri.[ReceiptItemQuantity])
7          END AS decimal(19,5)) / COUNT(r.ReceiptId)) AS 'Sales Average',
8      s.StaffId , s.StaffFirstName , s.StaffSurname
9  FROM Receipt r
10 INNER JOIN ReceiptItem ri ON r.ReceiptId = ri.ReceiptId
11 INNER JOIN Item i ON i.ItemId = ri.ItemId
12 INNER JOIN Price p ON p.PriceId = ri.PriceId
13 INNER JOIN Staff s ON s.StaffId = r.ReceiptStaffId
14 GROUP BY s.StaffId , s.StaffFirstName , s.StaffSurname
15 ORDER BY 'Sales Average' DESC;
```

With results as follows:

| Sales Average | StaffId | StaffFirstName | StaffSurname |
|----------------------|---------|----------------|--------------|
| 125.0260328014184397 | S3 | Amber | Hill |
| 116.2310872781065088 | S8 | Michelle | Miller |
| 113.9745675453047775 | S2 | Joseph | Reed |
| 112.6834835355285961 | S10 | Jonathan | Jenkins |
| 109.9688151840490797 | S14 | Noah | Brooks |
| 108.2802162756598240 | S19 | Kaitlyn | Ortiz |
| 107.6363862559241706 | S4 | Robert | Wood |
| 105.1450889121338912 | S12 | Leah | Harris |
| 104.9410258358662613 | S13 | Molly | Carter |
| 103.9498392857142857 | S17 | Daniel | Baker |
| 103.6623644578313253 | S5 | Stephanie | Watson |
| 99.7999849397590361 | S6 | Evan | Hill |
| 99.1612247278382581 | S16 | Jordan | Turner |
| 99.0901802884615384 | S15 | Bailey | Green |
| 97.9373355263157894 | S11 | Gavin | Thompson |
| 97.2679992138364779 | S9 | Mélissa | Garcia |
| 97.1199810318664643 | S1 | Lauren | Martin |
| 96.6041679748822605 | S20 | Dylan | Hall |
| 95.9287195121951219 | S7 | Molly | Jackson |
| 95.8533884615384615 | S18 | Megan | James |

We consider this to be a metric which weighs heavily in our analysis, as multiple factors would impact this result, the number of items on the sale (resulting in a lower total if discount was applied). Another consideration for this metric would be that it leans towards anyone who could sell a larger quantity of the same item, as this lends itself towards a higher receipt total.

We see that Ms Amber Hill (S3) has the highest average sale total and this is backed up by her high item sales count, indicating she is making more sales per receipt on average than any other sales office. However, Ms Hill has not made as much revenue as some other employees, approximately \$7,500 behind the sales leader who achieved \$78,572. Originally, we suspected that perhaps Ms Hill was a new employee but reviewing sale receipt dates we can confirm this was not a valid assumption and that she was working within the business throughout the entirety of 2017. After this revelation, we can rule Ms Hill out of our assessment for the best sales person. The data suggests that Ms Hill sells higher cost items which makes up for the lack of transactions she completes when compared to other staff members.

4.2 Notes on Analysis

As a group, we emphasised identifying and avoiding bad data as our top priority for the analysis outlined in this report. One of the major steps taken in this process included the removal of duplicate items on receipts and consolidating them into single line items. The reasoning behind this is that we discovered a number of receipts that included redundant entries which showed double or, in rare cases, triple the number of items expected which would falsely inflate the cost of items on a receipt.

5 Executive Summary

This report was created for the head sales executives of BIA Inc for the purpose of determining which sales staff member would be considered the best performer based on the data set provided by the firm. As no specific metrics or measurement requirements were provided, an analysis was performed by our team. This analysis included extracting, cleaning and loading the data using SQL scripts and a supplementary Python script which assisted in preparing the data provided by the firm for analysis and query in SQL Server Management Studio (SSMS).

The findings of our analysis resulted in ranking high achieving staff members determined by a number of key metrics selected by the team. Firstly, we ranked the sales officers by total number of sales and, from this group, identified the top staff member of sales, Mr Daniel Baker. Mr Baker had made the largest number of sales in the 12 months of data supplied with a total of 700 sales. Placed immediately after Mr Baker, with 18 fewer sales is Ms Kaitlyn Ortiz (682), then Ms Michelle Miller (676), followed by Ms Stephanie Watson (664) and Mr Evan Hill (664).

The next metric was total items sold. In this relation, the best sales officer is Ms Kaitlyn Ortiz, the second place sales officer in the first metric. In the 12 months of data supplied, it can be observed that of the 682 total sales, Ms Ortiz sold 4217 items with approximately six (6.18) items per sale. The second place sales officer, Mr Daniel Baker, sold 4212 items, only five items less than Ms Ortiz, and also sold approximately six (6.02) items per sale. The third placed employee in this relation is Ms Michelle Miller who resulted in 4414 total sales and approximately six (6.13) items per sale.

The third key metric is discounted sales ratio. As stated in the business rules document provided by the firm, any sale with five or more row items would be eligible for a 15% discount to the total sale. We identified this as an important factor because understanding how many items were discounted may offer insight into sales methods and techniques applied by sales officers which can then be used to enhance future performance. From the results, we aggregated the data by total percentage of sales which were discounted. Mr Robert Wood (84.57%) discounted the greatest share of his sales of all sales staff members. After Mr Wood comes Mr Dylan Hall (83.41%), then Ms Lauren Martin (83.31%), Mr Jordan Turner (82.74%), Mr Noah Brooks (82.72%) and Mr Daniel Baker (81.39%).

The final metric considered was total sales value per staff member. The purpose of this report is to find the “best” salesperson, therefore, we can assume that, along with other metrics, the firm would be interested to know which sales officer generates the most revenue. After examining the data, we can state that Ms Michelle Miller is the sales officer that generates the most value with a total of \$78,572.22 over 12 months. This equates to \$4,725.10 more than her nearest competitor, Mr Daniel Baker who generated \$72,764.89 who was followed by Mr Noah Brooks with \$71,699.67 and Ms Amber Hill with \$70,514.68.

After carefully considering the aforementioned key metrics and reviewing the results, we can conclude that Ms Michelle Miller should be considered the most valuable sales officer at BIA Inc. Our findings showed clearly that Ms Michelle Miller achieved the highest sales value when compared with other sales officer by a significant margin (\$4,725.10). While Ms Miller was not ranked first in total number of sales or total items sold, she did rank highly in both relations. Ms Miller did however score poorly in her discounted sales ratio. This could indicate that Ms Miller is not as effective at ‘upselling’ as her fellow sales officers and this could be an area for improvement. We can assert that Ms Miller should be considered for the reward (and possible cash prize) suggested in the original document outlining the firm requirements. If for any reason Ms Miller should not be applicable or eligible for the discount, we would recommend Ms Kaitlyn Ortiz as the alternative choice due to her high ranking in all metrics discussed in this summary.

References

- [1] Reasons against TSQL Money type: Stackoverflow User; *SQLMenace* <https://stackoverflow.com/questions/582797/should-you-choose-the-money-or-decimalx-y-datatypes-in-sql-server>
- [2] Microsoft TSQL documentation of Decimal/Numeric types <https://docs.microsoft.com/en-us/sql/t-sql/data-types/decimal-and-numeric-transact-sql?view=sql-server-2017>
- [3] Microsoft documentation: WITH common_table_expression (Transact-SQL) <https://docs.microsoft.com/en-us/sql/t-sql/queries/with-common-table-expression-transact-sql?view=sql-server-2017>
- [4] Upselling - Business Dictionary <http://www.businessdictionary.com/definition/upselling.html>

6 Appendix

6.1 CTE Raw Results

| Reciept_Id | Customer_Id | Staff_Id |
|------------|-------------|----------|
| 52137 | C27 | S4 |
| 52137 | C59 | S2 |
| 52138 | C29 | S13 |
| 52138 | C30 | S19 |
| 52139 | C3 | S5 |
| 52139 | C31 | S20 |
| 52140 | C38 | S4 |
| 52140 | C52 | S10 |
| 52141 | C24 | S19 |
| 52141 | C42 | S7 |
| 52142 | C46 | S8 |
| 52142 | C47 | S6 |
| 52143 | C51 | S17 |
| 52143 | C8 | S13 |
| 52144 | C11 | S10 |
| 52144 | C50 | S4 |
| 52145 | C21 | S8 |
| 52145 | C40 | S15 |
| 52146 | C38 | S16 |
| 52146 | C38 | S5 |
| 52147 | C40 | S18 |
| 52147 | C9 | S19 |
| 52148 | C26 | S8 |
| 52148 | C43 | S16 |
| 52149 | C10 | S19 |
| 52149 | C45 | S11 |
| 52150 | C15 | S10 |
| 52150 | C57 | S7 |

6.2 Python Script

```
1  #!/usr/bin/env python3.7
2  import classes as Classes
3  import os
4  import sys
5  import csv
6  import re
7  import openpyxl
8  import traceback
9
10 # This is courtesy of: https://stackoverflow.com/questions/1323364/in-python-how-to-check-if-a-string-only-contains-certain-characters
11 # Required to determine pesky dates Peter was nice enough to put in
    Receipt_Id column.
12 def special_match(strg, search=re.compile(r'^a-zA-Z').search):
13     return bool(search(strg))
14
15 # Function to parse all receipts once populated and add to employee
    totals
16 def populate_receipt_totals(sales, employees, customers, items):
17     for receipt_id, sale in sales.sales.items():
18         total = 0
19         for item_id, item in sale.receipt.items.items():
20             total = total + (item.quantity * item.price)
21             employees.employees[sale.receipt.staff.id].item_count +=
                item.quantity
22
23         if(len(sale.receipt.items.items()) > 4):
24             print("Total was adjusted from {} to {} due to business
                rules related to number\nof items in a sale.".format(
                    total, total * 0.85))
25             total *= 0.85
26             employees.employees[sale.receipt.staff.id].
                discounted_sales += 1
27
28             print("Total calculated for receipt {} is: {}, Items count
                was: {}".format(receipt_id, total, len(sale.receipt.items.
                    items()))
29             employees.employees[sale.receipt.staff.id].sales_count += 1
30             employees.employees[sale.receipt.staff.id].sales_total +=
                total
31
32     populate_customer_totals(sales, customers)
33     populate_item_totals(sales, items)
34     generate_employee_report(employees)
35
36 # Function to parse all receipts once populated and add to customer
    totals
37 def populate_item_totals(sales, items):
38     for receipt_id, sale in sales.sales.items():
39         total = 0
40         for item_id, item in sale.receipt.items.items():
41             total = total + (item.quantity * item.price)
42             items.items[sale.receipt.items[item_id].id].item_count +=
                item.quantity
```

```

43
44     if(len(sale.receipt.items.items()) > 4):
45         total *= 0.85
46         items.items[sale.receipt.items[item_id].id].
            discounted_sales += 1
47
48         items.items[sale.receipt.items[item_id].id].sales_count += 1
49         items.items[sale.receipt.items[item_id].id].sales_total +=
            total
50
51     generate_items_report(items)
52
53 # Function to parse all receipts once populated and add to customer
    totals
54 def populate_customer_totals(sales,customers):
55     for receipt_id,sale in sales.sales.items():
56         total = 0
57         for item_id,item in sale.receipt.items.items():
58             total = total + (item.quantity * item.price)
59             customers.customers[sale.receipt.customer.id].item_count
                += item.quantity
60
61         if(len(sale.receipt.items.items()) > 4):
62             total *= 0.85
63             customers.customers[sale.receipt.customer.id].
                discounted_sales += 1
64
65             customers.customers[sale.receipt.customer.id].sales_count +=
                1
66             customers.customers[sale.receipt.customer.id].sales_total +=
                total
67
68     generate_customer_report(customers)
69
70 # Generation of required output files
71 def generate_results_structures():
72     try:
73         if not os.path.exists('Results'):
74             os.makedirs('Results')
75
76         open('Results/Employee_Results.txt','w+').close()
77         open('Results/Item_Results.txt','w+').close()
78         open('Results/Customer_Results.txt','w+').close()
79
80     except Exception:
81         print("An error occurred: {}".format(traceback.format_exc()))
82
83 # Main branch of code to parse rows in excel file
84 def parse_rows(rows,logged_errors):
85     for row in rows:
86         receipt_id = row[1].value
87         if (isinstance(receipt_id,int) or special_match(receipt_id)):
88             if receipt_id in sales.sales:
89                 staff_id = row[5].value
90                 customer_id = row[2].value
91                 item_id = row[11].value

```



```

92         item_quantity = row[13].value
93
94         for item in sales.sales[receipt_id].receipt.items.items():
95             if item_id == item[0]:
96                 if staff_id == sales.sales[receipt_id].receipt.staff.id:
97                     print("Error in data row; {} is the same as {}".format(item_id, item_id))
98                     logged_errors.add_error(receipt_id, "Error in data row id: {}; {} is the same as {}".format(receipt_id, item_id, item_id), "Item.Id Duplicate", customer_id, staff_id, item_id, item_quantity, sales.sales[receipt_id].receipt.items[item_id].quantity)
99
100                 if sales.sales[receipt_id].receipt.staff.id != staff_id:
101                     print("Error in data row; {} is not the same as {}".format(sales.sales[receipt_id].receipt.staff.id, staff_id))
102                     logged_errors.add_error(receipt_id, "Error in data row id: {}; {} is not the same as {}".format(receipt_id, sales.sales[receipt_id].receipt.staff.id, staff_id), "Staff.Id Mismatch", customer_id, staff_id, item_id, item_quantity, None)
103
104                 if sales.sales[receipt_id].receipt.customer.id != customer_id:
105                     print("Error in data row; {} is not the same as {}".format(sales.sales[receipt_id].receipt.customer.id, customer_id))
106                     logged_errors.add_error(receipt_id, "Error in data row id: {}; {} is not the same as {}".format(receipt_id, sales.sales[receipt_id].receipt.customer.id, customer_id), "Customer.Id Mismatch", customer_id, staff_id, item_id, item_quantity, None)
107
108                 print("Found existing receipt {}, adding items instead".format(receipt_id))
109                 sales.add_items_to_sale(row, receipt_id)
110             else:
111                 sales.parse_row(row, employees, customers, items)
112         else:
113             raise ValueError('Non int receipt id.')
114
115 # Clear the current errors.txt file
116 def clear_error_log():
117     open('Results/Errors.txt', 'w+').close()
118     open('Results/SQL.txt', 'w+').close()
119
120 # Function to generate employee report and output to disk
121 def generate_employee_report(employees):

```

```

122 employee_output = ""
123 header = "Results for Employee analysis:"
124 for employee_id, employee in employees.items():
125     employee_output += ""Employee: {}, {} {} \n
126     Metrics: #####
127     Sales Count = {}
128     Total Discounted Sales: {}
129     Discounted Sales Ratio: {}
130     Total Items Sold: {} \n
131     Financials: #####
132     Sales Total = ${}
133     Average Sale Value: ${}
134     Average Item Sold Value: ${}
135     \n"".format(
136         employee_id,
137         employee.first_name,
138         employee.surname,
139         employee.sales_count,
140         employee.discounted_sales,
141         employee.discounted_sales / employee.sales_count,
142         employee.item_count,
143         employee.sales_total,
144         employee.sales_total / employee.sales_count,
145         employee.sales_total / employee.item_count)
146 write_report_results('Employee Results', header, employee_output)
147
148 # Function to generate customer report and output to disk
149 def generate_customer_report(customers):
150     customer_output = ""
151     header = "Results for Customer analysis:"
152     for customer_id, customer in customers.items():
153         customer_output += ""Customer: {}, {} {} \n
154         Metrics: #####
155         Sales Count = {}
156         Total Discounted Sales: {}
157         Discounted Sales Ratio: {}
158         Total Items Sold: {} \n
159         Financials: #####
160         Sales Total = ${}
161         Average Sale Value: ${}
162         Average Item Sold Value: ${}
163         \n"".format(
164             customer_id,
165             customer.first_name,
166             customer.surname,
167             customer.sales_count,
168             customer.discounted_sales,
169             customer.discounted_sales / customer.sales_count,
170             customer.item_count,
171             customer.sales_total,
172             customer.sales_total / customer.sales_count,
173             customer.sales_total / customer.item_count)
174     write_report_results('Customer Results', header, customer_output)
175
176 # Function to generate item report and output to disk
177 def generate_items_report(items):

```

```

178 items_output = ""
179 header = "Results for Item analysis:"
180 for item_id,item in items.items.items():
181     items_output += ""Item: {}\\n
182     Metrics: #####
183     Sales Count = {}
184     Total Discounted Sales: {}
185     Discounted Sales Ratio: {}
186     Total Items Sold: {}\\n
187     Financials: #####
188     Sales Total = ${}
189     Average Sale Value: ${}
190     Average Item Sold Value: ${}
191     \\n"".format(
192         item_id,
193         item.sales_count,
194         item.discounted_sales,
195         item.discounted_sales / item.sales_count,
196         item.item_count,
197         item.sales_total,
198         item.sales_total / item.sales_count,
199         item.sales_total / item.item_count)
200 write_report_results('Item Results',header,items_output)
201
202 # Function to generate error report and output to disk
203 def generate_error_report(logged_errors):
204     header = "Error Report:"
205     error_output = ""
206     for error_log_id,error_log in logged_errors.logged_errors.items():
207         error_output += "ErrorId = {}\\nErrorType = {}\\nReceiptId =
208             {}\\nError = {}\\n\\n"".format(
209             error_log_id,
210             error_log.error_type,
211             error_log.receipt_id,
212             error_log.trace)
213     write_report_results('Errors',header,error_output)
214
215 def generate_sql_move_items(logged_errors):
216     header = "USE EBUS3030;"
217     sql_output = ""
218     parsed_receipt_ids = []
219     for error_log_id,error_log in logged_errors.logged_errors.items():
220         if error_log.receipt_id not in parsed_receipt_ids and
221             error_log.error_type != "Item.Id Duplicate":
222             sql_output += ""
223             -- Auto-generated query to fix error of type: {}
224             -- Resolved error identified by UUID: {}
225             UPDATE Assignment1Data
226             SET Reciept_Id=(
227             SELECT MAX(Reciept_Id)+1
228             FROM Assignment1Data)
229             WHERE Reciept_Id={}
230             AND "".format(error_log.error_type,error_log_id,error_log.receipt_id
231 )

```

```

229         if error_log.customer_id is not None and error_log.
230            staff_id is not None:
231             sql_output += "Customer_Id = '{}{}' AND Staff_Id =
232                '{}{}'\nGO\n".format(error_log.customer_id,error_log
233                .staff_id)
234             elif error_log.customer_id is not None:
235              sql_output += "Customer_Id = '{}{}'\nGO\n".format(
236                error_log.customer_id)
237             elif error_log.staff_id is not None:
238              sql_output += "Staff_Id = '{}{}'\nGO\n".format(
239                error_log.staff_id)
240             else:
241              sql_output = None
242              parsed_receipt_ids.append(error_log.receipt_id)
243
244     write_report_results('SQL',header,sql_output)
245
246 def generate_sql_fix_duplicate_items(logged_errors):
247     header = "USE EBUS3030;"
248     sql_output = ""
249     for error_log_id,error_log in logged_errors.logged_errors.items():
250         print(error_log.error_type)
251         if error_log.error_type == "Item.Id Duplicate":
252             if error_log.item_quantity == error_log.
253                duplicate_item_quantity:
254                 new_quantity = error_log.item_quantity * 2
255                 print(error_log.receipt_id)
256                 # error_log.item_quantity += error_log.
257                    duplicate_item_quantity
258                 sql_output += ""
259
260     -- Auto-generated query to fix error of type: {}
261     -- Resolved error identified by UUID: {}
262     UPDATE Assignment1Data
263     SET [Item_Quantity]={}
264     WHERE Reciept_Id={}
265     AND Item_ID = {}
266     AND Item_Quantity = {}'\nGO\n"".format(error_log.error_type,
267     error_log_id,
268     new_quantity,
269     error_log.receipt_id,
270     error_log.item_id,
271     error_log.item_quantity)
272
273     else:
274         sql_output += ""
275
276     -- Auto-generated query to fix error of type: {}
277     -- Resolved error identified by UUID: {}
278     UPDATE Assignment1Data
279     SET [Item_Quantity]=(
280     SELECT SUM([Item_Quantity])
281     FROM Assignment1Data
282     WHERE Reciept_Id={}
283     AND Item_ID = {})
284     WHERE Reciept_Id={}
285     AND Item_ID = {}

```

```

277 AND Item_Quantity = {} \nGO\n"".format(error_log.error_type,
278                                         error_log_id,
279                                         error_log.receipt_id,
280                                         error_log.item_id,
281                                         error_log.receipt_id,
282                                         error_log.item_id,
283                                         error_log.item_quantity)
284         sql_output += ""
285 -- Auto-generated query to fix error of type: {}
286 -- Resolved error identified by UUID: {}
287 DELETE FROM Assignment1Data
288 WHERE Reciept_Id={}
289 AND Item_ID = {}
290 AND Item_Quantity < (
291     SELECT MAX([Item_Quantity])
292     FROM Assignment1Data
293     WHERE Reciept_Id={}
294     AND Item_ID = {}
295 ) \nGO\n"".format(error_log.error_type, error_log_id,
296                   error_log.receipt_id,
297                   error_log.item_id,
298                   error_log.receipt_id,
299                   error_log.item_id)
300
301     write_report_results('SQL', header, sql_output)
302
303
304 # Generalised function to write a report to disk
305 def write_report_results(report_name, header, report_body):
306     with open('Results/{}.txt'.format(report_name), 'a+') as report:
307         report.write(header + 2*'\n')
308         report.write(report_body)
309
310 # Main hook
311 if __name__ == '__main__':
312     # Open excel file stored in child folder
313     excel_file = openpyxl.load_workbook('Data/Assignment1Data.xlsx')
314     data = excel_file['Asgn1 Data']
315     sales = Classes.Sales()
316     employees = Classes.Employees()
317     customers = Classes.Customers()
318     items = Classes.Items()
319     logged_errors = Classes.LoggedErrors()
320
321     clear_error_log()
322
323     # Main branch of code to parse excel file.
324     parse_rows(data.rows, logged_errors)
325
326     # If results folder and required text files don't exist, create
    them
327     generate_results_structures()
328
329     # Output error report to disk.
330     generate_error_report(logged_errors)
331

```

```
332 # Output sql to disk to fix errors found
333 generate_sql_fix_duplicate_items(logged_errors)
334 generate_sql_move_items(logged_errors)
335
336 # Iterate over sales and employees to generate reports
337 # populate_receipt_totals(sales,employees,customers,items)
```

```

1  #!/usr/bin/env python3.7
2  import hashlib
3
4  # Item class, to imitate item entries in receipt
5  class Item:
6      def __init__(self, item_id, item_description, item_price,
7                    item_quantity):
8          self.id = item_id
9          self.description = item_description
10         self.price = item_price
11         self.quantity = item_quantity
12         self.sales_count = 0
13         self.sales_total = 0
14         self.item_count = 0
15         self.discounted_sales = 0
16
17  # Office class, to imitate office entries in staff
18  class Office:
19      def __init__(self, office_id, office_location):
20          self.id = office_id
21          self.location = office_location
22
23  # Staff class, to emulate staff
24  class Staff:
25      def __init__(self, staff_id, staff_first_name, staff_surname, office)
26          :
27          self.id = staff_id
28          self.first_name = staff_first_name
29          self.surname = staff_surname
30          self.office = office
31          self.sales_count = 0
32          self.sales_total = 0
33          self.item_count = 0
34          self.discounted_sales = 0
35
36  # Customer class to emulate customers
37  class Customer:
38      def __init__(self, customer_id, customer_first_name,
39                    customer_surname):
40          self.id = customer_id
41          self.first_name = customer_first_name
42          self.surname = customer_surname
43          self.sales_count = 0
44          self.sales_total = 0
45          self.item_count = 0
46          self.discounted_sales = 0
47
48  # Receipt class to hold data for a sale
49  class Receipt:
50      def __init__(self, receipt_id, customer, staff):
51          self.id = receipt_id
52          self.customer = customer
53          self.staff = staff
54          self.items = {}
55          self.item_count = 0

```

```

53         self.total = 0
54
55     # Function to add items to receipt
56     def add_item(self, item):
57         self.items[item.id] = item
58         self.item_count += 1
59
60 # Sale class to hold one receipt (Kinda redundant)
61 class Sale:
62     def __init__(self, date, receipt):
63         self.date = date
64         self.receipt = receipt
65
66 # Sales class to hold record of all sales
67 class Sales:
68     def __init__(self):
69         self.sales = {}
70
71     # Parse row function, intended to determine if row is a header
72     # row or contains formula
73     def parse_row(self, row, employees, customers, items):
74         if row[0].value != 'Sale Date' and isinstance(row[1].value,
75             int):
76             item = Item(row[11].value, row[12].value, row[14].value, row
77                 [13].value)
78             customer = Customer(row[2].value, row[3].value, row[4].
79                 value)
80             office = Office(row[8].value, row[9].value)
81             staff = Staff(row[5].value, row[6].value, row[7].value,
82                 office)
83             receipt = Receipt(row[1].value, customer, staff)
84             receipt.add_item(item)
85             sale = Sale(row[0].value, receipt)
86             self.sales[sale.receipt.id] = sale
87             print("Added sale: {}".format(sale.receipt.id))
88
89             if staff.id in employees.items():
90                 print("Duplicate employee: {}".format(staff.id))
91             else:
92                 employees.add_employee(staff.id, staff)
93
94             if customer.id in customers.items():
95                 print("Duplicate customer: {}".format(customer.id))
96             else:
97                 customers.add_customer(customer.id, customer)
98
99             # We itemed your items so you can .items() your items
100             if item.id in items.items():
101                 print("Duplicate item: {}".format(item.id))
102             else:
103                 items.add_item(item.id, item)
104
105         else:
106             print("Skipped row, either it was a row header: {} or it
107                 was a formula: {}".format(row[0].value, row[1].value))

```



```

103 # Add items to sale if the receipt already exists
104 def add_items_to_sale(self,row,existing_sale_identiflier):
105     item = Item(row[11].value,row[12].value,row[14].value,row
106                 [13].value)
107     self.sales[existing_sale_identiflier].receipt.add_item(item)
108     print("Added items to receipt {} : ID: {}, Desc: {}, Price:
109           {}, Quantity: {}".format(existing_sale_identiflier,item.id,
110                                     item.description,item.price,item.quantity))
111
112 # Employees class to hold all staff
113 class Employees:
114     def __init__(self):
115         self.employees = {}
116
117     # Function to add new employees if they currently don't exist
118     def add_employee(self,employee_id,employee):
119         self.employees[employee_id] = employee
120
121 # Customers class to hold all customers
122 class Customers:
123     def __init__(self):
124         self.customers = {}
125
126     # Function to add new customer if they currently don't exist
127     def add_customer(self,customer_id,customer):
128         self.customers[customer_id] = customer
129
130 # Items class to hold all items
131 class Items:
132     def __init__(self):
133         self.items = {}
134
135     # Function to add new items if they currently don't exist
136     def add_item(self,item_id,item):
137         self.items[item_id] = item
138
139 class Error_Log:
140     def __init__(self,trace,error_type,receipt_id,customer_id = None,
141                 staff_id = None,item_id = None,item.quantity = None,
142                 duplicate_item_quantity = None):
143         self.trace = trace
144         self.receipt_id = receipt_id
145         self.error_type = error_type
146         self.customer_id = None
147         self.staff_id = None
148         if customer_id is not None:
149             self.customer_id = customer_id
150         if staff_id is not None:
151             self.staff_id = staff_id
152         if item_id is not None:
153             self.item_id = item_id
154         if item.quantity is not None:
155             self.item_quantity = item.quantity
156         if duplicate_item_quantity is not None:
157             self.duplicate_item_quantity = duplicate_item_quantity
158         self.hash = self.generate_hash(trace + error_type + str(

```

```

        receipt_id))
155
156     def generate_hash(self, hashcontent):
157         return str(hashlib.shal(hashcontent.encode(encoding='UTF-8',
158                                     errors='strict')).hexdigest())
159
160 # Logged errors class to avoid logging the same error multiple times
161 class LoggedErrors:
162     def __init__(self):
163         self.logged_errors = {}
164         self.error_count = 0
165
166     # Determine if error related to receipt is already logged
167     def add_error(self, receipt_id, trace, error_type, customer_id = None,
168                 staff_id = None, item_id = None, item_quantity = None,
169                 duplicate_item_quantity = None):
170         error_log = Error_Log(trace, error_type, receipt_id, customer_id,
171                               staff_id, item_id, item_quantity, duplicate_item_quantity)
172         if error_log.hash not in self.logged_errors:
173             self.logged_errors[error_log.hash] = error_log

```
