

EBUS3030 Assignment 1

Steven Karmaniolos c3160280@uon.edu.au

Jay Rovacsek c3146220@uon.edu.au

Jacob Litherland c3263482@uon.edu.au

Edward Lonsdale c3252144@uon.edu.au

August 31, 2018

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.

- 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.

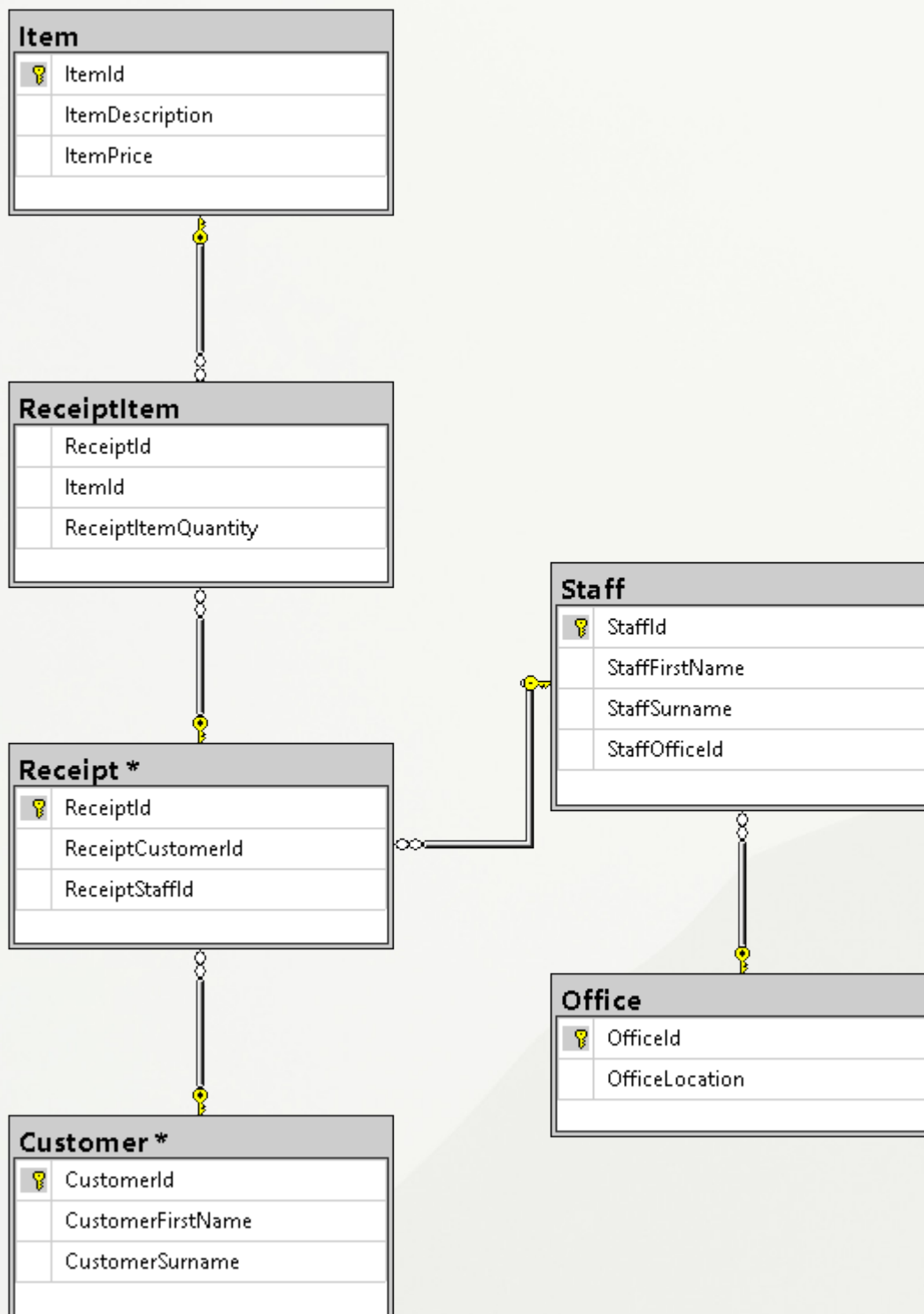
1 Datamart Business Notes

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.

2 Data Model

The below data model is only a suggestion and is still subject to change into the future.



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 analyse 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"/>

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.

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:

```
INSERT INTO Receipt([ReceiptId], [ReceiptCustomerId],[ReceiptStaffId])
SELECT DISTINCT([Reciept_Id]), [Customer_ID], [Staff_ID]
FROM [Assignment1Data]
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:

```
SELECT * FROM Assignment1Data WHERE Customer_ID='C32' AND Staff_ID='S15' AND
Sale_Date='2017-11-12 00:00:00.0000000';
```

```
SELECT * FROM Assignment1Data WHERE Customer_ID='C13' AND Staff_ID='S4' 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:

```
UPDATE Assignment1Data
SET Reciept_Id=51585,
Reciept_Transaction_Row_ID=(
    SELECT MAX(Reciept_Transaction_Row_ID)+1
    FROM Assignment1Data
    WHERE Reciept_Id=51585)
WHERE Customer_ID='C32'
AND Staff_ID='S15'
AND Sale_Date='2017-11-12 00:00:00.0000000'
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:

```
UPDATE Assignment1Data SET Reciept_Id=(SELECT MAX(Reciept_Id)+1 FROM Assignment1Data)
WHERE Customer_ID='C27' AND Staff_ID='S4' AND Sale_Date='2017-12-30 00:00:00.0000000';
```

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

```
UPDATE Assignment1Data SET Reciept_Id=(SELECT MAX(Reciept_Id)+1 FROM Assignment1Data)
WHERE Customer_ID='C30' AND Staff_ID='S19' 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:

```
SELECT Reciept_Id, Customer_ID,Staff_ID FROM Assignment1Data
WHERE Reciept_Id BETWEEN 52137 AND 52150
GROUP BY Reciept_Id, Customer_ID,Staff_ID
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.

```
WITH CTE AS
(
    SELECT ROW_NUMBER() OVER (ORDER BY Reciept_Id) AS RowNumber,
           Reciept_Id,
           Customer_ID,
           Staff_ID
    FROM Assignment1Data
    WHERE Reciept_Id BETWEEN 52137 AND 52150
    GROUP BY Reciept_Id, Customer_ID,Staff_ID
)
SELECT Reciept_Id,Customer_ID,Staff_ID FROM CTE WHERE (RowNumber % 2 = 0)
```


Results of the above command 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:

```
-- Auto-generated query to fix error of type: Staff.Id Mismatch
-- Resolved error identified by UUID: dcf16fba08c63ecc85556c385204d9524ec359cf
UPDATE Assignment1Data
SET Reciept_Id=(
SELECT MAX(Reciept_Id)+1
FROM Assignment1Data)
WHERE Reciept_Id=52136
AND Customer_Id = 'C13' AND Staff_Id = 'S4'
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:

```
-- Verify that no receipt has duplicate ItemIds and all are unique per order
SELECT *
FROM
(
    SELECT [ReceiptItem].[ReceiptId],
    COUNT([ReceiptItem].[ReceiptId]) AS 'ItemCount',
    COUNT(DISTINCT [ReceiptItem].[ItemId]) AS 'ItemIdCount'
    FROM [ReceiptItem]
    GROUP BY [ReceiptItem].[ReceiptId]) AS SubQuery
WHERE [SubQuery].[ItemIdCount] != [SubQuery].[ItemCount]
ORDER BY [SubQuery].[ReceiptId]
```

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:

```
-- Auto-generated query to fix error of type: Item.Id Duplicate
-- Resolved error identified by UUID: 0ee74976129cce87fb1558eb5586b1511f5c8d8f
UPDATE Assignment1Data
SET [Item_Quantity]=(
SELECT SUM([Item_Quantity])
FROM Assignment1Data
WHERE Reciept_Id=51500
AND Item_ID = 20)
WHERE Reciept_Id=51500
AND Item_ID = 20
AND Item_Quantity = 1
GO

-- Auto-generated query to fix error of type: Item.Id Duplicate
-- Resolved error identified by UUID: 0ee74976129cce87fb1558eb5586b1511f5c8d8f
DELETE FROM Assignment1Data
WHERE Reciept_Id=51500
AND Item_ID = 20
AND Item_Quantity < 1
GO
```

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:

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

With the results:

```
Max Item_ID
30
```

ItemDescription underwent some size optimisation, as the max datalength 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:

```
-- Determine current max varchar used in Item_Description
SELECT MAX(DATALength(Item_Description))
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.

ItemPrice while imported as float type was considered too precise for the usecase of a monetary value. While MONEY and derivatives exist in the TSQL ecosphere, there are real concerns of accuracy of the datatype [1], and therefore we decided for a decimal(19,5) typing [2].

The final Item table structure is reflected as:

Item			
	Column Name	Data Type	Allow Nulls
	ItemId	smallint	<input type="checkbox"/>
	ItemDescription	varchar(255)	<input type="checkbox"/>
	ItemPrice	decimal(19, 5)	<input type="checkbox"/>
			<input type="checkbox"/>

4 Base Analysis

4.1 Raw Results

5 Executive Summary

6 Assumptions

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>

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

7 Appendix