**WIE3008**

**BUSINESS ANALYTICS AND INTELLIGENCE**

**ASSIGNMENT 2**

**Name: Amir Firdaus Bin Abdul Hadi**

**Matric Number: 17204620/2**

**Group Members**

Imran Asyraaf Bin Mohd Ikhalan

Aiman Wafiq Bin Appandi
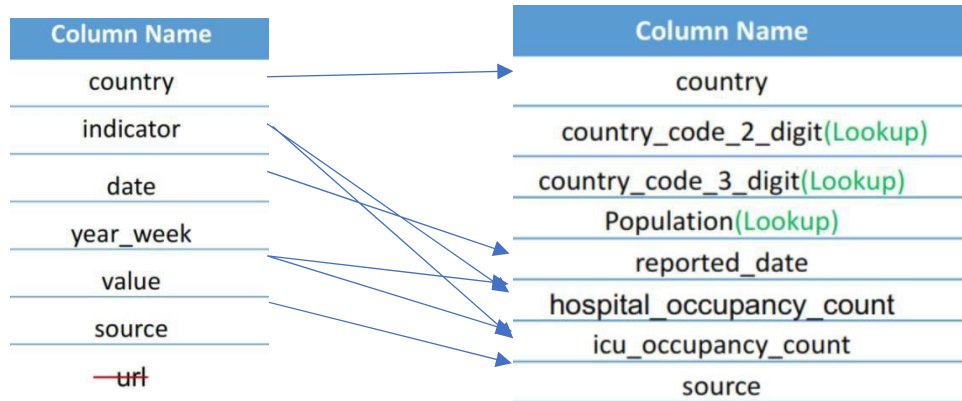
Irfan Abidin As Salik Bin Noor Riza Al Jeffery

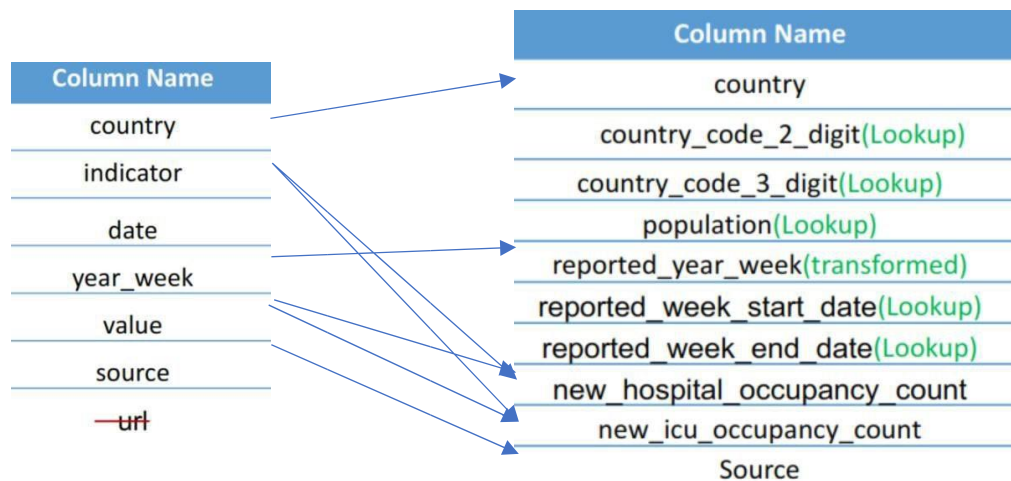Mohamad Nizar Mustaqeem Bin Mazlan

**Part A (5 marks)**

Building the Azure Data Factory pipeline.

1. Complete the pipeline by preparing the following sink datasets based on provided mapping and transformation requirements.
   a. Transformed Daily File

| Column Name | Column Name |
|---|---|
| country | country |
| indicator | country_code_2_digit(Lookup) |
| date | country_code_3_digit(Lookup) |
| year_week | Population(Lookup) |
| value | reported_date |
| source | hospital_occupancy_count |
| ~~url~~ | icu_occupancy_count |
|  | source |

   b. Transformed Weekly File

| Column Name | Column Name |
|---|---|
| country | country |
| indicator | country_code_2_digit(Lookup) |
| date | country_code_3_digit(Lookup) |
| year_week | population(Lookup) |
| value | reported_year_week(transformed) |
| source | reported_week_start_date(Lookup) |
| ~~url~~ | reported_week_end_date(Lookup) |
|  | new_hospital_occupancy_count |
|  | new_icu_occupancy_count |
|  | Source |

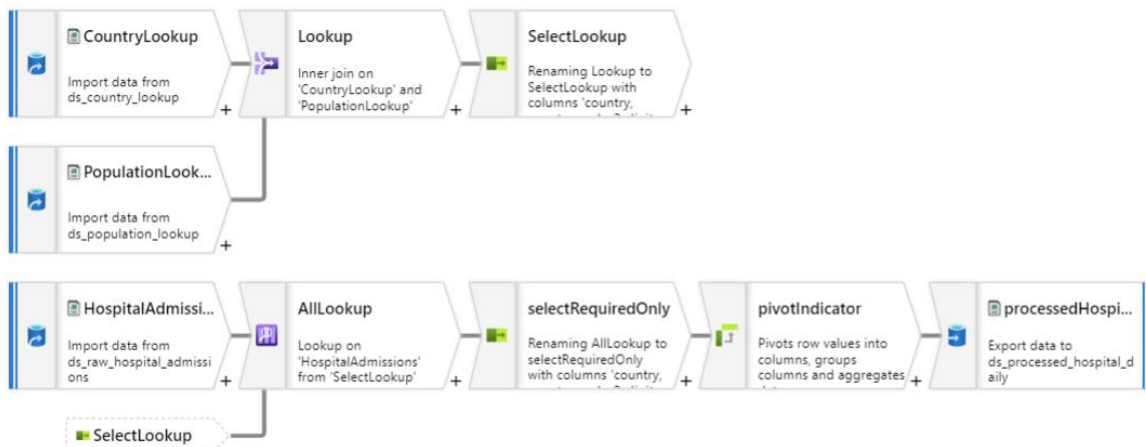a. Pipeline for Transformed Daily File:



Figure 1: Pipeline for Hospital Admission Daily File

For this pipeline, we used 3 datasets to achieve the required sink result which are Country Lookup, Population Lookup and Hospital Admission. There are 3 lookup columns needed which are "country_code_2_digit", "country_code_3_digit" and "population". Country Lookup provides the lookup for "country_code_2_digit" and "country_code_3_digit" while Population Lookup provides the lookup for "population". Hospital Admission dataset is obtained through http file and then stored into data lake.

Hospital Admission Http link:
https://raw.githubusercontent.com/cloudboxacademy/covid19/main/ecdc_data/hospital_admissions.csv

All Lookup file can be found in this drive:
https://drive.google.com/drive/folders/1nuXV8XIHCw6jrDM496jTmmCf8BkyMt6I?usp=share_link

First, we joined all lookup columns needed from Country Lookup and Population Lookup. Then, we select only required field after the lookup joining to remove any redundant and unnecessary columns. After that, to create "hospital_occupancy_count" and "icu_occupancy_count", we pivoted the "indicator" column into 2 new columns based on the "value" column. Lastly, the processed data is sink as a transformed file. The final transformation output can be seen on figure 2 while the sample exported excel sink file on figure 3.

| Order | Column | Type | | Updated | Input column |
|---|---|---|---|---|---|
| | Schema | | Input / **Output** | | |
| | Number of columns **Updated** 0 | | Dropped 0 | Unchanged 8 | Total 8 |
| 1 | country | abc string | | | country |
| 2 | country_code_2_digit | abc string | | | country_code_2_digit |
| 3 | country_code_3_digit | abc string | | | country_code_3_digit |
| 4 | population | abc string | | | population |
| 5 | reported_date | 📅 date | | | reported_date |
| 6 | source | abc string | | | source |
| 7 | Daily hospital occupancy_count | 12l long | | | Daily hospital occupancy_count |
| 8 | Daily ICU occupancy_count | 12l long | | | Daily ICU occupancy_count |

Figure 2: Final Transformation Output for Daily File

| | A | B | C | D | E | F | G | H |
|---|---|---|---|---|---|---|---|---|
| 1 | country | country_code_2_digit | country_code_3_digit | population | reported_date | source | hospital_occupancy_count | icu_occupancy_count |
| 2 | Austria | AT | AUT | 9015361 | 2/4/2020 | Surveillance | 1057 | 196 |
| 3 | Austria | AT | AUT | 9015361 | 8/4/2020 | Surveillance | 1096 | 176 |
| 4 | Austria | AT | AUT | 9015361 | 15/4/2020 | Surveillance | 1001 | 169 |
| 5 | Austria | AT | AUT | 9015361 | 16/4/2020 | Surveillance | 967 | 156 |
| 6 | Austria | AT | AUT | 9015361 | 17/4/2020 | Surveillance | 909 | 140 |
| 7 | Austria | AT | AUT | 9015361 | 19/4/2020 | Surveillance | 817 | 136 |
| 8 | Austria | AT | AUT | 9015361 | 21/4/2020 | Surveillance | 756 | 196 |
| 9 | Austria | AT | AUT | 9015361 | 22/4/2020 | Surveillance | 700 | 176 |
| 10 | Austria | AT | AUT | 9015361 | 23/4/2020 | Surveillance | 677 | 169 |

Figure 3: Sample of Exported Sink Daily File

b. Pipeline for Transformed Weekly File:



Figure 4: Pipeline for Hospital Admission Weekly File

The process for this weekly pipeline is similar with daily pipeline with removal of reported date and addition of reported year week and its start and end date. For the week date, we add new dataset to the data lake called "yearWeekLookup" to provide the dates for the start and end date of the week. The rest is the same and the processed data is sink as a transformed weekly file. The final transformation output can be seen on figure 5 while the sample exported excel sink file on figure 6.



| Order | Column | Type | Updated | Input column |
|---|---|---|---|---|
| 1 | country | abc string | | country |
| 2 | country_code_2_digit | abc string | | country_code_2_digit |
| 3 | country_code_3_digit | abc string | | country_code_3_digit |
| 4 | population | abc string | | population |
| 5 | reported_year_week | abc string | | reported_year_week |
| 6 | reported_week_start_date | date | | reported_week_start_date |
| 7 | reported_week_end_date | date | | reported_week_end_date |
| 8 | new_hospital_occupancy_count | long | | new_hospital_occupancy_count |
| 9 | new_icu_occupancy_count | long | | new_icu_occupancy_count |
| 10 | source | abc string | | source |

Figure 5: Final Transformation Output for Weekly File

| | A | B | C | D | E | F | G | H | I | J |
|---|---|---|---|---|---|---|---|---|---|---|
| 1 | country | country_code_2_digit | country_code_3_digit | population | reported_year_week | reported_wee | reported_wee | new_hospital_occupancy_count | new_icu_occupancy_count | source |
| 2 | Austria | AT | AUT | 9015361 | 2020-W14 | 30/3/2020 | 5/4/2020 | 1057 | 697 | Surveillance |
| 3 | Austria | AT | AUT | 9015361 | 2020-W15 | 6/4/2020 | 12/4/2020 | 1096 | 407 | Surveillance |
| 4 | Austria | AT | AUT | 9015361 | 2020-W16 | 13/4/2020 | 19/4/2020 | 3694 | 356 | Surveillance |
| 5 | Austria | AT | AUT | 9015361 | 2020-W17 | 20/4/2020 | 26/4/2020 | 2784 | 697 | Surveillance |
| 6 | Austria | AT | AUT | 9015361 | 2020-W18 | 27/4/2020 | 3/5/2020 | 1657 | 407 | Surveillance |
| 7 | Austria | AT | AUT | 9015361 | 2020-W18 | 27/4/2020 | 3/5/2020 | 1350 | 356 | Country_Website |
| 8 | Austria | AT | AUT | 9015361 | 2020-W19 | 4/5/2020 | 10/5/2020 | 2555 | 636 | Country_Website |
| 9 | Austria | AT | AUT | 9015361 | 2020-W20 | 11/5/2020 | 17/5/2020 | 1637 | 381 | Country_Website |
| 10 | Austria | AT | AUT | 9015361 | 2020-W21 | 18/5/2020 | 24/5/2020 | 998 | 202 | Country_Website |

Figure 6: Sample of Exported Sink Weekly File

**Part B (20 marks)**

Build a new project to ingest e-commerce and currency exchange rate data for localized reporting. Source or dataset:

a. E-Commerce Business Transaction from Kaggle
(https://www.kaggle.com/datasets/gabrielramos87/an-online-shop-business)
b. Exchange Rates by Bank Negara Malaysia
(https://apikijangportal.bnm.gov.my/openapi?category=Rates%20and%20Volumes)
c. Country and Currency Code (to find relevant dataset from a trusted source)

**High Level Requirements**

1. To ingest both E-Commerce Business Transaction and Exchange Rates dataset into data lake (5 marks)
2. Map Country field from transaction to respective currency code and retrieve equivalent price in MYR according to date of exchange rate and transaction. (note: assume Price shown is in local currency of purchaser) (10 marks)
3. Get total price for each transaction in local currency and MYR (5 marks)

Expected sink datasets structure:

**Daily Transaction**

| Field | Data Type | Description |
|---|---|---|
| TransactionNo | Integer | Unique number that defines each transaction |
| Date | Date | Transaction date |
| ProductNo | String | Unique number used to identify product |
| ProductName | String | Product name |
| PriceLocal | Decimal | Price in local currency of purchaser's country |
| CurrencyLocal | Decimal | Local 3-digit currency code of purchaser's country |
| PriceConverted | Decimal | MYR equivalent of transacted price |
| CurrencyConverted | Decimal | Default to 'MYR' |
| Quantity | Integer | Quantity of each product per transaction |
| ExchangeRate | Decimal | Currency exchange rate for transacted currency and date |
| TotalPriceLocal | Decimal | Total price = priceLocal * Quantity |
| TotalPriceConverted | Decimal | MYR equivalent of TotalPriceLocal |
| CustomerNo | Integer | Unique number to identify customer |
| Country | String | Country where customer reside |

**Weekly Summary**

Sum of transaction by product and country.

| Field | Data Type | Description |
|---|---|---|
| ProductNo | String | Unique number used to identify product |
| ProductName | String | Product name |
| PriceLocal | Decimal | Price in local currency of purchaser's country |
| TotalPriceLocal | Decimal | Total price = priceLocal * Quantity |
| TotalPriceConverted | Decimal | MYR equivalent of TotalPriceLocal |
| Country | String | Country where customer reside |
| ReportedMonth | Integer | Month of transaction |
| ReportedYear | Integer | Year of transaction |

**Daily Transaction:**

Ingestion

For this project, we are going to need 3 datasets; business transaction from Kaggle, exchange rates from API link and currency code lookup that can be found online. Figure 1 shows how business transaction data ingested into data lake link service that has been created in Part A. For exchange rates, the API link given from BNM do not work as intended, so we are using other exchange rates API (https://api.exchangerate.host/latest). To ingest API dataset, we created a REST link service and provides the base URL of the API (Figure 2). Then, we created a pipeline (Figure 3) to copy data using the link service into data lake in the form of JSON file. Figure 4 display the preview of the JSON file that have been successfully ingested.



Figure 1: Ingesting Business Transaction CSV Dataset to Data Lake





Figure 2: Setting up a REST Linked Service to the Exchange Rates Dataset API URL

Figure 3: Copy data from REST Linked Service into Data Lake



Figure 4: Successfully Ingested Exchange Rates JSON Dataset into Data Lake

Data Flow Pipeline

Figure 5 display the data flow of daily transaction. First, we add 3 dataset sources from ingested data lake into the data flow ("Transaction", "CurrencyCodeLookup" and "Rates"). We join currency code lookup with rates ("ExchangeRates") using left outer join to map country field to the respective currency code. Then, for transaction source, we selected only the needed columns ("selectTransaction") and joined with the lookup that we created before ("joinRate"). From there we need to create 5 new columns; PriceLocal, PriceConverted, ExchangeRate, TotalPriceLocal, TotalPriceConverted. By using Derived Column function, we calculated for each new columns (Figure 6) to be derived from exchange rates information available using the expression builder. Last but not least, we selected only required columns, rearrange them in order as shown in Figure 7 and successfully sink them into a new data lake (Figure 8) as a processed dataset. Figure 9 shows a snippet of the processed dataset after exported into a csv file.

Figure 5: Overall Data Flow of Daily Transaction



Figure 6: Creating Derived Columns According to Requirements



Figure 7: Final Transformation Output for Daily Transaction

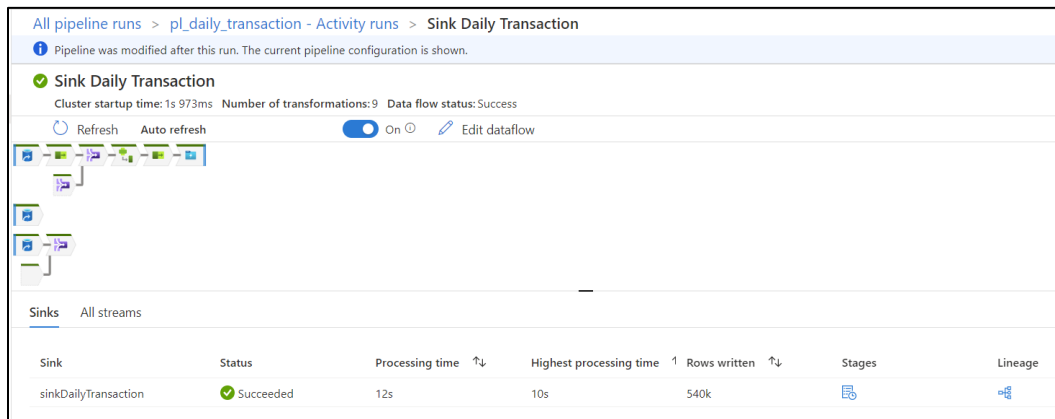Figure 8: Successfully Sink the Daily Transaction Pipeline

| | A | B | C | D | E | F | G | H | I | J | K | L | M | N |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | TransactionNo | Date | ProductNo | ProductName | PriceLocal | CurrencyL | PriceConverted | CurrencyConverted | Quantity | ExchangeRate | TotalPriceLocal | TotalPriceConverted | CustomerNo | Country |
| 2 | 581493 | 9/12/2019 | 22807 | Set Of 6 T-Lights T | 7.06 | EUR | 32.5 | 0.217106376 | 6 | 1.14 | 42.34 | 195 | 12423 | Belgium |
| 3 | 581493 | 9/12/2019 | 22252 | Birdcage Decorati | 8.25 | EUR | 38.01 | 0.217106376 | 12 | 1.14 | 99.04 | 456.16 | 12423 | Belgium |
| 4 | 581493 | 9/12/2019 | 21108 | Fairy Cake Flannel | 8.25 | EUR | 38.01 | 0.217106376 | 18 | 1.14 | 148.56 | 684.24 | 12423 | Belgium |
| 5 | 581493 | 9/12/2019 | 23204 | Charlotte Bag App | 8.25 | EUR | 38.01 | 0.217106376 | 10 | 1.14 | 82.53 | 380.13 | 12423 | Belgium |
| 6 | 581493 | 9/12/2019 | 20724 | Red Retrospot Cha | 8.25 | EUR | 38.01 | 0.217106376 | 10 | 1.14 | 82.53 | 380.13 | 12423 | Belgium |
| 7 | 581493 | 9/12/2019 | 22356 | Charlotte Bag Pink | 8.25 | EUR | 38.01 | 0.217106376 | 10 | 1.14 | 82.53 | 380.13 | 12423 | Belgium |
| 8 | 581493 | 9/12/2019 | 84945 | Multi Colour Silver | 8.25 | EUR | 38.01 | 0.217106376 | 12 | 1.14 | 99.04 | 456.16 | 12423 | Belgium |

Figure 9: Sample of Exported Sink Daily Transaction File

**Weekly Transaction:**

For weekly transaction, the process is similar with the daily transaction with added derived columns (reportedMonth and reportedYear) and removal of some columns. Figure 10 shows the overall data flow of weekly transaction. We created new calculations in the derivedPrices to derive reportedMonth and reportedYear using the expression builder to get the month and the year of transaction respectively. The rest of the process flow is similar to daily transaction where we selected only the required columns (Figure 12) and successfully sink the processed dataset into a new data lake (Figure 13). Figure 14 shows a snippet of the processed dataset for weekly transaction after exported into a csv file.
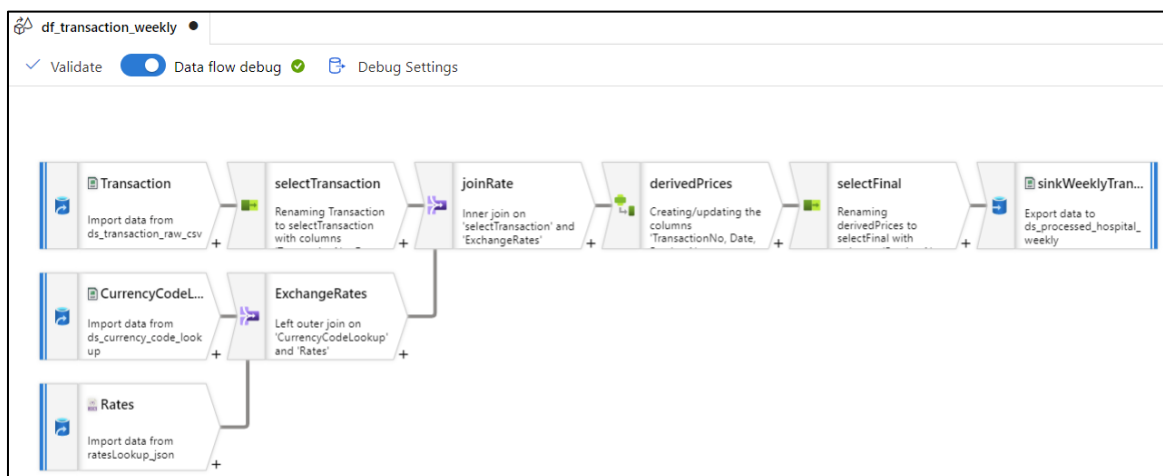

Figure 10: Overall Data Flow of Weekly Transaction

Figure 11: Creating Derived Columns According to Weekly Transaction Requirements



Figure 12: Final Transformation Output for Weekly Transaction



Figure 13: Successfully Sink the Weekly Transaction Pipeline

| | A | B | C | D | E | F | G | H |
|---|---|---|---|---|---|---|---|---|
| 1 | ProductNo | ProductName | PriceLocal | TotalPriceLocal | TotalPriceConverted | Country | ReportedMonth | ReportedYear |
| 2 | 22807 | Set Of 6 T-Lights | 7.06 | 42.34 | 195 | Belgium | 12 | 2019 |
| 3 | 22252 | Birdcage Decorat | 8.25 | 99.04 | 456.16 | Belgium | 12 | 2019 |
| 4 | 21108 | Fairy Cake Flanne | 8.25 | 148.56 | 684.24 | Belgium | 12 | 2019 |
| 5 | 23204 | Charlotte Bag Ap | 8.25 | 82.53 | 380.13 | Belgium | 12 | 2019 |
| 6 | 20724 | Red Retrospot Ch | 8.25 | 82.53 | 380.13 | Belgium | 12 | 2019 |
| 7 | 22356 | Charlotte Bag Pir | 8.25 | 82.53 | 380.13 | Belgium | 12 | 2019 |
| 8 | 84945 | Multi Colour Silve | 8.25 | 99.04 | 456.16 | Belgium | 12 | 2019 |

Figure 14: Sample of Exported Sink Weekly Transaction File

Summary

We are able to ingest both E-Commerce Business Transaction and Exchange Rates dataset into data lake although using different exchange rates API. Next, we are able to map Country field from transaction to respective currency code and retrieve equivalent price in MYR according to date of exchange rate and transaction. Moreover, we are able to calculate total price for each transaction in local currency and MYR and others required fields from the expected sink datasets structure for both daily and weekly transactions. Last but not least, both daily and weekly transactions successfully sink into new data lake using the explained pipelines and data flow.