



UJIAN TENGAH SEMESTER (UTS)

Untuk Memenuhi Salah Satu Syarat Penyelesaian

Mata Kuliah Data Warehouse

Program S-1 Teknik Informatika

Oleh

Toni Hari Wibowo

NIM. 20200801314

PROGRAM STUDI S-1 TEKNIK INFORMATIKA

FAKULTAS ILMU KOMPUTER

UNIVERSITAS ESA UNGGUL

NOVEMBER 2022

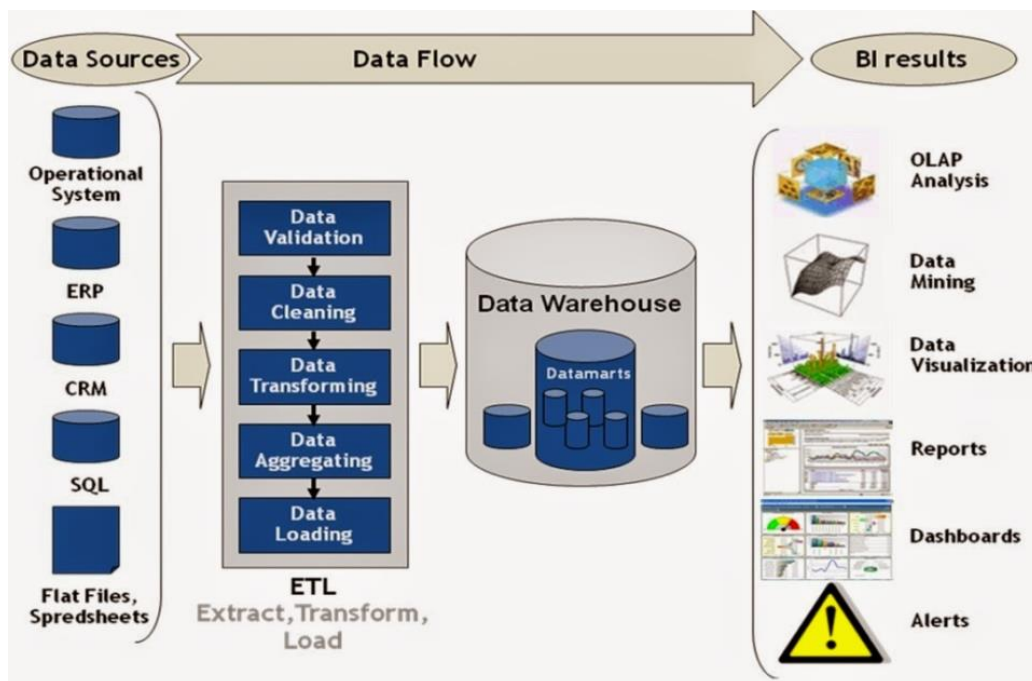
SOAL UTS

Compile lah semua penugasan yang sudah diberikan (dalam group) selama ini sehingga memenuhi ketentuan sebagai berikut:

- Gunakan data riil dari suatu bisnis (industri nya bebas) minimal 5 tahun
- Lakukan analisis kebutuhan dengan minimal mengacu kepada 3 driven untuk mencerminkan kondisi riil dari data yang dimiliki
- Buatlah dokumentasi dari mulai analisis kebutuhan dan disain konseptual dari data yang sudah didapatkan selengkap mungkin

Jawab :

Dari soal diatas maka penulis akan menjawab soal tersebut sebagai berikut, dimana penulis menggunakan dataset yang di ambil dari dataset publik E-commerce Brazilian by Olist. Dataset ini memiliki informasi 100 ribu pesanan dari tahun 2016 hingga 2018 yang dibuat di beberapa pasar di brasil. Fitur-fiturnya memungkinkan melihat pesanan dari berbagai dimensi: dari status pesanan, harga, kinerja pembayaran dan pengiriman hingga lokasi pelanggan, atribut produk, dan akhirnya ulasan yang ditulis oleh pelanggan. Selain itu di dalam dataset ini, olist juga merilis kumpulan data geolokasi yang menghubungkan kode pos brasil dengan koordinat lat/Ing. Dalam UTS ini penulis memproses dataset tersebut melalui alur sebagai berikut :



Gambar 1. Flow

1. Data Sources

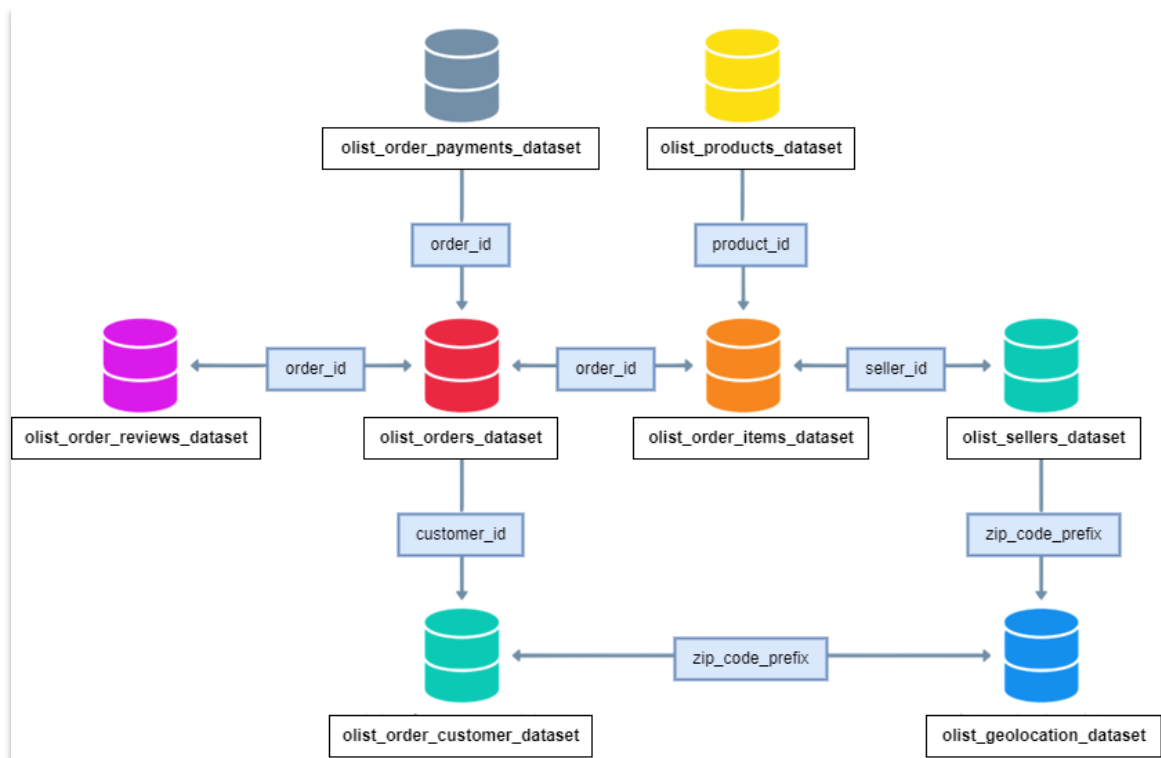
Data ini di ambil dapat dari berbagai sumber data, dengan contoh lain pada pengejaan UTS ini data di ambil dari situs Kaggle dataset publik dalam bentuk Flat Files yang berformatkan **.csv**

2. Data Flow

Dalam proses ini dataset yang sudah didapatkan akan di proses dengan beberapa tahapan layer sebelum dilakukannya load atau dimasukan kedalam data warehouse, tahap-tahap tersebut diantaranya *Extract, Transform, Load* (ETL) :

- **Extract**

Pada layer ini merupakan tahapan *Extracting* data, sebelum data digunakan terlebih dahulu data yang berformatkan file **.csv** tersebut akan dibaca menggunakan program python sehingga terbentuknya dataframe, berikut adalah gambaran schema dataset yang digunakan:

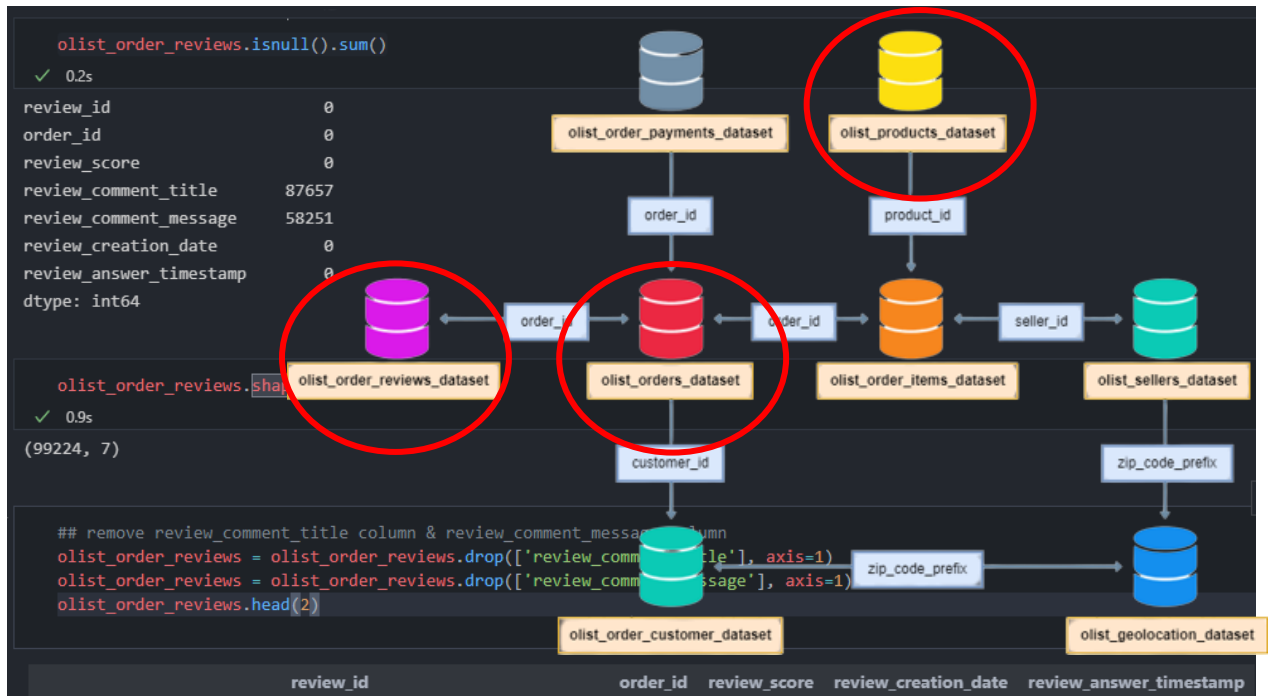


Gambar 2. Schema Dataset

Pada gambar 2 diatas merupakan gambaran dari schema dataset yang berasal dari file e-commerce Brazilian by olist dengan berjumlah 9 file. File-file tersebut nantinya akan diintegrasikan atau digabungkan menjadi satu kesatuan yang utuh melalui foreign key yang terdapat pada dataset tersebut.

- **Transform**

Sebelum proses penggabungan data, dilakukannya transformasi data, hal tersebut berguna untuk melakukan pembersihan dari baris yang kosong pada data yang akan dipakai untuk divisualisasi dan di analisa.



Gambar 3. Cleaning Data

Pada gambar 3 diatas merupakan gambaran pada saat tahap cleaning dataframe, pada saat dilakukannya pengecekan disetiap dataframe didapatkan 3 dataframe yang memiliki baris yang kosong diantaranya yaitu : `olist_order_reviews`, `olist_orders` dan `olist_products`. Dalam program python terdapat beberapa solusi untuk melakukan data cleaning, dan metode tersebut nantinya dapat dipilih dengan menyesuaikan kebutuhan data yang akan dipakai untuk divisualisasi dan di analisa. Berikut adalah metode-metode data cleaning yang terdapat pada program python :

- 1) `Dropna(how='all')`

jika semua kolom tidak mengandung nilai maka metode ini akan menghapus dari kolom tersebut.

- 2) `Fillna(0)`

Lain halnya dengan dropna, jika hanya beberapa nilai yang hilang di kolom, maka metode ini digunakan untuk mengisi nilai NA/NaN menggunakan nilai yang ditentukan.

3) fillna(method='ffill')

Metode ini merupakan metode duplikasi data dengan mengisi nilai yang hilang di kolom dengan mengisi nilai di baris **sebelumnya**.

4) fillna(method='bfill')

Metode ini merupakan metode duplikasi data dengan mengisi nilai yang hilang di kolom dengan mengisi nilai di baris **setelahnya**.

5) interpolate()

Metode ini merupakan metode duplikasi data dengan mengisi nilai yang hilang di kolom dengan mengisi nilai diantara baris **sebelumnya** dan baris **setelahnya**.

- **Load**

Pada layer ini merupakan proses melakukan pengiriman data kedalam postgresQL. Data yang sudah dilakukan pembersihan dan telah di integrasi menjadi satu kesatuan akan di load kedalam database yang ada di postgresQL. Berikut adalah gambaran pada saat melakukan push kedalam database yang ada di postgresQL

The screenshot displays the pgAdmin interface with a SQL query executed in the Query Editor. The query is as follows:

```
## Test the code by select the complete table

drop_column='''ALTER TABLE result_olist_marged DROP COLUMN index''';
engine.execute(drop_column)

sql='''
Select * from result_olist_marged''';

df_sql = pd.read_sql_query(sql,con=engine)
df_sql.head()
```

The results of the query are shown in the Data Output tab, displaying a table with 5 rows and 40 columns. The first few rows are visible:

	order_id		
0	e481f51cbdc54678b7cc49136f2d6af7	9ef432eb	
1	128e10d95713541c87cd1a2e48201934	a20e810	
2	0e7e841ddf8f8f2de2bad69267ecfbcf	26c7ac16	
3	bfc39df4f36c3693f3b63fcbca9e90a	53904dd	
4	5f49f31e537f8f1a496454b48edbe34d	a7260a	

Gambar 4. Push Data into PostgreSQL

3. Business intelligence (BI) Results

Pada tahap ini merupakan penerapan teknik transformasi data dari data mentah menjadi informasi yang berguna dan bermakna untuk tujuan analisis, salah satu teknik yang digunakan untuk melakukan pengerjaan SOAL UTS ini yaitu membuat dashboard di PowerBI dengan cara menghubungkan database server PostgreSQL dengan PowerBI. Data yang telah disimpan di dalam PostgreSQL nantinya akan di panggil kembali kedalam PowerBI, sehingga perlu dihubungkan antara PostgreSQL dengan PowerBI yang nantinya di dalam PowerBI data tersebut akan di visualisasikan agar dapat dengan mudah di analisis. Berikut adalah gambaran proses hasil load data dari PostgreSQL kedalam PowerBI :

The screenshot shows the PowerBI interface during a data load process. The Navigator pane on the left lists the following tables:

- localhost: Tugas03_DataWareHou...
- ☒ public.geolocation_customers
- ☒ public.geolocation_sellers
- ☒ public.result_olist_marged

The main area displays a preview of the 'public.result_olist_marged' table. The table has the following columns: order_id, customer_id, order_status, and order_purc. The data is truncated, showing 15 rows of sample data.

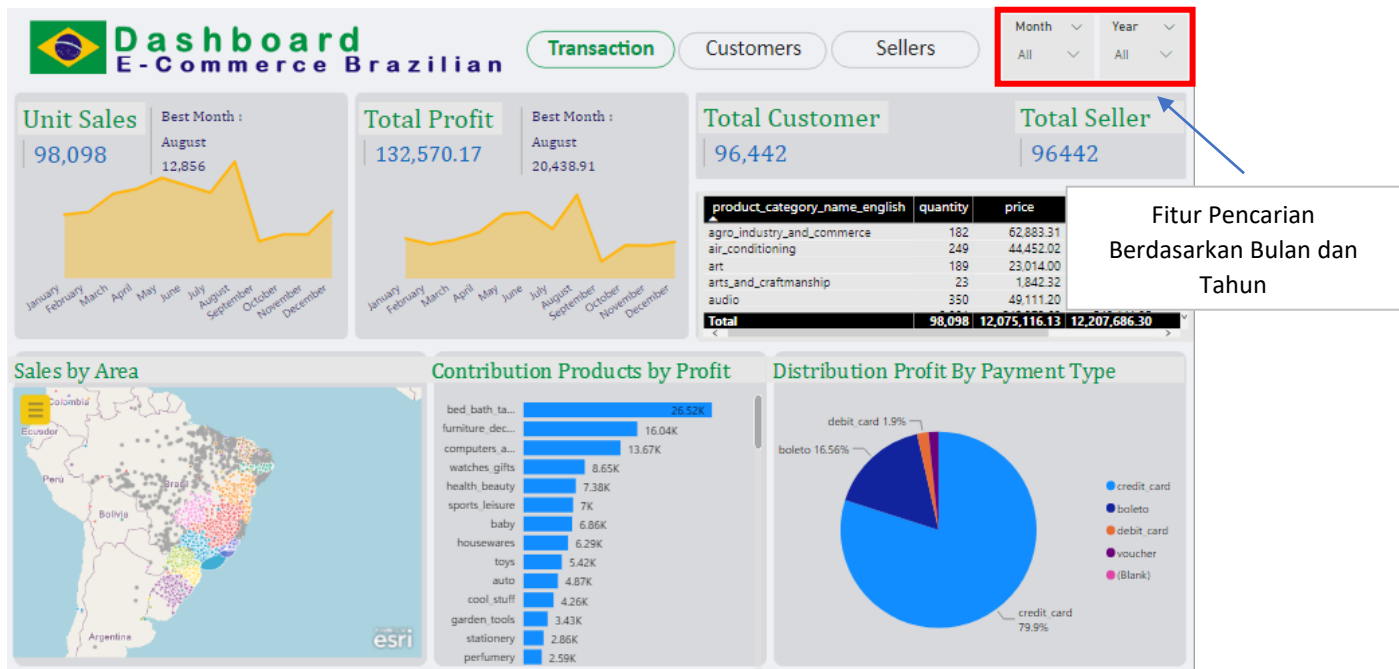
order_id	customer_id	order_status	order_purc
517b4b7e7f723df3693ae71473078d73	13c14842e78da8bb1e408509619f291f	delivered	2017-10-07
fa534840c9b5b3a7052a37a863ea1790	d14fd2f32ddd1cab26e66bbd544bd17a	delivered	2018-01-25
f4f8eefd039e279c8fac5a307cb469c3	c5156853ded5199521864606a85ff104	delivered	2017-06-01
60d5461dab300c5e833a29ac2f534ac4	971690d4a9c0044682ddf18ae5b32f1e	delivered	2017-07-18
4df371955120c03af473cf4160cb6c65	6a0b64aac6b11f3b988161dc7dee946f	delivered	2018-08-17
a1f24dcf6e5fc6b0f0cd0b265381767e	15f73cd40ee94734ce2f3ebce168f1ba	delivered	2017-09-27
1710c56069303e33d5014ad19ff1d402	6950e00305cacdb75bd497302c430d68	delivered	2017-11-15
2ab3c22d797ae5a4199d53021ac9e78b	084276e19165eb374ea209756511156a	delivered	2017-06-08
2306aeccf5b70ab26762993fdf2f37b9	33cdf7dd400e36954b11995afd7f0aa7	delivered	2017-05-07
27ee3c5633d367f617c3dd7c5037edd7	0d68f4d8de943d4d32ec8f72f1d88230	delivered	2018-03-31
15565710182cd35455ee0abf4b701ea9	5e7bb90d4342e8086e9a9acc1360a1e0	delivered	2018-02-17
e1a9f566dc38eac4f78556976596a43f	93dfd5dcfc590abda99d2e9f3960ea5	delivered	2018-07-01

Below the table, a message states: "The data in the preview has been truncated due to size limits." At the bottom, there are buttons for "Load", "Transform Data", and "Cancel".

Gambar 5. Load Data Into PowerBI

Hasil

A. Halaman Transaction



Gambar 6. Dashboard PowerBI, Halaman Transaction

Pada gambar 6 diatas merupakan tampilan dari Dashboard yang telah dibuat dan pada gambar tersebut merupakan tampilan informasi dari Halaman Transaction seperti gambar 6 diatas, penulis membuat beberapa fitur salah satunya untuk pencarian total transaksi berdasarkan bulan dan tahun yang ingin di tentukan, fitur ini penulis buat berdasarkan bulan dan tahun yang di ambil dari kolom **order_delivered_customer** yang ada pada database. Selain itu masih ada beberapa informasi yang penulis tampilkan pada halaman dashboard Transaksi seperti yang ada pada gambar Pada **Gambar 6. Dashboard PowerBI, Halaman Transaction**, yaitu berupa :

- **Unit Sales**

Informasi Unit Sales yang penulis buat menggunakan visualisasi multi_row card dan juga Area Chart, dari visualisasi multi_row card menggunakan penjumlahan keseluruhan data dari variable **order_item_id** untuk menghasilkan informasi, selain itu dari visualisasi Area Chart penulis juga menggunakan **order_item_id** yang di bandingkan dengan variable **order_delivered_date** dan di filter berdasarkan top 1 untuk pencarian bulan terbaik pada Unit Sales.

- **Total Profit**

Informasi Total Profit yang penulis buat menggunakan visualisasi multi_row card dan juga Area Chart, dari visualisasi multi_row card menggunakan penjumlahan keseluruhan data dari variable profit, dimana variable profit dibuat setelah pembuatan variable revenue, variable tersebut penulis buat terlebih dahulu agar nantinya dapat digunakan untuk memberikan informasi berupa total profit, berikut adalah rumus yang penulis gunakan :

$\text{revenue} = \text{SUMX}(\text{'public result_olist_marged'}, \text{'public result_olist_marged'}[\text{order_item_id}] * \text{'public result_olist_marged'}[\text{price}])$
--

Table 1. Rumus penghitungan revenue

$\text{profit} = \text{SUMX}(\text{'public result_olist_marged'}, \text{'public result_olist_marged'}[\text{revenue}] - \text{'public result_olist_marged'}[\text{price}])$

Table 2. Rumus penghitungan profit

- **Total Customer**

Informasi Total Customer yang penulis buat menggunakan visualisasi **multi_row card** yang menggunakan penjumlahan keseluruhan data dari variable **customer_unique_id**, keputusan variable yang penulis gunakan untuk pembuatan informasi tersebut dikarenakan bahwa setiap customer memiliki id yang unique, sehingga tidak ada duplikasi atau redundan pada penghitungan data.

- **Total Seller's**

Informasi Total Seller's yang penulis buat menggunakan visualisasi **multi_row card** yang menggunakan penjumlahan keseluruhan data dari variable seller_id, dari data yang ada pada kolom seller_id penulis filter terlebih dahulu agar nantinya tidak ada id yang redundan untuk di jumlahkan.

- **Tabel**

Informasi pada Tabel yang penulis buat terisi dari **product name, quantity, price, revenue, profit hingga gross profit margin(GPM)**. Untuk informasi quantity penulis ambil dari order_item_id dengan alasan penulis berasumsi bahwa dari setiap item yang di

order, maka dari adanya pemesanan barang dan jumlah barang yang harus di sediakan atau dengan kata lain bahwa data dari E-commerce ini melakukan pemesanan seperti **Purches Order**, dimana memesan terlebih dahulu baru adanya penghitungan dari quantity product yang harus di sediakan. penulis mengambil keputusan berikut dikarenakan pada dataset yang penulis gunakan tidak terdapat quantity barang, sehingga penulis mengambil keputusan tersebut. Selain itu Informasi yang penulis tampilkan yaitu quantity, dan revenue penulis gunakan berdasarkan pembuatan variable pada **Table 1**. Rumus penghitungan revenue

Dan **Table 2**. Rumus penghitungan profit, dan ada juga variable baru yang penulis buat yaitu **Gross Profit margin(GPM)**, beriku penghitungan rumus dari GPM :

$\text{GPM (\%)} = \text{DIVIDE}([\text{profit}], [\text{revenue}])$
--

Table 3. Rumus penghitungan Gross Profit Margin(GPM)

- **Sales by Area**

Pada informasi Sales by Area penulis menampilkan persebaran peta dari product yang di order, Informasi tersebut penulis buat menggunakan visualisasi ArcGIS Maps yang berisikan value dari **order_item_id** dan di bandingkan dengan **customer_city** lalu di filter berdasarkan **customer_state** agar peta dari persebaran penjualan product dapat terlihat dengan jelas pada setiap daerah nya.

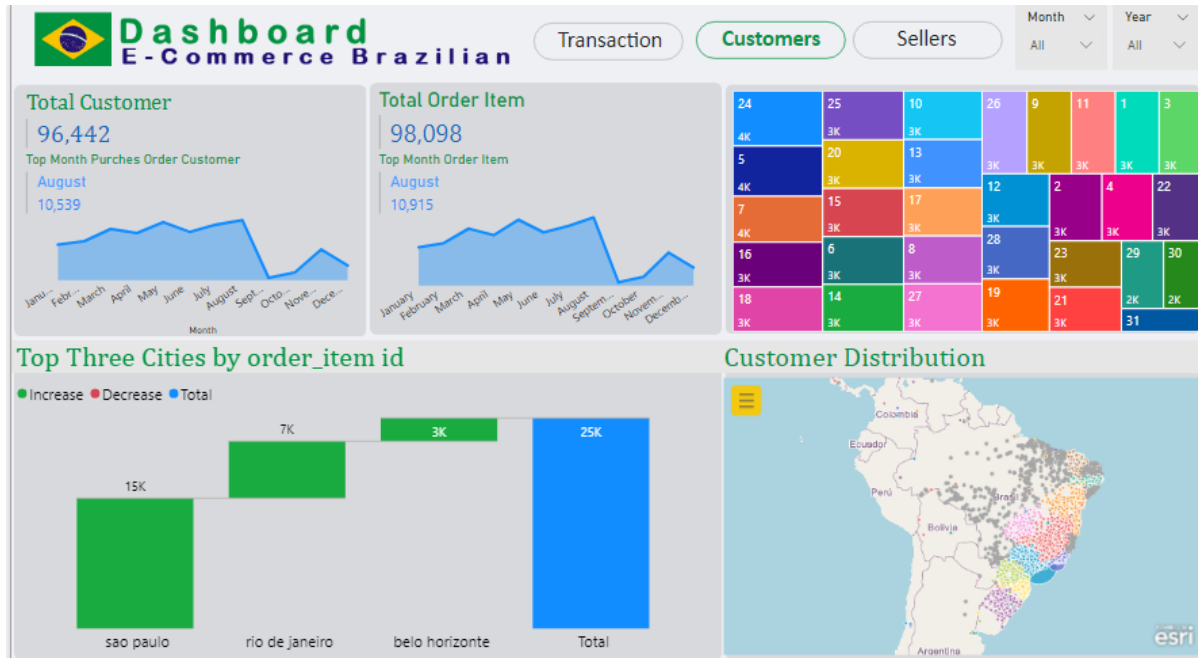
- **Contribution Products by Profit**

Informasi yang di tampilkan pada Contribution Products by Profit penulis buat menggunakan visualisasi Bar Chart, yang berisikan value berdasarkan perbandingan dari **product_category** dengan **profit**.

- **Distribution Profit by Payment Type**

Pada Informasi Distributin Profit by Payment Type penulis buat menggunakan visualisasi pie chart yang berisikan value dari **profit** dan di bandingkan dengan **payment_type**

B. Halaman Customer



Gambar 7. Dashboard PowerBI, Halaman Customer

Pada Halaman Customer seperti gambar 7 diatas, juga penulis membuat beberapa fitur salah satunya untuk pencarian total customer berdasarkan bulan dan tahun yang ingin di tentukan, selain itu masih ada beberapa informasi yang penulis tampilkan pada halaman dashboard Transaksi seperti yang ada pada gambar Pada Gambar 7. Dashboard PowerBI, Halaman Customer, yaitu berupa :

- **Total Customer**

Informasi Total Customer yang penulis buat menggunakan visualisasi multi_row card yang menggunakan penjumlahan keseluruhan data dari variable customer_unique_id, keputusan variable yang penulis gunakan untuk pembuatan informasi tersebut dikarenakan bahwa setiap customer memiliki id yang unique, sehingga tidak ada duplikasi atau redundan pada penghitungan data. Selain itu informasi yang penulis berikan pada Total Customer yaitu berupa:

- **Grafik**

Visualisasi yang penulis gunakan untuk pembuatan grafik yaitu menggunakan Area Chart dimana pada informasi tersebut dapat menampilkan kenaikan dan penurunan dari bulan pada saat customer melakukan pemesanan pembelian.

- **Top Month Purches Order Customer**

Agar lebih spesifik pada visualisas grafik yang di tampilkan sebelumnya, penulis menampilkan Informasi berupa visualisasi multi-row card di mana nilai dari visualisasi tersebut di isi dari variable `top_month_purches_order` yang di filter berdasarkan tingkat pertama pada keseluruhan bulan pemesanan pembelian.

- **Total Order Item**

Informasi `Total_order_item` yang penulis buat menggunakan visualisasi multi_row card yang menggunakan penjumlahan keseluruhan data dari variable `order_item_id`, , selain itu dari visualisasi Area Chart penulis juga menggunakan `order_item_id` yang di bandingkan dengan variable `order_approved_at` dan di filter berdasarkan top 1 untuk pencarian bulan terbaik pada Total Order Item.

- **Jumlah pemesanan berdasarkan tanggal**

Pada informasi jumlah pemesanan berdasarkan tanggal, dibuat dengan menggunakan visualisasi treemap, yang berisi dari nilai penjumlahan total customer dan di bandingkan dengan hari pada kolom `order_approved_at`.

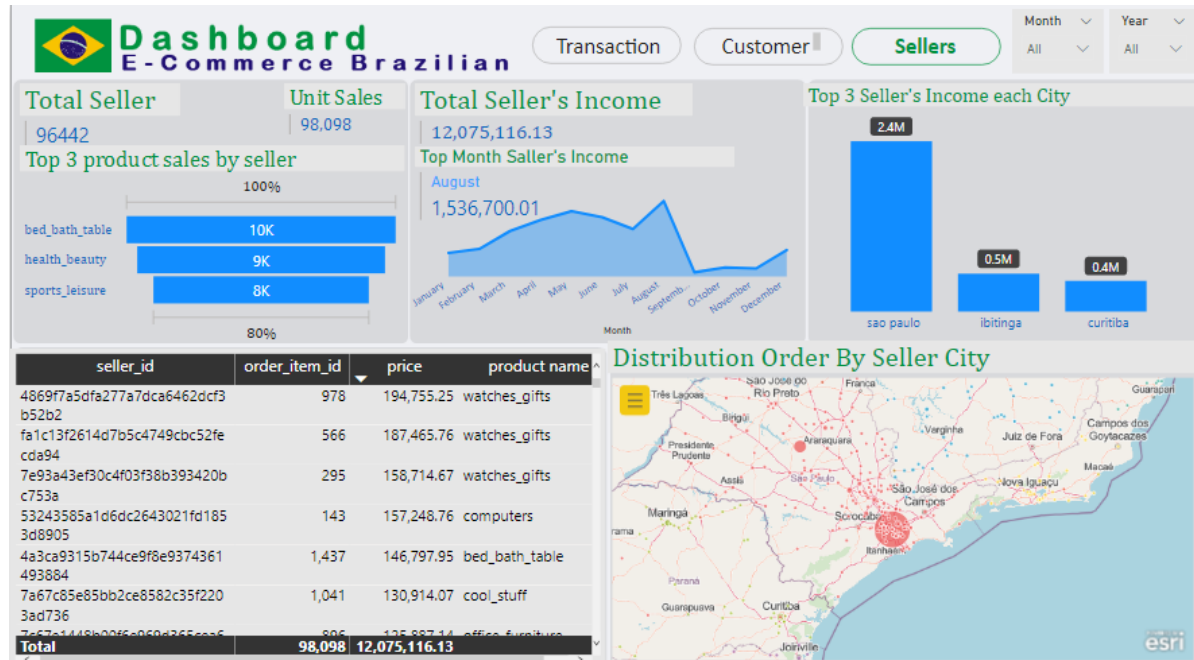
- **Top Three Cities by Order_item_id**

Dari Informasi Top Three Cities adalah informasi yang menampilkan 3 kota teratas berdasarkan banyaknya jumlah order yang terjadi di dalam database, informasi ini dibuat menggunakan visualisasi waterfall chart yang berisi nilai dari `order_item_id` yang di bandingkan dengan `customer_city` lalu di filter berdasarkan 3 kota teratas yang memiliki banyak order.

- **Customer Distribution**

Pada informasi Customer Distribution penulis menampilkan persebaran peta dari customer, dimana Informasi tersebut penulis buat menggunakan visualisasi ArcGIS Maps yang berisikan value dari jumlah customer lalu di cari berdasarkan `customer_city` menggunakan latitude dan longitude pada setiap customer.

C. Halaman Sellers



Gambar 8. Dashboard PowerBI, Halaman Sellers

Pada Halaman Seller's seperti gambar 8 diatas, juga penulis membuat beberapa fitur dihalaman ini yaitu berupa :

- **Total Seller's**

Informasi Total Seller's yang penulis buat menggunakan visualisasi multi_row card yang menggunakan penjumlahan keseluruhan data dari variable seller_id, dari data yang ada pada kolom seller_id penulis filter terlebih dahulu agar nantinya tidak ada id yang redundan untuk di jumlahkan.

- **Unit Sales**

Informasi Unit Sales yang penulis buat menggunakan visualisasi multi_row card dan juga Area Chart, dari visualisasi multi_row card menggunakan penjumlahan keseluruhan data dari variable order_item_id untuk menghasilkan informasi total dari item yang di order.

- **Top 3 Product sales by sellers**

Informasi Tope 3 Product sales by sellers yang penulis buat menggunakan visualisasi funnel, dimana visualisasi tersebut menggambarkan tingkatan dari product yang banyak

di jual oleh seller. Dari visualisasi tersebut penulis menggunakan nilai `order_item_id` yang di bandingkan berdasarkan `product_category_name`

- **Total Seller's Income**

Informasi Total Seller's Income yang penulis buat menggunakan visualisasi `multi_row card` yang menggunakan penjumlahan keseluruhan data dari variable `price`, dimana keputusan dari penggunaan variable tersebut, penulis berasumsi bahwa `price` adalah harga awal yang di jual dari seller's maka dari itu E-Commerce ini menjual dengan harga lebih dari harga awal yang dijual dari sellers.

- **Grafik**

Pada Informasi Grafik ini penulis menampilkan Informasi tersebut berupa grafik penaikan dan penurunan dari setiap bulan yang di hasilkan oleh seller's, untuk membuat visualisasi ini menggunakan `Area Chart` yang berisi nilai dari `price` yang di bandingkan dengan bulan pada setaip `order_delivered`.

- **Top Month Saller's Income**

Untuk menghasilkan informasi yang lebih spesifik pada grafik tersebut maka dari itu penulis membuat visualisasi `Top Month Saller's Income`, dimana informasi ini memberikan berupa total angka `Income` dari bulan yang tertinggi.

- **Top 3 Seller's Income each City**

Dari Informasi `Top 3 Seller's Income each City` adalah informasi yang menampilkan 3 kota teratas berdasarkan banyaknya jumlah `Income` yang terjadi di dalam database, informasi ini dibuat menggunakan visualisasi `Line and Cluster Column` atau sama dengan diagram batang, dari visualisasi ini berisi nilai dari total keseluruhan `price` yang di bandingkan dengan `seller's_city` lalu di filter berdasarkan 3 kota teratas yang memiliki banyak `Income`.

- **Table**

Informasi pada Tabel yang penulis buat pada halaman Seller's terisi dari `Seller_id`, `Order_item_id`, `Price`, dan `Product_name`, penulis memilih menampilkan `seller_id` dikarenakan di dataset yang penulis gunakan tidak terdapat dari nama seller, maka dari itu penulis membuat keputusan untuk menampilkan kolom tersebut, dari informasi tersebut nantinya dapat di cari lebih lanjut untuk penelusuran nama dari seller berdasarkan pencarian `seller_id`, dengan catatan jika dataset yang digunakan sudah benar-benar lengkap atau sesuai dengan database pada E-Commerce Brazilian.

- **Distribution Order by Seller City**

Pada informasi Distribution Order by Seller City penulis menampilkan persebaran peta, dimana Informasi tersebut penulis buat menggunakan visualisasi ArcGIS Maps yang berisikan value dari order_item_id lalu di cari berdasarkan seller_city menggunakan latitude dan longitude pada setiap customer dan di filter berdasarkan seller_state agar informasi dapat terpapar dengan jelas dari setiap perbedaan daerah yang ada pada di peta tersebut.

Kesimpulan

Dari penjelasan pada saat proses pengambilan data hingga pembuatan informasi data dalam bentuk Visual di dalam powerBI penulis menyimpulkan, Dashboard Visual dari dataset E-Commerce Brazilian By Olist telah dibuat, dashboard tersebut dibuat guna untuk memenuhi syarat dalam pengerjaan SOAL UTS ini dan juga berguna untuk menampilkan transformasi data dari data mentah menjadi informasi yang berguna dan bermakna untuk tujuan analisis.

Dari hasil analisis yang penulis lakukan, dimana pada dataset E-Commerce Brazilian By Olist, jumlah setiap tahun seperti customer, product, seller dan order barang terus meningkat, dari hasil peningkatannya tersebut banyak terjadi penjualan di bulan **Agustus** dari tahun **2016 s/d 2018**, dan untuk tanggal yang paling unggul pada saat pemesanan yaitu pada tanggal **24 di bulan Agustus**. Dari situ penulis berasumsi bahwa ada **hari penting** di dalam bulan Agustus sehingga dapat menghasilkan jumlah penjualan yang sangat tinggi, selain itu category product yang paling laris di jual yaitu category product **bed bath table** yang berada di tingkatan pertama dengan jumlah **9,593** dari total keseluruhan category_product **98,098**.