# Assignment 1

## Connect to postgre local db

Melakukan import pandas dan psycopg2 library, dan membuat koneksi db ke posgresql pada ubuntu local.

```python
import pandas as pd
import psycopg2

CONNECT_DB = "host=localhost port=5432 dbname=assignment1 user=cloud_user password=cloud_user"
```

✓ 0.0s                                                                                                          Python

## Create table for both csv file (customer_detail and transformation)

Melakukan create table query dengan menambahkan primary key pada column customer dan merchant pada table customer_detail, dan primary key source dan target pada table transformation. Primary key pada column customer / source saja akan menimbulkan error (duplicate value) sehingga memilih untuk membuat 2 primary key.

```python
create_table_query = '''CREATE TABLE IF NOT EXISTS customer_detail (
        step INT ,
        customer VARCHAR(15),
        age VARCHAR(10),
        gender VARCHAR(6),
        zipcodeOri VARCHAR(10),
        merchant VARCHAR(15),
        zipMerchant VARCHAR(10),
        category VARCHAR(30),
        amount FLOAT,
        fraud INT,
        PRIMARY KEY (customer,merchant)
    ); '''
try:
    # Make connection to db
    cxn = psycopg2.connect(CONNECT_DB)

    # Create a cursor to db
    cur = cxn.cursor()

    # Send sql query to request
    cur.execute(create_table_query)
    records = cxn.commit()

except (Exception, psycopg2.Error) as error :
    print ("Error while connecting to PostgreSQL", error)

finally:
    #closing database connection.
    if(cxn):
        cur.close()
        cxn.close()
        print("PostgreSQL connection is closed")

print(f'Records:\n {records}')
```

✓ 2.0s

```python
create_table_query = '''CREATE TABLE IF NOT EXISTS transformation (
        Source VARCHAR(15),
        Target VARCHAR(15),
        Weight FLOAT,
        typeTrans VARCHAR(30),
        fraud INT,
        PRIMARY KEY (source,Target)
    ); '''
try:
    # Make connection to db
    cxn = psycopg2.connect(CONNECT_DB)

    # Create a cursor to db
    cur = cxn.cursor()

    # Send sql query to request
    cur.execute(create_table_query)
    records = cxn.commit()

except (Exception, psycopg2.Error) as error :
    print ("Error while connecting to PostgreSQL", error)

finally:
    #closing database connection.
    if(cxn):
        cur.close()
        cxn.close()
        print("PostgreSQL connection is closed")

print(f'Records:\n {records}')
```

✓ 2.0s

## Add data to table

Menambahkan data pada kedua table dengan data csv yang telah diberikan. Inserting data dilakukan dengan membuat temporary table terlebih dahulu kemudian meng copy data dari temporary table ke main table dan Ketika ada conflik akan di skip.



**SQL fetching data function**



**Data fetching to dataframe**

**Data pada dataframe**



**Drop duplicate and null data**

Menghapus data duplikat dan kosong. Dilakukan karena terdapat beberapa data duplikat pada primary key(customer_detail dan source).

```
df_customer_detail.drop_duplicates(inplace=True)
df_transformation.drop_duplicates(inplace=True)
```
✓ 0.1s

```
df_customer_detail.isnull().sum()
```
✓ 0.0s

```
step          0
customer      0
age           0
gender        0
zipcodeori    0
merchant      0
zipmerchant   0
category      0
amount        0
fraud         0
dtype: int64
```

```
df_transformation.isnull().sum()
```
✓ 0.0s

```
source       0
target       0
weight       0
typetrans    0
fraud        0
dtype: int64
```

**Removing whitespace and singlequote**

Menghapus data yang memiliki spasi kosong (` `) dan menghapus tanda kutip satu , sehingga data pada kolom terlihat lebih rapih dan mudah di identifikasi.

```
                                          (property) str: StringMethods[Series, DataFrame]
df_customer_detail['customer'] = df_customer_detail['c
df_customer_detail['age'] = df_customer_detail['age'].str.strip().str[1:-1]
df_customer_detail['gender'] = df_customer_detail['gender'].str.strip().str[1:-1]
df_customer_detail['zipcodeori'] = df_customer_detail['zipcodeori'].str.strip().str[1:-1]
df_customer_detail['merchant'] = df_customer_detail['merchant'].str.strip().str[1:-1]
df_customer_detail['zipmerchant'] = df_customer_detail['zipmerchant'].str.strip().str[1:-1]
df_customer_detail['category'] = df_customer_detail['category'].str.strip().str[1:-1]

df_customer_detail.head()
```
✓ 0.2s

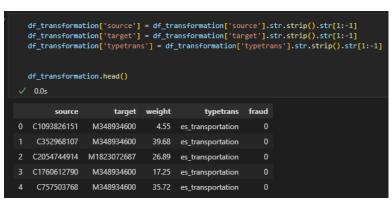| | step | customer | age | gender | zipcodeori | merchant | zipmerchant | category | amount | fraud |
|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 0 | C1093826151 | 4 | M | 28007 | M348934600 | 28007 | es_transportation | 4.55 | 0 |
| 1 | 0 | C352968107 | 2 | M | 28007 | M348934600 | 28007 | es_transportation | 39.68 | 0 |
| 2 | 0 | C2054744914 | 4 | F | 28007 | M1823072687 | 28007 | es_transportation | 26.89 | 0 |
| 3 | 0 | C1760612790 | 3 | M | 28007 | M348934600 | 28007 | es_transportation | 17.25 | 0 |
| 4 | 0 | C757503768 | 5 | M | 28007 | M348934600 | 28007 | es_transportation | 35.72 | 0 |

```
df_transformation['source'] = df_transformation['source'].str.strip().str[1:-1]
df_transformation['target'] = df_transformation['target'].str.strip().str[1:-1]
df_transformation['typetrans'] = df_transformation['typetrans'].str.strip().str[1:-1]

df_transformation.head()
```
✓ 0.0s

| | source | target | weight | typetrans | fraud |
|---|---|---|---|---|---|
| 0 | C1093826151 | M348934600 | 4.55 | es_transportation | 0 |
| 1 | C352968107 | M348934600 | 39.68 | es_transportation | 0 |
| 2 | C2054744914 | M1823072687 | 26.89 | es_transportation | 0 |
| 3 | C1760612790 | M348934600 | 17.25 | es_transportation | 0 |
| 4 | C757503768 | M348934600 | 35.72 | es_transportation | 0 |

**Github repo :** https://github.com/ChristopherRsl/df_assignment1