# Exporting data to Sql

```python
In [1]:  import pandas as pd
         import mysql.connector
         import os

         # List of CSV files and their corresponding table names
         csv_files = [
             ('customers.csv', 'customers'),
             ('orders.csv', 'orders'),
             ('sellers.csv', 'sellers'),
             ('products.csv', 'products'),
             ('geolocation.csv', 'geolocation'),
             ('payments.csv', 'payments'),
             ('order_items.csv', 'order_items') # Added payments.csv for specific har
         ]

         # Connect to the MySQL database
         conn = mysql.connector.connect(
             host='localhost',
             user='root',
             password='rushil',
             database='ecommerce'
         )
         cursor = conn.cursor()

         # Folder containing the CSV files
         folder_path = 'C:/Users/RUSHIL/OneDrive/Desktop/PROJECT/E-COMMERCE/NEWPROJEC

         def get_sql_type(dtype):
             if pd.api.types.is_integer_dtype(dtype):
                 return 'INT'
             elif pd.api.types.is_float_dtype(dtype):
                 return 'FLOAT'
             elif pd.api.types.is_bool_dtype(dtype):
                 return 'BOOLEAN'
             elif pd.api.types.is_datetime64_any_dtype(dtype):
                 return 'DATETIME'
             else:
                 return 'TEXT'

         for csv_file, table_name in csv_files:
             file_path = os.path.join(folder_path, csv_file)

             # Read the CSV file into a pandas DataFrame
             df = pd.read_csv(file_path)

             # Replace NaN with None to handle SQL NULL
             df = df.where(pd.notnull(df), None)

             # Debugging: Check for NaN values
             print(f"Processing {csv_file}")
```

```python
    print(f"NaN values before replacement:\n{df.isnull().sum()}\n")

    # Clean column names
    df.columns = [col.replace(' ', '_').replace('-', '_').replace('.', '_')

    # Generate the CREATE TABLE statement with appropriate data types
    columns = ', '.join([f'`{col}` {get_sql_type(df[col].dtype)}' for col in
    create_table_query = f'CREATE TABLE IF NOT EXISTS `{table_name}` ({colum
    cursor.execute(create_table_query)

    # Insert DataFrame data into the MySQL table
    for _, row in df.iterrows():
        # Convert row to tuple and handle NaN/None explicitly
        values = tuple(None if pd.isna(x) else x for x in row)
        sql = f"INSERT INTO `{table_name}` ({', '.join(['`' + col + '`' for
        cursor.execute(sql, values)

    # Commit the transaction for the current CSV file
    conn.commit()

# Close the connection
conn.close()
```

```
Processing customers.csv
NaN values before replacement:
customer_id                0
customer_unique_id         0
customer_zip_code_prefix   0
customer_city              0
customer_state             0
dtype: int64

Processing orders.csv
NaN values before replacement:
order_id                         0
customer_id                      0
order_status                     0
order_purchase_timestamp         0
order_approved_at              160
order_delivered_carrier_date    1783
order_delivered_customer_date   2965
order_estimated_delivery_date    0
dtype: int64

Processing sellers.csv
NaN values before replacement:
seller_id               0
seller_zip_code_prefix  0
seller_city             0
seller_state            0
dtype: int64

Processing products.csv
NaN values before replacement:
product_id                    0
product category            610
product_name_length         610
product_description_length  610
product_photos_qty          610
product_weight_g              2
product_length_cm             2
product_height_cm             2
product_width_cm              2
dtype: int64

Processing geolocation.csv
NaN values before replacement:
geolocation_zip_code_prefix  0
geolocation_lat              0
geolocation_lng              0
geolocation_city             0
geolocation_state            0
dtype: int64

Processing payments.csv
NaN values before replacement:
order_id             0
payment_sequential   0
payment_type         0
```

```
payment_installments       0
payment_value              0
dtype: int64

Processing order_items.csv
NaN values before replacement:
order_id                   0
order_item_id              0
product_id                 0
seller_id                  0
shipping_limit_date        0
price                      0
freight_value              0
dtype: int64
```

# Importing Libraries

```python
In [5]:  import pandas as pd
         import matplotlib.pyplot as plt
         import seaborn as sns
         import mysql.connector


         db = mysql.connector.connect(host = "localhost",
                                      username = "root",
                                      password = "rushil",
                                      database = "ecommerce")

         cur = db.cursor()
```

# List all unique cities where customers are located.

```python
In [7]:  query = """ select distinct customer_city from customers """

         cur.execute(query)

         data = cur.fetchall()

         df = pd.DataFrame(data)
         df.head()
```

|   | 0 |
|---|---|
| **0** | franca |
| **1** | sao bernardo do campo |
| **2** | sao paulo |
| **3** | mogi das cruzes |
| **4** | campinas |

# The number of orders placed in 2023

In [9]:
```python
query = """ select count(order_id) from orders where year(order_purchase_tim

cur.execute(query)

data = cur.fetchall()

"total orders placed in 2017 are", data[0][0]
```

Out[9]:  ('total orders placed in 2017 are', 45101)

# The total sales per category.

In [11]:
```python
query = """ select upper(products.product_category) category,
round(sum(payments.payment_value),2) sales
from products join order_items
on products.product_id = order_items.product_id
join payments
on payments.order_id = order_items.order_id
group by category
"""

cur.execute(query)

data = cur.fetchall()

df = pd.DataFrame(data, columns = ["Category", "Sales"])
df
```

| | Category | Sales |
|---|---|---|
| **0** | PERFUMERY | 506738.66 |
| **1** | FURNITURE DECORATION | 1430176.39 |
| **2** | TELEPHONY | 486882.05 |
| **3** | BED TABLE BATH | 1712553.67 |
| **4** | AUTOMOTIVE | 852294.33 |
| **...** | ... | ... |
| **69** | CDS MUSIC DVDS | 1199.43 |
| **70** | LA CUISINE | 2913.53 |
| **71** | FASHION CHILDREN'S CLOTHING | 785.67 |
| **72** | PC GAMER | 2174.43 |
| **73** | INSURANCE AND SERVICES | 324.51 |

74 rows × 2 columns

## Percentage of orders that were paid in installments.

```python
query = """ select ((sum(case when payment_installments >= 1 then 1
else 0 end))/count(*))*100 from payments
"""

cur.execute(query)

data = cur.fetchall()

"the percentage of orders that were paid in installments is", data[0][0]
```

('the percentage of orders that were paid in installments is',
Decimal('99.9981'))

## The year-over-year growth rate of total sales.

```python
query = """with a as(select year(orders.order_purchase_timestamp) as years,
round(sum(payments.payment_value),2) as payment from orders join payments
on orders.order_id = payments.order_id
group by years order by years)

select years, ((payment - lag(payment, 1) over(order by years))/
lag(payment, 1) over(order by years)) * 100 from a"""

cur.execute(query)
data = cur.fetchall()
```

```
df = pd.DataFrame(data, columns = ["years", "yoy % growth"])
df
```

Out[51]:

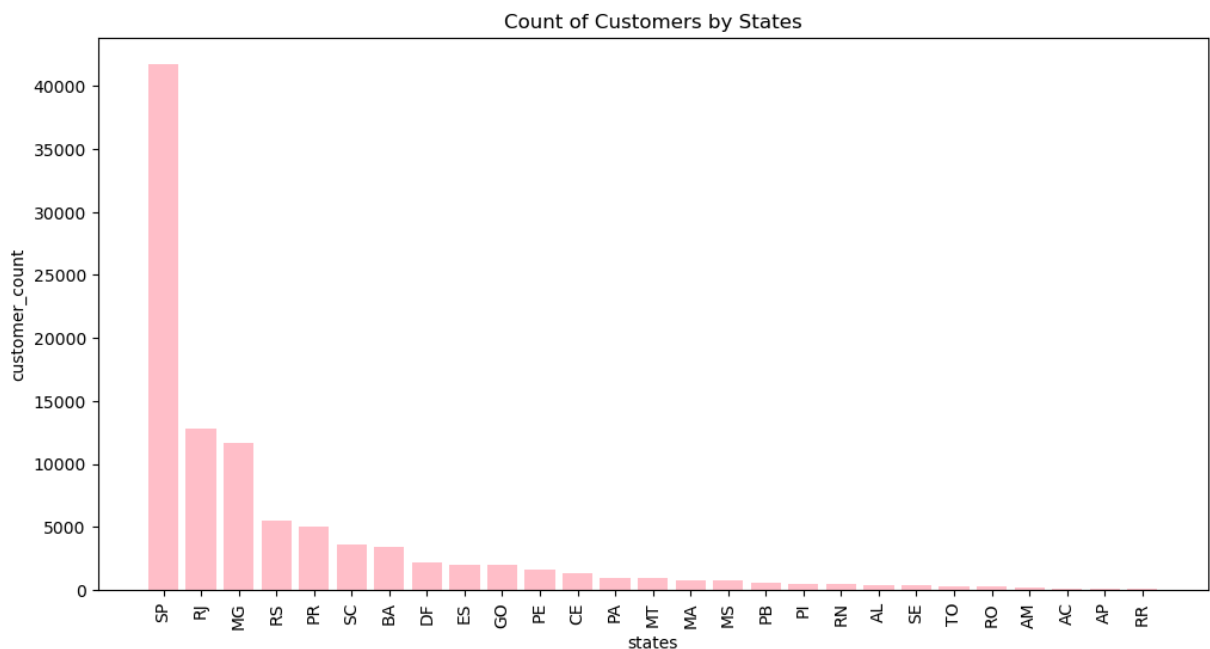|   | years | yoy % growth |
|---|-------|--------------|
| 0 | 2016  | NaN          |
| 1 | 2017  | 12112.703761 |
| 2 | 2018  | 20.000924    |

## Number of customers from each state.

In [43]:
```
query = """ select customer_state ,count(customer_id)
from customers group by customer_state
"""

cur.execute(query)

data = cur.fetchall()
df = pd.DataFrame(data, columns = ["state", "customer_count" ])
df = df.sort_values(by = "customer_count", ascending= False)

plt.figure(figsize = (12,6))
plt.bar(df["state"], df["customer_count"] , color = 'pink')
plt.xticks(rotation = 90)
plt.xlabel("states")
plt.ylabel("customer_count")
plt.title("Count of Customers by States")
plt.show()
```
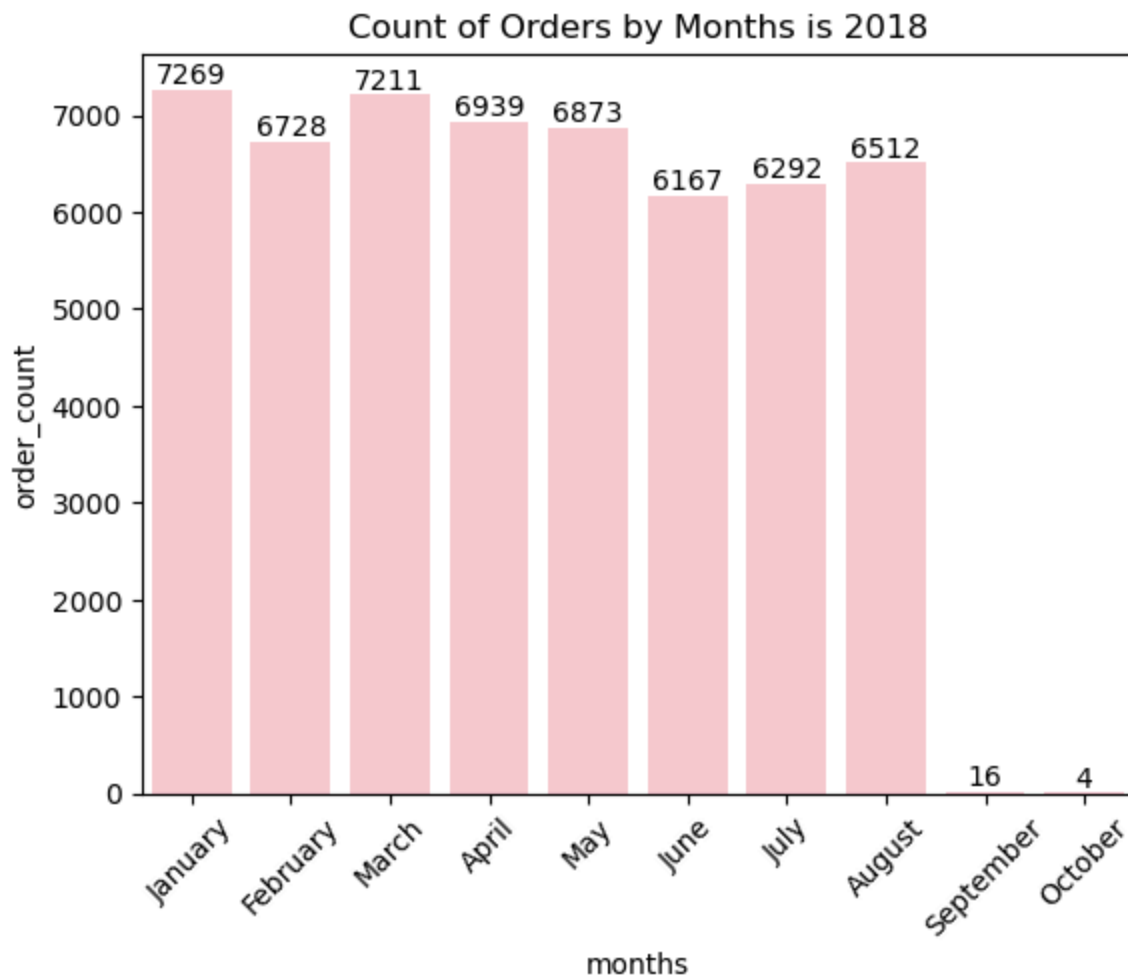


## The number of orders per month in 2018.

```
In [39]:  query = """ select monthname(order_purchase_timestamp) months, count(order_i
          from orders where year(order_purchase_timestamp) = 2018
          group by months
          """

          cur.execute(query)

          data = cur.fetchall()
          df = pd.DataFrame(data, columns = ["months", "order_count"])
          o = ["January", "February","March","April","May","June","July","August","Sep

          ax = sns.barplot(x = df["months"],y =  df["order_count"], data = df, order =
          plt.xticks(rotation = 45)
          ax.bar_label(ax.containers[0])
          plt.title("Count of Orders by Months is 2018")

          plt.show()
```



## Total revenue generated by each seller, and rank them by revenue.

```
In [47]:  query = """ select *, dense_rank() over(order by revenue desc) as rn from
          (select order_items.seller_id, sum(payments.payment_value)
```
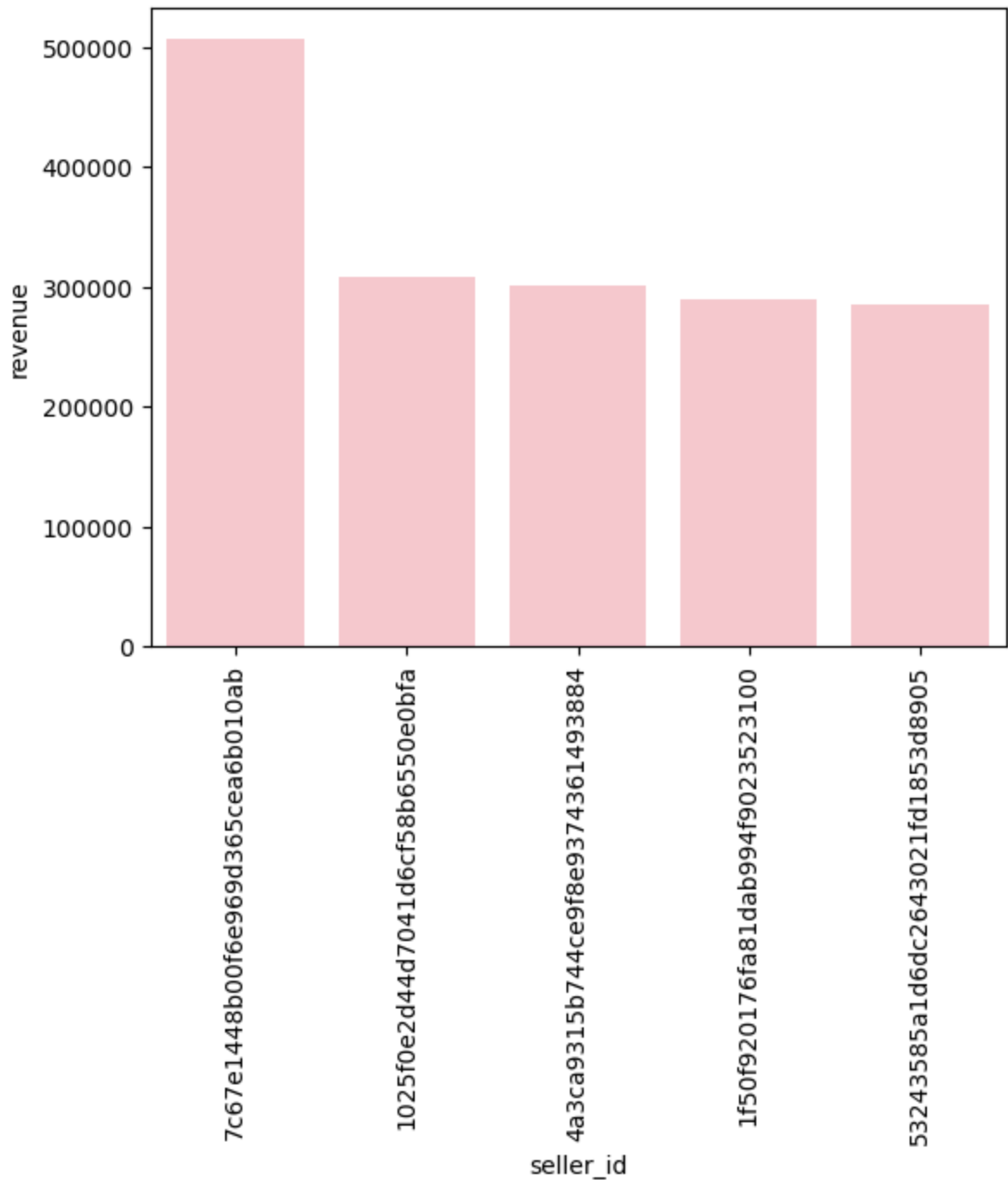
```
revenue from order_items join payments
on order_items.order_id = payments.order_id
group by order_items.seller_id) as a """

cur.execute(query)
data = cur.fetchall()
df = pd.DataFrame(data, columns = ["seller_id", "revenue", "rank"])
df = df.head()
sns.barplot(x = "seller_id", y = "revenue", data = df , color = 'pink')
plt.xticks(rotation = 90)
plt.show()
```
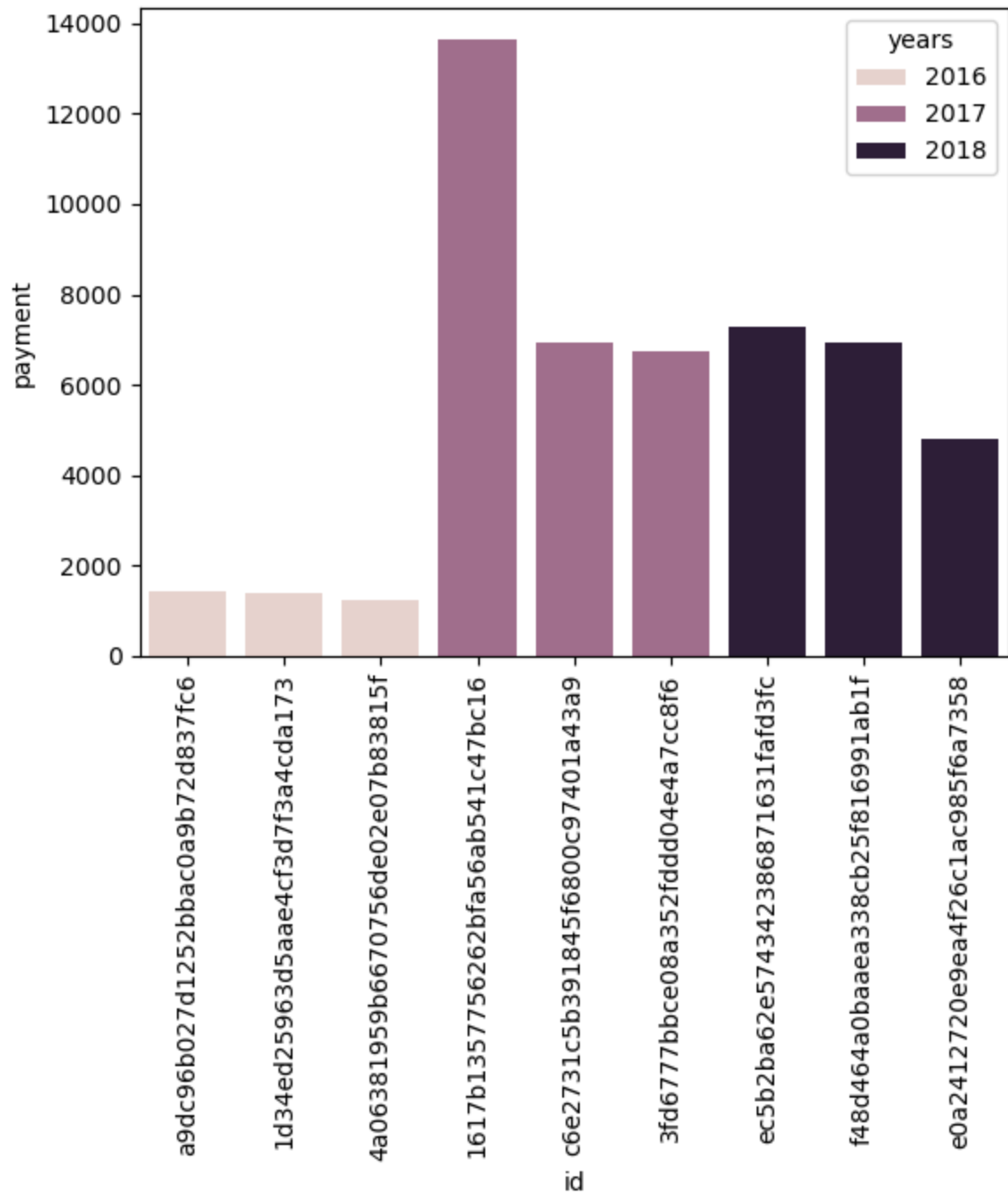
# Top 3 customers who spent the most money in each year.

In [27]:
```python
query = """select years, customer_id, payment, d_rank
from
(select year(orders.order_purchase_timestamp) years,
orders.customer_id,
sum(payments.payment_value) payment,
dense_rank() over(partition by year(orders.order_purchase_timestamp)
order by sum(payments.payment_value) desc) d_rank
from orders join payments
on payments.order_id = orders.order_id
group by year(orders.order_purchase_timestamp),
orders.customer_id) as a
where d_rank <= 3 ;"""

cur.execute(query)
data = cur.fetchall()
df = pd.DataFrame(data, columns = ["years","id","payment","rank"])
sns.barplot(x = "id", y = "payment", data = df, hue = "years")
plt.xticks(rotation = 90)
plt.show()
```

In [ ]: