

```
In [30]: 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')
]

# Connect to the MySQL database
db = mysql.connector.connect(
    host="localhost",
    user="root",
    password="root@123456",
    database="ecommerce"
)
cursor = db.cursor()

# Folder containing the CSV files
folder_path = "C:/Users/Pratibha/Desktop/Ecommerce"

def get_col_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.replace(np.nan, None, inplace=True)

    # Determine column types for MySQL
    columns = {}
    for col in df.columns:
        dtype = df[col].dtype
        columns[col] = get_col_type(dtype)

    # Generate the CREATE TABLE statement with appropriate data types
    create_table_query = f"CREATE TABLE {table_name} ({columns})"
    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(row if pd.isna(x) else x for x in row)
        sql = f"INSERT INTO {table_name} ({', '.join('' + col + '' for col in df.columns)}) VALUES ({', '.join('?' * len(row))})"
        cursor.execute(sql, values)

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

# Close the connection
conn.close()

# Processing customers.csv
df = pd.read_csv('customers.csv')
df.replace(np.nan, None, inplace=True)
df.columns = [col.replace(' ', '_') for col in df.columns]
df.to_sql('customers', db, if_exists='replace', index=False)

# Processing orders.csv
df = pd.read_csv('orders.csv')
df.replace(np.nan, None, inplace=True)
df.columns = [col.replace(' ', '_') for col in df.columns]
df.to_sql('orders', db, if_exists='replace', index=False)

# Processing products.csv
df = pd.read_csv('products.csv')
df.replace(np.nan, None, inplace=True)
df.columns = [col.replace(' ', '_') for col in df.columns]
df.to_sql('products', db, if_exists='replace', index=False)

# Processing geolocation.csv
df = pd.read_csv('geolocation.csv')
df.replace(np.nan, None, inplace=True)
df.columns = [col.replace(' ', '_') for col in df.columns]
df.to_sql('geolocation', db, if_exists='replace', index=False)

# Processing payments.csv
df = pd.read_csv('payments.csv')
df.replace(np.nan, None, inplace=True)
df.columns = [col.replace(' ', '_') for col in df.columns]
df.to_sql('payments', db, if_exists='replace', index=False)

# Processing order_items.csv
df = pd.read_csv('order_items.csv')
df.replace(np.nan, None, inplace=True)
df.columns = [col.replace(' ', '_') for col in df.columns]
df.to_sql('order_items', db, if_exists='replace', index=False)
```

In [31]: pip install mysql-connector-python

Requirement already satisfied: mysql-connector-python in c:\users\pratiha\appdata\local\programs\python\python32\lib\site-packages (9.0.0)

Note: you may need to restart the kernel to use updated packages.

[Notice] A new release of pip is available: 24.0 -> 24.2

[Notice] To update, run: python.exe -m pip install --upgrade pip

```
In [61]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
import mysql.connector

db = mysql.connector.connect(host="localhost",
                             username="root",
                             password="root@123456",
                             database="ecommerce")
cur = db.cursor()
```

List all unique cities where customers are located.

```
In [56]: query = """select distinct(customer_city) from customers """
cur.execute(query)
data = cur.fetchall()

df = pd.DataFrame(data, columns = ["customer_city"])
df.head()
```

```
Out[56]:
```

	customer_city
0	franca
1	sao bernardo do campo
2	sao paulo
3	mogi das cruzes
4	campinas

Count the number of orders placed in 2017

```
In [39]: query = """select count(order_id) from orders where year(order_purchase_timestamp) = 2017 """
cur.execute(query)
data = cur.fetchall()

print("Total orders placed in 2017 are:", data[0][0])

Out[39]:
```

(Total orders placed in 2017 are, 45181)

Find the total sales per category.

```
In [34]: 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
```

```
Out[34]:
```

	Category	Sales
0	PERFUMERY	106738.66
1	FURNITURE DECORATION	1430176.39
2	TELEPHONY	486682.05
3	BED TABLE BATH	1712503.67
4	AUTOMOTIVE	852294.33
...
69	CDS MUSIC DVDS	1199.43
70	LA CUISINE	2913.53
71	FASHION CHILDRENS CLOTHING	785.67
72	PC GAMER	2174.43
73	INSURANCE AND SERVICES	324.51
74 rows x 2 columns		

Calculate the percentage of orders that were paid in installments.

```
In [37]: query = """select count(case when payment_installments >= 3 then 1
else 0 end)/count(*)*100 from payments
"""
cur.execute(query)
data = cur.fetchall()

data

Out[37]:
```

(Decimal('99.9991'),)

Count the number of customers from each state.

```
In [44]: 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=(15,6))
plt.bar(df['state'], df['customer_count'])
plt.xticks(rotation = 90)
plt.title('Count of customers by state')
plt.show()
```

Calculate the number of orders per month in 2018.

```
In [40]: query = """select monthname(order_purchase_timestamp) months, count(order_id) order_count
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"])
df = df.sort_values(by = "order_count", ascending = False)
plt.figure(figsize=(15,6))
plt.bar(df['months'], df['order_count'])
plt.xticks(rotation = 45)
plt.title('Count of orders by month in 2018')
plt.show()
```

Find the average number of products per order, grouped by customer city.

```
In [57]: query = """with count_per_order as
(select orders.order_id, orders.customer_id, count(order_items.order_id) as oc
from orders join order_items
on orders.order_id = order_items.order_id
group by orders.order_id, orders.customer_id)
select customers.customer_city, round(avg(count_per_order.oc),2) average_orders
from customers join count_per_order
on customers.customer_id = count_per_order.customer_id
group by customers.customer_city order by average_orders desc
"""
cur.execute(query)
data = cur.fetchall()

df = pd.DataFrame(data, columns = ["customer_city", "average_products/order"])
df.head()
```

```
Out[57]:
```

	customer_city	average_products/order
0	padre carvalho	7.00
1	celso ramos	6.50
2	dadas	6.00
3	candido godoi	6.00
4	matias dimpio	5.00
5	cideblanda	4.00
6	picoara	4.00
7	monte de sao paulo	4.00
8	lexeira Soares	4.00
9	cunheiro	4.00

Calculate the percentage of total revenue contributed by each product category.

```
In [58]: query = """select upper(products.product_category) category,
round(sum(payments.payment_value)/sum(payment_value) from payments)*100,2) sales_percentage
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 order by sales_percentage desc
"""
cur.execute(query)
data = cur.fetchall()

df = pd.DataFrame(data, columns = ["Category", "percentage distribution"])
df.head()
```

```
Out[58]:
```

	Category	percentage distribution
0	BED TABLE BATH	10.70
1	HEALTH BEAUTY	10.36
2	COMPUTER ACCESSORIES	9.90
3	FURNITURE DECORATION	8.93
4	WATCHES PRESENT	8.93

Identify the correlation between product price and the number of times a product has been purchased.

```
In [60]: cur = db.cursor()
query = """select products.product_category,
count(order_items.product_id)
from products join order_items
on products.product_id = order_items.product_id
group by products.product_category
"""
cur.execute(query)
data = cur.fetchall()
df = pd.DataFrame(data, columns = ["Category", "order_count", "price"])
arr1 = df["order_count"]
arr2 = df["price"]
import numpy as np
a = np.corrcoef(arr1, arr2)
print("The correlation is:", a[0][1])

The correlation is -0.186325146757562
```

Calculate the total revenue generated by each seller, and rank them by revenue.

```
In [62]: 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)
plt.xticks(rotation = 90)
plt.show()
```

Calculate the moving average of order values for each customer over their order history.

```
In [68]: query = """select customer_id, order_purchase_timestamp, payment,
avg(payment) over(partition by customer_id order by order_purchase_timestamp
rows between 3 preceding and current row) as mov_avg
from
(select orders.customer_id, orders.order_purchase_timestamp,
payments.payment_value as payment
from payments join orders
on payments.order_id = orders.order_id) as a
"""
cur.execute(query)
data = cur.fetchall()
df = pd.DataFrame(data, columns = ["customer_id", "order_purchase_timestamp", "payment", "payment_value"])
df
```

```
Out[68]:
```

	customer_id	order_purchase_timestamp	payment	payment_value
0	00012a2e0ff8dc20a59a06491703	2017-11-14 16:08:26	114.74	114.739998
1	00016105f8020a09016078e4c21703	2017-07-16 09:40:32	67.41	67.410004
2	001106190a0a0308a0a340e079	2017-02-28 11:08:43	195.42	195.419998
3	0002414953443074040a0a7a28f4b	2017-08-16 13:09:20	179.38	179.380006
4	00037a0e02552490c315e7028f4b	2018-04-02 13:42:17	157.01	157.010002
...
103881	ff6c9f7968b764843a2951b11341	2018-03-29 16:59:26	71.23	27.120001
103882	ff6a2b1a0a0a0308a0a340e079	2018-05-22 13:36:02	63.13	63.130001
103883	ff6a2b1a0a0a0308a0a340e079	2018-06-13 16:57:05	214.13	214.130005
103884	ff6a2b1a0a0a0308a0a340e079	2017-09-02 11:53:32	45.50	45.500000
103885	ff6a2b1a0a0a0308a0a340e079	2017-09-29 14:07:03	18.37	18.370001
103886 rows x 4 columns				

Calculate the cumulative sales per month for each year.

```
In [70]: query = """select years, months, payment, sum(payment)
over(order by years, months) cumulative_sales
from
(select year(orders.order_purchase_timestamp) as years,
month(orders.order_purchase_timestamp) as months,
round(sum(payments.payment_value),2) as payment
from orders join payments
on orders.order_id = payments.order_id
group by years, months) as a
"""
cur.execute(query)
data = cur.fetchall()
df = pd.DataFrame(data, columns = ["years", "months", "payment", "payment_value"])
df.head()
```

```
Out[70]:
```

	years	months	payment	payment_value
0	2016	9	252.24	252.24
1	2016	10	59090.48	59342.72
2	2016	12	19.62	59362.34
3	2017	1	138468.04	157850.38
4	2017	2	291908.01	449758.39

Calculate the year-over-year growth rate of total sales.

```
In [71]: query = """with a as(select year(orders.order_purchase_timestamp) as years,
sum(payments.payment_value) as payment
from orders join payments
on orders.order_id = payments.order_id
group by year 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[71]:
```

	years	yoy % growth
0	2016	NaN
1	2017	1212.703761
2	2018	20.00904

Calculate the retention rate of customers, defined as the percentage of customers who make another purchase within 6 months of their first purchase.

```
In [72]: query = """with a as (select customers.customer_id,
min(orders.order_purchase_timestamp) first_order
from customers join orders
on customers.customer_id = orders.customer_id
group by customers.customer_id),
b as (select a.customer_id, count(distinct orders.order_purchase_timestamp) next_order
from a join orders
on orders.customer_id = a.customer_id
and orders.order_purchase_timestamp > first_order
and orders.order_purchase_timestamp <
date_add(first_order, interval 6 month)
group by a.customer_id)
select 100 * (count( distinct a.customer_id)/ count(distinct b.customer_id))
from a left join b
on a.customer_id = b.customer_id
"""
cur.execute(query)
data = cur.fetchall()

data

Out[72]:
```

(None,)

Identify the top 3 customers who spent the most money in each year.

```
In [73]: query = """select years, customer_id, payment, d.rank
from
(select year(orders.order_purchase_timestamp) years,
orders.customer_id,
sum(payments.payment_value) as 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 orders.order_id = payments.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()
```


