seller id seller_zip_code_prefix seller_city seller_state 0 dtype: int64 Processing products.csv NaN values before replacement: product_id product category 610 product_name_length 610 product_description_length 610 610 product_photos_qty 2 product_weight_g 2 product length cm product_height_cm product width cm dtype: int64 Processing geolocation.csv NaN values before replacement: geolocation_zip_code_prefix geolocation_lat geolocation_lng geolocation_city geolocation_state dtype: int64 Processing payments.csv NaN values before replacement: order id payment_sequential payment_type payment installments payment value dtype: int64 Processing order_items.csv NaN values before replacement: order_id order_item_id product id seller_id shipping limit date price freight_value dtype: int64 In [3]: **import** pandas **as** pd import numpy as np import matplotlib.pyplot as plt import mysql.connector db = mysql.connector.connect(host = "localhost", username = "root", password = "abhay55555", database = "ecommerce") cur = db.cursor() List all unique cities where customers are located. In [4]: query = """ select distinct customer_city from customers""" cur.execute(query) data = cur.fetchall() df = pd.DataFrame(data) df.head() Out[4]: 0 franca 1 sao bernardo do campo sao paulo mogi das cruzes campinas Count the number of orders placed in 2017. In [5]: query = """ select count(order_id) from orders where year(order_purchase_timestamp) = 2017 """ cur.execute(query) data = cur.fetchall() "total orders placed in 2017 are", data[0][0] ('total orders placed in 2017 are', 135303) Find the total sales per category. In [6]: 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']) Out[6]: Sales Category PERFUMERY 4053909.28 FURNITURE DECORATION 11441411.13 2 TELEPHONY 3895056.41 **3** FASHION BAGS AND ACCESSORIES 1745266.24 BED TABLE BATH 13700429.37 69 CDS MUSIC DVDS 9595.44 LA CUISINE 23308.24 70 FASHION CHILDREN'S CLOTHING 6285.36 72 PC GAMER 17395.44 73 INSURANCE AND SERVICES 2596.08 74 rows × 2 columns Calculate the percentage of orders that were paid in installments. In [7]: 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')) Count the number of customers from each state. In [8]: 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"], color="g") plt.xticks(rotation = 90) plt.xlabel("states") plt.ylabel("customer count") plt.title("Count of Customers by States") plt.show() Count of Customers by States 120000 100000 80000 60000 40000 20000 Calculate the number of orders per month in 2018. In [9]: import seaborn as sns query = """ select monthname(order_purchase_timestamp) months, count(order_id) order_count from orders where year(order_purchase_timestamp) = 2018 cur.execute(query) data = cur.fetchall() df = pd.DataFrame(data, columns = ["months", "order_count"]) o = ["January", "February", "March", "April", "May", "June", "July", "August", "September", "October"] $ax = sns.barplot(x = df["months"], y = df["order_count"], data = df, order = o, color = "#F02BC0")$ plt.xticks(rotation = 45) ax.bar_label(ax.containers[0]) plt.title("Count of Orders by Months is 2018") plt.show() Count of Orders by Months is 2018 20000 15000 10000 5000

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

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

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

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

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

round((sum(payments.payment_value)/(select sum(payment_value) from payments))*100,2) sales_percentage

(select orders.order_id, orders.customer_id, count(order_items.order_id) as oc

select customers.customer_city, round(avg(count_per_order.oc),2) average_orders

df = pd.DataFrame(data,columns = ["customer city", "average products/order"])

In [10]: query = """with count_per_order as

cur.execute(query)

Out[10]:

data = cur.fetchall()

padre carvalho

candido godoi

matias olimpio

8 morro de sao paulo

join payments

df.head()

In [12]: cur = db.cursor()

Out[11]:

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

cidelandia

curralinho

teixeira soares

picarra

from products join order_items

BED TABLE BATH

HEALTH BEAUTY

WATCHES PRESENT

count(order_items.product_id),
round(avg(order_items.price),2)
from products join order_items

group by products.product_category""'

print("the correlation is", a[0][-1])

the correlation is -0.10631514167157562

revenue from order_items join payments
on order_items.order_id = payments.order_id
group by order_items.seller_id) as a """

2 COMPUTER ACCESSORIES

3 FURNITURE DECORATION

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

arr2 = df["price"]

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

plt.show()

2.00 -1.75 -1.50 -1.25 -1.00 -0.75 -0.50 -0.25 -

plt.xticks(rotation = 90)

arr1 = df["order_count"]

a = np.corrcoef([arr1,arr2])

celso ramos

from orders join order_items

on orders.order_id = order_items.order_id
group by orders.order_id, orders.customer_id)

customer city average products/order

on customers.customer_id = count_per_order.customer_id

group by customers.customer_city order by average_orders desc

39.00

36.00

36.00

30.00

24.00

24.00

24.00

24.00

24.00

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

42.79

41.41

39.61

35.73

35.71

df = pd.DataFrame(data,columns = ["Category", "order_count","price"])

In [14]: query = """ select *, dense_rank() over(order by revenue desc) as rn from

df = pd.DataFrame(data, columns = ["seller_id", "revenue", "rank"])

(select order_items.seller_id, sum(payments.payment_value)

sns.barplot(x = "seller_id", y = "revenue", data = df)

In [15]: query = """select customer_id, order_purchase_timestamp, payment,

(select orders.customer_id, orders.order_purchase_timestamp,

0 00012a2ce6f8dcda20d059ce98491703 2017-11-14 16:08:26 114.74 114.739998
 1 00012a2ce6f8dcda20d059ce98491703 2017-11-14 16:08:26 114.74 114.739998

2 00012a2ce6f8dcda20d059ce98491703 2017-11-14 16:08:26 114.74 114.739998

3 00012a2ce6f8dcda20d059ce98491703 2017-11-14 16:08:26 114.74 114.739998

4 00012a2ce6f8dcda20d059ce98491703 2017-11-14 16:08:26 114.74 114.739998

623311 ffffe8b65bbe3087b653a978c870db99 2017-09-29 14:07:03 18.37 18.370001

623312 ffffe8b65bbe3087b653a978c870db99 2017-09-29 14:07:03 18.37 18.370001

623313 ffffe8b65bbe3087b653a978c870db99 2017-09-29 14:07:03 18.37 18.370001

623314 ffffe8b65bbe3087b653a978c870db99 2017-09-29 14:07:03 18.37 18.370001 **623315** ffffe8b65bbe3087b653a978c870db99 2017-09-29 14:07:03 18.37 18.370001

round(sum(payments.payment_value),2) as payment from orders join payments

In [17]: query = """select years, months , payment, sum(payment)
 over(order by years, months) cumulative_sales from
 (select year(orders.order_purchase_timestamp) as years,

on orders.order_id = payments.order_id

2

117.72 356174.04

1513.44

1 2016 10 354542.88 356056.32

3 2017 1 830928.24 1187102.28

4 2017 2 1751448.06 2938550.34

5 2017 3 2699181.60 5637731.94

6 2017 4 2506728.18 8144460.12

7 2017 5 3557512.92 11701973.04

8 2017 6 3067658.28 14769631.32

9 2017 7 3554297.52 18323928.84

10 2017 8 4046377.92 22370306.76

11 2017 9 4366574.70 26736881.46

12 2017 10 4678067.28 31414948.74

13 2017 11 7169296.80 38584245.54

14 2017 12 5270408.88 43854654.42

15 2018 1 6690025.07 50544679.49 **16** 2018 2 5954780.04 56499459.53

17 2018 3 6957912.72 63457372.25

18 2018 4 6964712.88 70422085.13

19 2018 5 6923892.91 77345978.04 **20** 2018 6 6143283.00 83489261.04

21 2018 7 6399244.49 89888505.53

22 2018 8 6134551.93 96023057.46

23 2018 9 26637.24 96049694.70

24 2018 10 3538.02 96053232.72

on orders.order_id = payments.order_id

group by years order by years)

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

years yoy % growth

1 2017 12112.703759

20.000924

In [19]: query = """with a as (select customers.customer_id,

min(orders.order_purchase_timestamp) first_order

and orders.order_purchase_timestamp > first_order

query = """select years, customer_id, payment, d_rank

dense_rank() over(partition by year(orders.order_purchase_timestamp)

df = pd.DataFrame(data, columns = ["years","id","payment","rank"])
sns.barplot(x = "id", y = "payment", data = df, hue = "years")

2016

2018

(select year(orders.order_purchase_timestamp) years,

order by sum(payments.payment_value) desc) d_rank

group by year(orders.order purchase timestamp),

on customers.customer_id = orders.customer_id

first purchase.

from customers join orders

from a join orders

from a left join b

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

orders.customer_id,

from orders join payments

orders.customer_id) as a
where d_rank <= 3 ;"""</pre>

plt.xticks(rotation = 90)

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

plt.show()

80000

70000

data

Out[19]: [(None,)]

group by a.customer_id)

group by customers.customer_id),

on orders.customer_id = a.customer_id

and orders.order_purchase_timestamp <
date_add(first_order, interval 6 month)</pre>

on a.customer_id = b.customer_id ;"""

sum(payments.payment_value) payment,

on payments.order id = orders.order id

0 2016

2 2018

Out[18]:

month(orders.order_purchase_timestamp) as months,

group by years, months order by years, months) as a

3

1513.44

Calculate the cumulative sales per month for each year.

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

b as (select a.customer_id, count(distinct orders.order_purchase_timestamp) next_order

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

select 100 * (count(distinct a.customer_id)/ count(distinct b.customer_id))

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

In [18]: query = """with a as(select year(orders.order_purchase_timestamp) as years,

select years, ((payment - lag(payment, 1) over(order by years))/

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

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

round(sum(payments.payment_value),2) as payment from orders join payments

rows between 2 preceding and current row) as mov avg

on payments.order_id = orders.order_id) as a"""

payments.payment_value as payment

from payments join orders

cur.execute(query)
data = cur.fetchall()
df = pd.DataFrame(data)

623316 rows × 4 columns

cur.execute(query)
data = cur.fetchall()
df = pd.DataFrame(data)

0 1

0 2016 9

2 2016 12

df

Out[17]:

df

avg(payment) over(partition by customer_id order by order_purchase_timestamp

In [11]: query = """select upper(products.product_category) category,

group by category order by sales_percentage desc"""

Category percentage distribution

query = """select products.product_category,

on products.product_id = order_items.product_id

on products.product_id = order_items.product_id

on payments.order_id = order_items.order_id

from customers join count_per_order

In [1]: import pandas as pd

csv_files = [

import os

import mysql.connector

List of CSV files and their corresponding table names

('customers.csv', 'customers'),

('payments.csv', 'payments'),
('order_items.csv','order_items')

('geolocation.csv', 'geolocation'),

Added payments.csv for specific handling

if pd.api.types.is_integer_dtype(dtype):

elif pd.api.types.is_float_dtype(dtype):

elif pd.api.types.is_bool_dtype(dtype):

elif pd.api.types.is_datetime64_any_dtype(dtype):

file_path = os.path.join(folder_path, csv_file)

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

Generate the CREATE TABLE statement with appropriate data types

Convert row to tuple and handle NaN/None explicitly
values = tuple(None if pd.isna(x) else x for x in row)

df.columns = [col.replace(' ', '_').replace('-', '_').replace('.', '_') for col in df.columns]

 $sql = f"INSERT INTO `{table_name}` ({', '.join(['`' + col + '`' for col in df.columns])}) VALUES ({', '.join(['%s'] * len(row))})"$

columns = ', '.join([f'`{col}` {get_sql_type(df[col].dtype)}' for col in df.columns])

create_table_query = f'CREATE TABLE IF NOT EXISTS `{table_name}` ({columns})'

Read the CSV file into a pandas DataFrame

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}")

cursor.execute(create_table_query)

cursor.execute(sql, values)

for , row in df.iterrows():

Insert DataFrame data into the MySQL table

Commit the transaction for the current CSV file

0

0

160

1783

2965

0

folder path = "C:\\Users\\abhay\\OneDrive\\Desktop\\ecommerce"

('orders.csv', 'orders'),
('sellers.csv', 'sellers'),
('products.csv', 'products'),

Connect to the MySQL database
conn = mysql.connector.connect(

password='abhay55555',
database='ecommerce'

Folder containing the CSV files

host='localhost',
user='root',

cursor = conn.cursor()

def get_sql_type(dtype):

else:

return 'INT'

return 'FLOAT'

return 'BOOLEAN'

return 'DATETIME'

for csv_file, table_name in csv_files:

df = pd.read_csv(file_path)

return 'TEXT'

Clean column names

conn.commit()

Close the connection

Processing customers.csv

customer zip code prefix

Processing orders.csv

order_purchase_timestamp

Processing sellers.csv

order_delivered_carrier_date

order_delivered_customer_date

order_estimated_delivery_date

NaN values before replacement:

order_approved_at

dtype: int64

NaN values before replacement:

NaN values before replacement:

conn.close()

customer_id

customer_city
customer state

dtype: int64

order_id
customer_id
order_status

customer_unique_id