

Data Engineer Intern Practical Exam

Name:P.H.H.Chamika Deemantha

Email:chamikade01@gmail.com

Table of Contents

1	Part 1: Video Demonstration & Connect to GitHub	2
2	Part 2: Data Preparation	3
2.1	Data Import:	3
2.2	Database Connection:	4
3	Part 3: Flask App Setup	7
3.1	Sidebar Filters:	7
3.2	Main Dashboard:	8
4	Part 4: Data Analysis	9
4.1	Machine Learning Model	9

1 Part 1: Video Demonstration & Connection of GitHub

Below is a video demonstration link of the completed task.

url: <https://www.youtube.com/watch?v=VSg1cmsiJ5Y>

Below is the GitHub link for codebase.

url: [https://github.com/chamikadeemantha/Delivergate Practicle Test](https://github.com/chamikadeemantha/Delivergate_Practicle_Test)

2 Part 2: Data Preparation

2.1 Data Import:

```
C:\Users\User\Desktop\Delivergate Pvt Ltd > xampp.py > ...
1  import pandas as pd      Import "pandas" could not be resolved from source
2  from sqlalchemy import create_engine, Table, Column, Integer, String, Float, MetaData, DateTime      Import "sqlalchemy" could not be resolved
3  import pymysql          Import "pymysql" could not be resolved from source
4
5  # MySQL Database Connection Parameters
6  MYSQL_USER = 'root'
7  MYSQL_PASSWORD = ''      # Empty password as per XAMPP config
8  MYSQL_HOST = 'localhost'
9  MYSQL_PORT = 3306
10 MYSQL_DB = 'delivergatedb'
11
12 # Create a connection string
13 db_connection_str = f'mysql+pymysql://{MYSQL_USER}:{MYSQL_PASSWORD}@{MYSQL_HOST}:{MYSQL_PORT}/{MYSQL_DB}'
14 engine = create_engine(db_connection_str)
15
16 # File paths for the CSVs
17 customers_file_path = r'C:\Users\User\Desktop\Delivergate Pvt Ltd\customers.csv'
18 orders_file_path = r'C:\Users\User\Desktop\Delivergate Pvt Ltd\order.csv'
19
20 # Load CSV files into DataFrames
21 customers_df = pd.read_csv(customers_file_path, usecols=['customer_id', 'name']) # Only keep relevant columns
22 orders_df = pd.read_csv(orders_file_path, usecols=['id', 'customer_id', 'total_amount', 'created_at']) # Only keep relevant columns
23
24 # Rename columns to match the database schema
25 customers_df = customers_df.rename(columns={'name': 'customer_name'})
26 orders_df = orders_df.rename(columns={'id': 'order_id', 'created_at': 'order_date'})
27
28 # Prepare the metadata for creating tables
29 metadata = MetaData()
30
31 # Define the Customers table
32 customers_table = Table(
33     'customers', metadata,
34     Column('customer_id', Integer, primary_key=True),
35     Column('customer_name', String(255))
36 )
37
38 # Define the Orders table
39 orders_table = Table(
40     'orders', metadata,
41     Column('order_id', Integer, primary_key=True),
42     Column('customer_id', Integer),
43     Column('total_amount', Float),
44     Column('order_date', DateTime)
45 )
46
47 # Create tables in the MySQL database
48 metadata.create_all(engine)
49
50 # Insert data into customers table
51 customers_df.to_sql('customers', con=engine, if_exists='append', index=False)
52
53 # Insert data into orders table
54 orders_df.to_sql('orders', con=engine, if_exists='append', index=False)
55
56 print("Data has been imported into the MySQL database successfully!")
57
```

2.2 Database Connection:

```
C:\Users\User\Desktop> Flask App > app.py > ...
1  import matplotlib      Import "matplotlib" could not be resolved from source
2  matplotlib.use('Agg') # Use non-interactive backend
3  from flask import Flask, render_template, request, redirect, url_for      Import "flask" could not be resolved
4  from sqlalchemy import create_engine      Import "sqlalchemy" could not be resolved
5  import pandas as pd      Import "pandas" could not be resolved from source
6  import matplotlib.pyplot as plt      Import "matplotlib.pyplot" could not be resolved from source
7  from io import BytesIO
8  import base64
9
10 app = Flask(__name__)
11
12 # Database connection
13 MYSQL_USER = 'root'
14 MYSQL_PASSWORD = '' # Empty password as per XAMPP config
15 MYSQL_HOST = 'localhost'
16 MYSQL_PORT = 3306
17 MYSQL_DB = 'delivertdb'
18
19 db_connection_str = f'mysql+pymysql://{MYSQL_USER}:{MYSQL_PASSWORD}@{MYSQL_HOST}:{MYSQL_PORT}/{MYSQL_DB}'
20 engine = create_engine(db_connection_str)
21
22 @app.route("/", methods=["GET", "POST"])
23 def index():
24     toast_message = None
25     try:
26         start_date = request.form.get("start_date")
27         end_date = request.form.get("end_date")
28         min_total_amount = request.form.get("min_total_amount", type=float, default=0)
29         min_order_count = request.form.get("min_order_count", type=int, default=1)
30     except ValueError:
31         toast_message = "Please enter valid values for all filters."
32         return render_template("index.html", toast_message=toast_message)
33
34     # Query to get customers with enough orders
35     customer_query = """
36         SELECT o.customer_id, c.customer_name, COUNT(o.order_id) AS order_count
37         FROM orders o
38         JOIN customers c ON o.customer_id = c.customer_id
39         WHERE o.total_amount >= %s
40     """
41
42     filters = [min_total_amount]
43     if start_date and end_date:
44         customer_query += " AND o.order_date BETWEEN %s AND %s"
45         filters.extend([start_date, end_date])
46     customer_query += " GROUP BY o.customer_id HAVING COUNT(o.order_id) >= %s"
47     filters.append(min_order_count)
48
49     connection = engine.raw_connection()
50     try:
51         customers_df = pd.read_sql(customer_query, connection, params=filters)
52         if customers_df.empty:
53             toast_message = "No customers found with the specified order count and total amount filters."
54             return render_template("index.html", toast_message=toast_message)
55         customer_ids = customers_df['customer_id'].tolist()
56     finally:
57         connection.close()
```

```

58 # Query to get orders of the filtered customers
59 if customer_ids:
60     orders_query = """
61         SELECT o.*, c.customer_name
62         FROM orders o
63         JOIN customers c ON o.customer_id = c.customer_id
64         WHERE o.customer_id IN (%s)
65     """ % ','.join(map(str, customer_ids))
66
67     connection = engine.raw_connection()
68     try:
69         orders_df = pd.read_sql(orders_query, connection)
70     finally:
71         connection.close()
72 else:
73     orders_df = pd.DataFrame()
74
75 # Ensure total_amount is numeric, setting non-numeric values to NaN and then dropping them
76 orders_df['total_amount'] = pd.to_numeric(orders_df['total_amount'], errors='coerce')
77 orders_df = orders_df.dropna(subset=['total_amount'])
78
79 # Convert order_date to datetime and handle errors
80 try:
81     orders_df['order_date'] = pd.to_datetime(orders_df['order_date'], errors='coerce')
82     orders_df = orders_df.dropna(subset=['order_date'])
83     orders_df = orders_df.set_index('order_date')
84 except Exception as e:
85     toast_message = "Error processing dates. Please check the date format."
86     return render_template("index.html", toast_message=toast_message)
87
88 # Calculate Total Spent by each customer
89 orders_df['Total Spent'] = orders_df.groupby('customer_id')['total_amount'].transform('sum')
90
91 # Apply search filter for customer name if provided
92 customer_name_filter = request.form.get("customer_name")
93 if customer_name_filter:
94     orders_df = orders_df[orders_df['customer_name'].str.contains(customer_name_filter, case=False, na=False)]
95
96 # Key metrics
97 total_revenue = orders_df['total_amount'].sum()
98 unique_customers = orders_df['customer_id'].nunique()
99 total_orders = len(orders_df)
100

```

```

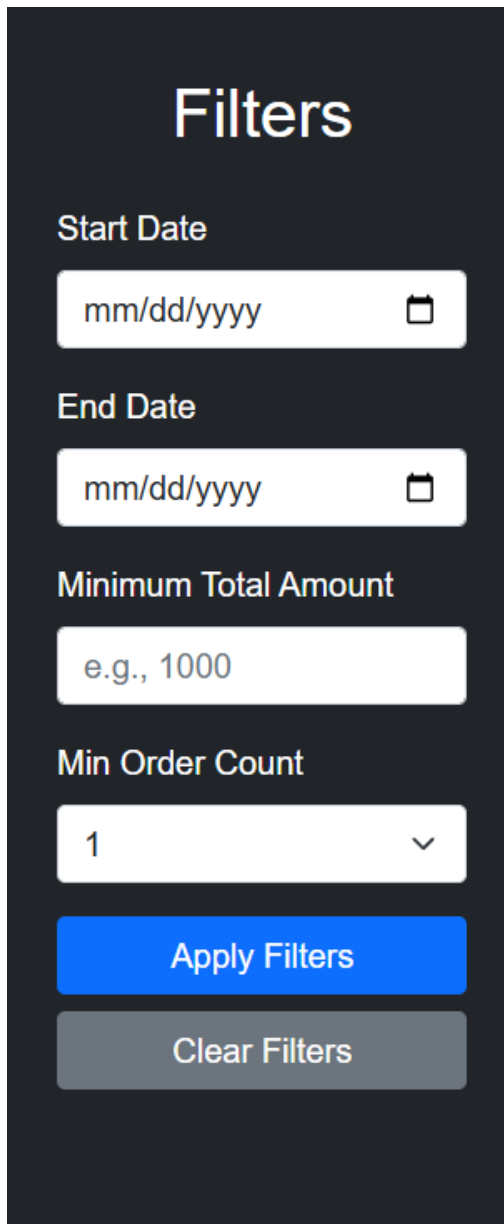
101 # Visualizations
102 top_customers_img = None
103 revenue_over_time_img = None
104
105 def plot_to_img(plot_func):
106     fig, ax = plt.subplots(figsize=(8, 6))
107     plot_func(ax)
108     plt.xticks(rotation=45, ha='right')
109     plt.tight_layout()
110     img = BytesIO()
111     fig.savefig(img, format='png')
112     img.seek(0)
113     plt.close(fig)
114     return base64.b64encode(img.getvalue()).decode()
115
116 if not orders_df.empty:
117     # Group by customer and get top 10 customers by revenue
118     top_customers = orders_df.groupby(['customer_name', 'customer_id']).total_amount.sum().nlargest(10)
119
120     # Generate top customers plot if data exists
121     if not top_customers.empty:
122         top_customers_img = plot_to_img(lambda ax: top_customers.plot(kind='bar', title="Top 10 Customers by Revenue", ax=ax, xlabel="Customer Name"))
123
124     # Generate revenue over time plot if there is data
125     revenue_over_time = orders_df.resample("M").total_amount.sum()
126     if not revenue_over_time.empty:
127         revenue_over_time_img = plot_to_img(lambda ax: revenue_over_time.plot(kind='line', title="Revenue Over Time", ax=ax, xlabel="Order Date"))
128
129 # Rename columns for a more readable table display
130 orders_df = orders_df.rename(columns={
131     'order_id': 'Order ID',
132     'customer_id': 'Customer ID',
133     'total_amount': 'Order Amount',
134     'customer_name': 'Customer Name',
135     'order_count': 'Order Count'
136 })
137
138 # Show only the first 10 rows of the filtered data
139 orders_table_html = orders_df.head(10).to_html(classes='table table-striped', index=False)
140
141 return render_template("index.html", total_revenue=total_revenue,
142     unique_customers=unique_customers,
143     total_orders=total_orders,
144     top_customers_img=top_customers_img,
145     revenue_over_time_img=revenue_over_time_img,
146     orders_table_html=orders_table_html,
147     toast_message=toast_message)
148
149 @app.route("/reset_filters", methods=["GET"])
150 def reset_filters():
151     # Redirect to the main index page without any filters
152     return redirect(url_for('index'))
153
154 if __name__ == "__main__":
155     app.run(debug=True)

```

3 Part 3: Flask App Setup


After experiencing persistent MySQL connection issues with the Streamlit application and multiple unsuccessful attempts to resolve them, I've decided to transition to using Flask. Flask provides a more stable connection handling mechanism for MySQL, which will improve application reliability and user experience.

3.1 Sidebar Filters:




Filters

Start Date

mm/dd/yyyy 


End Date

mm/dd/yyyy 

Minimum Total Amount

e.g., 1000

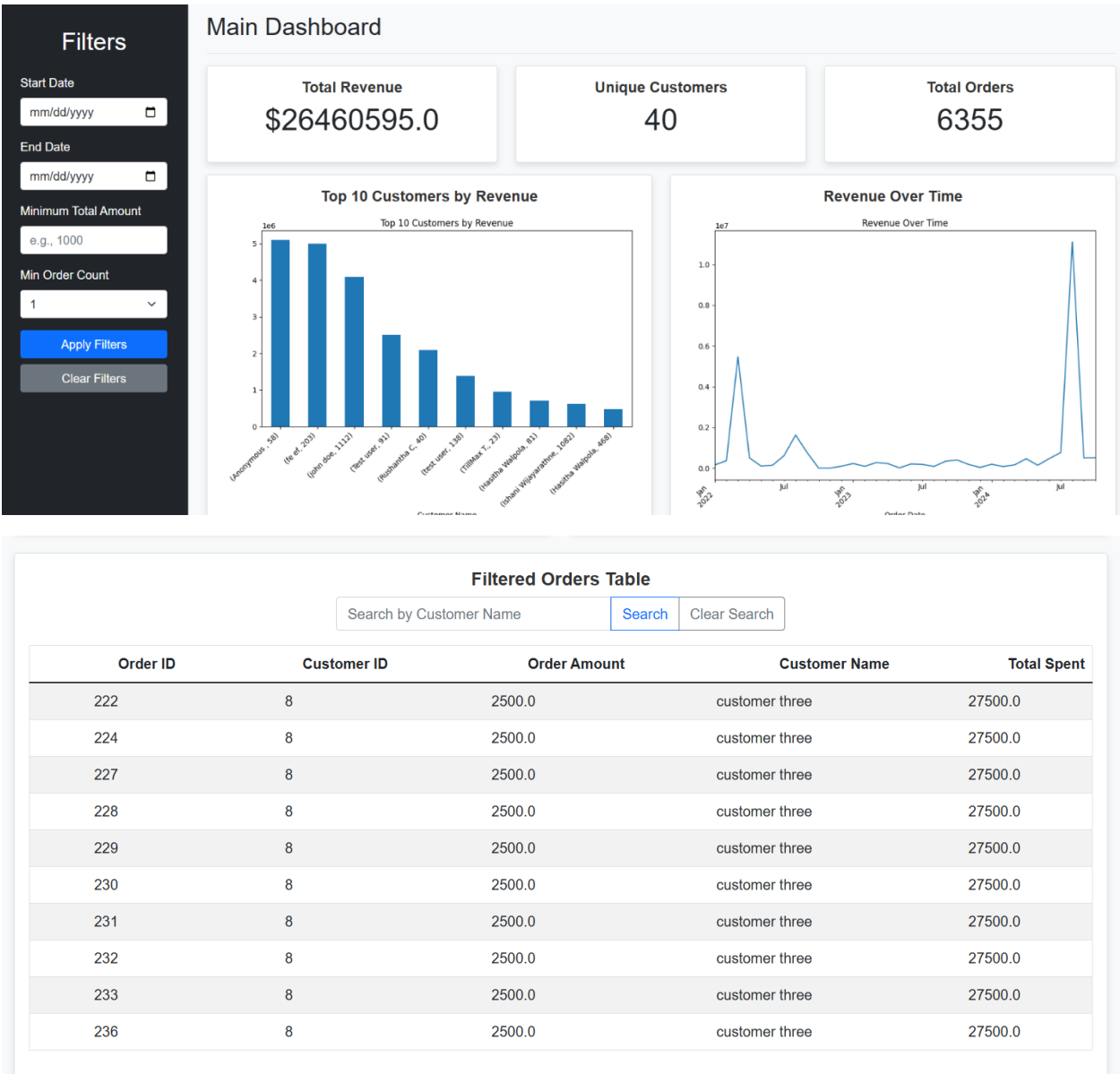
Min Order Count

1 

Apply Filters

Clear Filters

3.2 Main Dashboard:



4 Part 4: Data Analysis

4.1 Machine Learning Model

A simple logistic regression model was successfully implemented to predict repeat purchasers based on total orders and revenue, with validation checks for data sufficiency and accuracy metrics included. Relevant images illustrating the model and results are provided.

Inspect the the datasets

```
#First few rows of orders_df  
orders_df.head()
```

	id	display_order_id	total_amount	created_at	customer_id
0	13392	YTFA	425	2024-10-14 15:12:43	1251.0
1	13393	N1U7	1650	2024-10-14 15:17:25	1251.0
2	13394	PADV	1365	2024-10-14 17:02:16	468.0
3	13395	OKVW	525	2024-10-14 17:03:36	1251.0
4	13396	2G3Q	525	2024-10-14 17:04:49	468.0

Data Preprocessing

Checking Missing Values

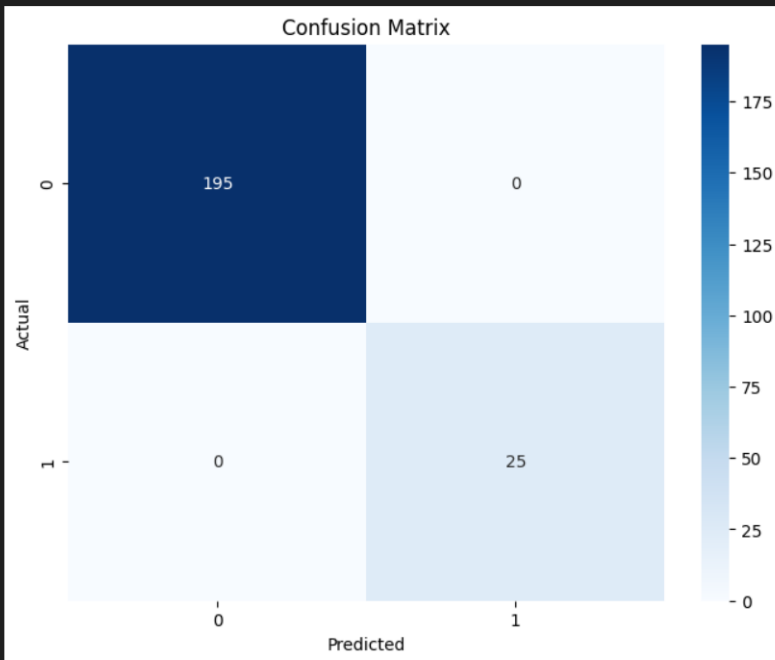
```
print(customers_df.isnull().sum())
```

```
customer_id    0  
name           6  
email          18  
dtype: int64
```

```
print(orders_df.isnull().sum())
```

```
id              0  
display_order_id  36  
total_amount    0  
created_at      0  
customer_id    125  
dtype: int64
```

Model Accuracy: 100.00%



Cross-Validation Accuracy Scores: [1. 1. 1. 1. 1.]

Mean CV Accuracy: 100.00%

