# ASSIGNMENT

### ~Data and Cloud Analytics internship~

**MUȘAT ANDRA**

# CUPRINS

# A) Provision an Azure SQL Database

## As a student account holder, I took the following steps:

* Before provisioning any resources, I created a new resource group to organize my resources logically. I named the resource group "InternshipResources-rg" and selected the East US location;
* Provisioning a single SQL database: From the search bar in Azure, I looked for Azure SQL, clicked on SQL databases, and left the Resource type set to Single database, then selected Create. I selected my student subscription and proceeded to give my database a name, "andra_internship-db". For Server, I selected Create new and filled out the New server form with the following values: Server name: "andra-internship-srv", Location: East US, Authentication method: Use SQL authentication, created a Server admin login and password below. After that, I selected the Workload Environment to be Development and clicked on Configure Database. I selected the Service Tier to be General Purpose and Compute tier to be Serverless, and applied the changes, ensuring the SQL elastic pool is set to No and Backup storage redundancy option is locally redundant storage. Next, I clicked on Networking and set the Connectivity method to Public endpoint. For Firewall rules, I set Add current client IP address to Yes, and left Allow Azure services and resources to access this server set to No. Then, I clicked on security to proceed and kept Enable Microsoft Defender for SQL set to No. Next, I clicked on additional settings and kept Use existing data set to None to create an empty database. Finally, I clicked on Review and Create and then on Create.
* Waited for Deployment: Azure deployed my SQL Server instance and database. This process took a few minutes. I monitored the deployment progress from the notification area in the upper right area of the portal;
* Verified Deployment: Once the deployment was complete, I visited the Azure SQL area again to ensure that both the SQL Server instance and the database were provisioned successfully and were in a running state;
* As an additional task, I went to my newly created database, navigated to the Query editor, and created a simple table (flowers table) using an SQL query;
* At the end, I took screenshots.

# B) In debt analysis of a dataset

### 1. I used QuickDBD



### 2. I expect the type of relationship between tables to be:

* The relationship between "pizzas" and "pizza_types" is a one to many relationship. Each pizza belongs to only one pizza type, but each pizza type can have multiple pizzas associated with it.
* The relationship between "pizzas" and "order_details" is a one to many relationship. Each pizza can have multiple order details associated with it, but each order detail is related to only one pizza.
* The relationship between "orders" and "order_details" is a one to many relationship. Each order can have multiple order details associated with it, but each order detail is related to only one order.

## 3. For this requirement:

* I split the sheets with tables from the given Excel file and saved them in CSV format. Then, using MS SQL from the online SQL server, I imported the CSV files with tables, and for each column in the tables, I adjusted the data type and added the necessary constraints.
* The code:

```
SELECT * FROM pizza_types;

ALTER TABLE pizza_types
ALTER COLUMN pizza_type_id VARCHAR(30) NOT NULL;
ALTER TABLE pizza_types
ALTER COLUMN name VARCHAR(50);
ALTER TABLE pizza_types
ALTER COLUMN category VARCHAR(30);
ALTER TABLE pizza_types
ALTER COLUMN ingredients VARCHAR(100);
ALTER TABLE pizza_types
ADD CONSTRAINT pk_pizza_type_id PRIMARY KEY (pizza_type_id);

SELECT * FROM pizzas;

ALTER TABLE pizzas
ALTER COLUMN pizza_id VARCHAR(30) NOT NULL;
ALTER TABLE pizzas
ALTER COLUMN pizza_type_id VARCHAR(30);
ALTER TABLE pizzas
ALTER COLUMN size VARCHAR(3);
ALTER TABLE pizzas
ALTER COLUMN price DECIMAL(6,2);
ALTER TABLE pizzas
ADD CONSTRAINT pk_pizza_id PRIMARY KEY (pizza_id);
ALTER TABLE pizzas
add CONSTRAINT fk_pizza_type_id FOREIGN KEY(pizza_type_id) REFERENCES
pizza_types(pizza_type_id);

SELECT TOP 100 * FROM orders;

ALTER TABLE orders
ALTER COLUMN order_id INT NOT NULL;
ALTER TABLE orders
ALTER COLUMN date DATE;
ALTER TABLE orders
ALTER COLUMN time TIME;
alter TABLE orders
```

ADD CONSTRAINT pk_order_id PRIMARY KEY (order_id);


SELECT TOP 100 * FROM order_details;

ALTER TABLE order_details
ALTER COLUMN order_details_id INT NOT NULL;
ALTER TABLE order_details
ALTER COLUMN order_id INT NOT NULL;
ALTER TABLE order_details
ALTER COLUMN pizza_id VARCHAR(30);
ALTER TABLE order_details
ALTER COLUMN quantity INT;
ALTER TABLE order_details
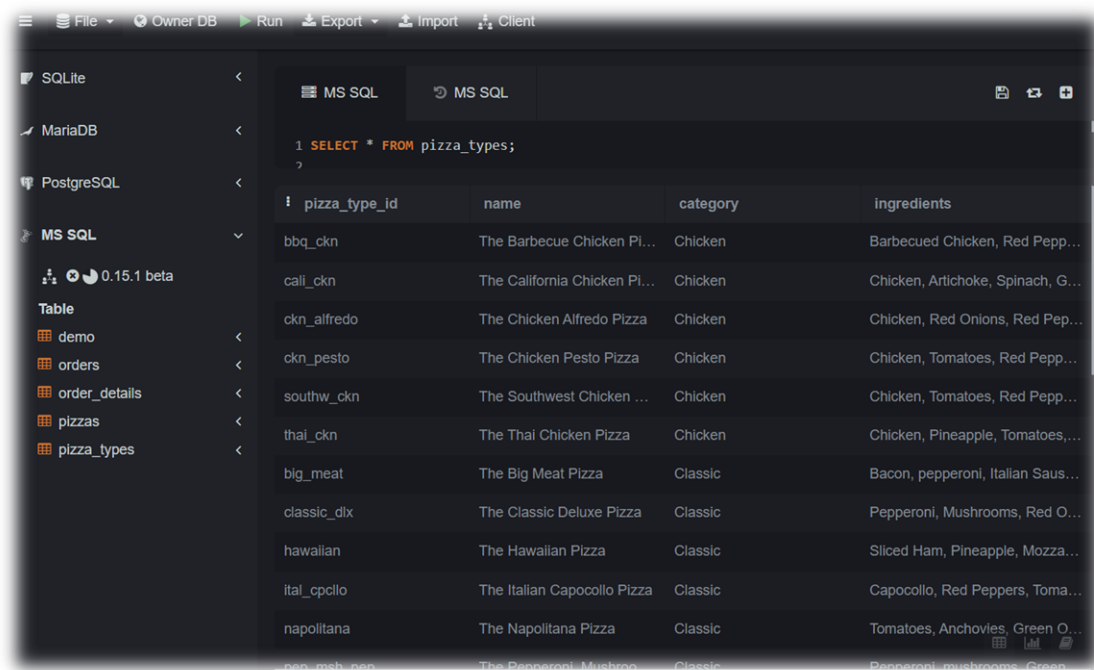ADD CONSTRAINT pk_order_details_id PRIMARY KEY (order_details_id);
ALTER TABLE order_details
ADD CONSTRAINT fk_order_id FOREIGN KEY(order_id) REFERENCES orders(order_id);
ALTER TABLE order_details
ADD CONSTRAINT fk_pizza_id FOREIGN KEY(pizza_id) REFERENCES pizzas(pizza_id);


   &ast;   The results:



**pizza_types table**

**pizzas table**



**orders table**

**order_details table**

## 4. I used python

∗ The code:

```python
import pandas as pd

#load the CSV file specifying the encoding(specifying the encoding parameter
#when reading the CSV file helps Pandas interpret the file's contents correctly,
#especially when dealing with characters that fall outside the standard ASCII range.)
df=pd.read_csv(r"D:\internship_diconium\tables\pizza_types.csv", encoding='latin1')

#initialize a dictionary to count ingredient occurrences
ingredient_counts={}

#iterate over each row in the DataFrame and parse the ingredients
for ingredients_str in df['ingredients']:
    ingredients=ingredients_str.split(", ")
    for ingredient in ingredients:
        if ingredient in ingredient_counts:
            ingredient_counts[ingredient]+=1
        else:
            ingredient_counts[ingredient]=1
```

#find the ingredient with the highest occurrence count
most_used_ingredient=max(ingredient_counts, key=ingredient_counts.get)
frequency_of_most_used_ingredient=ingredient_counts[most_used_ingredient]

print("The most used ingredient is:", most_used_ingredient)
print("It appears", frequency_of_most_used_ingredient, "times.")

* The result:

## 5. I used python and the jupyter notebook extension to display the result more clearly

* The code:

```python
import pandas as pd

#load the CSV file into a DataFrame
df=pd.read_csv(r"D:\internship_diconium\tables\pizza_types.csv", encoding='latin1')

#initialize a set to store unique types of cheese
cheese_types=set()

#function to split ingredients by both commas and semicolons
def split_ingredients(ingredients_str):
    return [item.strip() for item in ingredients_str.replace(';', ',').split(',')]

#iterate over each row in the DataFrame and parse the ingredients
for ingredients_str in df['ingredients']:
    ingredients=split_ingredients(ingredients_str)
    for ingredient in ingredients:
        if "cheese" in ingredient.lower():  #check if the word 'cheese' is in the ingredient
            cheese_types.add(ingredient.strip())

#convert the set of unique cheeses to a sorted list
unique_cheeses=sorted(cheese_types)

# Print the results
print(f"The restaurant uses {len(unique_cheeses)} types of cheese.")
print("The types of cheese used are:")
for cheese in unique_cheeses:
    print(cheese)

#create a DataFrame with an ID column and a Cheese column
cheese_df=pd.DataFrame({'ID': range(1, len(unique_cheeses)+1), 'Cheese': unique_cheeses})

#print the DataFrame with a title
print("List of Cheese Types Used by the Restaurant")
cheese_df
```

* The result:

## 6. I used python and the jupyter notebook extension to display the result more clearly

* The code:

import pandas as pd

#load the CSV file into a pandas DataFrame
pizzas_df=pd.read_csv(r"D:\internship_diconium\tables\pizzas.csv")

#group by pizza_type_id and size, then count the number of pizzas sold for each group
sales_summary=pizzas_df.groupby(['pizza_type_id', 'size']).size().reset_index(name='sold_quantity')

#display the sales summary DataFrame as a table
sales_summary

* The result:

```
✓ import pandas as pd ...

      pizza_type_id   size   sold_quantity
0         bbq_ckn       L            1
1         bbq_ckn       M            1
2         bbq_ckn       S            1
3         big_meat      L            1
4         big_meat      M            1
...          ...       ...          ...
91        the_greek     XL           1
92        the_greek     XXL          1
93        veggie_veg    L            1
94        veggie_veg    M            1
95        veggie_veg    S            1

96 rows × 3 columns
```

## 7. I used python and the jupyter notebook extension to display the result more clearly

∗ The code:

import pandas as pd

#load the order details CSV file into a pandas DataFrame
order_details_df=pd.read_csv(r"D:\internship_diconium\tables\order_details.csv")

#group by pizza_id and sum the quantities to get the total quantity of each pizza ordered
pizza_orders=order_details_df.groupby('pizza_id')['quantity'].sum().reset_index()

#sort the pizzas by total quantity in descending order and select the top 5
top_5_pizzas=pizza_orders.sort_values(by='quantity', ascending=False).head(5)

#display the top 5 most ordered pizzas as a table
top_5_pizzas

∗ The result:

## 8. I used python and the jupyter notebook extension to display the result more clearly

* The code:

```python
import pandas as pd
from pandasql import sqldf

#reading data from CSV files
order_details_df=pd.read_csv(r"D:\internship_diconium\tables\order_details.csv")
orders_df=pd.read_csv(r"D:\internship_diconium\tables\orders.csv")
pizzas_df=pd.read_csv(r"D:\internship_diconium\tables\pizzas.csv")

#converting the date column to datetime format
orders_df['date']=pd.to_datetime(orders_df['date'])

#defining the SQL query for LEFT JOIN
right_join_query="""
        SELECT DISTINCT p.pizza_id
        FROM pizzas_df p
        LEFT JOIN order_details_df od ON p.pizza_id=od.pizza_id
        LEFT JOIN orders_df o ON od.order_id=o.order_id
        WHERE o.date NOT LIKE '2015%'
          OR o.date IS NULL
        """

#executing the query
right_result=sqldf(right_join_query, locals())

#displaying the result as a DataFrame
right_result
```

* The result:



16

## 9. I used python and the jupyter notebook extension to display the result more clearly

* The code:

```python
import pandas as pd

#read the pizzas CSV file into a DataFrame
pizzas_df = pd.read_csv(r"D:\internship_diconium\tables\pizzas.csv")

#define a function to categorize the sizes
def categorize_size(size):
    if 's' in size.lower():
        return 'Small'
    elif 'm' in size.lower():
        return 'Medium'
    elif 'l' in size.lower():
        return 'Large'
    else:
        return 'Unknown'

#apply the function to create a new column 'size_category'
pizzas_df['size_category']=pizzas_df['size'].apply(categorize_size)

#group the pizzas by 'size_category' and calculate summary statistics for 'price'
summary_statistics=pizzas_df.groupby('size_category')['price'].agg(['mean', 'min', 'max'])

#display the summary statistics table
summary_statistics
```
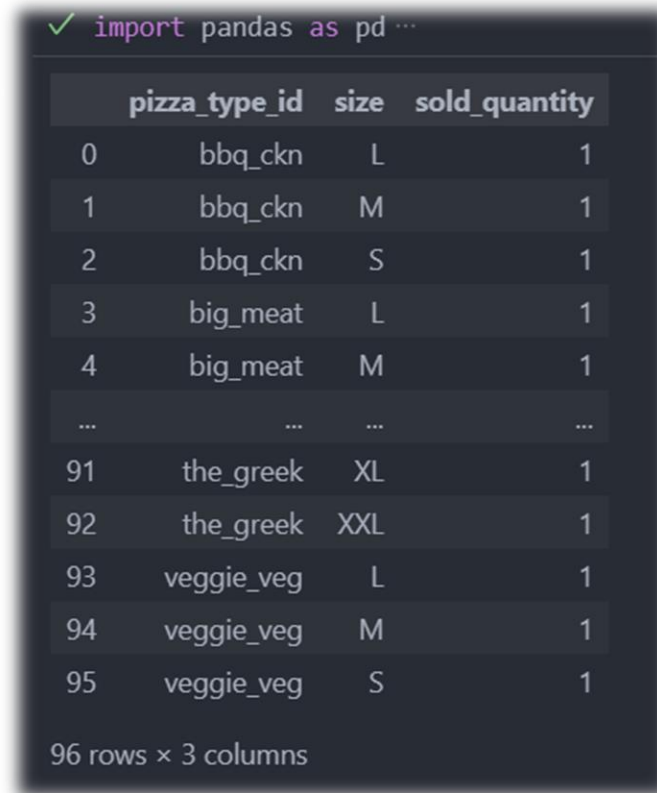
* The result:

## 10. I used python and the jupyter notebook extension to display the result more clearly

* The code:

```python
import pandas as pd
import matplotlib.pyplot as plt

#load the pizzas CSV file into a pandas DataFrame
pizzas_df = pd.read_csv(r"D:\internship_diconium\tables\pizzas.csv")

#plot the distribution of prices using a histogram
plt.figure(figsize=(10, 6))
plt.hist(pizzas_df['price'], bins=20, color='purple', edgecolor='black')
plt.title('Distribution of Pizza Prices')
plt.xlabel('Price')
plt.ylabel('Frequency')
plt.grid(True)
plt.show()
```

* The result:

## 11. I used python and the jupyter notebook extension to display the result more clearly

* The code:

```
import pandas as pd
from pandasql import sqldf

# Load the CSV files into pandas DataFrames
pizzas_df=pd.read_csv(r"D:\internship_diconium\tables\pizzas.csv")
orders_df=pd.read_csv(r"D:\internship_diconium\tables\orders.csv")
order_details_df=pd.read_csv(r"D:\internship_diconium\tables\order_details.csv")

#converting the date column to datetime format
orders_df['date']=pd.to_datetime(orders_df['date'])

#define the SQL query to join the tables and filter transactions in 2015
query = """
    SELECT COALESCE(ROUND(SUM(od.quantity*p.price)), 0) AS total_revenue_2015
    FROM pizzas_df p
    LEFT JOIN order_details_df od ON p.pizza_id=od.pizza_id
    LEFT JOIN orders_df o ON od.order_id=o.order_id
    WHERE o.date LIKE '2015%'
    """

#execute the query using pandasql
result=sqldf(query, locals())

#print the total revenue generated in 2015
print("Total revenue generated in 2015:", result['total_revenue_2015'].values[0])
```
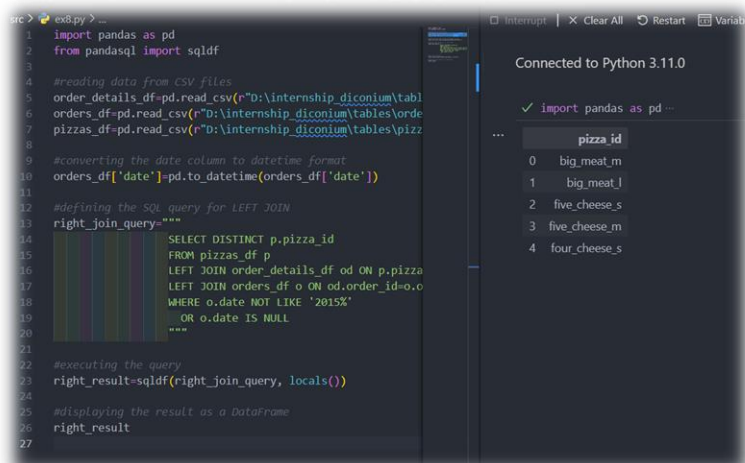
* The result:



19

## 12. I used python and the jupyter notebook extension to display the result more clearly

* The code:

```python
import pandas as pd
import matplotlib.pyplot as plt
from pandasql import sqldf

#load the CSV files into pandas DataFrames
pizzas_df=pd.read_csv(r"D:\internship_diconium\tables\pizzas.csv")
orders_df=pd.read_csv(r"D:\internship_diconium\tables\orders.csv")
order_details_df=pd.read_csv(r"D:\internship_diconium\tables\order_details.csv")

#converting the date column to datetime format
orders_df['date']=pd.to_datetime(orders_df['date'])

#define the SQL query to join the tables and filter transactions in 2015, grouped by month
query = """
    SELECT SUBSTR(o.date, INSTR(o.date, '/')+1) AS month,
        COALESCE(SUM(od.quantity*p.price), 0) AS revenue
    FROM pizzas_df p
    LEFT JOIN order_details_df od ON p.pizza_id=od.pizza_id
    LEFT JOIN orders_df o ON od.order_id=o.order_id
    WHERE o.date LIKE '2015%'
    GROUP BY month
    """

#execute the query using pandasql
result=sqldf(query, locals())

#plot the revenue for each month
plt.figure(figsize=(10, 6))
plt.bar(result['month'], result['revenue'], color='purple')
plt.xlabel('Month')
plt.ylabel('Revenue')
plt.title('Revenue Generated Each Month in 2015')
plt.grid(axis='y', linestyle='--', alpha=0.7)
plt.show()
```

* The result:



* Conclusion: I observe that revenues seem to vary depending on the month. For example, we can see that revenues appear to be higher during the warm months and lower during the cold months. This variation may be influenced by seasonal factors such as vacations or holidays.

## 13. I used python and the jupyter notebook extension to display the result more clearly

* The code:

```
import pandas as pd
import matplotlib.pyplot as plt

#load the orders CSV into a pandas DataFrame
orders_df=pd.read_csv(r"D:\internship_diconium\tables\orders.csv")

#convert the 'time' column to datetime with specified format
orders_df['time']=pd.to_datetime(orders_df['time'], format='%H:%M:%S')
```

```
#extract the hour, minute, and second components from the 'time' column
orders_df['hour']=orders_df['time'].dt.hour
orders_df['minute']=orders_df['time'].dt.minute
orders_df['second']=orders_df['time'].dt.second

#group the orders by hour, minute, and second and count the number of orders for each timestamp
order_frequency=orders_df.groupby(['hour', 'minute', 'second']).size()

#plot the order frequency over time
plt.figure(figsize=(12, 6))
order_frequency.plot(kind='line', color='purple')
plt.title('Order Frequency Over Time')
plt.xlabel('Time')
plt.ylabel('Number of Orders')
plt.grid(True)
plt.show()
```

* The result:



* Conclusion: The graph shows order frequency throughout the day, revealing busy peak hours and quieter times. This insight aids in resource planning for better customer service.

## 14. I used python and the jupyter notebook extension to display the result more clearly

* The code:

import pandas as pd

#load the orders CSV into a pandas DataFrame
orders_df=pd.read_csv(r"D:\internship_diconium\tables\orders.csv")

#extract the date component from the 'date' column
orders_df['date']=pd.to_datetime(orders_df['date'])

#group the orders by date and count the number of orders for each date
daily_orders=orders_df.groupby(orders_df['date'].dt.date).size()

#calculate the average number of orders per day
average_orders_per_day=round(daily_orders.mean())

print("Average orders per day:", average_orders_per_day)

* The result:

## 15. I used python and the jupyter notebook extension to display the result more clearly

* The code:

```python
import pandas as pd
import matplotlib.pyplot as plt

#load the orders CSV into a pandas DataFrame
orders_df=pd.read_csv(r"D:\internship_diconium\tables\orders.csv")

#convert the 'date' column to datetime format
orders_df['date']=pd.to_datetime(orders_df['date'])

#extract the day of the week from the 'date' column
orders_df['day_of_week']=orders_df['date'].dt.day_name()

#group the orders by the day of the week and count the number of orders for each day
orders_by_day=orders_df.groupby('day_of_week').size()

#plot the order frequency for each day of the week
plt.figure(figsize=(10, 6))
orders_by_day.plot(kind='bar', color='purple')
plt.title('Order Frequency by Day of the Week')
plt.xlabel('Day of the Week')
plt.ylabel('Number of Orders')
plt.xticks(rotation=45)
plt.grid(axis='y')
plt.show()
```
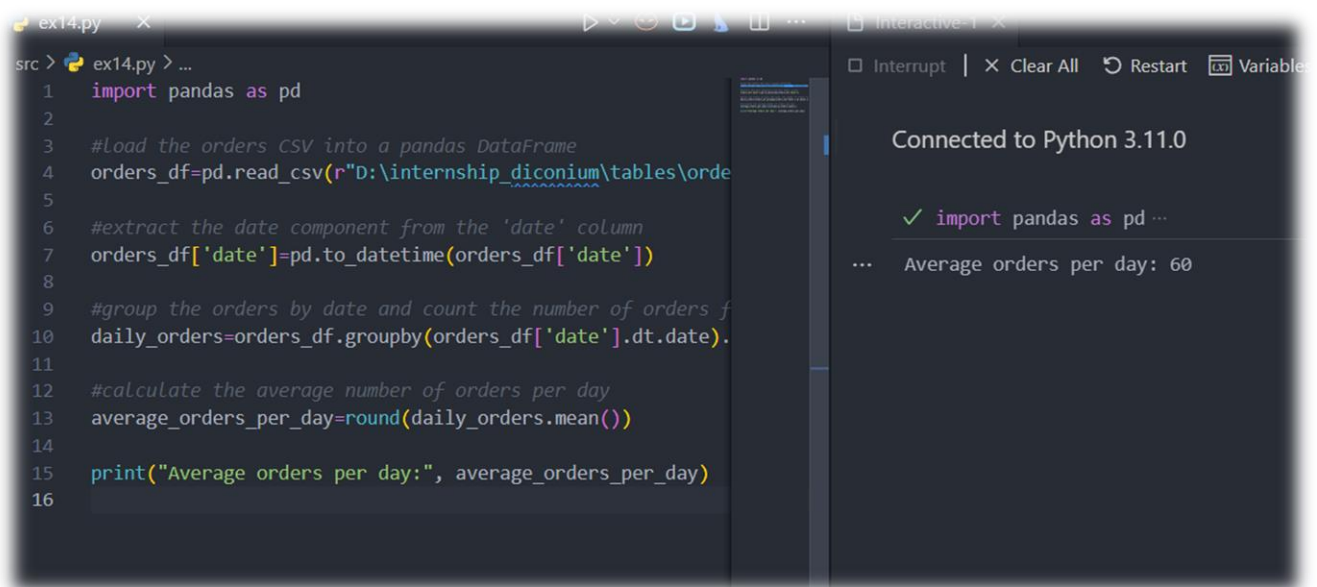
* The result:

* Conclusion: Based on the graph, it seems like Fridays are the most crowded days of the week. This could be because people tend to dine out more often as the workweek comes to a close, or perhaps it's influenced by social customs and traditions.

## 16. I used python and the jupyter notebook extension to display the result more clearly

* The code:

```
import pandas as pd
from pandasql import sqldf
import matplotlib.pyplot as plt

#load the CSV files into pandas DataFrames
pizzas_df=pd.read_csv(r"D:\internship_diconium\tables\pizzas.csv")
orders_df=pd.read_csv(r"D:\internship_diconium\tables\orders.csv")
order_details_df=pd.read_csv(r"D:\internship_diconium\tables\order_details.csv")

#clean the date column in orders_df and handle potential issues
orders_df['date']=pd.to_datetime(orders_df['date'], format='%m/%d/%Y', errors='coerce')
orders_df=orders_df.dropna(subset=['date'])  # Drop rows with invalid dates

#add a column for the day of the week to the orders_df
orders_df['day_of_week']=orders_df['date'].dt.day_name()

#define the SQL query to calculate revenue for each day of the week
query = """
    SELECT o.day_of_week, COALESCE(SUM(od.quantity*p.price), 0) AS revenue
    FROM order_details_df od
    LEFT JOIN pizzas_df p ON p.pizza_id=od.pizza_id
    LEFT JOIN orders_df o ON od.order_id=o.order_id
    GROUP BY o.day_of_week
    ORDER BY revenue DESC
    """

#execute the query using pandasql
result=sqldf(query, locals())

#plot the results
plt.figure(figsize=(10, 6))
plt.bar(result['day_of_week'], result['revenue'], color='purple')
plt.title('Revenue Generated by Day of the Week')
plt.xlabel('Day of the Week')
```
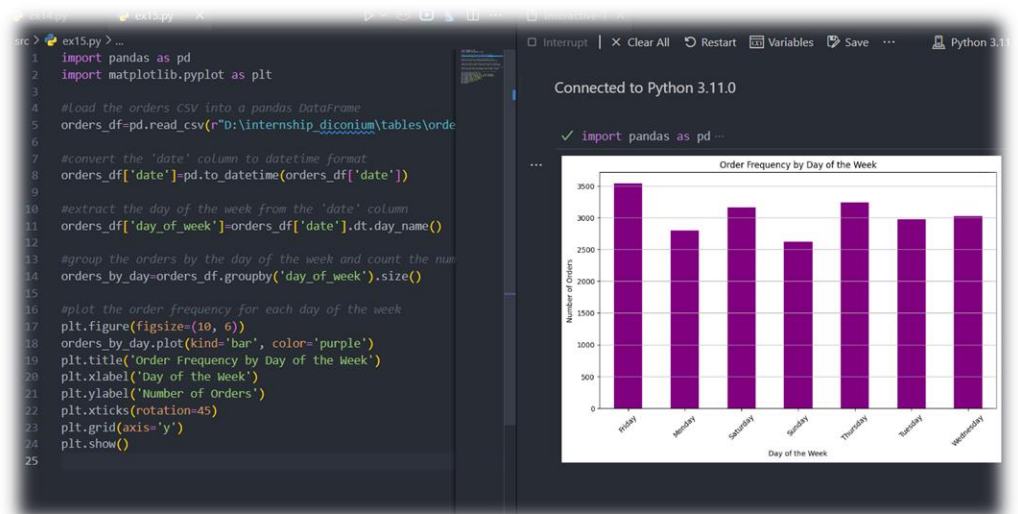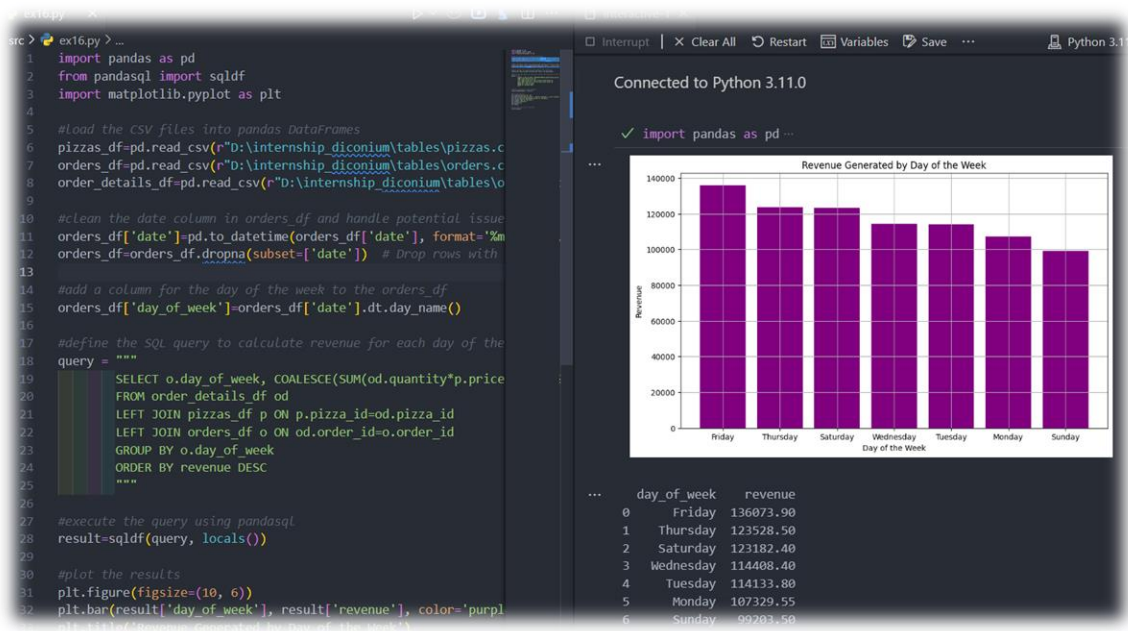
plt.ylabel('Revenue')
plt.grid(True)
plt.show()

#print the result for reference
print(result)

* The result:



* Conclusion: Based on the graph, you can observe which day of the week generates the highest revenue. This insight helps understand customer behavior and optimize restaurant operations. For instance, if Fridays show the highest revenue, it could be due to increased dining out as people celebrate the end of the workweek.

## 17. I used python and the jupyter notebook extension to display the result more clearly

* The code:

```
import pandas as pd
from pandasql import sqldf

#load the CSV files into pandas DataFrames
pizzas_df=pd.read_csv(r"D:\internship_diconium\tables\pizzas.csv")
orders_df=pd.read_csv(r"D:\internship_diconium\tables\orders.csv")
order_details_df=pd.read_csv(r"D:\internship_diconium\tables\order_details.csv")
```

#convert the 'date' column to datetime
orders_df['date']=pd.to_datetime(orders_df['date'], format='%m/%d/%Y')

#create a DataFrame with a full range of dates in 2015
full_date_range=pd.date_range(start='2015-01-01', end='2015-12-31')
full_date_df=pd.DataFrame(full_date_range, columns=['date'])

#define the SQL queries to join the tables and calculate daily revenue
query="""
    SELECT fd.date, COALESCE(SUM(od.quantity*p.price), 0) AS daily_revenue
    FROM full_date_df fd
    LEFT JOIN orders_df o ON fd.date=o.date
    LEFT JOIN order_details_df od ON o.order_id=od.order_id
    LEFT JOIN pizzas_df p ON od.pizza_id=p.pizza_id
    GROUP BY fd.date
    """

#execute the query using pandasql
full_revenue_df=sqldf(query, globals())

#identify the dates with no sales (daily_revenue is 0)
no_sales_dates=full_revenue_df[full_revenue_df['daily_revenue']==0]

#display the periods with no pizza sales
print("Periods with no pizza sales:")
no_sales_dates

* The result:



* Conclusion: The absence of pizza sales during certain periods could be attributed to factors like restaurant closures, changes in customer behavior, or external influences like holidays or events.