# Crime Data Exploration

Mini Capstone Project: Crime Data Analysis With Mysql And Python
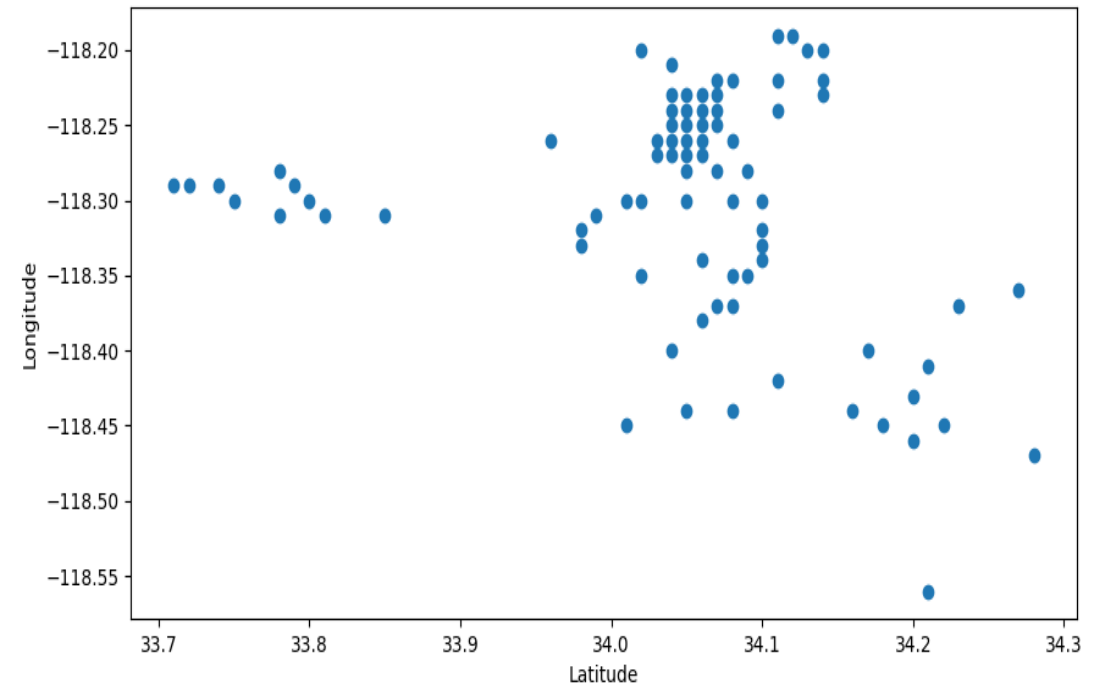
# Contents

# Project Overview

We're diving deep into crime data to understand it better. By analyzing spatial patterns,

victim demographics, location-specific occurrences, and crime code data, we aim to

unveil hidden trends, correlations, and hotspots within the dataset.

By finding these patterns, we can help make neighbourhood safer by focusing on the

right areas and issues. Our goal is to use facts and numbers to guide decisions and

actions, making our communities safer and reducing crime.

# Spatial Analysis

In our spatial analysis, we pinpoint geographical hotspots for reported crimes using a scatterplot. Focusing on latitude and longitude, we'll map crime incidents. This helps identify areas with high crime rates, aiding in targeted law enforcement and crime prevention strategies.
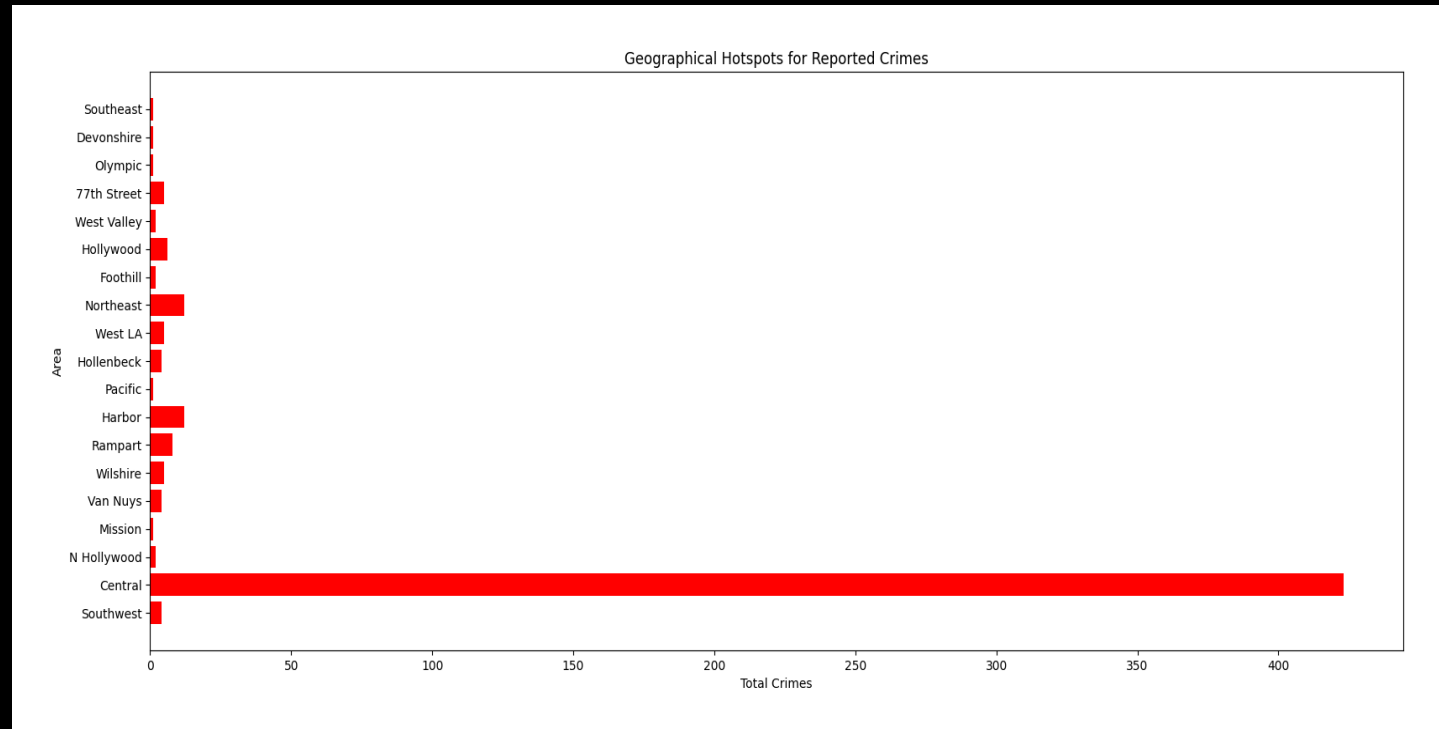
# Spatial Analysis

In another spatial analysis, we employ a bar

graph to pinpoint specific hotspots

of reported crimes. By analyzing crime

frequency across different areas,

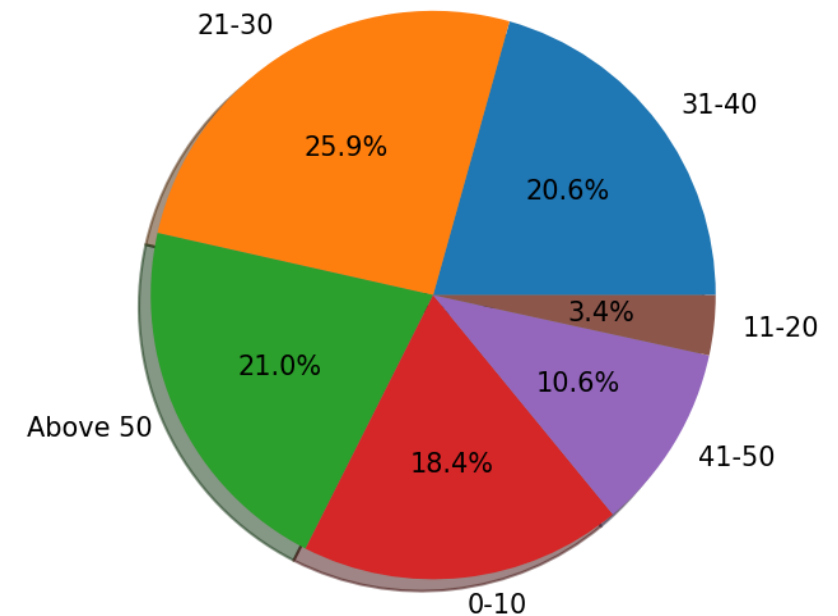we identify and highlight specific locations

with the highest crime rates.

This targeted approach aids in allocating

resources effectively for crime mitigation

strategies.

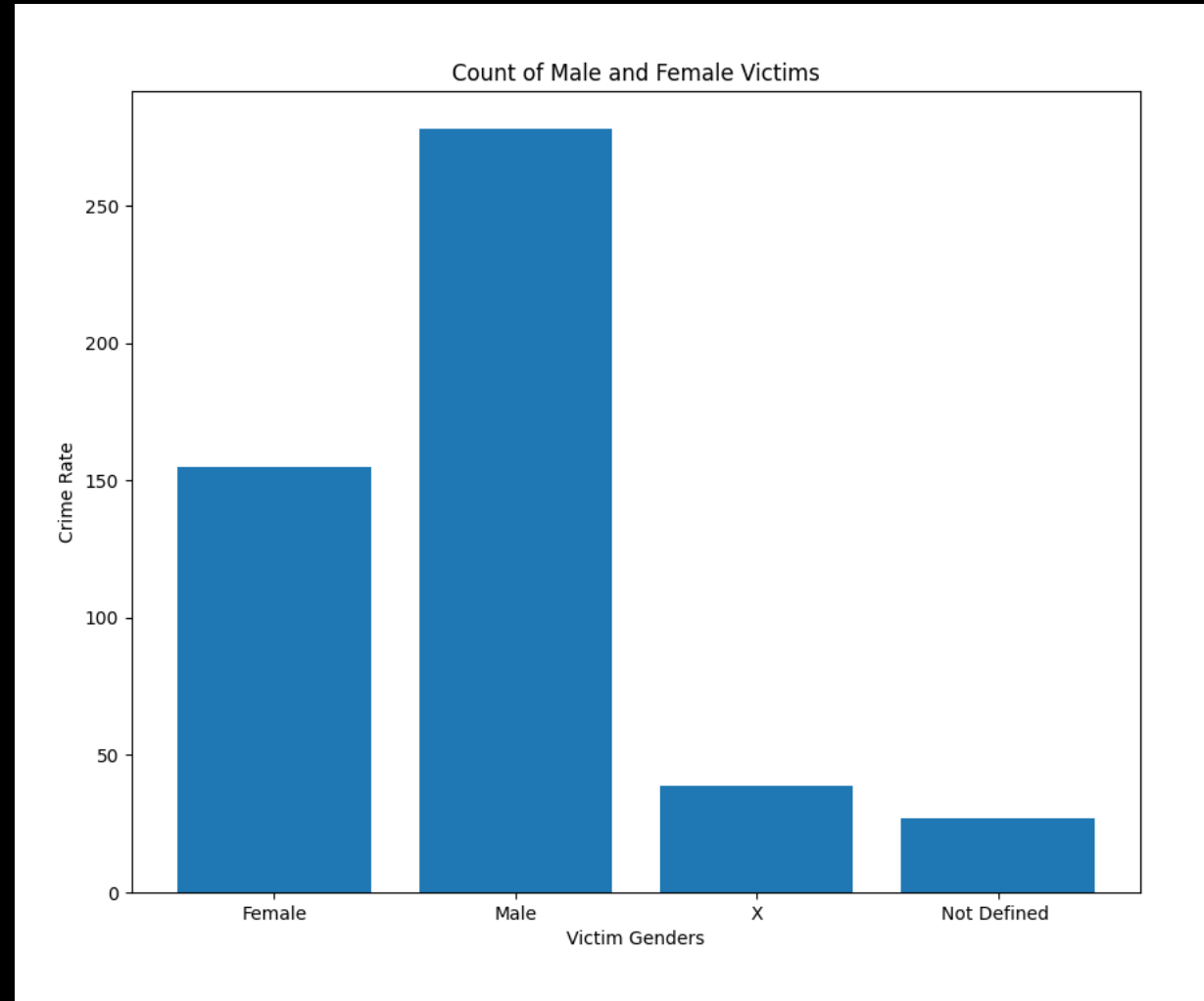# Victim Demographics

In Victim Demographics, we examine the distribution of victim ages in reported crimes using a pie chart. Predominantly, victims fall within the 21 to 30 age group, indicating a significant proportion of young individuals affected. This insight informs targeted interventions to address vulnerabilities within this demographic.



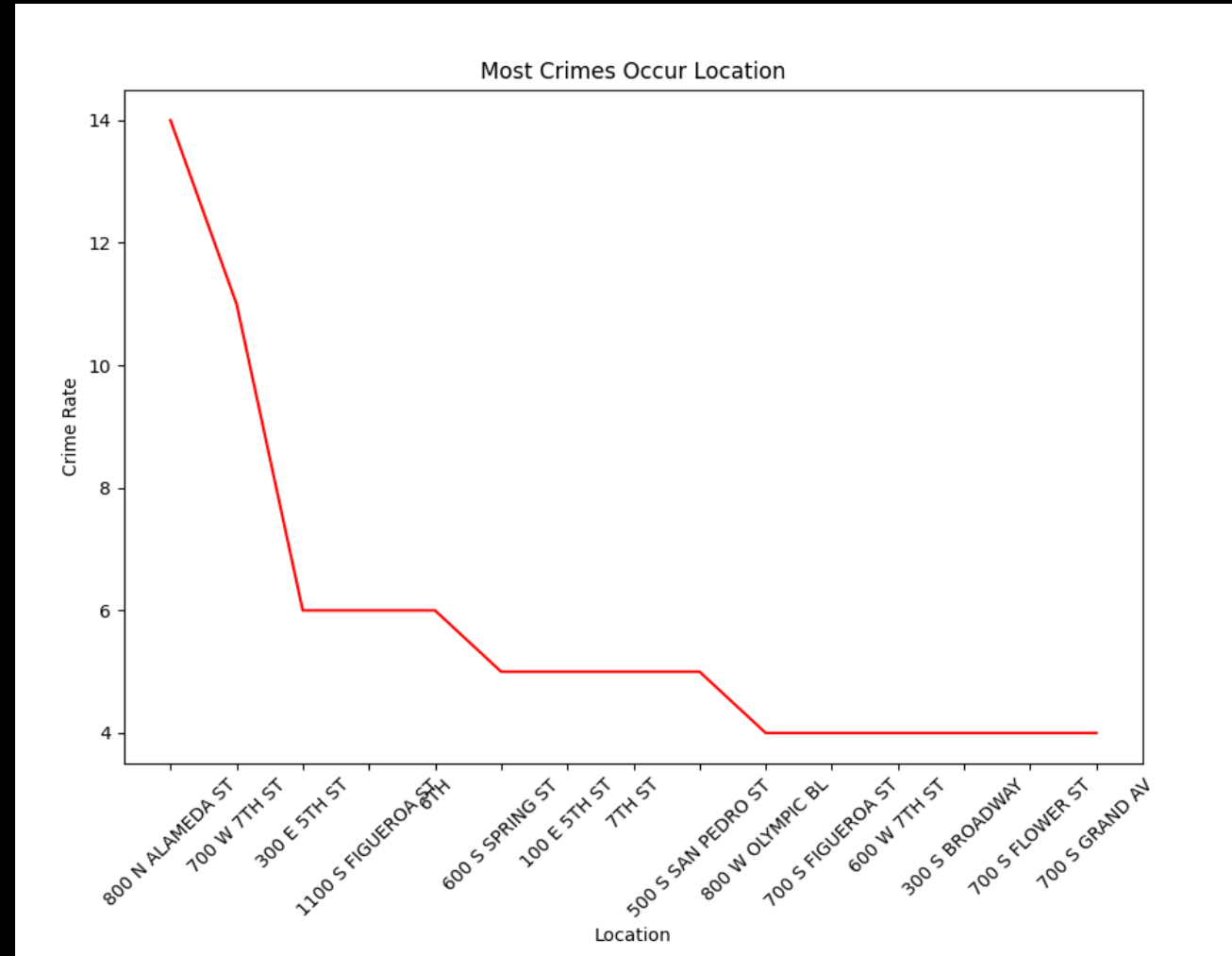Distribution of Victim Ages in Reported Crimes

# Victim Demographics

In another study of victim demographics, we used a bar graph to find difference in crime rates between males and females. The data showed that males experience more crimes, making up the majority of victims. This highlights the importance of gender-aware crime prevention efforts to tackle this issue effectively.



Count of Male and Female Victims

# Location Analysis

In Location Analysis, we identify locations with the highest crime rates. Using a line graph, we pinpoint these crime hotspots. Here, we can see '800 N Alameda ST' is the prime location of crimes. This data-driven approach helps focus law enforcement resources on key areas, improving crime management and public safety strategies.

# Crime Code Analysis

In Crime Code Analysis, we analyze the distribution of reported crimes based on Crime Code using a count plot. This visualization reveals the types of crimes and their respective frequencies by code. Understanding these patterns aids in developing targeted crime prevention measures and allocating resources more effectively.



Distribution of Reported Crimes by Crime Code

# Tools And Libraries

I used Visual Studio Code (VS Code) for writing

Python code. It's great because it has lots of

useful features and is easy to use. The

debugging tools and version control features

helped me a lot, and I could customize it to suit

my needs, making coding more efficient and

enjoyable.

# PyMySQL

The PyMySQL library in Python enables seamless interaction with MySQL databases, allowing you to execute SQL queries, retrieve data, and perform database operations directly from your Python code. It simplifies database connectivity tasks, making it easier to work with MySQL databases in various Python applications such as data analysis, web development, and automation scripts.



PYTHON AND MYSQL WITH
PYMYSQL

# Matplotlib and Seaborn

Matplotlib and Seaborn are powerful Python libraries for data visualization. Matplotlib offers a wide range of customizable plots like line charts, histograms, and scatter plots, ideal for exploring data and presenting insights. Seaborn, built on top of Matplotlib, provides additional functionalities and beautiful default styles for creating appealing statistical graphics, making it a favorite for data analysts and researchers worldwide.

# Python Scripts

database setup, data import

```
PROJECT q0.py > ...
 1   import pymysql
 2   import mysql.connector
 3   import pandas as pd
 4   import numpy as np
 5   import matplotlib.pyplot as plt
 6   import warnings
 7   warnings.filterwarnings("ignore")
 8
 9
10   connection = pymysql.connect(host= "localhost",
11                                user= "root",
12                                password= "tyjkl89",
13                                database= "crimedata"
14   )
15
16   print(connection)
17
18   qry1 = "select * from crime_data"
19   qry2 = "select count(*) from crime_data"
20   qry3 = "select distinct(crm_cd) from crime_data"
21
22   df1 = pd.read_sql(qry1, connection)
23   df2 = pd.read_sql(qry2, connection)
24   df3 = pd.read_sql(qry3, connection)
25
26
27
28
29   pd.set_option("display.max_rows", None)
30
31   print(df1)
32   print(df2)
33   print(df3)
```

## Spatial Analysis (bar graph)

```
PROJECT q1.py > ...
  1   # Spatial Analysis:
  2
  3   # Where are the geographical hotspots for reported crimes?
  4
  5
  6   import pymysql
  7   import mysql.connector
  8   import pandas as pd
  9   import numpy as np
 10   import matplotlib.pyplot as plt
 11   import seaborn as sns
 12   import plotly.express as px
 13   import geopandas as gpd
 14   import warnings
 15   warnings.filterwarnings("ignore")
 16
 17
 18   connection = pymysql.connect(host= "localhost",
 19                                user= "root",
 20                                password= "tyjkl89",
 21                                database= "crimedata"
 22   )
 23
 24   print(connection)
 25
 26   qry = "select AREA_NAME as area, count(DR_NO) as total_crimes from crime_data group by area"
 27
 28   df = pd.read_sql(qry, connection)
 29
 30
 31   pd.set_option("display.max_rows", None)
 32
 33   print(df)
 34
 35   connection.close()
 36
 37
 38   plt.figure(figsize=(20,8))
 39   plt.barh(df["area"],df["total_crimes"], color= "red")
 40   plt.xlabel("Total Crimes")
 41   plt.ylabel("Area")
 42   plt.title("Geographical Hotspots for Reported Crimes")
 43   plt.savefig("question1.png")
 44   plt.show()
```

## Spatial Analysis (scatterplot)

```
import pymysql
import mysql.connector
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
import plotly.express as px
import geopandas as gpd
import warnings
warnings.filterwarnings("ignore")


connection = pymysql.connect(host= "localhost",
                             user= "root",
                             password= "tyjkl89",
                             database= "crimedata"
)

print(connection)

qry = """SELECT LAT, LON, COUNT(*) AS CrimeCount
FROM crime_data
GROUP BY LAT, LON"""
df = pd.read_sql(qry, connection)
pd.set_option("display.max_rows", None)
print(df)
connection.close()

plt.figure(figsize=(10,5))
plt.scatter(df["LAT"], df["LON"])
plt.xlabel("Latitude")
plt.ylabel("Longitude")
plt.show()
```

## Victim Demographics (pie chart)

```
PROJECT q2.py > ...
1   # Victim Demographics:
2
3   # What is the distribution of victim ages in reported crimes?
4
5
6   import pymysql
7   import mysql.connector
8   import pandas as pd
9   import numpy as np
10  import matplotlib.pyplot as plt
11  import seaborn as sns
12  import plotly.express as px
13  import geopandas as gpd
14  import warnings
15  warnings.filterwarnings("ignore")
16
17  connection = pymysql.connect(host= "localhost",
18                               user= "root",
19                               password= "tyjk189",
20                               database= "crimedata"
21  )
22
23  print(connection)
24
25
26  qry = """SELECT
27      CASE
28          WHEN Vict_Age BETWEEN 0 AND 10 THEN '0-10'
29          WHEN Vict_Age BETWEEN 11 AND 20 THEN '11-20'
30          WHEN Vict_Age BETWEEN 21 AND 30 THEN '21-30'
31          WHEN Vict_Age BETWEEN 31 AND 40 THEN '31-40'
32          WHEN Vict_Age BETWEEN 41 AND 50 THEN '41-50'
33          ELSE 'Above 50'
34      END AS 'age_dist',
35      COUNT(*) AS reported_crime
36  FROM
37      crime_data
38  GROUP BY age_dist"""
39
40  df = pd.read_sql(qry, connection)
41
42  pd.set_option("display.max_rows", None)
43
44  print(df)
45
46  connection.close()
47
48  plt.figure(dpi= 150)
49
50  df['percentage'] = (df['reported_crime'] / df['reported_crime'].sum()) * 100
51
52  plt.pie(df["percentage"], labels= df["age_dist"],autopct='%1.1f%%', shadow= True)
53
54  plt.title("Distribution of Victim Ages in Reported Crimes")
55
56  plt.savefig("question2.png")
57
58  plt.show()
```

## Victim Demographics (bar graph)

```
PROJECT q3.py > ...
1   # Victim Demographics:
2
3   # Is there a significant difference in crime rates between male and female victims?
4
5   import pymysql
6   import mysql.connector
7   import pandas as pd
8   import numpy as np
9   import matplotlib.pyplot as plt
10  import seaborn as sns
11  import plotly.express as px
12  import geopandas as gpd
13  import warnings
14  warnings.filterwarnings("ignore")
15
16  connection = pymysql.connect(host= "localhost",
17                               user= "root",
18                               password= "tyjk189",
19                               database= "crimedata"
20  )
21
22  print(connection)
23
24
25  qry = """ SELECT
26      Vict_Sex, COUNT(*) AS total_crime
27  FROM
28      crime_data
29  GROUP BY Vict_Sex"""
30
31  df = pd.read_sql(qry, connection)
32
33  pd.set_option("display.max_rows", None)
34
35  print(df)
36
37  connection.close()
38
39  plt.figure(figsize=(10,8))
40
41  plt.bar(df["Vict_Sex"], df["total_crime"])
42
43  plt.xlabel("Victim Genders")
44
45  plt.ylabel("Crime Rate")
46
47  plt.title("Count of Male and Female Victims")
48
49  plt.xticks(df["Vict_Sex"], ["Female", "Male", "X", "Not Defined"])
50
51  plt.savefig("question3.png")
52
53  plt.show()
```

# Location Analysis

```
PROJECT q4.py > ...
1    # Location Analysis:
2
3    # Where do most crimes occur based on the "Location" column?
4
5    import pymysql
6    import mysql.connector
7    import pandas as pd
8    import numpy as np
9    import matplotlib.pyplot as plt
10   import seaborn as sns
11   import plotly.express as px
12   import geopandas as gpd
13   import warnings
14   warnings.filterwarnings("ignore")
15
16   connection = pymysql.connect(host= "localhost",
17                                user= "root",
18                                password= "tyjkl89",
19                                database= "crimedata"
20   )
21
22   print(connection)
23
24   qry = """ SELECT
25       Location, COUNT(DR_NO) AS crimes
26   FROM
27       crime_data
28   GROUP BY Location
29   ORDER BY crimes Desc
30   LIMIT 15"""
31
32   df = pd.read_sql(qry, connection)
33
34   pd.set_option("display.max_rows", None)
35
36   print(df)
37
38   connection.close()
39
40
41   plt.figure(figsize=(10, 8))
42
43   sns.lineplot(x= df["Location"], y= df["crimes"], color= "red")
44
45   plt.xlabel("Location")
46
47   plt.ylabel("Crime Rate")
48
49   plt.title("Most Crimes Occur Location")
50
51   plt.xticks(rotation= 45)
52
53   plt.savefig("question4.png")
54
55   plt.show()
```

# Crime Code Analysis

```python
1    # Crime Code Analysis:
2
3    # What is the distribution of reported crimes based on Crime Code?
4
5    import pymysql
6    import mysql.connector
7    import pandas as pd
8    import numpy as np
9    import seaborn as sns
10   import matplotlib.pyplot as plt
11   import warnings
12   warnings.filterwarnings("ignore")
13
14   connection = pymysql.connect(host= "localhost",
15                                user= "root",
16                                password= "tyjkl89",
17                                database= "crimedata"
18   )
19
20   print(connection)
21
22   qry = "select crm_cd, DR_No from crime_data"
23
24   df = pd.read_sql(qry, connection)
25
26   pd.set_option("display.max_rows", None)
27
28   print(df)
29
30   # plt.figure(figsize=(15, 10))
31
32   plt.figure(figsize=(10, 6))
33
34   sns.countplot(data=df, x="crm_cd", order=df["crm_cd"].value_counts().index[:499], palette="viridis", saturation=1)
35
36   plt.xlabel("Crime Code")
37
38   plt.ylabel("Number of Crimes")
39
40   plt.title("Distribution of Reported Crimes by Crime Code")
41
42   plt.xticks(rotation=45)
43
44   plt.tight_layout()
45
46   plt.savefig("question5.png")
47
48   plt.show()
```

# Insights

Based on our analysis, we recommend that law enforcement concentrate

their efforts on areas where crimes occur most frequently. We also suggest

allocating resources based on the ages of most victims and the types of

crimes common in different areas. This approach can help in better crime

prevention and control.