



## **Town Recommendation System and Report on an Individual Data Science Project**

Sandesh Sapkota

Softwarica College of IT and E-Commerce, Coventry University

ST5014CEM Data Science for Developers

Siddhartha Neupane

August 2024

## Table of Contents

Table of Contents .....	2
Table of Figures .....	5
Introduction .....	6
Cleaning Data.....	7
Cleaning Housing Data .....	7
Cleaning Broadband Speed.....	7
Cleaning Crime Data.....	8
Cleaning School Data .....	8
Exploratory Data Analysis .....	10
Visualization of Housing Price .....	10
Bar Chart: Average House Price by Town in 2023 .....	10
Bar Chart: Average House Price in 2023 .....	11
Line Chart: Average House Prices from 2020 to 2023 .....	12
Visualization of Broadband Speed Data .....	13
Box plot: Average Download Speed by County .....	13
Bar Chart: Average and Maximum Download Speeds in Cornwall .....	14
Bar Chart: Average and Maximum Download Speeds Bristol .....	15
Visualization of Crime Data .....	16
Radar Chart: Vehicle Crime Rate per 10,000 people .....	16

	3
Pie Chart: Robbery in 2023 by Month.....	17
Boxplot: Drug Offence Rate per 10,000 People in 2023 .....	18
Visualization of School Data .....	19
Line Graph: Average Attainment 8 scores for 2021-2022 in Cornwall.....	19
Line Graph: Average Attainment 8 scores for 2021-2022 in Bristol.....	20
Linear Modeling.....	22
Housing Price Vs. Average Download Speed .....	22
Average Attainment 8 Score Vs Housing Price.....	22
Housing Price Vs Drug Rates .....	23
Average Download Speed vs Attainment 8 scores .....	24
Town Recommendation System .....	25
Overview.....	25
Results.....	25
Based on Housing Price Ranking .....	25
Based on Crime Rates Ranking .....	26
Based on School Performance .....	27
Based on Broadband .....	27
Overall Ranking .....	28
Reflection.....	28
Legal and Ethical Considerations .....	29

Conclusion .....	30
References .....	31
Appendix .....	32

## Table of Figures

Figure 1: Barchart: Average House Price by Town.....	11
Figure 2: Bar Chart: Average House Price in 2023 .....	12
Figure 3: Line Chart of Average House Price .....	13
Figure 4: Boxplot: Average Download Speed by County .....	14
Figure 5: Barchart of Average and Max Download Speed in Cornwall.....	15
Figure 6: Barchart of Average and Max Downlaod Speed in Bristol.....	16
Figure 7: Radar Chart of Vehicle Crime Rate .....	17
Figure 8: Pie Chart: Robberies by Month in 2023 .....	18
Figure 9: Boxplot : Drug Offence per 10000.....	19
Figure 10: Linegraph of Attainment score in Cornwall.....	20
Figure 11: Linegraph of Attainment score in Bristol.....	21
Figure 12: House Price Vs Average Download Speed .....	22
Figure 13: Average Attainment Vs House price.....	23
Figure 14: House Price vs Drug Rate .....	24
Figure 15: Average Downlaod Speed Vs Attainment 8.....	24
Figure 16: Top 10 towns based on House Price .....	25
Figure 17: Top 10 Towns based on Crime Rates\.....	26
Figure 18:Top 10 based on Schools.....	27
Figure 19 Top 10 based on Broadband .....	27
Figure 20 Overall Top 10.....	28

## **Introduction**

This assignment shows the process and outcomes of using data analysis to recommend the appropriate city in the United Kingdom to purchase a property in the United Kingdom and particularly in Cornwall or Bristol counties. Each city has their own advantages but buying a house requires detailed analysis of several factors influencing property valuation and quality of life. The primary factors that are usually considered while buying a property include the housing price, crime rate in the neighborhood, internet connectivity, education and so on. This project utilizes the data made available by the government of United Kingdom and public sources to facilitate buyers to make informed decisions.

The objective of this project is to develop a recommendation system that uses data to analyze and compare the cities of Bristol and Cornwall based on several factors. The system provides a score to each city of the two counties based on the given criteria and rank the towns accordingly. The final goal is to suggest the top three cities to purchase a house to the friend in the context. This report outlines the process of obtaining the data, cleaning and preprocessing, performing exploratory data analysis and linear modeling to develop the recommendation system.

## **Cleaning Data**

Data cleaning is an important step in the data science lifecycle which comes right after obtaining the datasets. This step ensures that the datasets are accurate and consistent and prepares them appropriately to be used for analysis. The datasets obtained from the UK government and other public institutions also had several inconsistencies and inaccuracies. Each dataset were carefully cleaned and prepared to be used for the further analysis and development of the recommendation system.

### **Cleaning Housing Data**

The housing price data from 2020 and 2023 were combined into one using the `bind_rows()` function in R. Then the data was filtered to include the housing data of Cornwall and Bristol only. A new column was added named “Year” by extracting the year from the existing column “Transaction\_Date”. After that only the required columns were selected. Finally, the null values and redundant entries were removed using `na.omit()` and `distinct()` respectively. Now, the dataset was saved using the `write_csv()` function.

### **Cleaning Broadband Speed**

The broadband speed dataset was loaded using `read_csv()` function selecting the relevant columns. The the columns were renamed to simpler names for easier operation. For example, “Median download speed (Mbit/s)” was renamed to “MedianDownSpeed”. The final columns were ‘Postcode’, ‘MedianDownSpeed’, ‘MedianUpSpeed’, ‘AvgUpSpeed’, ‘MaxUpSpeed’, ‘AvgDownSpeed’ and ‘MaxDownSpeed’. Then the null values were removed using the `na.omit()` function. After that, the broadband data was merged with housing dataset using `inner_join()` function in R. The two datasets were joined based on the common field “Postcode”. Finally, the redundant rows were removed using the `distinct()` function. In this way, the

broadband speed dataset was cleaned and processed for further analysis. The cleaned dataset was then saved as a CSV using the `write_csv()` function.

### **Cleaning Crime Data**

The crime datasets from Bristol and Devon & Cornwall in the years 2020 to 2024 were loaded using the `read_csv()`. Then all the datasets were combined together using the `rbind()` functions and converted into a tibble. Since the crime dataset did not include postcodes but includes the LSOA codes, another dataset postcodes to LSOA was also loaded and cleaned to join with the crime dataset. Only the required columns were selected from the LSOA and crime data. From the combined crime dataset, month, LSOA code, crime type, and falls within columns were selected. The columns were renamed for simplicity in further processing. Similarly, from the LSOA to postcode dataset, only 'lsoa11cd', 'lsoa11nm', 'ladnm', 'pcds' were selected. The data frame was then filtered to include the counties Bristol and Cornwall only. The duplicate values for LSOA code were checked and removed from both crime and lsoa datasets. Finally, the selected crime dataset was merged with the lsoa dataset using left join by the column LSOA code common in both datasets. Then two new columns 'Year' and 'Month' were created using the `mutate` function from the original Month column by trimming from 1 to 4 and 6 to 7 respectively. Now, the population data was also merged using left join. Finally, `distinct()` and `na.omit()` functions were used to remove the redundant rows and null values to clean the final data. Now, the dataset was saved using the `write_csv()` function.

### **Cleaning School Data**

The performance school datasets from the academic years 2021 to 2022 and 2022 to 2023 were loaded using the `read_csv()` function from the `readR` library. Then the datasets were filtered to only include the relevant columns 'SCHNAME', 'PCODE', 'ATT8SCR', 'TOWN'. The



columns 'Year' and 'County' were added to each dataset using the `mutate()` function. Then the data from both academic years were combined for each county using the `rbind()` function. The combined dataset was now cleaned by filtering out the non-numeric values 'NE' and 'SUPP' in the ATT8SCR column. The null values and redundant rows were removed using the `na.omit()` and `distinct()` functions respectively. The final dataset was then saved as comma separated values (CSV) using the `write_csv()` function.

## **Exploratory Data Analysis**

Exploratory Data Analysis (EDA) is an approach in data analysis which uses graphical representation and data visualization to summarize the features of data. After cleaning and preprocessing the datasets, the trends were analyzed by visualizing the data. Several visualization techniques such as box plots, line graph, radar chart, bar chart and pie chart were used for better understanding of distribution, relationships and main features of the data.

### **Visualization of Housing Price**

#### **Bar Chart: Average House Price by Town in 2023**

The data was filtered to only include the year 2023 then the average housing price is calculated using the `summarize()` method. Then the data is visualized in a bar chart using the `ggplot2` library. The cities were plotted in the x-axis and the average price was plotted in the y-axis in the bar chart. The bar chart shows the average housing price in various towns in Bristol and Cornwall for the year 2023. Boscastle town is seen to have a significantly higher average housing price among the cities. Boscastle is followed by Port Isaac, Padstow, and Waderbridge as the most expensive cities to own a house. Camborne, Callington, Redruth are some of the least expensive cities to consider.

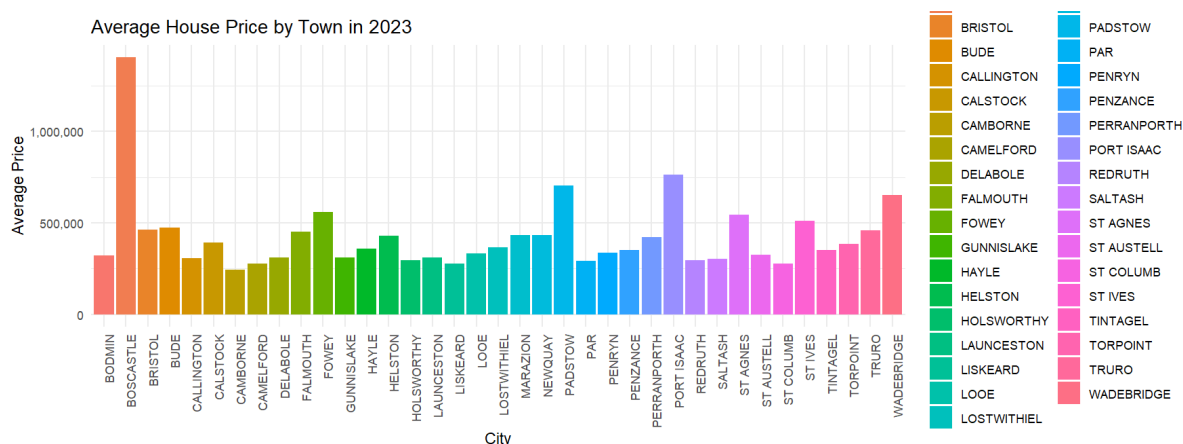
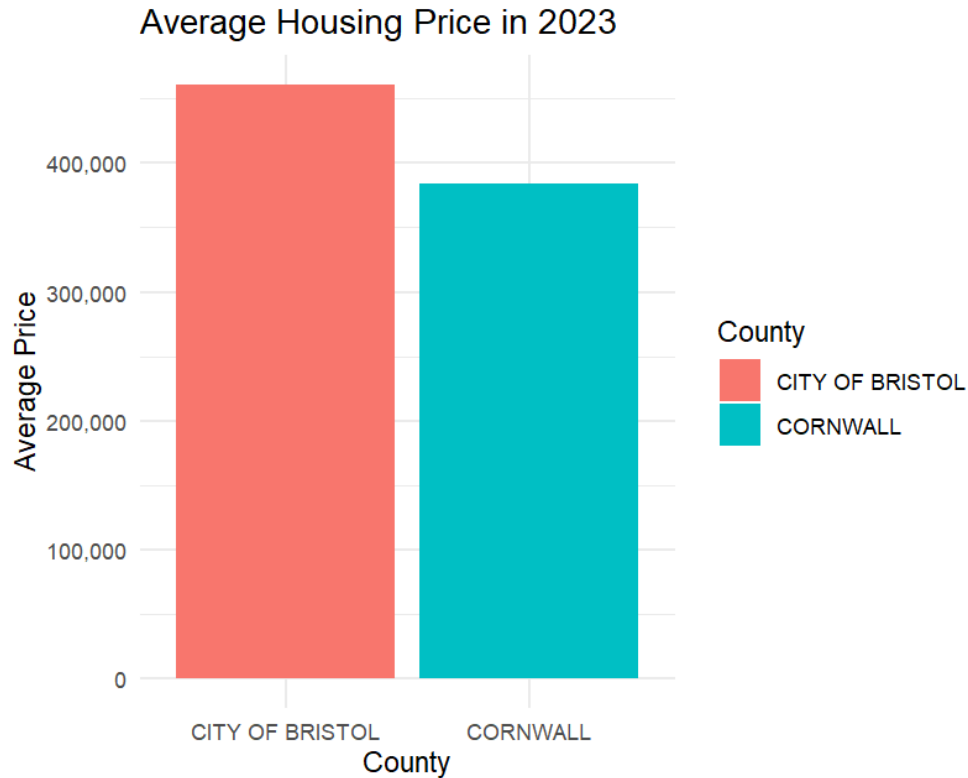


Figure 1: Barchart: Average House Price by Town

### Bar Chart: Average House Price in 2023

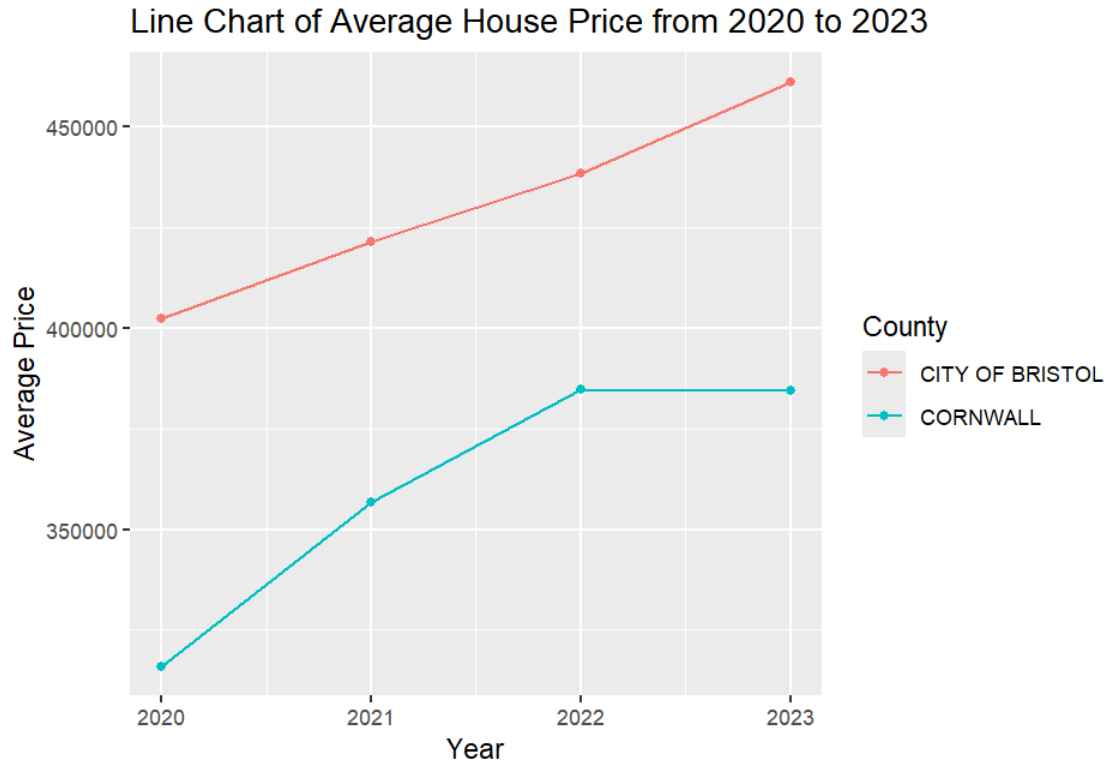
In this chart, the average house prices in 2023 were visualized by county wise in a bar chart. The data was first filtered to include only from the year 2023 and then grouped by County using `group_by()` function. Then the average house price was summarized by calculating the mean of 'Price column'. The data was then visualized in a bar chart using `ggplot2` to visualize the County to the x-axis and average house price to the y-axis. The chart shows that average house prices are noticeably higher in Bristol compared to Cornwall.



*Figure 2: Bar Chart: Average House Price in 2023*

### **Line Chart: Average House Prices from 2020 to 2023**

The data was filtered so that the year lies between 2020 to 2023 and grouped by Year and County. Then the average price was calculated using summarise() method. Then the line chart was plotted using the ggplot2 library with year on the x-axis and average price price on the y-axis. Looking at the line chart, the city of Bristol has consistently had higher prices compared to Cornwall throughout the year. After 2022, the housing market in Cornwall has seen more stability while the gap between the two counties has widened more in recent years.

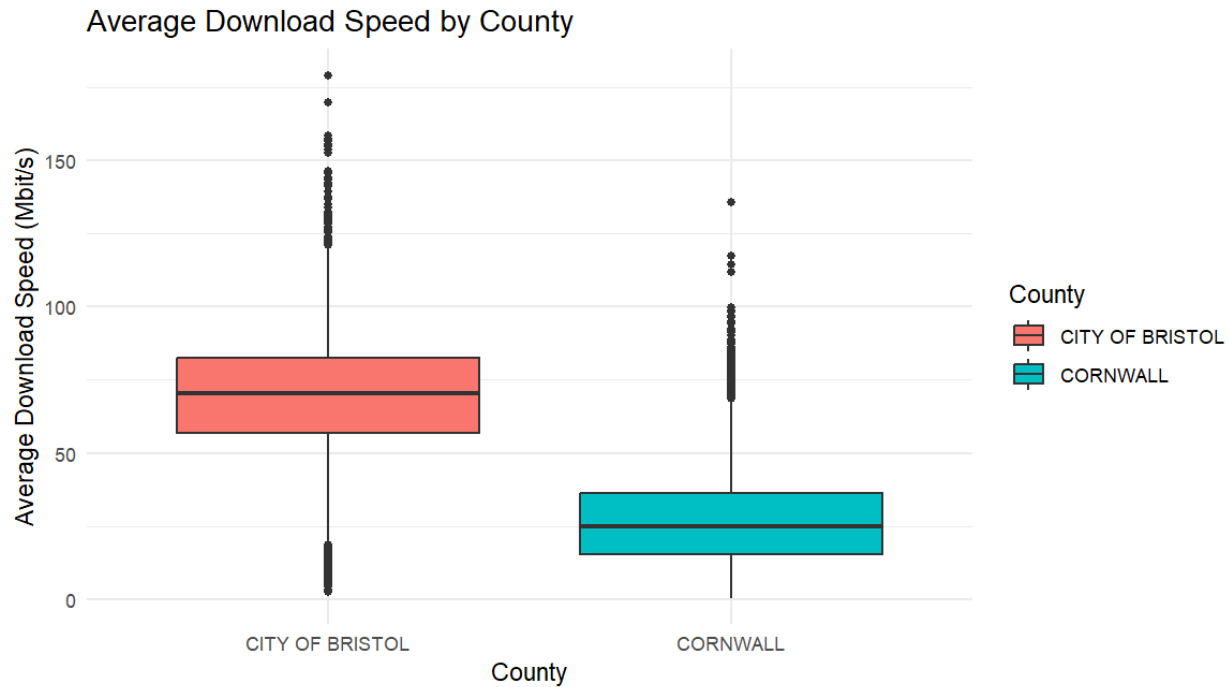


*Figure 3: Line Chart of Average House Price*

## Visualization of Broadband Speed Data

### Box plot: Average Download Speed by County

The cleaned broadband data is loaded from the csv file. Then the box plot is created using ggplot2 by grouping the average download speed and county. The box plot then visualizes the distribution of download speeds across the two counties. The box plot shows that Bristol has higher and more consistent download speeds with median around 80-90 Mbit/s while Cornwall has lower speed with median of 50 Mbits/s. The outliers denote the areas with significantly higher or lower speeds in both counties.



*Figure 4: Boxplot: Average Download Speed by County*

### **Bar Chart: Average and Maximum Download Speeds in Cornwall**

The data was filtered to only include the county Cornwall and grouped by Town/City. Then the average and maximum download speeds were calculated using the `summarize()` method. The bar chart was plotted using `ggplot2` which reveals significant differences in internet quality across Cornwall. The towns Saltash and St Austell have very high maximum speeds while the average among most towns is moderate.

### Average and Maximum Download Speeds by Town/City in Cornwall

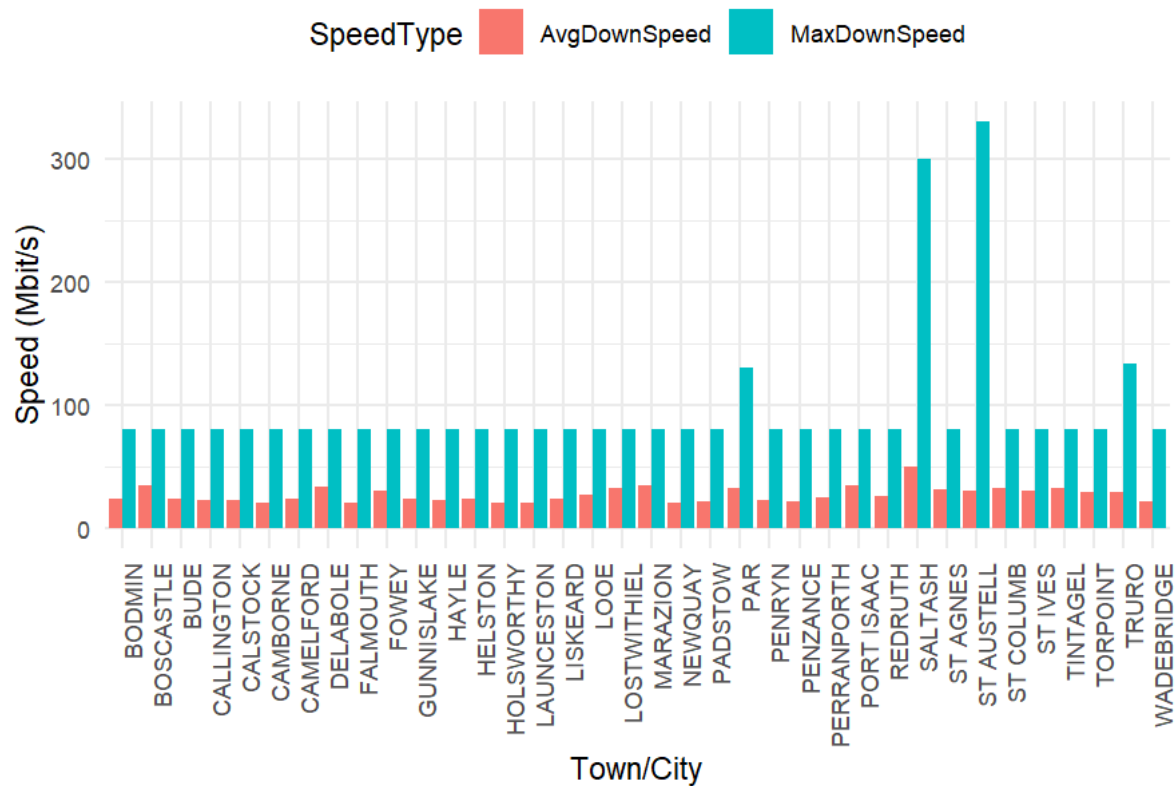
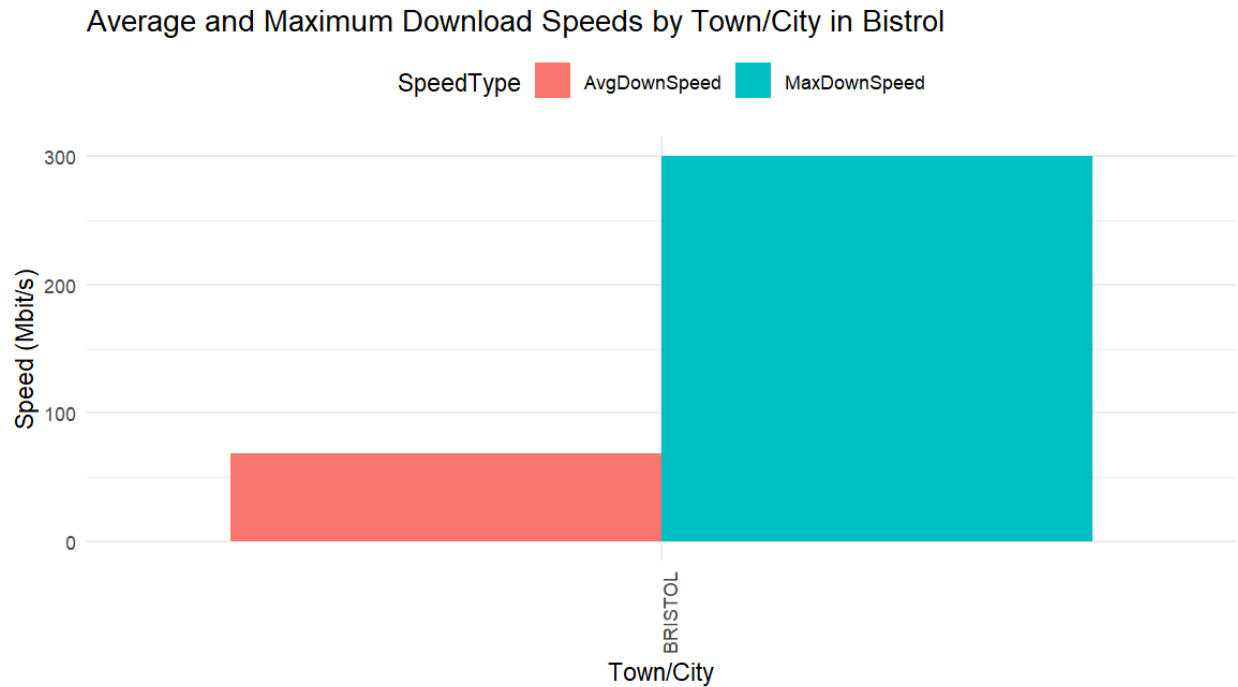


Figure 5: Barchart of Average and Max Download Speed in Cornwall

### Bar Chart: Average and Maximum Download Speeds Bristol

The data was filtered to include the data of Bristol only and grouped by Town/City. Then the average and maximum speeds were calculated using summarize() function before visualizing the bar chart using ggplot2. From the bar chart, the city's internet service is not uniformly distributed as there is a big gap between the maximum possible speed and the average speed most users get in the City of Bristol.



*Figure 6: Barchart of Average and Max Downlaod Speed in Bristol*

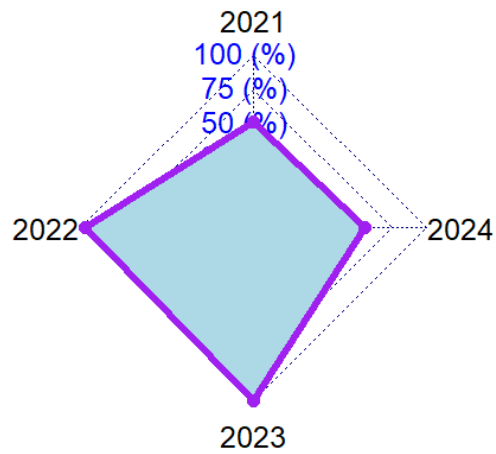
## Visualization of Crime Data

### Radar Chart: Vehicle Crime Rate per 10,000 people

The cleaned dataset is filtered to only include vehicle crimes and grouped by Year. Then the total crime is calculated by using the summarise() method. Then the data frame is prepared for the radar chart and plotted using the radarchart() function from 'fmsb' library in R. The chart shows the vehicle crimes were consistent from 2021 to 2024 with a few variations. The highest was in 2021 and the lowest was in 2023. Overall, the crime rates were stable across the years.



### Vehicle Crime Rate from 2021 to 2024



*Figure 7: Radar Chart of Vehicle Crime Rate*

### Pie Chart: Robbery in 2023 by Month

The data was filtered and summarized for robbery rates in each month in 2023. The data is then visualized in a pie chart using ggplot2. The chart shows that there was highest percentage of robberies at 13.9% in January, while April and September had the lowest at 5.1%. Other than that, the rate of robberies is even throughout the year.

## Robberies by Month in 2023

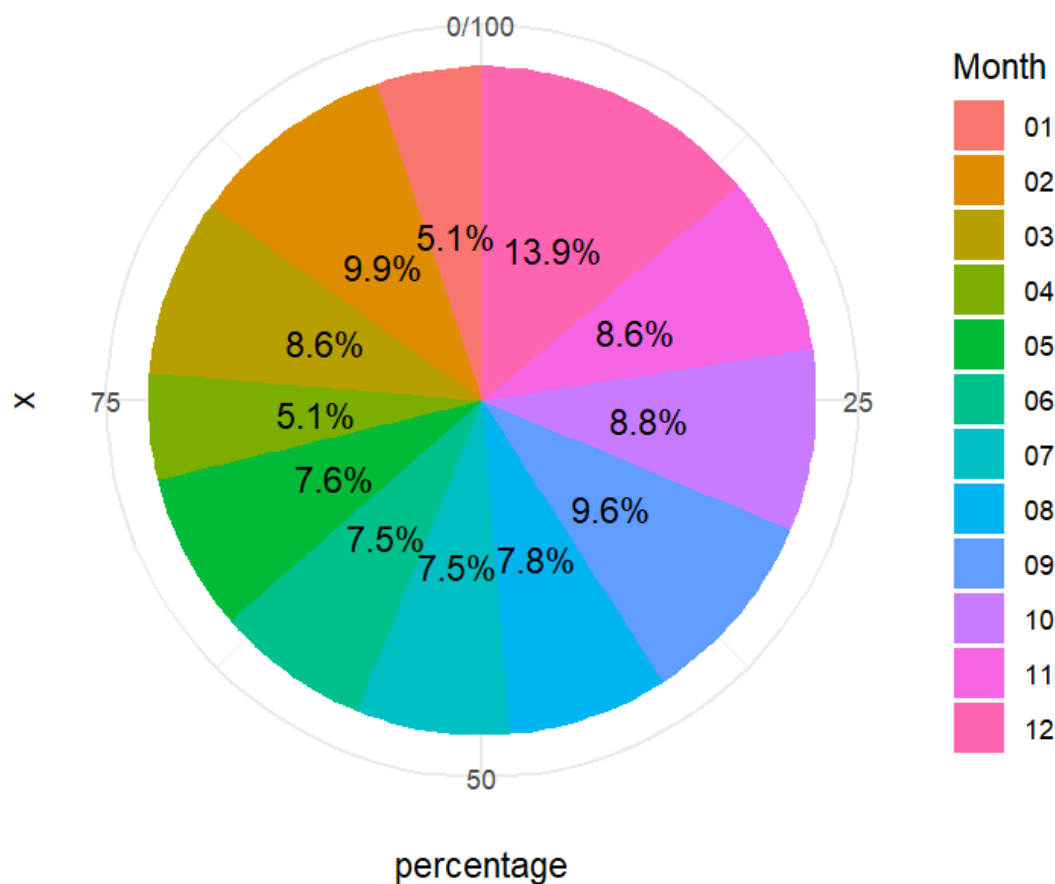


Figure 8: Pie Chart: Robberies by Month in 2023

## Boxplot: Drug Offence Rate per 10,000 People in 2023

The boxplot visualizes the distribution of drug related crimes in 2023 per 10,000 people in Bristol and Cornwall. The boxplot shows that the drug rates are higher in Cornwall compared to Bristol. Similarly, Cornwall has a wider spread indicating differences in offences in different regions while Bristol shows uniform distribution of data.

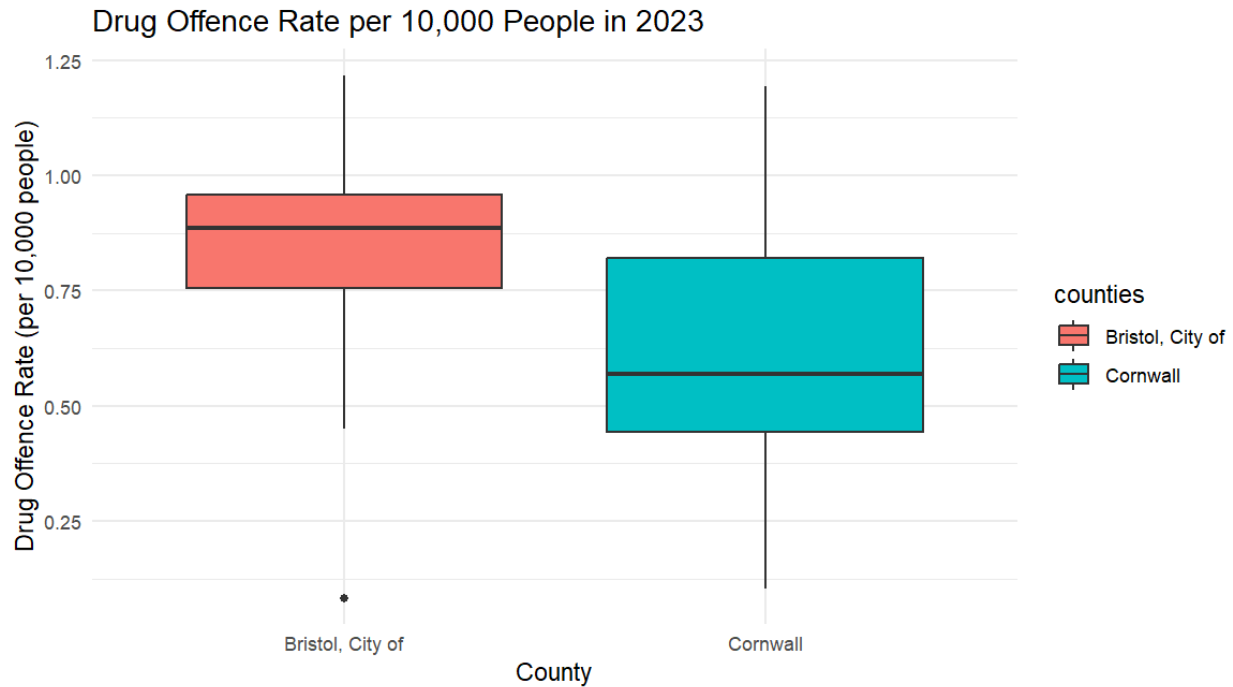
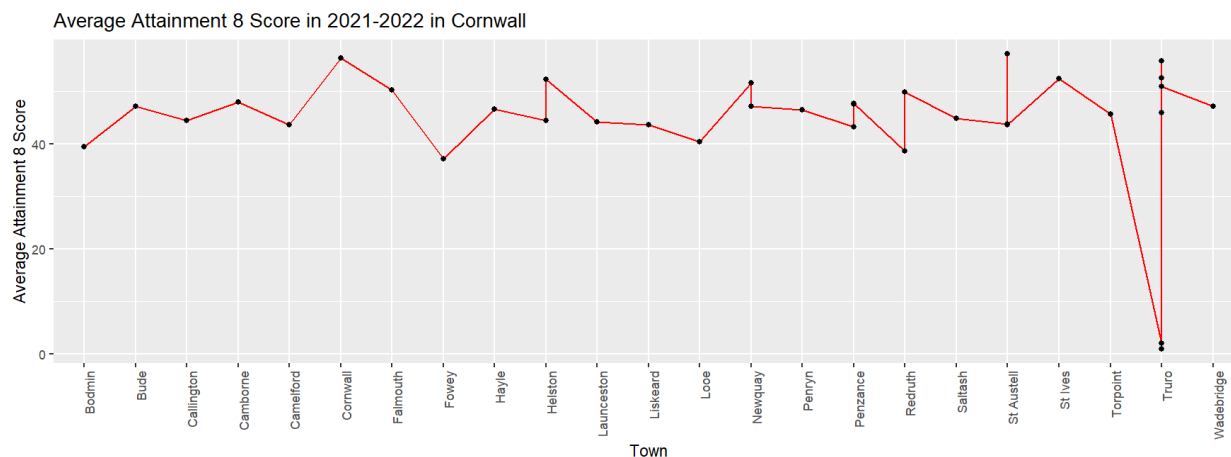


Figure 9: Boxplot : Drug Offence per 10000

## Visualization of School Data

### Line Graph: Average Attainment 8 scores for 2021-2022 in Cornwall

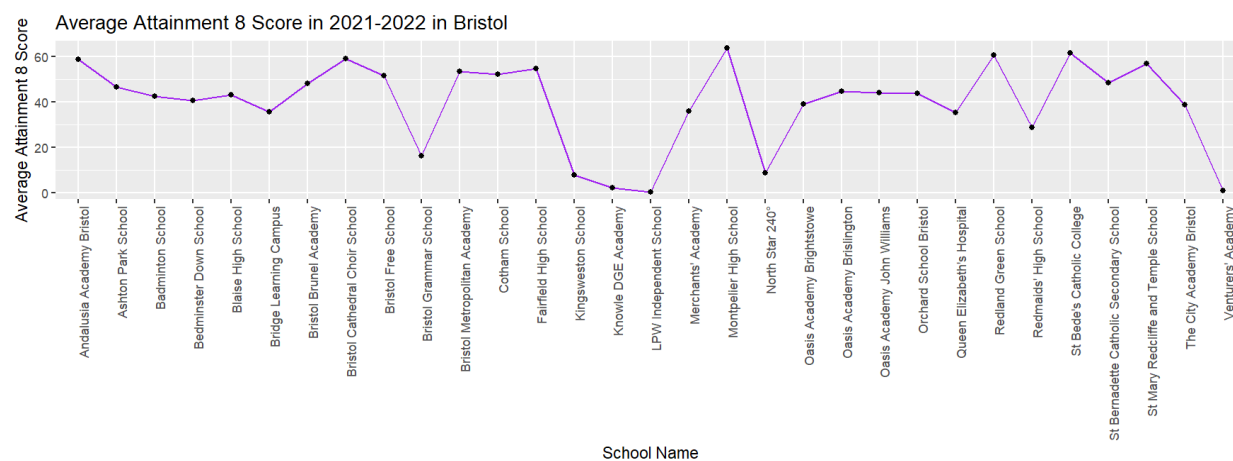
Attainment 8 scores is a way to measure the average performance of school students in the UK. The line graph is created to show the average attainment 8 scores for different towns in Cornwall in academic year 2021-2022. Town is taken in x-axis and scores is taken in y-axis of the graph. It shows a significant variation across the towns. Towns like Cornwall, St Ives performed well while others like Truro have performed extremely poor. The others had an above average performance.



*Figure 10: Linegraph of Attainment score in Cornwall*

### **Line Graph: Average Attainment 8 scores for 2021-2022 in Bristol**

For Bristol, the visualization for attainment 8 scores was done using ggplot2 based on Schools as there weren't significant data for the towns. In the line graph, x-axis represented Schools and y-axis contained their attainment scores. The graph shows schools like Montpelier High School peak the graph while schools like Venturers Academy show lower scores. This shows high academic disparities across the Bristol County.



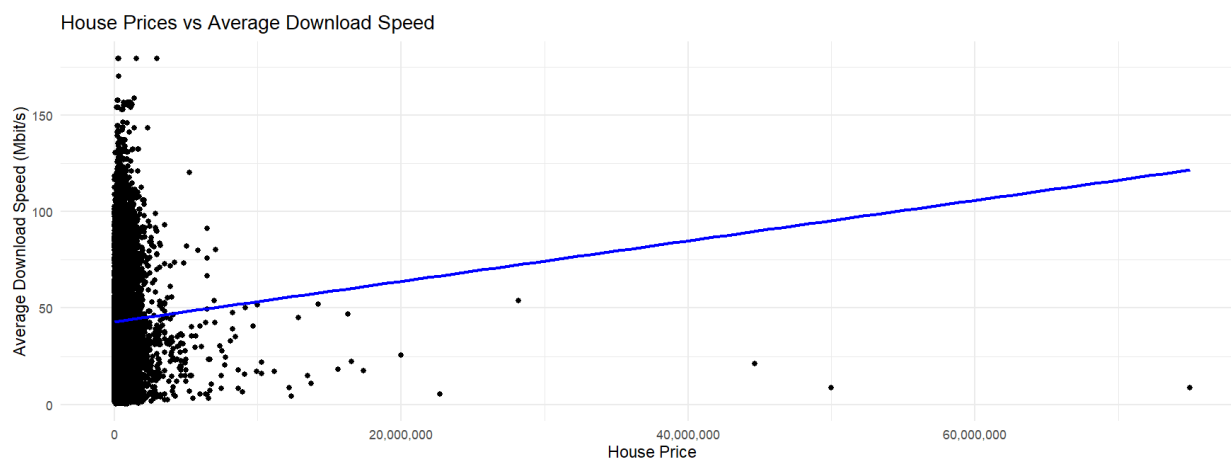
*Figure 11: Linegraph of Attainment score in Bristol*

## Linear Modeling

Linear Modeling is a statistical method used in data analysis to analyze the relationship between two or more variables. It shows the linear relationship between a dependent variable and one or more independent variables expressed in the form of a linear equation. This approach was used to identify the trends, patterns and make predictions based on the data of Cornwall and Bristol.

### Housing Price Vs. Average Download Speed

The shows a positive correlation between house price and average download speed. This suggests that higher the internet bandwidth, higher will be the housing price. This is likely because people with nicer house often want a good internet connection. However, there is a lot of outliers especially at the lower end of the house prices.

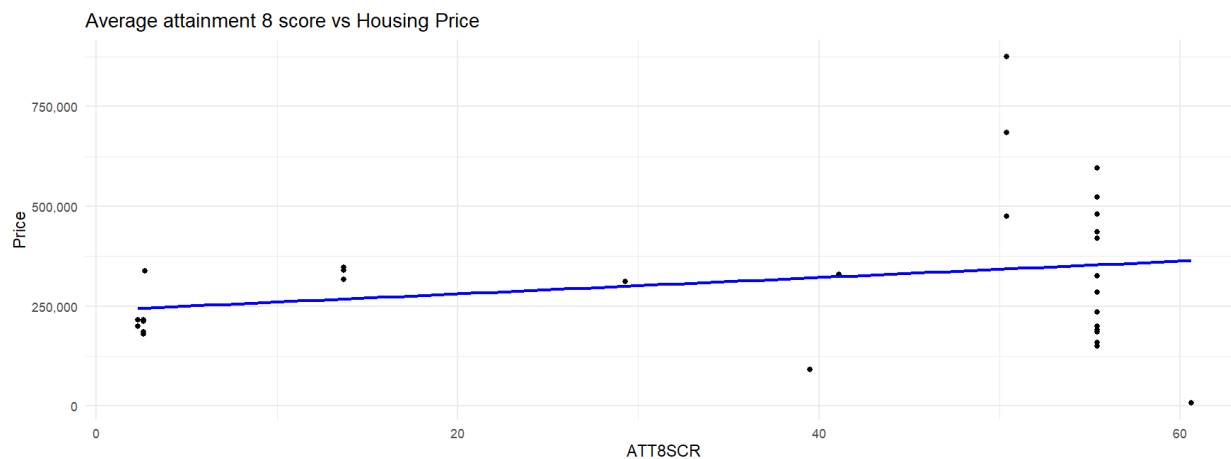


*Figure 12: House Price Vs Average Download Speed*

### Average Attainment 8 Score Vs Housing Price

The scatterplot shows a weak positive relationship between the attainment 8 score and housing prices in the two counties. This shows that the areas with good educational performance

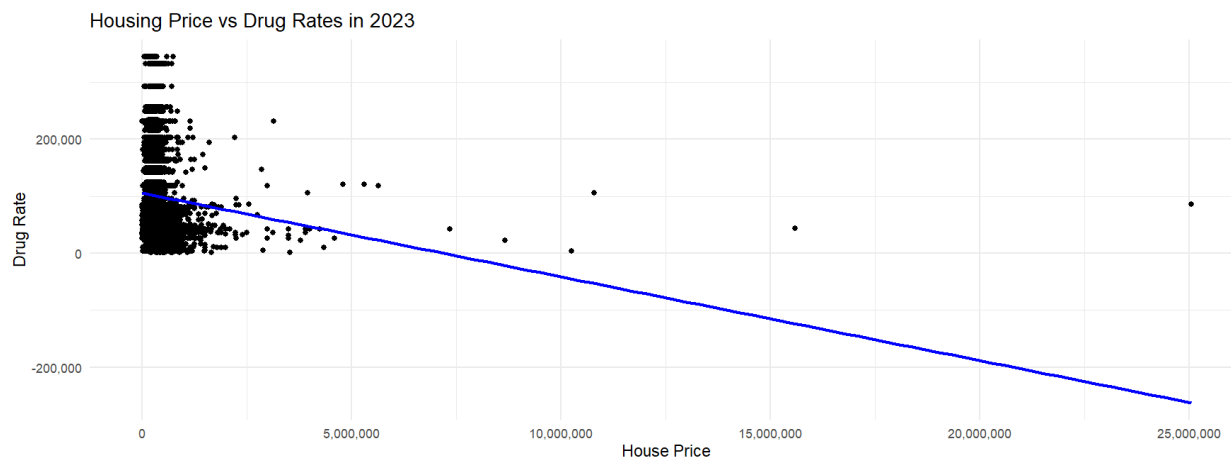
have higher housing price. The relationship is weak because the data points are widely scattered indicating that attainment 8 score might not have significant impact on the price.



*Figure 13: Average Attainment Vs House price*

### **Housing Price Vs Drug Rates**

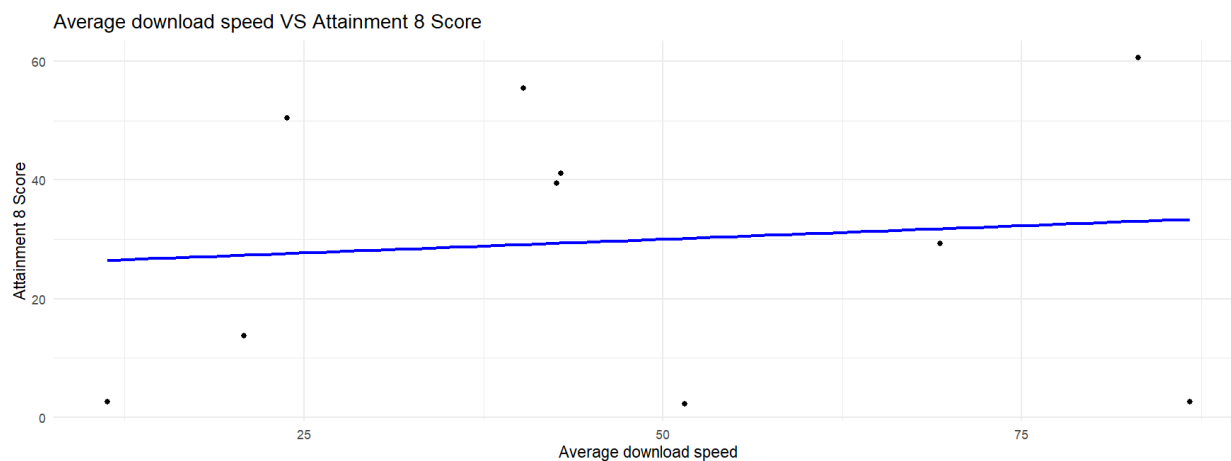
The plot shows a negative correlation between the housing prices and drug crime rates in the area. This suggests that the areas with expensive houses have lower drug offence rate while cheaper neighborhoods have higher drug crime rates. This might be because social factors like this affect the real estate market.



*Figure 14: House Price vs Drug Rate*

### **Average Download Speed vs Attainment 8 scores**

A weak positive correlation is seen between the average download speed and attainment 8 scores of a region. Even though the scattered data points suggest the relationship was not too strong, it could be said that the better internet has a slight impact in the academic performance of schools.



*Figure 15: Average Downlaod Speed Vs Attainment 8*



## Town Recommendation System

### Overview

After cleaning the obtained datasets, the data were prepared and visualized for exploratory data analysis, prepared linear models to compare the relationships between the different parameters in deciding the house purchase. Now, the task was to identify the most suitable town to purchase the house considering affordable housing, low crime rates, good internet connectivity, and good school performance measured by the attainment 8 scores. The preprocessed datasets were loaded and filtered and processed to calculate average for each factors by town. The data was then merged into a single data frame and normalized to ensure they can be compared later on. Finally, the towns are ranked based on the scores. Higher scores indicate more favorable to purchase a house or property in that location.

### Results

#### Based on Housing Price Ranking

	TOWN	avgPrice	normHousingPrice
1	CAMBORNE	244939.2	1.0000000
2	LISKEARD	276915.0	0.9214842
3	CAMELFORD	277289.9	0.9205635
4	REDRUTH	295858.8	0.8749681
5	SALTASH	301707.2	0.8606076
6	CALLINGTON	305634.0	0.8509653
7	LAUNCESTON	310922.9	0.8379786
8	BODMIN	323267.0	0.8076681
9	ST AUSTELL	325660.4	0.8017912
10	LOOE	333198.0	0.7832828

Figure 16: Top 10 towns based on House Price

### Based on Crime Rates Ranking

	TOWN	normCrimeRate	crimerate
1	FOWEY	1.0000000	98
2	CAMELFORD	0.9913297	199
3	LOOE	0.9800841	330
4	CALLINGTON	0.9683235	467
5	WADEBRIDGE	0.9663490	490
6	TORPOINT	0.9634303	524
7	PENRYN	0.9592240	573
8	ST IVES	0.9545884	627
9	HAYLE	0.9478067	706
10	BUDE	0.9420551	773

*Figure 17: Top 10 Towns based on Crime Rates\*

### Based on School Performance

	TOWN	avgAtt8	normAtt8
1	ST IVES	49.90000	1.000000000
2	PENRYN	47.75000	0.816112567
3	HELSTON	47.47500	0.792592082
4	ST AUSTELL	47.26667	0.774773532
5	WADEBRIDGE	46.60000	0.717754173
6	BUDE	46.45000	0.704924817
7	FALMOUTH	46.45000	0.704924817
8	CAMBORNE	46.25000	0.687819010
9	LAUNCESTON	44.65000	0.550972548
10	PENZANCE	44.31667	0.522462868

Figure 18:Top 10 based on Schools

### Based on Broadband

	TOWN	avg_down_speed	normDownSpeed
1	BRISTOL	68.48118	1.000000000
2	SALTASH	49.86904	0.609798505
3	FOWEY	30.87476	0.211585278
4	ST IVES	30.63361	0.206529673
5	ST AUSTELL	30.58067	0.205419759
6	TORPOINT	29.01877	0.172674778
7	LOOE	27.68960	0.144808782
8	REDRUTH	26.06246	0.110696055
9	LISKEARD	24.47380	0.077389865
10	HELSTON	24.44340	0.076752519

Figure 19 Top 10 based on Broadband

## Overall Ranking

	TOWN	avgPrice	crimerate	avgAtt8	avg_down_speed	finalPoints
1	SALTASH	301707.2	836	42.65000	49.86904	3.786967
2	ST AUSTELL	325660.4	2744	47.26667	30.58067	3.531017
3	ST IVES	510935.4	627	49.90000	30.63361	3.463236
4	TORPOINT	385197.4	524	43.60000	29.01877	3.093698
5	LOOE	333198.0	330	40.80000	27.68960	2.843493
6	HELSTON	431010.9	1361	47.47500	24.44340	2.830370
7	REDRUTH	295858.8	2542	42.87500	26.06246	2.783760
8	PENRYN	335319.3	573	47.75000	23.28255	2.728404
9	BUDE	475156.9	773	46.45000	24.13062	2.684772
10	LISKEARD	276915.0	1110	42.25000	24.47380	2.675947

*Figure 20 Overall Top 10*

Ergo, after analyzing the towns on different parameters, the towns SALTASH, ST AUSTELL and ST IVES can be recommended for purchasing a house since these cities provide the most balanced parameters optimal housing price, less crime rates, good internet and good education.

## Reflection

The relevant datasets were obtained from reliable sources of the official websites UK government and other public organizations to maintain the credibility of the data. Then the analysis was done using the R programming language to clean, process, visualize, construct linear models and finally rank the towns based on different parameters. The data cleaning included filtering required columns, removing null values and redundant fields. Then the data were visualized using different visualization methods such as box plots, bar charts, line graphs

and pie charts to find the patterns and trends. Similarly, linear modeling was done to analyse and compare the relationships between different parameters. This helped to classify the counties based on their performance in various factors. Finally, all the parameters were analyzed to identify the most favorable characteristics and recommendations were made on the basis of different comparative results.

### **Legal and Ethical Considerations**

The project has been done by considering all the ethical and legal compliances practiced in the industry. The datasets were obtained from the official websites of the UK government that are publicly available and approved by the government. The data weren't manipulated to push a certain political ideology and are completely based on the UK government approved sources.

## **Conclusion**

This project demonstrated the use of data analysis to develop a recommendation system to figure out the optimal town in the Bristol and Cornwall counties in the United Kingdom to purchase a house. By cleaning and preprocessing the datasets obtained from credible sources, the parameters were compared to identify Saltash, St Austell and St Ives as the top three most favorable cities to buy a house considering affordable pricing, low crime rates, good internet and strong academic performance of local schools. It was a great learning experience to explore the capacities of data analysis in recommendation systems.

## References

Sisense. (2018, June 12). What is Data Cleaning? | Sisense.

<https://www.sisense.com/glossary/data-cleaning/>

Stedman, C. (2022, January 28). data cleansing (data cleaning, data scrubbing). Data Management. <https://www.techtarget.com/searchdatamanagement/definition/data-scrubbing>

Data downloads | data.police.uk. (n.d.). <https://data.police.uk/data/>

Connected Nations 2018: Data downloads. (2019, January 2). www.ofcom.org.uk. <https://www.ofcom.org.uk/phones-and-broadband/coverage-and-speeds/data-downloads>

*RPUBS - Introduction to Linear Modeling in R.* (n.d.).

<https://rpubs.com/AnthonyCorbisieri/1100671>

*Radar.* (n.d.). <https://plotly.com/r/radar-chart/>

Chatterjee, S. (2024, February 26). *5 Ethical aspects for data science professionals to consider.*

Emeritus Online Courses. <https://emeritus.org/blog/data-science-and-analytics-data-science-course-curriculum/>

## Appendix

**Github Link:** <https://github.com/sandeshs0/data-science-assignment>

**Google Drive Link:**

[https://drive.google.com/drive/folders/11tJCS5Z5XRncfJ\\_cnh9a0\\_XBJJ31dKoX?usp=sharing](https://drive.google.com/drive/folders/11tJCS5Z5XRncfJ_cnh9a0_XBJJ31dKoX?usp=sharing)

### Cleaning.R

```

1 # Libraries
2 library(dplyr)
3 library(readr)
4 library(tidyverse)
5 # Defining column names
6 column_names <- c("Transaction_ID", "Price", "Transaction_Date", "Postcode",
7   "Property_Type", "old_New", "Duration", "PAON", "SAON",
8   "Street", "Locality", "Town_City", "District", "County",
9   "PPD_Category_Type", "Record_Status")
10 # Loading the datasets with column names
11 housing_2020 <- read_csv("D:/Academics/Fourth Semester/Data Science/Assignment/Obtained Datasets/Housing/pp-2020.csv", col_names = column_names)
12 housing_2021 <- read_csv("D:/Academics/Fourth Semester/Data Science/Assignment/Obtained Datasets/Housing/pp-2021.csv", col_names = column_names)
13 housing_2022 <- read_csv("D:/Academics/Fourth Semester/Data Science/Assignment/Obtained Datasets/Housing/pp-2022.csv", col_names = column_names)
14 housing_2023 <- read_csv("D:/Academics/Fourth Semester/Data Science/Assignment/Obtained Datasets/Housing/pp-2023.csv", col_names = column_names)
15 # Combining the datasets
16 combined_housing = bind_rows(housing_2020, housing_2021, housing_2022, housing_2023)
17 #Data Cleaning
18 cleaned_housing <- combined_housing %>%
19   filter(County %in% c('CORNWALL', 'CITY OF BRISTOL')) %>%
20   mutate(Year = year(Transaction_Date)) %>%
21   select(Price, Postcode, Year, Town_City, County) %>%
22   na.omit() %>%
23   distinct()
24 write_csv(cleaned_housing, 'D:/Academics/Fourth Semester/Data Science/Assignment/Cleaned Datasets/housing_cleaned.csv')
25
26 #-----CLEANING BROADBAND SPEED DATASET-----
27 broadband <- read_csv("D:/Academics/Fourth Semester/Data Science/Assignment/Obtained Datasets/Broadband Speed/201805_fixed_pc_performance_r0")
28 #Selecting Relevant Columns
29 broadband_selected <- broadband %>%
30   select("postcode_space",
31     "Median download speed (Mbit/s)",
32     "Median upload speed (Mbit/s)",
33     "Average upload speed (Mbit/s)",
34     "Maximum upload speed (Mbit/s)",
35     "Average download speed (Mbit/s)",
36     "Maximum download speed (Mbit/s)",
37   )
38 #Checking Null Values
39 na_summary <- sapply(broadband_selected, function(x) sum(is.na(x)))
40 print(na_summary)
41 #Cleaning
42 broadband_clean <- broadband_selected %>%
43   rename(
44     Postcode = postcode_space,
45     MedianDownSpeed = "Median download speed (Mbit/s)",
46     MedianUpSpeed = "Median upload speed (Mbit/s)",
47     AvgUpSpeed = "Average upload speed (Mbit/s)",
48     MaxUpSpeed = "Maximum upload speed (Mbit/s)",
49     AvgDownSpeed = "Average download speed (Mbit/s)",
50     MaxDownSpeed = "Maximum download speed (Mbit/s)"
51   ) %>%
52   na.omit() %>%
53   distinct()

```



```

53 distinct()
54 #Checking for null values after cleaning
55 na_summary_clean <- apply(broadband_clean, function(x) sum(is.na(x)))
56 print(na_summary_clean)
57 #Selecting only 3 rows from housing
58 housing_selected <- cleaned_housing %>%
59   select(Postcode,Town_City,County)
60
61 #Using Inner Join to Merge the broadband with housing data
62 broadband_w_housing = inner_join(housing_selected,broadband_clean,by="Postcode");
63 View(broadband_w_housing)
64
65 na_merged <- apply(broadband_w_housing, function(x) sum(is.na(x)))
66 print(na_merged)
67 dim(broadband_w_housing)
68
69
70 #Removing redundant rows
71 broadband_final <- broadband_w_housing %>%
72   distinct()
73 dim(broadband_final)
74 View(broadband_final)
75
76 #Saving the dataset
77 write_csv(broadband_final, "D:/Academics/Fourth Semester/Data Science/Assignment/Cleaned Datasets/broadband_cleaned.csv")
78

```

```

159 b1st31=read_csv("D:/Academics/Fourth Semester/Data Science/Assignment/obtained Datasets/Crime/2024-01/2024-01-avon-and-somerset-street.csv")
160 b1st32=read_csv("D:/Academics/Fourth Semester/Data Science/Assignment/obtained Datasets/Crime/2024-02/2024-02-avon-and-somerset-street.csv")
161 b1st33=read_csv("D:/Academics/Fourth Semester/Data Science/Assignment/obtained Datasets/Crime/2024-03/2024-03-avon-and-somerset-street.csv")
162 b1st34=read_csv("D:/Academics/Fourth Semester/Data Science/Assignment/obtained Datasets/Crime/2024-04/2024-04-avon-and-somerset-street.csv")
163 b1st35=read_csv("D:/Academics/Fourth Semester/Data Science/Assignment/obtained Datasets/Crime/2024-05/2024-05-avon-and-somerset-street.csv")
164 b1st36=read_csv("D:/Academics/Fourth Semester/Data Science/Assignment/obtained Datasets/Crime/2024-06/2024-06-avon-and-somerset-street.csv")
165
166 #Combining the Datasets
167 crime_combined=rbind(
168   b1st1,b1st2,b1st3,b1st4,b1st5,b1st6,b1st7,b1st8,b1st9,b1st10,b1st11,b1st12,b1st13,b1st14,b1st15,b1st16,b1st17,b1st18,b1st19,b1st20,b1st21,
169   b1st22,b1st23,b1st24,b1st25,b1st26,b1st27,b1st28,b1st29,b1st30,b1st31,b1st32,b1st33,b1st34,b1st35,b1st36,
170
171   corn1,corn2,corn3,corn4,corn5,corn6,corn7,corn8,corn9,corn10,corn11,corn12,corn13,corn14,corn15,corn16,corn17,corn18,corn19,corn20,corn21,
172   corn23,corn24,corn25,corn26,corn27,corn28,corn29,corn30,corn31,corn32,corn33,corn34,corn35,corn36
173 )
174
175 head(crime_combined)
176 View(crime_combined)
177 dim(crime_combined)
178
179 #Converting it to a tibble
180 crime_combined <- crime_combined %>%
181   as_tibble()
182
183
184 #POST CODE TO LSOA
185
186 postcode_to_lsoa <- read_csv("D:/Academics/Fourth Semester/Data Science/Assignment/Obtained Datasets/Postcode to LSOA.csv")
187 population <- read_csv("D:/Academics/Fourth Semester/Data Science/Assignment/Obtained Datasets/Population.csv")
188 #Selecting relevant columns
189 selected_lsoa <- postcode_to_lsoa %>%
190

```

```

186 postcode_to_lsoa <- read_csv("D:/Academics/Fourth Semester/Data Science/Assignment/Obtained Datasets/Postcode to LSOA.csv")
187 population <- read_csv("D:/Academics/Fourth Semester/Data Science/Assignment/Obtained Datasets/Population.csv")
188 #Selecting relevant columns
189 selected_lsoa <- postcode_to_lsoa %>%
190   select('lsoalld', 'lsoallnm', 'ladnm', 'pcds')
191
192 selected_crime <- crime_combined %>%
193   select('Month', 'LSOA code', 'Crime type', 'Falls within')
194
195 #Renaming the columns
196 colnames(selected_lsoa) = c('LSOA code', 'street', 'counties', "postcode")
197 colnames(population) = c("postcode", "population")
198 View(selected_crime)
199 View(selected_lsoa)
200 #Cleaning LSOA
201 clean_lsoa <- selected_lsoa %>%
202   filter(counties %in% c("Bristol", "City of", "Cornwall")) %>%
203   mutate(postcode=str_trim(substring(postcode,1,6)))
204
205 #Checking for duplicates of lsoa codes
206 any(duplicated(selected_crime$`LSOA code`))
207 any(duplicated(clean_lsoa$`LSOA code`))
208
209 #Since the duplicates are there, removing them
210 clean_lsoa=unique(clean_lsoa,by="LSOA code")
211 selected_crime=unique(selected_crime,by="LSOA code")
212
213 #Final Cleaning and merging
214 finalcrime= selected_crime %>%
215   left_join(clean_lsoa, by=("LSOA code"),relationship = "many-to-many") %>%
216   mutate(Year=str_trim(substring(Month,1,4))) %>%
217

```

```

213 #Final Cleaning and merging
214 finalCrime= selected_crime %>%
215   left_join(clean_lsoa, by=("LSOA code"),relationship = "many-to-many") %>%
216   mutate(Year=str_trim(substring(Month,1,4))) %>%
217   mutate(Month=str_trim(substring(Month,6,7))) %>%
218   left_join(population,by="postcode") %>%
219   distinct() %>%
220   na.omit()
221
222 dim(finalCrime)
223 dim(clean_lsoa)
224 view(finalCrime)
225 dim
226 write_csv(finalCrime, "D:/Academics/Fourth Semester/Data Science/Assignment/Cleaned Datasets/crime_cleaned.csv")
227
228 #-----SCHOOL-----
229 bschool121=read_csv("D:/Academics/Fourth Semester/Data Science/Assignment/Obtained Datasets/School/Bristol/2021-2022/801_ks4final.csv")
230 bschool122=read_csv("D:/Academics/Fourth Semester/Data Science/Assignment/Obtained Datasets/School/Bristol/2022-2023/801_ks4final.csv")
231 cschool121=read_csv("D:/Academics/Fourth Semester/Data Science/Assignment/Obtained Datasets/School/Cornwall/2021-2022/908_ks4final.csv")
232 cschool122=read_csv("D:/Academics/Fourth Semester/Data Science/Assignment/Obtained Datasets/School/Cornwall/2022-2023/908_ks4final.csv")
233
234 view(bschool121)
235 head(bschool121)
236 view(cschool121)
237
238 bschool121 <- bschool121 %>%
239   select(SCHNAME,PCODE,ATT8SCR,TOWN) %>%
240   mutate(YEAR=2021,COUNTY="Bristol")
241
242 bschool122 <- bschool122 %>%
243   select(SCHNAME,PCODE,ATT8SCR,TOWN) %>%
244   mutate(YEAR=2022,COUNTY="Bristol")
245
246 cschool121 <- cschool121 %>%
247   select(SCHNAME,PCODE,ATT8SCR,TOWN) %>%
248   mutate(YEAR=2021,COUNTY="Cornwall")
249
250 cschool122 <- cschool122 %>%
251   select(SCHNAME,PCODE,ATT8SCR,TOWN) %>%
252   mutate(YEAR=2022,COUNTY="Cornwall")
253
254 #Combining
255 bCombined=rbind(bschool121,bschool122)
256 cCombined=rbind(cschool121,cschool122)
257
258 View(bCombined)
259 View(cCombined)
260
261 schoolCombined=rbind(bCombined,cCombined)
262 View(schoolCombined)
263
264 schoolCombined <- schoolCombined %>%
265   filter(ATT8SCR!="NE"&ATT8SCR!="SUPP") %>%
266   na.omit() %>%
267   distinct()
268
269 dim(schoolCombined)
270 View(schoolCombined)
271
272 write_csv(schoolCombined,"D:/Academics/Fourth Semester/Data Science/Assignment/Cleaned Datasets/schoolcleaned.csv")
273

```

## Graphs.R

```

1 library(tidyverse)
2 library(ggplot2)
3 library(fmsb)
4 library(scales)
5
6 #Housing
7 housing_data=read_csv("D:/Academics/Fourth Semester/Data Science/Assignment/Cleaned Datasets/housing_cleaned.csv")
8
9 #For Average House Price of 2023 for Bristol and Cornwall
10
11 housing_2023<- housing_data %>%
12   filter(Year==2023) %>%
13   group_by(County) %>%
14   summarise(average_house_price_2023=mean(Price))
15
16 colnames(housing_data)
17
18 #Average house price for each county for 2020-2023
19 housing_20_23=housing_data %>%
20   filter(Year>=2020) %>%
21   group_by(Year,County) %>%
22   summarise(avg_house_price=mean(Price))
23
24 head(housing_2023)
25 head(housing_20_23)
26
27 #BarChart
28 ggplot(housing_2023,aes(x=County,y=average_house_price_2023,fill=County))+
29   geom_bar(stat="identity")+
30   ggtitle("Average Housing Price in 2023")+
31   ylab("Average Price")+
32   xlab("County")+
33   scale_y_continuous(labels = scales::comma)

```

```

27 #BarChart
28 ggplot(housing_2023,aes(x=County,y=average_house_price_2023,fill=County))+
29   geom_bar(stat="identity")+
30   ggtitle("Average Housing Price in 2023")+
31   ylab("Average Price")+
32   xlab("County")+
33   scale_y_continuous(labels = scales::comma) +
34   theme_minimal()
35
36 #Box Plot
37 ggplot(data = housing_2023, aes(x = County, y = average_house_price_2023, fill = County)) +
38   geom_boxplot() +
39   labs(title = "Boxplot for Average House Price in Year 2023", x = "County", y = "Price") +
40   scale_y_continuous(limits = c(0, 300000), labels = scales::comma) +
41   theme_minimal()
42
43 #Line Chart
44 housing_20to23 <- housing_data %>%
45   filter(Year >= 2020 & Year <= 2023)
46 average_per_year <- housing_20to23 %>%
47   group_by(Year, County) %>%
48   summarise(Average_Price = mean(Price))
49
50
51 ggplot(average_price_per_year, aes(x = Year, y = Average_Price, color = County)) +
52   geom_line() +
53   geom_point() +
54   ggtitle("Line chart of Average House Price from 2020 to 2023") +
55   xlab("Year") +
56   ylab("Average Price")+
57   scale_y_continuous(labels = scales::comma) +
58   scale_x_continuous(labels = scales::comma)
59

```

```

59
60 #housing by years
61 housing_by_year = housing_data %>%
62   filter(Year >= 2020) %>%
63   group_by(Year, Town_City) %>%
64   summarise(avg_price = mean(Price))
65
66 # Filtering data for 2023
67 housing_data_2023 = housing_by_year %>%
68   filter(Year == 2023)
69
70 ggplot(housing_data_2023, aes(x = Town_City, y = avg_price, fill = Town_City)) +
71   geom_bar(stat = "identity") +
72   ggtitle("Average House Price by Town in 2023") +
73   ylab("Average Price") +
74   xlab("City") +
75   scale_y_continuous(labels = scales::comma) +
76   theme_minimal() +
77   theme(axis.text.x = element_text(angle = 90, hjust = 1))
78
79
80 #-----Broadband Data
81 broadband=read_csv("D:/Academics/Fourth Semester/Data Science/Assignment/Cleaned Datasets/broadband_cleaned.csv")
82
83 colnames(broadband)
84 View(broadband)
85 #Average download speed
86 ggplot(broadband, aes(x = County, y = AvgDownSpeed, fill = County)) +
87   geom_boxplot() +
88   labs(title = "Average Download Speed by County",
89        x = "County",
90        y = "Average Download Speed (Mbit/s)") +
91   theme_minimal()
92
48:41 (Top Level)
R Script

```

```

86 ggplot(broadband, aes(x = County, y = AvgDownSpeed, fill = County)) +
87   geom_boxplot() +
88   labs(title = "Average Download Speed by County",
89        x = "County",
90        y = "Average Download Speed (Mbit/s)") +
91   theme_minimal()
92
93 #Average and maximum speed for cornwall
94 cornwall_speed <-broadband %>%
95   filter(County=="CORNWALL") %>%
96   group_by(Town_City) %>%
97   summarize(
98     AvgDownSpeed=mean(AvgDownSpeed),
99     MaxDownSpeed=max(MaxDownSpeed)
100   ) %>%
101   pivot_longer(cols=c(AvgDownSpeed,MaxDownSpeed),names_to="SpeedType",values_to = "Speed")
102
103 View(cornwall_speed)
104
105 #Visualization
106 ggplot(cornwall_speed, aes(x = Town_City, y = Speed, fill = SpeedType)) +
107   geom_bar(stat = "identity", position = "dodge") +
108   labs(title = "Average and Maximum Download Speeds by Town/City in Cornwall",
109        x = "Town/City",
110        y = "Speed (Mbit/s)") +
111   theme_minimal() +
112   theme(axis.text.x = element_text(angle = 90, hjust = 1),
113         legend.position = "top")
114
115 #Average and maximum speed for Bistol
116 bistol_speed <-broadband %>%
117   filter(County=="CTY OF BISTOL") %>%
118   group_by(Town_City) %>%
119   summarize(
120     AvgDownSpeed=mean(AvgDownSpeed),
121     MaxDownSpeed=max(MaxDownSpeed)
122   ) %>%
123   pivot_longer(cols=c(AvgDownSpeed,MaxDownSpeed),names_to="SpeedType",values_to = "Speed")
124
48:41 (Top Level)
R Script

```

```

116 #Average and maximum speed for Bristol
117 bistro1_speed <- broadband %>%
118   filter(County=="CITY OF BRISTOL") %>%
119   group_by(Town_City) %>%
120   summarize(
121     AvgDownSpeed=mean(AvgDownSpeed),
122     MaxDownSpeed=max(MaxDownSpeed)
123   ) %>%
124   pivot_longer(cols=c(AvgDownSpeed,MaxDownSpeed),names_to="SpeedType",values_to = "Speed")
125
126 View(bistro1_speed)
127
128 #Visualization
129 ggplot(bistro1_speed, aes(x = Town_City, y = Speed, fill = SpeedType)) +
130   geom_bar(stat = "identity", position = "dodge") +
131   labs(title = "Average and Maximum Download Speeds by Town/City in Bristol",
132        x = "Town/City",
133        y = "Speed (Mbit/s)") +
134   theme_minimal() +
135   theme(axis.text.x = element_text(angle = 90, hjust = 1),
136         legend.position = "top")
137
138 #-----CRIME-----
139 crime_data = read_csv("D:/Academics/Fourth Semester/Data Science/Assignment/Cleaned Datasets/crime_cleaned.csv")
140 view(crime_data)
141 colnames(crime_data)
142 dim(crime_data)
143
144 # Vehicle Crime per 10000 people in April 2022
145 vehicleCrime <- crime_data %>%
146   filter('Crime type' == "Vehicle crime")
147 colnames(vehicleCrime)
148

```

```

149 vehicle_sum <- vehicleCrime %>%
150   group_by(Year) %>%
151   summarise(total_crime = sum(population, na.rm = TRUE))
152
153 View(vehicle_sum)
154
155 for_radar = as.data.frame(t(vehicle_sum$total_crime))
156 colnames(for_radar) <- vehicle_sum$Year
157 view(for_radar)
158
159 for_radar <- rbind(rep(max(vehicle_sum$total_crime), length(years)),
160                  rep(0, length(years)),
161                  for_radar)
162
163 par(mar = c(2, 2, 2, 2))
164
165 #Plotting Radar Chart for Vehicle Crime
166 radarchart(for_radar,
167            axistype = 1,
168            pcol = "purple",
169            pfcol = "lightblue",
170            plwd = 4,
171            title = "Vehicle Crime Rate from 2021 to 2024"
172          )
173
174
175 #Pie Chart for Robbery in 2023 by month
176
177 robbery_data = crime_data %>%
178   filter('Crime type' == "Robbery" & Year == "2023") %>%
179   group_by(Month) %>%
180   summarise(count=n()) %>%
181   mutate(percentage=count/sum(count)*100)
182
183 #Pie Chart for Robbery in 2023 by month
184
185 robbery_data = crime_data %>%
186   filter('Crime type' == "Robbery" & Year == "2023") %>%
187   group_by(Month) %>%
188   summarise(count=n()) %>%
189   mutate(percentage=count/sum(count)*100)
190
191 #Plotting the Pie Chart
192 ggplot(robbery_data, aes(x = "", y = percentage, fill = as.factor(Month))) +
193   geom_bar(width = 1, stat = "identity") +
194   coord_polar("y") +
195   geom_text(aes(label = paste0(round(percentage, 1), "%")),
196            position = position_stack(vjust = 0.5)) +
197   labs(title = "Robberies by Month in 2023", fill = "Month") +
198   theme_minimal()
199
200 #Drugs
201 drugCrime <- subset(crime, 'Crime type' == "Drugs" & Year == 2023 & counties %in% c("Bristol, city of", "Cornwall"))
202
203 # create the boxplot
204 ggplot(drugCrime, aes(x = counties, y = population / 10000, fill = counties)) +
205   geom_boxplot() +
206   labs(title = "Drug Offence Rate per 10,000 People in 2023",
207        x = "County",
208        y = "Drug Offence Rate (per 10,000 people)") +
209   theme_minimal()
210

```



```

207
208 # Filtering data for drug offences
209 drug2023 = crime %>%
210   filter('crime type' == "Drugs" & Year == "2023")
211
212 # For Cornwall
213 cornDrug <- drugCrime %>%
214   filter(counties == "Cornwall") %>%
215   distinct('LSOA code', .keep_all = TRUE) %>%
216   summarise(total_population = sum(population),
217             total_drug_offences = n())
218 cornwall_data
219
220
221 # Summarize data for Bristol
222 bristolDrug <- drugCrime %>%
223   filter(counties == "Bristol, City of") %>%
224   distinct('LSOA code', .keep_all = TRUE) %>%
225   summarise(total_population = sum(population),
226             total_drug_offences = n())
227
228 cornDrug <- cornDrug %>%
229   mutate(offence_rate = (total_drug_offences / total_population) * 10000)
230
231 bristolDrug <- bristolDrug %>%
232   mutate(offence_rate = (total_drug_offences / total_population) * 10000)
233
234
235 combined_data <- bind_rows(
236   cornDrug %>% mutate(county = "Cornwall"),
237   bristolDrug %>% mutate(county = "Bristol, City of")
238 )
239

```

```

239
240 ggplot(combined_data, aes(x = county, y = offence_rate, fill = county)) +
241   geom_boxplot() +
242   labs(title = "Distribution of Drug Offence Rates (2023)",
243        x = "Location",
244        y = "Offence Rate (per 10,000)") +
245   theme_minimal()
246
247 # ----- Schools
248
249 school_data = read_csv("D:/Academics/Fourth Semester/Data Science/Assignment/Cleaned Datasets/schoolcleaned.csv")
250
251 head(school_data)
252 school_filtered = school_data %>%
253   filter(YEAR==2022)
254
255
256 #BOXPLOT: Average Attainment 8 scores by County in 2022
257 ggplot(school_filtered, aes(x = COUNTY, y = ATT8SCR, fill = COUNTY)) +
258   geom_boxplot() +
259   labs(title = "Average Attainment 8 Scores by County in 2022",
260        x = "County",
261        y = "Attainment 8 Score") +
262   theme_minimal()
263
264 #Line Graph : Average Attainment 8 Score in Academic Year 2021-2022 in BISTROL
265 bristol_21=school_data %>%
266   filter(YEAR==2021) %>%
267   filter(COUNTY=="Bristol")
268
269 ggplot(bristol_21, aes(x = SCHNAME, y = ATT8SCR, group = 1)) +
270   geom_line(color = "purple") +

```

```

271   theme_minimal()
272
273 #Line Graph : Average Attainment 8 Score in Academic Year 2021-2022 in BISTROL
274 bristol_21=school_data %>%
275   filter(YEAR==2021) %>%
276   filter(COUNTY=="Bristol")
277
278 ggplot(bristol_21, aes(x = SCHNAME, y = ATT8SCR, group = 1)) +
279   geom_line(color = "purple") +
280   geom_point() +
281   labs(title = "Average Attainment 8 Score in 2021-2022 in Bristol",
282        x = "School Name",
283        y = "Average Attainment 8 Score") +
284   theme(axis.text.x = element_text(angle = 90, hjust = 1))
285
286 #Line Graph : Average Attainment 8 Score in Academic Year 2021-2022 in CORNWALL
287 cornwall_21=school_data %>%
288   filter(YEAR==2021) %>%
289   filter(COUNTY=="Cornwall")
290 colnames(cornwall_21)
291 head(cornwall_21)
292
293 ggplot(cornwall_21, aes(x = TOWN, y = ATT8SCR, group = 1)) +
294   geom_line(color = "red") +
295   geom_point() +
296   labs(title = "Average Attainment 8 Score in 2021-2022 in Cornwall",
297        x = "Town",
298        y = "Average Attainment 8 Score") +
299   theme(axis.text.x = element_text(angle = 90, hjust = 1))
300
301

```

## LinearModel.R

```

1 library(tidyverse)
2 library(ggplot2)
3
4 #Importing Datasets
5 crimeData=read_csv("D:/Academics/Fourth Semester/Data Science/Assignment/Cleaned Datasets/crime_cleaned.csv")
6 schoolData=read_csv("D:/Academics/Fourth Semester/Data Science/Assignment/Cleaned Datasets/school_cleaned.csv")
7 housingData=read_csv("D:/Academics/Fourth Semester/Data Science/Assignment/Cleaned Datasets/housing_cleaned.csv")
8 broadbandData=read_csv("D:/Academics/Fourth Semester/Data Science/Assignment/Cleaned Datasets/broadband_cleaned.csv")
9
10 #Linear Model for : Housing Price vs Average Download Speed
11 housingPrice <- housingData %>%
12   select(Postcode,Price)
13
14 downSpeed <- broadbandData %>%
15   select(Postcode,AvgDownSpeed)
16
17 housing_speed <- merge(housingPrice, downSpeed, by="Postcode")
18
19 View(housing_speed)
20
21 ggplot(data = housing_speed, aes(x = Price, y = AvgDownSpeed)) +
22   geom_point() +
23   geom_smooth(method = "lm", se = FALSE, color = "blue") +
24   labs(title = "House Prices vs Average Download Speed",
25        x = "House Price",
26        y = "Average Download Speed (Mbit/s)") +
27   scale_y_continuous(labels = scales::comma) +
28   scale_x_continuous(labels = scales::comma) +
29   theme_minimal()
30
31 #Linear Model for Housing Price Vs Drug Offence
32 housingPrice23=housingData %>%
33   filter(Year==2023) %>%
34   mutate(Postcode=str_trim(substring(Postcode,1,6))) %>%
35   distinct()
36
37 drug_offence_23 <-crimeData %>%
38   filter(Year==2023, `Crime type`=="Drugs") %>%
39   group_by(postcode) %>%
40   summarise(DrugRate=sum(population)) %>%
41   distinct() %>%
42   na.omit()
43
44 housing_drugs <- housingPrice23 %>%
45   left_join(drug_offence_23,by=c("Postcode"="postcode")) %>%
46   distinct() %>%
47   na.omit()
48
49 ggplot(data = housing_drugs, aes(x = Price, y = DrugRate)) +
50   geom_point() +
51   geom_smooth(method = "lm", se = FALSE, color = "blue") +
52   labs(title = "Housing Price vs Drug Rates in 2023",
53        x = "House Price",
54        y = "Drug Rate") +
55   scale_y_continuous(labels = scales::comma) +
56   scale_x_continuous(labels = scales::comma) +
57   theme_minimal()
58
59 #Linear Model for Average Download Speed Vs Attainment 8 Score
60 school <- schoolData %>%
61   filter(YEAR==2022) %>%

```

```

59 #Linear Model for Average Download Speed Vs Attainment 8 Score
60 school <- schoolData %>%
61   filter(YEAR==2022) %>%
62   select(PCODE, ATT8SCR) %>%
63   rename(Postcode="PCODE")
64
65 school_broadband <-merge(downSpeed,school, by= "Postcode")
66
67 ggplot(data = school_broadband, aes(x = AvgDownSpeed, y = ATT8SCR)) +
68   geom_point() +
69   geom_smooth(method = "lm", se = FALSE, color = "blue") +
70   labs(title = "Average download speed VS Attainment 8 Score",
71        x = "Average download speed",
72        y = "Attainment 8 Score") +
73   scale_y_continuous(labels = scales::comma) +
74   scale_x_continuous(labels = scales::comma) +
75   theme_minimal()
76
77
78 #Linear Model for Attainment 8 Score Vs Housing Price
79 school_housing = merge(school, housingPrice, by = "Postcode")
80
81 ggplot(data = school_housing, aes(x = ATT8SCR, y = Price)) +
82   geom_point() +
83   geom_smooth(method = "lm", se = FALSE, color = "blue") +
84   labs(title = "Average attainment 8 score vs Housing Price",
85        x = "ATT8SCR",
86        y = "Price") +
87   scale_y_continuous(labels = scales::comma) +
88   scale_x_continuous(labels = scales::comma) +
89   theme_minimal()
90

```

## Ranking.R

```

1 library(tidyverse)
2
3 #Importing Libraries
4 crime=read_csv("D:/Academics/Fourth Semester/Data Science/Assignment/Cleaned Datasets/crime_cleaned.csv")
5 school=read_csv("D:/Academics/Fourth Semester/Data Science/Assignment/Cleaned Datasets/schoolcleaned.csv")
6 housingPrice=read_csv("D:/Academics/Fourth Semester/Data Science/Assignment/Cleaned Datasets/housing_cleaned.csv")
7 broadband=read_csv("D:/Academics/Fourth Semester/Data Science/Assignment/Cleaned Datasets/broadband_cleaned.csv")
8
9 dim(crime)
10
11 #Selecting Relevent Rows from Broadband
12 selected_broadband= broadband %>%
13   group_by(Town_City) %>%
14   summarise(
15     avg_upl_speed = mean(AvgUpSpeed),
16     avg_down_speed = mean(AvgDownSpeed)
17   )%>%
18   mutate(TOWN = str_trim(toupper(Town_City)))%>%
19   select(TOWN, avg_upl_speed, avg_down_speed)
20
21 selected_attainment8= school %>%
22   group_by(TOWN) %>%
23   summarise(avgAtt8 = mean(ATT8SCR))%>%
24   select(TOWN, avgAtt8)%>%
25   distinct()%>%
26   mutate(TOWN= str_trim(toupper(TOWN)))
27
28 selected_house = housingPrice %>%
29   filter(Year == 2023)%>%
30   mutate(TOWN = str_trim(toupper(Town_city)))%>%
31   group_by(TOWN) %>%
32   summarise(avgPrice = mean(Price))%>%
33   select(avgPrice, TOWN)%>%
34   distinct()
35
127:1 (Top Level) :

```



```

28 selected_house = housingPrice %>%
29   filter(Year == 2023)%>%
30   mutate(TOWN = str_trim(toupper(Town_City)))%>%
31   group_by(TOWN) %>%
32   summarise(avgPrice = mean(Price))%>%
33   select(avgPrice, TOWN)%>%
34   na.omit()%>%
35   distinct()
36
37 town = housingPrice %>%
38   mutate(postcode = str_trim(substring(Postcode, 1, 6))) %>%
39   mutate(TOWN = str_trim(toupper(Town_City)))%>%
40   select(postcode, TOWN)%>%
41   distinct()
42
43 selected_crime = crime %>%
44   filter(Year==2023)%>%
45   group_by(postcode)%>%
46   summarise(crimeno = n())%>%
47   arrange(desc(crimeno))%>%
48   select(postcode, crimeno)
49
50 final_crime = selected_crime%>%
51   left_join(town, by= "postcode")%>%
52   na.omit()%>%
53   distinct()
54
55 final_crime = final_crime%>%
56   group_by(TOWN)%>%
57   summarise(crimerate = sum(crimeno))%>%
58   select(TOWN, crimerate)
59
18:50 (Top Level) R Script

```

```

60 ranking = selected_house %>%
61   left_join(selected_attainment8, by = "TOWN") %>%
62   left_join(selected_broadband, by = "TOWN") %>%
63   left_join(final_crime, by = "TOWN") %>%
64   na.omit()
65 view(ranking)
66
67
68 # calculating the minimum and Maximum for each column
69 Extremes <- ranking_points %>%
70   summarise(
71     minDownSpeed= min(avg_down_speed),
72     maxDownSpeed = max(avg_down_speed),
73     minUpSpeed = min(avg_upl_speed),
74     maxUpSpeed = max(avg_upl_speed),
75     minAtt8 = min(avgAtt8),
76     maxAtt8 = max(avgAtt8),
77     minHousingPrice = min(avgPrice),
78     maxHousingPrice = max(avgPrice),
79     minCrimeRate = min(crimerate),
80     maxCrimeRate = max(crimerate)
81   )
82
83 # Normalizing and calculating the final points
84 finalRanking <- ranking_points %>%
85   mutate(
86     normDownSpeed = (avg_down_speed - Extremes$minDownSpeed) / (Extremes$maxDownSpeed - Extremes$minDownSpeed),
87     normUpSpeed = (avg_upl_speed - Extremes$minUpSpeed) / (Extremes$maxUpSpeed - Extremes$minUpSpeed),
88     normAtt8 = (avgAtt8 - Extremes$minAtt8) / (Extremes$maxAtt8 - Extremes$minAtt8),
89     normHousingPrice = 1 - (avgPrice - Extremes$minHousingPrice) / (Extremes$maxHousingPrice - Extremes$minHousingPrice),
90     normCrimeRate = 1 - (crimerate - Extremes$minCrimeRate) / (Extremes$maxCrimeRate - Extremes$minCrimeRate),
91     finalPoints = normDownSpeed + normUpSpeed + normAtt8 + normHousingPrice + normCrimeRate
92   )
93
18:50 (Top Level) R Script

```

```

94 # view the data with the final score
95 finalRanking= finalRanking%>%
96   arrange(desc(finalPoints))
97
98 view(finalRanking)
99 write_csv(finalRanking, "D:/Academics/Fourth Semester/Data Science/Assignment/ranking.csv")
100
101 houserank = finalRanking %>%
102   select(TOWN, avgPrice, normHousingPrice)%>%
103   arrange(desc(normHousingPrice))
104
105 view(houserank)
106
107 crimerank = finalRanking %>%
108   select(TOWN, normCrimeRate, crimerate) %>%
109   arrange(desc(normCrimeRate))
110 view(crimerank)
111
112 schoolrank = finalRanking %>%
113   select(TOWN, avgAtt8, normAtt8)%>%
114   arrange(desc(normAtt8))
115 view(schoolrank)
116
117 broadbandrank = finalRanking %>%
118   select(TOWN, avg_down_speed, normDownSpeed)%>%
119   arrange(desc(normDownSpeed))
120
121 view(broadbandrank)
122
123 final_rank = finalRanking %>%
124   select(TOWN, avgPrice, crimerate, avgAtt8, avg_down_speed, finalPoints) %>%
125   arrange(desc(finalPoints))
126
18:50 (Top Level) R Script

```