



Individual Report (Data Science for Developers)

BSc. (Hons.) Computing, Softwarica College of IT and E-commerce, University of Coventry

ST5014CEM (Data Science for Developers)

Rohan Chaulagain

19th August, 2024

GitHub link: <https://github.com/0to1e/DS-Assignment.git>

Table of Contents

Introduction.....	5
Data Cleaning.....	6
1. House Pricing Data:	7
Raw Dataset:	7
Cleaning Process:.....	7
Final Outcome:.....	7
2. Broadband Data:	8
Raw Datasets:.....	8
Issues Identified:	8
Cleaning Process:.....	8
Final Outcome:.....	9
3. Crime Data:	9
Raw Dataset:	9
Issues Identified:	9
Cleaning Process:.....	9
Final Outcome:.....	10
4. School Data:	10
Raw Datasets:.....	10
Issues Identified:	10
Cleaning Process:.....	10
Final Outcome:.....	11
EDA	11
1. EDA of House Pricing:	11
i. Boxplot of House Prices in 2023 for CITY OF BRISTOL and CORNWALL	12
ii. Total House Prices by Town within County (2023).....	12
2. EDA of Broad band.....	12
i. Average Download Speed by County2020-2023	12
ii. Bristol: Maximum vs Average Download Speed by City.....	13
iii. Cornwall: Maximum vs Average Download Speed by City.....	13
3. EDA of Crime	13
i. Yearly Drug Offence Rate per 10000 Population (2020-2023).....	13
ii. Drug Offence Rate (2023)	14
iii. Robbery Offences Distribution (October 2022)	14
4. EDA of School:	14

i.	Boxplot of Attainment 8 score by County	14
ii.	Average Attainment 8 score by City in Cornwall	15
iii.	Average Attainment 8 score by City in Bristol	15
	Linear Modeling:.....	15
i.	School ATT8 Scores and House Prices	15
ii.	Average Download Speed and Attainment 8 score.....	16
iii.	Average Download Speed and Drug Offence Rates per 10000 people	16
iv.	House Price vs Average Download Speed.....	17
v.	House Price vs Drug Offence Rate (2023).....	17
	Ranking:	18
	Overview.....	18
	House Pricing.....	19
	Broad band	19
	Crime.....	19
	Results	19
	Reflection.....	21
	Overall Score:	21
	Legal and Ethical Issue	22
	Conclusion:	22
	References	23
	Appendix	23

Table of Figures

Figure 1: Final Ranking	21
Figure 2: cleaning code for house pricing.....	23
Figure 3: Cleaning code for broadband.....	24
Figure 4: Cleaning code for crime	24
Figure 5: Cleaning code for school	25
Figure 6: cleaning post code	26
Figure 7: Cleaning Population dataset	27
Figure 8: Code to generate House Price Boxplot.....	27
Figure 9: House Price BoxPlot.....	28
Figure 10: Code to generate Stacked BarChart for House Prices 2023	28
Figure 11: Stacked Barchart House Price (2023).....	29
Figure 12:Code to generate Line graph for house price by town.....	29
Figure 13: Line Graph for house prices by town	29
Figure 14: Code to generate Boxplot for Average download speed by county.....	30
Figure 15: Box plot for average internet speed by County	31
Figure 16: Dodge Barchart for Average vs Maximum Download Speed in Bristol	31
Figure 17:Code to generate dodge bar chart of maximum and average download speed by City of Bristol.....	31
Figure 18: Figure 16: Code to generate Dodge Barchart for Average vs Maximum Download Speed in Cornwall.....	32
Figure 19:Dodge bar chart for Maximum vs Average Download Speed in Cornwall.....	33
Figure 20: Code to generate boxplot of Drug Offence rate per 10000 people 2023.....	33
Figure 21: Box plot for drug offence (2023).....	34
Figure 22: Code to generate Line graph for drug offence rate 2020-2023	34
Figure 23:Line graph for Yearly drug offence rate 2020-2023	35
Figure 24: Pie chart for robbery rate on October 2022	36
Figure 25: Code to generate Piechart for robberies on October 2022	36
Figure 26: Radar chart for vehicle crime rates from 2020-2023.....	37
Figure 27: Radar chart for yearly vehicle crimes 2020-2023	38
Figure 28: Code to generate Boxplot of attainment 8 score 2023	39
Figure 29: Attainment 8 score by county 2023	39
Figure 30: Code to lime graph of attainment 8 score city wise in Bristol and Cornwall.....	40
Figure 31: Line graph to see average attainment 8 score in Cornwall.....	41
Figure 32: Line graph to see average attainment 8 score in Bristol.....	41
Figure 33: Code to perform linear modelling of attainment 8 score and house price.....	42
Figure 34: Code to perform Linear Modelling between Attainment 8 score and average download speed	43
Figure 35: Code to perform linear modeling between drug offence rate per 1000 people and average download speed.....	43
Figure 36: Code to perform linear modelling between house prices and drug offence rate in 2023	45
Figure 37: Code to perform linear modelling of house price and Average download speed	46
Figure 38: Code to generate line of best fit between download speed and House Price	46
Figure 39: Line of best fit between download speed and house price	46
Figure 40: Code to generate line of best fit between drug offence rate and house price 2023	47
Figure 41: line of best fit between drug offence rate and house price 2023	47
Figure 42: Code to generate line of best fit for house price and attainment 8 score.....	48
Figure 43: line of best fit for house price and attainment 8 score.....	48

Figure 44: Code to generate line of best fit for download speed and drug offence rate per 1000 people	49
Figure 45:line of best fit for download speed and drug offence rate per 10000 people	49
Figure 46:Code to generate line of best fit for average download speed and attainment 8 score.....	50
Figure 47: line of best fit for average download speed and attainment 8 score.....	50
Figure 48: Code to perform city ranking for houses.....	51
Figure 49: Code to perform city ranking by broadband.....	51
Figure 50: Code to perform city ranking via schools.....	52
Figure 51:Code to perform city ranking via crime	53
Figure 52: Code to perform final ranking of cities	54

Introduction

This project involves analyzing datasets sourced from the UK government to generate data-driven recommendations for property investments in Bristol and Cornwall. Key factors such as housing prices, internet connectivity, and crime rates are examined. The datasets used include house price data from the housing and communities' section of the UK government website, crime statistics from Avon and Somerset Constabulary and Devon & Cornwall Police, school performance data from the UK government's "Compare school and college performance in England" section, and broadband data from OFCOM. Additionally, postcode-to-LSOA and population data based on postcodes were utilized, with the census data from 2011 being updated to reflect the 2023 population.

These datasets were deemed suitable as they were all publicly available and obtained from reliable sources, including official UK government websites, police websites, and OFCOM. This project provides an opportunity to enhance skills in data cleaning, normalization, exploratory data analysis, and statistical investigation. By developing a recommendation system that integrates various data attributes, practical experience in deriving actionable insights is gained. This process contributes to a deeper understanding of the data science life-cycle and improves the ability to make informed decisions based on complex datasets.

Data Cleaning

Before performing analytics and building the recommendation system, cleaning the data obtained from different sources is the foremost crucial process which involves identifying and correcting errors, inconsistencies, and missing values within a dataset. Cleaning the data supports ensuring data accuracy, consistency, and completeness, cleaning significantly enhances the reliability of subsequent analyses. This foundational step prevents misleading results, improves model performance, and ultimately empowers data-driven decision-making.

Talking about this task, datasets related to four entities, i.e. House Pricing data, Broadband Connection, Crime and Schools of UK has been provided to us to work on throughout the analysis process. Further discussion involves the cleaning process of each raw dataset obtained.

1. House Pricing Data:

Raw Dataset:

The house pricing data that was obtained wasn't usable at all when first received for following reasons:

- i. No header name, i.e. unable to identify composition of datasets.
- ii. Data separated into four different files yearwise, making it difficult to access the dataset.
- iii. Presence of Null values.
- iv. Presence of unwanted columns in the dataset.
- v. Improper date format.

Cleaning Process:

- i. Assigned standardized column names.
- ii. Selected columns only which are necessary to perform the analysis.
I.e.(Price, Date_of_Transfer, Postal.Code, Property_Type, Town.City, District, County)
- iii. Mutated date format appropriately.
- iv. Removed null values present.

Final Outcome:

- i. Only necessary columns in the cleaned dataset.
- ii. No null values
- iii. Appropriate datatypes handled

For detailed code implement, please refer to Appendix.

2. Broadband Data:

Raw Datasets:

The initial datasets obtained included:

- i. Broadband Coverage Data: Contains information on broadband availability.
- ii. Broadband Performance Data: Includes data on broadband performance metrics.

Issues Identified:

- i. No Standardized Column Names: Column names were inconsistent or unclear.
- ii. Irrelevant Columns: Some columns were not necessary for analysis.
- iii. Inconsistent Data: Postcodes had inconsistent formatting and were not aligned with other datasets.
- iv. Missing Values: Various columns had missing values that needed to be handled.
- v. Lack of City and County Information: Needed to enrich data with additional geographic information.

Cleaning Process:

- i. Load and Clean Broadband Coverage Data:
 - 1. Columns Selected: Focused on postcode, SFBB availability (% premises), UFBB availability (% premises), and FTTP availability (% premises).
 - 2. Postcode Formatting: Removed spaces from postcodes for consistency.
 - 3. Renaming: Standardized the column name for postcode to PostalCode.

- ii. Load and Clean Broadband Performance Data:
 - 1. Columns Selected: Focused on postcode, Maximum download speed (Mbit/s), and Average download speed (Mbit/s).
 - 2. Renaming: Standardized the column name for postcode to PostalCode.

- iii. Load and Clean Postcode Data:
 - 1. Postcode Data: Used to merge and enrich broadband datasets with city and county information.

- iv. Merge Data with Postcode Information:
 - 1. Broadband Coverage Data: Joined with postcode data to add city and county information. Assigned "NA" for unmatched entries.
 - 2. Broadband Performance Data: Similarly joined with postcode data, handled missing values by imputing the median for download speeds.

Final Outcome:

- i. Cleaned Datasets: Only relevant columns retained, and unnecessary ones removed.
- ii. Consistent Formatting: Postcodes standardized and merged with geographic information.
- iii. Handled Missing Values: Missing values addressed with median or mean imputation as appropriate.

For detailed code implementation, please refer to the Appendix.

3. Crime Data:

Raw Dataset:

- i. Crime Data Files: Multiple CSV files containing crime data.

Issues Identified:

- i. Multiple Files: Data split across several CSV files.
- ii. Irrelevant Columns: Not all columns were necessary.
- iii. Formatting Issues: Inconsistent date formats and missing county information.

Cleaning Process:

- i. Combine CSV Files: Merged all CSV files into a single dataset.
- ii. Select and Clean Data:
 - 1.Retained relevant columns and formatted Month to date.
 - 2.Filtered for crimes in Bristol or Cornwall.
 - 3.Renamed columns and assigned County based on LSOA name.

Final Outcome:

- iii. Consolidated Data: Combined multiple files into a single dataset.
- iv. Relevant Columns: Only necessary columns retained.
- v. Data Consistency: Dates formatted and county information assigned.

4. School Data:**Raw Datasets:**

- i. School Data Files: Key Stage 4 final data from Bristol and Cornwall for two academic years.

Issues Identified:

- i. Multiple Files: Data spread across multiple files for different years and regions.
- ii. Irrelevant Columns: Not all columns were necessary.
- iii. Missing Values: Missing school names and values in key columns.
- iv. Inconsistent Formatting: Postcode formatting needed standardization.

Cleaning Process:

- i. Load and Prepare Data:
 1. Loaded datasets for Bristol and Cornwall for the academic years 2021-2022 and 2022-2023.
 2. Selected relevant columns and added a Year column to each dataset.
- ii. Combine and Clean Data:
 1. Merged all datasets into one.
 2. Filtered out rows with missing school names.
 3. Converted and imputed missing values in ATT8SCR.
 4. Standardized PCODE, joined with postcode data, and handled missing city and county values.

Final Outcome:

- i. Consolidated Data: Combined multiple datasets into a single, comprehensive dataset.
- ii. Relevant Columns: Focused on necessary columns and handled missing values.
- iii. Data Consistency: Postcode formatting standardized, city and county information added.

EDA

1. EDA of House Pricing:

i. Boxplot of House Prices in 2023 for CITY OF BRISTOL and CORNWALL

- Graph: Boxplot
- Datasets Used: Cleaned House Pricing Data
- X-axis: County – Categorical
- Y-axis: House Price – Quantitative
- Summary: This boxplot compares the distribution of house prices in Bristol and Cornwall for the year 2023, showcasing the variation and central tendency of house prices in these areas.
- Result:

ii. Total House Prices by Town within County (2023)

- Graph: Stacked Bar Chart:
- Datasets Used: Cleaned House Pricing Data
- X-axis: County – Categorical
- Y-axis: Total Price – Quantitative
- Summary: This stacked bar chart illustrates the total house prices by town within each county for the year 2023, allowing comparison of house price contributions from different towns within Bristol and Cornwall.

2. EDA of Broad band

i. Average Download Speed by County 2020-2023

- Graph: Boxplot
- Datasets Used: Cleaned Broadband Performance Data
- X-axis: County - Categorical
- Y-axis: Average Download Speed (Mbit/s) – Quantitative
- Summary: This boxplot visualizes the distribution of average download speeds across different counties, providing insights into broadband performance variability within each county.

ii. Bristol: Maximum vs Average Download Speed by City

- Graph: Side-by-Side Bar Chart
- Datasets Used: Cleaned Broadband Performance Data
- X-axis: City – Categorical
- Y-axis: Download Speed (Mbit/s) – Quantitative
- Summary: This side-by-side bar chart compares maximum and average download speeds for different cities within Bristol, highlighting variations in broadband performance across these cities.

iii. Cornwall: Maximum vs Average Download Speed by City

- Graph: Side-by-Side Bar Chart
- Datasets Used: Cleaned Broadband Performance Data
- X-axis: City - Categorical
- Y-axis: Download Speed (Mbit/s) - Quantitative
- Summary: This side-by-side bar chart compares maximum and average download speeds for different cities within Cornwall, illustrating the variation in broadband performance across these cities.

3. EDA of Crime

- i. Yearly Drug Offence Rate per 10000 Population (2020-2023)
 - Graph: Line Graph
 - Datasets Used: Cleaned Crime Data
 - Related Datasets: Cleaned Population Data
 - X-axis: Year – Quantitative
 - Y-axis: Offence Rate per 10,000 Population - Quantitative
 - Summary: This line graph visualizes the yearly drug offence rates per 10,000 population for Bristol and Cornwall from 2020 to 2023, highlighting trends and comparisons between these regions.
- ii. Drug Offence Rate (2023)
 - Graph: Boxplot
 - Datasets Used: Cleaned Crime Data
 - Related Datasets: Cleaned Population Data
 - X-axis: County – Categorical
 - Y-axis: Offence Rate - Quantitative
 - Summary: This boxplot displays the distribution of drug offence rates in 2023 across different counties, highlighting variations in the rate of drug-related crimes.
- iii. Robbery Offences Distribution (October 2022)
 - Graph: Pie Chart
 - Datasets Used: Cleaned Crime Data
 - Related Datasets: Cleaned Population Data
 - X-axis: N/A (Pie chart does not have an X-axis)
 - Y-axis: Percentage - Quantitative
 - Summary: This pie chart visualizes the distribution of robbery offences across different counties for October 2022, showing the proportion of total offences represented by each county.

4. EDA of School:

- i. Boxplot of Attainment 8 score by County
 - Graph: Boxplot
 - Datasets Used: Cleaned Schools Data
 - X-axis: County - Categorical
 - Y-axis: ATT8SCR - Quantitative
 - Summary: This boxplot displays the distribution of Attainment 8 Scores (ATT8SCR) across various counties for the year 2023, providing insight into the educational performance variations between counties.
- ii. Average Attainment 8 score by City in Cornwall
 - Graph: Line Graph
 - Datasets Used: Cleaned Schools Data
 - X-axis: Year - Quantitative
 - Y-axis: Average Attainment 8 score - Quantitative
 - Summary: This line graph visualizes the trend of average Attainment 8 scores by city within Cornwall over the years, highlighting how educational performance has changed over time in different cities.
- iii. Average Attainment 8 score by City in Bristol
 - Graph: Line Graph
 - Datasets Used: Cleaned Schools Data
 - X-axis: Year - Quantitative
 - Y-axis: Average Attainment 8 score - Quantitative
 - Summary: This line graph shows the trend of average Attainment 8 Scores by city within Bristol over the years, providing insights into the changes in educational performance across Bristol's cities.

Linear Modeling:

i. School ATT8 Scores and House Prices

- Independent variable: Price
- Dependent variable: ATT8SCR
- Equation: $\text{ATT8SCR} = 37.12 + 0.000003557 * \text{Price}$
- P-value : 0.257
- R squared: 0.0095
- RSD: 16.41
- Summary: The model's fit is poor as the R-squared value is very low (0.0095), indicating that Price explains only a small fraction of the variability in ATT8SCR. The residual standard error (16.41) is relatively high compared to the range of ATT8SCR values, suggesting substantial unexplained variation.

ii. Average Download Speed and Attainment 8 score

- Independent variable: Average download speed (Mbit/s)
- Dependent variable: ATT8SCR
- Equation: $\text{ATT8SCR} = 43.10 - 0.0471 * \text{Average download speed (Mbit/s)}$
- P-value : <2e-16
- R squared: 0.0104
- RSD: 12.6
- Summary: The model's fit is poor as the R-squared value is very low (0.0104), indicating that Average download speed explains only a minimal fraction of the variability in ATT8SCR. Despite the statistically significant coefficient, the residual standard error (12.6) is relatively high, suggesting considerable unexplained variation in ATT8SCR.

iii. Average Download Speed and Drug Offence Rates per 10000 people

- Independent variable: Average download speed (Mbit/s)
- Dependent variable: offence_rate
- Equation: $\text{offence_rate} = 0.0345 - 0.00001973 * \text{Average download speed (Mbit/s)}$
- R squared variable: 0.0000672
- RSD: 0.06649
- P-value: 0.065
- Summary: The model's fit is very poor as the R-squared value is extremely low (0.0000672), indicating that Average download speed explains only an almost negligible fraction of the variability in offence_rate. The residual standard error (0.06649) suggests significant unexplained variation, and the p-value (0.065) indicates that the coefficient is marginally significant.

iv. House Price vs Average Download Speed

- Independent variable: Average download speed (Mbit/s)
- Dependent variable: Price
- Equation: $\text{Price} = 347015.29 + 489.17 * \text{Average download speed (Mbit/s)}$
- R squared: 0.0005203
- RSD: 576300
- P-value: 5.2e-09
- Summary: The model's fit is very poor as the R-squared value is extremely low (0.0005203), indicating that Average download speed explains only a negligible fraction of the variability in Price. Although the coefficient is statistically significant (p-value < 2e-16), the residual standard error (576300) is quite large, suggesting substantial unexplained variation in Price.

v. House Price vs Drug Offence Rate (2023)

- Independent variable: Offence rate (per 10,000 population)
- Dependent variable: House Price
- Equation: Price = 274,700 - 43,240,000,000 * Offence rate (per 10,000 population)
- P-value: < 2e-16
- R squared: 0.0135
- RSD: 1,309,000
- Summary: The R-squared value of 0.0135 indicates that the offence rate explains only a small fraction (1.35%) of the variability in house prices. Despite the coefficient for the offence rate being statistically significant (with a p-value less than 0.001), the model's fit is quite limited. The residual standard deviation (RSD) of 1,309,000 suggests considerable unexplained variation in house prices, highlighting that the offence rate may not be a strong predictor of house prices.

In comparing the models, none demonstrate a strong fit to their respective data. The models for School ATT8 Scores and Average Download Speed have low R-squared values, indicating minimal explanatory power for Price and ATT8SCR, with substantial unexplained variation. The model relating Average Download Speed to Drug Offence Rates shows an almost negligible R-squared, suggesting that average download speed has little impact on offence rates. The model for House Price versus Average Download Speed and Drug Offence Rate both shows very low R-squared values, with the latter having a high residual standard deviation, reflecting significant unexplained variability. Overall, these results highlight that neither average download speed nor drug offence rates are strong predictors of house prices or educational attainment, emphasizing the need for alternative explanatory variables or models.

Ranking:

Overview

The process used to rank cities within each county based on four criteria: school performance, house pricing, broadband availability, and crime rates. The process involves several steps, including data loading, ranking within individual categories, and combining the results to produce a final ranking of cities.

1. Data Loading

1.1. School Rankings: The school_dataset was loaded from the file Cleaned_datasets/Schools/Schools_Clean.csv. This dataset contains information on school performance metrics.

House Pricing

1.2. House Pricing: The house_pricing_dataset was loaded from Cleaned_datasets/House_Pricing_Data/clean_house_pricing_data.csv. It includes data on house prices across different towns and cities.

Broad band

1.3. Broadband Coverage and Performance:

Broadband data was sourced from two files:

Cleaned_datasets/Broadband/clean_broadband_coverage.csv (coverage data)
 Cleaned_datasets/Broadband/clean_broadband_performance.csv (performance data)

Crime

1.4. Crime Data: The crime data was loaded from Cleaned_datasets/Crime/Crime_Data_Combined.csv.

This dataset includes information on various crime types and their frequencies.

2. Ranking Cities

2.1. Ranking Schools

Objective: Rank schools based on their ATT8SCR (Attainment 8 score).

Method:

- Grouped data by County.
- Arranged schools in descending order of ATT8SCR.
- Assigned ranks to schools using row_number().
- Filtered to include only the top 10 schools per county.

Results

2.2. Ranking House Pricing

Objective: Rank towns and cities based on a calculated score from house pricing data.

Method:

- Calculated percentiles (p25, p50, p75) and skewness for house prices.
- Computed a score using a weighted formula combining these percentiles and skewness.
- Ranked cities within each county based on the score.
- Filtered to include only the top 10 cities per county.

2.3. Ranking Broadband

Objective: Rank cities based on the weighted average score of various broadband metrics.

Method:

- Merged broadband coverage and performance datasets.
- Normalized the data for several broadband metrics.
- Calculated a weighted score for each city based on normalized metrics.
- Aggregated the score at the city level.
- Ranked cities within each county based on the aggregated score.
- Filtered to include only the top 10 cities per county.

2.4. Ranking Crime

Objective: Rank cities based on the total number of crime offences.

Method:

- Aggregated total offences for each crime type.
- Calculated the total number of offences for each city.
- Ranked cities within each county based on total crime offences.
- Filtered to include only the top 10 cities with the fewest total crime offences.

3. Combining Rankings

Objective: Create an overall ranking by combining rankings from all categories.

Method:

- Merged all individual rankings using full_join to include all cities, even if they are missing from some rankings.
- Replaced NA values with a high rank to ensure they are ranked lower.
- Calculated the average rank across all categories for each city.
- Sorted cities within each county based on the average rank.
- Selected the top 3 cities per county based on the overall average rank.

Conclusion: This methodology allows for a comprehensive ranking of cities within each county based on multiple factors, providing a holistic view of city performance across schools, housing, broadband, and crime.

Reflection

This project underscored the importance of a structured approach to data analysis. Starting with thorough data cleaning ensured that the analysis was based on reliable data. EDA offered critical insights and guided the modeling process. Linear modeling provided a quantitative understanding of factors influencing city rankings. Finally, combining and ranking the results offered a holistic view of city performance.

The iterative nature of the project—from cleaning to ranking—demonstrates the complexity and interconnectedness of data analysis tasks. Each phase was crucial in building towards a robust final product. The experience highlighted the value of an integrated approach to data science, where each step builds upon the previous one, culminating in actionable insights and informed decision-making.

Overall, this project was a valuable exercise in applying data science methodologies to real-world problems, reinforcing the importance of careful planning, rigorous analysis, and thoughtful interpretation.

Overall Score:

Based on all the calculations below, here is the table of final ranking of the cities countywise.

(row)	County	City
1	Bristol	BRISTOL
2	Bristol	Bristol 040
3	Bristol	Bristol 060E
4	Cornwall	Cornwall 005
5	Cornwall	Cornwall 040F
6	Cornwall	Cornwall 068

Figure 1: Final Ranking

The data used in this project was sourced entirely from open-access and legally compliant channels, ensuring adherence to all ethical and legal standards. The datasets were obtained from reputable and official sources, including the UK government's housing and community's portal, Avon and Somerset Constabulary, Devon & Cornwall Police, the UK government's school performance database, and OFCOM for broadband information. Additionally, postcode-to-LSOA and census data were leveraged, with updates made to reflect current population estimates. All sources are publicly available and are provided under transparent data-sharing agreements, eliminating any concerns regarding data privacy or unauthorized use. By relying on these legitimate and ethical data sources, the project not only adheres to legal requirements but also upholds high standards of integrity and transparency in data handling and analysis. This approach reinforces the credibility of the analysis and ensures that the insights generated are both reliable and compliant with ethical norms.

Conclusion:

This project vividly illustrated the significance of a methodical approach to data analysis. Commencing with meticulous data cleaning ensured that our analysis was grounded in accurate and reliable data. Exploratory Data Analysis (EDA) played a pivotal role in uncovering critical insights and steering the subsequent modeling efforts. Linear modeling then offered a quantitative perspective on the factors impacting city rankings, enhancing our understanding of their relative importance.

The iterative process—from cleaning through to final ranking—highlighted the complexity and interdependency inherent in data analysis tasks. Each phase, from initial data preparation to the synthesis of final rankings, was integral to achieving a comprehensive and insightful final product. This project underscored the value of an integrated data science approach, where each step builds upon the last, leading to actionable insights and well-informed decision-making.

In conclusion, this project served as a practical application of data science methodologies to real-world issues, reaffirming the importance of careful planning, rigorous analysis, and thoughtful interpretation.

References

1. Dataquest. (2020, July 24). Loading and cleaning data with R and the tidyverse. Retrieved from <https://www.dataquest.io/blog/load-clean-data-r-tidyverse/>
2. Lamethods. (n.d.). An R approach to data cleaning and wrangling for education. Retrieved from <https://lamethods.org/chapters/ch04-data-cleaning/ch4-datacleaning.html>
3. Wickham, H., & Grolemund, G. (n.d.). Exploratory data analysis. In R for data science. Retrieved from <https://r4ds.had.co.nz/exploratory-data-analysis.html>
4. Wickham, H. (n.d.). Tidy data. In *tidyR*. Retrieved from <https://tidyR.tidyverse.org/articles/tidy-data.html>
5. Chang, W. (n.d.). R graphics cookbook. Retrieved from <https://www.cookbook-r.com/Graphs/>

Appendix

Figure 2: cleaning code for house pricing

```

1 # Defining column names as original data does'nt contain any
2 column_names <- c("Transaction_ID", "Price", "Date_of_Transfer", "Postal Code", "Property_Type", "Old/New", "Duration", "PAON", "SAON", "Street", "Locality", "Town/City", "District", "County", "PPD_Category_Type", "Record_Status")
3
4 postcode_data <- read_csv("Cleaned_datasets/Postcode_clean.csv") %>%
5   select(-pcd7, -Street)
6 # Read CSV files for house price data from 2020 to 2023 , each file containing a year's data
7 data0 <- read.csv("Obtained_Data/House_Pricing/pp-2020.csv", header = FALSE, col.names = column_names)
8 data1 <- read.csv("Obtained_Data/House_Pricing/pp-2021.csv", header = FALSE, col.names = column_names)
9 data2 <- read.csv("Obtained_Data/House_Pricing/pp-2022.csv", header = FALSE, col.names = column_names)
10 data3 <- read.csv("Obtained_Data/House_Pricing/pp-2023.csv", header = FALSE, col.names = column_names)
11
12 # Combine the datasets into one data frame
13 combined_data <- bind_rows(data0, data1, data2, data3)
14
15 cleaned_house_data <- combined_data %>%
16   # Select only the relevant columns
17   select(Price, Date_of_Transfer, Postal_Code, Town_City, County) %>%
18   mutate(County = ifelse(County == "CITY OF BRISTOL", "Bristol", County)) %>%
19   mutate(County = ifelse(County == "CORNWALL", "Cornwall", County)) %>%
20   rename("PostalCode" = "Postal.Code", "Date" = "Date_of_Transfer") %>%
21   mutate(PostalCode = gsub(" ", "", PostalCode)) %>%
22   # Filter for Bristol and Cornwall
23   filter((County %in% c("Bristol", "Cornwall")) %>%
24     # Convert Date_of_Transfer to Date type
25     mutate(Date = as.Date(Date)))
26
27 #join post_code dataset with combined house dataset to add street location data
28 final_house_data <- merge(
29   cleaned_house_data,
30   postcode_data,
31   by.x = "PostalCode",
32   by.y = "pcd7",
33   all.x = TRUE, # Keeps all rows from cleaned_house_data
34   all.y = FALSE # Excludes rows from postcode_data that do not match
35 ) %>%
36 na.omit()
37
38 write_csv(final_house_data, file = paste0(getwd(), "/Cleaned_datasets/House_Pricing_Data/clean_house_pricing_data.csv"))
39

```

Figure 3: Cleaning code for broadband

```

1 # Load and clean the broadband coverage data
2 broadband_coverage <- read_csv("Obtained_Data/broadband/201809_fixed_pc_coverage_r01.csv") %>%
3   # Select only relevant columns related to broadband coverage
4   select(postcode, `SFBB availability (% premises)`, `UFBB availability (% premises)`, `FTTP availability (% premises)` ) %>%
5   # Remove spaces from postcodes to ensure consistency
6   mutate(postcode = gsub(" ", "", postcode)) %>%
7   # Rename the postcode column to match other datasets
8   rename("PostalCode" = postcode)
9
10 # Load and clean the broadband performance data
11 broadband_performance <- read_csv("Obtained_Data/broadband/201805_fixed_pc_performance_r03.csv") %>%
12   # Select only relevant columns related to broadband performance
13   select(postcode, `Maximum download speed (Mbit/s)`, `Average download speed (Mbit/s)` ) %>%
14   # Rename the postcode column to match other datasets
15   rename("PostalCode" = postcode)
16
17 #Load postcode data to merge with street and County locations based on the postcode data available
18 postcode_clean <- read_csv("Cleaned_datasets/Postcode_clean.csv")
19
20 # Step 3: Join the broadband coverage data with postcode data to add city and county information
21 clean_coverage <- broadband_coverage %>%
22   inner_join(postcode_clean, by = c("PostalCode" = "pcd7")) %>%
23   # Handle any unmatched postcodes by assigning "NA" to City and County
24   mutate(City = ifelse(is.na(City), "NA", City), County = ifelse(is.na(County), "NA", County))
25
26 # Join the broadband performance data with postcode data to add city and county information
27 clean_performance <- broadband_performance %>%
28   inner_join(postcode_clean, by = c("PostalCode" = "pcd7")) %>%
29   mutate(City = ifelse(is.na(City), "NA", City), County = ifelse(is.na(County), "NA", County)) %>%
30   # Handle missing values in the Maximum download speed column by imputing the median
31   mutate(`Maximum download speed (Mbit/s)` = ifelse(is.na(`Maximum download speed (Mbit/s)`), median(`Maximum download speed (Mbit/s)`), na.rm = TRUE), `Maximum download speed (Mbit/s)` ) %>%
32   # Handle missing values in the Average download speed column by imputing the mean
33   mutate(`Average download speed (Mbit/s)` = ifelse(is.na(`Average download speed (Mbit/s)`), mean(`Average download speed (Mbit/s)`), na.rm = TRUE), `Average download speed (Mbit/s)` )
34
35 # Save the cleaned broadband coverage data to a CSV file
36 write_csv(clean_coverage, file = paste0(getwd(), "/Cleaned_datasets/Broadband/clean_broadband_coverage.csv"))
37
38 # Save the cleaned broadband performance data to a CSV file
39 write_csv(clean_performance, file = paste0(getwd(), "/Cleaned_datasets/Broadband/clean_broadband_performance.csv"))

```

Figure 4: Cleaning code for crime



```

1 # Define the directory containing the crime data files
2 crime_data_dir <- "Obtained_Data/Crime"
3 # List all CSV files in the directory recursively
4 all_files <- list.files(crime_data_dir, recursive = TRUE, full.names = TRUE, pattern = "\\.csv$")
5
6 # Function to read and combine multiple CSV files into one data frame
7 combine_csv_files <- function(files) {
8   map_df(files, read_csv, .id = "file_id")
9 }
10
11 # Combine all crime data files into one dataset
12 combined_data <- combine_csv_files(all_files) %>%
13   # Select relevant columns for crime analysis
14   select(`Crime ID`, `Month`, `Falls within`, `LSOA code`, `LSOA name`, `Crime type`)
15
16 # Clean the crime data
17 clean_data <- combined_data %>%
18   # Replace Crime ID with a unique identifier
19   mutate(`Crime ID` = row_number(), Month = ymd(paste0(Month, "-01"))) %>%
20   # Filter data to include only crimes in Bristol or Cornwall
21   filter(str_detect(`LSOA name`, "Bristol|Cornwall")) %>%
22   # Rename the Month column to date and assign County based on LSOA name
23   rename(date = Month) %>%
24   mutate(
25     `Falls within` = ifelse(
26       startsWith(`LSOA name`, "Bristol"),
27       "Bristol",
28       "Cornwall"
29     )
30   ) %>%
31   rename("County" = "Falls within", lsoa_code = `LSOA code`, City = `LSOA name`, Crime = `Crime type`)
32
33 # Save the cleaned crime data to a CSV file
34 write_csv(clean_data, file = paste0(getwd(), "/Cleaned_datasets/Crime/Crime_Data_Combined.csv"))

```

Figure 5: Cleaning code for school

```

1 # Load postcode dataset to enable joining tables and fetching street location information
2 postcode_dataset <- read_csv("Cleaned_datasets/Postcode_clean.csv")
3
4 # Load the 2021-2022 Bristol Key Stage 4 final data
5 bristol_key_stage_4_final_2021_2022 <- read_csv("Obtained_Data/schools/Bristol/2021-2022/801_ks4final.csv") %>%
6   # Select relevant columns
7   select(SCHNAME, PCODE, TOWN, ATT8SCR) %>%
8   # Add a Year column indicating the year of the data
9   mutate(Year = 2022)
10
11 # Load the 2022-2023 Bristol Key Stage 4 final data
12 bristol_key_stage_4_final_2022_2023 <- read_csv("Obtained_Data/schools/Bristol/2022-2023/801_ks4final.csv") %>%
13   # Select relevant columns
14   select(SCHNAME, PCODE, TOWN, ATT8SCR) %>%
15   # Add a Year column indicating the year of the data
16   mutate(Year = 2023)
17
18 # Load the 2021-2022 Cornwall Key Stage 4 final data
19 cornwall_key_stage_4_final_2021_2022 <- read_csv("Obtained_Data/schools/Cornwall/2021-2022/908_ks4final.csv") %>%
20   # Select relevant columns
21   select(SCHNAME, PCODE, TOWN, ATT8SCR) %>%
22   # Add a Year column indicating the year of the data
23   mutate(Year = 2022)
24
25 # Load the 2022-2023 Cornwall Key Stage 4 final data
26 cornwall_key_stage_4_final_2022_2023 <- read_csv("Obtained_Data/schools/Cornwall/2022-2023/908_ks4final.csv") %>%
27   # Select relevant columns
28   select(SCHNAME, PCODE, TOWN, ATT8SCR) %>%
29   # Add a Year column indicating the year of the data
30   mutate(Year = 2023)
31
32 # Combine all the school data into one dataset
33 combined_ks4_data <- bind_rows(
34   bristol_key_stage_4_final_2021_2022,
35   bristol_key_stage_4_final_2022_2023,
36   cornwall_key_stage_4_final_2021_2022,
37   cornwall_key_stage_4_final_2022_2023
38 )
39
40 # Filter out rows with missing school names (NA values)
41 clean_school_data <- combined_ks4_data %>%
42   filter(!is.na(SCHNAME)) %>% # Filter for rows where SCHNAME is not NA
43   # Convert ATT8SCR column to numeric (assuming it's a character vector)
44   mutate(ATT8SCR = as.numeric(ATT8SCR)) %>% # Convert ATT8SCR to numeric
45   # Impute missing values in ATT8SCR with the mean (excluding NAs)
46   mutate(ATT8SCR = ifelse(is.na(ATT8SCR), mean(ATT8SCR, na.rm = TRUE), ATT8SCR)) %>% # Replace NA with mean
47   # Remove spaces from PCODE column
48   mutate(PCODE = gsub(" ", "", PCODE)) %>% # Remove spaces in PCODE using gsub
49   # Inner join with postcode_clean data based on PCODE and pcd7
50   inner_join(postcode_dataset, by = c("PCODE" = "pcd7")) %>% # Join with postcode data
51   # Rename PCODE column to PostalCode
52   rename("PostalCode" = "PCODE") %>% # Rename PCODE to PostalCode
53   # Impute missing City and County with "NA"
54   mutate(City = ifelse(is.na(City), "NA", City), County = ifelse(is.na(County), "NA", County)) %>% # Replace NA with "NA"
55   # Select all columns except TOWN
56   select(-TOWN)
57
58 # Save the cleaned school data to a CSV file
59 write_csv(clean_school_data, file = paste0(getwd(), "/Cleaned_datasets/Schools/school_cleaned_data.csv"))
60

```

Figure 6: cleaning post code



```

1 #Read the CSV file containing postcode to LSOA mapping and select relevant columns
2 postcode_clean <- read_csv("Obtained_Data/Postcode to LSOA.csv") %>%
3   # Select only the necessary columns from the dataset
4   select(pcd7, oa11cd, lsoa11nm, msoa11nm, ladnm) %>%
5   # Rename the columns for clarity
6   rename(County = ladnm, City = msoa11nm, Street = lsoa11nm) %>%
7   # Remove spaces from the 'pcd7' column values for consistency
8   mutate(pcd7 = gsub(" ", "", pcd7)) %>%
9   # Standardize the 'County' column to ensure uniform naming
10  mutate(County = ifelse(County == "Bristol, City of", "Bristol", County)) %>%
11  # Filter the data to include only rows where the 'County' is either 'Bristol' or 'Cornwall'
12  filter((County %in% c("Bristol", "Cornwall")))
13
14 # Write the cleaned data to a new CSV file in the specified directory
15 write_csv(postcode_clean, file = paste0(getwd(), "/Cleaned_datasets/Postcode_clean.csv"))

```

Figure 7: Cleaning Population dataset



```

1 # Read the population data from the CSV file
2 population_dataset <- read_csv("Obtained_Data/Population2011_1656567141570.csv") %>%
3
4   # Extract the first two characters from the 'Postcode' column to match the format in the postcode dataset
5   mutate(Postcode = substr(Postcode, start = 1, stop = 2))
6
7 # Read the cleaned postcode data and select relevant columns
8 postcode_dataset <- read_csv("Cleaned_datasets/Postcode_clean.csv") %>%
9
10  # Select only the columns 'pcd7' (postcode) and 'County'
11  select(pcd7, County) %>%
12
13  # Extract the first two characters from the 'pcd7' column to match the format in the population dataset
14  mutate(pcd7 = substr(pcd7, start = 1, stop = 2)) %>%
15
16  # Remove duplicate rows based on 'pcd7' and 'County'
17  distinct()
18
19 # Merge the population data with the postcode data and calculate the total population by county
20 population_2011 <- population_dataset %>%
21
22  # Perform a left join to combine the population data with postcode data based on matching 'Postcode' and 'pcd7'
23  left_join(postcode_dataset, by = c("Postcode" = "pcd7")) %>%
24
25  # Group the data by 'County'
26  group_by(County) %>%
27
28  # Summarize the data by calculating the total population for each county
29  summarize(total_population = sum(Population, na.rm = TRUE))
30
31 # Write the cleaned and summarized population data to a new CSV file
32 write_csv(population_2011, file = paste0(getwd(), "/Cleaned_datasets/Population_clean.csv"))

```

Figure 8: Code to generate House Price Boxplot

```

1 # Load the cleaned house price dataset
2 house_dataset <- read_csv("Cleaned_datasets/House_Pricing_Data/clean_house_pricing_data.csv")
3
4 # Filter the dataset for transactions that occurred in 2023
5 house_data_2023 <- subset(house_dataset, format(as.Date(Date_of_Transfer), "%Y") == "2023")
6
7 # Generate a boxplot comparing house prices in Bristol and Cornwall for 2023
8 ggplot(house_data_2023, aes(x = County, y = log(Price), fill = County)) +
9   geom_boxplot() +
10  labs(
11    title = "Boxplot of House Prices in 2023 for CITY OF BRISTOL and CORNWALL",
12    x = "County",
13    y = "House Price"
14  )
15
16 # Save the boxplot as a PNG file
17 ggsave(paste0(getwd(), "/Graphs/House_Pricing/BoxP_BL_CL_2023.png"))

```

Figure 9: House Price BoxPlot

Boxplot of House Prices in 2023 for CITY OF BRISTOL and CORNWALL

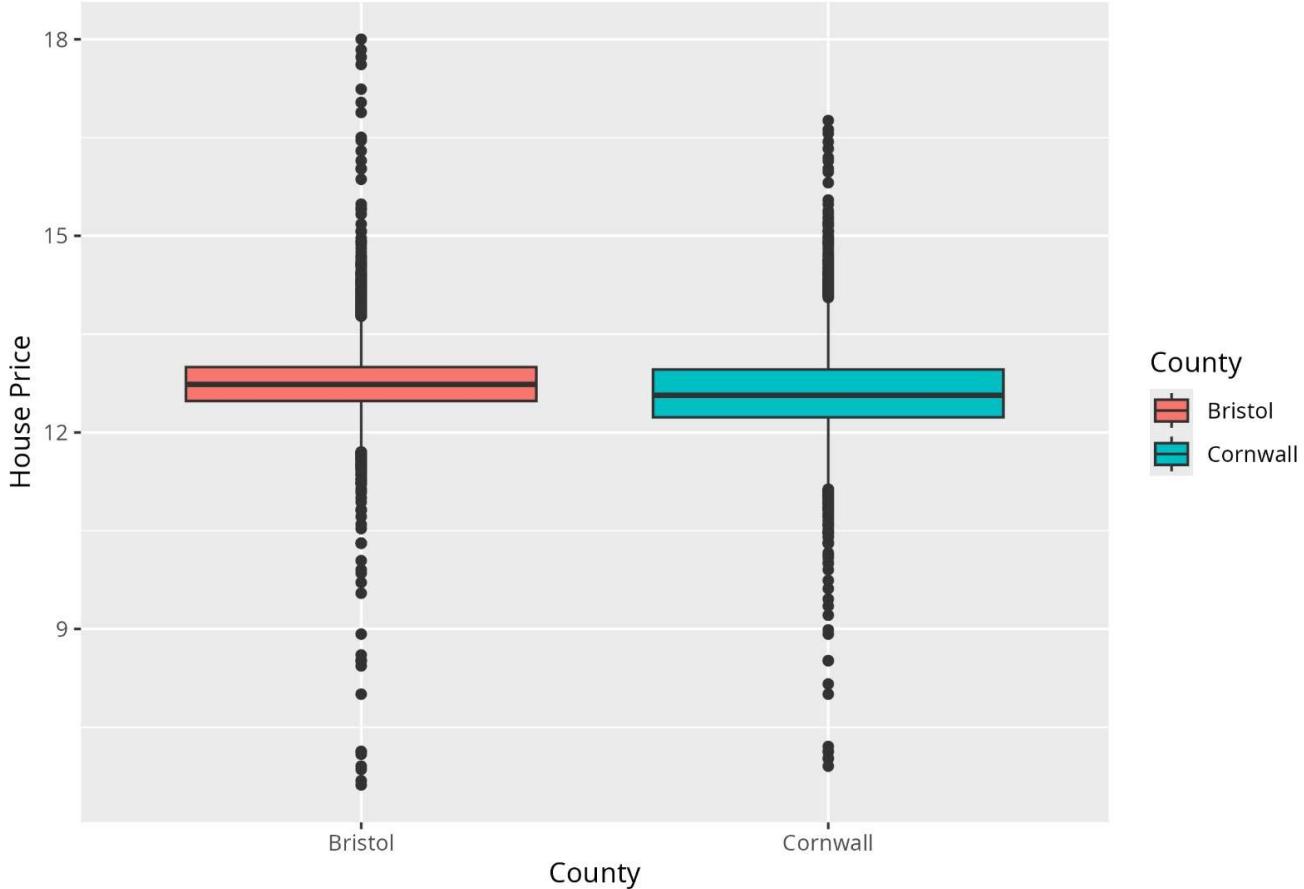


Figure 10: Code to generate Stacked BarChart for House Prices 2023



```

1 # Calculate the total house price for each town within each county in 2023
2 town_price_summary <- house_data_2023 %>%
3   group_by(County, Town.City) %>%
4   summarise(Total_Price = sum(Price))
5
6 # Create a stacked bar graph showing the total house prices by town within each county
7 ggplot(town_price_summary, aes(x = County, y = Total_Price, fill = Town.City)) +
8   geom_bar(stat = "identity") + # Stack bars based on the sum
9   labs(
10     title = "Total House Prices by Town within County (2023)",
11     x = "County",
12     y = "Total Price"
13   ) +
14   theme_classic() + # Optional: adjust plot aesthetics
15   guides(fill = guide_legend(title = "Town")) # Label the legend as "Town"
16
17 # Save the stacked bar graph as a PNG file
18 ggsave(paste0(getwd(), "/Graphs/House_Pricing/BarC_Stack_Town_BL_CL_2023.png"))

```

Figure 11: Stacked Barchart House Price (2023)

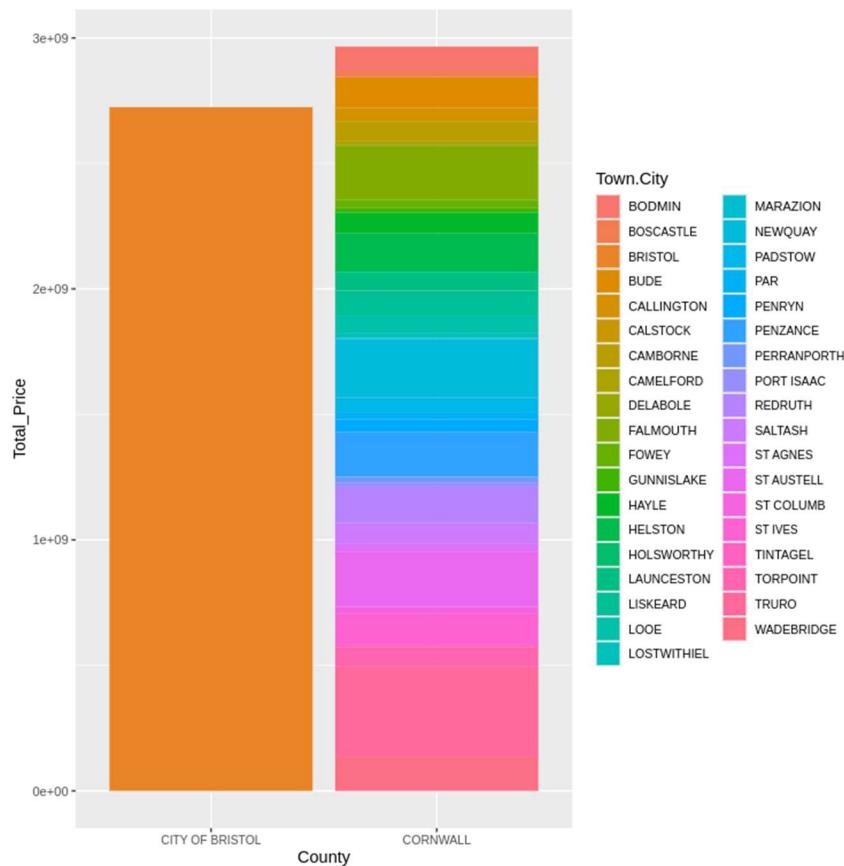


Figure 12: Code to generate Line graph for house price by town

```

1 avg_prices <- house_dataset %>%
2   group_by(Town.City) %>%
3   summarize(avg_price = mean(Price, na.rm = TRUE)) %>%
4   arrange(desc(avg_price))
5
6 # Create a line graph showing the average house price by town
7 ggplot(avg_prices, aes(x = reorder(Town.City, -avg_price), y = avg_price)) +
8   geom_line(group = 1) +
9   geom_point() +
10  theme_minimal() +
11  theme(axis.text.x = element_text(angle = 90, hjust = 1)) +
12  labs(
13    title = "Average House Prices by Town",
14    x = "Town/City",
15    y = "Average Price"
16  )

```

Figure 13: Line Graph for house prices by town

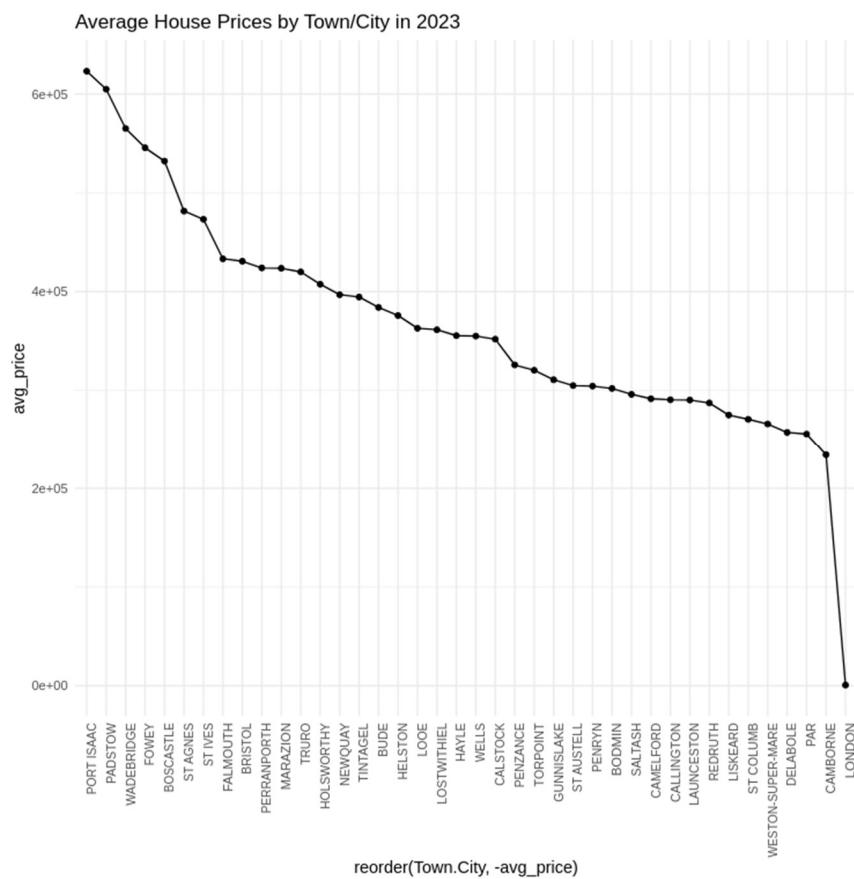


Figure 14: Code to generate Boxplot for Average download speed by county

```
1 # Load the cleaned broadband performance dataset
2 broadband_performance_dataset <- read_csv("Cleaned_datasets/Broadband/clean_broadband_performance.csv")
3
4 # Generate a boxplot of average download speeds by county
5 ggplot(broadband_performance_dataset, aes(x = county, y = `Average download speed (Mbit/s)`, fill = county)) +
6   geom_boxplot() + # Create a boxplot to show distribution of average download speeds
7   labs(
8     title = "Average Download Speed by County", # Title of the plot
9     x = "County", # Label for the x-axis
10    y = "Average Download Speed (Mbit/s)" # Label for the y-axis
11  )
12
13 # Save the boxplot as a PNG file in the specified directory
14 ggsave(paste0(getwd(), "/Graphs/Broadband/average_down_speed_boxplots.png"))
15
```

Figure 15: Box plot for average internet speed by County

Average Download Speed by County

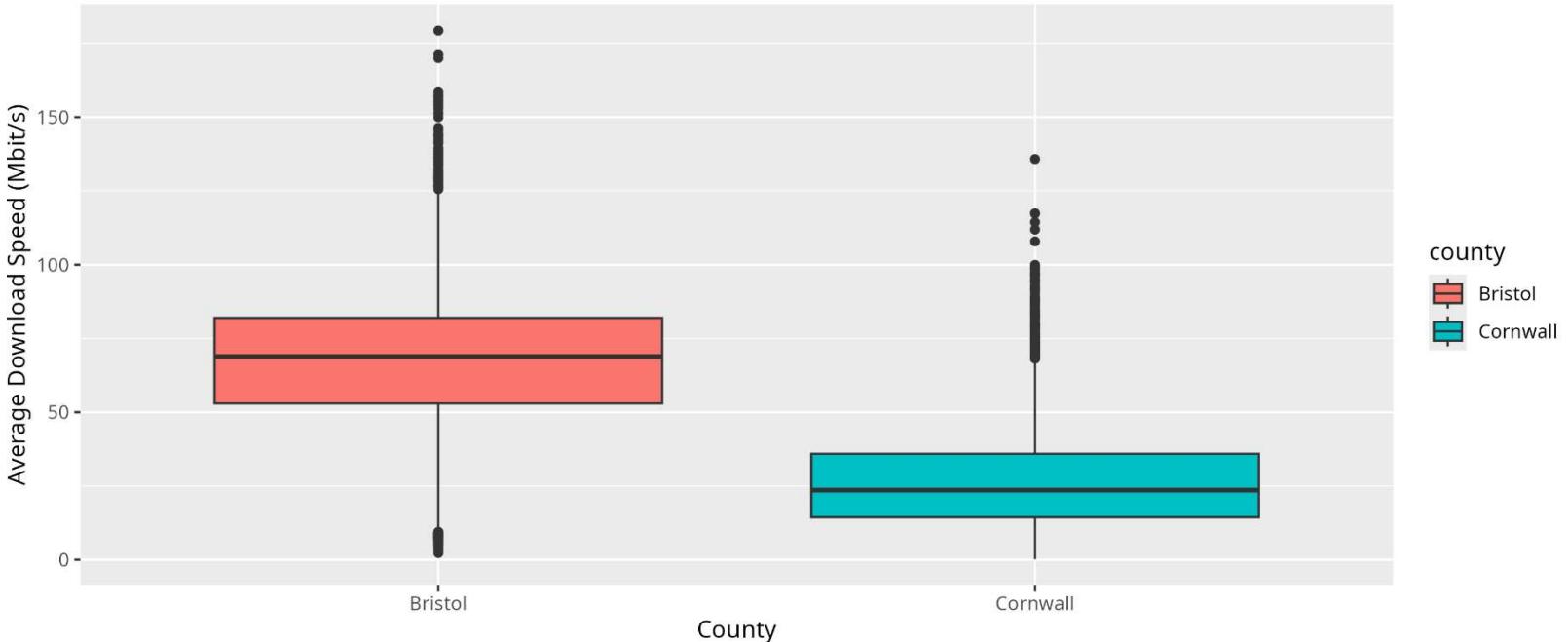


Figure 17: Code to generate dodge bar chart of maximum and average download speed by City of Bristol

```

● ● ●
1 # Filter the dataset to include only data for Bristol
2 bristol_data <- broadband_performance_dataset %>%
3   filter(county == "Bristol")
4
5 # Filter the dataset to include only data for Cornwall
6 cornwall_data <- broadband_performance_dataset %>%
7   filter(county == "Cornwall")
8
9 # Reshape the Bristol data into a long format for easier plotting
10 long_bristol <- tidyverse::gather(bristol_data, key = "speed_type", value = "speed", `Maximum download speed (Mbit/s)`, `Average download speed (Mbit/s)` )
11
12 # Reshape the Cornwall data into a long format for easier plotting
13 long_cornwall <- tidyverse::gather(cornwall_data, key = "speed_type", value = "speed", `Maximum download speed (Mbit/s)`, `Average download speed (Mbit/s)` )
14
15 # Create a side-by-side bar chart for Bristol showing maximum vs average download speeds by city
16 ggplot(long_bristol, aes(x = city, y = speed, fill = speed_type)) +
17   geom_bar(stat = "identity", position = "dodge") + # Use bars to show the values and dodge to place bars side by side
18   labs(
19     title = "Bristol: Maximum vs Average Download Speed by City", # Title of the plot
20     x = "City", # Label for the x-axis
21     y = "Download Speed (Mbit/s)" # Label for the y-axis
22   ) +
23   theme(axis.text.x = element_text(angle = 45, hjust = 1)) # Rotate x-axis labels for better readability
24
25 # Save the Bristol bar chart as a PNG file in the specified directory
26 ggsave(paste0(getwd(), "/Graphs/Broadband/bristol_down_median_bar_chart.png"))

```

Figure 16: Dodge Barchart for Average vs Maximum Download Speed in Bristol

Bristol: Maximum vs Average Download Speed by City

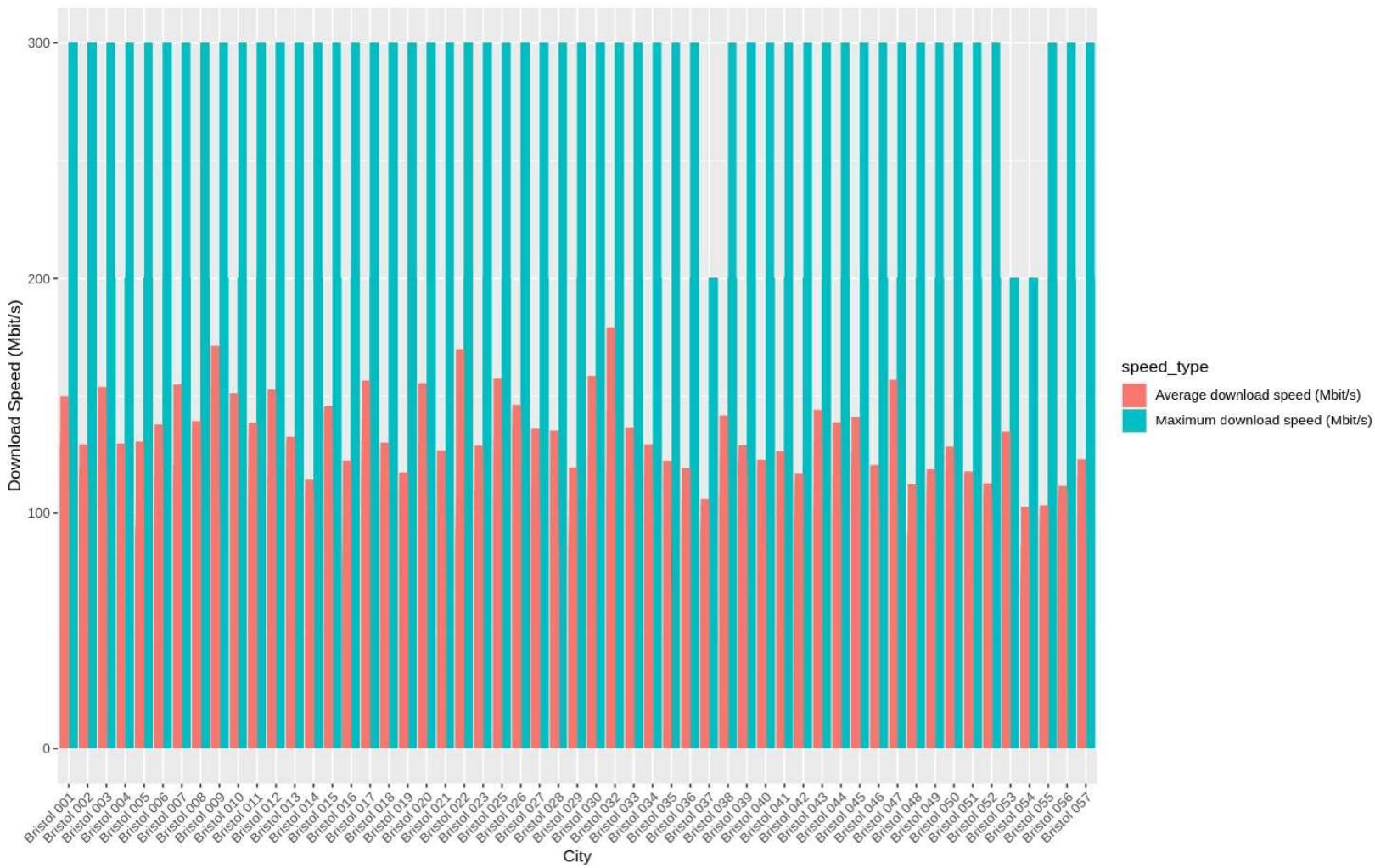


Figure 18: Figure 16: Code to generate Dodge Barchart for Average vs Maximum Download Speed in Cornwall

```

1 # Create a side-by-side bar chart for Cornwall showing maximum vs average download speeds by city
2 ggplot(long_cornwall, aes(x = city, y = speed, fill = speed_type)) +
3   geom_bar(stat = "identity", position = "dodge") + # Use bars to show the values and dodge to place bars side by side
4   labs(
5     title = "Cornwall: Maximum vs Average Download Speed by City", # Title of the plot
6     x = "City", # Label for the x-axis
7     y = "Download Speed (Mbit/s)" # Label for the y-axis
8   ) +
9   theme(axis.text.x = element_text(angle = 45, hjust = 1)) # Rotate x-axis labels for better readability
10
11 # Save the Cornwall bar chart as a PNG file in the specified directory
12 ggsave(paste0(getwd(), "/Graphs/Broadband/cornwall_down_median_bar_chart.png"))
13

```

Figure 19:Dodge bar chart for Maximum vs Average Download Speed in Cornwall

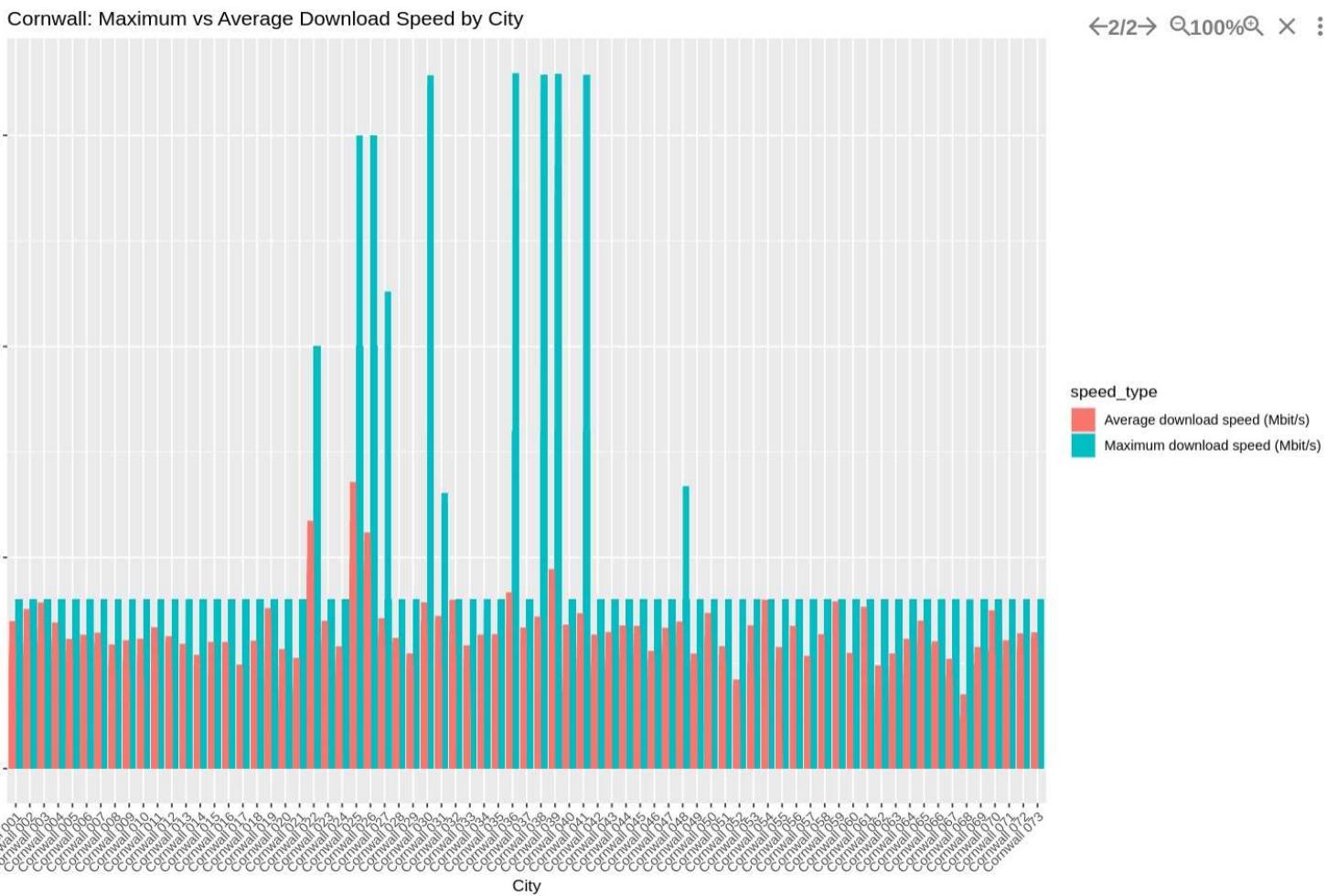


Figure 20: Code to generate boxplot of Drug Offence rate per 10000 people 2023



```

1  Filter drug offences data for the year 2023
2  drug_offences_2023 <- subset(drug_offences_2020_2023, format(as.Date(date), "%Y") == "2023")
3
4  # Calculate drug offence rates for 2023
5  drug_offences_rate_2023 <- drug_offences_2023 %>%
6    filter(Crime == "Drugs") %>% # Ensure the crime type is "Drugs"
7    mutate(month = floor_date(ymd(date), "month")) %>% # Extract month from the date
8    group_by(month, County) %>% # Group by month and county
9    summarize(
10      total_offences = n(), # Count total drug offences
11      offence_rate = (total_offences / population_2023[County]) * 10000 # Calculate offence rate per 10,000 population
12    )
13
14 # Create a boxplot to visualize drug offence rates in 2023
15 ggplot(drug_offences_rate_2023, aes(x = County, y = offence_rate, fill = County)) +
16   geom_boxplot() + # Draw boxplots
17   labs(
18     title = "Drug Offence Rate (2023)", # Title of the plot
19     x = "Region", # Label for the x-axis
20     y = "Offence Rate " # Label for the y-axis
21   )
22
23 # Save the boxplot as a PNG file
24 ggsave(paste0(getwd(), "/Graphs/Crime/BoxP_Drug_BL_CL_2023.png"))
25

```

Figure 21: Box plot for drug offence (2023)

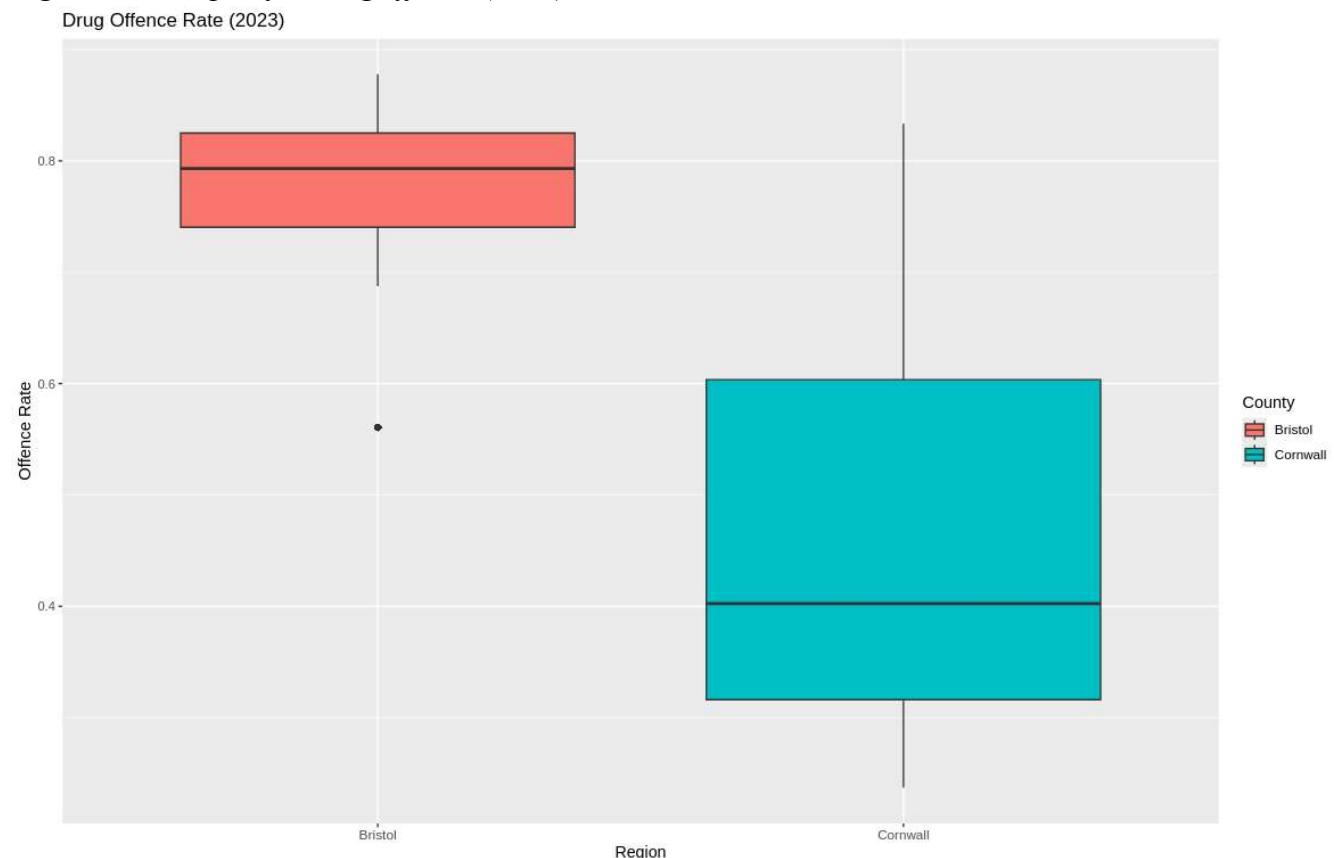


Figure 22: Code to generate Line graph for drug offence rate 2020-2023

```

1 #Load population dataset to access data on population around Bristol and Cornwall
2 population_2011 <- read_csv("Cleaned_datasets/Population_clean.csv") %>%
3   # Convert the data frame to a named vector
4   deframe()
5
6 # Calculate the total population for Bristol and Cornwall combined
7 total_population <- sum(population_2011[c("Bristol", "Cornwall")], na.rm = TRUE)
8
9 # Add the combined total to the named vector with the name "Total"
10 population_2011 <- c(population_2011, Total = total_population)
11
12 # Remove any NA entries from the named vector
13 population_2011 <- population_2011[!is.na(names(population_2011))]
14
15 # Estimate the population for 2023 using a growth factor
16 population_2023 <- floor((1.00561255390388033 * population_2011))
17
18 #####
19
20 # Load crime data
21 crime_dataset <- read_csv("Cleaned_datasets/Crime/Crime_Data_Combined.csv")
22
23 # Filter crime data for drug offences between 2020 and 2023
24 drug_offences_2020_2023 <- crime_dataset %>%
25   filter(Crime == "Drugs", year(date) >= 2020 & year(date) <= 2023)
26
27 # Calculate the drug offence rate per 10,000 population
28 drug_offences_rate_2020_2023 <- drug_offences_2020_2023 %>%
29   # Extract month and year from the date for grouping
30   mutate(
31     month = floor_date(ymd(date), "month"),
32     year = year(date)
33   ) %>%
34   # Group by year and county, then count the total number of offences
35   group_by(year, County) %>%
36   summarize(
37     total_offences = n(),
38     .groups = "drop"
39   ) %>%
40   # Join with population data to get population numbers for each county
41   left_join(
42     tibble(
43       County = names(population_2023),
44       population = as.numeric(population_2023)
45     ),
46     by = "County"
47   ) %>%
48   # Calculate the offence rate per 10,000 population
49   mutate(offence_rate = (total_offences / population) * 10000)
50
51 # Create a line graph to visualize drug offence rates by year
52 ggplot(drug_offences_rate_2020_2023, aes(x = year, y = offence_rate, color = County, group = County)) +
53   geom_line() + # Add lines to connect points
54   geom_point(size = 3) + # Add points to the line graph
55   labs(
56     title = "Yearly Drug Offence Rate per 10000 Population (2020-2023)", # Title of the plot
57     x = "Year", # Label for the x-axis
58     y = "Offence Rate per 10000 Population", # Label for the y-axis
59     color = "Region" # Label for the color legend
60   ) +
61   theme_minimal() + # Use a minimal theme for the plot
62   theme(
63     legend.position = "bottom" # Position the legend at the bottom
64   ) +
65   scale_x_continuous(breaks = 2020:2023) # Set x-axis breaks to display years 2020 to 2023
66
67 # Save the line graph as a PNG file in the specified directory
68 ggsave(paste0(getwd(), "/Graphs/Crime/LineG_Drug_10000_2023.png"))

```

Figure 23: Line graph for Yearly drug offence rate 2020-2023

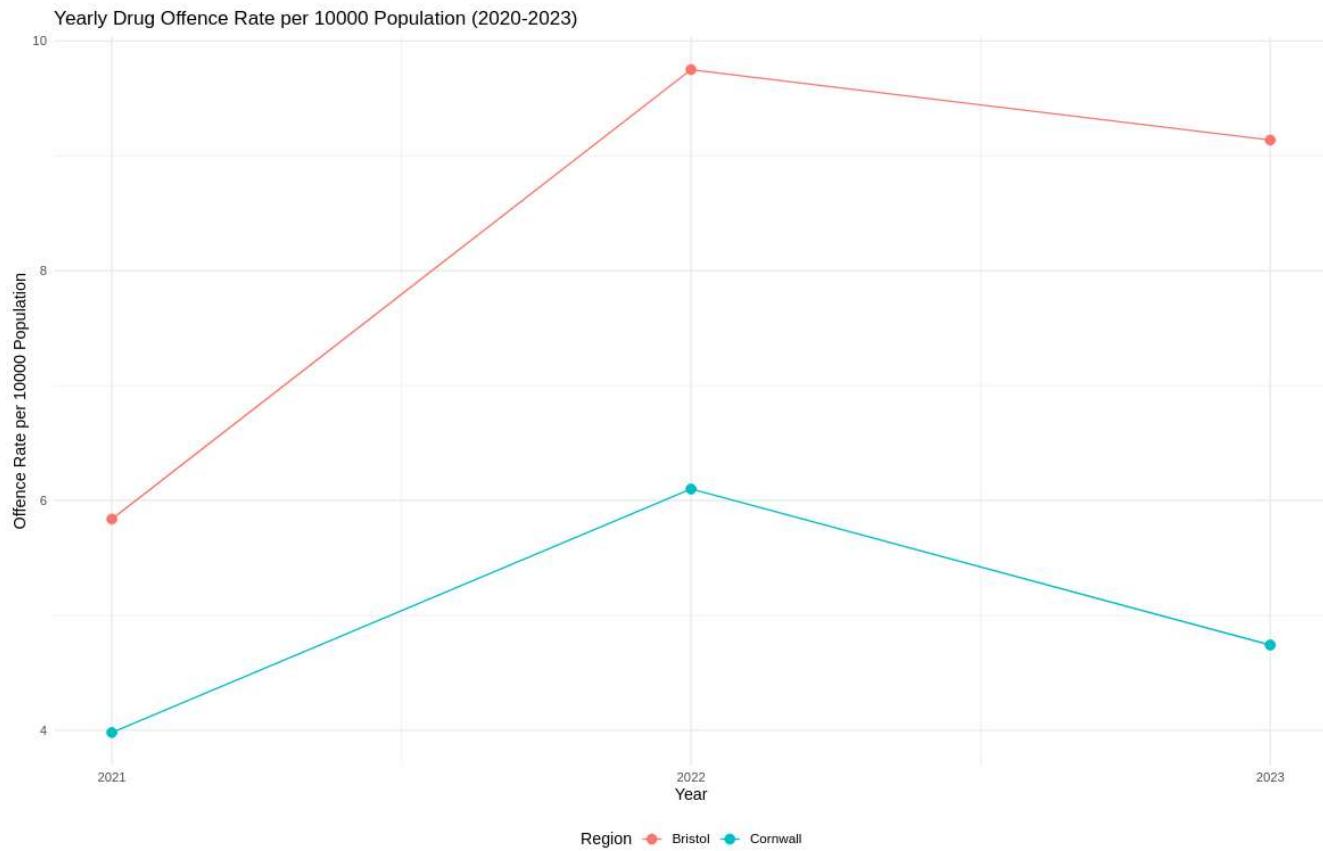


Figure 25: Code to generate Piechart for robberies on October 2022

```

1 robbery_offences_oct_2023 <- crime_dataset %>%
2   filter(
3     Crime == "Robbery", # Filter for robbery offences
4     year(date) == 2022, # Only data from the year 2022
5     month(date) == 10 # Only data from October
6   )
7
8 # Calculate robbery rates for October 2022
9 robbery_rates <- robbery_offences_oct_2023 %>%
10  group_by(County) %>% # Group by county
11  summarize(
12    total_offences = n(), # Count total robbery offences
13    .groups = "drop" # Drop the grouping after summarizing
14  ) %>%
15  # Join with population data for each county
16  left_join(
17    tibble(
18      County = names(population_2023),
19      population = as.numeric(population_2023)
20    ),
21    by = "County"
22  ) %>%
23  mutate(
24    offence_rate = (total_offences / population) * 10000, # Calculate offence rate per 10,000 population
25    percentage = total_offences / sum(total_offences) * 100
26  ) # Calculate percentage of total offences
27
28
29 # Create a pie chart to visualize the distribution of robbery offences in October 2022
30 ggplot(robbery_rates, aes(x = "", y = percentage, fill = County)) +
31  geom_bar(stat = "identity", width = 1) + # Use bars to represent proportions
32  coord_polar("y", start = 0) + # Convert bar chart to a pie chart
33  geom_text(aes(label = sprintf("%.1f%%\n(%d)", percentage, total_offences)), # Add labels with percentage and count
34    position = position_stack(vjust = 0.5)
35  ) +
36  labs(
37    title = "Robbery Offences Distribution (October 2022)", # Title of the plot
38    subtitle = paste("Total offences:", sum(robbery_rates$total_offences)), # Subtitle with total offences
39    fill = "Region" # Label for the fill legend
40  ) +
41  theme_void() + # Use a void theme to remove background and axes
42  theme(legend.position = "bottom") # Position the legend at the bottom
43
44 # Save the pie chart as a PNG file
45 ggsave(paste0(getwd(), "/Graphs/Crime/Pi_Robbery_oct_2022.png"))
46

```

Figure 24: Pie chart for robbery rate on October 2022

Robbery Offences Distribution (October 2022)
Total offences: 100

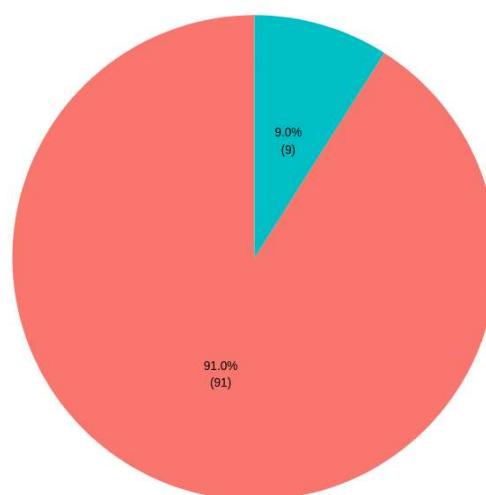


Figure 26: Radar chart for vehicle crime rates from 2020-2023

```

● ● ●
1 population_2011 <- read_csv("Cleaned_datasets/Population_clean.csv") %>%
2   deframe()
3
4 total_population <- sum(population_2011[c("Bristol", "Cornwall")], na.rm = TRUE)
5 # Add the total to the list with the name "Total"
6 population_2011 <- c(population_2011, Total = total_population)
7 # Remove the NA entry if it exists
8 population_2011 <- population_2011[!is.na(names(population_2011))]
9
10 population_2023 <- floor((1.00561255390388033 * population_2011))
11
12
13
14 crime_dataset <- read_csv("Cleaned_datasets/Crime/Crime_Data_Combined.csv")
15
16 vehicle_crime_2020_2023 <- crime_dataset %>%
17   filter(Crime == "Vehicle crime", year(date) >= 2020 & year(date) <= 2023)
18
19 vehicle_crime_rate_2020_2023 <- vehicle_crime_2020_2023 %>%
20   mutate(
21     year = year(date)
22   ) %>%
23   group_by(year, County) %>%
24   summarize(
25     total_offences = n(),
26     .groups = "drop"
27   ) %>%
28   left_join(
29     tibble(
30       County = names(population_2023),
31       population = as.numeric(population_2023)
32     ),
33     by = "County"
34   ) %>%
35   mutate(offence_rate = (total_offences / population) * 10000) %>%
36   select(year, County, offence_rate)
37
38 # Reshape the data for ggradar
39 vehicle_crime_rate_wide <- vehicle_crime_rate_2020_2023 %>%
40   pivot_wider(names_from = year, values_from = offence_rate, names_prefix = "Year_") %>%
41   mutate(across(starts_with("Year_"), ~replace_na(., 0)))
42
43 # Create radar chart
44 ggradar(vehicle_crime_rate_wide,
45   values.radar = c("0", "5", "10"),
46   grid.min = 0,
47   grid.max = 10,
48   grid.mid = 5,
49   group.point.size = 3,
50   group.line.width = 1.5,
51   axis.label.size = 3,
52   legend.text.size = 10,
53   legend.position = "bottom",
54   background.circle.colour = "white",
55   axis.line.colour = "gray60",
56   gridline.mid.colour = "gray60") +
57   labs(title = "Vehicle Crime Rate per 10,000 Population (2020-2023)")
58

```

Figure 27: Radar chart for yearly vehicle crimes 2020-2023

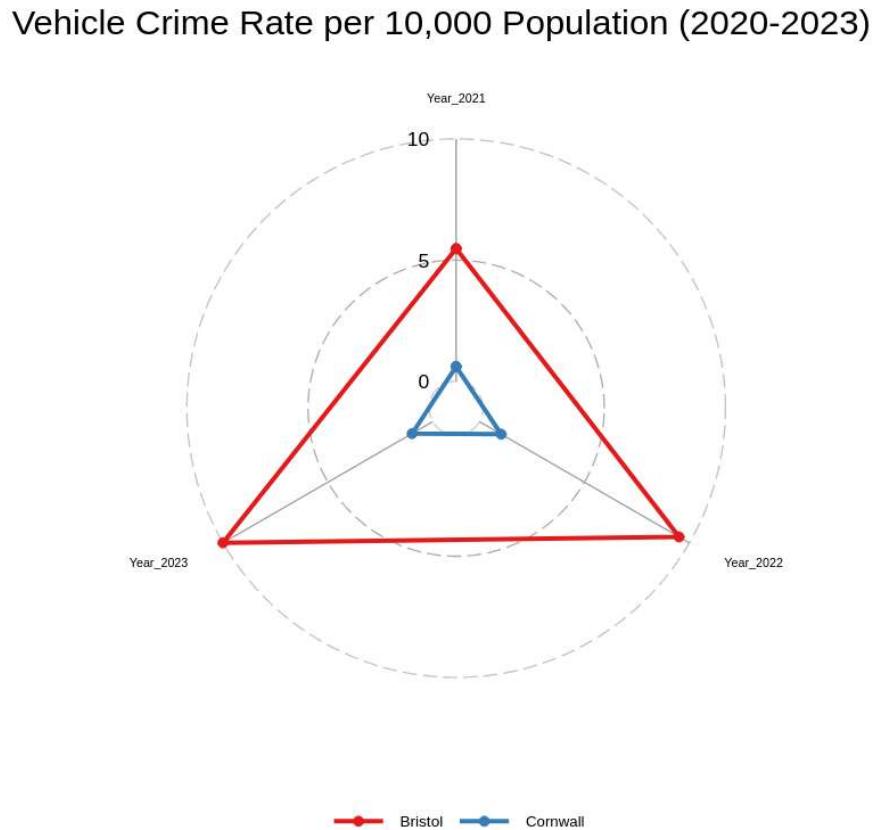


Figure 28: Code to generate Boxplot of attainment 8 score 2023

```

1  Read the cleaned school dataset and filter out data for Wiltshire
2  school_dataset <- read_csv("Cleaned_datasets/Schools/Schools_Clean.csv") %>%
3    filter(County != "Wiltshire")
4
5  # Subset the dataset to include only records for the year 2023
6  school_data_2023 <- subset(school_dataset, Year == 2023)
7
8  # Create a boxplot of ATT8SCR (Attainment 8 Score) by County for the year 2023
9  ggplot(school_data_2023, aes(x = County, y = ATT8SCR, fill = County)) +
10   geom_boxplot() + # Use a boxplot to visualize the distribution of ATT8SCR by County
11   labs(
12     title = "Boxplot of ATT8SCR by County",
13     x = "County",
14     y = "ATT8SCR"
15   ) + # Label the plot
16   theme(axis.text.x = element_text(hjust = 1)) # Adjust text alignment for x-axis labels
17
18 # Save the boxplot to a file
19 ggsave(paste0(getwd(), "/Graphs/Schools/BoxP_avg_att8_BL_CL_2023.png"))

```

Figure 29: Attainment 8 score by county 2023

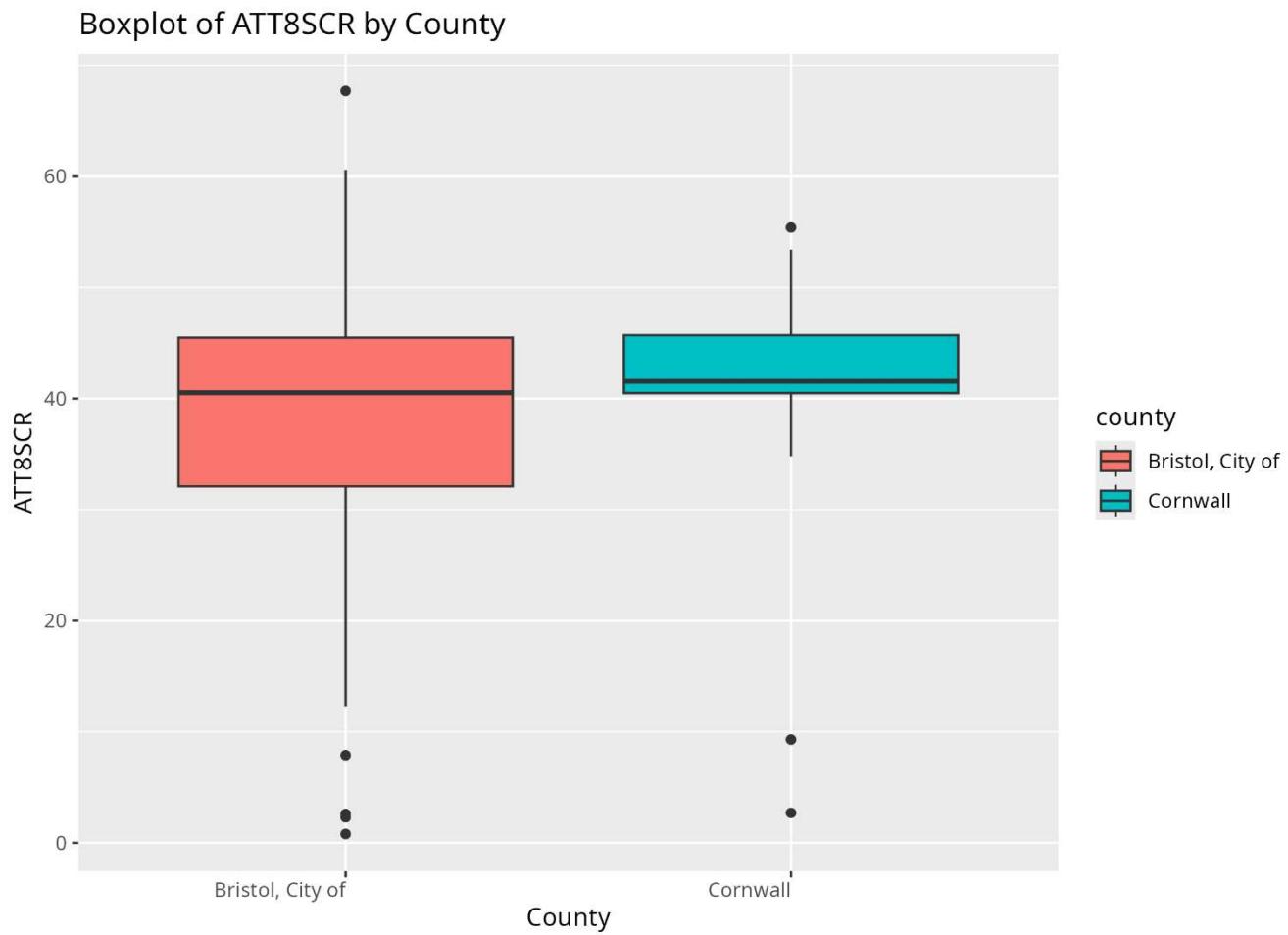


Figure 30: Code to lime graph of attainment 8 score city wise in Bristol and Cornwall

```
1 cornwall_data <- school_dataset %>%
2   filter(County == "Cornwall")
3
4 # Filter data for Bristol
5 bristol_data <- school_dataset %>%
6   filter(County == "Bristol")
7
8 # Create a line graph of average Attainment 8 Score by City in Cornwall over the years
9 ggplot(cornwall_data, aes(x = Year, y = ATT8SCR, color = City)) +
10   geom_line() + # Use a line graph to show the trend of ATT8SCR by City
11   labs(
12     title = "Average Attainment 8 Score by City in Cornwall",
13     x = "Year",
14     y = "Average Attainment 8 Score"
15   ) # Label the plot
16
17 # Save the boxplot to a file
18 ggsave(paste0(getwd(), "/Graphs/Schools/line_avg_att8_CL.png"))
19
20 # Create a line graph of average Attainment 8 Score by City in Bristol over the years
21 ggplot(bristol_data, aes(x = Year, y = ATT8SCR, color = City)) +
22   geom_line() + # Use a line graph to show the trend of ATT8SCR by City
23   labs(
24     title = "Average Attainment 8 Score by City in Bristol",
25     x = "Year",
26     y = "Average Attainment 8 Score"
27   ) # Label the plot
28
29 # Save the boxplot to a file
30 ggsave(paste0(getwd(), "/Graphs/Schools/line_avg_att8_BL.png"))
```

Figure 31: Line graph to see average attainment 8 score in Cornwall

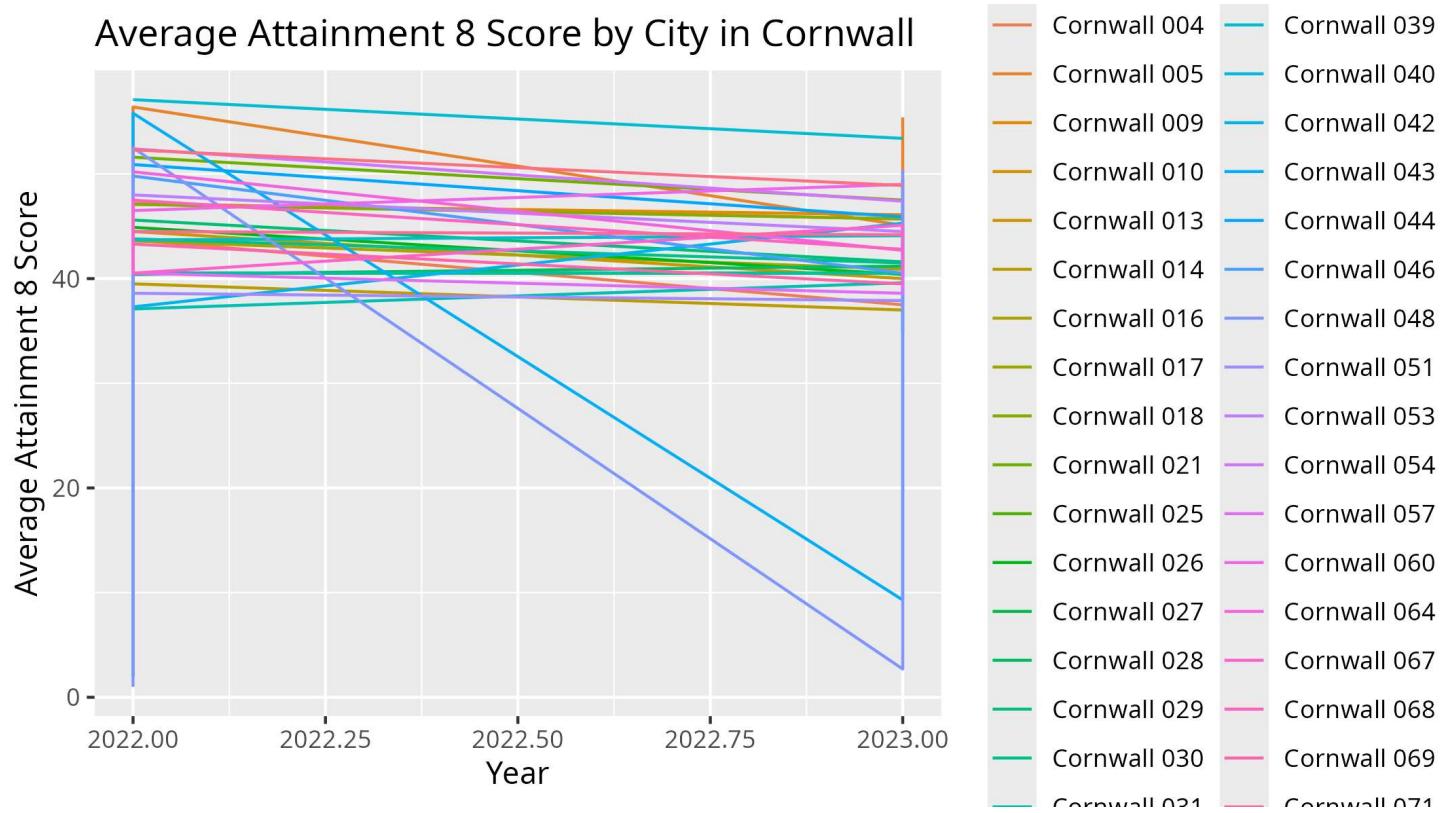


Figure 32: Line graph to see average attainment 8 score in Bristol

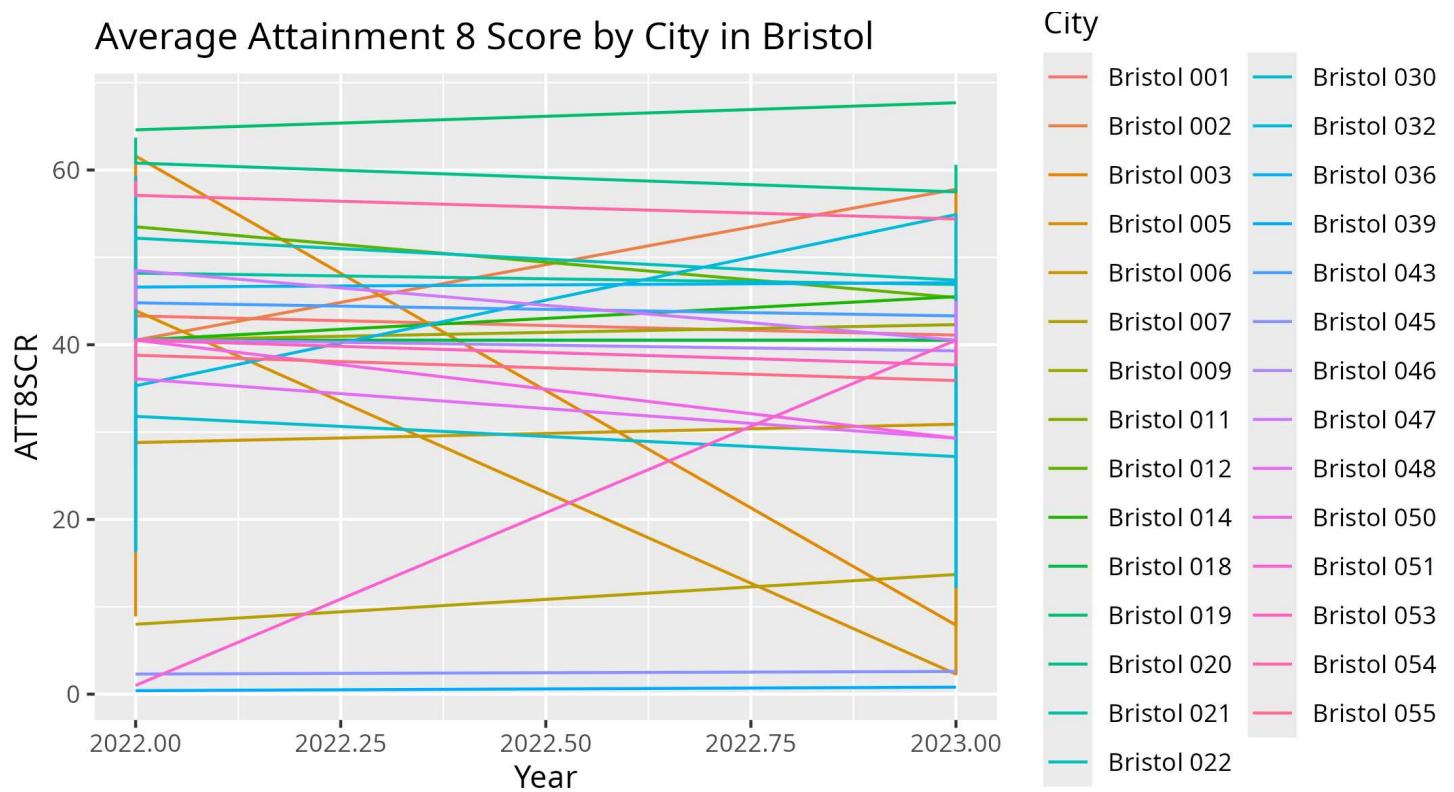


Figure 33: Code to perform linear modelling of attainment 8 score and house price

```

1 school_dataset <- read_csv("Cleaned_datasets/Schools/Schools_Clean.csv")
2 house_dataset <- read_csv("Cleaned_datasets/House_Pricing_Data/clean_house_pricing_data.csv")
3
4
5 population_2011 <- read_csv("Cleaned_datasets/Population_clean.csv") %>%
6   deframe()
7
8 total_population <- sum(population_2011[c("Bristol", "Cornwall")], na.rm = TRUE)
9 # Add the total to the list with the name "Total"
10 population_2011 <- c(population_2011, Total = total_population)
11 # Remove the NA entry if it exists
12 population_2011 <- population_2011[!is.na(names(population_2011))]
13
14
15 population_2023 <- floor((1.00561255390388033 * population_2011))
16
17 merged_dataset <- merge(
18   school_dataset,
19   house_dataset,
20   by = "PostalCode", # Use the updated column name for matching
21   all.x = TRUE, # Include all rows from the school_dataset even if there is no matching data in house_dataset
22   all.y = FALSE # Exclude rows from house_dataset if there is no matching row in school_dataset
23 ) %>%
24   filter(!is.na(ATT8SCR), !is.na(Price))
25
26 model <- lm(ATT8SCR ~ Price, data = merged_dataset)
27
28 # Check the model summary
29 summary(model)

```

Figure 34: Code to perform Linear Modelling between Attainment 8 score and average download speed

```

1 school_dataset <- read_csv("Cleaned_datasets/Schools/Schools_Clean.csv")
2 broadband_performance_dataset <- read_csv("Cleaned_datasets/Broadband/clean_broadband_performance.csv")
3
4 merged_dataset <- merge(
5   broadband_performance_dataset,
6   school_dataset,
7   by = "County", # Use the updated column name for matching
8   all.x = TRUE, # Include all rows from the school_dataset even if there is no matching data in house_dataset
9   all.y = FALSE # Exclude rows from house_dataset if there is no matching row in school_dataset
10 ) %>%
11   filter(!is.na(`Average download speed (Mbit/s)`), !is.na(ATT8SCR))
12
13 # Fit a linear model
14 model <- lm(ATT8SCR ~ `Average download speed (Mbit/s)`, data = merged_dataset)
15
16 # Check the model summary
17 summary(model)

```

Figure 35: Code to perform linear modeling between drug offence rate per 1000 people and average download speed

```

1 broadband_dataset <- read_csv("Cleaned_datasets/Broadband/clean_broadband_performance.csv")
2
3 drug_offences_2020_2023 <- read_csv("Cleaned_datasets/Crime/Crime_Data_Combined.csv") %>%
4   filter(Crime == "Drugs", year(date) >= 2020 & year(date) <= 2023)
5
6 population_dataset <- read_csv("Obtained_Data/Population2011_1656567141570.csv") %>%
7   mutate(Postcode = substr(Postcode, start = 1, stop = 2))
8
9 postcode_dataset <- read_csv("Cleaned_datasets/Postcode_clean.csv") %>%
10  select(pcd7, County) %>%
11  mutate(pcd7 = substr(pcd7, start = 1, stop = 2)) %>%
12  distinct()
13
14 population_2011 <- read_csv("Cleaned_datasets/Population_clean.csv") %>%
15  deframe()
16
17 total_population <- sum(population_2011[c("Bristol", "Cornwall")], na.rm = TRUE)
18 # Add the total to the list with the name "Total"
19 population_2011 <- c(population_2011, Total = total_population)
20 # Remove the NA entry if it exists
21 population_2011 <- population_2011[!is.na(names(population_2011))]
22
23
24 population_2023 <- floor((1.00561255390388033 * population_2011))
25
26 drug_offences_rate_2020_2023 <- drug_offences_2020_2023 %>%
27   # Extract month and year from the date for grouping
28   mutate(
29     month = floor_date(ymd(date), "month"),
30     year = year(date)
31   ) %>%
32   # Group by year and county, then count the total number of offences
33   group_by(year, City, County) %>%
34   summarize(
35     total_offences = n(),
36     .groups = "drop"
37   ) %>%
38   # Join with population data to get population numbers for each county
39   left_join(
40     tibble(
41       County = names(population_2023),
42       population = as.numeric(population_2023)
43     ),
44     by = "County"
45   ) %>%
46   # Calculate the offence rate per 10,000 population
47   mutate(offence_rate = (total_offences / population) * 10000) %>%
48   rename("Street" = "City")
49
50
51 # Merge datasets based on common columns, such as post codes or county
52 # Here, we'll merge on street as an example
53 merged_dataset <- merge(
54   broadband_dataset,
55   drug_offences_rate_2020_2023,
56   by = "Street",
57   all.x = TRUE,
58   all.y = FALSE
59 ) %>%
60   na.omit()
61
62 # Fit a linear model
63 model <- lm(offence_rate ~ `Average download speed (Mbit/s)`, data = merged_dataset)
64
65 # Check the model summary
66 summary(model)

```

Figure 36: Code to perform linear modelling between house prices and drug offence rate in 2023

```

1 house_dataset_2023 <- read_csv("Cleaned_datasets/House_Pricing_Data/clean_house_pricing_data.csv") %>%
2   filter(format(as.Date(Date), "%Y") == "2023")
3
4 drug_offences_2023 <- read_csv("Cleaned_datasets/Crime/Crime_Data_Combined.csv") %>%
5   filter(Crime == "Drugs", year(date) == 2023)
6
7 population_2011 <- read_csv("Cleaned_datasets/Population_clean.csv") %>%
8   deframe()
9
10 total_population <- sum(population_2011[c("Bristol", "Cornwall")], na.rm = TRUE)
11 # Add the total to the list with the name "Total"
12 population_2011 <- c(population_2011, Total = total_population)
13 # Remove the NA entry if it exists
14 population_2011 <- population_2011[!is.na(names(population_2011))]
15
16
17 population_2023 <- floor((1.00561255390388033 * population_2011))
18
19 drug_offences_rate_2023 <- drug_offences_2023 %>%
20   # Extract month and year from the date for grouping
21   mutate(
22     month = floor_date(ymd(date), "month"),
23     year = year(date)
24   ) %>%
25   # Group by year and county, then count the total number of offences
26   group_by(year, City, County) %>%
27   summarize(
28     total_offences = n(),
29     .groups = "drop"
30   ) %>%
31   # Join with population data to get population numbers for each county
32   left_join(
33     tibble(
34       County = names(population_2023),
35       population = as.numeric(population_2023)
36     ),
37     by = "County"
38   ) %>%
39   # Calculate the offence rate per 10,000 population
40   mutate(offence_rate = (total_offences / population)) %>%
41   rename("Street" = "City")
42
43 merged_data <- left_join(house_dataset_2023, drug_offences_rate_2023, by = c("Street", "County")) %>%
44   na.omit()
45
46
47
48 model <- lm(Price ~ offence_rate, data = merged_data)
49
50 # Summary of the model to view details
51 summary(model)

```

Figure 37: Code to perform linear modelling of house price and Average download speed

```

1 house_dataset <- read_csv("Cleaned_datasets/House_Pricing_Data/clean_house_pricing_data.csv")
2 broadband_performance_dataset <- read_csv("Cleaned_datasets/Broadband/clean_broadband_performance.csv")
3
4 merged_data <- left_join(house_dataset, broadband_performance_dataset, by = "PostalCode") %>%
5   na.omit()
6
7
8 model <- lm(Price ~ `Average download speed (Mbit/s)`, data = merged_data)
9
10 # Summary of the model to view details
11 summary(model)

```

Figure 38: Code to generate line of best fit between download speed and House Price

```

1 ggplot(merged_data, aes(x = `Average download speed (Mbit/s)`, y = Price)) +
2   geom_point(color = "blue") + # Scatter plot of the data
3   scale_y_continuous(trans = "log10") + # Plot the data points
4   geom_smooth(method = "lm", color = "red", se = FALSE) + # Line of best fit
5   labs(
6     title = "House Price vs Average Download Speed",
7     x = "Average Download Speed (Mbit/s)",
8     y = "House Price"
9   ) +
10   theme_minimal()

```

Figure 39: Line of best fit between download speed and house price

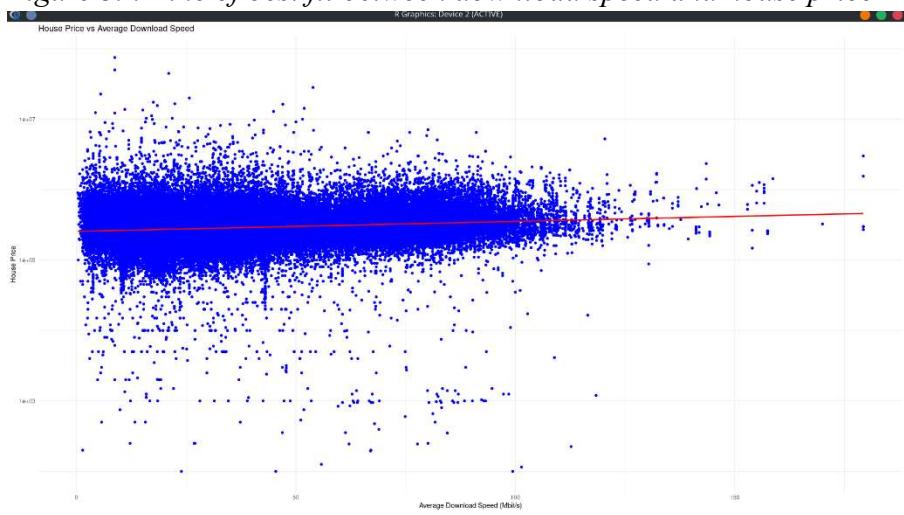


Figure 40: Code to generate line of best fit between drug offence rate and house price 2023

```

1 ggplot(merged_data, aes(x = offence_rate, y = Price)) +
2   geom_point(color = "blue") + # Scatter plot of the data
3   geom_smooth(method = "lm", color = "red", se = FALSE) +
4   scale_y_continuous(trans = "log10") + # Plot the data points# Line of best fit
5   labs(
6     title = "House Price vs Drug Offence Rate (2023)",
7     x = "Offence Rate",
8     y = "House Price"
9   ) +
10  theme_minimal()

```

Figure 41: line of best fit between drug offence rate and house price 2023

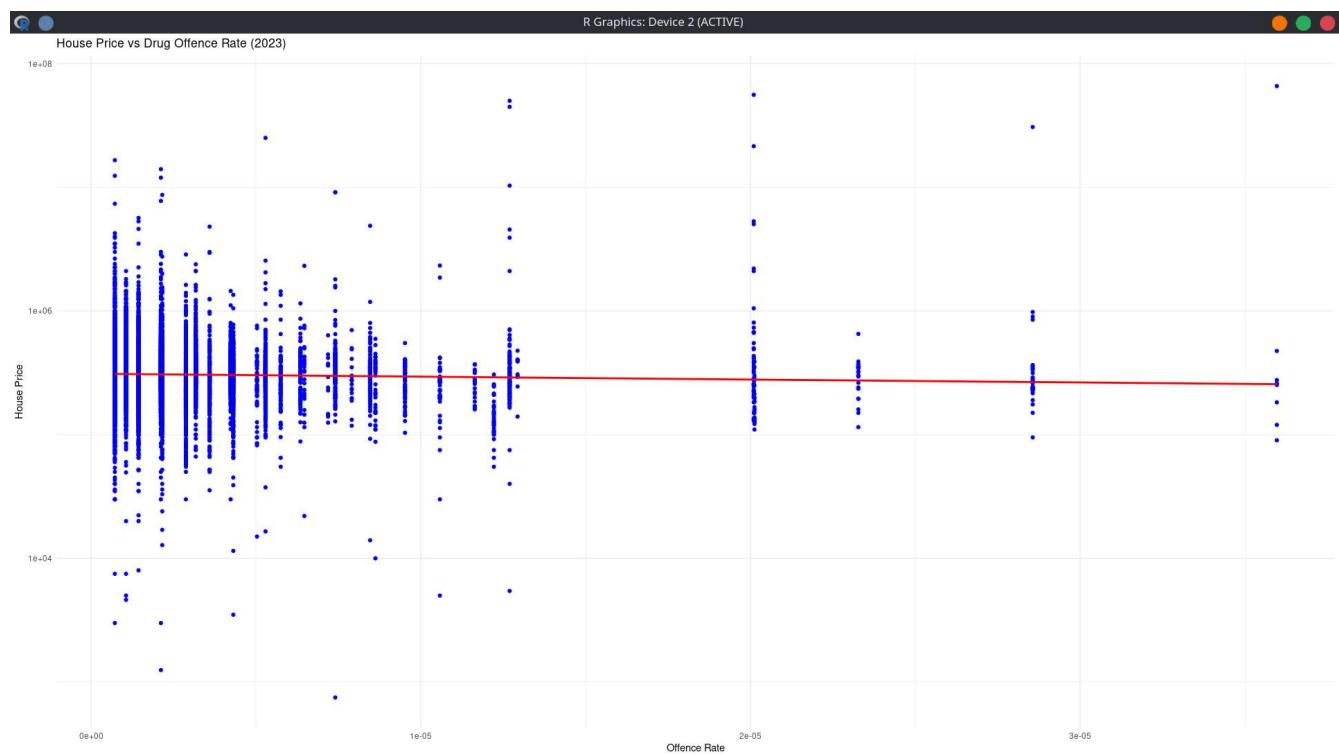


Figure 42: Code to generate line of best fit for house price and attainment 8 score

```
1 ggplot(merged_dataset, aes(x = Price, y = ATT8SCR)) +
2   geom_point(color = "blue", alpha = 0.5) + # Scatter plot of data points
3   geom_smooth(method = "lm", color = "red", se = FALSE) + # Line of best fit
4   labs(
5     title = "Relationship between School ATT8 Scores and House Prices",
6     x = "House Price",
7     y = "ATT8 Score"
8   ) +
9   theme_minimal()
```

Figure 43: line of best fit for house price and attainment 8 score

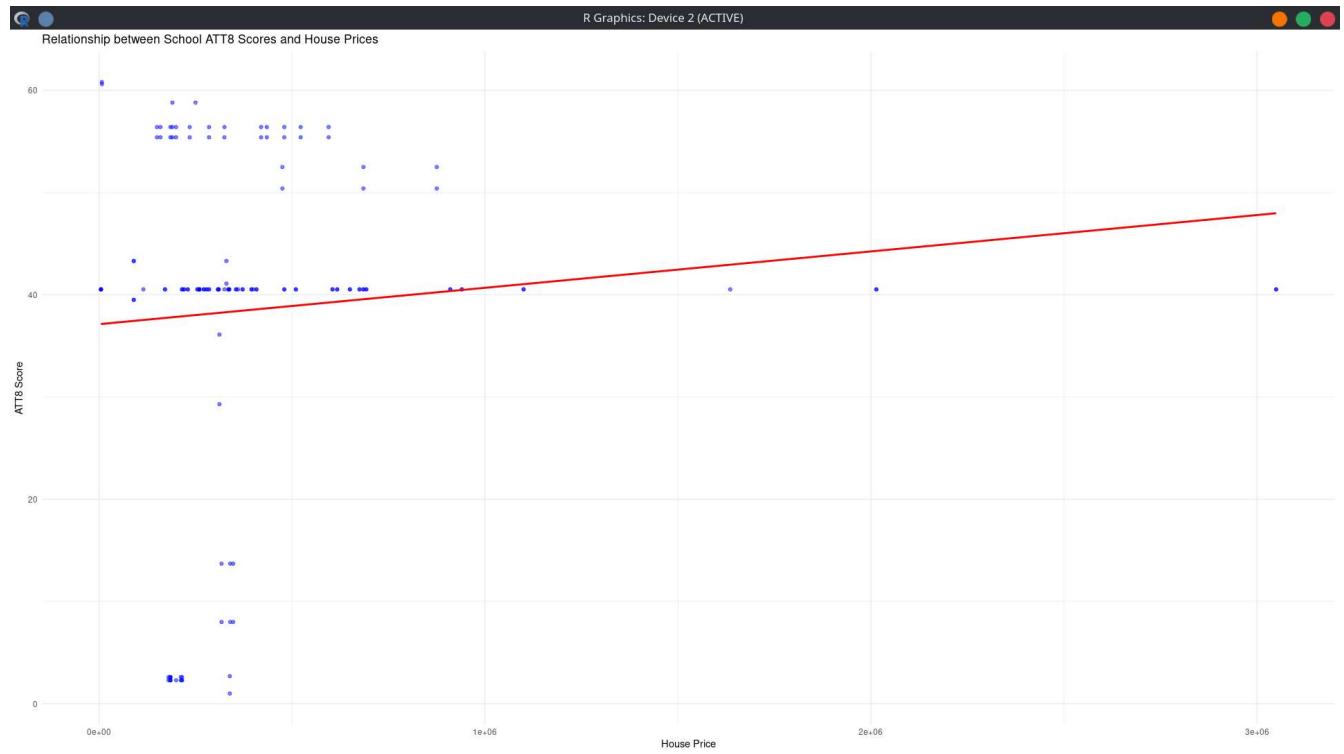


Figure 44: Code to generate line of best fit for download speed and drug offence rate per 1000 people

```

1 # Create a scatter plot with the line of best fit
2 ggplot(merged_dataset, aes(x = `Average download speed (Mbit/s)`, y = offence_rate)) +
3   geom_point(color = "blue", alpha = 0.5) + # Scatter plot of data points
4   geom_smooth(method = "lm", color = "red", se = FALSE) + # Line of best fit
5   labs(
6     title = "Relationship between Average Download Speed and Drug Offence Rates per 10000 people",
7     x = "Average Download Speed (Mbit/s)",
8     y = "Drug Offence Rate per 1000 Population"
9   ) +
10  theme_minimal()
11

```

Figure 45:line of best fit for download speed and drug offence rate per 10000 people

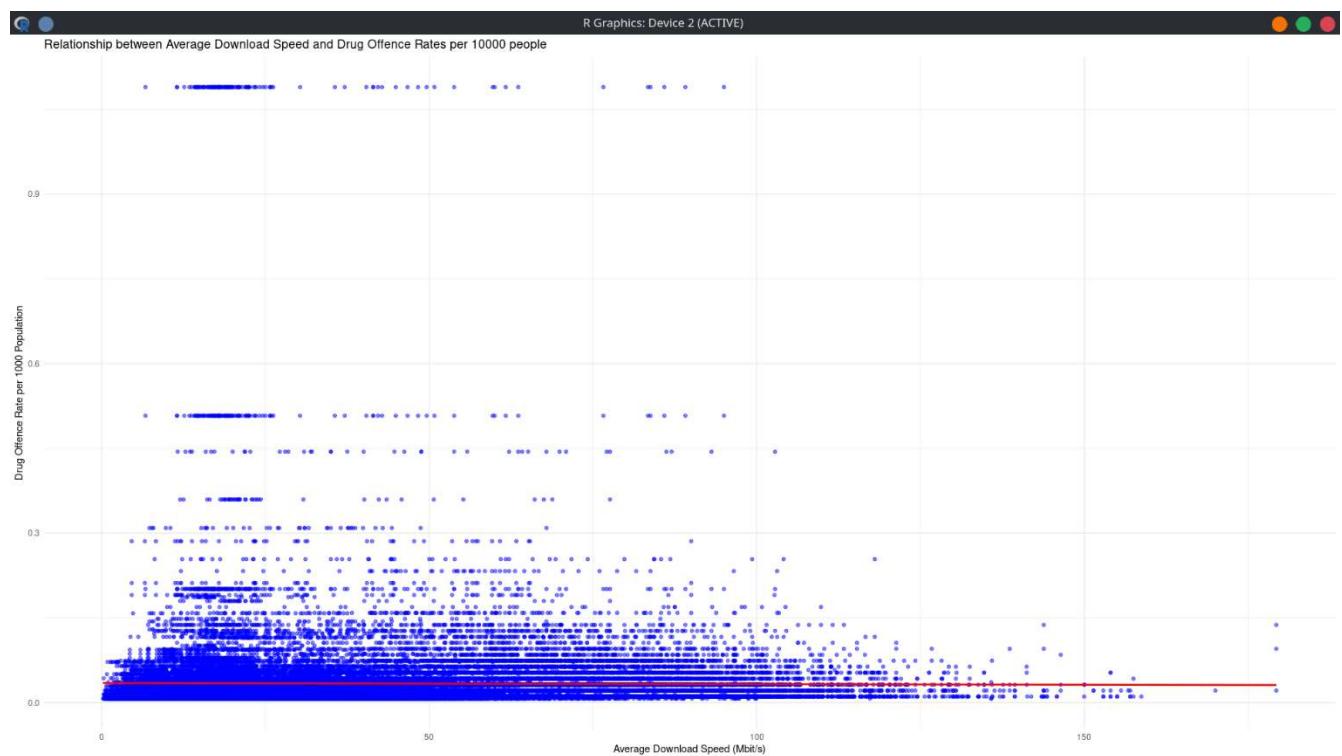
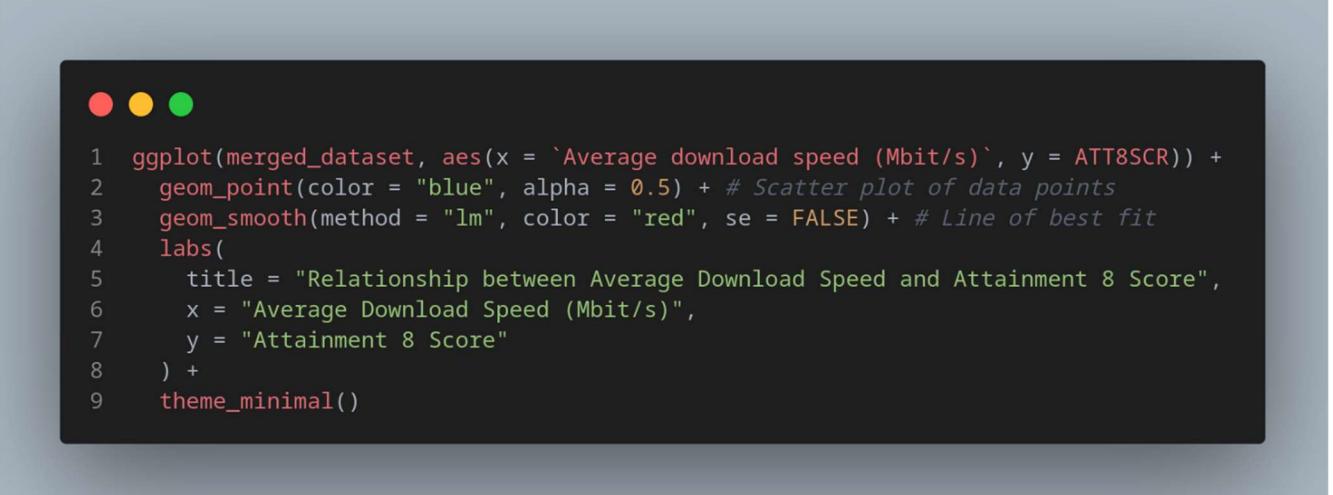


Figure 46: Code to generate line of best fit for average download speed and attainment 8 score



```

1 ggplot(merged_dataset, aes(x = `Average download speed (Mbit/s)`, y = ATT8SCR)) +
2   geom_point(color = "blue", alpha = 0.5) + # Scatter plot of data points
3   geom_smooth(method = "lm", color = "red", se = FALSE) + # Line of best fit
4   labs(
5     title = "Relationship between Average Download Speed and Attainment 8 Score",
6     x = "Average Download Speed (Mbit/s)",
7     y = "Attainment 8 Score"
8   ) +
9   theme_minimal()

```

Figure 47: line of best fit for average download speed and attainment 8 score

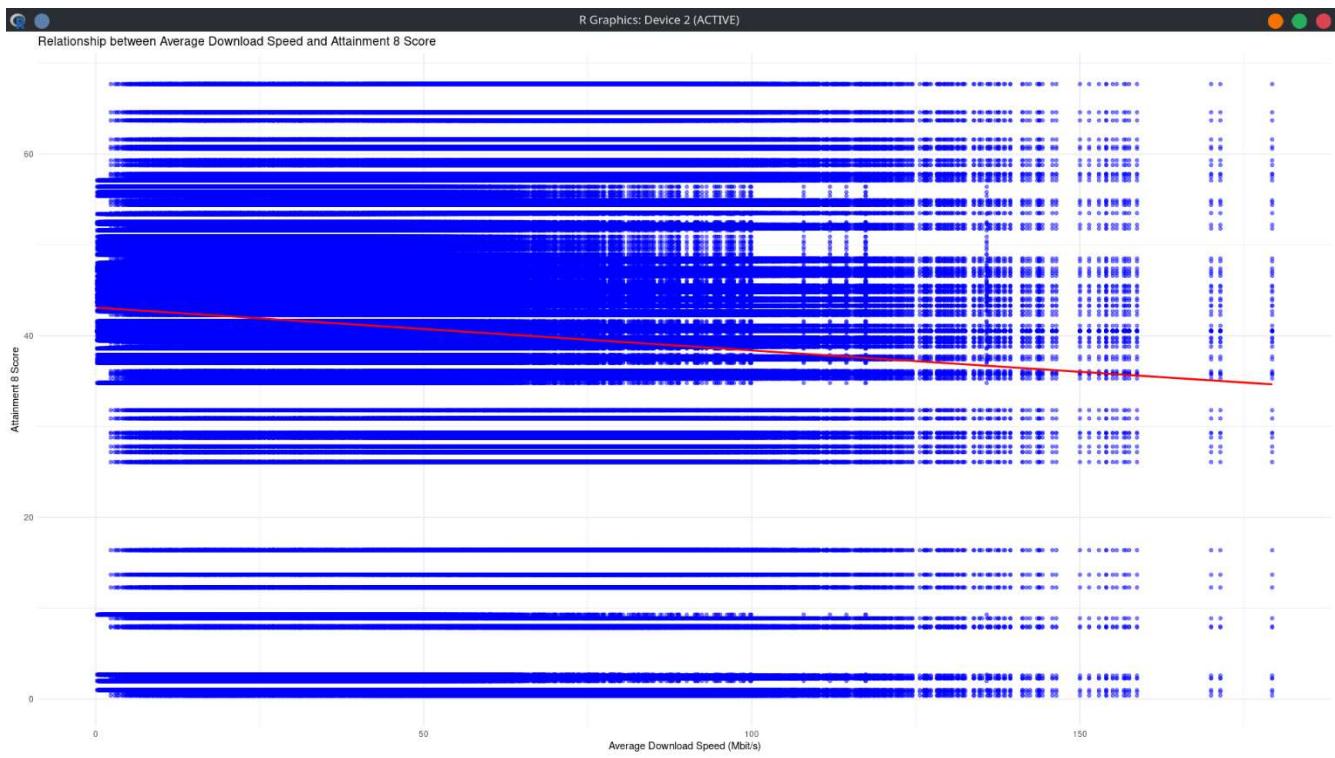


Figure 48: Code to perform city ranking for houses

```

1 house_pricing_dataset <- read_csv("Cleaned_datasets/House_Pricing_Data/clean_house_pricing_data.csv")
2 house_rankings <- house_pricing_dataset %>%
3   group_by(County, Town.City) %>%
4   # Calculate percentile values for price distribution
5   summarize(
6     p25 = quantile(Price, 0.25, na.rm = TRUE),
7     p50 = quantile(Price, 0.50, na.rm = TRUE),
8     p75 = quantile(Price, 0.75, na.rm = TRUE),
9     skewness = (mean(Price, na.rm = TRUE) - median(Price, na.rm = TRUE)) / sd(Price, na.rm = TRUE),
10    .groups = 'drop'
11  ) %>%
12  # Calculate a score based on your formula
13  mutate(score = p25 * 0.25 + p50 * 0.35 + p75 * 0.4 - skewness * 1000) %>%
14  # Rank cities within each county based on score
15  group_by(County) %>%
16  arrange(desc(score)) %>%
17  mutate(rank = row_number()) %>%
18  # Filter for top 10 cities in each county
19  filter(rank <= 10) %>%
20  ungroup()
21
22 # Print the result
23 print(house_rankings)

```

Figure 49: Code to perform city ranking by broadband

```

1 broadband_coverage_dataset <- read_csv("Cleaned_datasets/Broadband/clean_broadband_coverage.csv")
2 broadband_performance_dataset <- read_csv("Cleaned_datasets/Broadband/clean_broadband_performance.csv")
3
4
5
6 # Normalize the data
7 normalize <- function(x) {
8   (x - min(x, na.rm = TRUE)) / (max(x, na.rm = TRUE) - min(x, na.rm = TRUE))
9 }
10
11 # Merge the datasets by PostalCode
12 combined_dataset <- broadband_coverage_dataset %>%
13   inner_join(broadband_performance_dataset, by = c("PostalCode", "oa11cd", "Street", "City", "County"))
14
15
16 combined_dataset <- combined_dataset %>%
17   mutate(
18     SFBB_norm = normalize(`SFBB availability (% premises)`),
19     UFBB_norm = normalize(`UFBB availability (% premises)`),
20     FTTT_norm = normalize(`FTTP availability (% premises)`),
21     MaxDownload_norm = normalize(`Maximum download speed (Mbit/s)`),
22     AvgDownload_norm = normalize(`Average download speed (Mbit/s)`),
23     Weighted_Score = SFBB_norm * 0.2 + UFBB_norm * 0.25 + FTTT_norm * 0.25 +
24       MaxDownload_norm * 0.15 + AvgDownload_norm * 0.15
25   )
26

```

Figure 50: Code to perform city ranking via schools

Figure 51: Code to perform city ranking via crime

```
1 crime_dataset <- read_csv("Cleaned_datasets/Crime/Crime_Data_Combined.csv")
2
3 crime_totals <- crime_dataset %>%
4   group_by(County, City, Crime) %>%
5   summarize(total_offences = n(), .groups = 'drop') %>%
6   # Pivot the data to have total offences of each crime type as columns
7   pivot_wider(names_from = Crime, values_from = total_offences, values_fill = 0) %>%
8   # Calculate the total offences across all crime types for each city
9   rowwise() %>%
10  mutate(total_crime_offences = sum(c_across(where(is.numeric)))) %>%
11  ungroup() %>%
12  # Rank cities within each county based on total crime offences
13  group_by(County) %>%
14  arrange(total_crime_offences) %>%
15  mutate(rank = row_number()) %>%
16  # Select top 10 cities with the fewest total crime offences for each county
17  filter(rank <= 10) %>%
18  ungroup()
19
20 # Print the result
21 print(crime_totals)
22
23 crime_rank <- crime_totals %>%
24   select(County, City, rank)
```

Figure 52: Code to perform final ranking of cities

```

1 combined_rankings <- full_join(
2   broadband_rankings %>% select(County, City, broadband_rank = Rank),
3   crime_rank %>% select(County, City, crime_rank = rank),
4   by = c("County", "City")
5 ) %>%
6   full_join(house_rankings %>% select(County, City = Town.City, house_rank = rank), by = c("County", "City")) %>%
7   full_join(top_schools_by_county %>% select(County, City, school_rank = rank), by = c("County", "City"))
8
9 # Step 2: Replace NA values with a high number to ensure missing cities are ranked lower
10 combined_rankings <- combined_rankings %>%
11   mutate(
12     broadband_rank = ifelse(is.na(broadband_rank), max(combined_rankings$broadband_rank, na.rm = TRUE) + 1, broadband_rank),
13     crime_rank = ifelse(is.na(crime_rank), max(combined_rankings$crime_rank, na.rm = TRUE) + 1, crime_rank),
14     house_rank = ifelse(is.na(house_rank), max(combined_rankings$house_rank, na.rm = TRUE) + 1, house_rank),
15     school_rank = ifelse(is.na(school_rank), max(combined_rankings$school_rank, na.rm = TRUE) + 1, school_rank)
16   )
17
18 # Step 3: Calculate the average rank for each city
19 overall_rankings <- combined_rankings %>%
20   group_by(County, City) %>%
21   summarize(
22     avg_broadband_rank = mean(broadband_rank, na.rm = TRUE),
23     avg_crime_rank = mean(crime_rank, na.rm = TRUE),
24     avg_house_rank = mean(house_rank, na.rm = TRUE),
25     avg_school_rank = mean(school_rank, na.rm = TRUE),
26     overall_rank = mean(c(avg_broadband_rank, avg_crime_rank, avg_house_rank, avg_school_rank), na.rm = TRUE),
27     .groups = 'drop'
28   ) %>%
29   arrange(County, overall_rank)
30
31 # Step 4: Get top 3 entries from each county
32 top_3_by_county <- overall_rankings %>%
33   group_by(County) %>%
34   slice_head(n = 3) %>%
35   arrange(County, overall_rank)
36
37 # Print the result
38 print(top_3_by_county)

```