



in collaboration with



ST5014CEM

Data Science for Developers

Coursework: Individual Report

Nishan Bista

230348

Submitted To

Siddhartha Neupane

Table of Contents

Introduction	4
Cleaning Data.....	6
House Data	6
Town and Postcodes	7
Broadband Data.....	7
Crime Data	8
School Data.....	8
Exploratory Data Analysis	9
Visualization for house prices	9
<i>Line Chart: Average House Prices from 2021 to 2024</i>	9
<i>Bar Chart: Average House Price in 2023</i>	10
<i>Boxplots: Distribution of House Prices by County (District)</i>	11
Visualisation for Broadband speed data	13
<i>Boxplots: Average download speed by county</i>	13
Visualization for crime data.....	17
<i>Boxplot: Drug offense rate in the district of both counties</i>	17
<i>Radar chart: for Vehicle crime rate for West Yorkshire for May 2022</i>	18
<i>Pie chart: for Robbery rate for any one of two counties (for any specific month and year).....</i>	19
<i>Line chart for Drug offense rates per 10,000 people for both counties for all years.</i>	20
Visualization for school data.....	21
<i>Boxplot: Average attainment 8 score 2022</i>	21
<i>Line Graph: Relationship between attainment 8 score and years over multiple districts in South and West Yorkshire</i>	23
Linear Modeling	24
House Price vs Download Speed for both Counties	24
<i>Linear Model summary & correlation.....</i>	25

House price vs Drug rates (2023) per 10000 people for both counties	26
<i>Linear Model summary and correlation</i>	27
Attainment 8 score vs House Price for both counties	28
<i>Linear model summary and correlation</i>	29
Attainment 8 scores vs Drug Offense rates per 10000 people in 2023 for both counties	30
<i>Linear Model summary and correlation</i>	31
Average Download speed vs Drug Offense Rate per 10000 people for both counties	32
<i>Linear Model summary and correlation</i>	33
Average download speed vs Attainment 8 score for both counties.....	33
<i>Linear Model summary and correlation</i>	34
Recommendation System	35
Overview	35
Results	35
Based on Broadband speed	35
Based on School Attainment score.....	36
Based on House Prices	36
Based on Crime Data (safest)	37
Overall Score	37
Reflection.....	38
Legal and Ethical Issues	39
Conclusion.....	40
References	41
Appendix.....	42

Table of Figures

Figure 1, Average house prices Line chart	10
Figure 2, Average House price Bar chart	11
Figure 3, House by districts west Box plot.....	12
Figure 4, House by districts South Box plot	12
Figure 5, Average Download speed West Box plot	13
Figure 6, Average Download speed South box plot.....	14
Figure 7, West Towns by Average Download speed	15
Figure 8, South Towns by Average Download speed	16
Figure 9, Drug offense Rate South Districts	17
Figure 10, Drug offense Rate West districts	18
Figure 11, Vehicle Crime Radar chart	19
Figure 12, Robbery Distribution Pie chart.....	20
Figure 13, Drug offense trend Counties	21
Figure 14, Attainment 8 score South Box plot	22
Figure 15, Attainment 8 score West Box plot	22
Figure 16, Attainment 8 score South and West Line chart.....	23
Figure 17, Download speed vs House Price	25
Figure 18, LM 1 West Summary	25
Figure 19, LM 1 South Summary	26
Figure 20, House price vs Drug Offense	27
Figure 21, LM 2 Summary.....	27
Figure 22, Attainment 8 vs House price by county	28
Figure 23, LM 3 West Summary	29
Figure 24, LM 3 South Summary	29
Figure 25, Attainment 8 vs Drug offense rate	30
Figure 26, LM 4 West Summary	31
Figure 27, LM 4 South Summary	31
Figure 28, Drug Offense vs Average Download Speed.....	32
Figure 29, LM 5 Summary.....	33
Figure 30, Average Download Speed vs Attainment 8.....	34
Figure 31, LM 6 Summary.....	34
Figure 32, Broadband speed result.....	35
Figure 33, School Result	36
Figure 34, House prices Result	36
Figure 35, Crime Result	37
Figure 36, Overall Score - Final Results	37

Introduction

This assignment presents the process and outcomes of using data analysis to recommend the most suitable location to purchase a property in the United Kingdom, specifically comparing South Yorkshire and West Yorkshire. Both regions offer their own benefits, but making a property investment requires a detailed analysis of several factors that affect housing value and quality of life. Key aspects considered in this project include house prices, crime rates, broadband connectivity, education, and access to public amenities.

The main goal of this project is to build a simple recommendation system that uses publicly available UK government data to evaluate and compare towns in South and West Yorkshire. Based on selected criteria, each town receives a score which helps identify the most promising areas for investment. This report outlines the steps followed throughout the project, including data collection, cleaning, integration, exploratory data analysis, and linear modeling to support the final recommendation.

Cleaning Data

Data cleaning is a crucial step in the data science lifecycle that comes right after collecting the datasets. It helps ensure the data is accurate, consistent, and ready for analysis. The datasets used in this project, which were sourced from UK government portals and public sources, contained various inconsistencies such as missing values, formatting issues, and mismatched fields.

Each dataset was carefully cleaned and standardized to make sure they could be merged and analyzed together. This included removing duplicates, handling missing or invalid entries, formatting columns like postcodes, and filtering the data to include only relevant information from South Yorkshire and West Yorkshire. Cleaning the data properly was important to ensure reliable results and support the development of the final recommendation system.

House Data

The house price data from 2021 to 2024 was combined into a single dataset using the `bind_rows()` function in R. Since the original files did not have headers, column names were manually assigned based on the UK Price Paid Data structure. The postcode column was cleaned by trimming white spaces and converting all values to uppercase for consistency. Then, the dataset was filtered to include records only from South Yorkshire and West Yorkshire. Finally, the cleaned dataset was saved using the `write_csv()` function for further analysis.

Town and Postcodes

To map postcodes to their corresponding counties and LSOA codes, two datasets were used: one containing postcode-to-LSOA mappings and another with county information from the UK geoportal. These datasets were merged using the `inner_join()` function in R based on the local authority district code. The merged data was then cleaned by removing unnecessary columns and formatting the postcode column by removing spaces. Finally, the dataset was filtered to include only entries from South Yorkshire and West Yorkshire, and the result was saved as a CSV file for later use.

Broadband Data

Two separate broadband datasets were used, one for performance and one for coverage. These were cleaned by removing missing or invalid postcodes and standardizing them to uppercase with no spaces. The two datasets were then joined using the postcode field to combine performance and coverage information.

To link broadband data with location, a `postcode_to_LSOA` mapping data was used. After standardizing the postcodes in both datasets, they were merged to include LSOA codes. Further, the dataset was joined with a postcode-county mapping to isolate data for South Yorkshire and West Yorkshire. Irrelevant columns were removed to keep only the necessary features for analysis. Finally, the cleaned and filtered broadband dataset was saved for further use.

Crime Data

The crime data was collected from multiple CSV files stored in different folders, each representing a specific month and county. These files were combined into one dataset using a custom function that extracted metadata such as the year, month, and county name from the file names. Only the necessary columns were selected, including crime type, location, and LSOA code.

To enrich the dataset, a postcode-to-LSOA mapping was used to assign a postcode to each LSOA. Then, population data was added by matching postcodes, allowing crime rates to later be normalized by population. The final cleaned dataset includes crime, location, and population data for South Yorkshire and West Yorkshire, and was saved as a CSV file for further analysis.

School Data

School performance data from 2021 to 2024 was collected and combined into a single dataset. For each year, only the relevant columns were selected, including school name, postcode, town, and the Attainment 8 score. A new column was added to indicate the academic year.

The Attainment 8 score column was cleaned by removing non-numeric values such as "NE" and "SUPP", and then converted to numeric for analysis. Postcodes were trimmed to ensure consistency, and the dataset was joined with a postcode-county mapping to add county information. Finally, the data was filtered to include only schools located in South Yorkshire and West Yorkshire, and saved as a cleaned CSV for further use.

Exploratory Data Analysis

Exploratory Data Analysis (EDA) is an important step in understanding the underlying patterns and trends in a dataset. After cleaning and preparing the data, different visualization techniques were used to explore the key features across housing prices, crime rates, broadband availability, and school performance in South Yorkshire and West Yorkshire. Graphs such as bar charts, line plots, and box plots were used to help compare distributions, identify relationships, and highlight differences between the two regions. These insights supported the development of a data-driven recommendation system.

Visualization for house prices

Line Chart: Average House Prices from 2021 to 2024

To understand how house prices have changed over time, the dataset was first filtered to include the years from 2021 to 2024. The records were then grouped by year and county using the group_by() function, and the average house price was calculated with summarise(). The line chart was created using the ggplot2 library where the x-axis represents the year, and the y-axis shows the average house price.

From the line chart, we can clearly see how prices have shifted in West Yorkshire and South Yorkshire during these years. In general, South Yorkshire house prices have increased gradually from the past 2021 to 2024 . Whereas West Yorkshire house prices have decreased from 2023. However, the rate of growth in South Yorkshire appears to be more stable, while West Yorkshire shows slight fluctuations across 2023 to 2024.

Figure 1, Average house prices Line chart

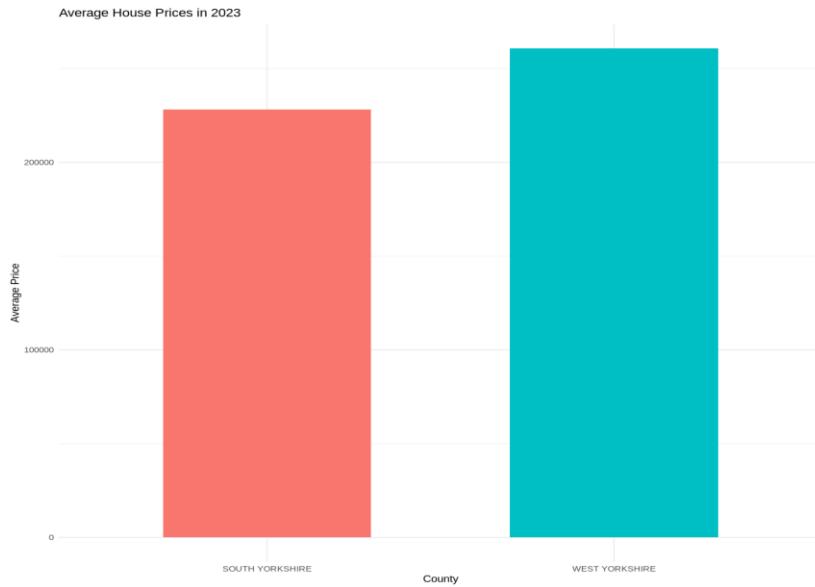


Bar Chart: Average House Price in 2023

For this chart, the dataset was filtered to include only transactions from the year 2023. Then, using group_by() and summarise(),. A bar chart was created using the ggplot2 package with the counties on the x-axis and the average price on the y-axis.

This bar chart gives a quick comparison of house prices in West Yorkshire and South Yorkshire for 2023. The chart shows that West Yorkshire had a slightly higher average price compared to South Yorkshire that year. This could be due to higher property values in certain towns within West Yorkshire or possibly due to different property types being more common in one area over the other.

Figure 2, Average House price Bar chart



Boxplots: Distribution of House Prices by County (District)

Boxplots was used to explore the spread and variation in house prices for both West Yorkshire and South Yorkshire. The dataset includes all years, and the prices were grouped by county. Separate boxplots were created for each county to better understand the distribution of housing prices.

The boxplots show a clear comparison of price ranges in both counties via district. South Yorkshire's prices seem to have a wider range, with some high outliers, suggesting that there might be certain towns or property types that are significantly more expensive. In contrast, West Yorkshire's distribution looks more compact, with most prices clustering around near the Q3 suggesting that the majority of house prices are relatively low compared to South Yorkshire, with only a few scattered outliers at much higher values. These visualizations help us understand not just the averages but also the variability and extremities in the housing market across both counties.

Figure 3, House by districts west Box plot



Figure 4, House by districts South Box plot



Visualisation for Broadband speed data

Boxplots: Average download speed by county

The clean broadband data is loaded from the csv file two boxplots are generated to visualize the variation in average download speeds across different districts within West Yorkshire and South Yorkshire. First, it filters the dataset for each county separately using the filter() function. Then, ggplot() is used to set up the plots, geom_boxplot() creates the boxplot. This visualization helps in identifying differences in internet speed distribution among districts within each county.

In West Yorkshire, All districts display a wide range of download speeds with numerous outliers, indicating the presence of faster connections in some areas.

In South Yorkshire, The median speeds appear similar across districts, all districts show a long right tail with numerous outliers indicating higher speeds, reaching up to 200 Mbps.

Figure 5, Average Download speed West Box plot

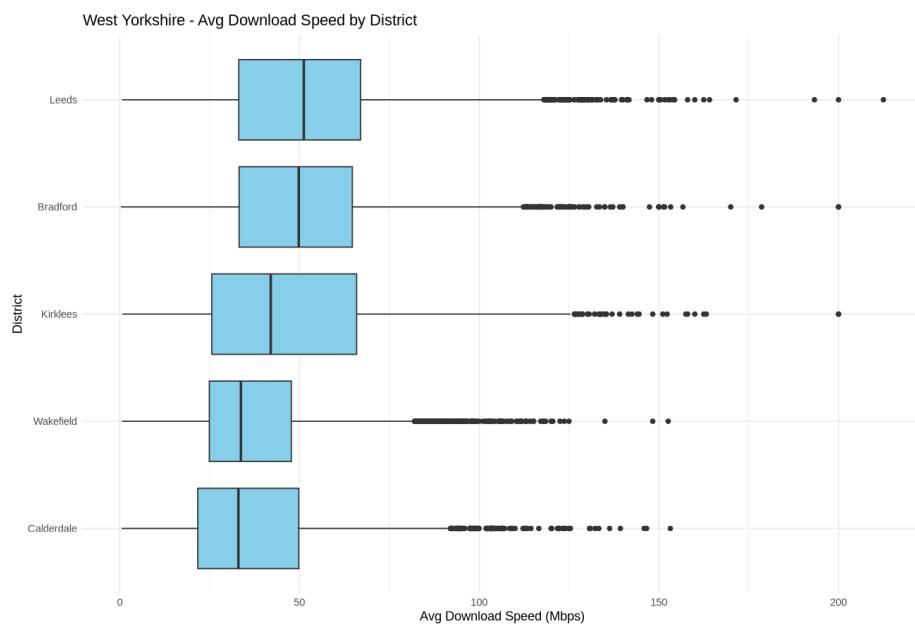
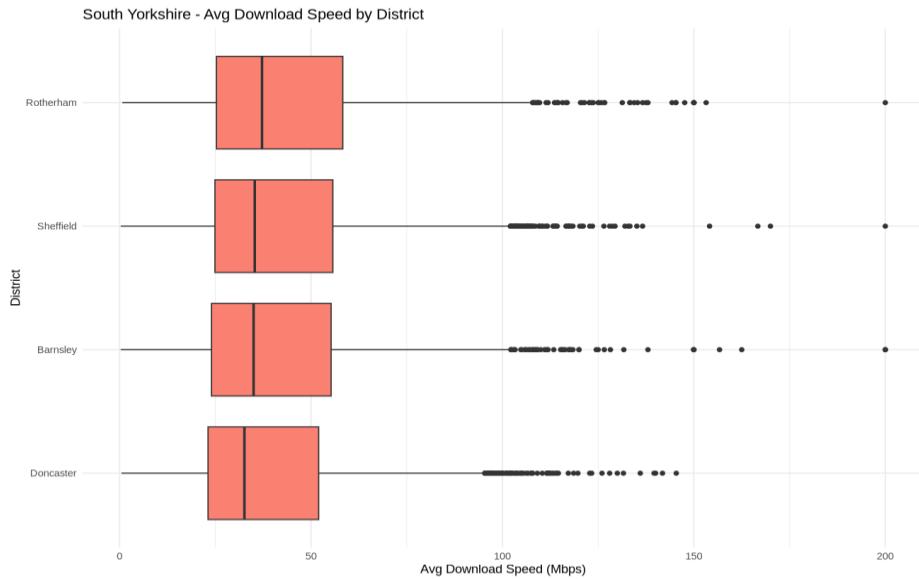


Figure 6, Average Download speed South box plot



Bar charts: Town vs Average Speed

This helps us to compare the cities in West Yorkshire and South Yorkshire based on average download speed. The data is grouped by district name (ladnm) using group_by(), and the mean download speed is calculated with summarise(), excluding any missing values. The results are sorted in descending order with arrange(). The ggplot() function then constructs a bar chart, where towns are reordered by average speed for clarity. These visuals help quickly identify which towns offer faster internet speeds within each region.

In West Yorkshire, Leeds is the town with the highest average broadband speed with around 52 Mbps followed by Bradford and Kirklees simultaneously.

In South Yorkshire, Rotherham Town has the highest average broadband speed with around 45 mbps. It is followed by Sheffield and Barnsley respectively with around average broadband speed of 42 Mbps

Figure 7, West Towns by Average Download speed

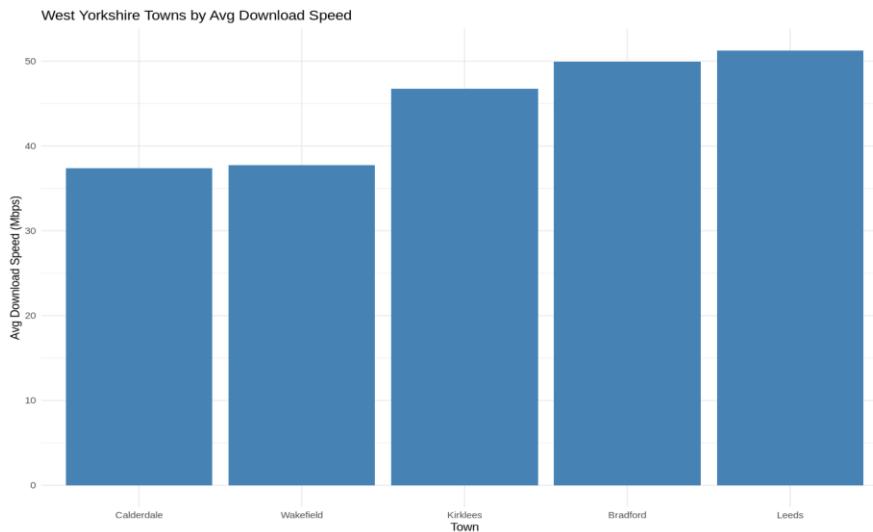
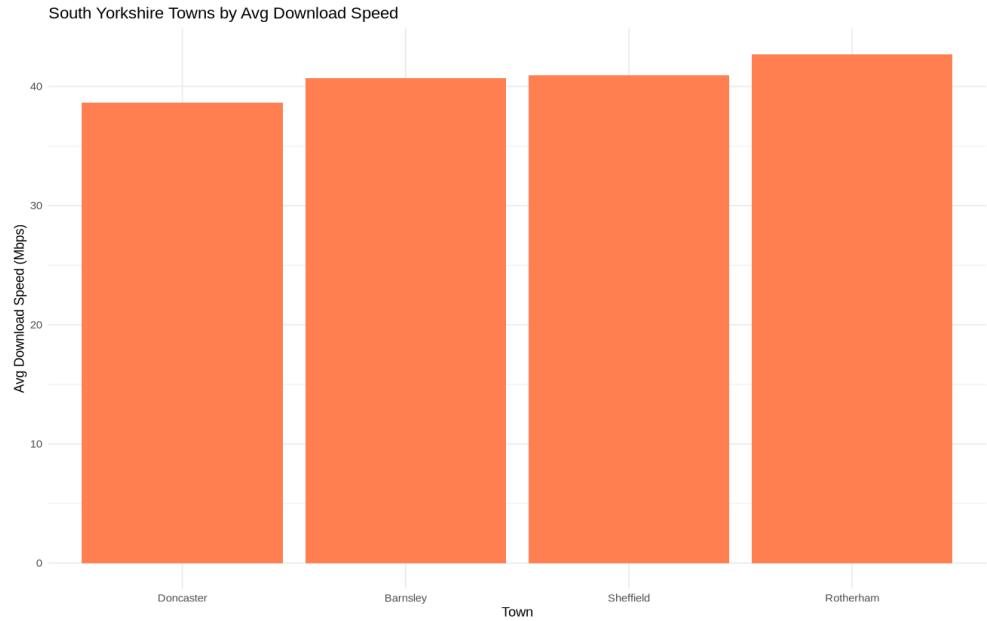


Figure 8, South Towns by Average Download speed



Visualization for crime data

Boxplot: Drug offense rate in the district of both counties

This boxplot compares drug offense rates in South and West Yorkshire. In South Yorkshire, Sheffield has the highest and most variable rates, while Barnsley shows the lowest and most stable. Rotherham and Doncaster fall in between. In West Yorkshire, Bradford leads with the highest and widest spread, followed by Leeds with a slightly more consistent rate. Wakefield stands out with the lowest rates, and Kirklees and Calderdale sit somewhere in the middle. Overall, Sheffield and Bradford seem to be the most affected, while Barnsley and Wakefield are the least.

Figure 9, Drug offense Rate South Districts

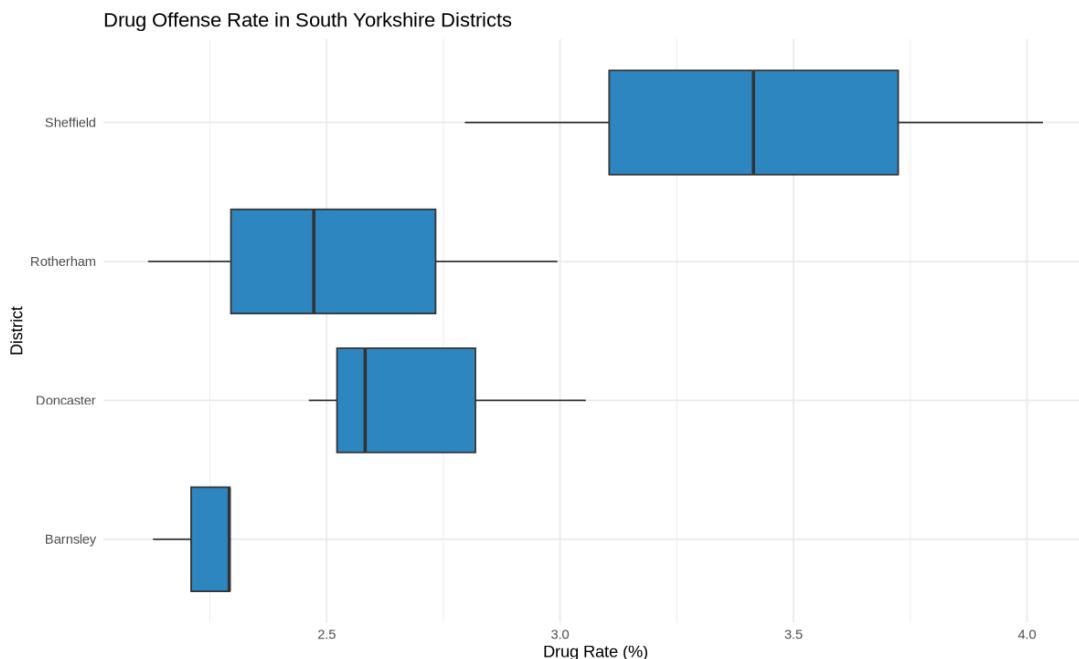
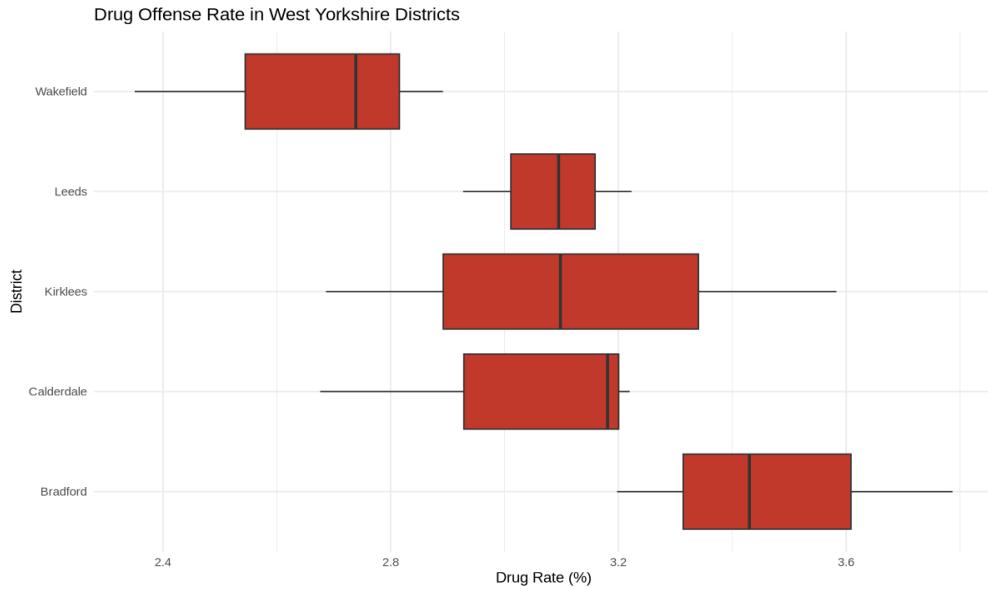


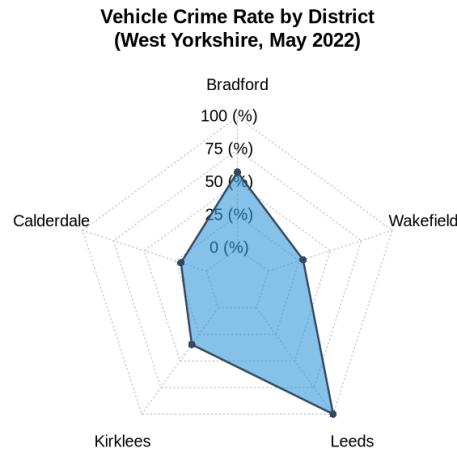
Figure 10, Drug offense Rate West districts



Radar chart: for Vehicle crime rate for West Yorkshire for May 2022

This radar chart visualizes vehicle crime rates across five districts in West Yorkshire for May 2022. Bradford and Leeds show the highest rates, indicating greater levels of vehicle-related crimes in these areas. In contrast, Kirklees, Calderdale, and Wakefield have noticeably lower rates, with Calderdale being the least affected. The chart highlights the regional disparity in vehicle crime, with Bradford standing out as the most impacted district.

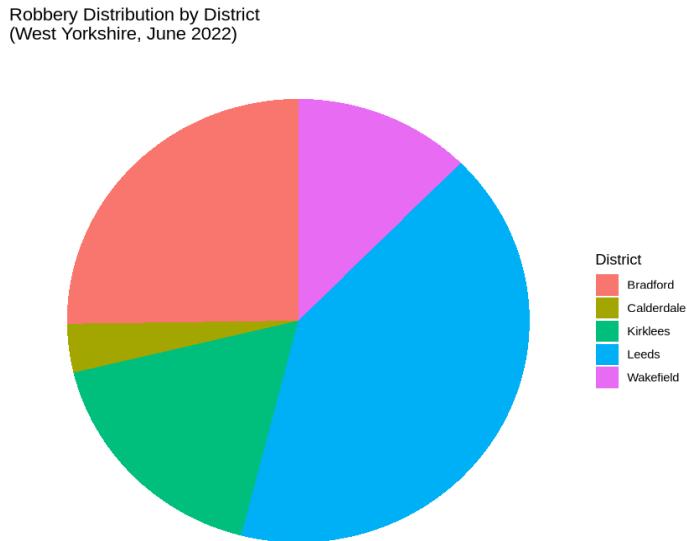
Figure 11, Vehicle Crime Radar chart



Pie chart: for Robbery rate for any one of two counties (for any specific month and year)

This pie chart shows how robbery cases were distributed across West Yorkshire districts in June 2022. Leeds had the largest share, followed by Bradford and Kirklees. Wakefield contributed a smaller portion, and Calderdale had the lowest share of robberies. The chart clearly highlights Leeds as the most affected district during this period.

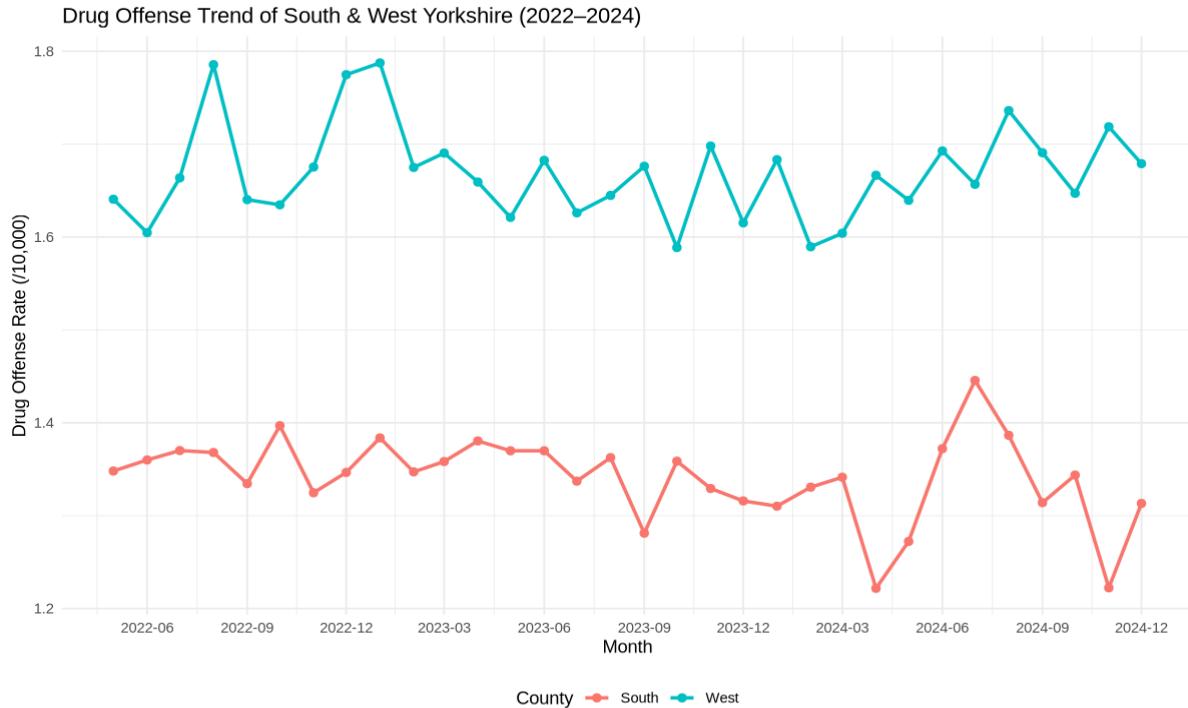
Figure 12, Robbery Distribution Pie chart



Line chart for Drug offense rates per 10,000 people for both counties for all years.

This line chart shows the monthly drug offense rates in South and West Yorkshire from 2022 to 2024. West Yorkshire consistently had higher rates than South Yorkshire throughout the period. While both regions experienced some fluctuations, West Yorkshire's trend remained more stable, whereas South Yorkshire showed more noticeable ups and downs over time.

Figure 13, Drug offense trend Counties



Visualization for school data

Boxplot: Average attainment 8 score 2022

The boxplot shows the distribution of the *Average Attainment 8 Score* for schools in 2022. This score reflects student performance across 8 key subjects. This will help us to know which district has better attainment scores in both counties.

Figure 14, Attainment 8 score South Box plot

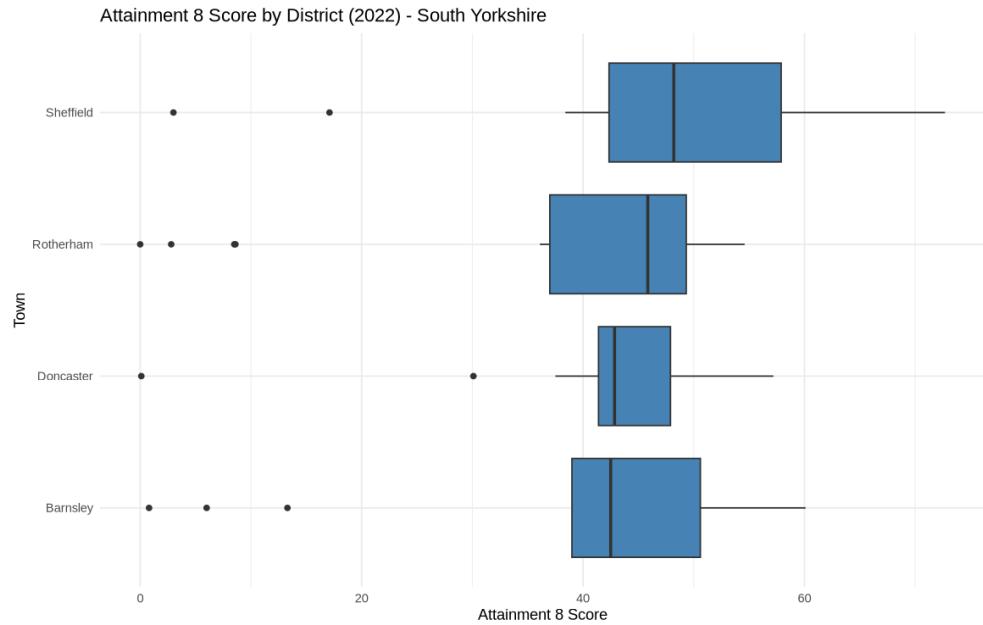
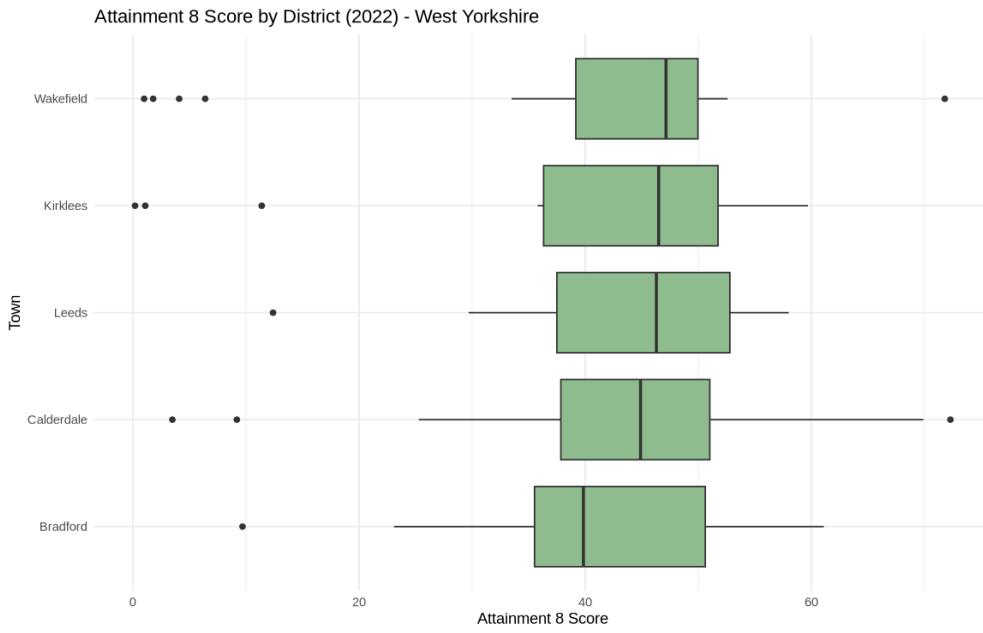


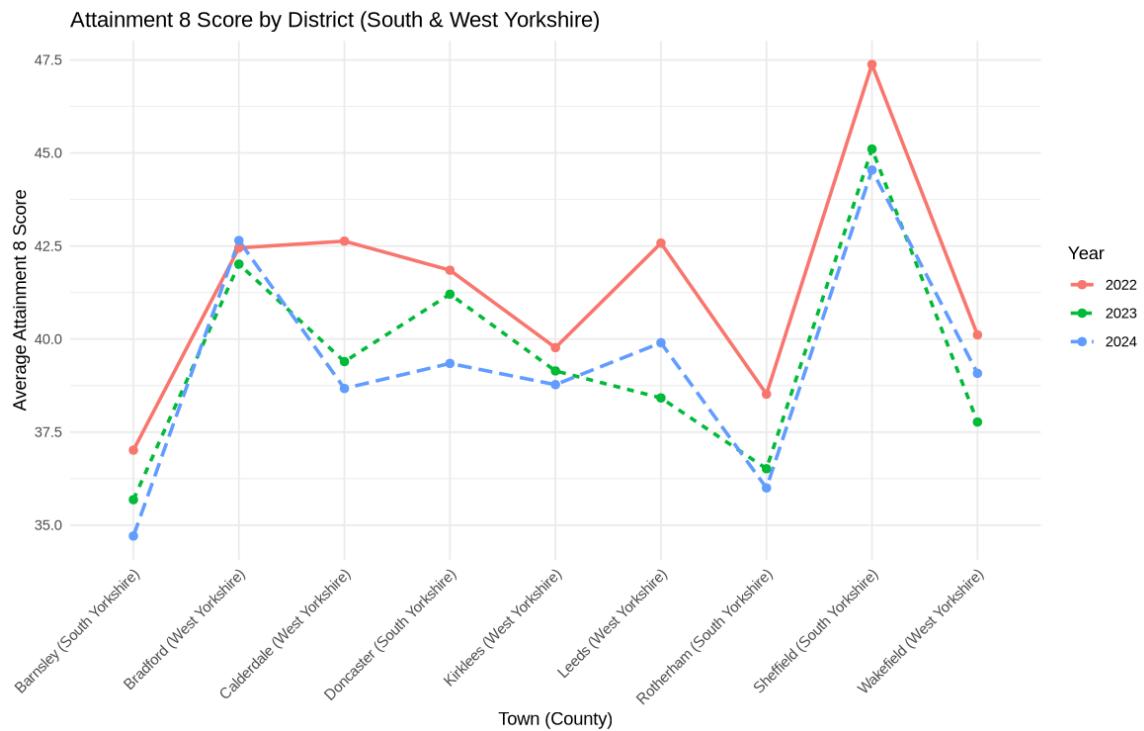
Figure 15, Attainment 8 score West Box plot



Line Graph: Relationship between attainment 8 score and years over multiple districts in South and West Yorkshire

From the line graph below, We can analyze that in three years Sheffield has a high attainment score whereas Barnsley has significantly low attainment score over the year. Sheffield also has the more consistent attainment score among other districts

Figure 16, Attainment 8 score South and West Line chart



Linear Modeling

Linear modeling is an easy to use yet potent statistical tool used for investigating the relationships between two or more variables. In this project, it was used to demonstrate trends and patterns between house price, broadband speed, crime rates, and school performance in both South Yorkshire and West Yorkshire.

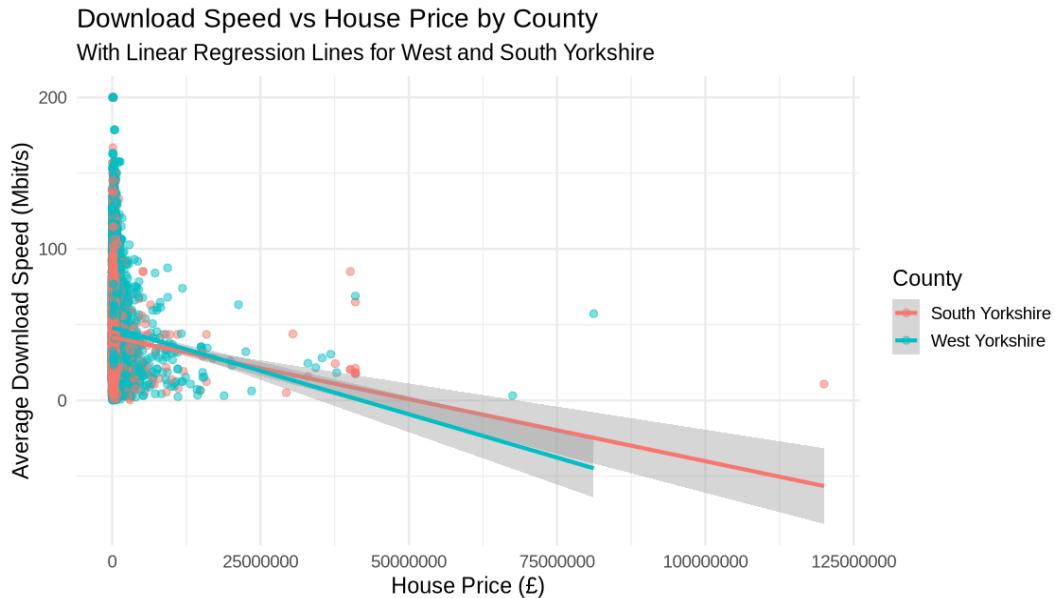
By Using a linear regression analysis allowed us to look at the shift when one variable (house price) changes with another (like average broadband speed or crime rate). For example, we assessed if areas with higher drug offense had lower house prices, or if the value of housing was linked with better school or college performance. We were able to go through these models and assess which were correlated and/or able to compare districts in an informed way, so we could address our top line conclusion on which district was a good place to buy a residential property.

House Price vs Download Speed for both Counties

Linear regression analysis carried out to analyze the relationship between house prices, and average download speed of broadband in West and South Yorkshire. The correlation in both regions was very weak and negative (-0.026 in West Yorkshire and -0.027 in South Yorkshire), which indicates there is almost no relationship at all.

While the regression models were statistically significant ($p < 0.001$), the effect size was negligible, download speed dropped by less than 0.000001 Mbps for every £1 increase in house price. The R-squared values were close to 0, suggesting that house prices explain almost no variation in download speed.

Figure 17, Download speed vs House Price



Linear Model summary & correlation

Figure 18, LM 1 West Summary

```

Call:
lm(formula = download_speed ~ price, data = west)

Residuals:
    Min      1Q  Median      3Q     Max 
-47.702 -17.848 -1.575  16.136 152.173 

Coefficients:
            Estimate Std. Error t value     Pr(>|t|)    
(Intercept) 48.0619745271  0.0679648863 707.159 <0.0000000000000002 *** 
price        -0.0000011440  0.0000001212  -9.442 <0.0000000000000002 *** 
...
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1 

Residual standard error: 22.46 on 131414 degrees of freedom
Multiple R-squared:  0.0006779, Adjusted R-squared:  0.0006703 
F-statistic: 89.15 on 1 and 131414 DF,  p-value: < 0.0000000000000022 

> cat("Correlation (West Yorkshire):", cor_west, "\n")
Correlation (West Yorkshire): -0.02603654

```

Figure 19, LM 1 South Summary

```

Call:
lm(formula = download_speed ~ price, data = south)

Residuals:
    Min      1Q  Median      3Q     Max 
-41.011 -15.936 -6.118  15.616 158.438 

Coefficients:
            Estimate Std. Error t value    Pr(>|t|)    
(Intercept) 41.7479820895  0.0794921759 525.184 < 0.000000000000002 *** 
price        -0.0000008189  0.0000001064   -7.695  0.000000000000143 *** 
---
Signif. codes:  0 ‘***’ 0.001 ‘**’ 0.01 ‘*’ 0.05 ‘.’ 0.1 ‘ ’ 1 

Residual standard error: 21.34 on 78406 degrees of freedom
Multiple R-squared:  0.0007547, Adjusted R-squared:  0.0007419 
F-statistic: 59.22 on 1 and 78406 DF,  p-value: 0.0000000000001429 

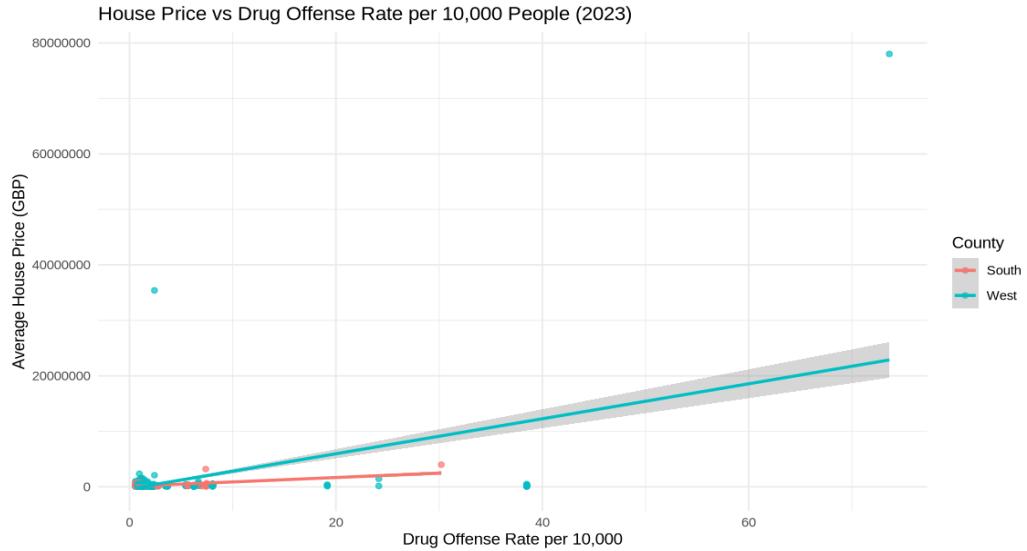
> cat("Correlation (South Yorkshire):", cor_south, "\n")
Correlation (South Yorkshire): -0.02747158

```

House price vs Drug rates (2023) per 10000 people for both counties

The line chart illustrates the relationship between average house prices and drug crime rates per 10,000 people between districts in West and South Yorkshire for 2023. A moderate positive correlation was observed ($r \approx 0.48$), suggesting that areas with higher drug crime rates may also have higher average house prices. The linear regression model further supports this, with a statistically significant positive coefficient for drug rate ($p < 0.001$). However, the R^2 value of 0.23 indicates that drug rates explain only 23% of the variation in house prices, suggesting other factors also play a major role.

Figure 20, House price vs Drug Offense



Linear Model summary and correlation

Figure 21, LM 2 Summary

```

Call:
lm(formula = AveragePrice ~ RatePer10k, data = joined_data)

Residuals:
    Min      1Q   Median      3Q     Max 
-10999066 -39040  104531  214628 56565648 

Coefficients:
            Estimate Std. Error t value    Pr(>|t|)    
(Intercept) -273653     88962  -3.076 0.00216 **  
RatePer10k    294885    17646  16.711 < 0.0000000000000002 *** 
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1 

Residual standard error: 2459000 on 933 degrees of freedom
Multiple R-squared:  0.2304,    Adjusted R-squared:  0.2295 
F-statistic: 279.3 on 1 and 933 DF,  p-value: < 0.0000000000000022 

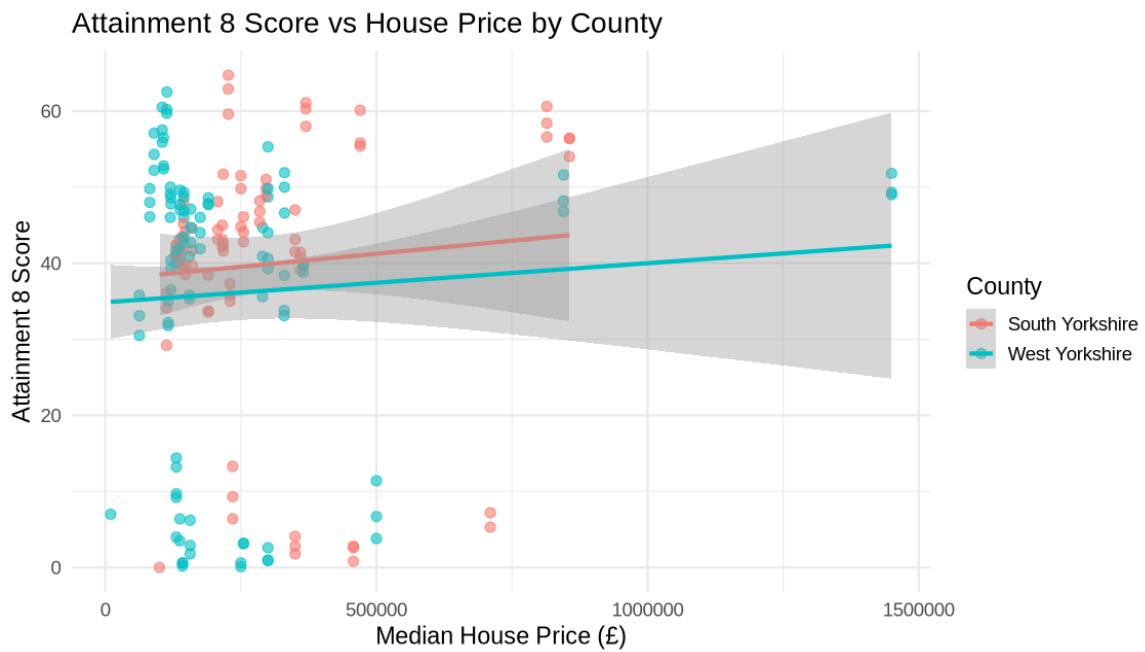
> cat("Correlation between Drug Rate and House Price:", correlation, "\n")
Correlation between Drug Rate and House Price: 0.4799685

```

Attainment 8 score vs House Price for both counties

The analysis explores the relationship between median house prices and average Attainment 8 scores between West and South Yorkshire. In both counties, the correlation is very weak ($r \approx 0.07\text{--}0.08$), and linear regression models showed no significant relationship ($p > 0.47$). This suggests that house prices in these regions do not have a meaningful impact on academic performance as measured by Attainment 8 scores.

Figure 22, Attainment 8 vs House price by county



Linear model summary and correlation

Figure 23, LM 3 West Summary

```

Call:
lm(formula = attainment ~ median_price, data = west)

Residuals:
    Min      1Q  Median      3Q     Max 
-36.046 -3.507  6.562 12.337 27.055 

Coefficients:
            Estimate Std. Error t value     Pr(>|t|)    
(Intercept) 34.861895345  2.500477359 13.94 <0.000000000000002 *** 
median_price  0.000005136  0.000007132   0.72      0.473  
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1 

Residual standard error: 18.57 on 102 degrees of freedom
Multiple R-squared:  0.00506,  Adjusted R-squared:  -0.004694 
F-statistic: 0.5188 on 1 and 102 DF,  p-value: 0.473 

> cat("Correlation (West Yorkshire):", cor_west, "\n")
Correlation (West Yorkshire): 0.07113407

```

Figure 24, LM 3 South Summary

```

Call:
lm(formula = attainment ~ median_price, data = south)

Residuals:
    Min      1Q  Median      3Q     Max 
-40.153 -1.411  3.639  9.459 25.329 

Coefficients:
            Estimate Std. Error t value     Pr(>|t|)    
(Intercept) 37.816226088  3.484159695 10.854 <0.000000000000002 *** 
median_price  0.000006856  0.000009706   0.706      0.482  
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1 

Residual standard error: 16.76 on 78 degrees of freedom
Multiple R-squared:  0.006357,  Adjusted R-squared:  -0.006382 
F-statistic: 0.499 on 1 and 78 DF,  p-value: 0.482 

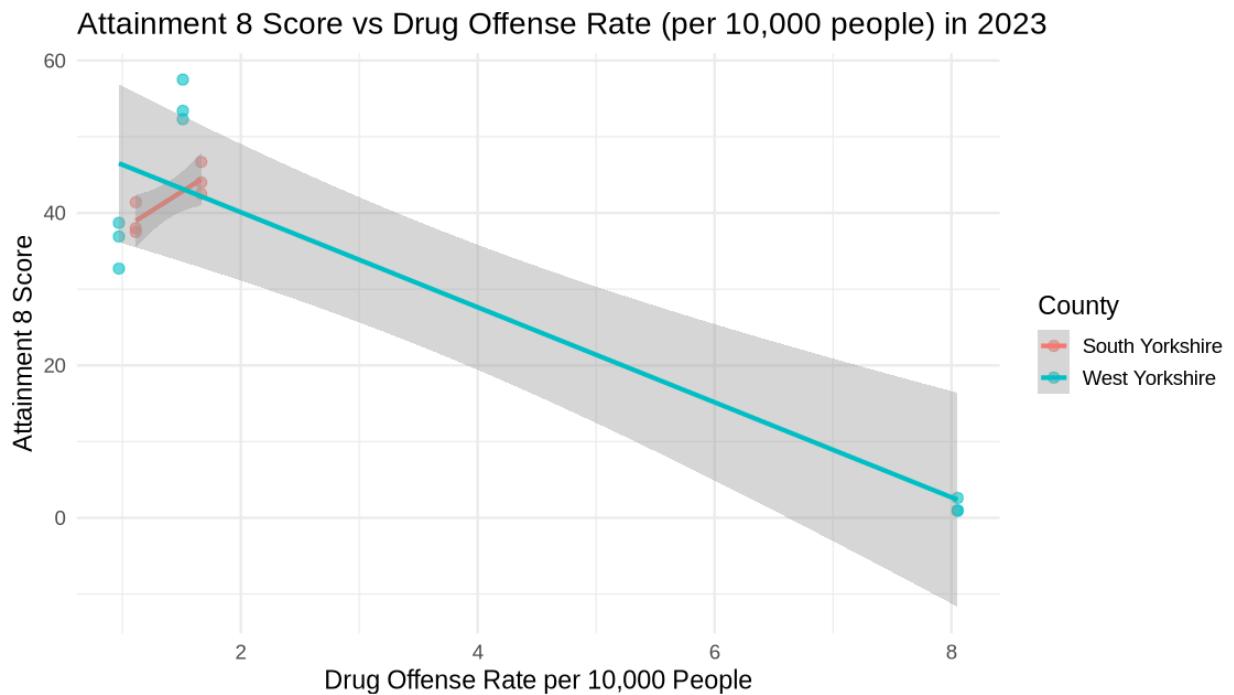
> cat("Correlation (South Yorkshire):", cor_south, "\n")
Correlation (South Yorkshire): 0.07973142

```

Attainment 8 scores vs Drug Offense rates per 10000 people in 2023 for both counties

The analysis shows a strong negative correlation in West Yorkshire ($r = -0.91$), indicating that areas with lower Attainment 8 score has higher drug offense rates. The regression model supports this with a statistically significant negative relationship ($p < 0.001$). In contrast, South Yorkshire shows a strong positive correlation ($r = 0.84$), and the model also finds a significant positive relationship ($p = 0.035$). This suggests that the relationship between education outcomes and drug offense rates varies notably between South and West Yorkshire.

Figure 25, Attainment 8 vs Drug offense rate



Linear Model summary and correlation

Figure 26, LM 4 West Summary

```

Call:
lm(formula = attainment ~ rate_per_10k, data = west)

Residuals:
    Min      1Q  Median      3Q     Max 
-13.803 -7.803 -1.358  9.161 14.361 

Coefficients:
            Estimate Std. Error t value Pr(>|t|)    
(Intercept) 52.557     5.083 10.340 0.0000172 ***
rate_per_10k -6.235     1.067 -5.842 0.000636 ***  
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 10.3 on 7 degrees of freedom
Multiple R-squared:  0.8298,   Adjusted R-squared:  0.8055 
F-statistic: 34.13 on 1 and 7 DF,  p-value: 0.0006358

> cat("Correlation (West Yorkshire):", cor_west, "\n")
Correlation (West Yorkshire): -0.9109349

```

Figure 27, LM 4 South Summary

```

Call:
lm(formula = attainment ~ rate_per_10k, data = south)

Residuals:
    1      2      3      4      5      6 
2.4333 2.3000 -0.9667 -1.9000 -1.4667 -0.4000 

Coefficients:
            Estimate Std. Error t value Pr(>|t|)    
(Intercept) 27.990     4.459  6.278 0.00329 ** 
rate_per_10k  9.855     3.147  3.131 0.03515 *  
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 2.125 on 4 degrees of freedom
Multiple R-squared:  0.7102,   Adjusted R-squared:  0.6378 
F-statistic: 9.804 on 1 and 4 DF,  p-value: 0.03515

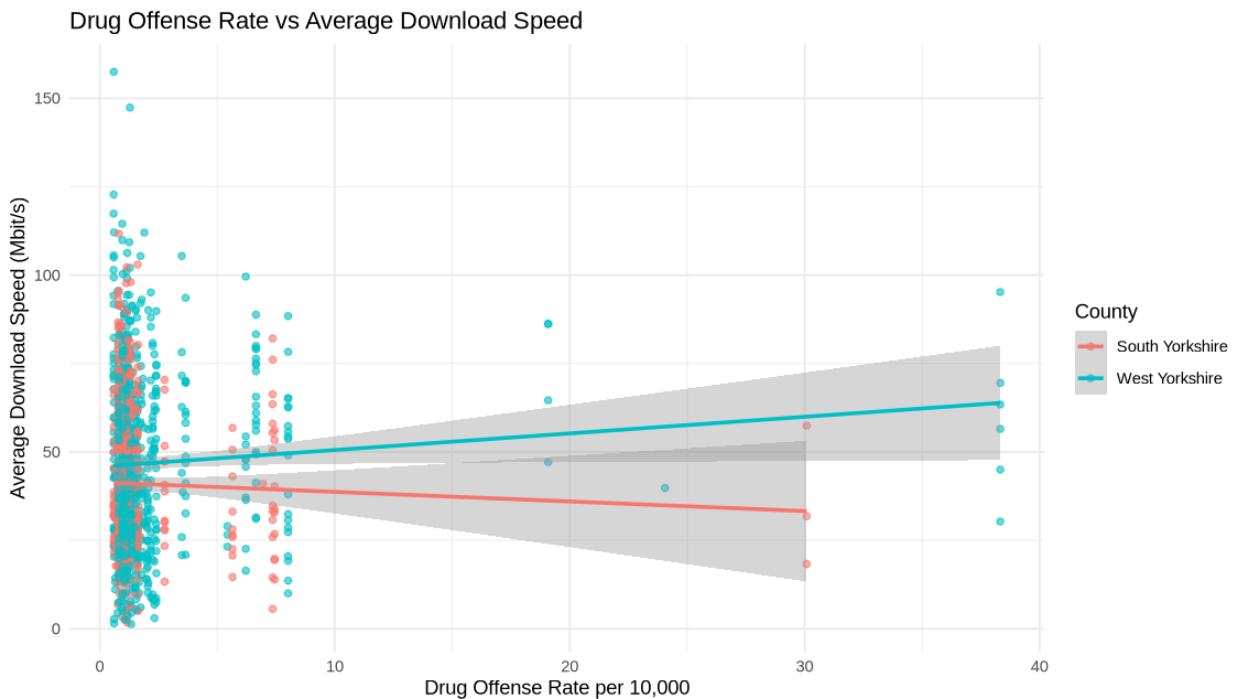
> cat("Correlation (South Yorkshire):", cor_south, "\n")
Correlation (South Yorkshire): 0.8427515

```

Average Download speed vs Drug Offense Rate per 10000 people for both counties

There is a very weak positive correlation ($r = 0.052$) between drug offense rate and average download speed, indicating almost no linear relationship. The regression model also supports this, with a small positive coefficient (0.35) and a p-value of 0.067, which is slightly above the conventional significance threshold of 0.05. This suggests that the relationship is not statistically significant, and drug offense rates do not meaningfully predict download speeds across the regions studied.

Figure 28, Drug Offense vs Average Download Speed



Linear Model summary and correlation

Figure 29, LM 5 Summary

```

Call:
lm(formula = `Average download speed (Mbit/s)` ~ RatePer10k,
    data = final_data)

Residuals:
    Min      1Q  Median      3Q     Max 
-43.102 -17.657 -4.081  15.355 113.357 

Coefficients:
            Estimate Std. Error t value     Pr(>|t|)    
(Intercept) 43.9351    0.7633  57.556 <0.0000000000000002 *** 
RatePer10k   0.3476    0.1899   1.831      0.0674 .  
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 23.3 on 1232 degrees of freedom
Multiple R-squared:  0.002713, Adjusted R-squared:  0.001904 
F-statistic: 3.352 on 1 and 1232 DF,  p-value: 0.06737

> cat("Correlation:", round(correlation, 3), "\n")
Correlation: 0.052

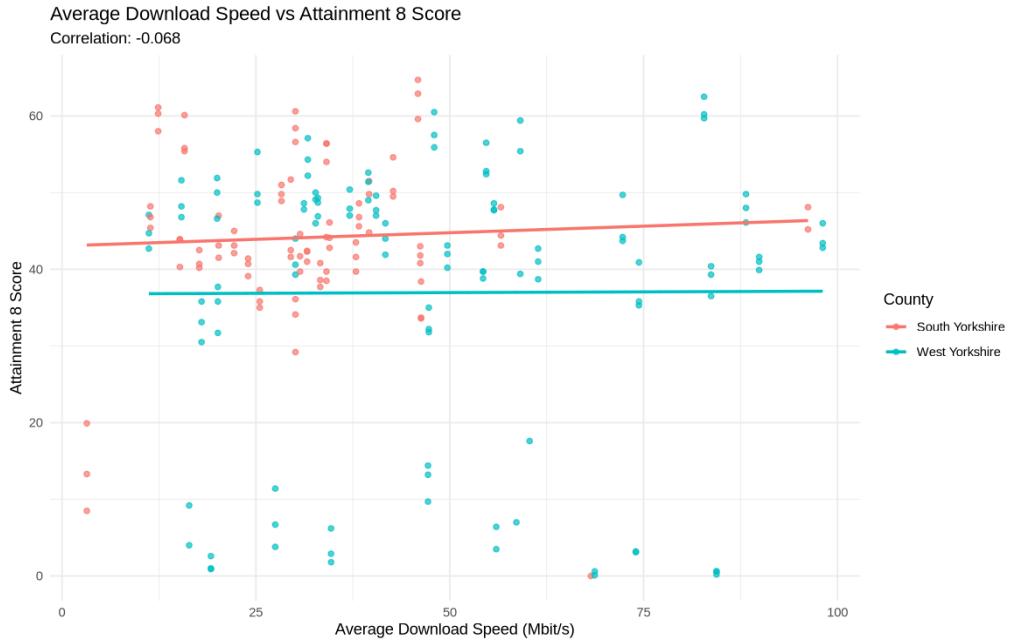
```

Average download speed vs Attainment 8 score for both counties

The correlation between Attainment 8 and Average Download Speed score was very weak and negative (correlation = -0.068), suggesting no meaningful linear relationship.

This was confirmed by Linear regression model. The slope of the model was -0.048 ($p = 0.335$), indicating the association is not statistically significant. The model suggests almost none of the variation in Attainment 8 scores ($R^2 \approx 0.005$), meaning download speed is not a good predictor of academic performance in this dataset.

Figure 30, Average Download Speed vs Attainment 8



Linear Model summary and correlation

Figure 31, LM 6 Summary

```

Call:
lm(formula = ATT8SCR ~ Average.download.speed..Mbit.s., data = filtered_data)

Residuals:
    Min      1Q  Median      3Q     Max 
-40.227 -1.845  3.400  8.963 24.864 

Coefficients:
            Estimate Std. Error t value    Pr(>|t|)    
(Intercept) 42.05521   2.33897 17.980 <0.000000000000002 ***
Average.download.speed..Mbit.s. -0.04835   0.05001 -0.967      0.335  
---
Signif. codes:  0 '****' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 15.76 on 199 degrees of freedom
Multiple R-squared:  0.004675, Adjusted R-squared:  -0.0003264 
F-statistic: 0.9347 on 1 and 199 DF,  p-value: 0.3348

> cat("Correlation: ", correlation, "\n")
Correlation: -0.06837595

```

Recommendation System

Overview

The recommendation system was developed to evaluate and compare towns in South Yorkshire and West Yorkshire based on multiple key factors: house prices, broadband speed, crime rates, and school performance. Each factor was cleaned, standardized, and analyzed to generate comparable scores across towns.

The system uses a weighted scoring method where each town receives a normalized score for each criterion. These scores are then combined to calculate an overall score, helping to identify the most attractive locations for property investment. By integrating data from multiple domains, the system provides a data-driven and balanced approach to support decision-making for potential home buyers or investors.

Results

Based on Broadband speed

Figure 32, Broadband speed result

▲ TOWN ▾	avg_down_speed	normDownSpeed
1 LEEDS	53.21395	10.000000
2 HUDDERSFIELD	52.17439	9.653007
3 ELLAND	52.04573	9.610060
4 BRADFORD	50.48136	9.087893
5 MEXBOROUGH	49.03646	8.605603
6 BRIGHOUSE	46.63752	7.804864
7 BARNSLEY	43.45569	6.742809
8 ROTHERHAM	42.95196	6.574668
9 HALIFAX	41.22150	5.997061
10 NORMANTON	40.54146	5.770073

Based on School Attainment score

Figure 33, School Result

▲	TOWN	avgAtt8	normAtt8
1	HOLMFIRTH	53.38333	10.000000
2	HEBDEN BRIDGE	48.90000	7.787222
3	ELLAND	46.20000	6.454620
4	NORMANTON	44.13333	5.434604
5	SHEFFIELD	44.09559	5.415974
6	HALIFAX	43.51905	5.131419
7	MEXBOROUGH	42.66667	4.710721
8	BRADFORD	42.52875	4.642651
9	SOWERBY BRIDGE	41.24444	4.008774
10	DONCASTER	40.88431	3.831030

Based on House Prices

Figure 34, House prices Result

▲	TOWN	avgPrice	normHousingPrice
1	MEXBOROUGH	156509.8	10.000000
2	BRADFORD	182448.5	8.517212
3	ROTHERHAM	193302.4	7.896743
4	BARNSLEY	193977.6	7.858147
5	HALIFAX	198690.6	7.588729
6	DONCASTER	200276.1	7.498092
7	BRIGHOUSE	204026.7	7.283688
8	PONTEFRACT	222419.6	6.232257
9	OSSETT	222863.5	6.206877
10	SOWERBY BRIDGE	229973.1	5.800456

Based on Crime Data (safest)

Figure 35, Crime Result

	TOWN	normCrimeRate	crimeno
1	SHEFFIELD	10.000000	37
2	ROTHERHAM	9.994925	49
3	HUDDERSFIELD	9.983930	75
4	DONCASTER	9.983084	77
5	HOLMFIRTH	9.982661	78
6	LEEDS	9.981392	81
7	BRIGHOUSE	9.975472	95
8	BRADFORD	9.974626	97
9	WAKEFIELD	9.965745	118
10	MEXBOROUGH	9.965322	119

Overall Score

Figure 36, Overall Score - Final Results

	TOWN	avgPrice	crimeno	avgAtt8	avg_down_speed	finalPoints
1	MEXBOROUGH	156509.8	119	42.66667	49.03646	8.320412
2	BRADFORD	182448.5	97	42.52875	50.48136	8.055596
3	ELLAND	272270.0	330	46.20000	52.04573	7.330829
4	HALIFAX	198690.6	159	43.51905	41.22150	7.166404
5	HUDDERSFIELD	233966.0	75	37.15610	52.17439	6.800021
6	ROTHERHAM	193302.4	49	37.85957	42.95196	6.701122
7	DONCASTER	200276.1	77	40.88431	38.10030	6.566862
8	HOLMFIRTH	308118.4	78	53.38333	36.49605	6.433922
9	BRIGHOUSE	204026.7	95	33.12222	46.63752	6.266006
10	BARNESLEY	193977.6	121	33.61818	43.45569	6.202554

After examining and ranking towns on house price, broadband speed, frequency and rate of crime, and school ratings I am recommending and found the three best towns to buy in are Mexborough, Bradford, and Elland. Mexborough was ranked highest on affordability and broadband reliability, Bradford was also ranked very high due to its excellent infrastructure and connectivity and Elland is a good option due to the quality of schools and lower crime. These three towns offer a good balance of price, livability and buying a house with potential upside from a long-term ownership perspective.

Reflection

This project provided valuable hands-on experience in applying data science techniques to solve a real-world problem. Through data cleaning, integration, analysis, and modeling, I developed a stronger understanding of the entire data science workflow. One key takeaway was the importance of high-quality, well-structured data. A significant portion of the time was spent cleaning and preparing datasets, which proved essential for producing reliable insights.

I additionally learned how different variables interact and how to use linear models to understand a trend or directionality even when the relationships were weak or contrary to expectations. Creating the recommendation system based on multiple weighted factors allowed me to include the analytical side of thinking with the praxis of decision-making. Overall, this assignment reinforced my confidence in using R for data analysis and helped bolster my ability to effectively communicate results visually and descriptively through statistical summaries.

Legal and Ethical Issues

This project relied entirely on publicly available datasets from UK government and open data portals. All data used was collected, processed, and presented in compliance with open data licenses. Personal information or sensitive data was not accessed or used at any point.

From an ethical standpoint, we ensured that we conducted an impartial analysis, with no bias. The recommendation system was based exclusively on a data-screening process and no personal beliefs. Likewise, while some of the neighborhoods might look unfavorable based on their higher crime rates or lower school scores, the report will not objectify a neighborhood, but rather will emphasize only measurable indicators that will vet property investment choices.

Conclusion

This analysis used house prices, broadband speed, crime, and school data to compare locations in South Yorkshire and West Yorkshire. Using data science methods and a scoring system, Mexborough, Bradford, and Elland were identified as the best locations to buy properties, as all provided an affordable, livable, and quality infrastructure combinations. Although not all factors provided strong linear relationships, the multi-criteria assessment provided a more holistic understanding of what makes certain place viable. Ultimately, this project illustrates how data science can contribute to evidence-based decision making and future buyers and investor understanding of the location.

References

- Open Geography Portal.* (n.d.). Open Geography Portal. <https://geoportal.statistics.gov.uk/>
- 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>
- 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>
- Data downloads | data.police.uk. (n.d.). <https://data.police.uk/data/>
- 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-sciencecourse-curriculum/>
- Compare school and college performance in England.* (n.d.). Compare School and College Performance in England. <https://www.compare-school-performance.service.gov.uk/download-data>
- RPUBS - Introduction to Linear Modeling in R. (n.d.). <https://rpubs.com/AnthonyCorbisieri/1100671>
- Radar. (n.d.). <https://plotly.com/r/radar-chart/>
- 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>

Appendix

GitHub Link: <https://github.com/nii-san/Datascience-Assignment>

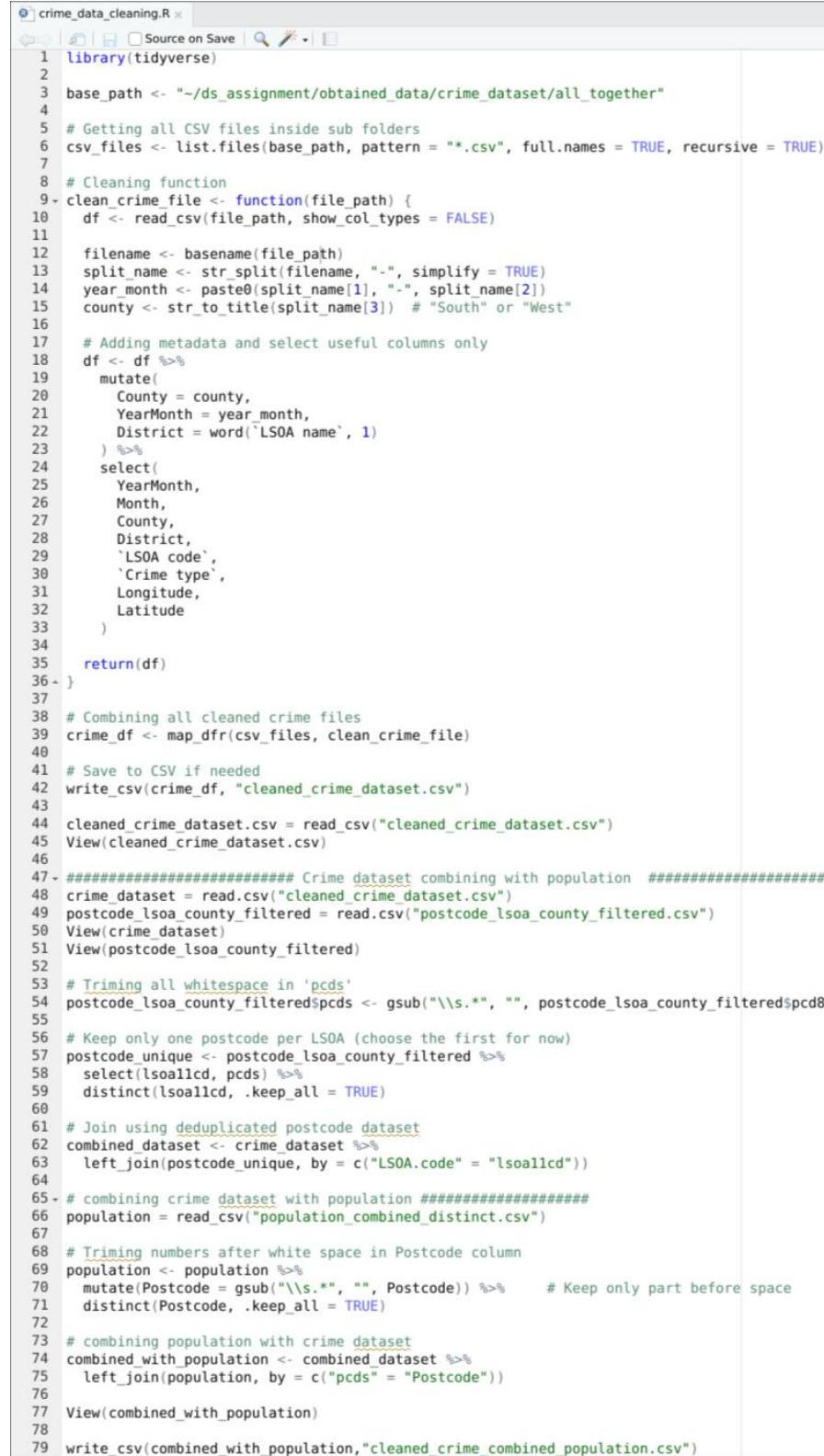
Cleaning.R



```

1 library(tidyverse)
2
3 # Reading Datasets
4 broadband_perf <- read_csv("201805_fixed_pc_performance_r03.csv")
5 broadband_cov <- read_csv("201809_fixed_pc_coverage_r01.csv")
6
7 # Standardizing and cleaning postcodes in broadband data sets
8 broadband_perf <- broadband_perf %>%
9   mutate(postcode = toupper(str_replace_all(postcode, " ", ""))) %>%
10  filter(!is.na(postcode) & postcode != "")
11 broadband_cov <- broadband_cov %>%
12  mutate(postcode = toupper(str_replace_all(postcode, " ", ""))) %>%
13  filter(!is.na(postcode) & postcode != "")
14
15 # Joining broadband performance and coverage data on postcode
16 broadband_combined <- left_join(broadband_perf, broadband_cov, by = "postcode")
17
18 # Reading postcode to LSOA data set
19 postcode_to_lsoa <- read_csv("PostcodeToLSOA.csv")
20 View(postcode_to_lsoa)
21
22 postcode_to_lsoa_clean <- postcode_to_lsoa %>%
23  mutate(pcds = toupper(str_replace_all(pcds, " ", ""))) %>%
24  filter(!is.na(pcds) & pcds != "")
25
26 broadband_with_lsoa <- broadband_combined %>%
27  left_join(postcode_to_lsoa_clean, by = c("postcode" = "pcds"))
28
29 broadband_final <- broadband_with_lsoa %>%
30  filter(!is.na(lsoalcid))
31
32 write_csv(broadband_final, "cleaned_broadband_speed_with_lsoa.csv")
33 View(broadband_final)
34
35 colnames(broadband_final)
36
37 postcode_lsoa_county = read.csv("postcode_lsoa_county.csv")
38 View(postcode_lsoa_county)
39
40 broadband_final <- broadband_final %>%
41  mutate(pcds = str_remove_all(postcode, "\\\s+"))
42
43 postcode_lsoa_county <- postcode_lsoa_county %>%
44  mutate(pcds = str_remove_all(pcds, "\s+"))
45
46 broadband_with_county <- broadband_final %>%
47  inner_join(
48    postcode_lsoa_county %>% select(pcds, County),
49    by = "pcds"
50  )
51 colnames(broadband_with_county)
52
53 # Filtering only South Yorkshire, West Yorkshire and needed columns
54 broadband_filtered <- broadband_with_county %>%
55  filter(County %in% c("South Yorkshire", "West Yorkshire")) %>%
56  select(-`% of premises unable to receive 2Mbit/s`, -`% of premises unable to receive 5Mbit/s`,
57  -`% of premises unable to receive 10Mbit/s`, -`% of premises unable to receive 30Mbit/s`,
58  -`% of premises unable meet USO`, -`% of premises able to receive decent broadband from FWA`,
59  -`% of premises able to receive SFBB from FWA`, -`% of premises able to receive NGA`,
60  -`Average data usage (GB) for lines < 10Mbit/s`, -`UFBB availability (% premises)`, -`SFBB availability (% premises)`,
61  -`FTTP availability (% premises)`, -`Number of connections < 2 Mbit/s (number of lines)`,
62  -`Number of connections 2<= 5 Mbit/s (number of lines)`, -`Number of connections 5<= 10 Mbit/s (number of lines)`,
63  -`Number of connections 10<=30 Mbit/s (number of lines)`, -`Number of connections 30<300 Mbit/s (number of lines)`,
64  -`Number of connections >= 300 Mbit/s (number of lines)`, -`Number of connections >= 30 Mbit/s (number of lines)`,
65  -`Average data usage (GB) for lines < 10Mbit/s`, -`Average data usage (GB) for Basic BB lines`,
66  -`Average data usage (GB) for SFBB lines`, -`Average data usage (GB) for UFBB lines`,
67  , -`All Premises`, -`All Matched Premises`, -`doterm`)
68
69 View(broadband_filtered)
70 write_csv(broadband_filtered, "broadband_filtered_yorkshire.csv")

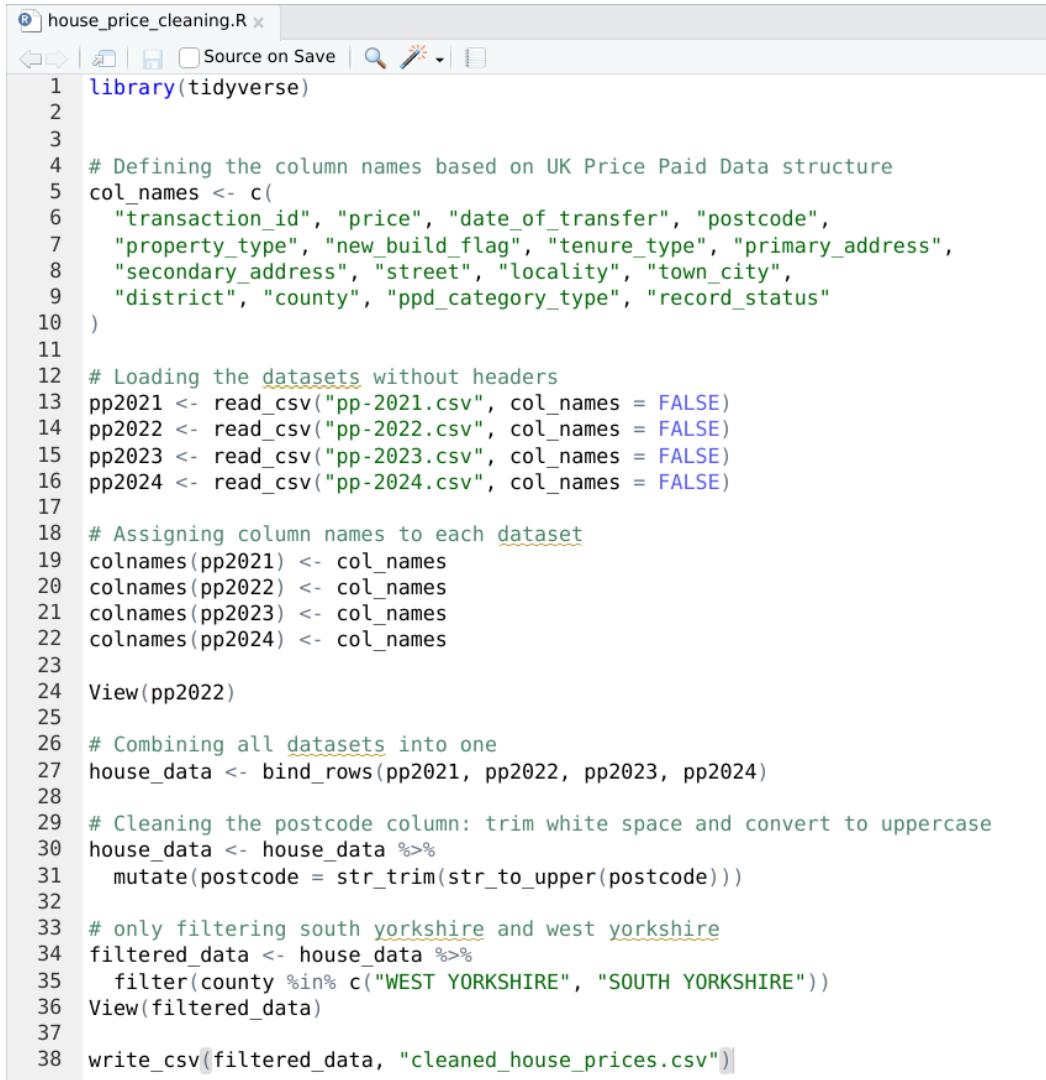
```



```

1 library(tidyverse)
2
3 base_path <- "~/ds_assignment/obtained_data/crime_dataset/all_together"
4
5 # Getting all CSV files inside sub folders
6 csv_files <- list.files(base_path, pattern = "*.csv", full.names = TRUE, recursive = TRUE)
7
8 # Cleaning function
9 clean_crime_file <- function(file_path) {
10   df <- read_csv(file_path, show_col_types = FALSE)
11
12   filename <- basename(file_path)
13   split_name <- str_split(filename, "-", simplify = TRUE)
14   year_month <- paste0(split_name[1], "-", split_name[2])
15   county <- str_to_title(split_name[3]) # "South" or "West"
16
17   # Adding metadata and select useful columns only
18   df <- df %>%
19     mutate(
20       County = county,
21       YearMonth = year_month,
22       District = word(`LSOA name`, 1)
23     ) %>%
24     select(
25       YearMonth,
26       Month,
27       County,
28       District,
29       `LSOA code`,
30       `Crime type`,
31       Longitude,
32       Latitude
33     )
34
35   return(df)
36 }
37
38 # Combining all cleaned crime files
39 crime_df <- map_dfr(csv_files, clean_crime_file)
40
41 # Save to CSV if needed
42 write_csv(crime_df, "cleaned_crime_dataset.csv")
43
44 cleaned_crime_dataset.csv = read_csv("cleaned_crime_dataset.csv")
45 View(cleaned_crime_dataset.csv)
46
47 ##### Crime dataset combining with population #####
48 crime_dataset = read.csv("cleaned_crime_dataset.csv")
49 postcode_lsoa_county_filtered = read.csv("postcode_lsoa_county_filtered.csv")
50 View(crime_dataset)
51 View(postcode_lsoa_county_filtered)
52
53 # Trimming all whitespace in 'pcds'
54 postcode_lsoa_county_filtered$pcds <- gsub("\\s.*", "", postcode_lsoa_county_filtered$pcd8)
55
56 # Keep only one postcode per LSOA (choose the first for now)
57 postcode_unique <- postcode_lsoa_county_filtered %>%
58   select(lsoal1cd, pcds) %>%
59   distinct(lsoal1cd, .keep_all = TRUE)
60
61 # Join using deduplicated postcode dataset
62 combined_dataset <- crime_dataset %>%
63   left_join(postcode_unique, by = c("LSOA.code" = "lsoal1cd"))
64
65 # combining crime dataset with population #####
66 population = read_csv("population_combined_distinct.csv")
67
68 # Trimming numbers after white space in Postcode column
69 population <- population %>%
70   mutate(Postcode = gsub("\\s.*", "", Postcode)) %>%      # Keep only part before space
71   distinct(Postcode, .keep_all = TRUE)
72
73 # combining population with crime dataset
74 combined_with_population <- combined_dataset %>%
75   left_join(population, by = c("pcds" = "Postcode"))
76
77 View(combined_with_population)
78
79 write_csv(combined_with_population, "cleaned_crime_combined_population.csv")

```



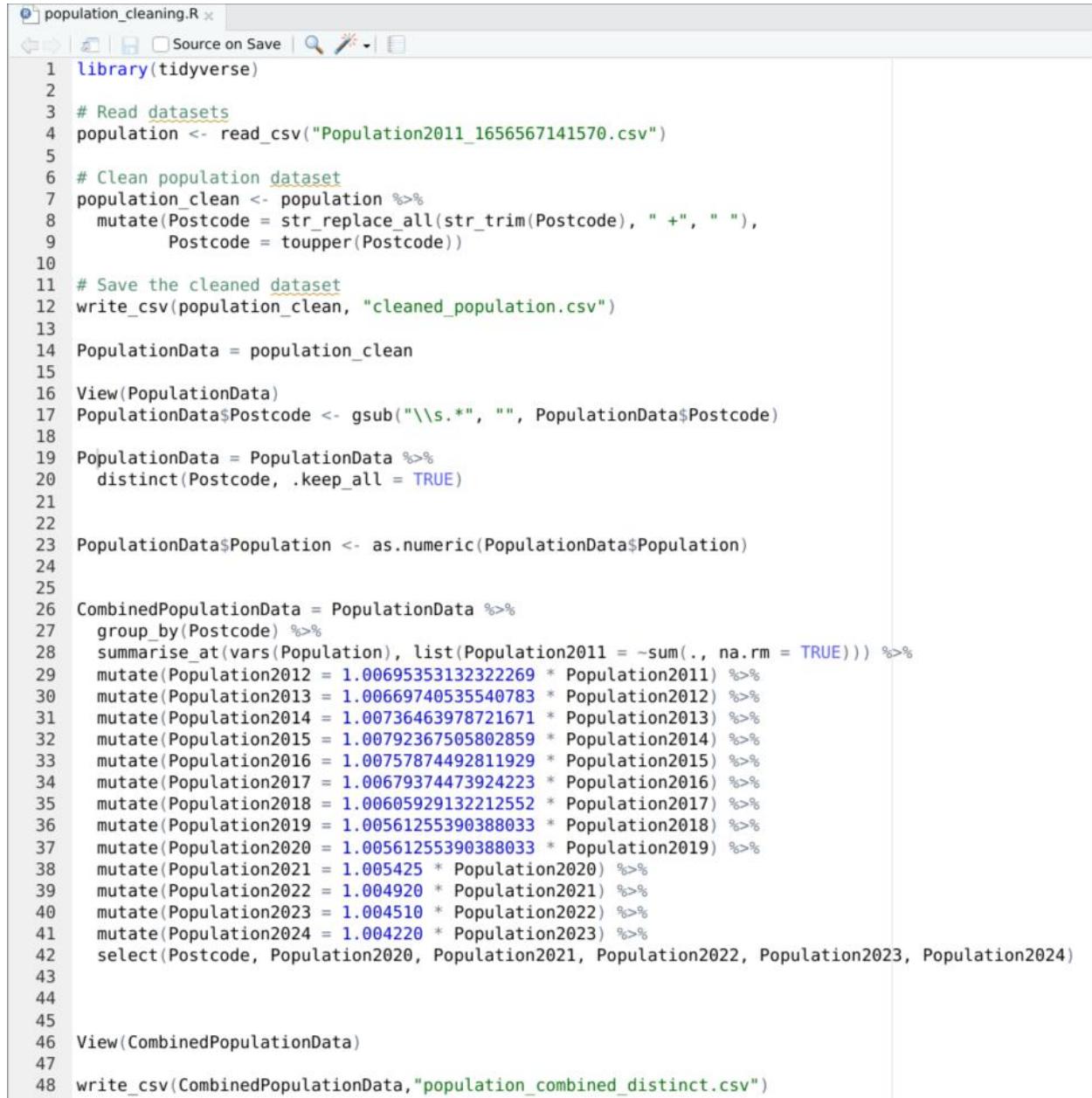
The screenshot shows an RStudio interface with the file 'house_price_cleaning.R' open. The code is as follows:

```
1 library(tidyverse)
2
3
4 # Defining the column names based on UK Price Paid Data structure
5 col_names <- c(
6   "transaction_id", "price", "date_of_transfer", "postcode",
7   "property_type", "new_build_flag", "tenure_type", "primary_address",
8   "secondary_address", "street", "locality", "town_city",
9   "district", "county", "ppd_category_type", "record_status"
10 )
11
12 # Loading the datasets without headers
13 pp2021 <- read_csv("pp-2021.csv", col_names = FALSE)
14 pp2022 <- read_csv("pp-2022.csv", col_names = FALSE)
15 pp2023 <- read_csv("pp-2023.csv", col_names = FALSE)
16 pp2024 <- read_csv("pp-2024.csv", col_names = FALSE)
17
18 # Assigning column names to each dataset
19 colnames(pp2021) <- col_names
20 colnames(pp2022) <- col_names
21 colnames(pp2023) <- col_names
22 colnames(pp2024) <- col_names
23
24 View(pp2022)
25
26 # Combining all datasets into one
27 house_data <- bind_rows(pp2021, pp2022, pp2023, pp2024)
28
29 # Cleaning the postcode column: trim white space and convert to uppercase
30 house_data <- house_data %>%
31   mutate(postcode = str_trim(str_to_upper(postcode)))
32
33 # only filtering south yorkshire and west yorkshire
34 filtered_data <- house_data %>%
35   filter(county %in% c("WEST YORKSHIRE", "SOUTH YORKSHIRE"))
36 View(filtered_data)
37
38 write_csv(filtered_data, "cleaned_house_prices.csv")
```

```

1 library(tidyverse)
2
3 ks4final_2021_2022 = read_csv("england_ks4final_2021_2022.csv")
4 ks4final_2022_2023 = read_csv("england_ks4final_2022_2023.csv")
5 ks4final_2023_2024 = read_csv("england_ks4final_2023_2024.csv")
6
7 postcode_lsoa_county = read_csv("postcode_lsoa_county_filtered.csv")
8
9 test = ks4final_2021_2022 %>%
10   distinct(TOWN, .keep_all = TRUE)
11
12 View(ks4final_2021_2022)
13 View(ks4final_2022_2023)
14 View(ks4final_2023_2024)
15
16 # Adding year and selecting relevant columns
17 ks4_21_22 <- ks4final_2021_2022 %>%
18   select(URN, SCHNAME, PCODE, TOWN, ATT8SCR) %>%
19   mutate(Year = "2022")
20 ks4_22_23 <- ks4final_2022_2023 %>%
21   select(URN, SCHNAME, PCODE, TOWN, ATT8SCR) %>%
22   mutate(Year = "2023")
23 ks4_23_24 <- ks4final_2023_2024 %>%
24   select(URN, SCHNAME, PCODE, TOWN, ATT8SCR) %>%
25   mutate(Year = "2024")
26
27 # Combining all three years ks4 dataset
28 ks4_combined <- bind_rows(ks4_21_22, ks4_22_23, ks4_23_24)
29
30 # Cleaning ATT8SCR by removing 'NE', 'SUPP', and non-numeric values, then remove NAs
31 ks4_clean <- ks4_combined %>%
32   filter(!ATT8SCR %in% c("NE", "SUPP")) %>%
33   filter(!is.na(ATT8SCR)) %>%
34   filter(str_detect(ATT8SCR, "[0-9]+"))
35
36 # Converting ATT8SCR to numeric for analysis/visualization
37 ks4_clean <- ks4_clean %>%
38   mutate(ATT8SCR = as.numeric(ATT8SCR))
39
40 postcode_lsoa_county = read_csv("postcode_lsoa_county.csv")
41
42 # Trimming whitespace in postcode columns
43 ks4_clean <- ks4_clean %>%
44   mutate(PCODE = str_trim(PCODE))
45 postcode_lsoa_county <- postcode_lsoa_county %>%
46   mutate(pcds = str_trim(pcds))
47
48 # Joining to add 'County'
49 ks4_with_county <- ks4_clean %>%
50   left_join(postcode_lsoa_county %>% select(pcds, County), by = c("PCODE" = "pcds"))
51
52 View(ks4_with_county)
53
54 # Filtering only West Yorkshire and South Yorkshire
55 ks4_filtered <- ks4_with_county %>%
56   filter(County %in% c("West Yorkshire", "South Yorkshire"))
57
58 View(ks4_filtered)
59
60 # Joining ladnm to school data set
61 school_data <- inner_join(
62   school_data,
63   postcode_lsoa_county %>% select(pcd7, ladnm),
64   by = c("PCODE" = "pcd7")
65 )
66
67 write.csv(ks4_filtered, "cleaned_filtered_school_dataset.csv", row.names = FALSE)

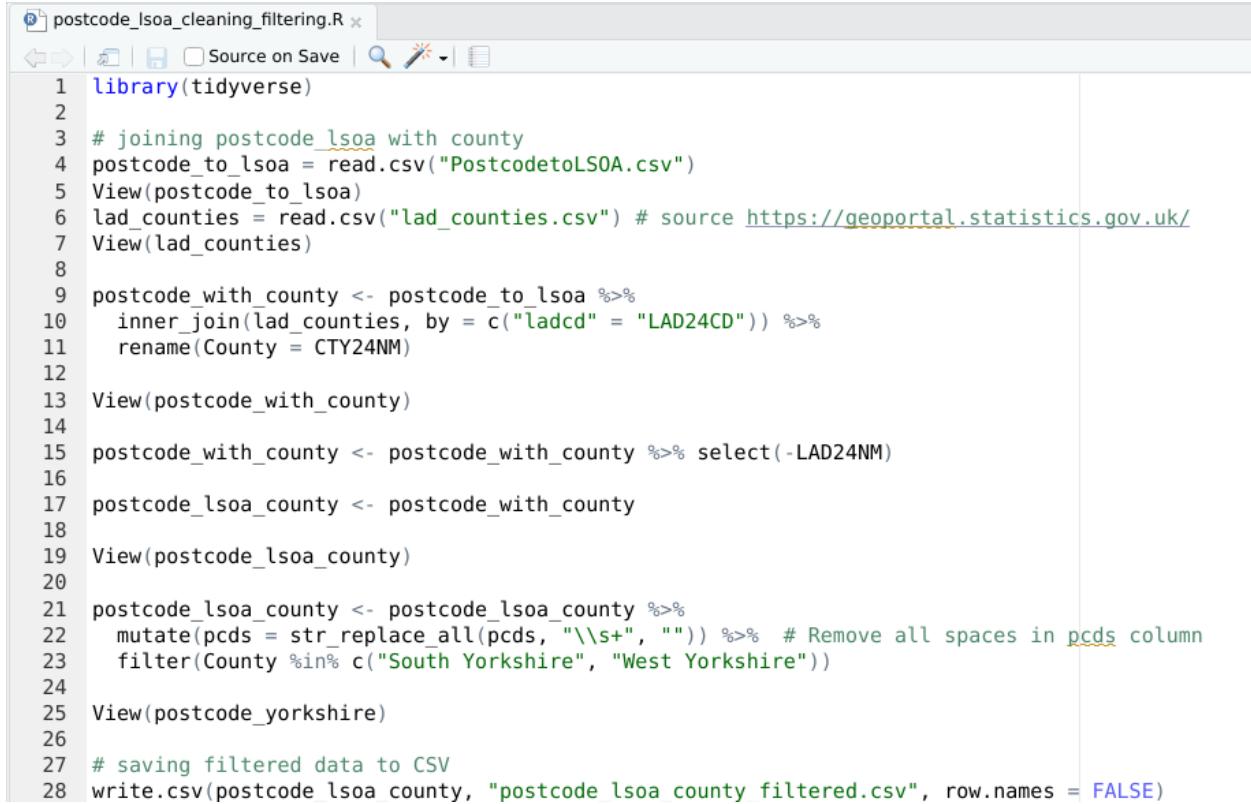
```



```

1 library(tidyverse)
2
3 # Read datasets
4 population <- read_csv("Population2011_1656567141570.csv")
5
6 # Clean population dataset
7 population_clean <- population %>%
8   mutate(Postcode = str_replace_all(str_trim(Postcode), " +", " "),
9         Postcode = toupper(Postcode))
10
11 # Save the cleaned dataset
12 write_csv(population_clean, "cleaned_population.csv")
13
14 PopulationData = population_clean
15
16 View(PopulationData)
17 PopulationData$Postcode <- gsub("\s.*", "", PopulationData$Postcode)
18
19 PopulationData = PopulationData %>%
20   distinct(Postcode, .keep_all = TRUE)
21
22
23 PopulationData$Population <- as.numeric(PopulationData$Population)
24
25
26 CombinedPopulationData = PopulationData %>%
27   group_by(Postcode) %>%
28   summarise_at(vars(Population), list(Population2011 = ~sum(., na.rm = TRUE))) %>%
29   mutate(Population2012 = 1.00695353132322269 * Population2011) %>%
30   mutate(Population2013 = 1.00669740535540783 * Population2012) %>%
31   mutate(Population2014 = 1.00736463978721671 * Population2013) %>%
32   mutate(Population2015 = 1.00792367505802859 * Population2014) %>%
33   mutate(Population2016 = 1.00757874492811929 * Population2015) %>%
34   mutate(Population2017 = 1.00679374473924223 * Population2016) %>%
35   mutate(Population2018 = 1.00605929132212552 * Population2017) %>%
36   mutate(Population2019 = 1.00561255390388033 * Population2018) %>%
37   mutate(Population2020 = 1.00561255390388033 * Population2019) %>%
38   mutate(Population2021 = 1.005425 * Population2020) %>%
39   mutate(Population2022 = 1.004920 * Population2021) %>%
40   mutate(Population2023 = 1.004510 * Population2022) %>%
41   mutate(Population2024 = 1.004220 * Population2023) %>%
42   select(Postcode, Population2020, Population2021, Population2022, Population2023, Population2024)
43
44
45
46 View(CombinedPopulationData)
47
48 write_csv(CombinedPopulationData, "population_combined_distinct.csv")

```



The screenshot shows an RStudio interface with the following details:

- Title Bar:** Shows the file name "postcode_lsoa_cleaning_filtering.R".
- Toolbar:** Includes standard icons for file operations (New, Open, Save, Print), a "Source on Save" checkbox, and search/filter tools.
- Code Editor:** Displays the R code for data cleaning and filtering. The code uses the tidyverse library and performs the following steps:
 - Imports the tidyverse library.
 - Reads "PostcodetoLSOA.csv" into "postcode_to_lsoa".
 - Views "postcode_to_lsoa".
 - Reads "lad_counties.csv" into "lad_counties".
 - Views "lad_counties".
 - Joins "postcode_to_lsoa" with "lad_counties" on the "ladcd" column, renaming the "ladcd" column to "CTY24NM".
 - Views the resulting "postcode_with_county" data.
 - Selects all columns except "CTY24NM" from "postcode_with_county".
 - Creates "postcode_lsoa_county" by selecting the first column of "postcode_with_county".
 - Views "postcode_lsoa_county".
 - Creates "postcode_lsoa_county" by selecting the first column of "postcode_lsoa_county".
 - Modifies "pcds" column by removing all spaces using str_replace_all and filters for "South Yorkshire" and "West Yorkshire".
 - Views the resulting "postcode_yorkshire" data.
 - Saves the filtered data to a CSV file named "postcode_lsoa_county_filtered.csv".

Eda.R

```

broadbandspeed.R x
Source on Save | 🔎 | 🖌️ | 📁

1 library(tidyverse)
2
3 broadband = read.csv("broadband_filtered_yorkshire.csv")
4 colnames(broadband)
5
6 # Remove rows with missing or invalid average download speed
7 broadband <- broadband %>%
8   filter(!is.na(Average.download.speed..Mbit.s.))
9
10
11 # Boxplot of Avg Download Speed by District
12 # West Yorkshire boxplot
13 broadband %>%
14   filter(County == "West Yorkshire") %>%
15   ggplot(aes(x = fct_reorder(ladnm.x, Average.download.speed..Mbit.s., .fun = median),
16             y = Average.download.speed..Mbit.s.)) +
17   geom_boxplot(fill = "skyblue") +
18   coord_flip() +
19   labs(title = "West Yorkshire - Avg Download Speed by District",
20        x = "District", y = "Avg Download Speed (Mbps)") +
21   theme_minimal()
22
23 # South Yorkshire boxplot
24 broadband %>%
25   filter(County == "South Yorkshire") %>%
26   ggplot(aes(x = fct_reorder(ladnm.x, Average.download.speed..Mbit.s., .fun = median),
27             y = Average.download.speed..Mbit.s.)) +
28   geom_boxplot(fill = "salmon") +
29   coord_flip() +
30   labs(title = "South Yorkshire - Avg Download Speed by District",
31        x = "District", y = "Avg Download Speed (Mbps)") +
32   theme_minimal()
33
34 # Bar Chart of Avg Download Speed by Town (from lsoallnm)
35 # West Yorkshire
36 broadband %>%
37   filter(County == "West Yorkshire") %>%
38   group_by(ladnm.x) %>%
39   summarise(avg_speed = mean(Average.download.speed..Mbit.s., na.rm = TRUE)) %>%
40   arrange(desc(avg_speed)) %>%
41   slice_head(n = 15) %>%
42   ggplot(aes(x = fct_reorder(ladnm.x, avg_speed), y = avg_speed)) +
43   geom_col(fill = "steelblue") +
44   labs(title = "West Yorkshire Towns by Avg Download Speed",
45        x = "Town", y = "Avg Download Speed (Mbps)") +
46   theme_minimal()
47
48 # South Yorkshire
49 broadband %>%
50   filter(County == "South Yorkshire") %>%
51   group_by(ladnm.x) %>%
52   summarise(avg_speed = mean(Average.download.speed..Mbit.s., na.rm = TRUE)) %>%
53   arrange(desc(avg_speed)) %>%
54   slice_head(n = 15) %>%
55   ggplot(aes(x = fct_reorder(ladnm.x, avg_speed), y = avg_speed)) +
56   geom_col(fill = "coral") +
57   labs(title = "South Yorkshire Towns by Avg Download Speed",
58        x = "Town", y = "Avg Download Speed (Mbps)") +
59   theme_minimal()

```

```

1 library(tidyverse)
2 library(fmsb)
3
4 # Read data
5 crime_data <- read.csv("cleaned_crime_combined_population.csv") %>%
6   filter(Year != 2025) %>%
7   filter(!(County == "South" & District == "Leeds")) %>%
8   filter(!(County == "West" & District == "Barnsley")) %>%
9   distinct()
10
11 # ---- 1. Boxplot for Drug Offense Rate by District (2 graphs) ----
12 # Calculate drug rate per 100 total crimes per District-Year
13 drug_data <- crime_data %>%
14   filter(Crime.type == "Drugs") %>%
15   group_by(County, District, Year) %>%
16   summarise(drug_count = n(), .groups = "drop")
17
18 total_data <- crime_data %>%
19   group_by(County, District, Year) %>%
20   summarise(total_crimes = n(), .groups = "drop")
21 drug_rate <- drug_data %>%
22   left_join(total_data, by = c("County", "District", "Year")) %>%
23   mutate(rate_percent = (drug_count / total_crimes) * 100)
24
25 # Separate data
26 drug_south <- drug_rate %>% filter(County == "South")
27 drug_west <- drug_rate %>% filter(County == "West")
28
29 # Plot 1 - South
30 ggplot(drug_south, aes(x = District, y = rate_percent)) +
31   geom_boxplot(fill = "#2E86C1") +
32   labs(title = "Drug Offense Rate in South Yorkshire Districts",
33        y = "Drug Rate (%)", x = "District") +
34   theme(axis.text.x = element_text(angle = 45, hjust = 1)) +
35   coord_flip() + theme_minimal()
36
37 # Plot 2 - West
38 ggplot(drug_west, aes(x = District, y = rate_percent)) +
39   geom_boxplot(fill = "#C0392B") +
40   labs(title = "Drug Offense Rate in West Yorkshire Districts",
41        y = "Drug Rate (%)", x = "District") +
42   theme(axis.text.x = element_text(angle = 45, hjust = 1)) +
43   coord_flip() + theme_minimal()
44
45 # ---- 2. Radar Chart for Vehicle Crime Rate (West, May 2022) ----
46 radar_data <- crime_data %>%
47   filter(County == "West", Year == 2022, Month == 5, Crime.type == "Vehicle crime") %>%
48   group_by(District) %>%
49   summarise(vehicle_crime = n(), .groups = "drop")
50
51 # Prepare for radar
52 radar_plot_data <- radar_data %>%
53   column_to_rownames("District") %>%
54   t() %>%
55   as.data.frame()
56
57 # Add max-min rows required by fmsb
58 radar_plot_data <- rbind(rep(max(radar_plot_data), ncol(radar_plot_data)),
59                           rep(0, ncol(radar_plot_data)),
60                           radar_plot_data)
61
62 # Plot radar
63 radarchart(radar_plot_data, axistype = 1,
64             pcol = "#34495E", pfcol = scales::alpha("#34495E", 0.6), plwd = 2,
65             cglcol = "grey", axislabcol = "black",
66             title = "Vehicle Crime Rate by District\n(West Yorkshire, May 2022)")
67
68 # ---- 3. Pie Chart for Robbery (South, June 2022) ----
69 robbery_data <- crime_data %>%
70   filter(County == "West", Year == 2022, Month == 6, Crime.type == "Robbery") %>%
71   group_by(District) %>%
72   summarise(robbery_count = n(), .groups = "drop")

```

```

71 # Plot pie chart
72 ggplot(robbery_data, aes(x = "", y = robbery_count, fill = District)) +
73   geom_bar(stat = "identity", width = 1) +
74   coord_polar("y") +
75   labs(title = "Robbery Distribution by District\n(West Yorkshire, June 2022)") +
76   theme_void()
77
78 # ---- 4. Line Chart for Drug Offense Rate per 10,000 People ----
79 # Filter for drug crimes only
80 drug_data <- crime_data %>%
81   filter(Crime.type == "Drugs")
82 # Ensure Year is numeric if not already
83 drug_data$Year <- as.numeric(drug_data$Year)
84 # Add population per year
85 drug_data <- drug_data %>%
86   mutate(Population = case_when(
87     Year == 2020 ~ Population2020,
88     Year == 2021 ~ Population2021,
89     Year == 2022 ~ Population2022,
90     Year == 2023 ~ Population2023,
91     Year == 2024 ~ Population2024,
92     TRUE ~ NA_real_
93   ))
94
95 # Group by County, Year, and Month
96 district_rate_yearly <- drug_data %>%
97   group_by(County, Year, Month) %>%
98   summarise(
99     drug_offense_count = n(),
100    Population = sum(Population, na.rm = TRUE),
101    .groups = "drop"
102  ) %>%
103  mutate(
104    drug_rate_per_10k = (drug_offense_count / Population) * 10000,
105    date = as.Date(paste(Year, Month, "01", sep = "-"))
106  )
107
108 # Plot the monthly drug offense rate with proper date formatting
109 ggplot(district_rate_yearly, aes(x = date, y = drug_rate_per_10k, color = County)) +
110   geom_line(size = 1) +
111   geom_point(size = 2) +
112   scale_x_date(
113     date_breaks = "3 months",
114     date_labels = "%Y-%m"
115   ) +
116   labs(
117     title = "Drug Offense Trend of South & West Yorkshire (2022–2024)",
118     x = "Month",
119     y = "Drug Offense Rate (/10,000)",
120     color = "County"
121   ) +
122   theme_minimal() +
123   theme(
124     legend.position = "bottom"  )

```

```

house_prices.R x
Source on Save | 🔍 | 🖌️ | 📁

1 library(tidyverse)
2
3 house_prices = read_csv("cleaned_house_prices.csv")
4 colnames(house_prices)
5
6 # Extract year
7 house_prices <- house_prices %>%
8   mutate(year = year(as.Date(date_of_transfer)))
9 # Filter years and counties
10 filtered_hp <- house_prices %>%
11   filter(year %in% 2021:2024,
12         county %in% c("WEST YORKSHIRE", "SOUTH YORKSHIRE"))
13 # Group and summarize
14 avg_prices_by_year_county <- filtered_hp %>%
15   group_by(year, county) %>%
16   summarise(avg_price = mean(price, na.rm = TRUE), .groups = "drop")
17 # Plot line graph
18 ggplot(avg_prices_by_year_county, aes(x = year, y = avg_price, color = county)) +
19   geom_line(size = 1.2) +
20   geom_point() +
21   labs(
22     title = "Average House Prices (2021–2024)",
23     x = "Year",
24     y = "Average Price",
25     color = "County"
26   ) + theme_minimal()
27
28 # bar chart
29 # Filter for 2023 only
30 avg_2023 <- filtered_hp %>%
31   filter(year == 2023) %>%
32   group_by(county) %>%
33   summarise(avg_price = mean(price, na.rm = TRUE), .groups = "drop")
34
35 # Plot bar chart
36 ggplot(avg_2023, aes(x = county, y = avg_price, fill = county)) +
37   geom_col(width = 0.6) +
38   labs(
39     title = "Average House Prices in 2023",
40     x = "County",
41     y = "Average Price"
42   ) +
43   theme_minimal() +
44   theme(legend.position = "none")
45
46 # Boxplot for average house prices for both counties in separate diagrams (take
47 # variables Price and District)
48 # Calculate IQR and filter out outliers
49 filtered_hp_clean <- filtered_hp %>%
50   group_by(county) %>%
51   mutate(
52     Q1 = quantile(price, 0.25, na.rm = TRUE),
53     Q3 = quantile(price, 0.75, na.rm = TRUE),
54     IQR = Q3 - Q1,
55     lower_bound = Q1 - 1.5 * IQR,
56     upper_bound = Q3 + 1.5 * IQR
57   ) %>%
58   filter(price >= lower_bound & price <= upper_bound) %>%
59   ungroup()
60
61 # Split datasets
62 west_yorkshire <- filtered_hp_clean %>% filter(county == "WEST YORKSHIRE")
63 south_yorkshire <- filtered_hp_clean %>% filter(county == "SOUTH YORKSHIRE")
64
65 # Plot for West Yorkshire
66 ggplot(west_yorkshire, aes(x = district, y = price, fill = district)) +
67   geom_boxplot() +
68   labs(
69     title = "House Prices by District – West Yorkshire",
70     x = "District",

```

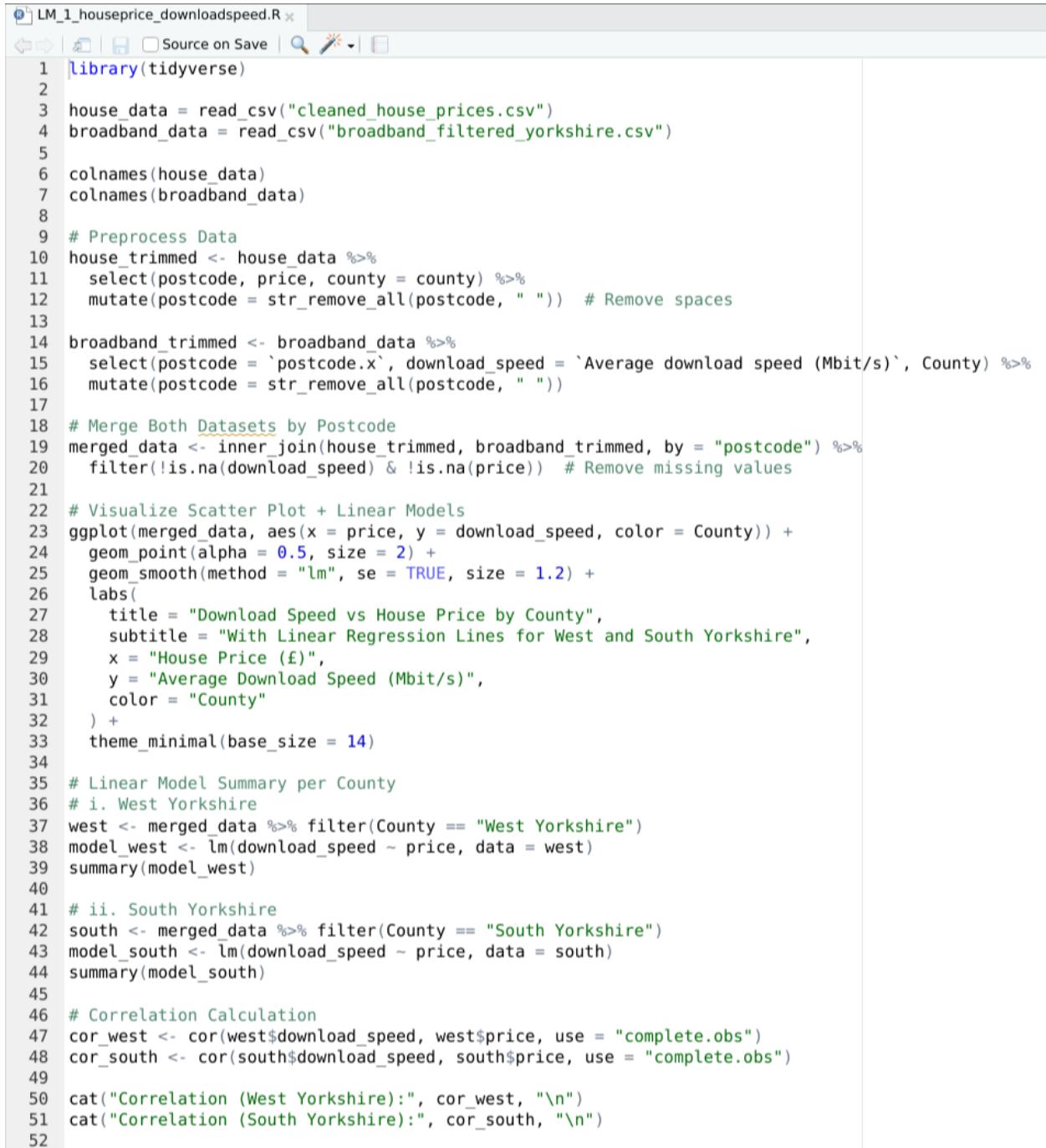
```
70     x = "District",
71     y = "Price"
72   ) +
73   theme_minimal() +
74   coord_flip() +
75   theme(
76     legend.position = "none"
77   )
78
79 # Plot for South Yorkshire
80 ggplot(south_yorkshire, aes(x = district, y = price, fill = district)) +
81   geom_boxplot() +
82   labs(
83     title = "House Prices by District – South Yorkshire",
84     x = "District",
85     y = "Price"
86   ) +
87   theme_minimal() +
88   coord_flip() +
89   theme(
90     legend.position = "none"
91   )
```

```

1 library(tidyverse)
2
3 school_data = read.csv("cleaned_filtered_school_dataset.csv")
4 View(school_data)
5 colnames(school_data)
6
7 # Boxplot: South Yorkshire – Attainment 8 Score by District (TOWN)
8 # South Yorkshire boxplot
9 school_data %>%
10   filter(County == "South Yorkshire", Year == 2022, !is.na(ladnm)) %>%
11   ggplot(aes(x = fct_reorder(ladnm, ATT8SCR), y = ATT8SCR)) +
12   geom_boxplot(fill = "steelblue") +
13   labs(
14     title = "Attainment 8 Score by District (2022) - South Yorkshire",
15     x = "Town", y = "Attainment 8 Score"
16   ) +
17   theme_minimal() +
18   coord_flip()
19
20 # Boxplot for West Yorkshire schools (2022 only)
21 school_data %>%
22   filter(County == "West Yorkshire", Year == 2022, !is.na(ladnm)) %>%
23   ggplot(aes(x = fct_reorder(ladnm, ATT8SCR), y = ATT8SCR)) +
24   geom_boxplot(fill = "darkseagreen") +
25   labs(
26     title = "Attainment 8 Score by District (2022) - West Yorkshire",
27     x = "Town", y = "Attainment 8 Score"
28   ) +
29   theme_minimal() +
30   coord_flip()
31
32 # Line Graph to show the relationship between attainment 8 score and years over
33 # multiple districts in South Yorkshire and west Yorkshire|
34 # West Yorkshire Line Graph
35 school_data %>%
36   filter(County %in% c("South Yorkshire", "West Yorkshire"), !is.na(ladnm)) %>%
37   group_by(County, ladnm, Year) %>%
38   summarise(avg_attainment = mean(ATT8SCR, na.rm = TRUE), .groups = "drop") %>%
39   mutate(ladnm_full = paste0(ladnm, " (", County, ")")) %>%
40   ggplot(aes(x = ladnm_full, y = avg_attainment, group = Year, color = as.factor(Year))) +
41   geom_line(aes(linetype = as.factor(Year)), size = 1) +
42   geom_point(size = 2) +
43   labs(
44     title = "Attainment 8 Score by District (South & West Yorkshire)",
45     x = "Town (County)",
46     y = "Average Attainment 8 Score",
47     color = "Year",
48     linetype = "Year"
49   ) +
50   theme_minimal() +
51   theme(axis.text.x = element_text(angle = 45, hjust = 1))

```

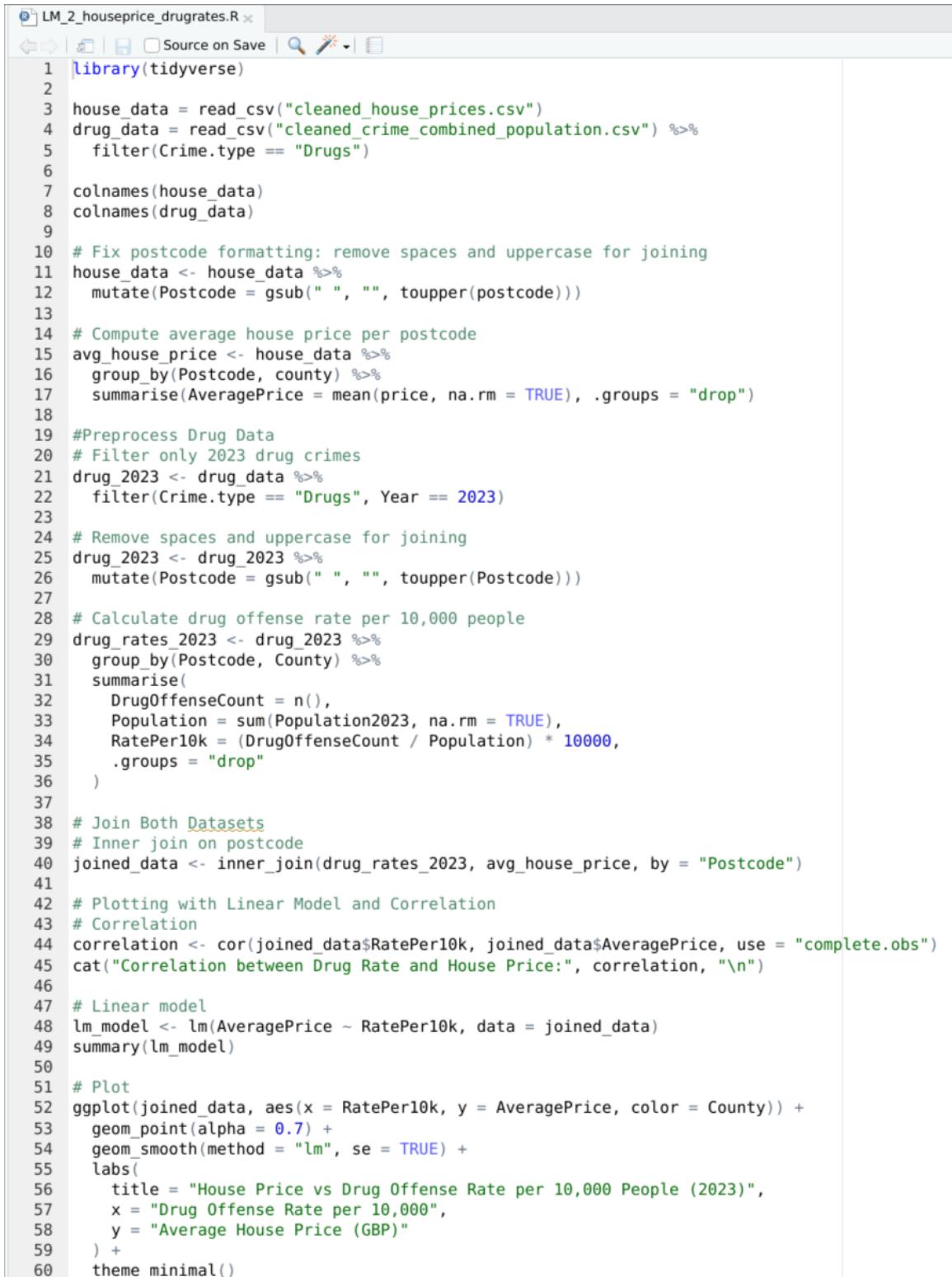
LinearModeling.R



```

1 library(tidyverse)
2
3 house_data = read_csv("cleaned_house_prices.csv")
4 broadband_data = read_csv("broadband_filtered_yorkshire.csv")
5
6 colnames(house_data)
7 colnames(broadband_data)
8
9 # Preprocess Data
10 house_trimmed <- house_data %>%
11   select(postcode, price, county = county) %>%
12   mutate(postcode = str_remove_all(postcode, " ")) # Remove spaces
13
14 broadband_trimmed <- broadband_data %>%
15   select(postcode = `postcode.x`, download_speed = `Average download speed (Mbit/s)`, County) %>%
16   mutate(postcode = str_remove_all(postcode, " "))
17
18 # Merge Both Datasets by Postcode
19 merged_data <- inner_join(house_trimmed, broadband_trimmed, by = "postcode") %>%
20   filter(!is.na(download_speed) & !is.na(price)) # Remove missing values
21
22 # Visualize Scatter Plot + Linear Models
23 ggplot(merged_data, aes(x = price, y = download_speed, color = County)) +
24   geom_point(alpha = 0.5, size = 2) +
25   geom_smooth(method = "lm", se = TRUE, size = 1.2) +
26   labs(
27     title = "Download Speed vs House Price by County",
28     subtitle = "With Linear Regression Lines for West and South Yorkshire",
29     x = "House Price (£)",
30     y = "Average Download Speed (Mbit/s)",
31     color = "County"
32   ) +
33   theme_minimal(base_size = 14)
34
35 # Linear Model Summary per County
36 # i. West Yorkshire
37 west <- merged_data %>% filter(County == "West Yorkshire")
38 model_west <- lm(download_speed ~ price, data = west)
39 summary(model_west)
40
41 # ii. South Yorkshire
42 south <- merged_data %>% filter(County == "South Yorkshire")
43 model_south <- lm(download_speed ~ price, data = south)
44 summary(model_south)
45
46 # Correlation Calculation
47 cor_west <- cor(west$download_speed, west$price, use = "complete.obs")
48 cor_south <- cor(south$download_speed, south$price, use = "complete.obs")
49
50 cat("Correlation (West Yorkshire):", cor_west, "\n")
51 cat("Correlation (South Yorkshire):", cor_south, "\n")
52

```



```

1 library(tidyverse)
2
3 house_data = read_csv("cleaned_house_prices.csv")
4 drug_data = read_csv("cleaned_crime_combined_population.csv") %>%
5   filter(Crime.type == "Drugs")
6
7 colnames(house_data)
8 colnames(drug_data)
9
10 # Fix postcode formatting: remove spaces and uppercase for joining
11 house_data <- house_data %>%
12   mutate(Postcode = gsub(" ", "", toupper(postcode)))
13
14 # Compute average house price per postcode
15 avg_house_price <- house_data %>%
16   group_by(Postcode, county) %>%
17   summarise(AveragePrice = mean(price, na.rm = TRUE), .groups = "drop")
18
19 #Preprocess Drug Data
20 # Filter only 2023 drug crimes
21 drug_2023 <- drug_data %>%
22   filter(Crime.type == "Drugs", Year == 2023)
23
24 # Remove spaces and uppercase for joining
25 drug_2023 <- drug_2023 %>%
26   mutate(Postcode = gsub(" ", "", toupper(Postcode)))
27
28 # Calculate drug offense rate per 10,000 people
29 drug_rates_2023 <- drug_2023 %>%
30   group_by(Postcode, County) %>%
31   summarise(
32     DrugOffenseCount = n(),
33     Population = sum(Population2023, na.rm = TRUE),
34     RatePer10k = (DrugOffenseCount / Population) * 10000,
35     .groups = "drop"
36   )
37
38 # Join Both Datasets
39 # Inner join on postcode
40 joined_data <- inner_join(drug_rates_2023, avg_house_price, by = "Postcode")
41
42 # Plotting with Linear Model and Correlation
43 # Correlation
44 correlation <- cor(joined_data$RatePer10k, joined_data$AveragePrice, use = "complete.obs")
45 cat("Correlation between Drug Rate and House Price:", correlation, "\n")
46
47 # Linear model
48 lm_model <- lm(AveragePrice ~ RatePer10k, data = joined_data)
49 summary(lm_model)
50
51 # Plot
52 ggplot(joined_data, aes(x = RatePer10k, y = AveragePrice, color = County)) +
53   geom_point(alpha = 0.7) +
54   geom_smooth(method = "lm", se = TRUE) +
55   labs(
56     title = "House Price vs Drug Offense Rate per 10,000 People (2023)",
57     x = "Drug Offense Rate per 10,000",
58     y = "Average House Price (GBP)"
59   ) +
60   theme_minimal()

```

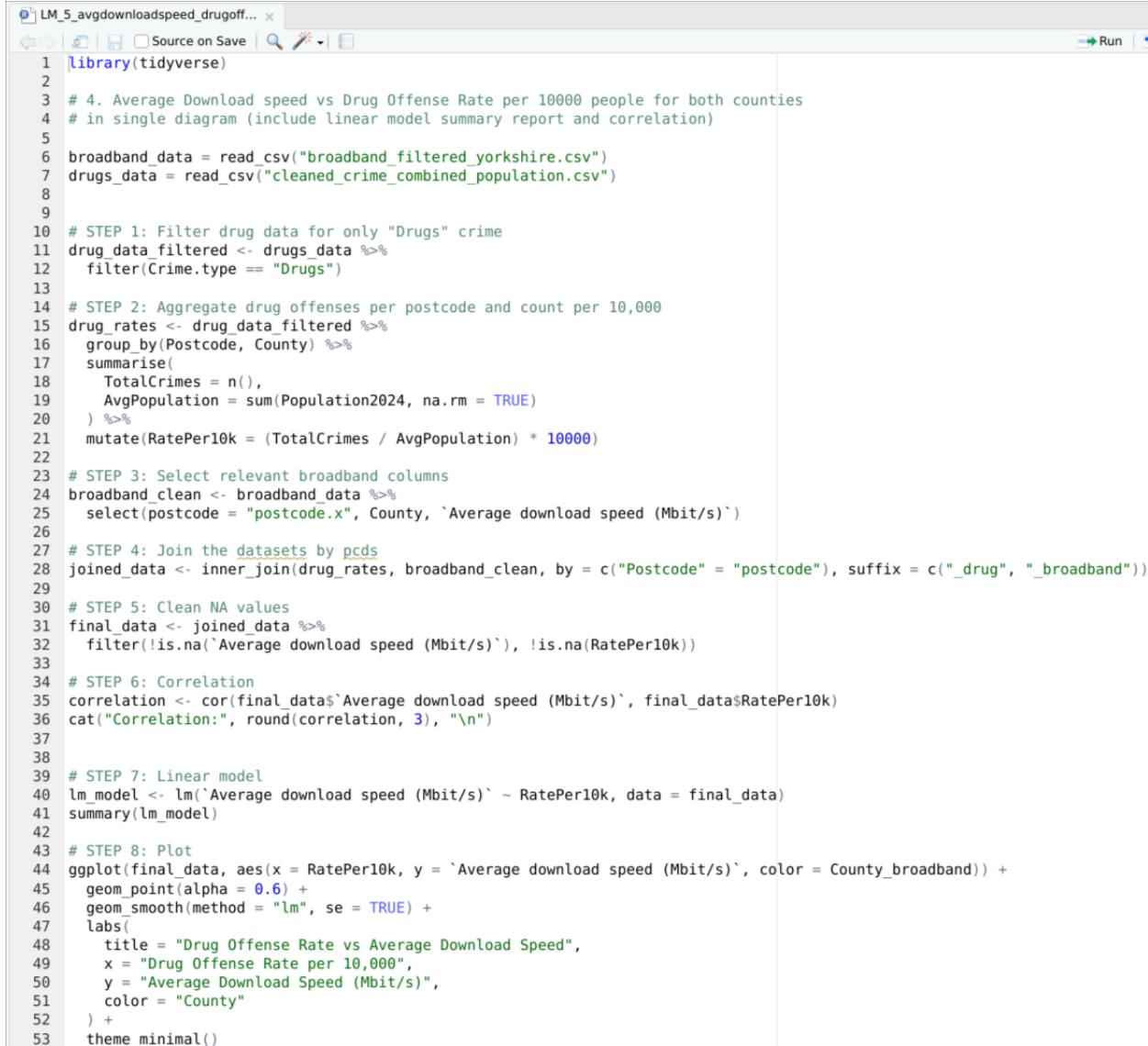
```
1 library(tidyverse)
2
3 school_dataset = read_csv("cleaned_filtered_school_dataset.csv")
4 house_dataset = read_csv("cleaned_house_prices.csv")
5 colnames(school_dataset)
6 colnames(house_dataset)
7
8
9 # Clean postcodes (remove spaces for consistent joining)
10 house_clean <- house_dataset %>%
11   select(postcode, price, county) %>%
12   mutate(postcode = str_remove_all(postcode, " ")) # e.g. "HX28SN"
13
14 school_clean <- school_dataset %>%
15   select(PCODE, ATT8SCR, County) %>%
16   rename(postcode = PCODE, attainment = ATT8SCR) %>%
17   mutate(postcode = str_remove_all(postcode, " ")) # e.g. "S639EW"
18
19 # Aggregate House Prices by Postcode
20 house_avg <- house_clean %>%
21   group_by(postcode, county) %>%
22   summarise(median_price = median(price, na.rm = TRUE), .groups = "drop")
23
24
25 # Merge School & House Datasets
26 merged_data <- inner_join(school_clean, house_avg, by = "postcode") %>%
27   filter(!is.na(attainment) & !is.na(median_price))
28
29 # Merge School & House Datasets
30 merged_data <- inner_join(school_clean, house_avg, by = "postcode") %>%
31   filter(!is.na(attainment) & !is.na(median_price))
32
33 # Plot (House Price vs Attainment Score)
34 ggplot(merged_data, aes(x = median_price, y = attainment, color = County)) +
35   geom_point(alpha = 0.6, size = 2.5) +
36   geom_smooth(method = "lm", se = TRUE, size = 1.2) +
37   labs(
38     title = "Attainment 8 Score vs House Price by County",
39     x = "Median House Price (£)",
40     y = "Attainment 8 Score",
41     color = "County"
42   ) +
43   theme_minimal(base_size = 14)
44
45 # Linear Model Summary (Per County)
46 # West Yorkshire model
47 west <- merged_data %>% filter(County == "West Yorkshire")
48 model_west <- lm(attainment ~ median_price, data = west)
49 summary(model_west)
50
51 # South Yorkshire model
52 south <- merged_data %>% filter(County == "South Yorkshire")
53 model_south <- lm(attainment ~ median_price, data = south)
54 summary(model_south)
55
56 # Correlation Coefficients
57 cor_west <- cor(west$median_price, west$attainment, use = "complete.obs")
58 cor_south <- cor(south$median_price, south$attainment, use = "complete.obs")
59
60 cat("Correlation (West Yorkshire):", cor_west, "\n")
61 cat("Correlation (South Yorkshire):", cor_south, "\n")
```

```

LM_4_att8_drugoffense.R x
Source on Save | 🔎 🖌️ | ⌂

1 library(tidyverse)
2
3 school_dataset = read_csv("cleaned_filtered_school_dataset.csv")
4 crime_data = read_csv("cleaned_crime_combined_population.csv")
5
6 # Filter for Drugs only and year 2023
7 drug_data <- crime_data %>%
8   filter(Crime.type == "Drugs", Year == 2023) %>%
9   distinct(LSOA.code, .keep_all = TRUE)
10
11 # Calculate Drug Offense Rate per 10,000 people
12 drug_rates <- drug_data %>%
13   group_by(Postcode, Year, Month) %>%
14   summarise(
15     drug_crimes = n(),
16     population = sum(Population2023, na.rm = TRUE),
17     .groups = "drop"
18   ) %>%
19   mutate(rate_per_10k = (drug_crimes / population) * 10000)
20
21 # Clean & Prepare School Dataset
22 school_clean <- school_dataset %>%
23   select(PCODE, ATT8SCR, County) %>%
24   rename(attainment = ATT8SCR) %>%
25   mutate(PCODE = str_remove_all(PCODE, " "))
26
27 colnames(school_clean)
28 colnames(drug_rates)
29
30
31 # Merge Drug Offense Rates with School Attainment
32 merged_data <- inner_join(school_clean, drug_rates, by = c("PCODE" = "Postcode")) %>%
33   filter(!is.na(attainment) & !is.na(rate_per_10k))
34
35 # Plot: Drug Crime Rate vs Attainment Score
36 ggplot(merged_data, aes(x = rate_per_10k, y = attainment, color = County)) +
37   geom_point(alpha = 0.6, size = 2.5) +
38   geom_smooth(method = "lm", se = TRUE, size = 1.2) +
39   labs(
40     title = "Attainment 8 Score vs Drug Offense Rate (per 10,000 people) in 2023",
41     x = "Drug Offense Rate per 10,000 People",
42     y = "Attainment 8 Score",
43     color = "County"
44   ) +
45   theme_minimal(base_size = 14)
46
47
48 # Linear Model Summary (Per County)
49 west <- merged_data %>% filter(County == "West Yorkshire")
50 south <- merged_data %>% filter(County == "South Yorkshire")
51
52 model_west <- lm(attainment ~ rate_per_10k, data = west)
53 model_south <- lm(attainment ~ rate_per_10k, data = south)
54
55 summary(model_west)
56 summary(model_south)
57
58 # Correlation
59 cor_west <- cor(west$rate_per_10k, west$attainment, use = "complete.obs")
60 cor_south <- cor(south$rate_per_10k, south$attainment, use = "complete.obs")
61
62 cat("Correlation (West Yorkshire):", cor_west, "\n")
63 cat("Correlation (South Yorkshire):", cor_south, "\n")

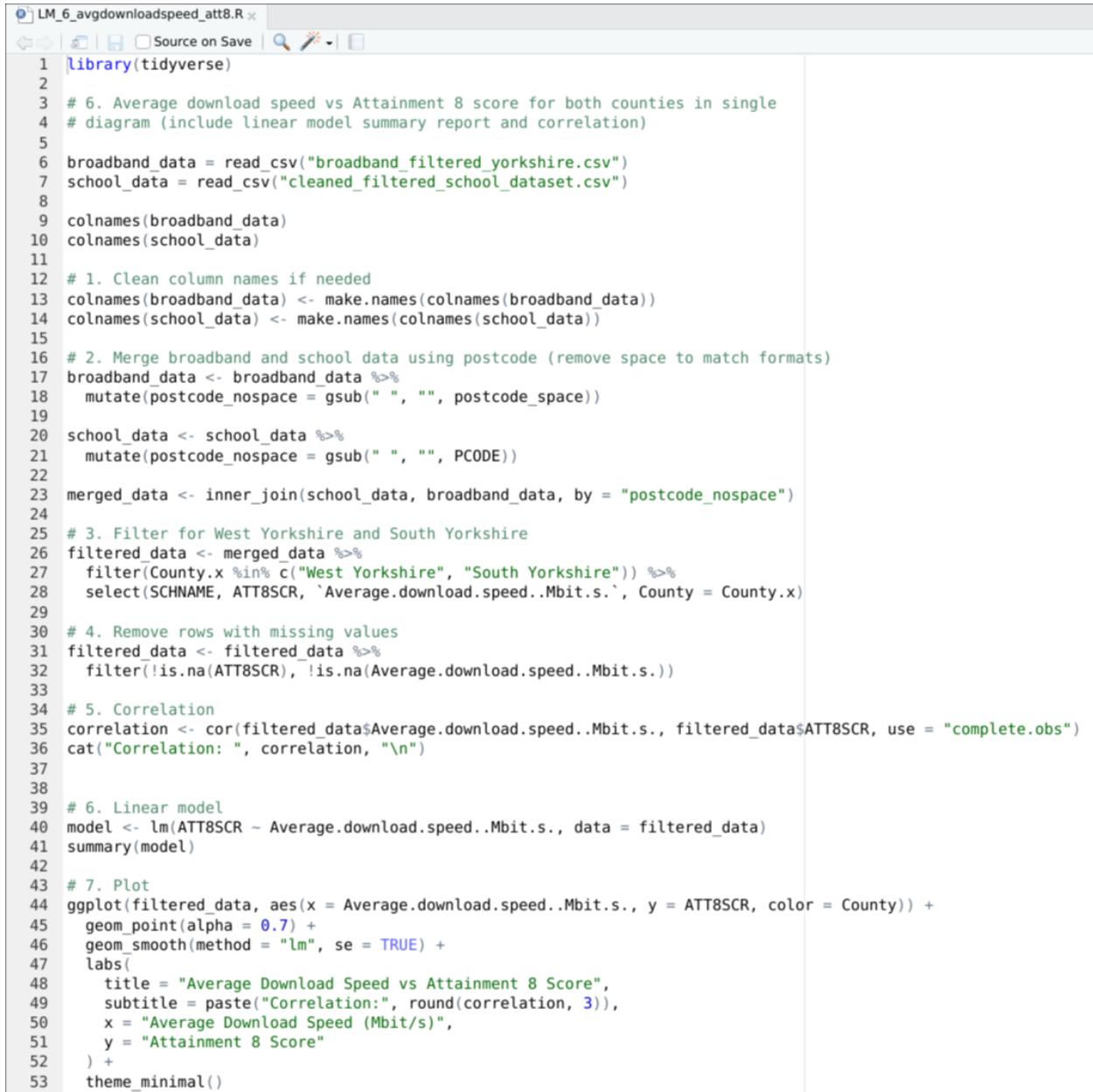
```



```

1 library(tidyverse)
2
3 # 4. Average Download speed vs Drug Offense Rate per 10000 people for both counties
4 # in single diagram (include linear model summary report and correlation)
5
6 broadband_data = read_csv("broadband_filtered_yorkshire.csv")
7 drugs_data = read_csv("cleaned_crime_combined_population.csv")
8
9
10 # STEP 1: Filter drug data for only "Drugs" crime
11 drug_data_filtered <- drugs_data %>%
12   filter(Crime.type == "Drugs")
13
14 # STEP 2: Aggregate drug offenses per postcode and count per 10,000
15 drug_rates <- drug_data_filtered %>%
16   group_by(Postcode, County) %>%
17   summarise(
18     TotalCrimes = n(),
19     AvgPopulation = sum(Population2024, na.rm = TRUE)
20   ) %>%
21   mutate(RatePer10k = (TotalCrimes / AvgPopulation) * 10000)
22
23 # STEP 3: Select relevant broadband columns
24 broadband_clean <- broadband_data %>%
25   select(postcode = "postcode.x", County, `Average download speed (Mbit/s)`)
26
27 # STEP 4: Join the datasets by pcds
28 joined_data <- inner_join(drug_rates, broadband_clean, by = c("Postcode" = "postcode"), suffix = c("_drug", "_broadband"))
29
30 # STEP 5: Clean NA values
31 final_data <- joined_data %>%
32   filter(!is.na(`Average download speed (Mbit/s)`), !is.na(RatePer10k))
33
34 # STEP 6: Correlation
35 correlation <- cor(final_data$`Average download speed (Mbit/s)`, final_data$RatePer10k)
36 cat("Correlation:", round(correlation, 3), "\n")
37
38
39 # STEP 7: Linear model
40 lm_model <- lm(`Average download speed (Mbit/s)` ~ RatePer10k, data = final_data)
41 summary(lm_model)
42
43 # STEP 8: Plot
44 ggplot(final_data, aes(x = RatePer10k, y = `Average download speed (Mbit/s)`, color = County_broadband)) +
45   geom_point(alpha = 0.6) +
46   geom_smooth(method = "lm", se = TRUE) +
47   labs(
48     title = "Drug Offense Rate vs Average Download Speed",
49     x = "Drug Offense Rate per 10,000",
50     y = "Average Download Speed (Mbit/s)",
51     color = "County"
52   ) +
53   theme_minimal()

```



```

1 library(tidyverse)
2
3 # 6. Average download speed vs Attainment 8 score for both counties in single
4 # diagram (include linear model summary report and correlation)
5
6 broadband_data = read_csv("broadband_filtered_yorkshire.csv")
7 school_data = read_csv("cleaned_filtered_school_dataset.csv")
8
9 colnames(broadband_data)
10 colnames(school_data)
11
12 # 1. Clean column names if needed
13 colnames(broadband_data) <- make.names(colnames(broadband_data))
14 colnames(school_data) <- make.names(colnames(school_data))
15
16 # 2. Merge broadband and school data using postcode (remove space to match formats)
17 broadband_data <- broadband_data %>%
18   mutate(postcode_nospace = gsub(" ", "", postcode_space))
19
20 school_data <- school_data %>%
21   mutate(postcode_nospace = gsub(" ", "", PCODE))
22
23 merged_data <- inner_join(school_data, broadband_data, by = "postcode_nospace")
24
25 # 3. Filter for West Yorkshire and South Yorkshire
26 filtered_data <- merged_data %>%
27   filter(County.x %in% c("West Yorkshire", "South Yorkshire")) %>%
28   select(SCHNAME, ATT8SCR, `Average.download.speed..Mbit.s.`, County = County.x)
29
30 # 4. Remove rows with missing values
31 filtered_data <- filtered_data %>%
32   filter(!is.na(ATT8SCR), !is.na(Average.download.speed..Mbit.s.))
33
34 # 5. Correlation
35 correlation <- cor(filtered_data$Average.download.speed..Mbit.s., filtered_data$ATT8SCR, use = "complete.obs")
36 cat("Correlation: ", correlation, "\n")
37
38 # 6. Linear model
39 model <- lm(ATT8SCR ~ Average.download.speed..Mbit.s., data = filtered_data)
40 summary(model)
41
42
43 # 7. Plot
44 ggplot(filtered_data, aes(x = Average.download.speed..Mbit.s., y = ATT8SCR, color = County)) +
45   geom_point(alpha = 0.7) +
46   geom_smooth(method = "lm", se = TRUE) +
47   labs(
48     title = "Average Download Speed vs Attainment 8 Score",
49     subtitle = paste("Correlation:", round(correlation, 3)),
50     x = "Average Download Speed (Mbit/s)",
51     y = "Attainment 8 Score"
52   ) +
53   theme_minimal()

```

RecommendationSystem.R

```

1 library(tidyverse)
2 library(lubridate)
3
4 house_data = read_csv("cleaned_house_prices.csv")
5 crime = read_csv("cleaned_crime_combined_population.csv")
6 broadband_data = read_csv("broadband_filtered_yorkshire.csv")
7 schoolData = read_csv("cleaned_filtered_school_dataset.csv")
8
9 # Check original data dimensions
10 cat("Original data dimensions:\n")
11 cat(sprintf("house_data: %d rows\n", nrow(house_data)))
12 cat(sprintf("crime: %d rows\n", nrow(crime)))
13 cat(sprintf("schoolData: %d rows\n", nrow(schoolData)))
14 cat(sprintf("broadband_data: %d rows\n", nrow(broadband_data)))
15
16 # Check unique towns in each dataset
17 cat("\nUnique towns count:\n")
18 cat(sprintf("house_data towns: %d\n", length(unique(house_data$town_city))))
19 cat(sprintf("crime towns: %d\n", length(unique(crime$town_city))))
20 cat(sprintf("schoolData towns: %d\n", length(unique(schoolData$TOWN))))
21 cat(sprintf("broadband_data towns: %d\n", length(unique(broadband_data$town_city))))
22
23
24 # Housing Data Processing
25 selected_house <- house_data %>%
26   mutate(TOWN = str_trim toupper(town_city)) %>%
27   group_by(TOWN) %>%
28   summarise(avgPrice = mean(price, na.rm = TRUE)) %>%
29   select(avgPrice, TOWN) %>%
30   na.omit() %>%
31   distinct()
32
33 # School Data Processing
34 selected_attainment8 <- schoolData %>%
35   mutate(
36     TOWN = str_trim toupper(TOWN),
37     ATT8SCR_clean = suppressWarnings(as.numeric(as.character(ATT8SCR)))
38   ) %>%
39   filter(!is.na(ATT8SCR_clean), !TOWN == "SOUTH YORKSHIRE", !TOWN == "WEST YORKSHIRE") %>%
40   group_by(TOWN) %>%
41   summarise(avgAtt8 = mean(ATT8SCR_clean, na.rm = TRUE)) %>%
42   ungroup()
43
44
45 # Broadband Data Processing
46 selected_broadband <- broadband_data %>%
47   group_by(town_city) %>%
48   mutate(
49     AvgUpSpeed = `Average upload speed (Mbit/s)`,
50     AvgDownSpeed = `Average download speed (Mbit/s)`
51   ) %>%
52   summarise(
53     avg_down_speed = mean(AvgDownSpeed, na.rm = TRUE)
54   ) %>%
55   mutate(TOWN = str_trim toupper(town_city)) %>%
56   select(TOWN, avg_down_speed)
57
58 # Crime Data processing
59 selected_crime <- crime %>%
60   group_by(postcode) %>%
61   summarise(
62     crimeno = n(),
63     TOWN = first(town_city) # or use unique(Town_City)[1] for safety
64   ) %>%
65   select(postcode, crimeno, TOWN) %>%
66   arrange(desc(crimeno))
67
68 cat(sprintf("crime towns: %d\n", length(unique(selected_crime$TOWN))))
69 colnames(selected_crime)
70

```

```

84 #Combine All Data (using left_join to keep as many as possible)
85 ranking <- selected_house %>%
86   left_join(selected_attainment8, by = "TOWN") %>%
87   left_join(selected_broadband, by = "TOWN") %>%
88   left_join(selected_crime, by = "TOWN") %>%
89   na.omit()
90
91
92 Extremes <- ranking %>%
93   summarise(
94     minDownSpeed = min(avg_down_speed, na.rm = TRUE),
95     maxDownSpeed = max(avg_down_speed, na.rm = TRUE),
96     minAtt8 = min(avgAtt8, na.rm = TRUE),
97     maxAtt8 = max(avgAtt8, na.rm = TRUE),
98     minHousingPrice = min(avgPrice, na.rm = TRUE),
99     maxHousingPrice = max(avgPrice, na.rm = TRUE),
100    minCrimeRate = min(crimeno, na.rm = TRUE),
101    maxCrimeRate = max(crimeno, na.rm = TRUE)
102  )
103
104 print(Extremes)
105
106
107
108 finalRanking <- ranking %>%
109   mutate(
110     normDownSpeed = 10 * (avg_down_speed - Extremes$minDownSpeed) / (Extremes$maxDownSpeed - Extremes$minDownSpeed),
111     normAtt8 = 10 * (avgAtt8 - Extremes$minAtt8) / (Extremes$maxAtt8 - Extremes$minAtt8),
112     normHousingPrice = 10 * (1 - (avgPrice - Extremes$minHousingPrice)) / (Extremes$maxHousingPrice - Extremes$minHousingPrice),
113     normCrimeRate = 10 * (1 - (crimeno - Extremes$minCrimeRate)) / (Extremes$maxCrimeRate - Extremes$minCrimeRate),
114     finalPoints = normDownSpeed + normAtt8 + normHousingPrice + normCrimeRate
115   )
116
117 colnames(finalRanking)
118
119
120 houserank <- finalRanking %>%
121   select(TOWN, avgPrice, normHousingPrice) %>%
122   arrange(desc(normHousingPrice)) %>%
123   distinct(TOWN,.keep_all = TRUE) %>%
124   slice_head(n=10)
125
126 crimerank <- finalRanking %>%
127   select(TOWN, normCrimeRate, crimeno) %>%
128   arrange(desc(normCrimeRate)) %>%
129   distinct(TOWN,.keep_all = TRUE) %>%
130   slice_head(n=10)
131
132 schoolrank <- finalRanking %>%
133   select(TOWN, avgAtt8, normAtt8) %>%
134   arrange(desc(normAtt8)) %>%
135   distinct(TOWN,.keep_all = TRUE) %>%
136   slice_head(n=10)
137
138 broadbandrank <- finalRanking %>%
139   select(TOWN, avg_down_speed, normDownSpeed) %>%
140   arrange(desc(normDownSpeed)) %>%
141   distinct(TOWN,.keep_all = TRUE) %>%
142   slice_head(n=10)
143
144
145 final_rank <- finalRanking %>%
146   select(TOWN, avgPrice, crimeno, avgAtt8, avg_down_speed, finalPoints) %>%
147   mutate(finalPoints = finalPoints / 4) %>%
148   arrange(desc(finalPoints)) %>%
149   distinct(TOWN,.keep_all = TRUE) %>%
150   slice_head(n=10)

```