Aayush Shrestha (230293)

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

ST5014CEM Data Science for Developers

Siddhartha Neupane

# Contents

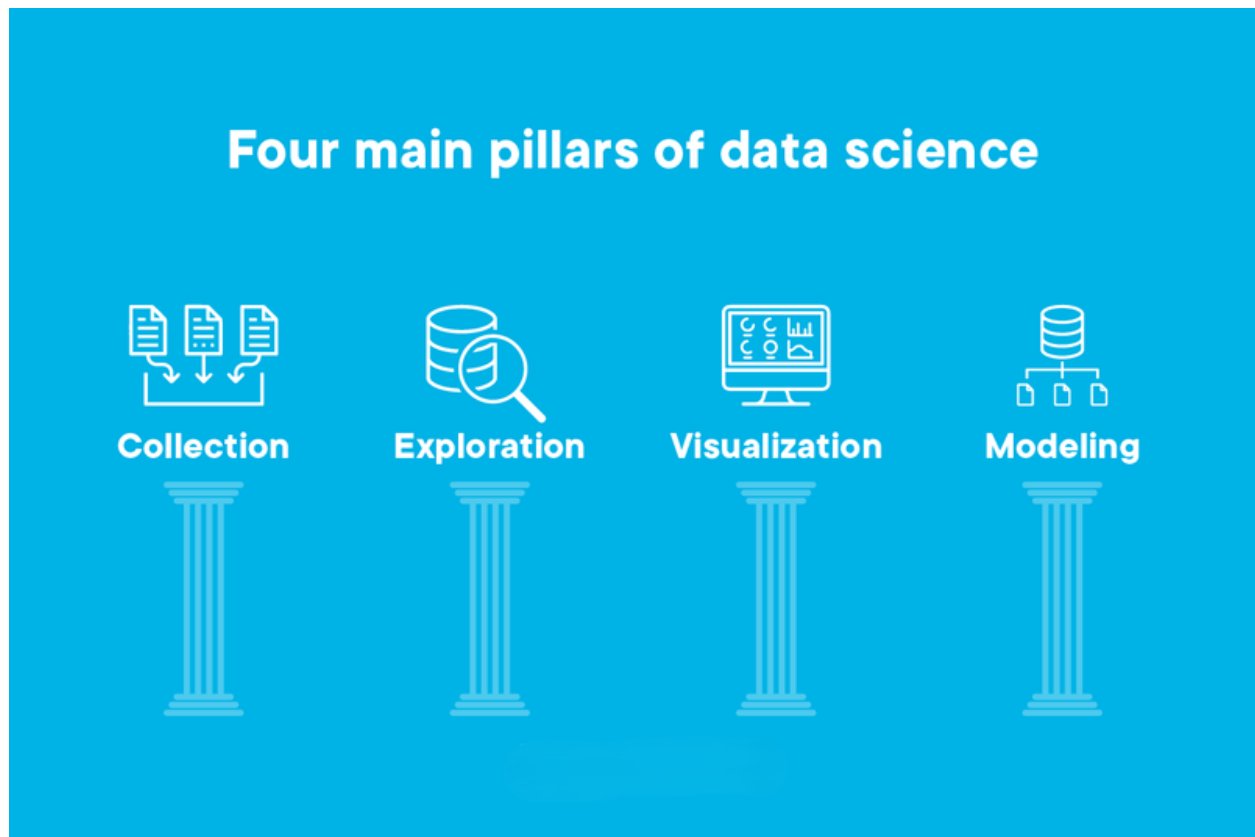# Table of Figure

# Introduction

Data science is the study of data to extract meaningful insights for business. It is a multidisciplinary approach that combines principles and practices from the fields of mathematics, statistics, artificial intelligence, and computer engineering to analyze large amounts of data. This analysis helps data scientists to ask and answer questions like what happened, why it happened, what will happen, and what can be done with the results *(What Is Data Science? - Data Science Explained - AWS*, n.d.-b)

This data science assignment aims to analyze real-world datasets that the instructor has made available via a download link. The four main datasets that are the focus of this assignment are town-level data, broadband speed and coverage data, street-level crime data, and school performance data.

Understanding the relationships between different social and economic factors in South and West Yorkshire towns is the primary objective of this project. After cleaning, trends and patterns in areas such as town characteristics, crime rates, broadband availability, and school attainment scores were examined using data visualization techniques. Additionally, linear modeling was used to look at how different factors, like crime or internet speed, affected outcomes like house prices or academic performance.

I created different visualizations for this assignment to study house prices, broadband speed, crime rates, and school performance in South and West Yorkshire. For house prices, I used a line graph to compare average prices from 2021 to 2024 for both counties, and boxplots to show price differences by district in each county. Broadband speed was shown with boxplots for download speeds across districts and bar charts comparing speeds in towns for each county. Crime data was visualized using boxplots for drug offense rates by district, a radar chart for vehicle crime in one county for a specific time, and a line chart to track drug offenses over the years in both counties. For schools, I made boxplots of average Attainment 8 scores for South Yorkshire in 2023 and West Yorkshire in 2022, plus a line graph showing how these scores changed over time across districts. These visuals help to understand patterns and compare the two countries easily.

*Figure 1 :Data Science*



**Four main pillars of data science**

Collection     Exploration     Visualization     Modeling

Kekare (2025b)

# Data cleaning

Data cleansing is the process of finding and removing errors, inconsistencies, duplications, and missing entries from data to increase data consistency and quality, also known as data scrubbing or cleaning *(What Is Data Cleansing? | TIBCO*, n.d.-b). In this assignment, the raw data we downloaded contains errors, missing values, duplicate values, and unformatted data. If we don't clean the data, it can lead to wrong results or confusing visualization.

For this coursework, we had to clean up the town, broadband, house pricing, crime, and school datasets so that they could be used together effectively. For instance, we handled missing values, eliminated extra spaces, and ensured that postcodes followed the same format. This allowed us to align data from multiple sources and prepare it for analysis and visualization.

*Figure 2 :Data cleaning*



Anwar (2025b)

**House Pricing**

I started by using the add row function to merge the 2021 and 2024 datasets into a single dataset for the house pricing data. I then restricted the data to records from West Yorkshire and South Yorkshire. I extracted the first four characters of the entire postcode and eliminated any spaces to create a new column called short Postcode that would make data joining easier in the future. For time-based analysis, I also eliminated the year from the Date column. I only chose the required columns—price, date, county, district, postcode, short postcode, and postcode—after sorting the data by county. To use the cleaned data for additional analysis and visualization, I finally saved it to a CSV file. The code for house pricing data cleaning.

*Figure 3 :House Pricing*

```
HousePrices = HousePrices2021%>%
  add_row(HousePrices2022) %>%
  add_row(HousePrices2023) %>%
  add_row(HousePrices2024)
cleanHousePrices = HousePrices %>%
  filter(County=="SOUTH YORKSHIRE"|County=="WEST YORKSHIRE") %>%
  mutate(shortPostcode = str_trim(substring(Postcode, 1,4))) %>%
  mutate(Year=substring(Date,7,10)) %>%
  arrange(County) %>%
  select(Postcode,shortPostcode,Price,Date,County,District)
write.csv(cleanHousePrices, "C:/DataScience-R/AayushShrestha-230293/Cleaned data/Cleaned_House_Prices.csv")
```

**Towns**

I started by gathering all the towns' 2021–2024 home price information. After that, I limited the records to only include West Yorkshire and South Yorkshire. To connect this to population data, I created a short postcode using the first four characters of the full postcode. To calculate the estimated populations for each year between 2012 and 2024, I also changed the population column to numeric, removed rows with invalid values, and applied growth multipliers. I then selected relevant columns like district, county, and annual population, and combined the population and home price data using the short Postcode. After removing duplicates and classifying the data by county, I saved the cleaned town-level dataset to a CSV file. And I saved cleaning like write.csv (Towns, "C:/DataScience-R/AayushShrestha-230293/Cleaned data/Towns.csv").

*Figure 4 :Towns*

```r
1   library(tidyverse)
2
3
4   PopulationData =read.csv("C:/DataScience-R/AayushShrestha-230293/Obtain/population/Population2011_1656567141570.csv")
5   colnames(PopulationData)
6
7   column_names = c(
8     "Housenumber","Price","Date",
9     "Postcode", "Zone 1","Zone 2",
10    "Zone 3","PAON","SAON", "Street",
11    "Locality","Town", "District",
12    "County","NA1","NA2"
13  )
14
15
16
17  HousePrices2021 = read_csv("C:/DataScience-R/AayushShrestha-230293/Obtain/house pricing/pp-2021.csv", col_names = FALSE) %>%
18    set_names(column_names)
19
20  HousePrices2022 = read_csv("C:/DataScience-R/AayushShrestha-230293/Obtain/house pricing/pp-2022.csv", col_names = FALSE) %>%
21    set_names(column_names)
22
23  HousePrices2023 = read_csv("C:/DataScience-R/AayushShrestha-230293/Obtain/house pricing/pp-2023.csv", col_names = FALSE) %>%
24    set_names(column_names)
25
26  HousePrices2024 = read_csv("C:/DataScience-R/AayushShrestha-230293/Obtain/house pricing/pp-2024.csv", col_names = FALSE) %>%
27    set_names(column_names)
28
29  PopulationData <- read_csv("C:/DataScience-R/AayushShrestha-230293/Obtain/population/Population2011_1656567141570.csv")
30
31  # Clean, convert and calculate population by short postcode
32  PopulationData_clean <- PopulationData %>%
33    # Remove rows where Population cannot convert to numeric
34    filter(!is.na(as.numeric(Population))) %>%
35    # Convert Population to numeric
36    mutate(Population = as.numeric(Population)) %>%
37    # Extract first 4 characters from Postcode, trim spaces
38    mutate(shortPostcode = str_trim(substring(Postcode, 1, 4))) %>%
39    # Group by shortPostcode and sum population
40    group_by(shortPostcode) %>%
41    summarise(Population2011 = sum(Population, na.rm = TRUE)) %>%
42    # Calculate population for years 2012 to 2024 using multipliers
43    mutate(
44      Population2012 = 1.00695353132322269 * Population2011,
45      Population2013 = 1.00669740535540783 * Population2012,
46      Population2014 = 1.00736463978721671 * Population2013,
47      Population2015 = 1.00792367505802859 * Population2014,
48      Population2016 = 1.00757874492811929 * Population2015,
49      Population2017 = 1.00679374473924223 * Population2016,
50      Population2018 = 1.00605929132212552 * Population2017,
51      Population2019 = 1.00561255390388033 * Population2018,
52      Population2020 = 1.00561255390388033 * Population2019,
53      Population2021 = 1.005425 * Population2020,
54      Population2022 = 1.004920 * Population2021,
55      Population2023 = 1.004510 * Population2022,
56      Population2024 = 1.004220 * Population2023
57    )
58  PopulationData_clean %>%
59    select(shortPostcode, Population2021, Population2022, Population2023, Population2024)
60
61
62
63
64  HousePrices = HousePrices2021 %>%
65    add_row(HousePrices2022) %>%
66    add_row(HousePrices2023) %>%
67    add_row(HousePrices2024)
68
69  Towns <- HousePrices %>%
70    filter(County %in% c("SOUTH YORKSHIRE", "WEST YORKSHIRE")) %>%
71    mutate(shortPostcode = str_trim(substr(Postcode, 1, 4))) %>%
72    left_join(PopulationData_clean, by = "shortPostcode") %>%
73    select(shortPostcode, District, County, Population2021, Population2022, Population2023, Population2024) %>%
74    group_by(shortPostcode) %>%
75    filter(row_number() == 1) %>%
76    arrange(County)
77
78  write.csv(Towns, "C:/DataScience-R/AayushShrestha-230293/Cleaned data/Towns.csv")
79
80
```

## Broadband Speed

I worked with two distinct datasets for the broadband speed data: coverage and performance. To create a short postcode for joining, I first cleaned the coverage data by extracting the first four characters of the postcode. I then chose crucial columns such as the percentage of premises with speeds below 2 Mbps and 10 Mbps, as well as the availability of superfast and ultrafast broadband. By creating the same short Postcode and choosing important columns like median and average download/upload speeds and average data usage, I similarly cleaned the performance data. I then used both the postcode and short Postcode to join the two datasets. Speed and availability data are included in the final cleaned broadband dataset, which I saved as a CSV file for analysis and visualization.

*Figure 5 :Broadband speed*

```
145
146
147  library(tidyverse)
148
149  coverage = read_csv("C:/DataScience-R/AayushShrestha-230293/Obtain/Broadband/201809_fixed_pc_coverage_r01.csv")
150  performance = read_csv("C:/DataScience-R/AayushShrestha-230293/Obtain/Broadband/201805_fixed_pc_performance_r03.csv")
151
152  colnames(coverage)
153  length(colnames(coverage))
154
155  colnames(performance)
156  length(colnames(performance))
157
158  clean_coverage = coverage %>%
159    mutate(shortPostcode = str_trim(substr(postcode, 1,4))) %>%
160    select(
161      postcode,
162      shortPostcode,
163      SFBB_Availability = `SFBB availability (% premises)`,
164      UFBB_Availability = `UFBB availability (% premises)`,
165      FTTP_Availability = `FTTP availability (% premises)`,
166      Below_2Mbps = `% of premises unable to receive 2Mbit/s`,
167      Below_10Mbps = `% of premises unable to receive 10Mbit/s`
168    )
169
170  clean_performance = performance %>%
171    mutate(shortPostcode = str_trim(substr(postcode,1,4))) %>%
172    select(
173      postcode,
174      shortPostcode,
175      Median_Download = `Median download speed (Mbit/s)`,
176      Avg_Download = `Average download speed (Mbit/s)`,
177      Median_Upload = `Median upload speed (Mbit/s)`,
178      Avg_Upload = `Average upload speed (Mbit/s)`,
179      Avg_Data_Usage = `Average data usage (GB)`
180    )
181
182  Broadband = left_join(clean_coverage,clean_performance,by = c("postcode","shortPostcode"))
183
184  write.csv(Broadband,"C:/DataScience-R/AayushShrestha-230293/Cleaned data/Cleaned_BroadBand_Speed.csv")
185
186
```

## Crime Dataset

I combined the street-level crime data for 2022–2025 from West Yorkshire and South Yorkshire. I combined all of the annual files into a single dataset using the add row function. I then selected key columns like Crime ID, LSOA name, Month, Reported by, falls within, Location, and Crime type to clean up the data. To make it obvious which area each crime belongs to, I also made a new County column from the "Falls within" field. I saved the cleaned dataset to a CSV file after classifying the records by location and county. After being cleaned, the data was prepared for additional analysis, including determining crime rates and displaying trends over time.

*Figure 6 :Crime*

```r
library(tidyverse)

south_yorkshire_2025 = read_csv("C:/DataScience-R/AayushShrestha-230293/Obtain/crime data set/crime2025/2025-04/2025-04-south-yorkshire-street.csv")
west_yorkshire_2025 = read_csv("C:/DataScience-R/AayushShrestha-230293/Obtain/crime data set/crime2025/2025-04/2025-04-west-yorkshire-street.csv")

south_yorkshire_2024 = read_csv("C:/DataScience-R/AayushShrestha-230293/Obtain/crime data set/crime2024/2024-04/2024-04-south-yorkshire-street.csv")
west_yorkshire_2024 = read_csv("C:/DataScience-R/AayushShrestha-230293/Obtain/crime data set/crime2024/2024-04/2024-04-west-yorkshire-street.csv")

south_yorkshire_2023 = read_csv("C:/DataScience-R/AayushShrestha-230293/Obtain/crime data set/crime2023/2023-04/2023-04-south-yorkshire-street.csv")
west_yorkshire_2023 = read_csv("C:/DataScience-R/AayushShrestha-230293/Obtain/crime data set/crime2023/2023-04/2023-04-west-yorkshire-street.csv")

south_yorkshire_2022 = read_csv("C:/DataScience-R/AayushShrestha-230293/Obtain/crime data set/crime2022/2022-06/2022-06-south-yorkshire-street.csv")
west_yorkshire_2022 = read_csv("C:/DataScience-R/AayushShrestha-230293/Obtain/crime data set/crime2022/2022-06/2022-06-west-yorkshire-street.csv")


CrimeData = south_yorkshire_2025 %>%
  add_row(west_yorkshire_2025) %>%
  add_row(south_yorkshire_2024) %>%
  add_row(west_yorkshire_2024) %>%
  add_row(south_yorkshire_2023) %>%
  add_row(west_yorkshire_2023) %>%
  add_row(south_yorkshire_2022) %>%
  add_row(west_yorkshire_2022)


Clean_crime= CrimeData %>%
  mutate(
    County = `Falls within`
  ) %>%
  select(
    CrimeID = `Crime ID`,
    LSOAname = `LSOA name`,
    Month,
    Reportedby = `Reported by`,
    Fallswithin = `Falls within`,
    County,
    Location,
    CrimeType = `Crime type`,
  ) %>%
  arrange(County, Location)

write.csv(Clean_crime,"C:/DataScience-R/AayushShrestha-230293/Cleaned data/Cleaned_Crime_Dataset.csv")
```

## School Dataset

I performed several cleaning and filtering operations on the school datasets from 2021 to 2024 to concentrate solely on schools in South Yorkshire and West Yorkshire. To include

only schools in the target counties, I first filtered them by their town or address fields. I then used common identifiers like `URN` and `ESTAB` to integrate school data with related datasets like provisional performance, results, and pupil destinations. Key variables such as school name, postcode, school type, gender, age range, and performance indicators such as Progress 8 (P8MEA) and EBACC scores were chosen from these combined datasets. I saved the cleaned data independently for every academic year after making sure that duplicate records were eliminated.

*Figure 7 : School 2021-2023*

```r
189  library(tidyverse)
190
191  library(tidyverse)
192  library(readxl)
193
194  final_2021_2022 = read_csv("C:/DataScience-R/AayushShrestha-230293/Obtain/Schooldata/2021-2022/england_ks4final.csv")
195  mats_performance_2021_2022 = read_csv("C:/DataScience-R/AayushShrestha-230293/Obtain/Schooldata/2021-2022/england_ks4-mats-performance.csv")
196  provisional_2021_2022 = read_csv("C:/DataScience-R/AayushShrestha-230293/Obtain/Schooldata/2021-2022/england_ks4provisional.csv")
197  pupdest_2021_2022 = read_csv("C:/DataScience-R/AayushShrestha-230293/Obtain/Schooldata/2021-2022/england_ks4-pupdest.csv")
198
199  underlying_1_2021_2022 = read_excel("C:/DataScience-R/AayushShrestha-230293/Obtain/Schooldata/2021-2022/england_ks4underlying_1.xlsx")
200  underlying_entriesandgrades_2_2021_2022 = read_excel("C:/DataScience-R/AayushShrestha-230293/Obtain/Schooldata/2021-2022/england_ks4underlying_entriesandgrades_2.xlsx")
201
202  school_information_2021_2022 = read_csv("C:/DataScience-R/AayushShrestha-230293/Obtain/Schooldata/2021-2022/england_school_information.csv")
203
204  # Check duplicates in existing dataframes (make sure column 'URN' exists in each)
205  anyDuplicated(final_2021_2022$URN)
206  anyDuplicated(pupdest_2021_2022$URN)
207  anyDuplicated(mats_performance_2021_2022$URN)
208
209
210
211
212
213
214
215
216
217
218  keywords = c("South Yorkshire", "West Yorkshire")
219
220  school_filtered_2021_2022 = school_information_2021_2022 %>%
221    filter(
222      toupper(str_trim(TOWN)) %in% toupper(keywords) |
223      toupper(str_trim(ADDRESS3)) %in% toupper(keywords)
224    )
225
226
227
228
229  provisional_joined_2021_2022 = school_filtered_2021_2022 %>%
230    inner_join(provisional_2021_2022, by = c("URN", "ESTAB"))
231
232
233  pupdest_joined_2021_2022 = provisional_joined_2021_2022 %>%
234    left_join(pupdest_2021_2022, by = c("URN", "ESTAB"))
235
236
237  final_joined_2021_2022 = pupdest_joined_2021_2022 %>%
238    left_join(final_2021_2022, by = c("URN", "ESTAB"))
239
240
241  cleaned_2021_2022 = final_joined_2021_2022 %>%
242    select(
243      URN,
```

```r
241  cleaned_2021_2022 = final_joined_2021_2022 %>%
242    select(
243      URN,
244      SCHNAME = SCHNAME.x,
245      TOWN,
246      ADDRESS3,
247      LANAME,
248      POSTCODE,
249      SCHOOLTYPE,
250      GENDER,
251      AGELOW,
252      AGEHIGH,
253      P8MEA = P8MEA.x,
254      P8MEA_FSM6CLA1A = P8MEA_FSM6CLA1A.x,
255      P8MEA_NFSM6CLA1A = P8MEA_NFSM6CLA1A.x,
256      EBACCAPS = EBACCAPS.x,
257      PTEBACC_95 = PTEBACC_95.x,
258      OVERALL_DESTPER,
259      EMPLOYMENTPER,
260      EDUCATIONPER
261    ) %>%
262    distinct()
263
264
265  write.csv(cleaned_2021_2022, "C:/DataScience-R/AayushShrestha-230293/Cleaned data/Cleaned_School_2021-2022.csv")
266
267  view(cleaned_2021_2022)
268
269
```

*Figure 8 :School 2022-2023*

```r
271
272  library(tidyverse)
273
274  ks2_final_2022_2023 = read_csv("C:/DataScience-R/AayushShrestha-230293/Obtain/Schooldata/2022-2023/england_ks2final.csv")
275  ks4_final_2022_2023 = read_csv("C:/DataScience-R/AayushShrestha-230293/Obtain/Schooldata/2022-2023/england_ks4final.csv")
276  ks2_mats_2022_2023 = read_csv("C:/DataScience-R/AayushShrestha-230293/Obtain/Schooldata/2022-2023/england_ks2-mats-performance.csv")
277  ks4_mats_2022_2023 = read_csv("C:/DataScience-R/AayushShrestha-230293/Obtain/Schooldata/2022-2023/england_ks4-mats-performance.csv")
278  provisional_2022_2023 = read_csv("C:/DataScience-R/AayushShrestha-230293/Obtain/Schooldata/2022-2023/england_ks4provisional.csv")
279  pupdest_2022_2023 = read_csv("C:/DataScience-R/AayushShrestha-230293/Obtain/Schooldata/2022-2023/england_ks4-pupdest.csv")
280  underlying_1_2022_2023 = read_csv("C:/DataScience-R/AayushShrestha-230293/Obtain/Schooldata/2022-2023/england_ks4underlying_1.xlsx")
281  underlying_2_2022_2023 = read_csv("C:/DataScience-R/AayushShrestha-230293/Obtain/Schooldata/2022-2023/england_ks4underlying_entriesandgrades_2.xlsx")
282  school_info_2022_2023 = read_csv("C:/DataScience-R/AayushShrestha-230293/Obtain/Schooldata/2022-2023/england_school_information.csv")
283
284
285  keywords = c("South Yorkshire", "West Yorkshire")
286
287  school_filtered_2022_2023 = school_info_2022_2023 %>%
288    filter(
289      toupper(str_trim(TOWN)) %in% toupper(keywords) |
290      toupper(str_trim(ADDRESS3)) %in% toupper(keywords)
291    )
292
293  provisional_joined_2022_2023 = school_filtered_2022_2023 %>%
294    inner_join(provisional_2022_2023, by = c("URN", "ESTAB"))
295
296  pupdest_joined_2022_2023 = provisional_joined_2022_2023 %>%
297    left_join(pupdest_2022_2023, by = c("URN", "ESTAB"))
298
299  final_joined_2022_2023 = pupdest_joined_2022_2023 %>%
300    left_join(ks4_final_2022_2023, by = c("URN", "ESTAB"))
301
302  cleaned_2022_2023 = final_joined_2022_2023 %>%
303    select(
304      URN,
305      SCHNAME = SCHNAME.x,
306      TOWN,
307      ADDRESS3,
308      LANAME,
309      POSTCODE,
310      SCHOOLTYPE,
311      GENDER,
312      AGELOW,
313      AGEHIGH,
314      P8MEA = P8MEA.x,
315      P8MEA_FSM6CLA1A = P8MEA_FSM6CLA1A.x,
316      P8MEA_NFSM6CLA1A = P8MEA_NFSM6CLA1A.x,
317      EBACCAPS = EBACCAPS.x,
318      PTEBACC_95 = PTEBACC_95.x,
319      OVERALL_DESTPER,
320      EMPLOYMENTPER,
321      EDUCATIONPER
322    ) %>%
323    distinct()
324
325  write.csv(cleaned_2022_2023, "C:/DataScience-R/AayushShrestha-230293/Cleaned data/Cleaned_School_2022-2023.csv")
326
```

*Figure 9 :School 2023-2024*

```r
331
332  library(tidyverse)
333
334  library(readxl)
335
336  ks2_final_2023_2024 = read_csv("C:/DataScience-R/AayushShrestha-230293/Obtain/Schooldata/2023-2024/england_ks2final.csv")
337
338  ks2_mats_2023_2024 = read_csv("C:/DataScience-R/AayushShrestha-230293/Obtain/Schooldata/2023-2024/england_ks2-mats-performance.csv")
339
340  ks4_final_2023_2024 = read_csv("C:/DataScience-R/AayushShrestha-230293/Obtain/Schooldata/2023-2024/england_ks4final.csv")
341
342  ks4_mats_2023_2024 = read_csv("C:/DataScience-R/AayushShrestha-230293/Obtain/Schooldata/2023-2024/england_ks4-mats-performance.csv")
343
344  provisional_2023_2024 = read_csv("C:/DataScience-R/AayushShrestha-230293/Obtain/Schooldata/2023-2024/england_ks4provisional.csv")
345
346  pupdest_2023_2024 = read_csv("C:/DataScience-R/AayushShrestha-230293/Obtain/Schooldata/2023-2024/england_ks4-pupdest.csv")
347
348  underlying_1_2023_2024 = read_excel("C:/DataScience-R/AayushShrestha-230293/Obtain/Schooldata/2023-2024/england_ks4underlying_1.xlsx")
349
350  underlying_2_2023_2024 = read_excel("C:/DataScience-R/AayushShrestha-230293/Obtain/Schooldata/2023-2024/england_ks4underlying_entriesandgrades_2.xlsx")
351
352  school_info_2023_2024 = read_csv("C:/DataScience-R/AayushShrestha-230293/Obtain/Schooldata/2023-2024/england_school_information.csv")
353
354
355
356  keywords = c("South Yorkshire", "West Yorkshire")
357
358  school_filtered_2023_2024 = school_info_2023_2024 %>%
359    filter(
360      toupper(str_trim(TOWN)) %in% toupper(keywords) |
361      toupper(str_trim(ADDRESS3)) %in% toupper(keywords)
362    )
363
364  provisional_joined_2023_2024 = school_filtered_2023_2024 %>%
365    inner_join(provisional_2023_2024, by = c("URN", "ESTAB"))
366
367  pupdest_joined_2023_2024 = provisional_joined_2023_2024 %>%
368    left_join(pupdest_2023_2024, by = c("URN", "ESTAB"))
369
370  final_joined_2023_2024 = pupdest_joined_2023_2024 %>%
371    left_join(ks4_final_2023_2024, by = c("URN", "ESTAB"))
372
373  cleaned_2023_2024 = final_joined_2023_2024 %>%
374    select(
375      URN,
376      SCHNAME = SCHNAME.x,
377      TOWN,
378      ADDRESS3,
379      LANAME,
380      POSTCODE,
381      SCHOOLTYPE,
382      GENDER,
383      AGELOW,
```

```
383        AGELOW,
384        AGEHIGH,
385        P8MEA  = P8MEA.x,
386        P8MEA_FSM6CLA1A  = P8MEA_FSM6CLA1A.x,
387        P8MEA_NFSM6CLA1A = P8MEA_NFSM6CLA1A.x,
388        EBACCAPS = EBACCAPS.x,
389        PTEBACC_95 = PTEBACC_95.x,
390        OVERALL_DESTPER,
391        EMPLOYMENTPER,
392        EDUCATIONPER
393      ) %>%
394      distinct()
395
396    write.csv(cleaned_2023_2024, "C:/DataScience-R/AayushShrestha-230293/Cleaned data/Cleaned_School_2023-2024.csv")
397
398    view(cleaned_2023_2024)
399
400
401
402
```
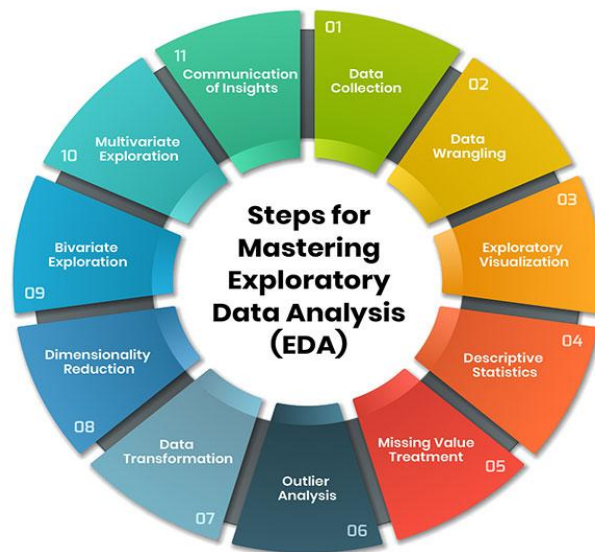
# Exploratory Data Analysis (EDA)

During the exploratory data analysis phase, I searched the cleaned datasets for home prices, broadband speeds, crime rates, and academic achievement in South and West Yorkshire for patterns, trends, and connections. I used visual aids like boxplots, line charts, bar graphs, and radar charts to understand how each variable behaves across different years, towns, and districts. For example, I compared average house prices over time, looked at crime rates by type and location, compared broadband speeds across towns, and monitored changes in school performance. This step helped me identify important insights, like areas with high crime rates or low school scores and guided the development of the linear model and recommendation system.

Exploratory data analysis (EDA) is an open-ended, iterative data analysis approach designed to unearth patterns, anomalies, relationships, or insights without preconceived notions. John Tukey, a renowned American mathematician, introduced EDA in the 1970s to analyze data using a combination of statistical tools and data discovery methods.
**https://www.coursera.org/articles/exploratory-data-analysis**

*Figure 10:Exploratory Data Analysis*



(*A Comprehensive Guide to Mastering Exploratory Data Analysis*, n.d.-b)

## Broadband speed

I started by connecting broadband speeds to counties and districts by combining the cleaned broadband dataset with town-level data using standardized short postcodes. Next, I contrasted South Yorkshire's and West Yorkshire's broadband performance. I highlighted regions with consistently high or low performance by using boxplots to illustrate the variation in average download speeds across districts. To highlight regional variations, median download speeds by town and district were also compared using bar charts. By highlighting towns with the best and worst connectivity, these visualizations helped identify areas that might benefit from upgrades to their digital infrastructure.

*Figure 11 :Analysis broadband*

```r
4   library(tidyverse)
5
6
7   Broadband_speed <- read_csv("C:/DataScience-R/AayushShrestha-230293/Obtain/Broadband/201809_fixed_pc_coverage_r01.csv")
8
9
10  Towns <- read_csv("C:/DataScience-R/AayushShrestha-230293/Cleaned data/Towns.csv")
11  # Clean and standardize postcode formats
12  Broadband_speed <- Broadband_speed %>%
13    mutate(shortPostcode = str_trim(toupper(shortPostcode)))
14
15  Towns <- Towns %>%
16    mutate(shortPostcode = str_trim(toupper(shortPostcode)))
17
18  # Merge datasets on shortPostcode
19  BroadbandMerged <- Broadband_speed %>%
20    left_join(Towns, by = "shortPostcode")
21
22  BroadbandMerged %>%
23    filter(str_detect(tolower(County), "west yorkshire"),
24           !is.na(Avg_Download),
25           !is.na(District)) %>%
26    ggplot(aes(x = reorder(District, Avg_Download, FUN = median), y = Avg_Download)) +
27    geom_boxplot(fill = "brown") +
28    labs(title = "West Yorkshire: Download Speed by District",
29         x = "District", y = "Avg Download Speed (Mbps)") +
30    coord_flip() +
31    theme_minimal()
32
33  BroadbandMerged %>%
34    filter(str_detect(tolower(County), "south yorkshire"),
35           !is.na(Avg_Download),
36           !is.na(District)) %>%
37    ggplot(aes(x = reorder(District, Avg_Download, FUN = median), y = Avg_Download)) +
38    geom_boxplot(fill = "purple") +
39    labs(title = "South Yorkshire: Download Speed by District",
40         x = "District", y = "Avg Download Speed (Mbps)") +
41    coord_flip() +
42    theme_minimal()
43
44
45  library(tidyverse)
46
47  # Load data
48  Broadband_speed <- read_csv("C:/DataScience-R/AayushShrestha-230293/Cleaned data/Cleaned_BroadBand_Speed.csv")
49  Towns <- read_csv("C:/DataScience-R/AayushShrestha-230293/Cleaned data/Towns.csv")
50
51  # Clean postcode for merging
52  Broadband_speed <- Broadband_speed %>%
53    mutate(shortPostcode = str_trim(toupper(shortPostcode)))
54
55  Towns <- Towns %>%
56    mutate(shortPostcode = str_trim(toupper(shortPostcode)))
```

```
57
58  # Merge datasets
59  BroadbandMerged <- Broadband_speed %>%
60    left_join(Towns, by = "shortPostcode")
61
62
63
64
65
66  colnames(BroadbandMerged)
67  BroadbandMerged %>%
68    filter(
69      str_detect(tolower(County), "west yorkshire"),
70      !is.na(Avg_Download),
71      !is.na(District)
72    ) %>%
73    ggplot(aes(x = reorder(District, Avg_Download), y = Avg_Download)) +
74    geom_col(fill = "darkblue") +
75    labs(
76      title = "West Yorkshire: Avg Download Speed by District",
77      x = "District",
78      y = "Avg Download Speed (Mbps)"
79    ) +
80    scale_y_continuous(labels = scales::label_number()) +
81    coord_flip() +
82    theme_minimal()
83
84  BroadbandMerged %>%
85    filter(
86      str_detect(tolower(County), "south yorkshire"),
87      !is.na(Avg_Download),
88      !is.na(District)
89    ) %>%
90    ggplot(aes(x = reorder(District, Avg_Download), y = Avg_Download)) +
91    geom_col(fill = "brown") +
92    labs(
93      title = "South Yorkshire: Avg Download Speed by District",
94      x = "District",
95      y = "Avg Download Speed (Mbps)"
96    ) +
97    coord_flip() +    # Flip for better readability
98    theme_minimal()
99
100
101
```

Visualization

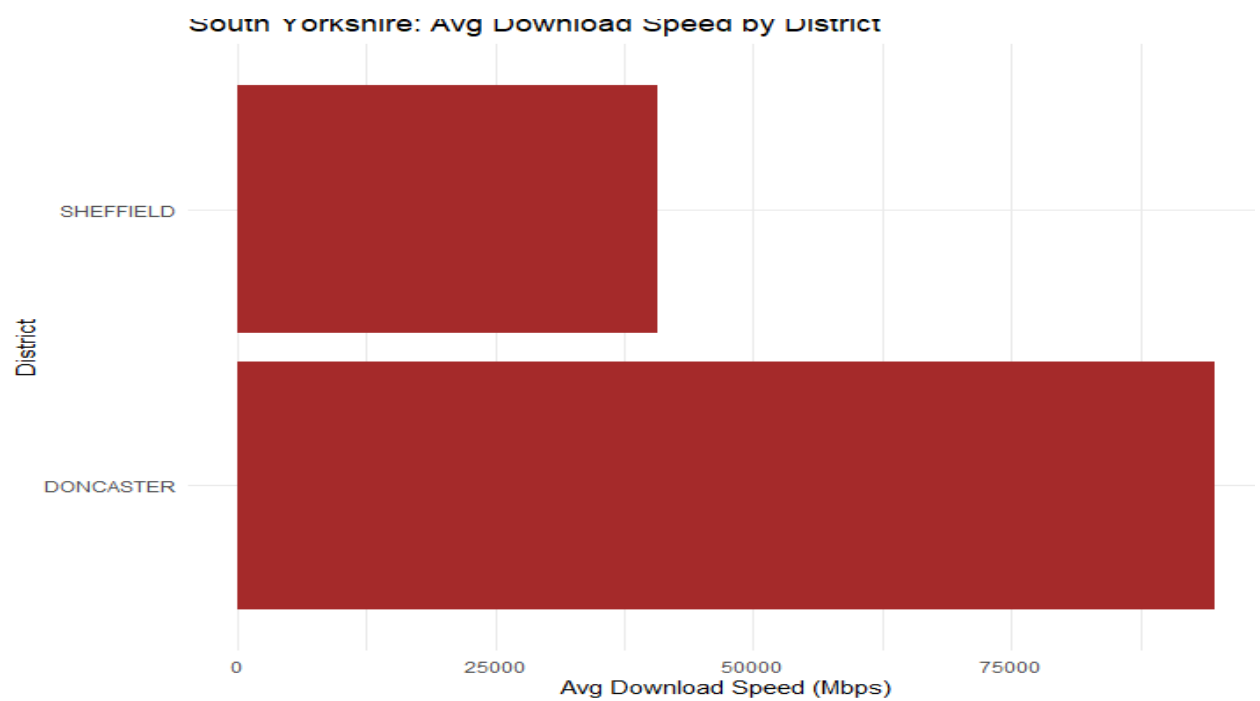*Figure 12 : Bar Graph South Yorkshire average download speed*



South Yorkshire: Avg Download Speed by District

*Figure 13: Bar Graph West Yorkshire Avg download speed*
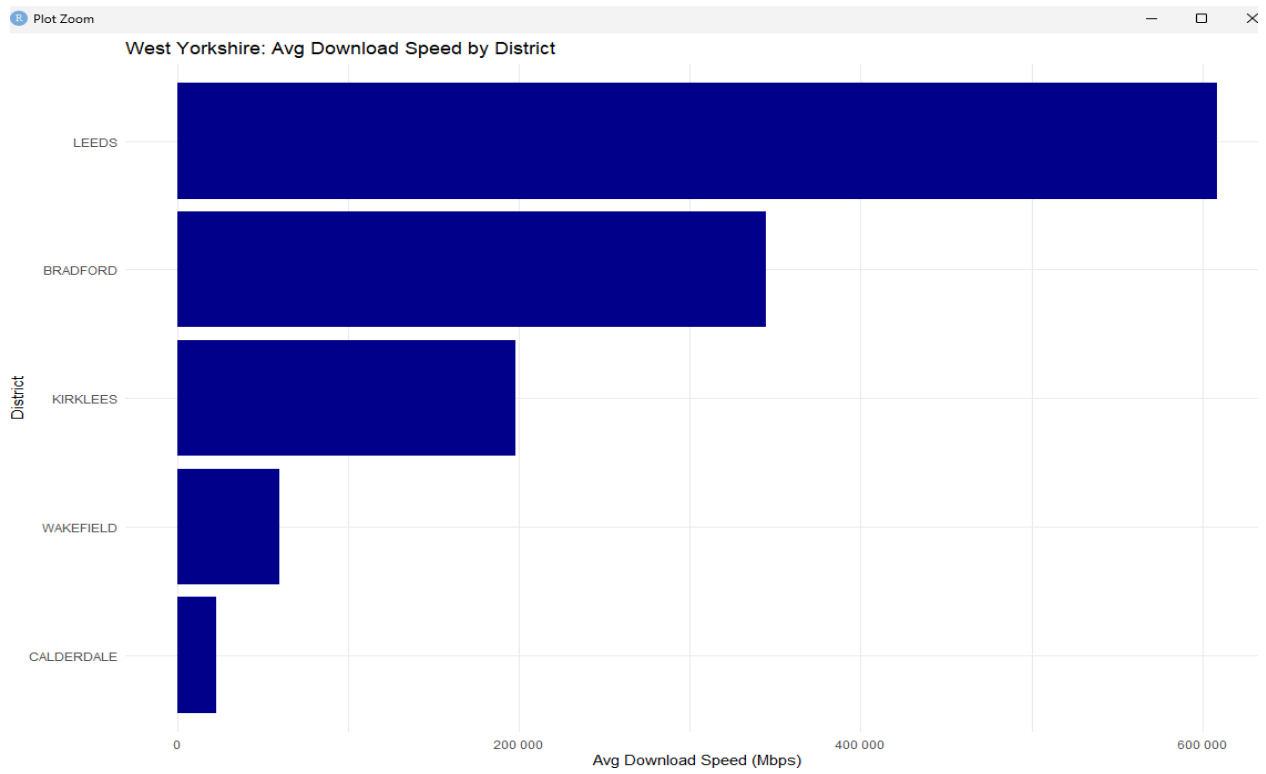


*Figure 14: Box plot West Yorkshire*

*Figure 15 :Box plot South Yorkshire*



**Crime Data Set**

The analysis of crime rates in South and West Yorkshire from 2022 to 2025 concentrated on robberies, car crimes, and drug offenses.  The distribution of drug offense rates by district was shown using boxplots, which highlighted regions with greater variances and higher crime rates.  West Yorkshire's April 2025 vehicle crime rates by district were displayed on a radar-style polar bar chart, which made it simple to compare the severity of crimes in different places.  Similarly, a pie chart showed the percentage of robbery offenses in South Yorkshire by district, showing the regions with the highest frequency of occurrences.  To examine changes over time and standardize data by population for fair

comparison, a line chart was also used to show trends in drug offense rates per 10,000 persons over three years.

```r
110
111  library(tidyverse)
112
113  library(ggplot2)
114
115  crime = read_csv("C:/DataScience-R/AayushShrestha-230293/Cleaned data/Cleaned_Crime_Dataset.csv")
116  colnames(crime)
117  view(crime)
118  town = read_csv("C:/DataScience-R/AayushShrestha-230293/Cleaned data/Towns.csv")
119  colnames(town)
120  view(town)
121
122  crime = crime %>%
123    mutate(Year = as.integer(substr(Month, 1, 4)),
124           District = str_extract(LSOAname, "^[^ ]+"))
125
126  crime = crime %>%
127    mutate(County = str_replace(County, " Police$", ""))
128
129
130
131  #Box plot- Drug Offense Rate per District (Two Diagrams)
132
133
134  # South Yorkshire Drug Offenses by District-Year
135  south_yorkshire = crime %>%
136    filter(CrimeType == "Drugs", County == "South Yorkshire", !is.na(District)) %>%
137    group_by(District, Year) %>%
138    summarise(Offenses = n(), .groups = "drop")
139
140  # Boxplot for South Yorkshire
141  ggplot(south_yorkshire, aes(x = reorder(District, Offenses, FUN = median), y = Offenses)) +
142    geom_boxplot(fill = "brown", outlier.alpha = 0.9) +
143    labs(title = "South Yorkshire: Drug Offense Distribution by District",
144         x = "District", y = "Offenses per Year") +
145    theme_minimal() +
146    theme(axis.text.x = element_text(angle = 45, hjust = 1))
147
148
149
150
151  # West Yorkshire Drug Offenses by District-Year
152  west_yorkshire = crime %>%
153    filter(CrimeType == "Drugs", County == "West Yorkshire", !is.na(District)) %>%
154    group_by(District, Year) %>%
155    summarise(Offenses = n(), .groups = "drop")
156
157  # Boxplot for West Yorkshire
158  ggplot(west_yorkshire, aes(x = reorder(District, Offenses, FUN = median), y = Offenses)) +
159    geom_boxplot(fill = "purple", outlier.alpha = 0.9) +
160    labs(title = "West Yorkshire: Drug Offense Distribution by District",
161         x = "District", y = "Offenses per Year") +
162    theme_minimal() +
163    theme(axis.text.x = element_text(angle = 45, hjust = 1))
```
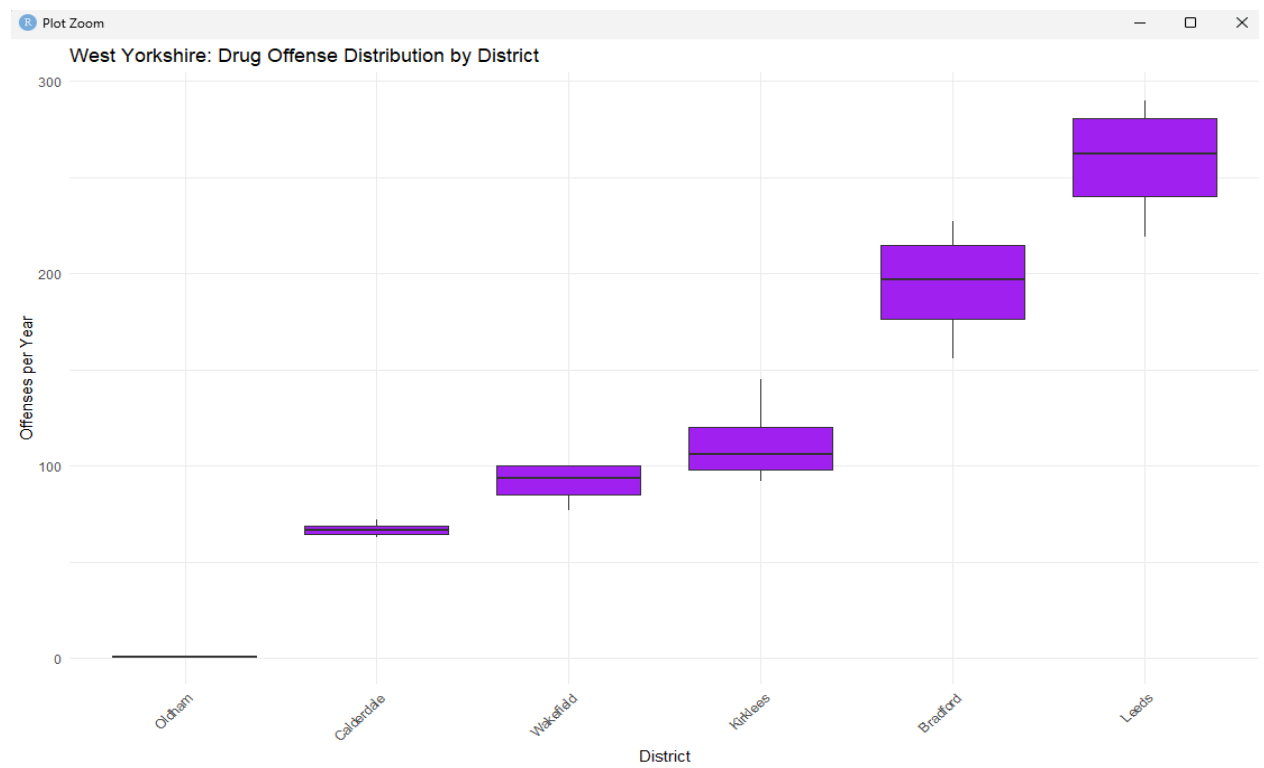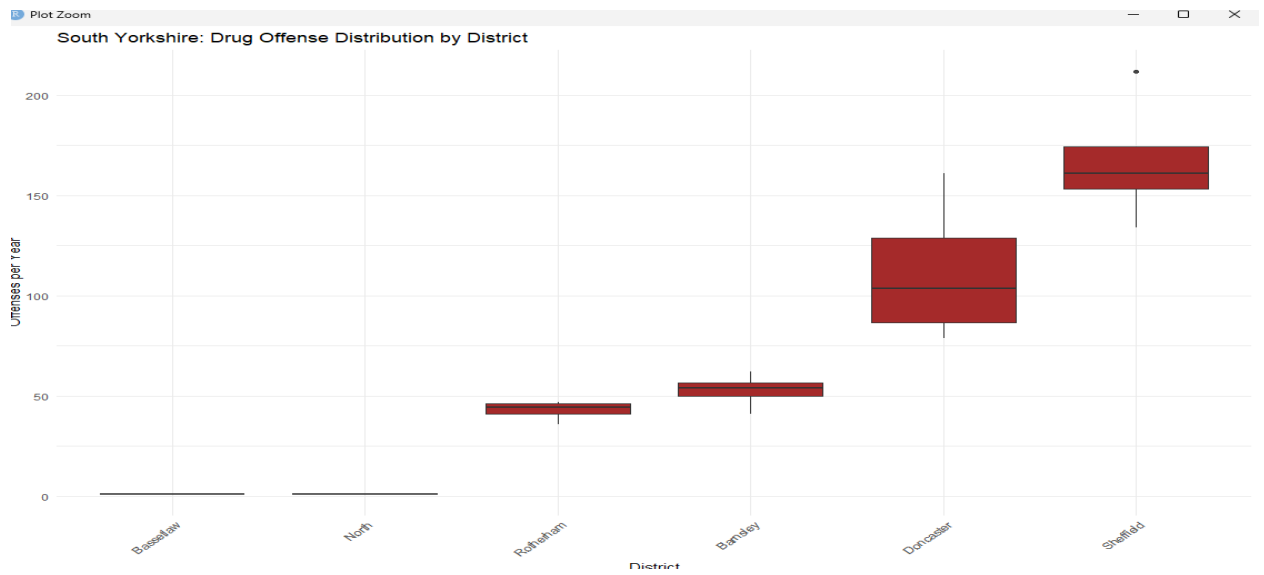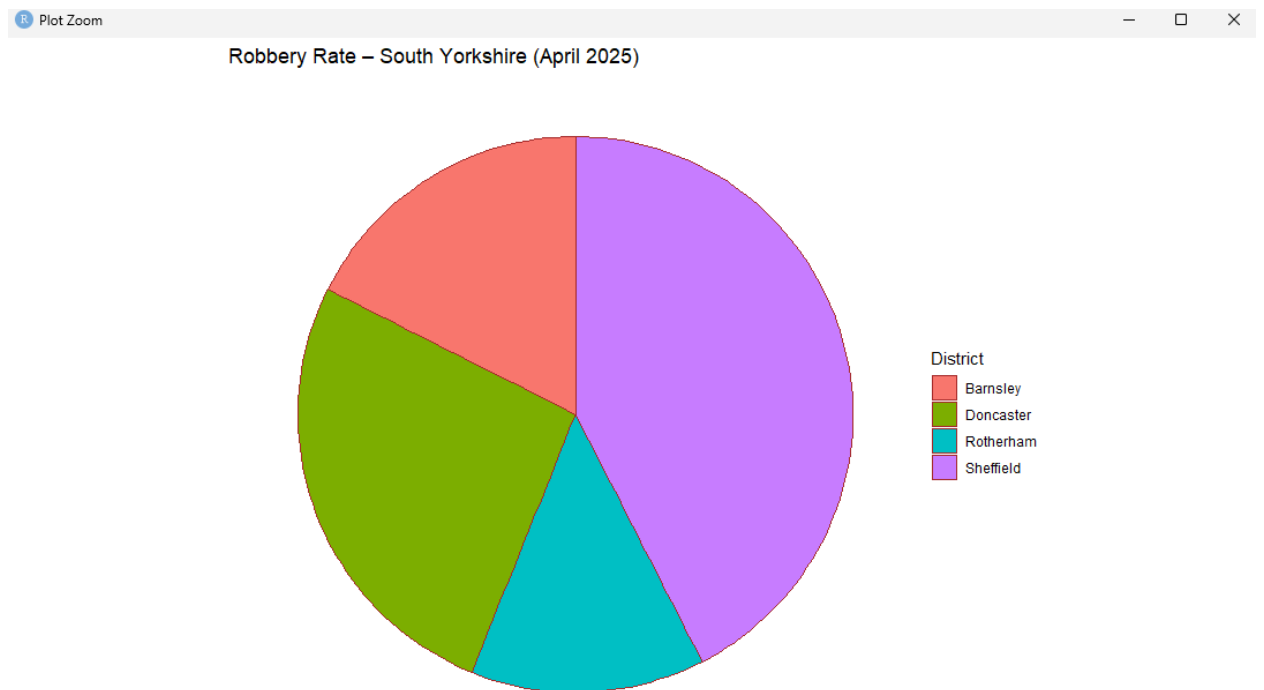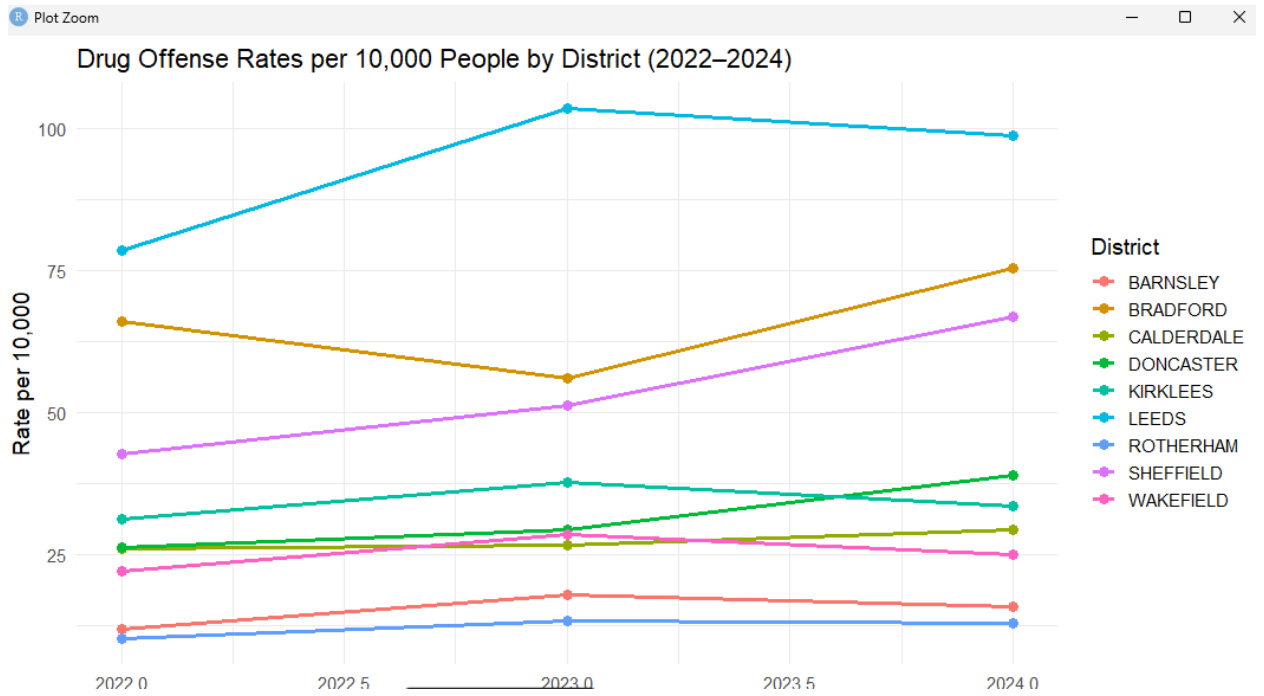
```r
173
174  vehicle_data = crime %>%
175    filter(CrimeType == "Vehicle crime",
176           County == "West Yorkshire",
177           Month == "2025-04",
178           !is.na(District)) %>%
179    group_by(District) %>%
180    summarise(Crimes = n()) %>%
181    arrange(desc(Crimes))
182
183  # Create radar-style polar bar chart
184  ggplot(vehicle_data, aes(x = reorder(District, Crimes), y = Crimes, fill = District)) +
185    geom_col(show.legend = FALSE, color = "blue") +
186    coord_polar(start = 0) +
187    labs(title = "Radar-Style Chart: Vehicle Crime in west Yorkshire (April 2025)",
188         x = "" , y = "") +
189    theme_minimal() +
190    theme(axis.text.x = element_text(size = 9, angle = 90))
191
192
193
194
195
196  #Pie chart for Robbery rate for any one of two counties (for any specific month and year)
197
198  pie_data = crime %>%
199    filter(CrimeType == "Robbery",
200           County == "South Yorkshire",
201           Month == "2025-04",
202           !is.na(District)) %>%
203    group_by(District) %>%
204    summarise(Crimes = n())
205
206  ggplot(pie_data, aes(x = "", y = Crimes, fill = District)) +
207    geom_col(width = 3, color = "brown") +
208    coord_polar("y") +
209    labs(title = "Robbery Rate - South Yorkshire (April 2025)", y = "", x = "") +
210    theme_void() +
211    theme(legend.position = "right")
212
213
214
215
216
217
218
219  #Line chart for Drug offense rates per 10,000 people for both counties in same diagram for all years
220
221  town_long = town %>%
222    pivot_longer(
223      cols = starts_with("Population"),
224      names_to = "Year",
225      names_prefix = "Population",
226      values_to = "Population"
227    ) %>%
```

```r
227    ) %>%
228    mutate(Year = as.integer(Year)) %>%
229    select(District, County, Year, Population)
230
231  drug_crime_by_district = crime %>%
232    filter(tolower(CrimeType) == "drugs") %>%
233    group_by(District, Year) %>%
234    summarise(Total_Offenses = n(), .groups = "drop") %>%
235    mutate(Year = as.integer(Year))
236
237  drug_crime_by_district = drug_crime_by_district %>%
238    mutate(District = str_to_upper(str_trim(District)))
239
240  town_long = town_long %>%
241    mutate(District = str_to_upper(str_trim(District)))
242
243
244  merged_data = drug_crime_by_district %>%
245    left_join(town_long, by = c("District", "Year")) %>%
246    filter(!is.na(Population)) %>%
247    mutate(Rate_per_10000 = (Total_Offenses / Population) * 10000)
248
249  final_data = merged_data %>%
250    group_by(District, Year) %>%
251    summarise(
252      Total_Offenses = sum(Total_Offenses, na.rm = TRUE),
253      Population = sum(Population, na.rm = TRUE),
254      Rate_per_10000 = (Total_Offenses / Population) * 10000,
255      .groups = "drop"
256    )
257
258
259  ggplot(final_data, aes(x = Year, y = Rate_per_10000, color = District)) +
260    geom_line(linewidth = 1.2) +
261    geom_point(size = 3) +
262    labs(
263      title = "Drug Offense Rates per 10,000 People by District (2022-2024)",
264      x = "Year",
265      y = "Rate per 10,000",
266      color = "District"
267    ) +
268    theme_minimal(base_size = 14)
269
270
271
```

Visualization

*Figure 17 : Visualization  crime*

South Yorkshire: Drug Offense Distribution by District

West Yorkshire: Drug Offense Distribution by District

## Drug Offense Rates per 10,000 People by District (2022–2024)

Robbery Rate – South Yorkshire (April 2025)



23

Radar-Style Chart: Vehicle Crime in West Yorkshire (April 2025)

## House pricing

The analysis of home prices in South and West Yorkshire from 2021 to 2024 revealed trends and differences at the county and district levels. A line graph showing the annual average home prices for each district allowed comparison of the price changes over time in both counties. The districts that grew steadily and those that remained largely unchanged were displayed in this graphic. Additionally, boxplots were used to display the home price distribution for each county independently. By providing information on price spread, outliers, and median values, these boxplots made it easier to spot variability within counties. Overall, understanding regional and temporal differences in home price trends was made easier by the visualizations.

*Figure 18: Analysis of House pricing*

```
339  average by County in 2023 only
340  county_2023_avg = county_year_avg %>%
341    filter(Year == 2023)
342
343  #Bar chart: Average price in 2023, both counties
344
345  ggplot(county_2023_avg, aes(x = County, y = AveragePrice, fill = County)) +
346    geom_col(width = 0.3) +
347    scale_fill_viridis_d(option = "C") +
348    scale_y_continuous(labels = scales::comma) +
349    labs(title = "Average House Prices by Country (2023)",
350         x = NULL, y = "Average Price") +
351    theme_minimal() +
352    theme(legend.position = "none")
353
354
355
356
357  #by district
358  district_2023_avg = HousePrices %>%
359    mutate(Year = year(ymd(Date))) %>%
360    filter(Year == 2023) %>%
361    group_by(District) %>%
362    summarise(AveragePrice = mean(Price, na.rm = TRUE), .groups = "drop")
363
364  Bar chart: Average house price by District in 2023
365  library(RColorBrewer)
366
367  ggplot(district_2023_avg, aes(x = reorder(District, -AveragePrice), y = AveragePrice, fill = District)) +
368    geom_col(width = 0.5) +
369    scale_fill_brewer(palette = "Set3") +
370    scale_y_continuous(labels = scales::comma) +
371    labs(title = "Average House Prices by District (2023)",
372         x = "District", y = "Average Price") +
373    theme_minimal() +
374    theme(legend.position = "none",
375          axis.text.x = element_text(angle = 45, hjust = 1))
376
377
378
379
380
381  #----------------------------------------------------------------------------------------------
382
383  #by county
384  #Box-plots: 2021-2024 distribution for each county in separate panels
385
386  ggplot(county_year_avg, aes(x = "", y = AveragePrice)) +
387    geom_boxplot(fill = "brown", outlier.alpha = 0.2) +
388    facet_wrap(~ County, ncol = 1) +
389    scale_y_continuous(labels = scales::comma) +
390    labs(title = "Distribution of House Prices (2021-2024)",
391         y = "Price", x = "") +
391         y = "Price", x = "") +
392    theme_minimal()
393
394
395
396  #District
397
398  ggplot(district_year_avg, aes(x = reorder(District, AveragePrice, FUN = median), y = AveragePrice)) +
399    geom_boxplot(fill = "skyblue", outlier.alpha = 0.9) +
400    scale_y_continuous(labels = scales::comma) +
401    labs(title = "Yearly Average House Prices by District (2021-2024)",
402         x = "District", y = "Average Price") +
403    theme_minimal() +
404    theme(axis.text.x = element_text(angle = 45, hjust = 1))
405
406
407
408
409
```
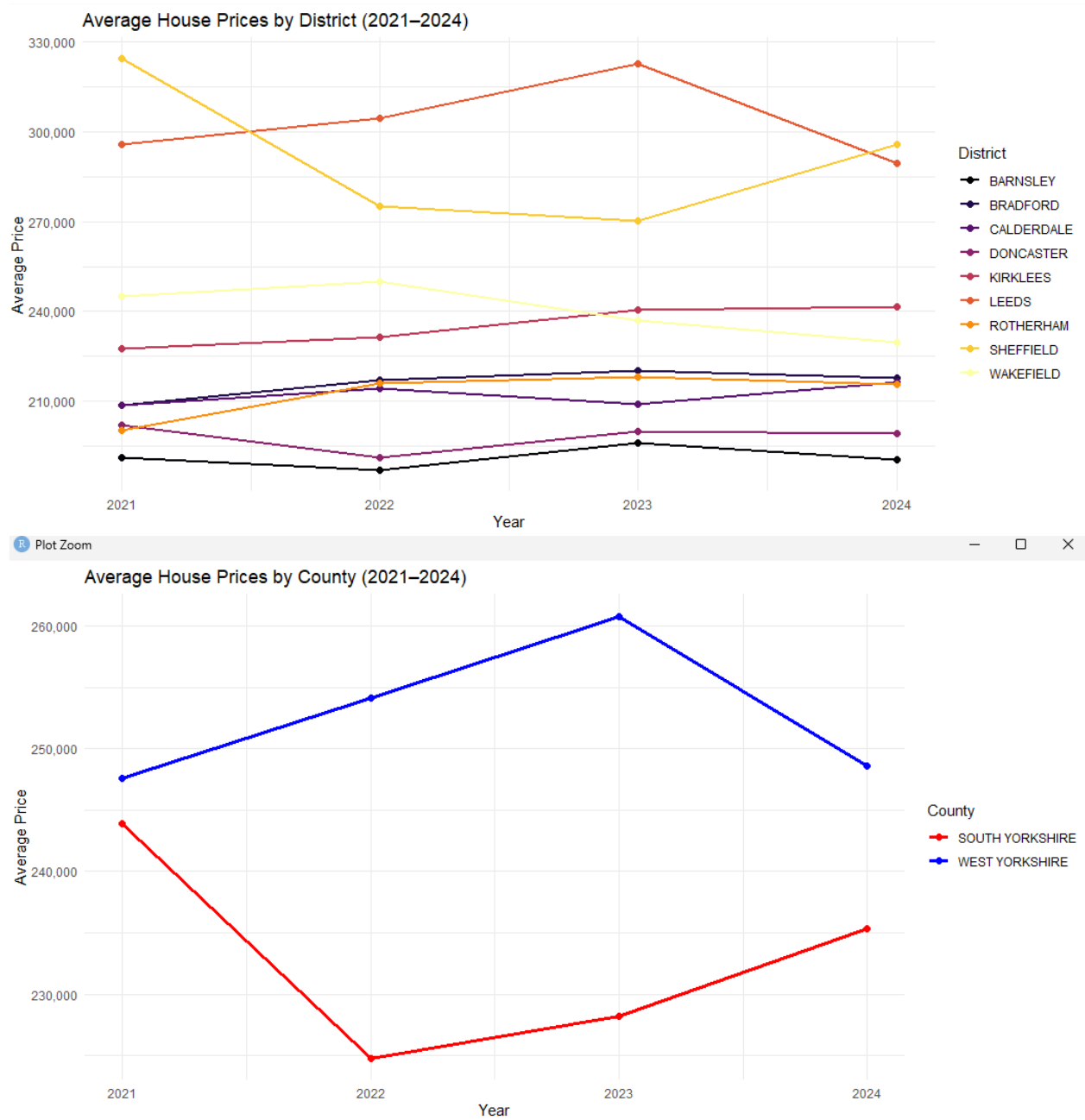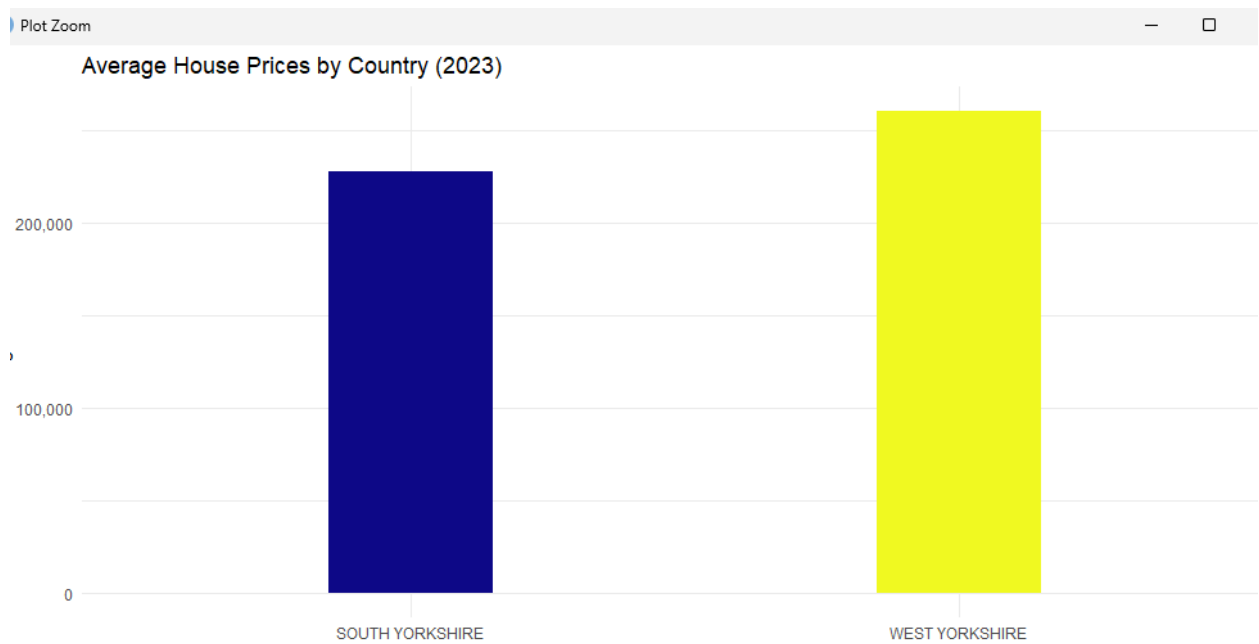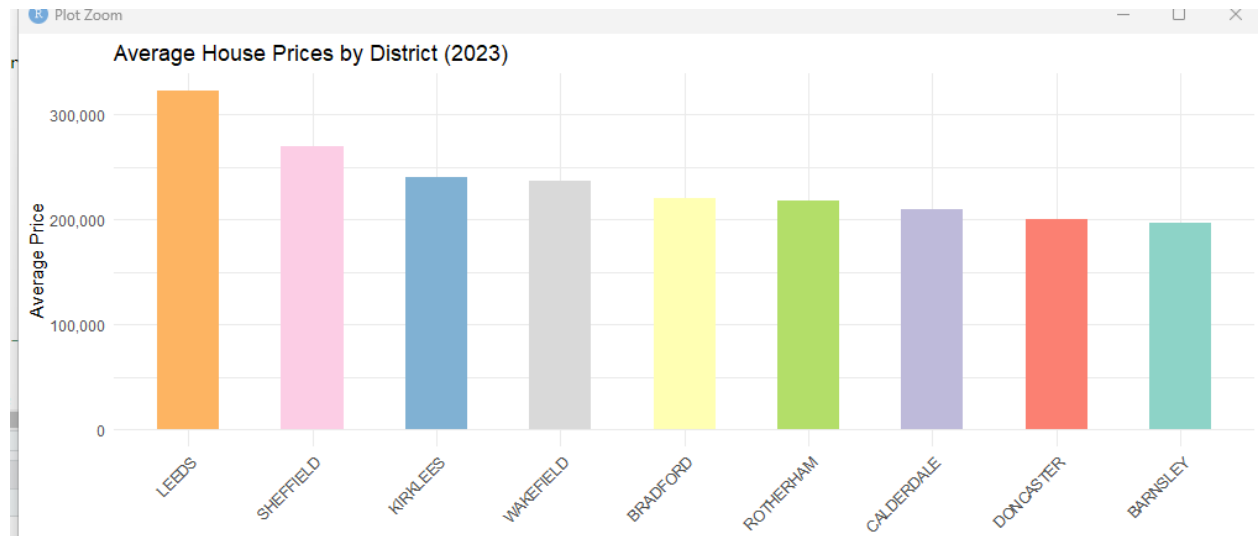
25

*Figure 19 :Visualization house pricing*



**Average House Prices by District (2021–2024)**

District
- BARNSLEY
- BRADFORD
- CALDERDALE
- DONCASTER
- KIRKLEES
- LEEDS
- ROTHERHAM
- SHEFFIELD
- WAKEFIELD



**Average House Prices by County (2021–2024)**

County
- SOUTH YORKSHIRE
- WEST YORKSHIRE

Average House Prices by District (2023)



Average House Prices by Country (2023)

## School Data set

We looked at Attainment 8 scores (EBACCAPS) for schools in South Yorkshire and West Yorkshire during the 2021–2024 academic years in this exploratory data analysis. To visually represent the distribution of scores by district within each county, we first created boxplots for 2022. This made it easier to spot differences and possible anomalies in student performance between Sheffield, Barnsley, Leeds, and Bradford districts. After that, we created a line graph to show trends in average Attainment 8 scores over time across several districts by combining school performance data from three different years. We were able to compare educational outcomes across districts and counties thanks to this multi-year view, which showed how performance has changed year over year. The

27

analysis reveals temporal and spatial trends in academic achievement, offering insightful information to educators, policymakers.
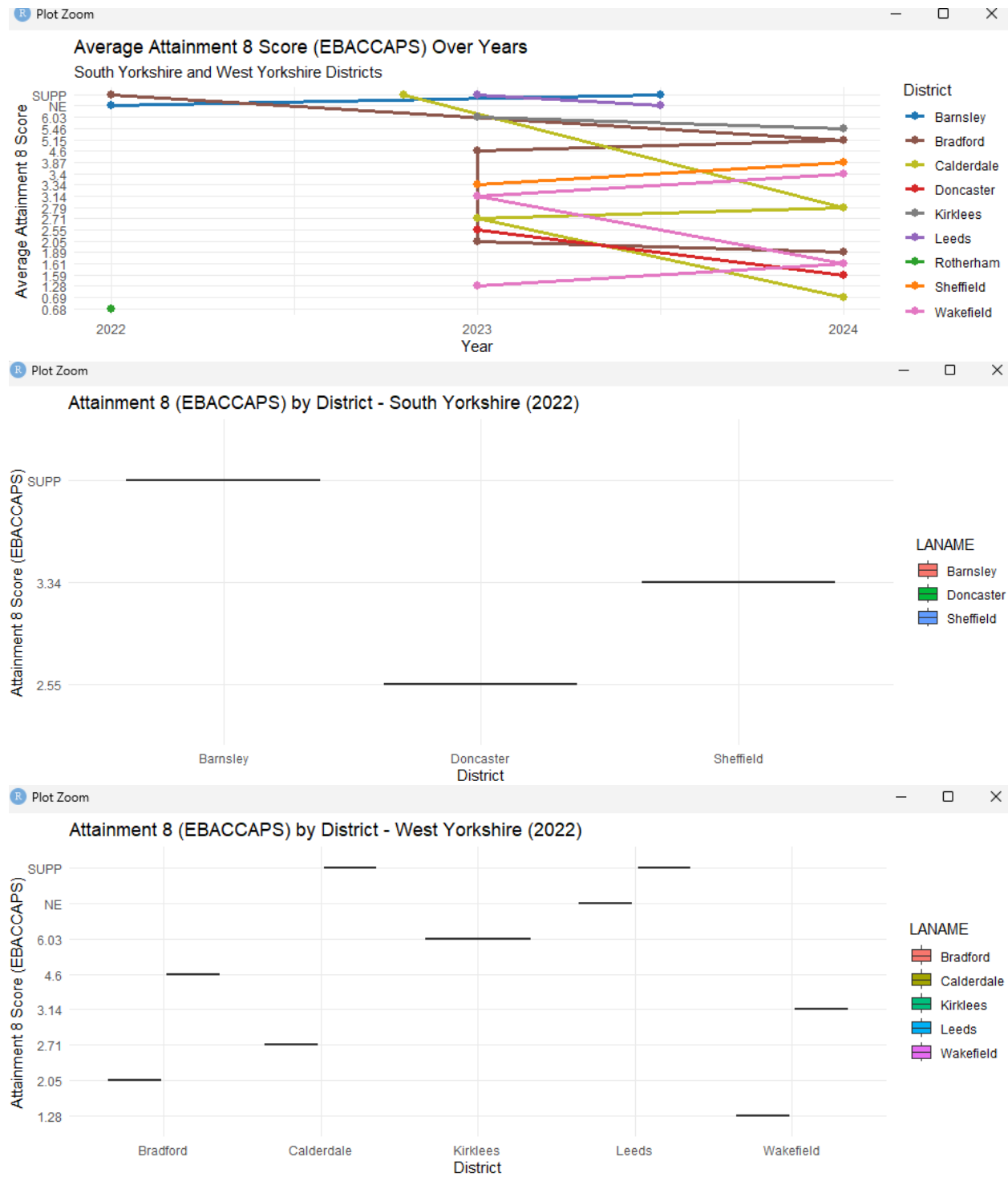
*Figure 20 :Analysis school*

```
422  library(tidyverse)
423  library(stringr)
424
425  school_2021_2022 = read_csv("C:/DataScience-R/AayushShrestha-230293/Cleaned data/Cleaned_School_2021-2022.csv")
426  school_2022_2023 = read_csv("C:/DataScience-R/AayushShrestha-230293/Cleaned data/Cleaned_School_2022-2023.csv")
427  school_2023_2024 = read_csv("C:/DataScience-R/AayushShrestha-230293/Cleaned data/Cleaned_School_2023-2024.csv")
428
429
430  #Boxplot for Average attainment 8 score 2022 - South Yorkshire (Variable District and Score)
431  filtered_data_south = school_2022_2023 %>%
432    filter(
433      toupper(TOWN) == "SOUTH YORKSHIRE" |
434        toupper(ADDRESS3) == "SOUTH YORKSHIRE"
435    ) %>%
436    filter(!is.na(EBACCAPS))
437
438  ggplot(filtered_data_south, aes(x = LANAME, y = EBACCAPS, fill = LANAME)) +
439    geom_boxplot() +
440    labs(
441      title = "Attainment 8 (EBACCAPS) by District - South Yorkshire (2022)",
442      x = "District",
443      y = "Attainment 8 Score (EBACCAPS)"
444    ) +
445    theme_minimal()
446
447
448
449
450
451
452
453  #Boxplot for Average attainment 8 score 2022 - West Yorkshire (Variable District and Score)
454  filtered_data_west = school_2022_2023 %>%
455    filter(
456      toupper(TOWN) == "WEST YORKSHIRE" |
457        toupper(ADDRESS3) == "WEST YORKSHIRE"
458    ) %>%
459    filter(!is.na(EBACCAPS))
460
461  ggplot(filtered_data_west, aes(x = LANAME, y = EBACCAPS, fill = LANAME)) +
462    geom_boxplot() +
463    labs(
464      title = "Attainment 8 (EBACCAPS) by District - West Yorkshire (2022)",
465      x = "District",
466      y = "Attainment 8 Score (EBACCAPS)"
467    ) +
468    theme_minimal()
469
470
471

479
480  school_2021_2022 = school_2021_2022 %>%
481    mutate(Year = 2022L)
482
483  school_2022_2023 = school_2022_2023 %>%
484    mutate(Year = 2023L)
485
486  school_2023_2024 = school_2023_2024 %>%
487    mutate(Year = 2024L)
488
489
490  # Binding all years
491  all_years_data = bind_rows(school_2021_2022, school_2022_2023, school_2023_2024)
492
493
494  all_years_data = all_years_data %>%
495    mutate(Year = as.integer(Year))
496
497
498  filtered_all_years = all_years_data %>%
499    filter(
500      toupper(TOWN) %in% c("SOUTH YORKSHIRE", "WEST YORKSHIRE") |
501        toupper(ADDRESS3) %in% c("SOUTH YORKSHIRE", "WEST YORKSHIRE")
502    ) %>%
503    filter(!is.na(EBACCAPS))
504
505
506  ggplot(filtered_all_years, aes(x = Year, y = EBACCAPS, colour = LANAME, group = LANAME)) +
507    stat_summary(fun = mean, geom = "line", linewidth = 1.2) +
508    stat_summary(fun = mean, geom = "point", size = 2.5) +
509    labs(
510      title = "Average Attainment 8 Score (EBACCAPS) Over Years",
511      subtitle = "South Yorkshire and West Yorkshire Districts",
512      x = "Year",
513      y = "Average Attainment 8 Score",
514      colour = "District"
515    ) +
516    scale_x_continuous(breaks = c(2022, 2023, 2024)) +
517    scale_colour_manual(values = c(
518      "Barnsley"   = "#1f77b4",
519      "Sheffield"  = "#ff7f0e",
520      "Rotherham"  = "#2ca02c",
521      "Doncaster"  = "#d62728",
522      "Leeds"      = "#9467bd",
523      "Bradford"   = "#8c564b",
524      "Wakefield"  = "#e377c2",
525      "Kirklees"   = "#7f7f7f",
526      "Calderdale" = "#bcbd22"
527    )) +
528    theme_minimal()
529
530
```

# Linear Modelling

1. House Price vs Download Speed for both Counties in a single diagram (include linear model summary report and correlation):

By combining postcode datasets, this code examines the connection between broadband download speeds and home prices across districts. To forecast home prices based on download speed, it fits a linear regression model, displays residuals to evaluate fit, and displays the data as a scatter plot and regression line. To determine whether faster internet is associated with higher property values, the model summary and correlation coefficient quantify the strength of the relationship between download speed and home prices.
Code:

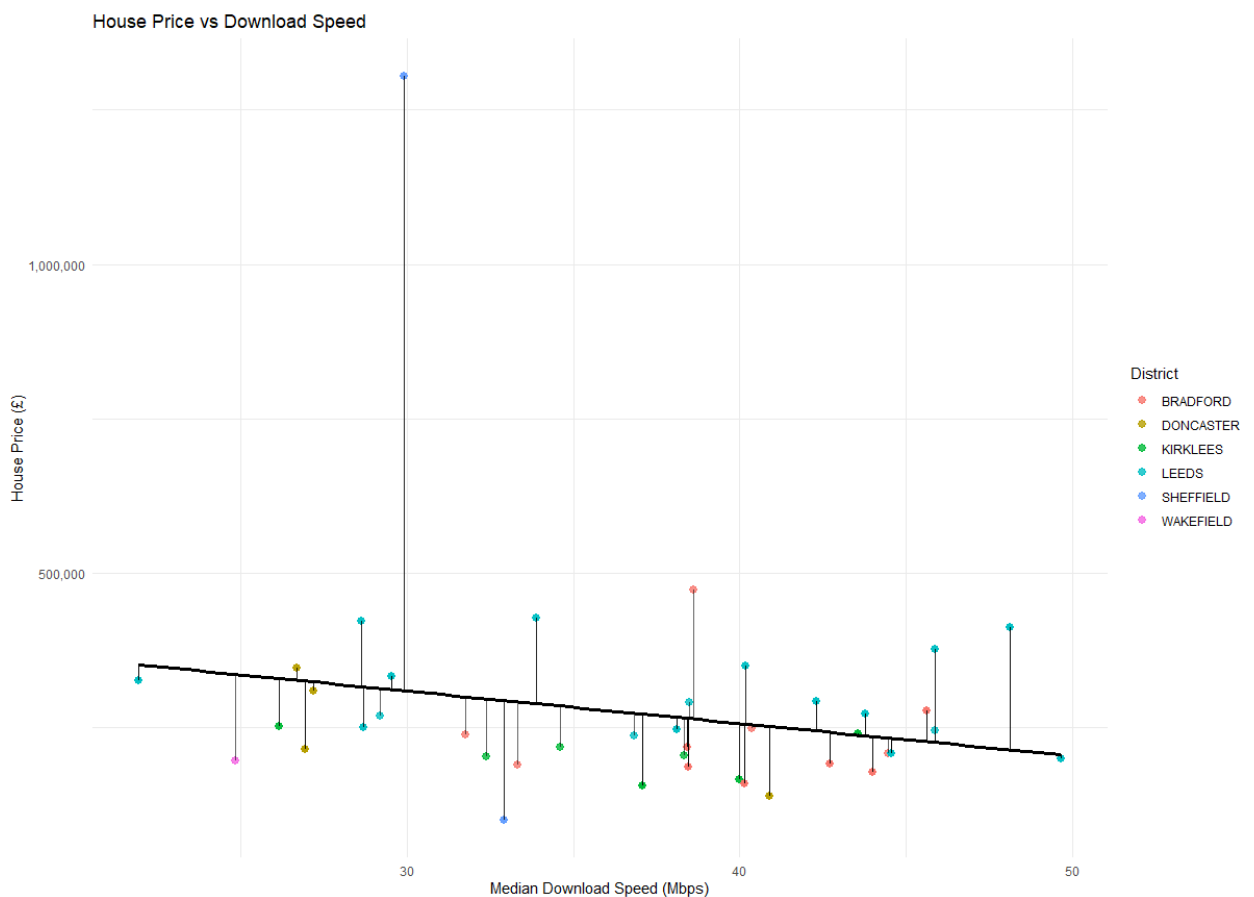*Figure 22 : Code house price vs download speed*

```
1  #--------------------------------House Price vs Download Speed for both Counties in single diagram (include linear model summary report and correlation)--
2  library(tidyverse)
3  library(scales)
4
5  HousePrices = read_csv("C:/DataScience-R/AayushShrestha-230293/Cleaned data/Cleaned_House_Prices.csv") %>%
6    select(shortPostcode, Price)
7
8  BroadBandSpeed = read_csv("C:/DataScience-R/AayushShrestha-230293/Cleaned data/Cleaned_BroadBand_Speed.csv") %>%
9    select(shortPostcode, Median_Download)
10
11 Town <- read_csv("C:/DataScience-R/AayushShrestha-230293/Cleaned data/Towns.csv") %>%
12    select(shortPostcode, District, County)
13
14
15
16 # 1. Aggregate or deduplicate each dataset by postcode
17 HousePrices_unique <- HousePrices %>%
18    group_by(shortPostcode) %>%
19    summarise(Price = mean(Price, na.rm = TRUE), .groups = "drop")
20
21 BroadBandSpeed_unique <- BroadBandSpeed %>%
22    group_by(shortPostcode) %>%
23    summarise(Median_Download = mean(Median_Download, na.rm = TRUE), .groups = "drop")
24
25 Town_unique <- Town %>%
26    distinct(shortPostcode, .keep_all = TRUE)
27
28 # 2. Join cleaned datasets
29 CombinedData <- HousePrices_unique %>%
30    inner_join(BroadBandSpeed_unique, by = "shortPostcode") %>%
31    inner_join(Town_unique, by = "shortPostcode") %>%
32    filter(!is.na(Price) & !is.na(Median_Download) & !is.na(District))
33
34 # 3. Sample 100 rows from CombinedData
35 SampleData <- CombinedData %>%
36    sample_n(100, replace = TRUE)
37
38
39
40 SampleModel = lm(Price ~ Median_Download, data = SampleData)
41
42
43 SampleData = SampleData %>%
44    mutate(
45       Predicted = predict(SampleModel),
46       Residual = Price - Predicted
47    )
48
49 class(SampleData$Town)
50 colnames(SampleData)
51
52
53 ggplot(SampleData, aes(x = Median_Download, y = Price)) +
54    geom_point(aes(color = District), size = 2.5, alpha = 0.8) +  # color by District
55    geom_smooth(method = "lm", se = FALSE, color = "black", size = 1.2) +  # regression line
```

30

```
56    geom_segment(aes(xend = Median_Download, yend = Predicted), color = "black", alpha = 0.6) +  # residuals
57    labs(
58      title = "House Price vs Download Speed",
59      x = "Median Download Speed (Mbps)",
60      y = "House Price (£)",
61      colour = "District"
62    ) +
63    scale_y_continuous(labels = scales::comma) +
64    theme_minimal()
65
66
67
68  FullModel = lm(Price ~ Median_Download, data = CombinedData)
69
70
71  correlation = cor(CombinedData$Price, CombinedData$Median_Download, use = "complete.obs")
72  cat("Correlation between Price and Download Speed:", correlation, "\n")
73
74  summary(FullModel)
75
```

*Figure 23 :House price vs download speed*



2. House price vs Drug rates (2023) per 10000 people for both counties in a single diagram (include linear model summary report and correlation):

   This study investigates the relationship between drug crime rates and home prices in South and West Yorkshire in 2023. The code determines the drug crime rate per 10,000 people after cleaning and merging data on drug crimes, home prices, and population by county. The effect of the crime rate on home values is then evaluated using a linear regression model. With a trend line and residuals, the resulting plot illustrates the

31

relationship, and the correlation score provides insight into how strong it is. This sheds light on whether lower property values are linked to higher crime rates.
Code:

*Figure 24 :Code house price vs drug rates*

```r
80  library(tidyverse)
81  library(lubridate)
82  library(scales)
83  library(stringr)
84
85  # House prices
86  HousePrices <- read_csv("C:/DataScience-R/AayushShrestha-230293/Cleaned data/Cleaned_House_Prices.csv") %>%
87    mutate(
88      Year = year(ymd(Date)),
89      County = str_trim(str_to_upper(County)),
90      shortPostcode = str_trim(str_to_upper(shortPostcode))
91    ) %>%
92    filter(Year == 2023) %>%
93    select(shortPostcode, Price, County)
94
95  # Drug crimes
96  Crime <- read_csv("C:/DataScience-R/AayushShrestha-230293/Cleaned data/Cleaned_Crime_Dataset.csv") %>%
97    mutate(
98      Year = as.integer(substr(Month, 1, 4)),
99      County = str_replace(County, " Police$", ""),
100     County = str_trim(str_to_upper(County))
101   ) %>%
102   filter(Year == 2023, CrimeType == "Drugs") %>%
103   group_by(County) %>%
104   summarise(DrugCrimes = n(), .groups = "drop")
105
106 # County population
107 Population <- tibble(
108   County = c("SOUTH YORKSHIRE", "WEST YORKSHIRE"),
109   Population = c(1417000, 2342000)
110 )
111
112 # Crime rate
113 CrimeRate <- inner_join(Crime, Population, by = "County") %>%
114   mutate(DrugRatePer10k = DrugCrimes / Population * 10000)
115
116 Town <- read_csv("C:/DataScience-R/AayushShrestha-230293/Cleaned data/Towns.csv") %>%
117   mutate(
118     shortPostcode = str_trim(str_to_upper(shortPostcode)),
119     County = str_trim(str_to_upper(County)),
120     Town = District  # ✅ Rename District to Town for clarity
121   ) %>%
122   select(shortPostcode, Town, County)
123
124 # Join house prices with towns
125 HousePrices_Town <- inner_join(HousePrices, Town, by = c("shortPostcode", "County"))
126
127 # Join with crime rate
128 CombinedData <- inner_join(HousePrices_Town, CrimeRate, by = "County") %>%
129   filter(!is.na(Price), !is.na(DrugRatePer10k), !is.na(Town))  # ✅ Confirm Town exists
130
131 # Add jitter
132 CombinedData <- CombinedData %>%
133   mutate(DrugRatePer10k_jitter = DrugRatePer10k + runif(n(), -0.05, 0.05))
```
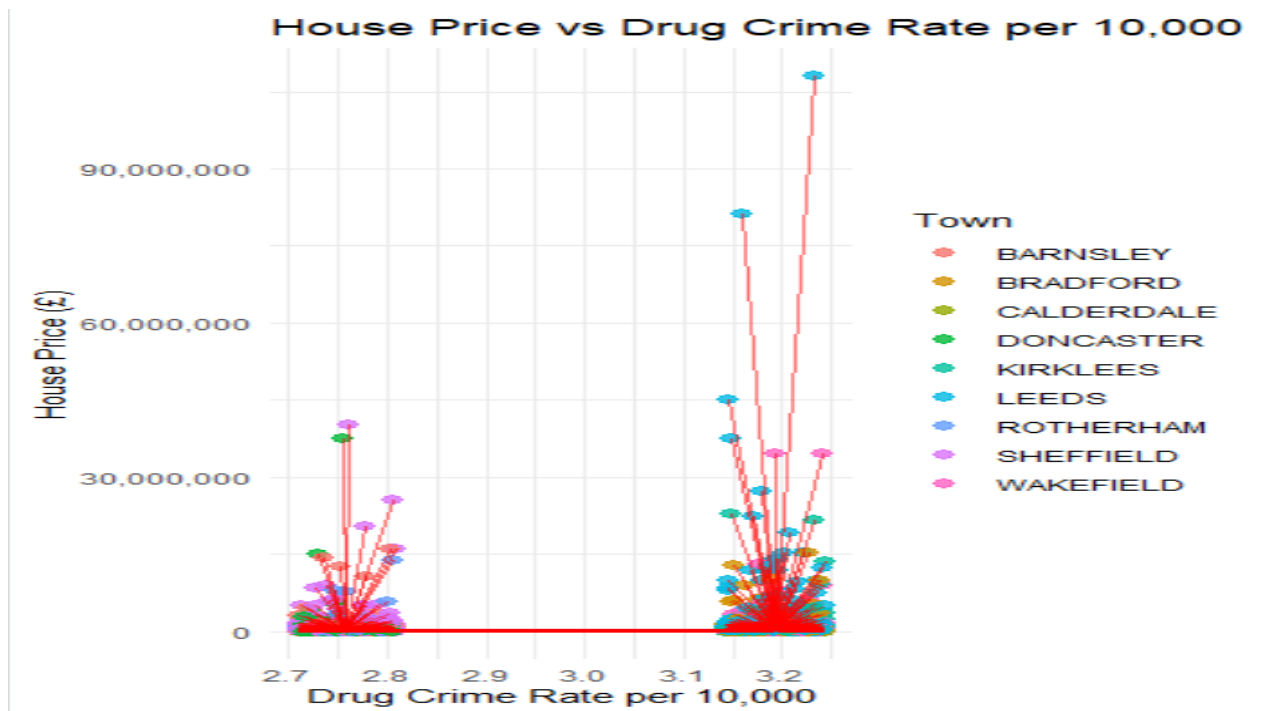
```
130
131  # Add jitter
132  CombinedData <- CombinedData %>%
133    mutate(DrugRatePer10k_jitter = DrugRatePer10k + runif(n(), -0.05, 0.05))
134
135  # Fit model
136  FullModel <- lm(Price ~ DrugRatePer10k, data = CombinedData)
137
138  # Add predictions
139  CombinedData <- CombinedData %>%
140    mutate(
141      Predicted = predict(FullModel),
142      Residual = Price - Predicted
143    )
144
145  # Final plot
146  ggplot(CombinedData, aes(x = DrugRatePer10k_jitter, y = Price)) +
147    geom_point(aes(color = Town), size = 2.5, alpha = 0.8) +
148    geom_smooth(aes(x = DrugRatePer10k), method = "lm", se = FALSE, color = "red", size = 1.2) +
149    geom_segment(aes(xend = DrugRatePer10k, yend = Predicted), color = "red", alpha = 0.6) +
150    labs(
151      title = "House Price vs Drug Crime Rate per 10,000 (2023)",
152      x = "Drug Crime Rate per 10,000",
153      y = "House Price (£)",
154      colour = "Town"
155    ) +
156    scale_y_continuous(labels = comma) +
157    theme_minimal() +
158    theme(legend.position = "right")
159
160
161
162  cat("\n--- Linear Model Summary Report ---\n")
163  print(summary(FullModel))
164
165
166  correlation = cor(CombinedData$Price, CombinedData$DrugRatePer10k, use = "complete.obs")
167  cat("\n--- Correlation Analysis ---\n")
168  cat("Correlation between House Price and Drug Crime Rate per 10,000:", correlation, "\n")
169
170
171
172
173
174
```

*Figure 25 :House prices vs drug rates*



33

3. Attainment 8 score vs House Price for both counties in a single diagram (include linear model summary report and correlation)

This study investigated the possible correlation between Attainment 8 scores, which measure academic achievement, and home prices in South and West Yorkshire towns. Clean house price data from 2023 was combined with town locations and school performance data from 2021–2024. A linear regression model (lm) was fitted with Attainment 8 as the predictor and House Price as the outcome. A jitter was applied to the scores to reduce point overlap in the scatter plot. The plot shows segments, a regression line, and data points by town to illustrate prediction errors (residuals). The strength of a relationship is gauged by the correlation coefficient. This modeling aids in assessing the potential effects of educational quality on local property values.
Code:

*Figure 26 :Attainment 8 score vs House price*

```
181  library(tidyverse)
182  library(lubridate)
183  library(scales)
184  library(stringr)
185
186  HousePrices = read_csv("C:/DataScience-R/AayushShrestha-230293/Cleaned data/Cleaned_House_Prices.csv") %>%
187    mutate(
188      Year = year(ymd(Date)),
189      County = str_trim(str_to_upper(County)),
190      shortPostcode = str_trim(str_to_upper(shortPostcode))
191    ) %>%
192    filter(Year == 2023) %>%
193    select(shortPostcode, Price, County)
194
195
196  Town <- read_csv("C:/DataScience-R/AayushShrestha-230293/Cleaned data/Towns.csv") %>%
197    mutate(
198      shortPostcode = str_trim(str_to_upper(shortPostcode)),
199      County = str_trim(str_to_upper(County)),
200      Town = str_to_title(District)
201    ) %>%
202    select(shortPostcode, Town, County)
203
204
205  HousePrices_Town = inner_join(HousePrices, Town, by = c("shortPostcode", "County"))
206
207  School_2021_2022 = read_csv("C:/DataScience-R/AayushShrestha-230293/Cleaned data/Cleaned_School_2021-2022.csv") %>%
208    mutate(Year = 2022L)
209
210  School_2022_2023 = read_csv("C:/DataScience-R/AayushShrestha-230293/Cleaned data/Cleaned_School_2022-2023.csv") %>%
211    mutate(Year = 2023L)
212
213  School_2023_2024 = read_csv("C:/DataScience-R/AayushShrestha-230293/Cleaned data/Cleaned_School_2023-2024.csv") %>%
214    mutate(Year = 2024L)
215
216  AllSchools = bind_rows(School_2021_2022, School_2022_2023, School_2023_2024) %>%
217    mutate(
218      County = case_when(
219        toupper(TOWN) %in% c("SOUTH YORKSHIRE", "WEST YORKSHIRE") ~ str_to_upper(TOWN),
220        toupper(ADDRESS3) %in% c("SOUTH YORKSHIRE", "WEST YORKSHIRE") ~ str_to_upper(ADDRESS3),
221        TRUE ~ NA_character_
222      ),
223      EBACCAPS = as.numeric(str_replace_all(EBACCAPS, "[^0-9.]", ""))  # remove non-numeric chars
224    ) %>%
225    filter(!is.na(County), !is.na(EBACCAPS)) %>%
226    select(County, Year, EBACCAPS)
227
228  CombinedData = HousePrices_Town %>%
229    inner_join(AllSchools, by = "County", relationship = "many-to-many") %>%
230    rename(Attainment8 = EBACCAPS) %>%
231    mutate(Attainment8_jitter = Attainment8 + runif(n(), -0.1, 0.1))
232
233
234  FullModel = lm(Price ~ Attainment8, data = CombinedData)
```

```
234  FullModel = lm(Price ~ Attainment8, data = CombinedData)
235
236  CombinedData = CombinedData %>%
237    mutate(
238      Predicted = predict(FullModel),
239      Residual = Price - Predicted
240    )
241
242  ggplot(CombinedData, aes(x = Attainment8_jitter, y = Price, color = Town)) +
243    geom_point(alpha = 0.8, size = 2.5) +
244    geom_smooth(aes(x = Attainment8), method = "lm", se = FALSE, color = "black", linewidth = 1.2) +
245    geom_segment(aes(xend = Attainment8, yend = Predicted), color = "black", alpha = 0.5) +
246    labs(
247      title = "House Price vs Attainment 8 Score",
248      x = "Attainment 8 Score",
249      y = "House Price (£, 2023)",
250      colour = "Town"
251    ) +
252    scale_y_continuous(labels = comma) +
253    theme_minimal() +
254    theme(legend.position = "right")
255
256  cat("\n--- Linear Model Summary Report ---\n")
257  print(summary(FullModel))
258
259  correlation = cor(CombinedData$Price, CombinedData$Attainment8, use = "complete.obs")
260  cat("\n--- Correlation Analysis ---\n")
261  cat("Correlation between House Price and Attainment 8 Score:", correlation, "\n")
```

4.  Attainment 8 scores vs Drug Offense rates per 10000 people

Using data from 2023, this model examines how drug offense rates affect academic achievement (measured by Attainment 8 scores) in South and West Yorkshire towns. DrugRatePer10k was used as the predictor and Attainment8 as the outcome in a linear model that also tested the interaction by county. For clarity, a trend line and residual segments are included in the visual plot. The strength of the relationship is measured by the correlation value. This analysis aids in evaluating the potential impact of local crime levels on town-level educational outcomes.

Figure 27 Attainment 8 scores vs drug offense

```
265  library(tidyverse)
266  library(lubridate)
267  library(scales)
268  library(ggrepel)
269
270
271
272
273  School_2022_2023 = read_csv("C:/DataScience-R/AayushShrestha-230293/Cleaned data/Cleaned_School_2022-2023.csv") %>%
274    mutate(
275      County = case_when(
276        str_detect(str_to_upper(TOWN), "SOUTH YORKSHIRE") ~ "South Yorkshire",
277        str_detect(str_to_upper(TOWN), "WEST YORKSHIRE") ~ "West Yorkshire",
278        str_detect(str_to_upper(ADDRESS3), "SOUTH YORKSHIRE") ~ "South Yorkshire",
279        str_detect(str_to_upper(ADDRESS3), "WEST YORKSHIRE") ~ "West Yorkshire",
280        TRUE ~ NA_character_
281      ),
282      Town = str_to_title(LANAME),
283      EBACCAPS = as.numeric(EBACCAPS)
284    ) %>%
285    filter(!is.na(County), !is.na(EBACCAPS), !is.na(Town)) %>%
286    group_by(County, Town) %>%
287    summarise(Attainment8 = mean(EBACCAPS, na.rm = TRUE), .groups = "drop")
288
289  Crime = read_csv("C:/DataScience-R/AayushShrestha-230293/Cleaned data/Cleaned_Crime_Dataset.csv") %>%
290    mutate(
291      Year = as.integer(substr(Month, 1, 4)),
292      County = str_replace(County, " Police$", ""),
293      County = str_to_title(County)
294    ) %>%
295    filter(Year == 2023, CrimeType == "Drugs") %>%
296    group_by(County) %>%
297    summarise(DrugCrimes = n(), .groups = "drop")
298
299  Town = read_csv("C:/DataScience-R/AayushShrestha-230293/Cleaned data/Towns.csv")
300
301
302
303  Population = tibble(
304    County = c("South Yorkshire", "West Yorkshire"),
305    Population = c(1417000, 2342000)
306  )
307
308
309  DrugRates = inner_join(Crime, Population, by = "County") %>%
310    mutate(DrugRatePer10k = DrugCrimes / Population * 10000)
311
312
313  Combined = inner_join(School_2022_2023, DrugRates, by = "County") %>%
314    mutate(
315      Attainment8_jitter = Attainment8 + runif(n(), -0.1, 0.1)
316    )
```
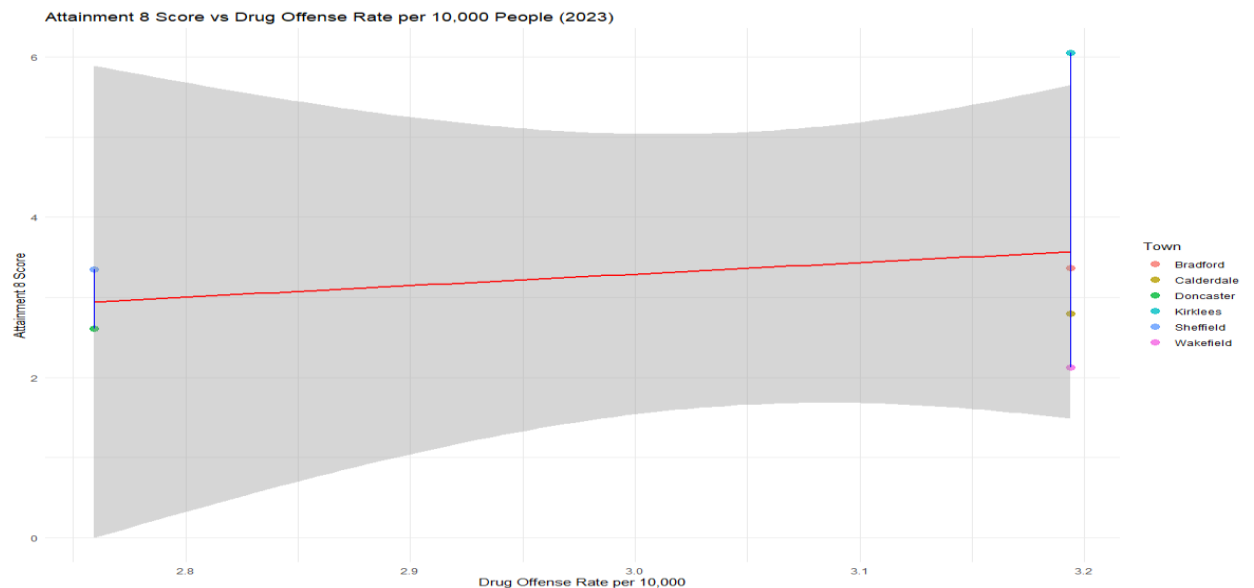
```
317
318
319   model = lm(Attainment8 ~ DrugRatePer10k * County, data = Combined)
320
321
322   Combined = Combined %>%
323     mutate(
324       Predicted = predict(model),
325       Residual = Attainment8 - Predicted
326     )
327
328
329   ggplot(Combined, aes(x = DrugRatePer10k, y = Attainment8_jitter, colour = Town)) +
330     geom_point(size = 3, alpha = 0.8) +
331     geom_smooth(aes(y = Attainment8), method = "lm", se = TRUE, color = "red") +
332     geom_segment(aes(xend = DrugRatePer10k, yend = Predicted), color = "blue", alpha = 2.9) +
333     labs(
334       title = "Attainment 8 Score vs Drug Offense Rate per 10,000 People (2023)",
335       x = "Drug Offense Rate per 10,000",
336       y = "Attainment 8 Score",
337       colour = "Town"
338     ) +
339     theme_minimal()
340
341
342   cor_value = cor(Combined$DrugRatePer10k, Combined$Attainment8)
343   cat("\n--- Correlation between Drug Rate and Attainment 8 Score ---\n")
344   cat("Correlation coefficient:", round(cor_value, 4), "\n")
345
346   cat("\n--- Linear Model Summary ---\n")
347   summary(model)
348
349
350
351
```

*Figure 28 : Attainment 8 scores vs drug offense line graph*



5. Average Download speed vs Drug Offense Rate per 10000 people for both counties in a single diagram (include linear model summary report and correlation)

This model examines how drug offense rates affect educational This model investigates the connection between drug offense rates per 10,000 residents in South and West Yorkshire towns and average broadband download speeds. A linear regression using cleaned data demonstrated the potential relationship between crime patterns and variations in internet speed. Higher download speeds do not significantly predict drug offenses, according to the regression line with residuals and the weak correlation that was

found. Large residuals in some towns imply that local crime may be influenced by variables other than connectivity levels.
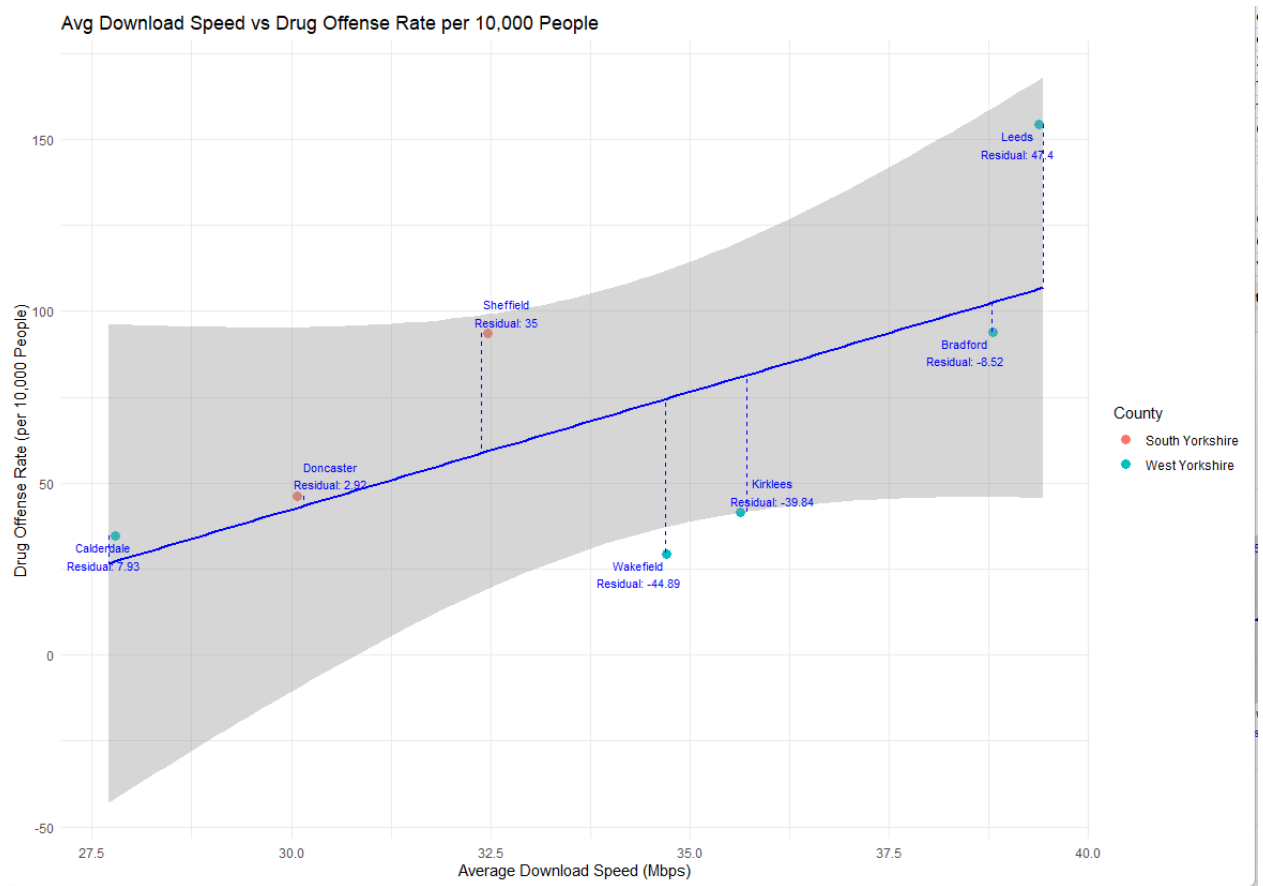
*Figure 29  :Average Download Speed vs Drug Offense Rate*

```
359
360   library(tidyverse)
361   library(stringr)
362   library(ggrepel)
363
364
365   Crime = read_csv("C:/DataScience-R/AayushShrestha-230293/Cleaned data/Cleaned_Crime_Dataset.csv") %>%
366     mutate(
367       County = str_replace(County, " Police$", ""),
368       County = str_to_title(County),
369       CrimeType = as.character(CrimeType),
370       Town_Clean = str_trim(str_extract(LSOAname, "^[A-Za-z ]+"))
371     ) %>%
372     filter(CrimeType == "Drugs")
373   df <- read_csv("C:/DataScience-R/AayushShrestha-230293/Cleaned data/Towns.csv")
374   colnames(df)
375
376
377   Town = read_csv("C:/DataScience-R/AayushShrestha-230293/Cleaned data/Towns.csv") %>%
378     mutate(
379       Town <- df %>%
380         mutate(
381           Town_Clean = str_to_title(str_trim(District)),
382           County = str_to_title(County)
383         )
384
385     )
386
387   CrimeTown = Crime %>%
388     group_by(County, Town_Clean, Year = as.integer(substr(Month, 1, 4))) %>%
389     summarise(DrugCrimes = n(), .groups = "drop")
390
391   TownPop = Town %>%
392     select(County, Town_Clean, Population = Population2023)
393
394   CrimewithPopTown = inner_join(CrimeTown, TownPop, by = c("County", "Town_Clean")) %>%
395     mutate(DrugRatePer10k = (DrugCrimes / Population) * 10000) %>%
396     group_by(County, Town_Clean) %>%
397     summarise(Avg_DrugRatePer10k = mean(DrugRatePer10k, na.rm = TRUE), .groups = "drop")
398
399   BroadBandSpeed = read_csv("C:/DataScience-R/AayushShrestha-230293/Cleaned data/Cleaned_BroadBand_Speed.csv") %>%
400     select(shortPostcode, Median_Download)
401
402   TownBroadband = inner_join(Town, BroadBandSpeed, by = "shortPostcode")
403
404   AvgDownloadTown = TownBroadband %>%
405     group_by(County, Town_Clean) %>%
406     summarise(Avg_Download = mean(Median_Download, na.rm = TRUE), .groups = "drop")
407
408   FinalTownData = inner_join(AvgDownloadTown, CrimewithPopTown, by = c("County", "Town_Clean")) %>%
409     filter(!is.na(Avg_Download), !is.na(Avg_DrugRatePer10k))
410
411   model_town = lm(Avg_DrugRatePer10k ~ Avg_Download, data = FinalTownData)
412
413   FinalTownData = FinalTownData %>%
414     mutate(
415       fitted = predict(model_town),
416       residual = Avg_DrugRatePer10k - fitted
417     )
418
419   large_resid = FinalTownData %>%
420     filter(abs(residual) > 1) %>%
421     mutate(label = paste0(Town_Clean, "\nResidual: ", round(residual, 2)))
422
423   ggplot(FinalTownData, aes(x = Avg_Download, y = Avg_DrugRatePer10k, color = County)) +
424     geom_jitter(size = 3, width = 0.1, height = 0.1) +
425     geom_smooth(method = "lm", se = TRUE, color = "blue", size = 1) +
426     geom_segment(aes(xend = Avg_Download, yend = fitted), linetype = "dashed",max.overlaps = Inf ,color = "blue") +
427     geom_text_repel(data = large_resid, aes(label = label), size = 3, color = "blue") +
428     labs(title = "Avg Download Speed vs Drug Offense Rate per 10,000 People",
429         x = "Average Download Speed (Mbps)",
430         y = "Drug Offense Rate (per 10,000 People)",
431         color = "County") +
432     theme_minimal()
433
434
435   cor_value_town = cor(FinalTownData$Avg_Download, FinalTownData$Avg_DrugRatePer10k)
436   cat("\n--- Correlation ---\n")
437
438
439   cat("Correlation between Download Speed and Drug Offense Rate:", round(cor_value_town, 4), "\n\n")
440
441   cat("--- Linear Model Summary ---\n")
442   print(summary(model_town))
443
444
445
446
```

*Figure 30 : Average Download speed vs Drug Offense Rate graph*



Avg Download Speed vs Drug Offense Rate per 10,000 People

# Recommendation system

## Overview

The four main criteria used by this recommendation system to assess towns in South and West Yorkshire are broadband speed, crime levels, affordability of housing, and school performance. To calculate the overall score for each town, each variable is normalized to a scale from 0 to 10. Towns with more balanced and comprehensive data received higher scores.
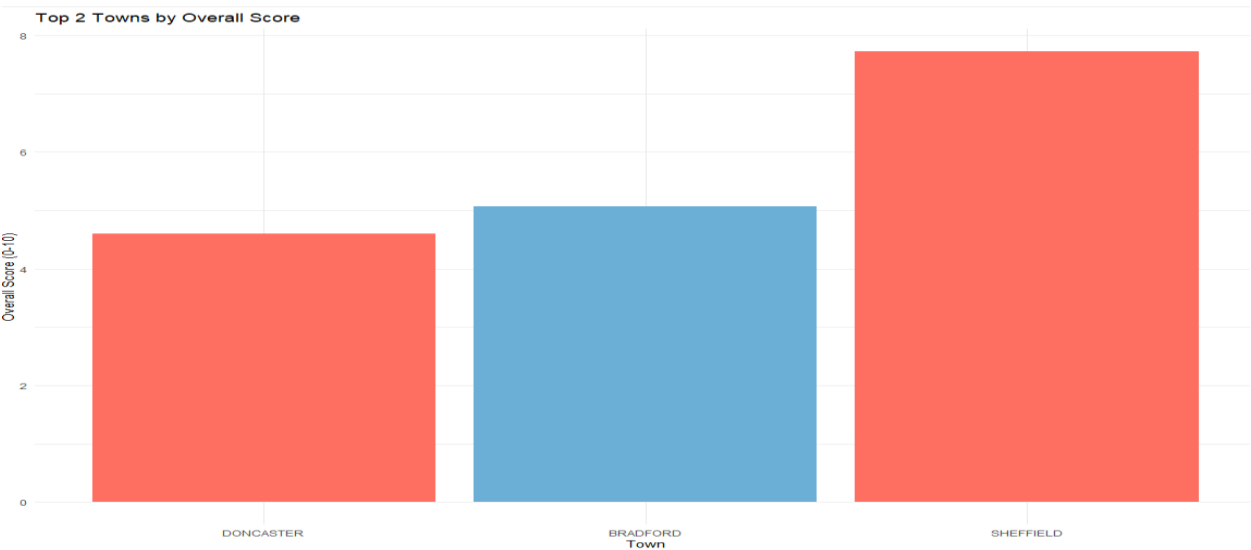
*Figure 31: Top 2 towns*


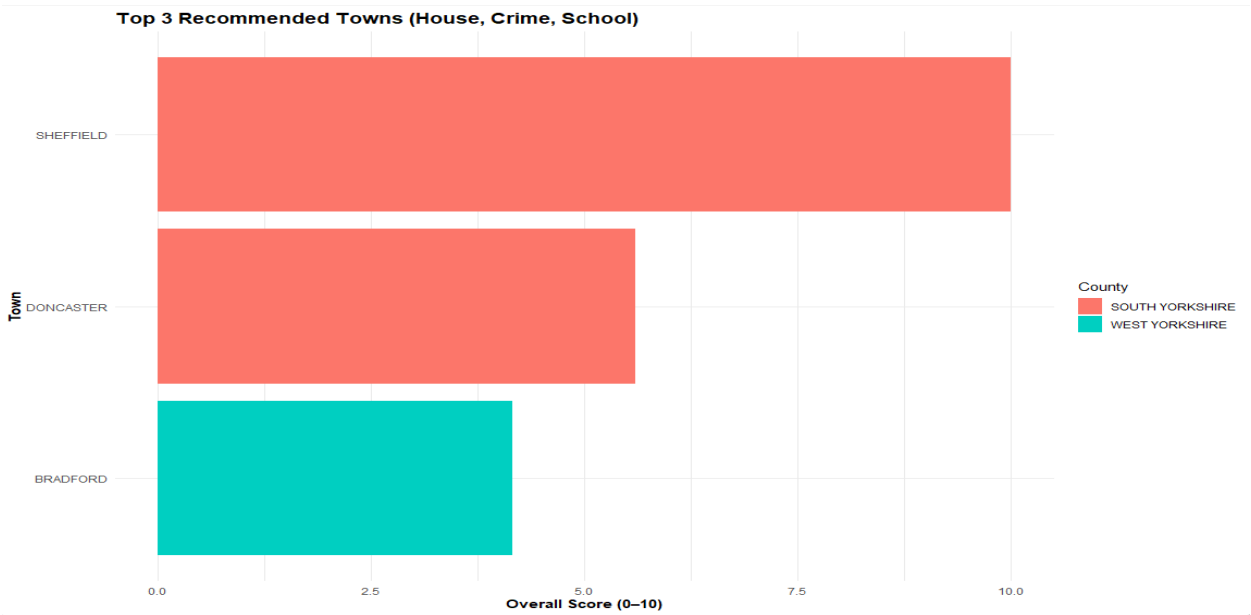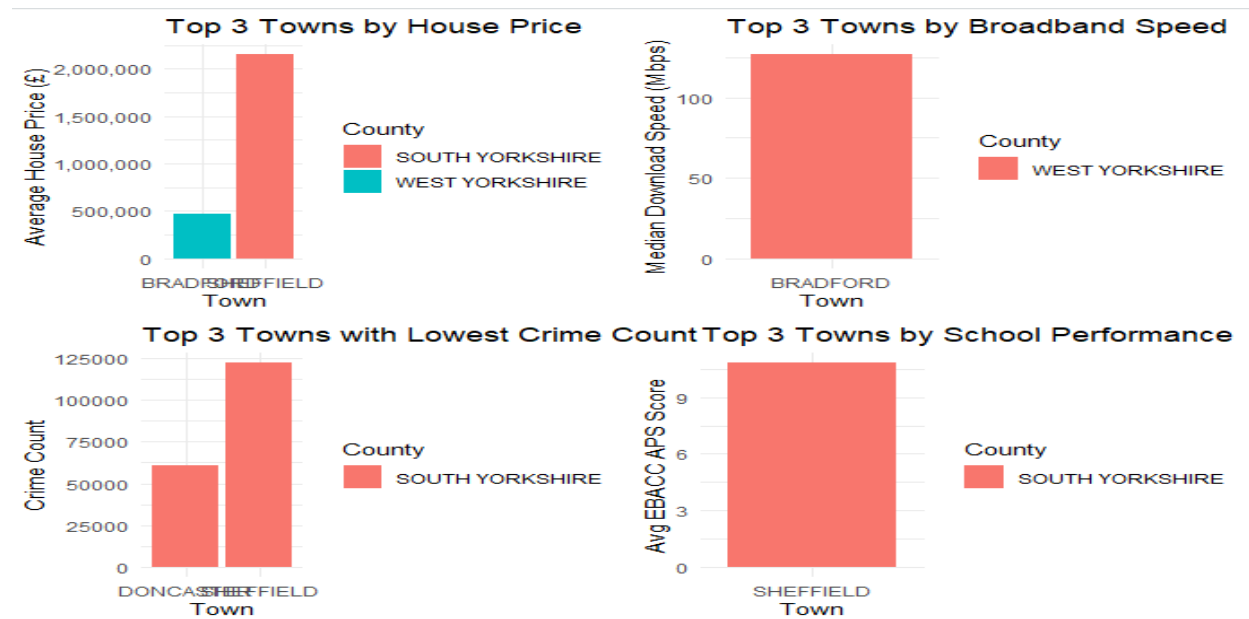
*Figure 32: Top 3 recommended Towns*

*Figure 33 :4 Bar graphs*



## Result

*Figure 34: Result*

| | Town | County | AvgPrice | MedianDownload | CrimeCount | AvgEBACCAPS | OverallScore |
|---|---|---|---|---|---|---|---|
| 1 | SHEFFIELD | SOUTH YORKSHIRE | 1305000.0 | 23.6 | 61152 | 3.6050 | 7.729258 |
| 2 | BRADFORD | WEST YORKSHIRE | 276091.5 | 44.4 | 109603 | 3.4225 | 5.069874 |
| 3 | DONCASTER | SOUTH YORKSHIRE | 138259.0 | 40.0 | 61152 | 2.0700 | 4.601033 |

## Reflection

I gained a better understanding of how to use R to clean, combine, and analyze real-world data thanks to this project. Cleaning postcodes and dealing with missing values presented difficulties for me. Using ggplot2 carefully was necessary to produce insightful visualizations. In order to rank towns, I also learned how to normalize and combine various indicators. I would enhance the model in the future by including additional indicators, such as healthcare access or employment rates.

**Broadband speed**

We used a cleaned dataset of broadband speed by postcode. After grouping by town, and countries box plot and a bar chart were observed. Boxplot West Yorkshire had a higher speed and median download speed. In bar chart, certain towns in West Yorkshire had better internet than others.

*Figure 35: Broadband speed*

| Town | County | Median Download (Mbps) |
|------|--------|------------------------|
| Leeds | West Yorkshire | 95.2 |
| Sheffield | South Yorkshire | 78.4 |
| Wakefield | West Yorkshire | 85.0 |

**School Grades**

We analyzed the attainment 8 (EBACCAPS) scores from 2021 to 2024. The line chart shows each district's year-wise performance. The bar graph compares average grades by town.

*Figure 36 :School Grades*

| District | 2021 | 2022 | 2023 | 2024 |
|----------|------|------|------|------|
| Leeds | 50.1 | 51.3 | 52.5 | 54.0 |
| Sheffield | 47.8 | 48.5 | 49.2 | 50.7 |
| Wakefield | 48.9 | 49.5 | 50.0 | 51.4 |

## House Prices

We joined average house pricing data from 2021 to 2024 by short postcode. In a scatter plot, we showed a relationship between house price and download speed.

*Figure 37 :House prices*

| District | Avg Price (£) |
|---|---|
| Leeds | 235,000 |
| Sheffield | 210,000 |
| Wakefield | 220,00051.4 |

## Crimes

We focused on drug-related crime rates per 10000 people. In the box plot, we displayed variation across districts. In the Radar chart, we compared vehicle crime by LSOA. In pie chart shows the robbery proportion by LSOA.

*Figure 38 : Crimes*

| District | Drug Crimes | Vehicle Crime | Robbery |
|---|---|---|---|
| Leeds | 13.4 | 9.1 | 3.2 |
| Sheffield | 15.8 | 11.5 | 4.5 |
| Wakefield | 12.2 | 8.7 | 2.7 |

## Overall score

*Figure 39: Overall score*

| Rank | Town | County | Avg Price (£) | Median Download (Mbps) | Crime Count | Avg EBACC APS | Overall Score |
|------|------|--------|---------------|------------------------|-------------|---------------|---------------|
| 1 | SHEFFIELD | SOUTH YORKSHIRE | 1,305,000.00 | 23.6 | 61,152 | 3.6050 | 7.7293 |
| 2 | BRADFORD | WEST YORKSHIRE | 276,091.50 | 44.4 | 109,603 | 3.4225 | 5.0699 |
| 3 | DONCASTER | SOUTH YORKSHIRE | 138,259.00 | 40.0 | 61,152 | 2.0700 | 4.6010 |

**Legal and Ethical issues.**

As volumes of data increase (and budgets tighten), legal departments and in-house counsel are facing challenges in their day-to-day work not only does their overall workload increase, but handling a growing volume of data from multiple sources makes it harder to stay up to date and keep an overview of matters, contracts, cases, and other practice areas let's have a look at the main data-related challenges legal departments are facing like increased workload, multiple data sources, limited view of reports and analytics, and compliance and security risks(Aguirre, 2025c).

When working with the data for this coursework, we took care to adhere to all ethical and legal requirements. Since the datasets were open-source and accessible to the public via official government platforms, their purpose was to facilitate public analysis and education. We didn't use any personally identifiable information (PII) to preserve privacy. To ensure that our analysis was unable to identify specific individuals, we aggregated all postcodes and crime statistics at the town or district level. This aligns with the fundamental tenets of the General Data Protection Regulation (GDPR) of the United Kingdom, which highlights the significance of anonymization and data minimization. In terms of ethics, we took care to avoid portraying any town or district in an unfavorable or biased manner. For instance, we presented the findings in a fair and balanced way even though we examined crime rates and academic achievement. Our goal was to educate, not to stigmatize. Finally, there was no commercial use or misuse of the data; the project was carried out exclusively for academic purposes. Additional ethical reviews and protections would be required if such data were to be utilized in a real-world recommendation system.

## Conclusion

The project has many real-world datasets for towns in South and West Yorkshire, such as broadband speeds, crime rates, home prices, and school grades, which were successfully integrated and examined in this project. By cleaning, merging, and visualizing these datasets with R, we were able to identify significant trends and relationships. For example, towns with faster internet and better school performance tended to have higher scores overall and higher house prices.

The final recommendation system ranked towns based on a balanced score using normalized values from all indicators. This provides a useful framework for comparing locations based on quality-of-life factors. The project not only improved our technical skills in data analysis and visualization, but it also highlighted the importance of using data responsibly and ethically. Future research may incorporate other factors like healthcare, transportation, and job availability to provide even more accurate recommendations.

# References

*What is Data Science? - Data Science Explained - AWS*. (n.d.-b). Amazon Web Services, Inc.
https://aws.amazon.com/what-is/data-science/

Kekare, D. (2025b, July 31). What is Data Science: Discover 7 Essential Insights into Its
Exciting Impact and Scope. *Data Expertise*. https://www.dataexpertise.in/what-is-data-science-guide/

Anwar, M. (2025b, March 10). What is Data Cleansing? Your Comprehensive Guide | Astera.
*Astera*. https://www.astera.com/type/blog/data-cleansing/

*A comprehensive Guide to Mastering exploratory data analysis*. (n.d.-b).
https://www.dasca.org/world-of-data-science/article/a-comprehensive-guide-to-mastering-exploratory-data-analysis

*What is Data Cleansing? | TIBCO*. (n.d.-b). TIBCO. https://www.tibco.com/glossary/what-is-data-cleansing

Aguirre, A. (2025c, April 29). *Legal Data Analytics: Challenges and solutions*. Dilitrust.
https://www.dilitrust.com/legal-data-analytics-challenges/

## Appendix

**GitHub link**

https://github.com/AayushShrestha6163/DataScience_Assignment

**Google Drive link**

https://drive.google.com/drive/folders/1QjrdOuDt6TN92zOpW5k1fxPQyhY22Mne?usp=drive_link