Report on Town Recommendation System

Aashista Karki

B.Sc. Computing Softwarica College of IT and E-commerce

ST5014CEM Data Science for Developers

Siddhartha Neupane

Jan 24 2024

Table of Contents

# Table of Figures

**Introduction**

This Individual Coursework endeavors to develop a town recommender system catering to the needs of international students planning a study exchange program in England. The primary goal is to assist in making informed decisions about towns in the counties of Kent and Surrey. Recommendations are based on a detailed analysis of various factors such as educational institutions, cost of living (exemplified by house prices), broadband speed, and safety ratings (illustrated by local crime statistics).

The exclusive use of datasets released by the UK government ensures the reliability of the analysis. These datasets are sourced from reputable entities, primarily accessed through the https://data.gov.uk/ website, aligning with ethical and legal considerations.

The steps involved in this project is first data gathering using official uk government data, secondly data cleaning using r language, thirdly to plot out various graphs as the need of project and lastly to recommend town based on the top scores.

The report unfolds with a comprehensive Exploratory Data Analysis (EDA) section, using graphical plots and summary statistics to comprehend the distribution of single-variable data, identify outliers, and investigate relationships between variables through plots and correlation coefficients.

The report concludes with a discussion of legal and ethical considerations associated with the data used and the recommendation system developed. Additionally, it reflects on the application of the data mining lifecycle to this problem, summarizes conclusions, and offers recommendations for future improvements or extensions.

**Cleaning data**

Cleaning data is the process of detecting and rectifying errors in datasets, including addressing issues like missing values, duplicates, and inconsistencies. This practice is essential for bolstering the accuracy and comprehensiveness of data, thereby improving the trustworthiness of analyses and decision-

making. Employing techniques like exploratory analysis aids in honing the quality of data. The importance of clean data lies in its ability to predict errors and the formation of deceptive conclusions, underscoring its pivotal role in generating meaningful insights.

The obtained data will be then cleaned to be used for plotting , summary, graphs etc. Various commands such as  omit, distinct, mean etc will be used to clean the data to make it duplicate free and to remove null values and make the data clean. After the data is cleaned it is written in a new csv file using write.csv command.

House pricing cleaning data

Upon acquiring the dataset, new names were assigned to the columns. Subsequently, the data for each individual year underwent filtering to eliminate any unnecessary variables. Finally, the data from all the years were merged into a unified tibble for comprehensive analysis.

*Figure 1- House price cleaning1*



```
13
14   #----------2019 Dataset Cleaning----------#
15
16   #Cleaning data through the use of pipe operator
17   houseprices_2019 <-read_csv("C:/Users/aasis/Desktop/DataScience-Assignment/Obtained-data/housing-price/pp-2019.csv", col_names = FALSE) %>%  #Importing CSV into R
18     setNames(c("Transaction unique identifier", "Price", "Date of Transfer", "Postcode", "Property Type", "Old/New", "Duration", "PAON",
19             "SAON", "Street", "Locality", "Town/City", "District", "County", "PPD Category type", "Record Status")) %>% #Changing Column name
20     as_tibble() %>% #Converting into tibble
21     na.omit() %>% #Removing rows with null value
22     select(Price, `Date of Transfer`, Postcode, `Town/City`, District, County) %>% #selecting only columns that are required
23     filter(County =="KENT" | County== "SURREY") %>% #Preserving rows with Kent and Surrey as county
24     mutate(`Date of Transfer` = year(as.Date(`Date of Transfer`, format = "%y/%m/%d"))) %>%  #modifying the date of transfer column to only show year
25     mutate(S_No = row_number()) %>% #Adding a new serial number column
26     select(S_No, everything()) #moving the serial number column at first
27
28
29   #----------2020 Dataset Cleaning----------#
30
31   houseprices_2020 <-read_csv("C:/Users/aasis/Desktop/DataScience-Assignment/Obtained-data/housing-price/pp-2020.csv", col_names = FALSE) %>%  #Importing CSV into R
32     setNames(c("Transaction unique identifier", "Price", "Date of Transfer", "Postcode", "Property Type", "Old/New", "Duration", "PAON",
33             "SAON", "Street", "Locality", "Town/City", "District", "County", "PPD Category type", "Record Status")) %>% #Changing Column name
34     as_tibble() %>% #Converting into tibble
35     na.omit() %>% #Removing rows with null value
36     select(Price, `Date of Transfer`, Postcode, `Town/City`, District, County) %>% #selecting only columns that are required
37     filter(County =="KENT" | County== "SURREY") %>% #Preserving rows with Kent and Surrey as county
38     mutate(`Date of Transfer` = year(as.Date(`Date of Transfer`, format = "%y/%m/%d"))) %>%  #modifying the date of transfer column to only show year
39     mutate(S_No = row_number()) %>% #Adding a new serial number column
40     select(S_No, everything()) #moving the serial number column at first
41
42   #----------2021 Dataset Cleaning----------#
43
44   houseprices_2021 <-read_csv("C:/Users/aasis/Desktop/DataScience-Assignment/Obtained-data/housing-price/pp-2021.csv", col_names = FALSE) %>%  #Importing CSV into R
45     setNames(c("Transaction unique identifier", "Price", "Date of Transfer", "Postcode", "Property Type", "Old/New", "Duration", "PAON",
46             "SAON", "Street", "Locality", "Town/City", "District", "County", "PPD Category type", "Record Status")) %>% #Changing Column name
47     as_tibble() %>% #Converting into tibble
48     na.omit() %>% #Removing rows with null value
49     select(Price, `Date of Transfer`, Postcode, `Town/City`, District, County) %>% #selecting only columns that are required
50     filter(County =="KENT" | County== "SURREY") %>% #Preserving rows with Kent and Surrey as county
51     mutate(`Date of Transfer` = year(as.Date(`Date of Transfer`, format = "%y/%m/%d"))) %>%  #modifying the date of transfer column to only show year
52     mutate(S_No = row_number()) %>% #Adding a new serial number column
53     select(S_No, everything()) #moving the serial number column at first
```

*Figure 2-House price cleaning 2*

```
68
69  #merging all the cleaned dataset into a single tibble
70
71  combined_houseprices<- bind_rows(houseprices_2019, houseprices_2020, houseprices_2021, houseprices_2022) %>%
72    mutate(`Short Postcode`= substr(Postcode, 1,5)) #adding another column to the combine dataset
73
74
75  #defining path to save the cleaned dataset
76  file_path <- "C:/Users/aasis/Desktop/DataScience-Assignment/Clean-data/Cleaned House Prices.csv"
77
78
```

Population Cleaning data

After obtaining the population data it was imported into r studio. The na values as well as duplicated values were to be removed using omit and distinct commands. But the population data was already clean so it was written into a csv file.

*Figure 3-Population cleaning data*

```
1  library(tidyverse)
2  library(dplyr)
3
4  population_data=read.csv("Obtained-data/Population2011_1656567141570.csv")
5  population_data
6
7  clean_pop= population_data %>%
8    na.omit() %>%
9    distinct()
10
11  #no na data or duplicate data found in this pop csv file
12
13  write_csv(clean_pop, "C:/Users/aasis/Desktop/DataScience-Assignment/Clean-data/Cleaned Population.csv")
14
```

Broadband Speed Cleaning data

This R code begins by loading essential libraries for data manipulation and setting the working directory. It then imports a cleaned dataset linking postcodes to LSOA codes. The code proceeds to read a broadband speed dataset, selecting relevant columns, renaming them, and joining the data with the cleaned postcode-to-LSOA dataset. Additional data manipulations include selecting specific columns, handling missing values, and adding a serial number column. The cleaned broadband speed dataset is

finally saved as a CSV file. Overall, the code ensures data consistency, prepares it for analysis, and stores the cleaned dataset for further use.

*Figure 4-Broadband speed cleaning code*

```r
1   # Load necessary libraries
2   library(tidyverse)
3   library(dplyr)
4
5   |
6   # Set the working directory
7   setwd("C:/Users/aasis/Desktop/DataScience-Assignment")
8   getwd()
9   #Importing cleaned postcode to LSOA csv into R
10  cleaned_postcode_to_LSOA<- read_csv("Clean-data/Cleaned Postcode To LSOA Code.csv")
11
12  #Cleaning and joining data through the use of pipe operator
13  broadband_speed<-read_csv("Obtained-data/broadband-data/201805_fixed_pc_performance_r03.csv") %>% #Importing broadband speed csv into R
14    as_tibble() %>% #converting into tibble
15    select('Average download speed (Mbit/s)', postcode_space) %>%  #only selecting columns that are required
16    rename(Postcode= `postcode_space`) %>% #renaming the post_space column to Postcode
17    right_join(cleaned_postcode_to_LSOA, by="Postcode") %>% #Joining with the cleaned house price dataset by matching Postcode
18    select('Average download speed (Mbit/s)',Postcode, `Short Postcode`, `Town/City`, District, County,) %>% #selecting only required columns
19    na.omit() %>%  #Removing rows with null value
20    mutate(`Short Postcode`= substr(Postcode, 1,5)) %>% #Filling missing short code values
21    mutate(S_No = row_number()) %>% #Adding a new serial number column
22    select(S_No, everything()) #moving the serial number column at first
23
24
25  #defining path to save the cleaned dataset
26  file_path <- "Clean-data/Cleaned Broadband Speed Dataset.csv"
27
28
29  #saving the cleaned dataset
30  write.csv(broadband_speed,file_path, row.names = FALSE)
31
32
33
34
```

Crime rate Cleaning data

This combines crime data from Kent and Surrey, loads the datasets using the tidyverse and dplyr libraries, merges them into a single dataset (merged data), and displays the result. Subsequently, it slices values in the 'LSOA name' and 'Falls within' columns, removing excess characters. The 'Context' column is then deleted to remove any null values in this column. Finally, the cleaned dataset is saved as a CSV file named "Cleaned Crime Data.csv" at the specified path.

*Figure 5-Crime data cleaning code*

```r
1  library(tidyverse)
2  library(dplyr)
3
4  #load the datasets
5  kentcrimedata = read_csv("C:/Users/aasis/Desktop/DataScience-Assignment/Obtained-data/police-data/2023-04/2023-04-kent-street.csv")
6  surreycrimedata = read_csv("C:/Users/aasis/Desktop/DataScience-Assignment/Obtained-data/police-data/2023-04/2023-04-surrey-street.csv")
7
8  #Merge the data
9  merged_data = rbind(kentcrimedata,surreycrimedata)
10 merged_data
11 |
12 #Slicing the values in the columns LSOA and Falls
13 merged_data$`LSOA name` = substr(merged_data$`LSOA name`,1,nchar(merged_data$`LSOA name`)-5)
14 merged_data$`Falls within`= substr(merged_data$`Falls within`,1,nchar(merged_data$`Falls within`)-nchar("Police "))
15
16 #Deleting all Null Column
17 merged_data$Context = NULL
18
19 merged_data
20 write_csv(merged_data, "C:/Users/aasis/Desktop/DataScience-Assignment/Clean-data/Cleaned Crime Data.csv")
21
22
23
```

School Cleaning data

This R code processes school data from Kent and Surrey. It reads  CSV files, each representing a specific region and academic year, loads the necessary libraries (tidyverse, dplyr), and manipulates the data. The code combines the datasets into a single tibble (**combine_data**) using **rbind**. Then, it cleans the data by handling empty strings, removing rows with any missing values, filtering out rows with "NE" or "SUPP" in the 'ATT8SCR' column, converting 'ATT8SCR' to numeric, and selecting specific columns. Finally, the cleaned school data is saved as a CSV file named "Cleaned School Data.csv" in the specified directory.

*Figure 6- School data cleaning code*

```r
1   library(tidyverse)
2   library(dplyr)
3
4   setwd("C:/Users/aasis/Desktop/DataScience-Assignment")
5
6   getwd()
7
8   kent2021_22= read.csv("C:/Users/aasis/Desktop/DataScience-Assignment/Obtained-data/school-data/2021-2022 kent/886_ks4provisional.csv",fill = TRUE) %>%
9     mutate(Year = 2021) %>%
10    select(Year, PCODE, SCHNAME, ATT8SCR) %>%
11    distinct()
12  kent2022_23= read.csv("C:/Users/aasis/Desktop/DataScience-Assignment/Obtained-data/school-data/2022-2023 kent/886_ks4provisional.csv",fill=TRUE) %>%
13    mutate(Year = 2022) %>%
14    select( Year,PCODE,SCHNAME, ATT8SCR,) %>%
15    na.omit() %>%
16    distinct()
17  surrey2021_22= read.csv("C:/Users/aasis/Desktop/DataScience-Assignment/Obtained-data/school-data/2021-2022 surrey/936_ks4provisional.csv",fill=TRUE) %>%
18    mutate(Year = 2021) %>%
19    select( Year,PCODE,SCHNAME, ATT8SCR) %>%
20    na.omit() %>%
21    distinct()
22  surrey2022_23= read.csv("C:/Users/aasis/Desktop/DataScience-Assignment/Obtained-data/school-data/2022-2023 surrey/936_ks4provisional.csv",fill=TRUE) %>%
23    mutate(Year = 2022) %>%
24    select( Year,PCODE,SCHNAME, ATT8SCR,) %>%
25    na.omit() %>%
26    distinct()
27
28  combine_data = rbind(kent2021_22,kent2022_23,surrey2021_22,surrey2022_23)
29
30  cleanSchooldata=combine_data %>%
31    mutate_all(~ifelse(. == "", NA, .)) %>%   # Replace empty strings with NA
32    filter_all(all_vars(!is.na(.))) %>%    # Remove rows with any NA values
33
34  # Remove rows where ATT8SCR contains "NE" or "SUPP"
35  filter(!grepl("NE|SUPP", ATT8SCR, ignore.case = TRUE)) %>%
36  # Convert ATT8SCR to numeric (assuming it's a numeric score column)
37  mutate(ATT8SCR = as.numeric(ATT8SCR)) %>%
38  filter(!is.na(ATT8SCR)) %>%
39  select(Year, PCODE, SCHNAME, ATT8SCR) %>%
40  distinct()
41
42
43  write.csv(cleanSchooldata,"Clean-data/Cleaned School Data.csv",row.names =FALSE)
44
```

Exploratory data analysis

Exploratory Data Analysis (EDA) is telling a story using pictures and graphs. It's about digging into data to find interesting insights without jumping to conclusions. EDA makes sure that the questions we ask about the data make sense and that the answers match what we already know.
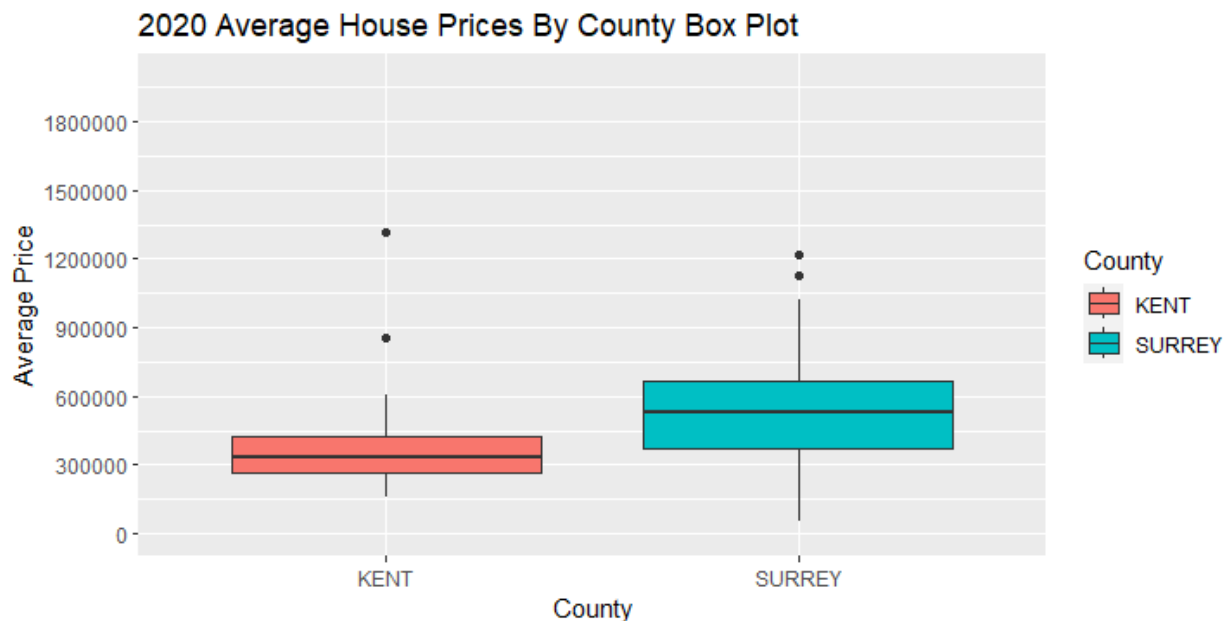
In the beginning, we decide whether to use visuals or other methods to explore the data. Then, we figure out if we want to focus on one thing or look at multiple aspects together. EDA acts like a spotlight, helping us spot unusual patterns in the data and giving us guidance on how to analyze it effectively. We start by looking at individual parts, then explore relationships between two things, and finally, we check how different factors interact with each other (Patil, What is exploratory data analysis? 2022).

While sometimes we use tables with numbers like averages, most people prefer pictures for better understanding. Different types of visuals and tools help us track and understand the data. If one way isn't clear, we try another to get a fuller picture of what the data is telling us. EDA lets us uncover hidden insights in the data, making it an exciting journey of discovery.

House price data representation

A box plot to visually compare the average house prices in Kent and Surrey for the year 2020. It starts by grouping cleaned house price data by town, district, county, and date of transfer, calculating the average price for each group. The code then filters this data to include only records from 2020. Using the ggplot library, it creates a box plot, with the x-axis representing counties (Kent and Surrey), the y-axis representing average prices, and the plot filled by county for clarity. The chart is titled "2020 Average House Prices By County Box Plot" and is customized to show prices up to 2,000,000 with breaks at intervals of 300,000 on the y-axis. The resulting visualization provides insights into the distribution of average house prices in the specified counties for the specified year.
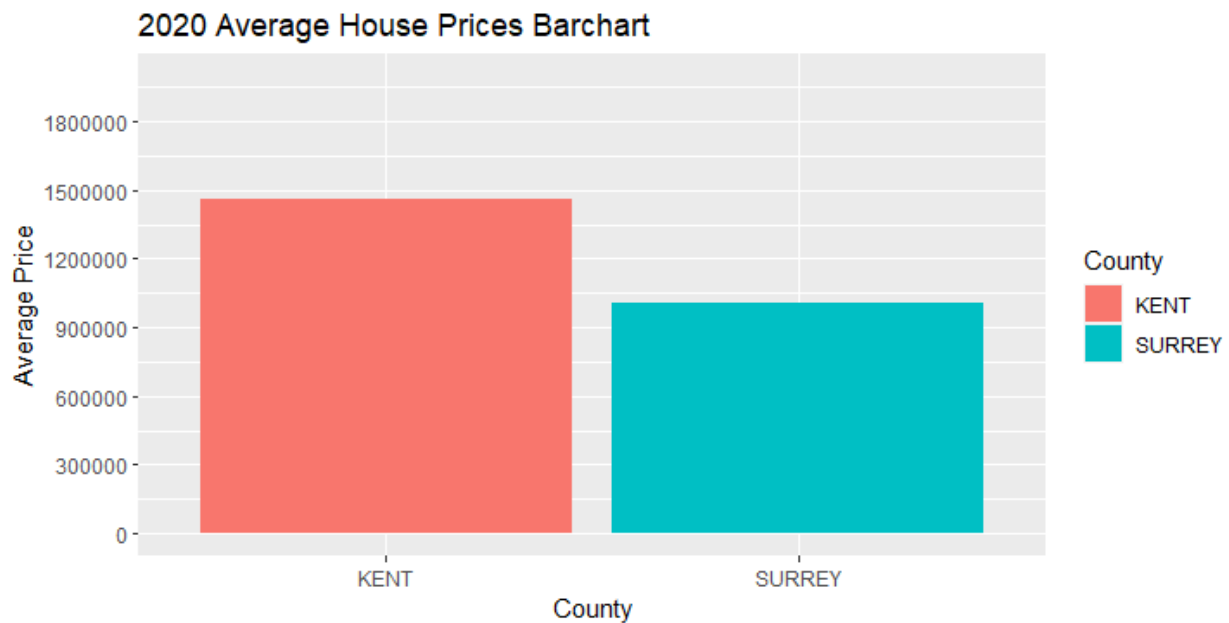
*Figure 7-Average house prices by county*



A bar chart to visually compare the average house prices in Kent and Surrey for the year 2020. It filters the grouped house price data to include only records from 2020, then uses ggplot to plot a bar chart with counties on the x-axis, average prices on the y-axis, and bars filled by county. The chart is titled
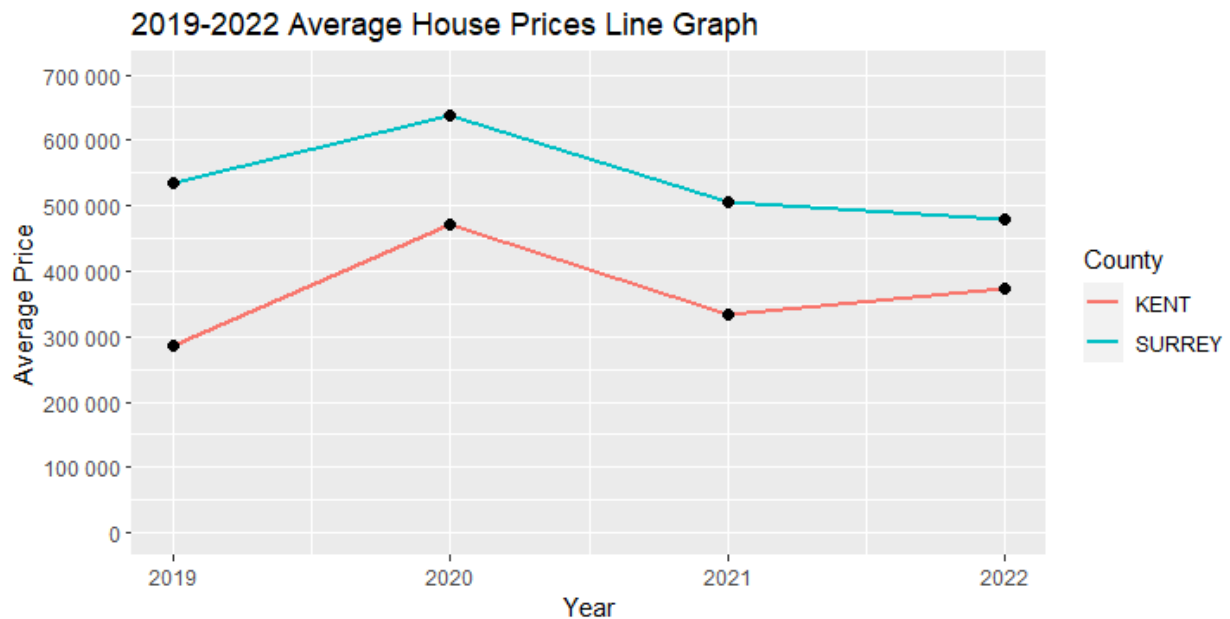
"2020 Average House Prices Barchart" and is customized to display prices up to 2,000,000 with breaks at intervals of 300,000 on the y-axis. The resulting visualization provides a clear comparison of average house prices between the two specified counties for the specified year.

*Figure 8-Average house prices barchart*



This graph groups cleaned house prices by county and year, calculating the average price for each group. It then creates a line graph to visualize average house prices from 2019 to 2022. The data is filtered to include only records from these years, and the graph compares county prices over this period. The resulting line graph is titled "2019-2022 Average House Prices Line Graph," with the x-axis representing years, the y-axis representing average prices, and different colors indicating different counties. Points on the lines highlight specific data points, and the graph is customized with specific limits, breaks, and labels for better clarity.

*Figure 9-Average house prices line graph*



2019-2022 Average House Prices Line Graph

Broadband Speed data representation

Ggplot was used to plot various graphs of broadband speed of both kent in surrey. At first average download speed of both county was compared in a box plot where surrey had more speed. Following that average download speed within kent was put in a bar chart. Also pie chart of robbery crimes committed is also plotted.

*Figure 10-Average Download speed by COunty*



Average Download Speed By County Box Plot

*Figure 11-Average download speed within Kent*



Average Download Speed Within Kent Bar Chart

*Figure 12-Average Download Speed within surrey*



Average Download Speed Within Surrey Bar Chart

- Crime rate data representation

Ggplot and fmsb library was used to plot various graphs of crime rate of both Kent and surrey. At first Drug offence rate by district is plotted in a box plot. Following that a radar chart of vehicle offence rate per 10k is plotted. Also average download speed of surrey was visualized in a bar chart along with appropriate legend and colors. At last line chart of drug offence rate of both kent and surrey is plotted.

*Figure 13- Drug offence rate by district*



2023 Drug Offence Rate By District Box Plot

*Figure 14-Radar chart*

*Figure 15-Robbery rate in march 2022*

**Robbery Crime Rate by District in March 2022**



District
- CANTERBURY
- MAIDSTONE
- REIGATE AND BANSTEAD
- SWALE
- TANDRIDGE
- THANET
- TUNBRIDGE WELLS

*Figure 16 Drug offence rate*



2020-2023 Drug Offence Rate

School data representation

In this section average attainment score is compared between both kent and surrey in a box plot. Then average attainment score line graph is drawn of kent county following that it is also drawn for surrey county.

*Figure 17-Average attainment score both county*



2021 Average Attainment Score By County Box

*Figure 18-Average attainment score kent line graph*



2020-2023 Average Attainment Score Line Graph For Kent's District

*Figure 19-Average attainment score surrey line graph*



Linear modeling

 Leveraging the principles of linear regression, a statistical technique employed to identify

connections between variables, streamlines the analysis of diverse datasets without repetitive efforts. In

this approach, data points are linked by a straight line, and the objective is to determine regression

coefficients that minimize errors, ensuring the best possible fit. Linear regression comes in two main

forms: simple linear regression, dealing with a single variable, and multiple linear regression, which

involves several independent variables and is more intricate (Mishra, Linear Modeling 2021).

 The analysis encompasses a broad spectrum of information, including housing prices, crime rates

(such as drug-related, robbery, and vehicle-related crimes), and significant educational metrics. By

amalgamating these diverse datasets, a comprehensive comparison is facilitated, aiding in the identification of countries that have performed exceptionally well.

House price vs Average Download Speed

This begins by importing cleaned house prices and broadband speed datasets. It then groups the house prices and broadband speed data by town and county, calculating the average price and download speed for each group. The two datasets are joined into a single table. The code proceeds to create a linear model predicting house prices based on average download speed. The summary of the linear model is displayed, providing insights into the relationship between house prices and broadband speed. Additionally, a graphical representation of the linear model is generated using ggplot, with data points colored differently for Kent and Surrey. The resulting plot visualizes the 2020 house prices against average download speed, including a linear regression line.

*Figure 20-House price vs Average download speed*



House price vs Drug rate

     This R code begins by grouping cleaned house prices for the year 2021 by town and county, calculating the average price for each group. It then modifies the crime dataset to focus on drug-related offenses, creating a new dataset showing drug offense rates for each town and county in 2021. The two datasets are joined into a single table, and a linear model is created to predict house prices based on drug offense rates. The summary of the linear model is displayed, offering insights into the relationship between house prices and drug offenses. A graphical representation of the linear model is generated using ggplot, with data points colored differently for Kent and Surrey. The resulting plot visualizes 2021 house prices against drug offense rates, including a linear regression line. Finally, null values are removed from the joined dataset.

*Figure 21-House price vs Drug offence rate*



Attainment 8 Score vs House price

First group cleaned house prices for the year 2021 by town and county, calculating the average price for each group. It then groups school data by town and county for the same year, calculating the average attainment score for each group. The two datasets are joined into a single table based on the town, converting town names to lowercase for consistency. A linear model is created to predict average attainment scores based on average house prices. The summary of the linear model is displayed, providing insights into the relationship between attainment scores and house prices. A graphical representation of the linear model is generated using ggplot, with data points colored differently for Kent and Surrey. The resulting plot visualizes the 2021 attainment scores against house prices, including a linear regression line. Null values are removed from the joined dataset.

*Figure 22-Attainment score vs House prices*



2021 Attainment Score vs House Prices

Attainment 8 score vs drug rate

This begins by grouping school data for the year 2021 by town and county, calculating the average attainment score for each group. It then modifies the crime dataset to focus on drug-related offenses, creating a new dataset showing drug offense rates for each town and county in 2021. The two datasets are joined into a single table based on the town, converting town names to lowercase for consistency. A linear model is created to predict average attainment scores based on drug offense rates. The summary of the linear model is displayed, providing insights into the relationship between attainment scores and drug offenses. A graphical representation of the linear model is generated using ggplot, with data points colored differently for Kent and Surrey. The resulting plot visualizes the 2021 attainment scores against drug offense rates, including a linear regression line.

*Figure 23-Attainment Score vs Drug offence rate*



2021 Attainment Score vs Drug Offence Rate

Average Download Speed vs Drug rate

This begins by grouping broadband speed data by town and county, calculating the average download speed for each group. It then modifies the crime dataset to focus on drug-related offenses, creating a new dataset showing drug offense rates for each town and county in 2022. The two datasets are joined into a single table based on the town. A linear model is created to predict average download speeds based on drug offense rates. The summary of the linear model is displayed, providing insights into the relationship between download speeds and drug offenses. A graphical representation of the linear model is generated using ggplot, with data points colored differently for Kent and Surrey. The resulting plot visualizes the 2022 average download speeds against drug offense rates, including a linear regression line.

*Figure 24-Average download speed vs Drug offence Rate*



2022 Average Download Speed vs Drug Offence Rate

Average Download Speed vs Attainment Score

This begins by grouping broadband speed data by town and county, converting town names to lowercase for consistency, and calculating the average download speed for each group. It then groups school data by town and county, converting town names to lowercase, and finding the average attainment score for each group. The two datasets are joined into a single table based on the town. A linear model is created to predict average download speeds based on attainment scores. The summary of the linear model is displayed, providing insights into the relationship between download speeds and attainment scores. A graphical representation of the linear model is generated using ggplot, with data points colored differently for Kent and Surrey. The resulting plot visualizes the average download speeds against attainment scores, including a linear regression line.

*Figure 25-Average Download speed vs Attainment Score*



Average Download Speed vs Attainment Score

## Recommendation System

House ranking

Margate came in first place from kent county with an average price of 125226. Also higher score is given to more affordable house.

*Figure 26-House ranking*

```
 7  #-----------House Price Ranking-----------#
 8
 9  # Clean the Price column (convert to numeric)
10  house_prices_data$Price <- as.numeric(gsub("[^0-9.]", "", house_prices_data$Price))
11
12  # Group by County and Town/City, calculate the average price for each group
13  grouped_data <- house_prices_data %>%
14    group_by(County, `Town.City`) %>%
15    summarise(`Average Price` = mean(Price))
16
17  # Assign a score to housing prices (higher score indicates affordability)
18  grouped_data <- grouped_data %>%
19    mutate(Score = 1 - scale(`Average Price`))
20
21  # Arrange the data in descending order based on the score
22  sorted_data <- grouped_data[order(-grouped_data$Score), ]
23
24  # Select the top 10 entries
25  top_10 <- head(sorted_data, 10)
26
27  # Print the top 10 best counties and town/cities with the least house prices and their scores
28  print(top_10[, c("County", "Town.City", "Average Price", "Score")])
29
```

```
29:1   (Top Level) ÷                                                                    R S
```

```
  R 4.3.1 · C:/Users/aasis/Desktop/DataScience-Assignment/
# A tibble: 10 x 4
# Groups:   County [2]
   County Town.City      `Average Price` Score[,1]
   <chr>  <chr>                    <dbl>     <dbl>
 1 KENT   MARGATE                145226.      2.93
 2 KENT   SHEERNESS              163832.      2.71
 3 KENT   NEW ROMNEY             168584.      2.66
 4 KENT   WHITSTABLE             177834.      2.55
 5 KENT   GRAVESEND              220276.      2.06
 6 KENT   DEAL                   223601.      2.02
 7 KENT   BIRCHINGTON            229328.      1.95
 8 KENT   ROMNEY MARSH           233316.      1.90
 9 SURREY EDENBRIDGE             165000.      1.83
10 KENT   WEST MALLING           247423.      1.74
```

## Broadband Ranking

The town of Chatham from kent county came in first with the average download speed of 106 .

*Figure 27-Broadband speed ranking*

```
30 ▾ #--------------------------------------------------------------------------------
31   broadband_speed_data <- read.csv("C:/Users/aasis/Desktop/DataScience-Assignment/Clean-data/Cleaned Broadband Speed Dat
32   broadband_speed_data
33   # Group by County and Town/City, calculate the average download speed for each group
34   grouped_data <- broadband_speed_data %>%
35     group_by(County, `Town.City`) %>%
36     summarise(`Average Download Speed` = mean(`Average.download.speed..Mbit.s.`))
37
38   # Assign a score to download speeds (higher score indicates better speed)
39   grouped_data <- grouped_data %>%
40     mutate(Score = scale(`Average Download Speed`))
41
42   # Arrange the data in descending order based on the score
43   sorted_data <- grouped_data[order(-grouped_data$Score), ]
44
45   # Select the top 10 entries
46   top_10 <- head(sorted_data, 10)
47
48   # Print the top 10 towns/cities and counties with the highest average download speed and their scores
49   print(top_10[, c("County", "Town.City", "Average Download Speed", "Score")])
50
```

```
30:39    # (Untitled) ≎                                                                          R Scri

  R 4.3.1 · C:/Users/aasis/Desktop/DataScience-Assignment/ ➔
 A tibble: 10 x 4
# Groups:   County [2]
   County Town.City        `Average Download Speed`   Score[,1]
   <chr>  <chr>                              <dbl>      <dbl>
 1 KENT   CHATHAM                            106.       3.79
 2 KENT   SWANSCOMBE                         106.       3.79
 3 SURREY SURBITON                           87.2       2.85
 4 SURREY WEYBRIDGE                          69.2       1.80
 5 SURREY EPSOM                              64.0       1.50
 6 SURREY ALDERSHOT                          63.4       1.46
 7 SURREY ADDLESTONE                         62.1       1.39
 8 SURREY EGHAM                              55.6       1.01
 9 SURREY CAMBERLEY                          53.7       0.894
10 SURREY WALTON-ON-THAMES                   52.9       0.850
```

## Crime ranking

Orpington town in surrey county came in 1st place with the least amount of crimes count.

```
51  #--------------------------------------------------------------------------------
52  crime_data <- read.csv("C:/Users/aasis/Desktop/DataScience-Assignment/Clean-data/Cleaned Crime Dataset.csv")
53
54  grouped_data <- crime_data %>%
55    group_by(Falls.within, `Town.City`) %>%
56    summarise(`Total Crime Count` = n())
57
58  # Assign a score to total crime counts (higher score indicates lower crime)
59  grouped_data <- grouped_data %>%
60    mutate(Score = rank(`Total Crime Count`))
61
62  # Arrange the data in ascending order based on the score
63  sorted_data <- grouped_data[order(grouped_data$Score), ]
64
65  # Select the top 10 entries
66  top_10 <- head(sorted_data, 10)
67
68  # Print the top 10 towns/cities and counties with the lowest total crime count and their scores
69  print(top_10[, c("Falls.within", "Town.City", "Total Crime Count", "Score")])
70
```

```
70:1    # (Untitled)                                                                                      R Script
```

```
R 4.3.1 · C:/Users/aasis/Desktop/DataScience-Assignment/
# A tibble: 10 x 4
# Groups:   Falls.within [2]
   Falls.within  Town.City   `Total Crime Count` Score
   <chr>         <chr>                     <int> <dbl>
 1 Kent Police   CATERHAM                      2     1
 2 Surrey Police ORPINGTON                     1     1
 3 Kent Police   OXTED                         7     2
 4 Surrey Police LONGFIELD                     3     2
 5 Kent Police   CRANLEIGH                    19     3
 6 Surrey Police DARTFORD                      7     3
 7 Kent Police   LINGFIELD                    21     4
 8 Surrey Police SWANLEY                      22     4
 9 Kent Police   HERNE BAY                   293     5
10 Surrey Police EDENBRIDGE                  344     5
>
```

School ranking

Dartford town from kent came in 1st place with highest average attainment score.

```
71
72  # Read the CSV file
73  school_data <- read.csv("C:/Users/aasis/Desktop/DataScience-Assignment/Clean-data/Cleaned School Dataset.csv")
74  school_data
75  # Group by County and Town/City, calculate the average attainment score for each group
76  grouped_data <- school_data %>%
77    group_by(County, Town) %>%
78    summarise(`Average Attainment Score` = mean(`Attainment.Score`))
79
80  # Assign a score to average attainment scores (higher score indicates higher attainment)
81  grouped_data <- grouped_data %>%
82    mutate(Score = rank(`Average Attainment Score`, na.last = "keep"))
83
84  # Arrange the data in descending order based on the score
85  sorted_data <- grouped_data[order(-grouped_data$Score), ]
86
87  # Select the top 10 entries
88  top_10 <- head(sorted_data, 10)
89
90  # Print the top 10 towns/cities and counties with the highest average attainment scores and their scores
91  print(top_10[, c("County", "Town", "Average Attainment Score", "Score")])
92
```

`92:1`  `(Untitled)`   `R Sc`

R 4.3.1 · C:/Users/aasis/Desktop/DataScience-Assignment/

```
# A tibble: 10 x 4
# Groups:   County [2]
   County Town       `Average Attainment Score` Score
   <chr>  <chr>                           <dbl> <dbl>
 1 Kent   Dartford                         59.2    39
 2 Kent   Sandwich                         52.0    38
 3 Kent   Faversham                        49.9    37
 4 Kent   Tonbridge                        49.7    36
 5 Kent   Wye                              49.2    35
 6 Kent   Rochester                        49      34
 7 Kent   Maidstone                        44.9    33
 8 Surrey Weybridge                        61.8    33
 9 Kent   Gravesend                        44.7    32
10 Surrey Esher                            57.1    32
>
```

Overall ranking

```
[1] Recommended city:
> print(recommended_city)
      Town/City County HouseScore Attainment.Score SchoolScore AvgDownloadSpeed BroadbandScore TotalCrimes CrimeScore TotalScore
133 WARLINGHAM SURREY       9.45                0           0             12.1       9.886171        1221   9.997334   29.33351
>
```

**Legal and ethical issues**

In the modern era, the utilization of data has resulted in a plethora of legal, ethical, and social dilemmas. One of the most crucial concerns is the protection of data. The collection, storage, and sharing of personal information can lead to infringement of individual rights. To safeguard privacy, it is imperative to implement measures against unauthorized access and misuse.

Ethical considerations are also of utmost importance. Striking a balance between utilizing data for gaining insights and respecting individuals' autonomy is a challenging task. When data is used to alter people's behavior, it exacerbates the ethical dilemma. This brings to light issues of transparency and

potential manipulation.

(*5 Principles of Data Ethics for Business*, 2021)

The aspect of safety is another crucial aspect in this scenario. Data breaches can result in theft of personal information, financial harm, and damage to reputation. Ensuring data security is not only a legal requirement, but also a moral obligation to prevent harm to individuals.

## Conclusion

In brief, the study aimed to assist a friend in selecting a suitable location in the UK for relocation. Key factors such as housing costs, broadband speed, school quality, and crime rates were examined to gather crucial information. The data underwent meticulous cleaning and was visually presented through graphs and charts for clarity. A specialized model was employed to explore correlations among various factors, leading to the creation of a ranking based on the analysis results.

## References

Elgabry, O. (2019, March 2). *The Ultimate Guide to data cleaning*. Medium.

https://towardsdatascience.com/the-ultimate-guide-to-data-cleaning-3969843991d4

*R tutorial: Learn R programming language tutorial - javatpoint*. www.javatpoint.com.

Police Data. (2022). Police Data Download. https://data.police.uk/data/

Registry, H. L. (2023, July 28). *Price paid data*. GOV.UK. https://www.gov.uk/government/statistical-data-sets/price-paid-data-downloads

*5 Principles of Data Ethics for Business*. (2021, March 16). Business Insights Blog.

https://online.hbs.edu/blog/post/data-ethics

Mishra, P. (2021, December 19). *Linear Modeling*. GeeksforGeeks.

https://www.geeksforgeeks.org/how-to-use-lm-function-in-r-to-fit-linear-models/

Ofcom. (2023, March 16). *Connected nations 2018: Data downloads*. Ofcom.

https://www.ofcom.org.uk/research-and-data/multi-sector-research/infrastructure-

research/connected-nations-2018/data-downloads

**Github Link**

**https://github.com/aasiskrk/Datasci-assign**

**Appendix**

```
11
12  #importing the cleaned school dataset
13  cleaned_school_dataset= read_csv('Clean-data/Cleaned School Dataset.csv')
14
15  #grouping broadband speed by town and county and finding average download speed for each group
16  grouped_broadband_speeds = cleaned_broadband_speed %>%
17    group_by(`Town/City`,County) %>%
18    mutate(`Town/City`= tolower(`Town/City`)) %>%   #converting the town from to all lowercase
19    summarise(`Average download speed (Mbit/s)`= mean(`Average download speed (Mbit/s)`))
20
21  #grouping school data by town and county and finding average score for each group
22  grouped_school_dataset = cleaned_school_dataset %>%
23    group_by(`Town`,County) %>%
24    mutate(Town= tolower(Town)) %>%   #converting the town from to all lowercase
25    summarise(`Attainment Score`=mean(`Attainment Score`))
26
27
28  #joining broadband data and school data in a single table
29  broadband_attainment_data = grouped_broadband_speeds %>%
30    left_join(grouped_school_dataset,by=c("Town/City"="Town")) %>%
31    na.omit #removing rows with null value
32
33  #creating a linear model
34  l_model = lm(data=broadband_attainment_data, `Average download speed (Mbit/s)`~`Attainment Score`) #t
35
36  #showing summary of the Linear Model
37  summary(l_model)
38
39
40  #creating the linear model graph
41  ggplot(broadband_attainment_data,aes(x=`Attainment Score`,y=`Average download speed (Mbit/s)`)) +
42    scale_y_continuous(limits=c(0,80), breaks = seq(0,80,5))+ #setting limits and breaks
43    geom_point(data = filter(broadband_attainment_data,County.x=="KENT"),aes(color=c("Red"="Kent")))+ #
44    geom_point(data = filter(broadband_attainment_data,County.x=="SURREY"), aes(color=c("Blue"="Surrey"
45    geom_smooth(method=lm,se=FALSE,color="lightgreen")+ #adding linear regression line and omitting err
46    labs(x="Attainment Score",
47        y="Average Download Speed (Mbit/s)",
48        title="Average Download Speed vs Attainment Score",color="County") #setting labels
49
```

```r
22
23  #modifying our crime dataset to show drug offence rate and crime count
24  crime_dataset_drugs2 <-cleaned_crime_dataset %>%
25    mutate(`Date of crime`= substr(`Date of crime`, 1, 4)) %>% #Mutating this column to only show year
26    group_by(`Short Postcode`,`Crime type`,`Date of crime`, `Falls within`) %>% #Grouping to show crime count in each postcode by year
27    select(`Short Postcode`,`Crime type`,`Date of crime`, `Falls within`) %>%
28    na.omit() %>%
29    tally() %>% #creating crime count column
30    rename(`Crime Count`=n) %>%  #renaming crime count column %>%
31    right_join(population_dataset, by = "Short Postcode") %>% #joining with population dataset to show district and population
32    select(`Short Postcode`,`Crime type`,`Crime Count`, `Population`, `Date of crime`, `Falls within`, `Town/City`, District) %>% #select the required co
33    na.omit() %>%
34    filter(`Crime type`== "Drugs") %>% #filtering to show only drug crimes of 2022
35    mutate(`Drug Offence Rate` = (`Crime Count` / Population)) #calculating drug offence rate
36
37  #grouping the drug crime dataset by county and town and showing the rate for each group for the year 2022
38  grouped_drug_crime <- crime_dataset_drugs2 %>%
39    filter(`Date of crime`=="2022") %>%
40    group_by(`Falls within`,`Town/City`) %>%
41    summarise(`Drug Offence Rate`= mean(`Drug Offence Rate`))
42
43
44  #joining broadband data and drug crime rate data in a single table
45  broadband_crime_data = grouped_broadband_speeds %>%
46    left_join(grouped_drug_crime,by="Town/City") %>%
47    na.omit #removing null values
48
49
50  #creating a linear model
51  l_model = lm(data=broadband_crime_data, `Average download speed (Mbit/s)`~`Drug Offence Rate`) #this model predicts Average download speed as a functio
52
53  #showing summary of the Linear Model
54  summary(l_model)
55
56
57  #creating the linear model graph
58  ggplot(broadband_crime_data,aes(x=`Drug Offence Rate`,y=`Average download speed (Mbit/s)`)) +
59    scale_y_continuous(limits=c(0,50), breaks = seq(0,50,5))+ #setting limits and breaks
60    geom_point(data = filter(broadband_crime_data,County=="KENT"),aes(color=c("Red"="Kent")))+ #setting color as red for Kent's data point
61    geom_point(data = filter(broadband_crime_data,County=="SURREY"), aes(color=c("Blue"="Surrey"))) + #setting color as blue for Surrey's data point
62    geom_smooth(method=lm,se=FALSE,color="lightgreen")+ #adding linear regression line and omitting error bands
63    labs(x="Drug Offence Rate",
64        y="Average Download Speed (Mbit/s)",
65        title="2022 Average Download Speed vs Drug Offence Rate",color="County") #setting labels
66
```

```r
20  #grouping school data by town and county and finding average score for each group
21  grouped_school_dataset = cleaned_school_dataset %>%
22    filter(`Year`=="2021") %>%
23    group_by(`Town`,County) %>%
24    mutate(Town= tolower(Town)) %>%  #converting the town from to all lowercase
25    summarise(`Attainment Score`=mean(`Attainment Score`))
26
27  #modifying our crime dataset to show drug offence rate and crime count
28  crime_dataset_drugs2 <-cleaned_crime_dataset %>%
29    mutate(`Date of crime`= substr(`Date of crime`, 1, 4)) %>% #Mutating this column to only show year
30    group_by(`Short Postcode`,`Crime type`,`Date of crime`, `Falls within`) %>% #Grouping to show crime count in each postcode by year
31    select(`Short Postcode`,`Crime type`,`Date of crime`, `Falls within`) %>%
32    na.omit() %>%
33    tally() %>% #creating crime count column
34    rename(`Crime Count`=n) %>%  #renaming crime count column %>%
35    right_join(population_dataset, by = "Short Postcode") %>% #joining with population dataset to show district and population
36    select(`Short Postcode`,`Crime type`,`Crime Count`, `Population`, `Date of crime`, `Falls within`, `Town/City`, District) %>% #select the required co
37    na.omit() %>%
38    filter(`Crime type`== "Drugs") %>% #filtering to show only drug crimes of 2022
39    mutate(`Drug Offence Rate` = (`Crime Count` / Population)) #calculating drug offence rate
40
41  #grouping the drug crime dataset by county and town and showing the rate for each group for the year 2021
42  grouped_drug_crime <- crime_dataset_drugs2 %>%
43    filter(`Date of crime`=="2021") %>%
44    group_by(`Falls within`,`Town/City`) %>%
45    mutate(`Town/City`= tolower(`Town/City`)) %>%  #converting the town from to all lowercase
46    summarise(`Drug Offence Rate`= mean(`Drug Offence Rate`))
47
48  #joining school data and house price data in a single table
49  school_drug_data = grouped_school_dataset %>%
50    left_join(grouped_drug_crime ,by=c("Town"="Town/City")) %>%
51    na.omit #removing rows with null value
52
53  #creating a linear model
54  l_model = lm(data=school_drug_data, `Attainment Score`~`Drug Offence Rate`) #this model predicts Average attainment score as a function of Drug offence
55
56  #showing summary of the Linear Model
57  summary(l_model)
58
59
60  #creating the linear model graph
61  ggplot(school_drug_data,aes(x=`Drug Offence Rate`,y= `Attainment Score`)) +
62    scale_y_continuous(limits=c(0,50), breaks = seq(0,50,5))+ #setting limits and breaks
63    geom_point(data = filter(school_drug_data,County=="Kent"),aes(color=c("Red"="Kent")))+ #setting color as red for Kent's data point
64    geom_point(data = filter(school_drug_data,County=="Surrey"), aes(color=c("Blue"="Surrey"))) + #setting color as blue for Surrey's data point
65    geom_smooth(method=lm,se=FALSE,color="lightgreen")+ #adding linear regression line and omitting error bands
66    labs(x="Drug Offence Rate",
67        y="Attainment Score",
68
```

```r
14
15  #grouping house prices by town and county and finding average price for each group
16  grouped_house_prices = cleaned_houseprices %>%
17    filter(`Date of Transfer`=="2021") %>%
18    group_by(`Town/City`,County) %>%
19    mutate(`Town/City` = tolower(`Town/City`)) %>% #converting the town from uppercase to all lowercase
20    summarise(Price=mean(Price))
21
22
23  #grouping school data by town and county and finding average score for each group
24  grouped_school_dataset = cleaned_school_dataset %>%
25    filter(`Year`=="2021") %>%
26    group_by(`Town`,County) %>%
27    mutate(Town= tolower(Town)) %>%  #converting the town from to all lowercase
28    summarise(`Attainment Score`=mean(`Attainment Score`))
29
30
31  #joining school data and house price data in a single table
32  school_houseprice_data = grouped_school_dataset %>%
33    left_join(grouped_house_prices,by=c("Town"="Town/City")) %>%
34    na.omit #removing rows with null value
35
36  #creating a linear model
37  l_model = lm(data=school_houseprice_data, `Attainment Score`~Price) #this model predicts Average attainment score as a function of Average house prices
38
39  #showing summary of the Linear Model
40  summary(l_model)
41
42  #creating the linear model graph
43  ggplot(school_houseprice_data,aes(x=Price,y= `Attainment Score`)) +
44    scale_y_continuous(limits=c(0,80), breaks = seq(0,80,5))+ #setting limits and breaks
45    geom_point(data = filter(school_houseprice_data,County.x=="Kent"),aes(color=c("Red"="Kent")))+ #setting color as red for Kent's data point
46    geom_point(data = filter(school_houseprice_data,County.x=="Surrey"), aes(color=c("Blue"="Surrey"))) + #setting color as blue for Surrey's data point
47    geom_smooth(method=lm,se=FALSE,color="lightgreen")+ #adding linear regression line and omitting error bands
48    labs(x="House Price",
49         y="Attainment Score",
50         title="2021 Attainment Score vs House Prices",color="County") #setting labels
51
```

```r
24  #modifying our crime dataset to show drug offence rate and crime count
25  crime_dataset_drugs2 <-cleaned_crime_dataset %>%
26    mutate(`Date of crime`= substr(`Date of crime`, 1, 4)) %>% #Mutating this column to only show year
27    group_by(`Short Postcode`,`Crime type`,`Date of crime`, `Falls within`) %>% #Grouping to show crime count in each postcode by year
28    select(`Short Postcode`,`Crime type`,`Date of crime`, `Falls within`) %>%
29    na.omit() %>%
30    tally() %>% #creating crime count column
31    rename(`Crime Count`=n) %>%  #renaming crime count column %>%
32    right_join(population_dataset, by = "Short Postcode") %>% #joining with population dataset to show district and population
33    select(`Short Postcode`,`Crime type`,`Crime Count`, `Population`, `Date of crime`, `Falls within`, `Town/City`, District) %>% #select the required co
34    na.omit() %>%
35    filter(`Crime type`== "Drugs") %>% #filtering to show only drug crimes of 2022
36    mutate(`Drug Offence Rate` = (`Crime Count` / Population)) #calculating drug offence rate
37
38  #grouping the drug crime dataset by county and town and showing the rate for each group for the year 2020
39  grouped_drug_crime <- crime_dataset_drugs2 %>%
40    filter(`Date of crime`=="2021") %>%
41    group_by(`Falls within`,`Town/City`) %>%
42    summarise(`Drug Offence Rate`= mean(`Drug Offence Rate`))
43
44
45  #joining house price data and drug crime rate data in a single table
46  house_price_drug_crime_data = grouped_house_prices %>%
47    left_join(grouped_drug_crime,by="Town/City") %>%
48    na.omit #removing null values
49  http://127.0.0.1:26089/graphics/5a4605a6-f469-49de-992e-e087dc02f32a.png
50
51  #creating a linear model
52  l_model = lm(data=house_price_drug_crime_data, Price~`Drug Offence Rate`) #this model predicts House Price as a function of Drug offence rate
53
54  #showing summary of the Linear Model
55  summary(l_model)
56
57  #creating the linear model graph
58  ggplot(house_price_drug_crime_data,aes(x=`Drug Offence Rate`,y=Price)) +
59    scale_y_continuous(limits=c(0,1000000), breaks = seq(0,1000000,200000))+ #setting limits and breaks
60    geom_point(data = filter(house_price_drug_crime_data,County=="KENT"),aes(color=c("Red"="Kent")))+ #setting color as red for Kent's data point
61    geom_point(data = filter(house_price_drug_crime_data,County=="SURREY"), aes(color=c("Blue"="Surrey"))) + #setting color as blue for Surrey's data poi
62    geom_smooth(method=lm,se=FALSE,color="lightgreen")+ #adding linear regression line and omitting error bands
63    labs(x="Drug Offence Rate",
64         y="Price",
65         title="2021 House Prices vs Drug Offence Rate",color="County") #setting labels
66
67
```