

CA2: NI Postcode and Crime data

Section-1:

Introduction:

The NI post code and crime data are given to work on the CA2. NI post code datasets consists of the postcode details of the Northern Ireland like Organisation name, Building name, number, Primary thorfare, locality, town etc. Crime data consist of the information on crimes of the Northern Ireland. Assignment has some task to do the for the analysis in order to obtain the final output.

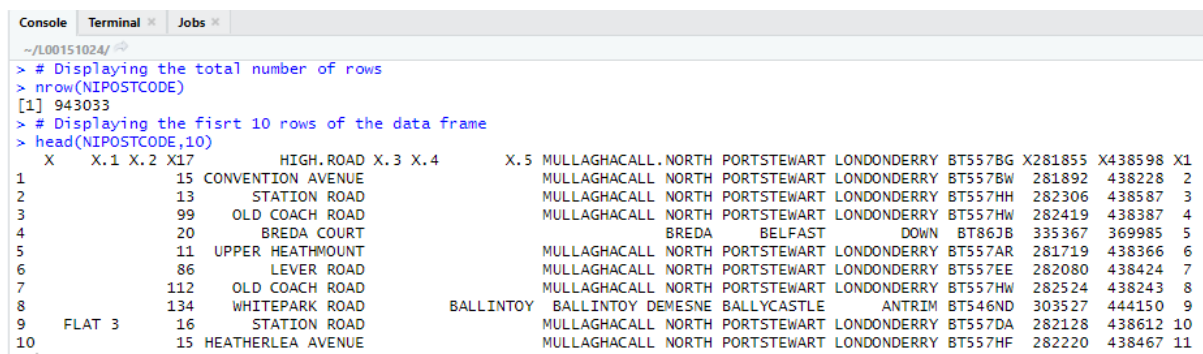
Section 1:

- A) The first task of the assignment is to show the total number of rows, the structure of the data and the first 10 rows of the data frame containing all of the NI postcode data. The code for the above task is shown in the fig.1.

```
# Loading the CSV file by placing in the working directory
NIPOSTCODE = read.csv("NIPostcodes.csv", header = FALSE)
# Displaying the total number of rows
nrow(NIPOSTCODE)
# Structure of the data
str(NIPOSTCODE)
# Displaying the first 10 rows of the data frame
head(NIPOSTCODE,10)
```

Fig.1: Loading the CSV file, showing total number of rows and first 10 rows of the data frame

The CSV file is placed in the working directory to access publicly. The dataset is loaded by using the read.csv function and number of rows are displayed using the function nrow. Head function gives you limited number of rows in this task we need only 10. The structure and outputs of the task are shown in the below figure. The Header = false in the first line of the code is used because in the figure.2 we can observe that the first column is treated as the header to ignore this header must be false.



	X	X.1	X.2	X17	HIGH. ROAD	X.3	X.4	X.5	MULLAGHACALL	NORTH	PORTSTEWART	LONDONDERRY	BT557BG	X281855	X438598	X1
1				15	CONVENTION AVENUE				MULLAGHACALL	NORTH	PORTSTEWART	LONDONDERRY	BT557BW	281892	438228	2
2				13	STATION ROAD				MULLAGHACALL	NORTH	PORTSTEWART	LONDONDERRY	BT557HH	282306	438587	3
3				99	OLD COACH ROAD				MULLAGHACALL	NORTH	PORTSTEWART	LONDONDERRY	BT557HW	282419	438387	4
4				20	BREDA COURT				BREDA	BELFAST	DOWN	BT86JB	335367	369985		5
5				11	UPPER HEATHMOUNT				MULLAGHACALL	NORTH	PORTSTEWART	LONDONDERRY	BT557AR	281719	438366	6
6				86	LEVER ROAD				MULLAGHACALL	NORTH	PORTSTEWART	LONDONDERRY	BT557EE	282080	438424	7
7				112	OLD COACH ROAD				MULLAGHACALL	NORTH	PORTSTEWART	LONDONDERRY	BT557HW	282524	438243	8
8				134	WHITEPARK ROAD			BALLINTOY	BALLINTOY	DEMESNE	BALLYCASTLE	ANTRIM	BT546ND	303527	444150	9
9	FLAT 3			16	STATION ROAD				MULLAGHACALL	NORTH	PORTSTEWART	LONDONDERRY	BT557DA	282128	438612	10
10				15	HEATHERLEA AVENUE				MULLAGHACALL	NORTH	PORTSTEWART	LONDONDERRY	BT557HF	282220	438467	11

Fig.3: Error in the output

```

> str(NIPOSTCODE)
'data.frame':   943033 obs. of  15 variables:
 $ X      : Factor w/ 40859 levels: "", "ASCERT", "BALLYMAC HOTEL",...: 1 1 1 1 1 1 1 1 1 1 ...
 $ X.1    : Factor w/ 6186 levels: "", "RETURN APARTMENT" B",...: 1 1 1 1 1 1 1 1 1 1 3166 1 ...
 $ X.2    : Factor w/ 12215 levels: "", "ARDEEVIN",...: 1 1 1 1 1 1 1 1 1 1 ...
 $ X17    : Factor w/ 5831 levels: "", "25C",...: 804 564 5076 1368 298 4701 348 612 926 804 ...
 $ HIGH.ROAD : Factor w/ 24540 levels: "", "ABBACY ROAD",...: 6072 21646 18024 3392 23334 14504 18024 23992 21646 11945 ...
 $ X.3    : Factor w/ 453 levels: "", "AN BEALACH LEATHAN",...: 1 1 1 1 1 1 1 1 1 1 ...
 $ X.4    : Factor w/ 288 levels: "", "ABBEY ROAD",...: 1 1 1 1 1 1 1 1 1 1 ...
 $ X.5    : Factor w/ 676 levels: "", "ABBEY BUSINESS PARK",...: 1 1 1 1 1 1 1 1 55 1 1 ...
 $ MULLAGHACALL.NORTH: Factor w/ 7705 levels: "ABBEY PARK", "ABOCURRAGH",...: 6145 6145 6145 1532 6145 6145 482 6145 6145 ...
 $ PORTSTEWART : Factor w/ 314 levels: "", "AGHAGALLON",...: 271 271 271 46 271 271 271 27 271 271 ...
 $ LONDONDERRY : Factor w/ 6 levels: "ANTRIM", "ARMAGH",...: 5 5 5 3 5 5 5 1 5 5 ...
 $ BT557BG    : Factor w/ 47931 levels: "", "BR925BN", "BT00BT",...: 30877 30962 30972 43705 30854 30904 30972 30557 30881 30960 ...
 $ X281855    : int  281892 282306 282419 335367 281719 282080 282524 303527 282128 282220 ...
 $ X438598    : int  438228 438587 438387 369985 438366 438424 438243 444150 438612 438467 ...
 $ X1         : int    2 3 4 5 6 7 8 9 10 11 ...

```

Fig.4: Structure of the data frame

Output:

```

> head(NIPOSTCODE,10)
  V1      V2 V3  V4      V5 V6 V7      V8      V9      V10      V11      V12      V13      V14 V15
1   17      HIGH ROAD      MULLAGHACALL NORTH PORTSTEWART LONDONDERRY BT557BG 281855 438598 1
2   15 CONVENTION AVENUE      MULLAGHACALL NORTH PORTSTEWART LONDONDERRY BT557BW 281892 438228 2
3   13 STATION ROAD      MULLAGHACALL NORTH PORTSTEWART LONDONDERRY BT557HH 282306 438587 3
4   99 OLD COACH ROAD      MULLAGHACALL NORTH PORTSTEWART LONDONDERRY BT557HW 282419 438387 4
5   20 BREDA COURT      MULLAGHACALL NORTH PORTSTEWART LONDONDERRY BT557AR 281719 438366 6
6   11 UPPER HEATHMOUNT      MULLAGHACALL NORTH PORTSTEWART LONDONDERRY BT557EE 282080 438424 7
7   86 LEVER ROAD      MULLAGHACALL NORTH PORTSTEWART LONDONDERRY BT557HW 282524 438243 8
8   112 OLD COACH ROAD      MULLAGHACALL NORTH PORTSTEWART LONDONDERRY BT557DA 282128 438612 10
9   134 WHITEPARK ROAD      BALLINTOY BALLINTOY DEMESNE BALLYCASTLE ANTRIM BT546ND 303527 444150 9
10  FLAT 3 16 STATION ROAD      MULLAGHACALL NORTH PORTSTEWART LONDONDERRY BT557DA 282128 438612 10

```

Fig.5: Output without error in the headers

CSV files has 943033 row in total and the first 10 rows of the data frame are showed in the above figure. Now we can see that the headers are changed by keeping false.

- B) The next task is to assign the relevant titles to each attribute in the data frame. Initially the header of the columns is named as x, x.1, x.2 and so on. We have changed with the relevant titles as shown in the below code and stored in the data frame col_names. Modified data is then moved to NIPOSTCODE and output can be displayed by using view(NIPOSTCODE).

```

# Assigning the relevant title for each attribute of the data
col_names = c("Organisation Name",
              "Sub-building Name",
              "Building Name",
              "Number",
              "Primary Thorfare",
              "Alt Thorfare",
              "Secondary Thorfare",
              "Locality",
              "Townland",
              "Town",
              "County",
              "Postcode",
              "x-coordinates",
              "y-coordinates",
              "Primary Key (identifier)")
colnames(NIPOSTCODE) <- col_names
View(NIPOSTCODE)

```

Fig.3: Changing the names of each attributes

Output:

	Organisation Name	Sub-building Name	Building Name	Number	Primary Thorfare	Alt Thorfare	Secondary Thorfare	Locality	Townland	Town	County
1				17	HIGH ROAD				MULLAGHACALL NORTH	PORTSTEWART	PK
2				15	CONVENTION AVENUE				MULLAGHACALL NORTH	PORTSTEWART	PK
3				13	STATION ROAD				MULLAGHACALL NORTH	PORTSTEWART	PK
4				99	OLD COACH ROAD				MULLAGHACALL NORTH	PORTSTEWART	PK
5				20	BREDA COURT				BREDA	BELFAST	BI
6				11	UPPER HEATHMOUNT				MULLAGHACALL NORTH	PORTSTEWART	PK
7				86	LEVER ROAD				MULLAGHACALL NORTH	PORTSTEWART	PK
8				112	OLD COACH ROAD				MULLAGHACALL NORTH	PORTSTEWART	PK
9				134	WHITEPARK ROAD			BALLINTOY	BALLINTOY DEMESNE	BALLYCASTLE	BI
10		FLAT 3		16	STATION ROAD				MULLAGHACALL NORTH	PORTSTEWART	PK
11				15	HEATHERLEA AVENUE				MULLAGHACALL NORTH	PORTSTEWART	PK
12				30	STATION ROAD				MULLAGHACALL NORTH	PORTSTEWART	PK
13				21	HIGH ROAD				MULLAGHACALL NORTH	PORTSTEWART	PK
14				1	PROSPECT ROAD				TULLAGHMURRY WEST	PORTSTEWART	PK

Fig.4: Change in the titles

- C) In this task we need to replace and recode all the missing entries with a suitable identifier. In order to do this, we need to check the count of missing entries in each attribute by writing the code as shown in below figure.5

	Organisation Name	Sub-building Name	Building Name	Number	Primary Thorfare	Alt Thorfare	Secondary Thorfare	Locality	Townland	Town	County
1				17	HIGH ROAD				MULLAGHACALL NORTH	PORTSTEWART	PK
2				15	CONVENTION AVENUE				MULLAGHACALL NORTH	PORTSTEWART	PK
3				13	STATION ROAD				MULLAGHACALL NORTH	PORTSTEWART	PK
4				99	OLD COACH ROAD				MULLAGHACALL NORTH	PORTSTEWART	PK
5				20	BREDA COURT				BREDA	BELFAST	BI
6				11	UPPER HEATHMOUNT				MULLAGHACALL NORTH	PORTSTEWART	PK
7				86	LEVER ROAD				MULLAGHACALL NORTH	PORTSTEWART	PK
8				112	OLD COACH ROAD				MULLAGHACALL NORTH	PORTSTEWART	PK
9				134	WHITEPARK ROAD			BALLINTOY	BALLINTOY DEMESNE	BALLYCASTLE	BI
10		FLAT 3		16	STATION ROAD				MULLAGHACALL NORTH	PORTSTEWART	PK
11				15	HEATHERLEA AVENUE				MULLAGHACALL NORTH	PORTSTEWART	PK
12				30	STATION ROAD				MULLAGHACALL NORTH	PORTSTEWART	PK
13				21	HIGH ROAD				MULLAGHACALL NORTH	PORTSTEWART	PK
14				1	PROSPECT ROAD				TULLAGHMURRY WEST	PORTSTEWART	PK
15				29	LEVER ROAD				MULLAGHACALL NORTH	PORTSTEWART	PK
16		APARTMENT B		27	PRINCESS STREET				PORT RUSH	PORTSTEWART	PK
17				94	OLD COACH ROAD				MULLAGHACALL NORTH	PORTSTEWART	PK
18				1	HARRYVILLE				MULLAGHACALL NORTH	PORTSTEWART	PK

Fig.6: Initial missing entries in the dataset

Now we need to check the count of missing entries of each attributes. The lapply is used in the code to iterate each column and check the missing entries.

```

29 # Checking the count of missing entries in each attribute
30 data.frame(lapply(NIPOSTCODE, function(x) sum(x == "")))
31

```

Fig.5:

Output:

```

> # Checking the count of missing entries in each attribute
> data.frame(lapply(NIPOSTCODE, function(x) sum(x == "")))
  Organisation.Name Sub.building.Name Building.Name Number Primary.Thorfare Alt.Thorfare Secondary.Thorfare Locality Townland Town
1      890537      884099      895540      28753      470      921788      938400      856789      0 19872
1 County Postcode x.coordinates y.coordinates Primary.Key..identifier.
1      0      8900      0      0      0
> |

```

Fig.6: Count of missing entries

Plotting a graph in order to justify my decision in choosing the best. The below code is used to plot the graph by using the summary function.

```
32 # Using library VIM
33 library(VIM)
34 # Plotting a graph to choose the best consideration for the missing entries
35 missing_values <- aggr(NIPOSTCODE, prop= FALSE, numbers = TRUE)
36 summary(missing_values)
```

Fig.7: using library vim

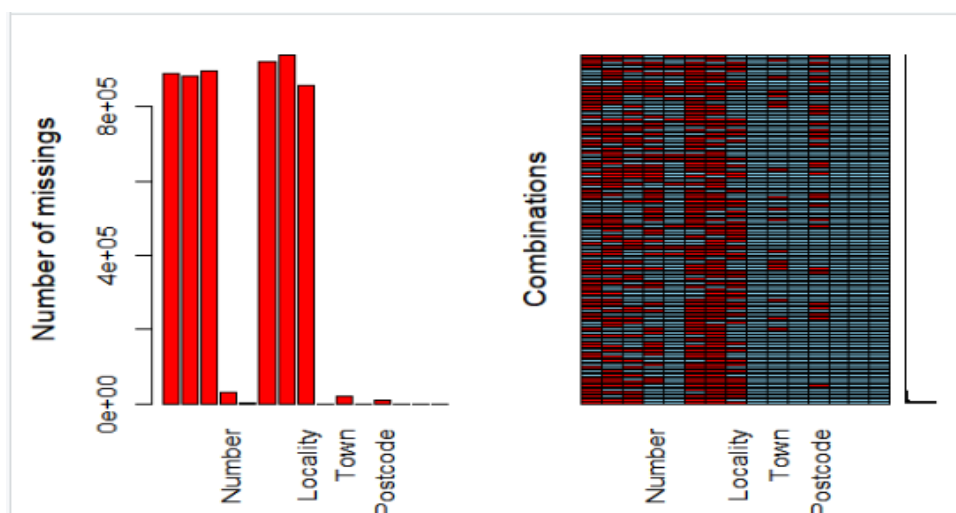


Fig.8: bar plot

While observing the graph I have decided to replace with NA to the missing entries because we cannot hardcode the data and some of the columns have missing values more than 50%. We cannot even remove the rows of the missing entries as they effect the other rows and lead to data loss.

```
32 # Replacing the missing values with the NA
33 NIPOSTCODE[NIPOSTCODE == ""] <- NA
34 # To view the data frame
35 View(NIPOSTCODE)
```

Fig.9: Recode and replace the missing entries

Output:

	Organisation Name	Sub-building Name	Building Name	Number	Primary Thorfare	Alt Thorfare	Secondary Thorfare	Locality	Townland
1	NA	NA	NA	17	HIGH ROAD	NA	NA	NA	MULLAGH
2	NA	NA	NA	15	CONVENTION AVENUE	NA	NA	NA	MULLAGH
3	NA	NA	NA	13	STATION ROAD	NA	NA	NA	MULLAGH
4	NA	NA	NA	99	OLD COACH ROAD	NA	NA	NA	MULLAGH
5	NA	NA	NA	20	BREDA COURT	NA	NA	NA	BREDA
6	NA	NA	NA	11	UPPER HEATHMOUNT	NA	NA	NA	MULLAGH
7	NA	NA	NA	86	LEVER ROAD	NA	NA	NA	MULLAGH
8	NA	NA	NA	112	OLD COACH ROAD	NA	NA	NA	MULLAGH
9	NA	NA	NA	134	WHITEPARK ROAD	NA	NA	BALLINTOY	BALLINTOY
10	NA	FLAT 3	NA	16	STATION ROAD	NA	NA	NA	MULLAGH
11	NA	NA	NA	15	HEATHERLEA AVENUE	NA	NA	NA	MULLAGH
12	NA	NA	NA	30	STATION ROAD	NA	NA	NA	MULLAGH
13	NA	NA	NA	21	HIGH ROAD	NA	NA	NA	MULLAGH

Fig.7: Dataset replaced with NA

- D) As I have replaced with NA the count of missing entries and after replacing with NA has same number of count. The function used in the code is used to iterate each column of the dataset and check the count of the NA values in each attribute.

```
43 # Checking the count of NA in each attribute
44 data.frame(sapply(NIPOSTCODE, function(x)sum(length(which(is.na(x))))))
45
```

Fig.10: Checking the count of NA

Output:

```
> data.frame(sapply(NIPOSTCODE, function(x)sum(length(which(is.na(x))))))
      sapply.NIPOSTCODE..function.x..sum.length.which.is.na.x.....
Organisation Name                                     890537
Sub-building Name                                     884099
Building Name                                          895540
Number                                                 28753
Primary Thorfare                                       470
Alt Thorfare                                           921788
Secondary Thorfare                                    938400
Locality                                              856789
Townland                                               0
Town                                                  19872
County                                                 0
Postcode                                              8900
x-coordinates                                         0
y-coordinates                                         0
Primary Key (identifier)                             0
```

Fig.11: structure of Nipostcode

- E) In this task we need to shift the primary key identifier to the first of the dataset. We can see in the below figure the primary key identifier is at the last.

Alt Thorfare	Secondary Thorfare	Locality	Townland	Town	County	Postcode	x-coordinates	y-coordinates	Primary Key (identifier)
			MULLAGHACALL NORTH	PORTSTEWART	LONDONDERRY	BT557BG	281855	438598	1
			MULLAGHACALL NORTH	PORTSTEWART	LONDONDERRY	BT557BW	281892	438228	2
			MULLAGHACALL NORTH	PORTSTEWART	LONDONDERRY	BT557HH	282306	438587	3
			MULLAGHACALL NORTH	PORTSTEWART	LONDONDERRY	BT557HW	282419	438387	4
			BREDA	BELFAST	DOWN	BT86JB	335367	369985	5
			MULLAGHACALL NORTH	PORTSTEWART	LONDONDERRY	BT557AR	281719	438366	6
			MULLAGHACALL NORTH	PORTSTEWART	LONDONDERRY	BT557EE	282080	438424	7
			MULLAGHACALL NORTH	PORTSTEWART	LONDONDERRY	BT557HW	282524	438243	8
		BALLINTOY	BALLINTOY DEMESNE	BALLYCASTLE	ANTRIM	BT546ND	303527	444150	9
			MULLAGHACALL NORTH	PORTSTEWART	LONDONDERRY	BT557DA	282128	438612	10

Fig.12: Data frame before shifting the primary key

Now we need to move the primary key identifier to the first place. In the code I have changed the column values by shifting the last value to the first. As we can see in the below code the number 15 is shifted to the first place. In output we can see that primary key identifier is shifted to the first place.

```
46 # Changing the position of the primary key to the first
47 NIPOSTCODE <-NIPOSTCODE[, c(15,1,2,3,4,5,6,7,8,9,10,11,12,13,14)]
48 # Displaying the first 10 rows
49 head(NIPOSTCODE,10)
```

Fig.13: changing the position of primary key identifier

Output:

	Primary Key (Identifier)	Organisation Name	Sub-building Name	Building Name	Number	Primary Thoroughfare	Alt Thoroughfare	Secondary Thoroughfare	Locality	Townland
1	1	NA	NA	NA	17	HIGH ROAD	NA	NA	NA	MULLAGHACALL NORTH
2	2	NA	NA	NA	15	CONVENTION AVENUE	NA	NA	NA	MULLAGHACALL NORTH
3	3	NA	NA	NA	13	STATION ROAD	NA	NA	NA	MULLAGHACALL NORTH
4	4	NA	NA	NA	99	OLD COACH ROAD	NA	NA	NA	MULLAGHACALL NORTH
5	5	NA	NA	NA	20	BREDA COURT	NA	NA	NA	BREDA
6	6	NA	NA	NA	11	UPPER HEATHMOUNT	NA	NA	NA	MULLAGHACALL NORTH
7	7	NA	NA	NA	86	LEVER ROAD	NA	NA	NA	MULLAGHACALL NORTH
8	8	NA	NA	NA	112	OLD COACH ROAD	NA	NA	NA	MULLAGHACALL NORTH
9	9	NA	NA	NA	134	WHITEPARK ROAD	NA	NA	BALLINTOY	BALLINTOY
10	10	NA	FLAT 3	NA	16	STATION ROAD	NA	NA	NA	MULLAGHACALL NORTH

Fig.14:

- F) In this task we need to create a new dataset called `Limavady_data` and store information only within it where locality, townland, town contain the name limavady. Count, display the rows and store the information in the csv file called `limavady`.

Primary Thoroughfare	Alt Thoroughfare	Secondary Thoroughfare	Locality	Townland	Town	County	Postcode	x-coordinates	y-coordinates
HIGH ROAD	NA	NA	NA	MULLAGHACALL NORTH	PORTSTEWART	LONDONDERRY	BT557BG	281855	438598
CONVENTION AVENUE	NA	NA	NA	MULLAGHACALL NORTH	PORTSTEWART	LONDONDERRY	BT557BW	281892	438228
STATION ROAD	NA	NA	NA	MULLAGHACALL NORTH	PORTSTEWART	LONDONDERRY	BT557HH	282306	438587
OLD COACH ROAD	NA	NA	NA	MULLAGHACALL NORTH	PORTSTEWART	LONDONDERRY	BT557HW	282419	438387
BREDA COURT	NA	NA	NA	BREDA	BELFAST	DOWN	BT86JB	335367	369985
UPPER HEATHMOUNT	NA	NA	NA	MULLAGHACALL NORTH	PORTSTEWART	LONDONDERRY	BT557AR	281719	438366
LEVER ROAD	NA	NA	NA	MULLAGHACALL NORTH	PORTSTEWART	LONDONDERRY	BT557EE	282080	438424
OLD COACH ROAD	NA	NA	NA	MULLAGHACALL NORTH	PORTSTEWART	LONDONDERRY	BT557HW	282524	438243
WHITEPARK ROAD	NA	NA	BALLINTOY	BALLINTOY DEMESNE	BALLYCASTLE	ANTRIM	BT546ND	303527	444150
STATION ROAD	NA	NA	NA	MULLAGHACALL NORTH	PORTSTEWART	LONDONDERRY	BT557DA	282128	438612

Fig.15: Dataset before applying the changes

Code: In this task we have used `subset` function to check the columns within the `NIPostcode` dataset which contains the name `limavady`. In the code `subset` function uses the `NI` postcode data and checks the columns `locality`, `townland`, `town` and store only the information containing the name `limavady`. After that check the number of rows and store the information in the csv file and name it as `limavady`.

```

51 # To store only the data within name containing limavady
52 attach(NIPOSTCODE)
53 limavady_data <- subset(NIPOSTCODE, Locality=="LIMAVADY"|Townland == "LIMAVADY"|Town=="LIMAVADY")
54 View(limavady_data)
55 # To check the rows in the Limavady data
56 nrow(limavady_data)
57 # The modified data is loaded in to a csv file
58 write.csv(limavady_data, "Limavady.csv")

```

Fig.16: creating a data frame

Output:

As we can see in the output the name containing `limavady` is stored in the assigned columns in the dataset `limavady`. Number of rows are displayed in the below figure.17 output


```
> # To check the rows in the Limavady data
> nrow(limavady_data)
[1] 8467
>
```

Fig.17: Number of rows of limavady data

Limavady.csv - Excel																		
Rampai Paul Noah Sujith																		
Tell me what you want to do...																		
Clipboard Font Alignment Number Styles Conditional Formatting Cell Styles Insert Delete Format AutoSum Sort & Find Filter Select																		
L1 Town																		
	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R
1		Primary Ki	Organisation	Sub-build	Building Number	Primary TI	Alt Thorfa	Secondary Locality	Townland	Town	County	Postcode	x-coordinates	y-coordinates				
2	73	73	NA	NA	NA	30	LEIGHRY	NA	NA	BALLYLEIGHRY UPPER	LIMAVADY	LONDONDERRY	BT490IG	270404	431604			
3	75	75	NA	NA	NA	38	CURRAGH	NA	NA	BALLYSCU DUNCRUN	LIMAVADY	LONDONDERRY	BT490JE	268543	432534			
4	76	76	NA	NA	NA	405	SEACOAST	NA	NA	CLOONEY	LIMAVADY	LONDONDERRY	BT490LW	268640	433777			
5	89	89	NA	NA	NA	370	SEACOAST	NA	NA	BALLYSCULLION	LIMAVADY	LONDONDERRY	BT490LA	267172	432608			
6	92	92	NA	TELEPHON	NA		SEACOAST	NA	NA	OUGHTYMOYLE	LIMAVADY	LONDONDERRY	BT490JY	266864	431456			
7	99	99	NA	NA	NA	21	DRUMAV	NA	NA	DRUMAVALLY	LIMAVADY	LONDONDERRY	BT490LT	267148	432852			
8	226	226	NA	NA	NA	254	SEACOAST	NA	NA	BELLARENA	LIMAVADY	LONDONDERRY	BT490HZ	266391	429886			
9	236	236	NA	NA	NA	100	DOWLANE	NA	NA	FRUITHILL CARBULLION	LIMAVADY	LONDONDERRY	BT490HR	267029	425935			
10	310	310	NA	NA	NA	47	DOWLANE	NA	NA	FRUITHILL ARTIKELLY	LIMAVADY	LONDONDERRY	BT490HR	268259	424813			
11	455	455	NA	APARTME	NA	5	LODGE CO	NA	NA	NEWTOWN LIMAVADY	LIMAVADY	LONDONDERRY	BT490EY	266899	423197			
12	456	456	NA	NA	NA	40	LINENHAL	NA	NA	NEWTOWN LIMAVADY	LIMAVADY	LONDONDERRY	BT490HQ	266936	423105			
13	457	457	DRY ARCH	NA	NA	55	CATHERIN	NA	NA	NEWTOWN LIMAVADY	LIMAVADY	LONDONDERRY	BT490DA	266943	422977			
14	458	458	NA	NA	NA	9	TYLER PAR	NA	NA	RATHBRADY MORE	LIMAVADY	LONDONDERRY	BT490DS	267857	423478			
15	461	461	NA	NA	NA	9	BALLYCLO	NA	NA	RATHBRADY MORE	LIMAVADY	LONDONDERRY	BT490BL	267359	423426			
16	462	462	NA	NA	NA	100	CONNELL	NA	NA	NEWTOWN LIMAVADY	LIMAVADY	LONDONDERRY	BT490EA	267364	422972			
17	463	463	NA	NA	NA	12	DUNMORI	NA	NA	NEWTOWN LIMAVADY	LIMAVADY	LONDONDERRY	BT490AN	267370	422688			
18	464	464	NA	NA	NA	2	GLENBEG	NA	NA	ENAGH	LIMAVADY	LONDONDERRY	BT490NL	267723	421964			
19	465	465	NA	NA	NA	74A	SCROGGY	NA	NA	ENAGH	LIMAVADY	LONDONDERRY	BT490NB	267986	422179			
20	466	466	NA	NA	NA	51	CARLARAC	NA	NA	TERRYDRUM	LIMAVADY	LONDONDERRY	BT490LF	265869	419219			
21	467	467	CRAIGS SE	NA	NA	6	BALLYCLO	NA	NA	KILLANE	LIMAVADY	LONDONDERRY	BT490BN	267310	423476			
22	468	468	NA	NA	NA	42	DRUMACH	NA	NA	ENAGH	LIMAVADY	LONDONDERRY	BT490NZ	268024	422250			
23	472	472	NA	NA	NA	13	PROTESTA	NA	NA	NEWTOWN LIMAVADY	LIMAVADY	LONDONDERRY	BT499BP	266982	422879			

Fig.18: Limavady.csv file

Now check the rows of the limavady data and

- G) Now save the modified NIPostCode data into a csv file and name it as CleanNIPostCodeData. Write.csv file saves the data into a csv file.

```
60 # Saving the modified data in new csv file
61 write.csv(NIPOSTCODE, "CleanNIPostCode.csv")
```

Fig.19: writing the NIpostcodedata

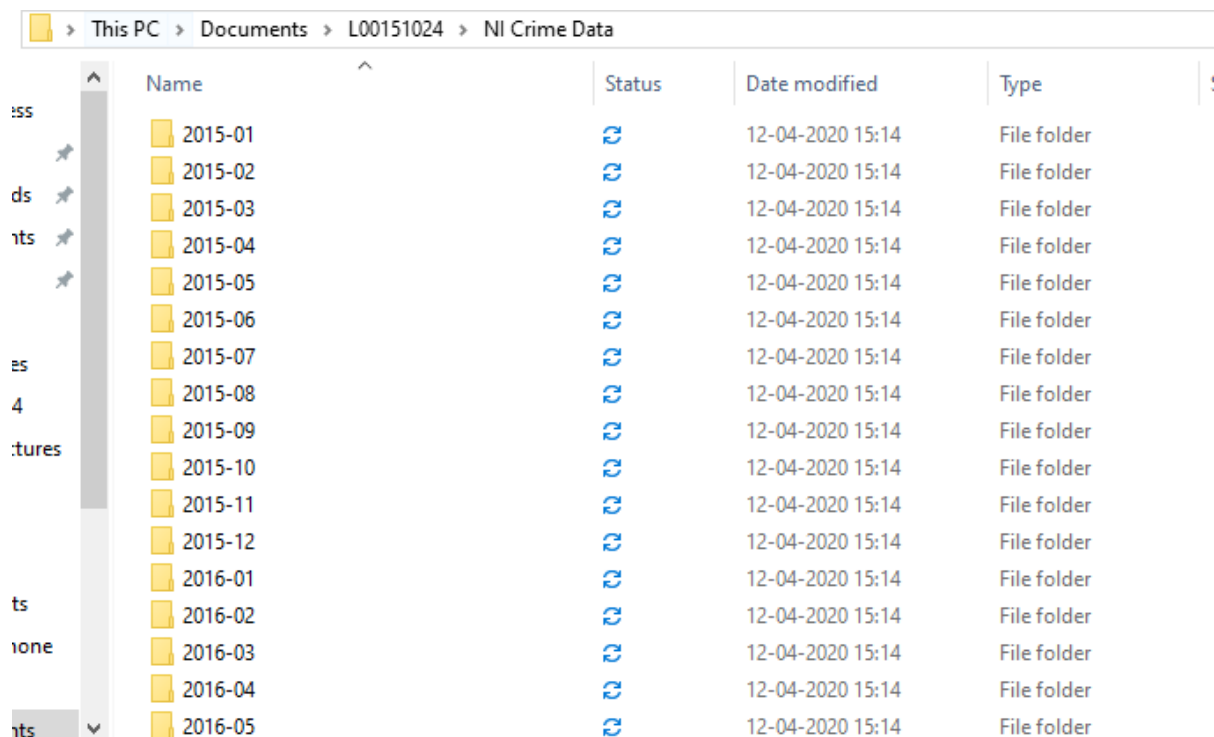
Output:

CleanNIPostCode.csv - Excel																			Rajmaji Paul Noah Sujith		Share																																													
File Home Insert Page Layout Formulas Data Review View Tell me what you want to do...																			Conditional Formatting				Format as Table		Cell Styles		Insert		Delete		Format		AutoSum		Sort & Find		Filter & Select																													
Clipboard			Font			Alignment			Number			Conditional Formatting		Format as Table		Cell Styles		Insert		Delete		Format		AutoSum		Sort & Find		Filter & Select																																						
Paste			Cut			Copy			Format Painter			General		Number		Text		Formulas		References		Styles		Cells		Editing																																								
A1																			X		Y		Z		AA		AB		AC		AD		AE		AF		AG		AH		AI		AJ		AK		AL		AM		AN		AO		AP		AQ		AR		AS		AT		AU	
	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U																																													
1		Primary Ki	Organisation	Sub-build	Building Number	Primary TI	Alt Thorfa	Secondary Locality	Townland	Town	County	Postcode	x-coordinates	y-coordinates																																																				
2	1	1	NA	NA	NA	17	HIGH ROA	NA	NA	NA	MULLAGH PORTSTEV	LONDONDERRY	BT557BG	281855	438598																																																			
3	2	2	NA	NA	NA	15	CONVENT	NA	NA	NA	MULLAGH PORTSTEV	LONDONDERRY	BT557BW	281892	438228																																																			
4	3	3	NA	NA	NA	13	STATION FNA	NA	NA	NA	MULLAGH PORTSTEV	LONDONDERRY	BT557HH	282306	438587																																																			
5	4	4	NA	NA	NA	99	OLD COACNA	NA	NA	NA	MULLAGH PORTSTEV	LONDONDERRY	BT557HW	282419	438387																																																			
6	5	5	NA	NA	NA	20	BREDA CO NA	NA	NA	NA	BREDA BELFAST DOWN	BT86JB	335367	369985																																																				
7	6	6	NA	NA	NA	11	UPPER HE,NA	NA	NA	NA	MULLAGH PORTSTEV	LONDONDERRY	BT557AR	281719	438366																																																			
8	7	7	NA	NA	NA	86	LEVER RO,NA	NA	NA	NA	MULLAGH PORTSTEV	LONDONDERRY	BT557EE	282080	438424																																																			
9	8	8	NA	NA	NA	112	OLD COACNA	NA	NA	NA	MULLAGH PORTSTEV	LONDONDERRY	BT557HW	282524	438243																																																			
10	9	9	NA	NA	NA	134	WHITEPAFNA	NA	BALLINTO	BALLINTO	BALLYCAS ANTRIM	BT546ND	303527	444150																																																				
11	10	10	NA	FLAT 3	NA	16	STATION FNA	NA	NA	NA	MULLAGH PORTSTEV	LONDONDERRY	BT557DA	282128	438612																																																			
12	11	11	NA	NA	NA	15	HEATHER,NA	NA	NA	NA	MULLAGH PORTSTEV	LONDONDERRY	BT557HF	282220	438467																																																			
13	12	12	NA	NA	NA	30	STATION FNA	NA	NA	NA	MULLAGH PORTSTEV	LONDONDERRY	BT557DA	282222	438603																																																			
14	13	13	NA	NA	NA	21	HIGH ROA NA	NA	NA	NA	MULLAGH PORTSTEV	LONDONDERRY	BT557BG	281900	438599																																																			
15	14	14	NA	NA	NA	1	PROSPECT,NA	NA	NA	NA	TULLAGH PORTSTEV	LONDONDERRY	BT557NF	281522	437730																																																			
16	15	15	NA	NA	NA	29	LEVER RO,NA	NA	NA	NA	MULLAGH PORTSTEV	LONDONDERRY	BT557EB	281947	438175																																																			
17	16	16	NA	APARTME	NA	27	PRINCESS,NA	NA	NA	NA	PORT RUS PORTRUS	ANTRIM	BT568AX	285529	441033																																																			
18	17	17	NA	NA	NA	94	OLD COAC NA	NA	NA	NA	MULLAGH PORTSTEV	LONDONDERRY	BT557HW	282358	438327																																																			
19	18	18	NA	NA	NA	1	HARRYVIL,NA	NA	NA	NA	MULLAGH PORTSTEV	LONDONDERRY	BT557AU	281730	438499																																																			
20	19	19	NA	NA	NA	13	CENTRAL,NA	NA	NA	NA	MULLAGH PORTSTEV	LONDONDERRY	BT557BP	282034	438521																																																			
21	20	20	NA	FLAT 29	BLOCK B	NA	COVEHILL	NA	NA	NA	GLENAMA PORTRUS	ANTRIM	BT568GL	284936	439713																																																			
22	21	21	NA	NA	NA	20	OLD COAC NA	NA	NA	NA	MULLAGH PORTSTEV	LONDONDERRY	BT557BX	281906	438507																																																			
23	22	22	NA	2A	NA	BLOCK 2	MILLFIELD	NA	MILL ROA	NA	GALVALLY PORTSTEV	LONDONDERRY	BT557PQ	282443	437855																																																			

Fig.20: CleanNIPostCode.csv file

Section:2

To complete this section, they have given a dataset which contains all the crime data of Northern Ireland from 2015 to 2016 in different folders as shown in below figure



Name	Status	Date modified	Type
2015-01		12-04-2020 15:14	File folder
2015-02		12-04-2020 15:14	File folder
2015-03		12-04-2020 15:14	File folder
2015-04		12-04-2020 15:14	File folder
2015-05		12-04-2020 15:14	File folder
2015-06		12-04-2020 15:14	File folder
2015-07		12-04-2020 15:14	File folder
2015-08		12-04-2020 15:14	File folder
2015-09		12-04-2020 15:14	File folder
2015-10		12-04-2020 15:14	File folder
2015-11		12-04-2020 15:14	File folder
2015-12		12-04-2020 15:14	File folder
2016-01		12-04-2020 15:14	File folder
2016-02		12-04-2020 15:14	File folder
2016-03		12-04-2020 15:14	File folder
2016-04		12-04-2020 15:14	File folder
2016-05		12-04-2020 15:14	File folder

Fig.1: Checking the dataset folders

- A) In this task we need to combine all the files in crime data from each csv file into one dataset and save this dataset into a csv file called AllNICrimeData. Show the count of number of rows in the AllNICrimeData.

Code: In this task we have used list function is used to extract the files within the folders in the NICrimeData by using the for loop. In this task cannot hardcode the data as it is a bad option to combine the files. Rbind is used to combine the data frames by rows. Finally saving the combined dataset into a csv file called AllNICrimeData.csv.


```

2 # A)
3 # combining all of the crime data from each csv file into one data set
4 # creating an empty dataframe to store the data after looping
5 final_crime_data <- data.frame()
6 #files_list contains all the csv files from NI Crime Data
7 files_list <- list.files(recursive = TRUE)
8 for(file_list in files_list){
9   crime_data <- read.csv(file_list)
10  final_crime_data <- rbind(final_crime_data, crime_data)
11 }
12 # Getting the current working directory
13 getwd()
14 setwd("NI Crime Data/")
15 # Saving the combined dataset into a csv file
16 write.csv(final_crime_data, "AllNICrimeData.csv")
17 AllNICrimeData <- read.csv("AllNICrimeData.csv")
18 str(AllNICrimeData)
19 # showing the number of rows
20 nrow(AllNICrimeData)

```

Fig.1: code for task A

Output:

In the output we can see that all the files from the NI crime folders have been combined and count of rows is shown in the below figure.

```

> # Saving the combined dataset into a csv file
> write.csv(final_crime_data, "AllNICrimeData.csv")
> AllNICrimeData <- read.csv("AllNICrimeData.csv")
> str(AllNICrimeData)
'data.frame': 477696 obs. of 13 variables:
 $ X      : int  1 2 3 4 5 6 7 8 9 10 ...
 $ Crime.ID : Factor w/ 11667 levels "", "0009d3218c478283888080303fed14c46c61e6b3b8f55963a2671dac3afb3907",...: 1 1 1
 1 1 1 1 1 1 ...
 $ Month   : Factor w/ 36 levels "2015-01","2015-02",...: 1 1 1 1 1 1 1 1 1 1 ...
 $ Reported.by : Factor w/ 1 level "Police Service of Northern Ireland": 1 1 1 1 1 1 1 1 1 1 ...
 $ Falls.within : Factor w/ 1 level "Police Service of Northern Ireland": 1 1 1 1 1 1 1 1 1 1 ...
 $ Longitude : num  -6 -5.71 -5.82 -6.39 -6.25 ...
 $ Latitude  : num  54.6 54.6 54.7 54.2 54.9 ...
 $ Location  : Factor w/ 14984 levels "No Location",...: 12359 2 9993 3624 13143 8121 11826 2 6913 7811 ...
 $ LSOA.code : logi  NA NA NA NA NA NA ...
 $ LSOA.name : logi  NA NA NA NA NA NA ...
 $ Crime.type : Factor w/ 14 levels "Anti-social behaviour",...: 1 1 1 1 1 1 1 1 1 1 ...
 $ Last.outcome.category: logi  NA NA NA NA NA NA ...
 $ Context   : logi  NA NA NA NA NA NA ...
> nrow(AllNICrimeData)
[1] 477696

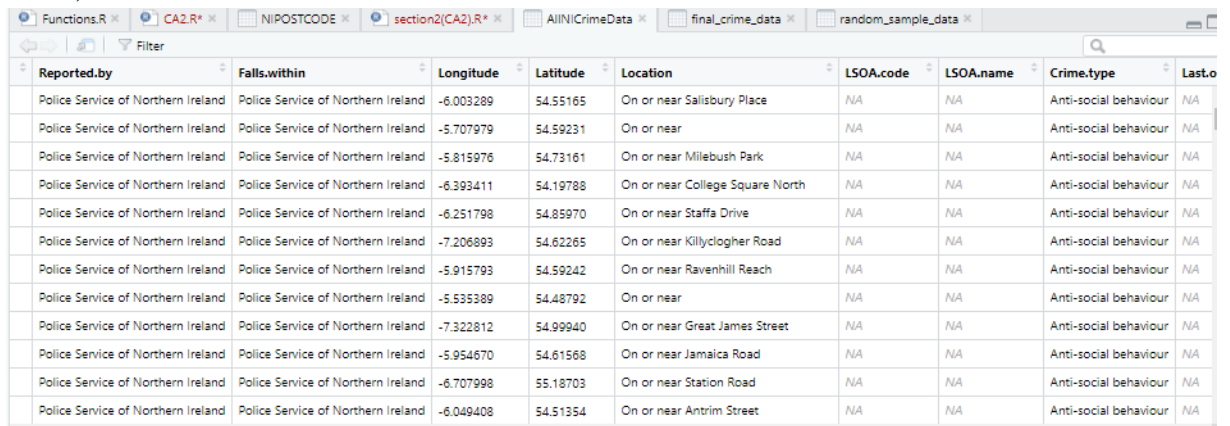
```

Fig.2: output for the structure and count of the AllNICrimeData

X	Crime.ID	Month	Reported.by	Falls.within	Longitude	Latitude	Location	LSOA.code	LSOA.name
1	1	2015-01	Police Service of Northern Ireland	Police Service of Northern Ireland	-6.003289	54.55165	On or near Salisbury Place	NA	NA
2	2	2015-01	Police Service of Northern Ireland	Police Service of Northern Ireland	-5.707979	54.59231	On or near	NA	NA
3	3	2015-01	Police Service of Northern Ireland	Police Service of Northern Ireland	-5.815976	54.73161	On or near Milebush Park	NA	NA
4	4	2015-01	Police Service of Northern Ireland	Police Service of Northern Ireland	-6.393411	54.19788	On or near College Square North	NA	NA
5	5	2015-01	Police Service of Northern Ireland	Police Service of Northern Ireland	-6.251798	54.85970	On or near Staffa Drive	NA	NA
6	6	2015-01	Police Service of Northern Ireland	Police Service of Northern Ireland	-7.206893	54.62265	On or near Killyclogher Road	NA	NA
7	7	2015-01	Police Service of Northern Ireland	Police Service of Northern Ireland	-5.915793	54.59242	On or near Ravenhill Reach	NA	NA
8	8	2015-01	Police Service of Northern Ireland	Police Service of Northern Ireland	-5.535389	54.48792	On or near	NA	NA
9	9	2015-01	Police Service of Northern Ireland	Police Service of Northern Ireland	-7.322812	54.99940	On or near Great James Street	NA	NA
10	10	2015-01	Police Service of Northern Ireland	Police Service of Northern Ireland	-5.954670	54.61568	On or near Jamaica Road	NA	NA
11	11	2015-01	Police Service of Northern Ireland	Police Service of Northern Ireland	-6.707998	55.18703	On or near Station Road	NA	NA
12	12	2015-01	Police Service of Northern Ireland	Police Service of Northern Ireland	-6.049408	54.51354	On or near Antrim Street	NA	NA

Fig.3: output for the combined files placed in the AllNICrimeData file

- B) In this task we need to modify newly created AllNICrimeData by removing the following attributed like Crime ID, Reported by, Falls within, LSOA code, LSOA name, last outcome and context.



The screenshot shows an RStudio window with several tabs. The active tab is 'AllNICrimeData', which displays a data table with the following columns: Reported.by, Falls.within, Longitude, Latitude, Location, LSOA.code, LSOA.name, Crime.type, and Last.o. The data rows show various police service locations in Northern Ireland, all reporting 'Anti-social behaviour'.

Reported.by	Falls.within	Longitude	Latitude	Location	LSOA.code	LSOA.name	Crime.type	Last.o
Police Service of Northern Ireland	Police Service of Northern Ireland	-6.003289	54.55165	On or near Salisbury Place	NA	NA	Anti-social behaviour	NA
Police Service of Northern Ireland	Police Service of Northern Ireland	-5.707979	54.59231	On or near	NA	NA	Anti-social behaviour	NA
Police Service of Northern Ireland	Police Service of Northern Ireland	-5.815976	54.73161	On or near Milebush Park	NA	NA	Anti-social behaviour	NA
Police Service of Northern Ireland	Police Service of Northern Ireland	-6.393411	54.19788	On or near College Square North	NA	NA	Anti-social behaviour	NA
Police Service of Northern Ireland	Police Service of Northern Ireland	-6.251798	54.85970	On or near Staffa Drive	NA	NA	Anti-social behaviour	NA
Police Service of Northern Ireland	Police Service of Northern Ireland	-7.206893	54.62265	On or near Killyclogher Road	NA	NA	Anti-social behaviour	NA
Police Service of Northern Ireland	Police Service of Northern Ireland	-5.915793	54.59242	On or near Ravenhill Reach	NA	NA	Anti-social behaviour	NA
Police Service of Northern Ireland	Police Service of Northern Ireland	-5.535389	54.48792	On or near	NA	NA	Anti-social behaviour	NA
Police Service of Northern Ireland	Police Service of Northern Ireland	-7.322812	54.99940	On or near Great James Street	NA	NA	Anti-social behaviour	NA
Police Service of Northern Ireland	Police Service of Northern Ireland	-5.954670	54.61568	On or near Jamaica Road	NA	NA	Anti-social behaviour	NA
Police Service of Northern Ireland	Police Service of Northern Ireland	-6.707998	55.18703	On or near Station Road	NA	NA	Anti-social behaviour	NA
Police Service of Northern Ireland	Police Service of Northern Ireland	-6.049408	54.51354	On or near Antrim Street	NA	NA	Anti-social behaviour	NA

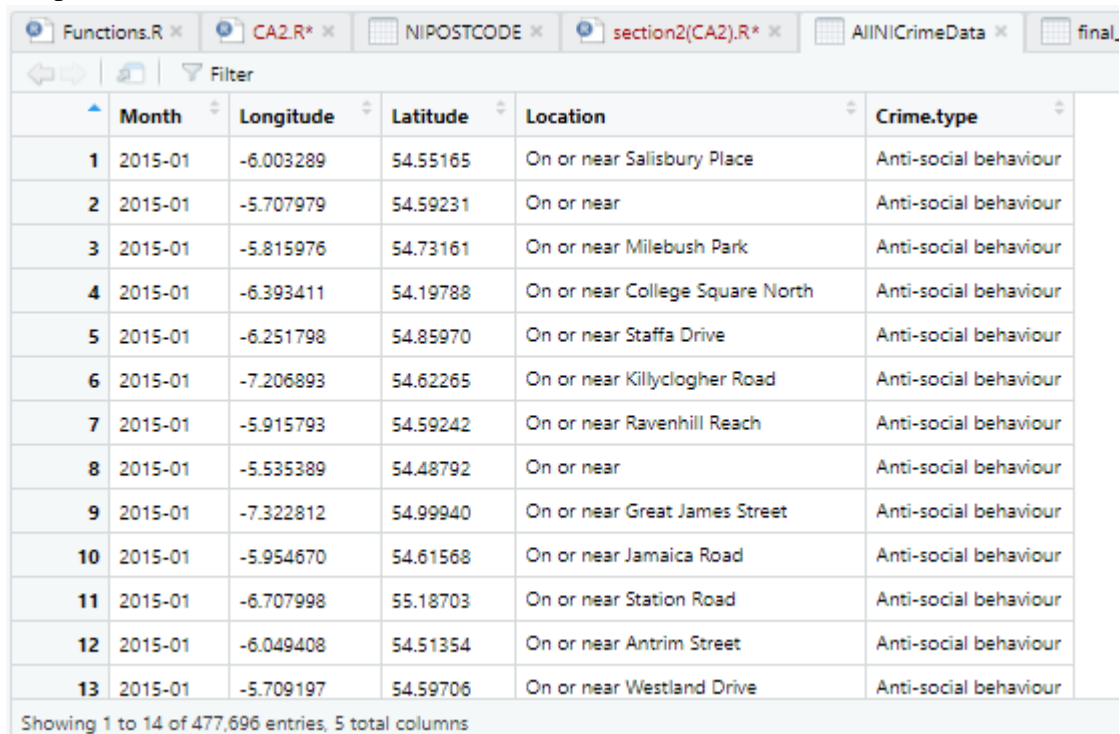
Fig.4: Before removing the attributes

Code: I have selected the required attributes and stored into the AllNICrimeData instead of removing the above given attributes.

```
21 # B)
22 # Selecting required attribute and removing other attributes according to the given task
23 AllNICrimeData <- final_crime_data[c(2,5,6,7,10)]
24 str(AllNICrimeData)
```

Fig.5:

Output:



The screenshot shows the RStudio window with the 'final_crime_data' tab active. The data table now has columns: Month, Longitude, Latitude, Location, and Crime.type. The data rows are numbered 1 to 13, showing the same locations as before but with the 'Month' attribute added and other attributes removed.

	Month	Longitude	Latitude	Location	Crime.type
1	2015-01	-6.003289	54.55165	On or near Salisbury Place	Anti-social behaviour
2	2015-01	-5.707979	54.59231	On or near	Anti-social behaviour
3	2015-01	-5.815976	54.73161	On or near Milebush Park	Anti-social behaviour
4	2015-01	-6.393411	54.19788	On or near College Square North	Anti-social behaviour
5	2015-01	-6.251798	54.85970	On or near Staffa Drive	Anti-social behaviour
6	2015-01	-7.206893	54.62265	On or near Killyclogher Road	Anti-social behaviour
7	2015-01	-5.915793	54.59242	On or near Ravenhill Reach	Anti-social behaviour
8	2015-01	-5.535389	54.48792	On or near	Anti-social behaviour
9	2015-01	-7.322812	54.99940	On or near Great James Street	Anti-social behaviour
10	2015-01	-5.954670	54.61568	On or near Jamaica Road	Anti-social behaviour
11	2015-01	-6.707998	55.18703	On or near Station Road	Anti-social behaviour
12	2015-01	-6.049408	54.51354	On or near Antrim Street	Anti-social behaviour
13	2015-01	-5.709197	54.59706	On or near Westland Drive	Anti-social behaviour

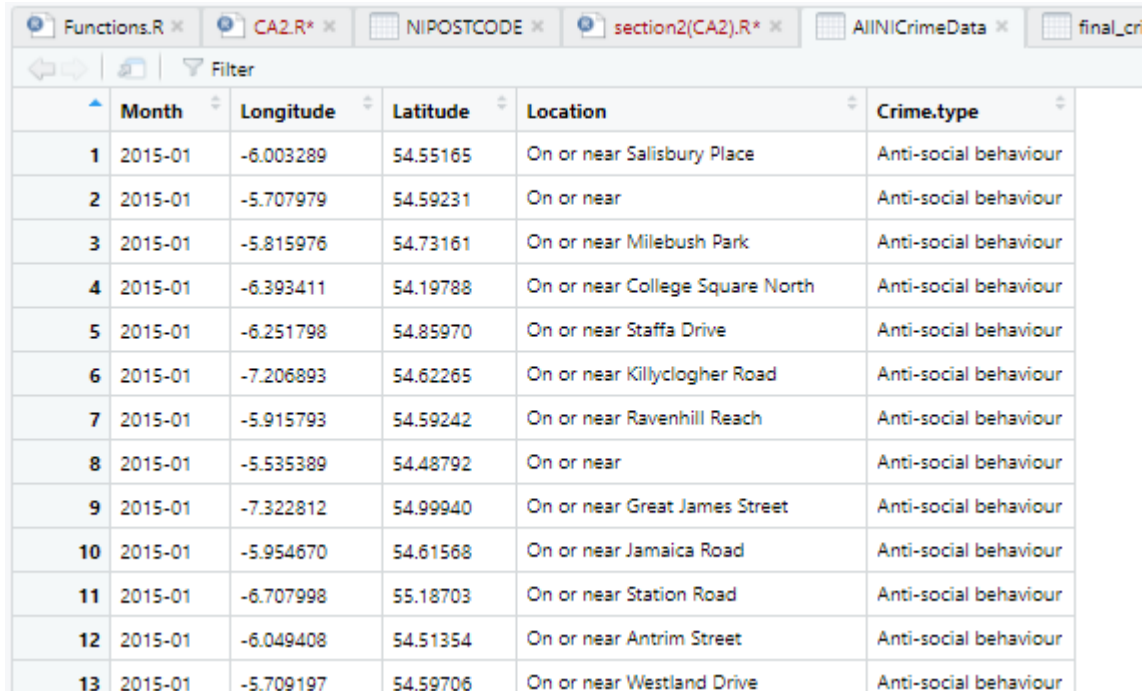
Showing 1 to 14 of 477,696 entries, 5 total columns

Fig.6: Output after removing certain attributes

```
> str(AllNICrimeData)
'data.frame': 477696 obs. of 5 variables:
 $ Month      : Factor w/ 36 levels "2015-01","2015-02",...: 1 1 1 1 1 1 1 1 1 1 ...
 $ Longitude  : num -6 -5.71 -5.82 -6.39 -6.25 ...
 $ Latitude   : num 54.6 54.6 54.7 54.2 54.9 ...
 $ Location   : Factor w/ 14984 levels "No Location",...: 3507 2 2790 1022 3724 2263 3363 2 1909 2179 ...
 $ Crime.type : Factor w/ 14 levels "Anti-social behaviour",...: 1 1 1 1 1 1 1 1 1 1 ...
```

Fig.7: Structure of the modified file

- C) In this task we need to change the crime type in AllNICrimeData by giving them short forms like for anti-social behavior be ASBO, Bicycle theft as BITH and so on as shown in the below code.



	Month	Longitude	Latitude	Location	Crime.type
1	2015-01	-6.003289	54.55165	On or near Salisbury Place	Anti-social behaviour
2	2015-01	-5.707979	54.59231	On or near	Anti-social behaviour
3	2015-01	-5.815976	54.73161	On or near Milebush Park	Anti-social behaviour
4	2015-01	-6.393411	54.19788	On or near College Square North	Anti-social behaviour
5	2015-01	-6.251798	54.85970	On or near Staffa Drive	Anti-social behaviour
6	2015-01	-7.206893	54.62265	On or near Killyclogher Road	Anti-social behaviour
7	2015-01	-5.915793	54.59242	On or near Ravenhill Reach	Anti-social behaviour
8	2015-01	-5.535389	54.48792	On or near	Anti-social behaviour
9	2015-01	-7.322812	54.99940	On or near Great James Street	Anti-social behaviour
10	2015-01	-5.954670	54.61568	On or near Jamaica Road	Anti-social behaviour
11	2015-01	-6.707998	55.18703	On or near Station Road	Anti-social behaviour
12	2015-01	-6.049408	54.51354	On or near Antrim Street	Anti-social behaviour
13	2015-01	-5.709197	54.59706	On or near Westland Drive	Anti-social behaviour

Fig.8: Data frame before shortening the crime type

Code: In the below code I have used function as.character to change the character and assign a short forms to the attributes.

```
26 # C)
27 # Giving short forms to the crime type attribute
28 AllNICrimeData$Crime.type <- as.character(AllNICrimeData$Crime.type)
29 AllNICrimeData$Crime.type[AllNICrimeData$Crime.type == "Anti-social behaviour"] <- "ASBO"
30 AllNICrimeData$Crime.type[AllNICrimeData$Crime.type == "Bicycle theft"] <- "BITH"
31 AllNICrimeData$Crime.type[AllNICrimeData$Crime.type == "Burglary"] <- "BURG"
32 AllNICrimeData$Crime.type[AllNICrimeData$Crime.type == "Criminal damage and arson"] <- "CDAR"
33 AllNICrimeData$Crime.type[AllNICrimeData$Crime.type == "DRUGS"] <- "DRUG"
34 AllNICrimeData$Crime.type[AllNICrimeData$Crime.type == "Other Theft"] <- "OTTH"
35 AllNICrimeData$Crime.type[AllNICrimeData$Crime.type == "Public order"] <- "PUBO"
36 AllNICrimeData$Crime.type[AllNICrimeData$Crime.type == "Robbery"] <- "ROBY"
37 AllNICrimeData$Crime.type[AllNICrimeData$Crime.type == "Shoplifting"] <- "SHOP"
38 AllNICrimeData$Crime.type[AllNICrimeData$Crime.type == "Theft from the person"] <- "THPR"
39 AllNICrimeData$Crime.type[AllNICrimeData$Crime.type == "Vehicle crime"] <- "VECR"
40 AllNICrimeData$Crime.type[AllNICrimeData$Crime.type == "Violence and sexual offences"] <- "VISO"
41 AllNICrimeData$Crime.type[AllNICrimeData$Crime.type == "Other crime"] <- "OTCR"
42 AllNICrimeData$Crime.type[AllNICrimeData$Crime.type == "Drugs"] <- "DRUG"
43 AllNICrimeData$Crime.type[AllNICrimeData$Crime.type == "Possession of weapons"] <- "POW"
```

Fig.8: code to shorten the crime type

Output:

	Month	Longitude	Latitude	Location	Crime.type
1	2015-01	-6.003289	54.55165	On or near Salisbury Place	ASBO
2	2015-01	-5.707979	54.59231	On or near	ASBO
3	2015-01	-5.815976	54.73161	On or near Milebush Park	ASBO
4	2015-01	-6.393411	54.19788	On or near College Square North	ASBO
5	2015-01	-6.251798	54.85970	On or near Staffa Drive	ASBO
6	2015-01	-7.206893	54.62265	On or near Killyclogher Road	ASBO
7	2015-01	-5.915793	54.59242	On or near Ravenhill Reach	ASBO
8	2015-01	-5.535389	54.48792	On or near	ASBO
9	2015-01	-7.322812	54.99940	On or near Great James Street	ASBO
10	2015-01	-5.954670	54.61568	On or near Jamaica Road	ASBO
11	2015-01	-6.707998	55.18703	On or near Station Road	ASBO
12	2015-01	-6.049408	54.51354	On or near Antrim Street	ASBO
13	2015-01	-5.709197	54.59706	On or near Westland Drive	ASBO

Showing 1 to 14 of 477,696 entries, 5 total columns

Fig.9: Output after shortening the crime type

- D) In this task we need to show a bar plot of the frequency of each crime from the crime type field and specify relevant options such as labels and bar colors.

Code: In the below I have used the library called ggplot2 which is a system for declaratively creating graphics. I have specified the X label as crime type and colors blue and black.

```
45 # D)
46 # Plotting a bar graph of frequency of crime type and setting the axe labels and bar colours
47 library(ggplot2)
48 ggplot(AIINICrimeData, aes(x=Crime.type)) + geom_bar(fill = "blue", color = "black")
49
```

Fig.10:

Output: By observing the below bar plot we can say that anti-social behavior (ASBO) has the highest frequency of crimes comparing with the other crimes. While Theft from a person(THPR) and Robbery(ROBY) record lowest frequency in the crime rate. Violence and sexual offence(VISO) is the second highest in the crime type.

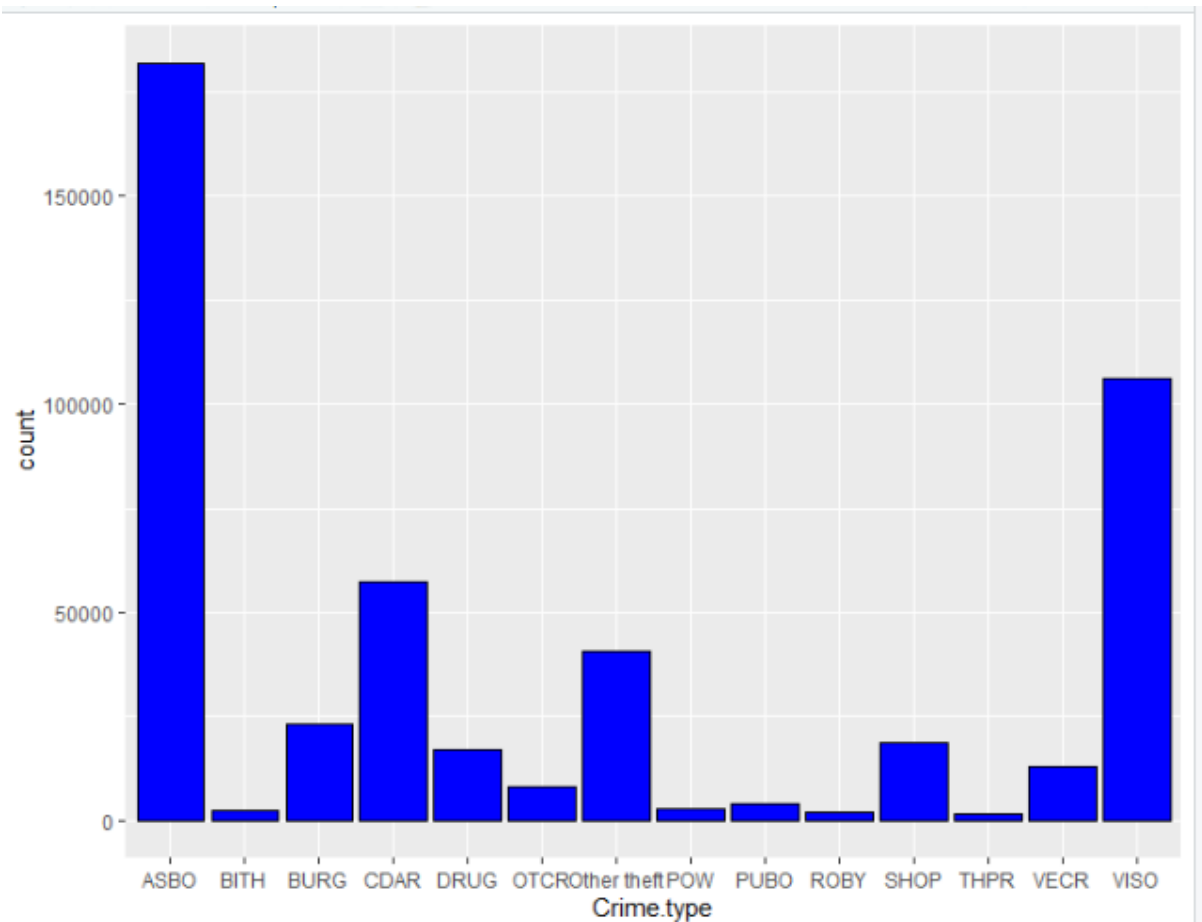


Fig.11: Bar plot for the crime data

- E) In this task we need to modify the AllNICrimeData such that the location attribute should contain only street name. we have observed that the attribute values are like on or near west rock square we need to remove the on or near and only assign street name to the attribute. And some attribute values are like no location.

	Month	Longitude	Latitude	Location	Crime.type
1	2015-01	-6.003289	54.55165	On or near Salisbury Place	ASBO
2	2015-01	-5.707979	54.59231	On or near	ASBO
3	2015-01	-5.815976	54.73161	On or near Milebush Park	ASBO
4	2015-01	-6.393411	54.19788	On or near College Square North	ASBO
5	2015-01	-6.251798	54.85970	On or near Staffa Drive	ASBO
6	2015-01	-7.206893	54.62265	On or near Killyclogher Road	ASBO
7	2015-01	-5.915793	54.59242	On or near Ravenhill Reach	ASBO
8	2015-01	-5.535389	54.48792	On or near	ASBO
9	2015-01	-7.322812	54.99940	On or near Great James Street	ASBO
10	2015-01	-5.954670	54.61568	On or near Jamaica Road	ASBO
11	2015-01	-6.707998	55.18703	On or near Station Road	ASBO
12	2015-01	-6.049408	54.51354	On or near Antrim Street	ASBO
13	2015-01	-5.709197	54.59706	On or near Westland Drive	ASBO

Showing 1 to 14 of 477,696 entries, 5 total columns

Fig.12: Data frame before changing the data

Code: we have used `str_remove_all` function by importing the library `stringr` and to remove the on or near in the location attribute and saved `AllNICrimeData`. For attribute value with no location we have assigned NA values as in the upcoming tasks we need to remove the NA values.

```
51 # E)
52 #install.packages("stringr")
53 library(stringr)
54 # Removing the extra strings in the location attribute and only keeping the street names
55 AllNICrimeData$Location <- str_remove_all(AllNICrimeData$Location, 'On or near ')
56 # Assigning NA values for no location attribute value
57 AllNICrimeData$Location[AllNICrimeData$Location == "no location" | AllNICrimeData$Location == ""] <- NA
```

Fig.13: Removing the NA values

Output:

	Month	Longitude	Latitude	Location	Crime.type
1	2015-01	-6.003289	54.55165	Salisbury Place	ASBO
2	2015-01	-5.707979	54.59231	NA	ASBO
3	2015-01	-5.815976	54.73161	Milebush Park	ASBO
4	2015-01	-6.393411	54.19788	College Square North	ASBO
5	2015-01	-6.251798	54.85970	Staffa Drive	ASBO
6	2015-01	-7.206893	54.62265	Killyclogher Road	ASBO
7	2015-01	-5.915793	54.59242	Ravenhill Reach	ASBO
8	2015-01	-5.535389	54.48792	NA	ASBO
9	2015-01	-7.322812	54.99940	Great James Street	ASBO
10	2015-01	-5.954670	54.61568	Jamaica Road	ASBO
11	2015-01	-6.707998	55.18703	Station Road	ASBO
12	2015-01	-6.049408	54.51354	Antrim Street	ASBO
13	2015-01	-5.709197	54.59706	Westland Drive	ASBO

Showing 1 to 14 of 477,696 entries, 5 total columns

Fig.14: Output after removing unwanted strings and assigning NA for no location

- F) This task assigns us to choose 5000 random samples from the crime data of AllNICrimeData dataset where location attribute should contain location name and there should be no NA values. We need to set seed value to 100 and create a new data frame called random_crime_sample. Now create a function called find_a_town which uses CleanNIPostCodeData dataset to find the correct town information for each location variable within the random_crime_sample dataset. Save each matched town into the random_crime_sample.

Code: In the first line of the code we have set the seed value to 100 to generate random number. After that by using subset function we have extracted the 5000 random samples and used !is.na function to ignore the NA values and saved in the variable random. Now save the extracted data into new dataset and name it as random_crime_sample.

```

61 # F)
62 # Choosing 5000 random samples of crime data from AllNICrimeData
63 # setting the seed value to 100
64 set.seed(100)
65 # Where location attribute should not contain any NA values
66 random <- subset(AllNICrimeData[, ], !is.na(AllNICrimeData$Location))
67 # saving the data into the new file random_crime_data
68 random_crime_sample <- random[sample(1:nrow(random), 5000, replace = FALSE), ]

```

Fig.15: code to extract 5000 random samples

Output:

	Month	Longitude	Latitude	Location	Crime.type
95810	2015-08	-5.933066	54.59109	Dublin Road	ASBO
384531	2017-06	-7.476470	54.78748	Daisy Park	ASBO
189080	2016-03	-5.935054	54.59857	King Street	ASBO
104404	2015-08	-5.963882	54.63147	Tyndale Gardens	VECR
241344	2016-07	-6.223332	54.48017	Meeting Street	ASBO
410624	2017-07	-5.839762	54.63121	Ardnagreena Gardens	VISO
103416	2015-08	-5.896897	54.21220	Shimna Mile	PUBO
206367	2016-04	-6.431276	54.43006	Grantham Park	CDAR
133178	2015-10	-5.930067	54.60018	Rosemary Street	SHOP
241430	2016-07	-6.739900	54.63720	Stewart Avenue	ASBO
219424	2016-05	-7.281825	54.60935	Knockshee Park	CDAR
97330	2015-08	-6.278889	54.86499	George Street	ASBO
146854	2015-11	-6.457149	54.43173	Ballyoran Heights	VISO
300797	2016-11	-6.338295	54.45736	Sloan Street	BURG
448064	2017-10	-5.970893	54.62690	Silverstream Parade	Other theft
314169	2016-12	-6.281757	54.85361	Montague Park	DRUG
265748	2016-08	-6.275444	54.87672	Kew Gardens	VISO
447649	2017-10	-6.213340	54.59704	Gobrana Road	DRUG
327654	2017-01	-6.068101	54.52843	The Oaks	Other theft
395384	2017-06	-5.948576	54.68843	Reverdev Road	VISO

Showing 1 to 21 of 5,000 entries, 5 total columns

Fig.16: Output of 5000 random samples

Now change the string to lower case by using `str_to_lower` of both the data sets `NIPOSTCODE` and `random_crime_sample` in the location and primary thorfare attribute respectively to match the attribute values.

```
70 # changing the string to lower case to match the attribute
71 NIPOSTCODE$'Primary Thorfare' <- str_to_lower(NIPOSTCODE$'Primary Thorfare')
72 random_crime_sample$Location <- str_to_lower(random_crime_sample$Location)
```

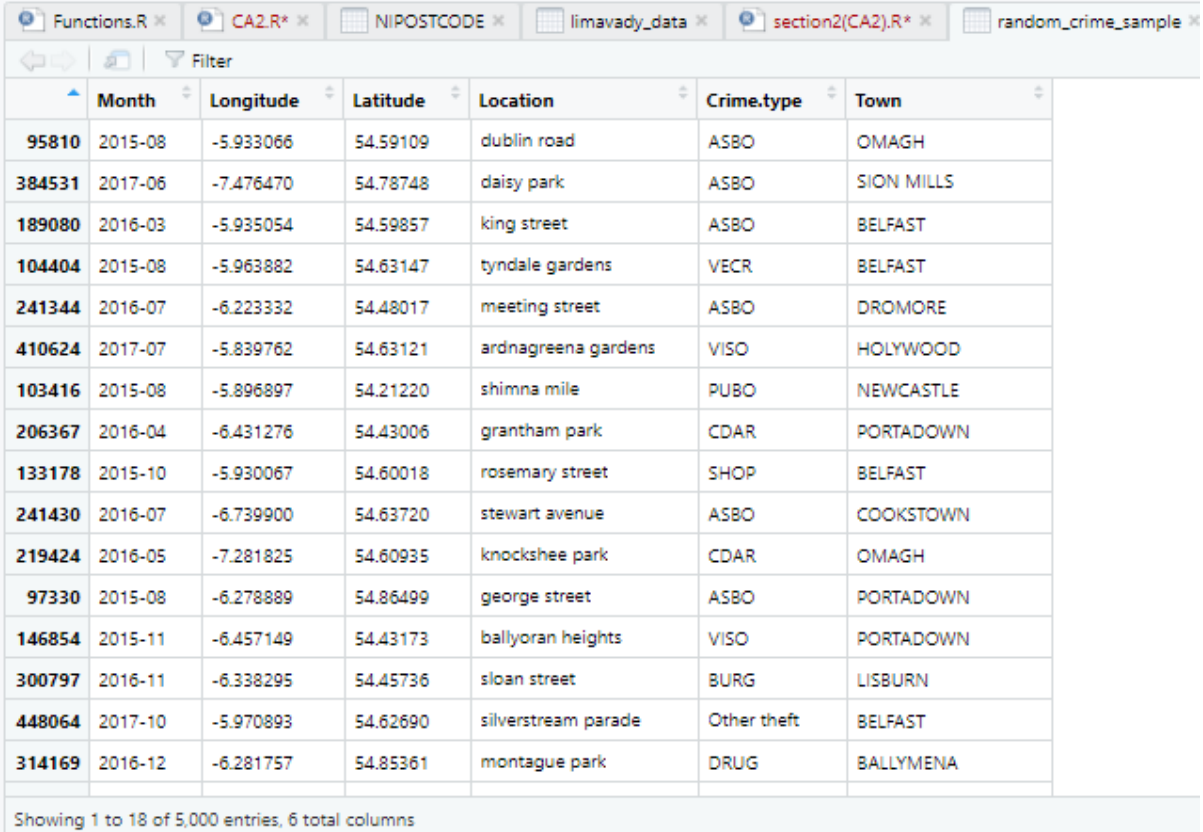
Fig.17: code for changing to lower case

After changing to lower now create a function `find_a_town` and return the match function to compare the location of two datasets. Store it into the dataset `random_crime_sample` and check the structure of the data frame

```
74 # creating a function to match to find correct town/city information for each location variable within the
75 find_a_town <- function(x,y){
76   return(NIPOSTCODE$Town[match(x,y)])
77 }
78 # Saving the match attributes into the town column of the random_crime_sample
79 random_crime_sample$Town <- find_a_town(random_crime_sample$Location, NIPOSTCODE$'Primary Thorfare')
81 str(random_crime_sample)
```

Fig.18: creating a function

Output:



	Month	Longitude	Latitude	Location	Crime.type	Town
95810	2015-08	-5.933066	54.59109	dublin road	ASBO	OMAGH
384531	2017-06	-7.476470	54.78748	daisy park	ASBO	SION MILLS
189080	2016-03	-5.935054	54.59857	king street	ASBO	BELFAST
104404	2015-08	-5.963882	54.63147	tyndale gardens	VECR	BELFAST
241344	2016-07	-6.223332	54.48017	meeting street	ASBO	DROMORE
410624	2017-07	-5.839762	54.63121	ardnagreena gardens	VISO	HOLYWOOD
103416	2015-08	-5.896897	54.21220	shimna mile	PUBO	NEWCASTLE
206367	2016-04	-6.431276	54.43006	grantham park	CDAR	PORTADOWN
133178	2015-10	-5.930067	54.60018	rosemary street	SHOP	BELFAST
241430	2016-07	-6.739900	54.63720	stewart avenue	ASBO	COOKSTOWN
219424	2016-05	-7.281825	54.60935	knockshee park	CDAR	OMAGH
97330	2015-08	-6.278889	54.86499	george street	ASBO	PORTADOWN
146854	2015-11	-6.457149	54.43173	ballyoran heights	VISO	PORTADOWN
300797	2016-11	-6.338295	54.45736	sloan street	BURG	LISBURN
448064	2017-10	-5.970893	54.62690	silverstream parade	Other theft	BELFAST
314169	2016-12	-6.281757	54.85361	montague park	DRUG	BALLYMENA

Showing 1 to 18 of 5,000 entries, 6 total columns

Fig.19: Output after the match and adding the town

```
> str(random_crime_sample)
'data.frame': 5000 obs. of 6 variables:
 $ Month      : Factor w/ 36 levels "2015-01","2015-02",...: 8 30 15 8 19 31 8 16 10 19 ...
 $ Longitude  : num -5.93 -7.48 -5.94 -5.96 -6.22 ...
 $ Latitude   : num 54.6 54.8 54.6 54.6 54.5 ...
 $ Location   : chr "dublin road" "daisy park" "king street" "tyndale gardens" ...
 $ Crime.type : chr "ASBO" "ASBO" "ASBO" "VECR" ...
 $ Town       : Factor w/ 314 levels "", "AGHAGALLON",...: 260 287 46 46 127 177 251 265 46 94 ...
```

Fig.20: structure of the random_crime_sample

- G) Now we need to create a function called `add_town_data` which examines each file of crime data in the `random_crime_sample` and matches the relevant information in the `village_list` dataset.

	A	B	C	D	E	F	G	H	I	J	K	L
1	CITY/TOWN/VILLAGE	POPULATION	STATUS	COUNTY								
2	Belfast	3,35,665	City	County Antrim								
3	Armagh	1,75,621	Town	County Armagh								
4	Lisburn	1,21,654	Town	County Antrim								
5	Newtown	1,18,261	Town	County Antrim								
6	Derry	87,269	City	County Londonderry								
7	Killeel	69,258	Town	County Down								
8	Bangor	62,789	Town	County Down								
9	Larne	33,620	Town	County Antrim								
10	Ballymena	30,265	Town	County Antrim								
11	Newtown	28,654	Town	County Down								
12	Carrickfergus	28,653	Town	County Antrim								
13	Lurgan	28,568	Town	County Armagh								
14	Newry	27,234	City	County Down								
15	Coleraine	24,694	Town	County Londonderry								
16	Acton	22,654	Town	County Armagh								
17	Portadown	22,611	Town	County Armagh								
18	Antrim	21,896	Town	County Antrim								
19	Omagh	21,325	Town	County Tyrone								
20	Ballyclare	18,675	Town	County Antrim								

Fig.21: Village_list.csv file

Code:

In the first line of the code we have read the village_list.csv file. Now we have changes the city/town information to upper case by using str_to_upper. From the above figure we can see that the city Londonderry is names as Derry as it will be unmatched with the random_crime_sample. So we have replaced Derry with Londonderry by using the function str_replace_all.

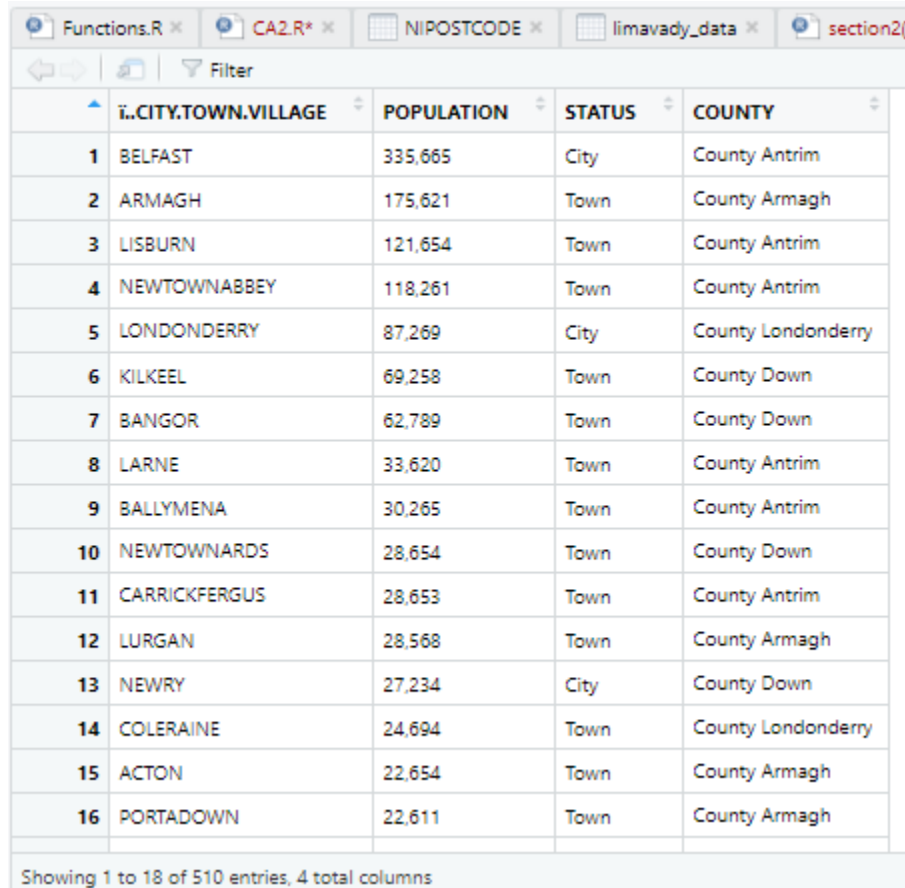
```

87 # Reading village data
88 village_data <- read.csv("VillageList.csv")
89 # Changing the string into upper case
90 village_data$CITY.TOWN.VILLAGE <- str_to_upper(village_data$CITY.TOWN.VILLAGE)
91 village_data$CITY.TOWN.VILLAGE <- str_replace_all(village_data$CITY.TOWN.VILLAGE, "DERRY", "LONDONDERRY")

```

Fig.22: changing the strings

Output:



	i..CITY.TOWN.VILLAGE	POPULATION	STATUS	COUNTY
1	BELFAST	335,665	City	County Antrim
2	ARMAGH	175,621	Town	County Armagh
3	LISBURN	121,654	Town	County Antrim
4	NEWTOWNABBEY	118,261	Town	County Antrim
5	LONDONDERY	87,269	City	County Londonderry
6	KILKEEL	69,258	Town	County Down
7	BANGOR	62,789	Town	County Down
8	LARNE	33,620	Town	County Antrim
9	BALLYMENA	30,265	Town	County Antrim
10	NEWTOWNARDS	28,654	Town	County Down
11	CARRICKFERGUS	28,653	Town	County Antrim
12	LURGAN	28,568	Town	County Armagh
13	NEWRY	27,234	City	County Down
14	COLERAINE	24,694	Town	County Londonderry
15	ACTON	22,654	Town	County Armagh
16	PORTADOWN	22,611	Town	County Armagh

Showing 1 to 18 of 510 entries, 4 total columns

Fig.23: change in string values to upper case and change Derry to Londonderry

Now we need to create a function `add_town_data` and match the town attribute within the `random_crime_sample` to the city/town/village in the village dataset. And store the matched population data from the two datasets into the population attribute in the `random_crime_sample`.

```

91 # Creating a function to match the each crime data to relevant information in village list
92 add_town_data <- function(x,y){
93
94   return(village_data$POPULATION[match(x,y)])
95 }
96 # storing it into the random_crime_sample
97 random_crime_sample$Population <- add_town_data(random_crime_sample$Town, village_data$i..CITY.TOWN.VILLAGE)

```

Fig.24: creating a function add_town_data

Output:

	Month	Longitude	Latitude	Location	Crime.type	Town	Population
95810	2015-08	-5.933066	54.59109	dublin road	ASBO	OMAGH	21,325
384531	2017-06	-7.476470	54.78748	daisy park	ASBO	SION MILLS	2,064
189080	2016-03	-5.935054	54.59857	king street	ASBO	BELFAST	335,665
104404	2015-08	-5.963882	54.63147	tyndale gardens	VECR	BELFAST	335,665
241344	2016-07	-6.223332	54.48017	meeting street	ASBO	DROMORE	5,254
410624	2017-07	-5.839762	54.63121	ardnagreena gardens	VISO	HOLYWOOD	12,657
103416	2015-08	-5.896897	54.21220	shimna mile	PUBO	NEWCASTLE	7,524
206367	2016-04	-6.431276	54.43006	grantham park	CDAR	PORTADOWN	22,611
133178	2015-10	-5.930067	54.60018	rosemary street	SHOP	BELFAST	335,665
241430	2016-07	-6.739900	54.63720	stewart avenue	ASBO	COOKSTOWN	10,718
219424	2016-05	-7.281825	54.60935	knockshee park	CDAR	OMAGH	21,325
97330	2015-08	-6.278889	54.86499	george street	ASBO	PORTADOWN	22,611
146854	2015-11	-6.457149	54.43173	ballyoran heights	VISO	PORTADOWN	22,611
300797	2016-11	-6.338295	54.45736	sloan street	BURG	LISBURN	121,654
448064	2017-10	-5.970893	54.62690	silverstream parade	Other theft	BELFAST	335,665
314169	2016-12	-6.281757	54.85361	montague park	DRUG	BALLYMENA	30,265

Showing 1 to 18 of 5,000 entries, 7 total columns

Fig.25: output after adding the population attribute

```
> str(random_crime_sample)
'data.frame': 5000 obs. of 7 variables:
 $ Month      : Factor w/ 36 levels "2015-01","2015-02",...: 8 30 15 8 19 31 8 16 10 19 ...
 $ Longitude  : num -5.93 -7.48 -5.94 -5.96 -6.22 ...
 $ Latitude   : num 54.6 54.8 54.6 54.6 54.5 ...
 $ Location   : chr "dublin road" "daisy park" "king street" "tyndale gardens" ...
 $ Crime.type : chr "ASBO" "ASBO" "ASBO" "VECR" ...
 $ Town       : Factor w/ 314 levels "", "AGHAGALLON",...: 260 287 46 46 127 177 251 265 46 94 ...
 $ Population : Factor w/ 387 levels "1,021","1,062",...: 167 127 232 232 275 70 333 174 232 52 ...
```

Fig.26: structure of random_crime_sample

H) Now we need to update the random_crime_sample so that it only contains the attributes Month, Longitude, Latitude, Location, Crime type, City-Town-Village, Population.

```
100 # changing the column name of town to city-town-village
101 colnames(random_crime_sample)[6] <- c("City-Town-Village")
```

Fig.27: changing the column name

As we can see in the above structure the town attribute is in the 6th index position we need to change the name to city-town-village by creating the attribute name and change in the random_crime_sample.

Output:

	Month	Longitude	Latitude	Location	Crime.type	City-Town-Village	Population
95810	2015-08	-5.933066	54.59109	dublin road	ASBO	OMAGH	21,325
384531	2017-06	-7.476470	54.78748	daisy park	ASBO	SION MILLS	2,064
189080	2016-03	-5.935054	54.59857	king street	ASBO	BELFAST	335,665
104404	2015-08	-5.963882	54.63147	tyndale gardens	VECR	BELFAST	335,665
241344	2016-07	-6.223332	54.48017	meeting street	ASBO	DROMORE	5,254
410624	2017-07	-5.839762	54.63121	ardnagreena gardens	VISO	HOLYWOOD	12,657
103416	2015-08	-5.896897	54.21220	shimna mile	PUBO	NEWCASTLE	7,524
206367	2016-04	-6.431276	54.43006	grantham park	CDAR	PORTADOWN	22,611
133178	2015-10	-5.930067	54.60018	rosemary street	SHOP	BELFAST	335,665
241430	2016-07	-6.739900	54.63720	stewart avenue	ASBO	COOKSTOWN	10,718
219424	2016-05	-7.281825	54.60935	knockshee park	CDAR	OMAGH	21,325
97330	2015-08	-6.278889	54.86499	george street	ASBO	PORTADOWN	22,611
146854	2015-11	-6.457149	54.43173	ballyoran heights	VISO	PORTADOWN	22,611
300797	2016-11	-6.338295	54.45736	sloan street	BURG	LISBURN	121,654
448064	2017-10	-5.970893	54.62690	silverstream parade	Other theft	BELFAST	335,665
314169	2016-12	-6.281757	54.85361	montague park	DRUG	BALLYMENNA	30,265

Showing 1 to 17 of 5,000 entries, 7 total columns

Fig.28: Output after changing the attribute name

- I) We need to plot the crimes in the cities Belfast and Londonderry. We need to show the plot data side by side.

Code: first we need create to variables and use the subset function to group the crimes in the cities Belfast and Londonderry. And use the libraries dplyr and ggplot for graphics. In the code we can see that we have grouped by crime type and summarise the frequency counts and change the label names.

```

106 # I)
107 # Plotting the crimes in cities Belfast and Londonderry
108 random_crime_belfast <- subset(random_crime_sample, random_crime_sample$'City-Town-Village' == "BELFAST")
109 random_crime_derry <- subset(random_crime_sample, random_crime_sample$'City-Town-Village' == "LONDONDERRY")
110
111
112 library(dplyr)
113 library(ggplot2)
114 # Step 1
115 plot1 <- random_crime_belfast %>%
116   #Step 2
117   group_by(Crime.type) %>%
118   #Step 3
119   summarise(count_fre = n()) %>%
120   #Step 4
121   ggplot(aes(x = Crime.type, y = sort(count_fre, decreasing = TRUE), fill = Crime.type)) +
122   geom_bar(stat = "identity") +
123   theme_classic() +
124   labs(
125     x = "Crime.type",
126     y = "Count",
127     title = paste(
128       "Crimes in Belfast"
129     )
130 )

```

Fig.29: code for plotting the crimes in Belfast

```

131 plot2 <- random_crime_derry %>%
132   #Step 2
133   group_by(Crime.type) %>%
134   #Step 3
135   summarise(count_fre = n()) %>%
136   #Step 4
137   ggplot(aes(x = Crime.type, y = sort(count_fre, decreasing = TRUE), fill = Crime.type)) +
138   geom_bar(stat = "identity") +
139   theme_classic() +
140   labs(
141     x = "Crime.type",
142     y = "Count",
143     title = paste(
144       "Crimes in Londonderry"
145     )
146   )
147
148 install.packages("gridExtra")
149 require(gridExtra)
150 grid.arrange(plot1,plot2,ncol=2)

```

Fig.30: code for plotting the crimes in Londonderry

Output: From the below output we can see that anti-social behavior is the highest in the both cities. when compared to Bicycle theft Londonderry has slightly highest frequency than Belfast. Violence and sexual offence is the least in both the cities. Burglary and criminal damage and arson are vice versa in both the cities.

