

# COMP5721M: Programming for Data Science

## Coursework 3: Data Analysis Project

### UK Accidents and Casualties

- Kriti Garg, [mm22kg@leeds.ac.uk](mailto:mm22kg@leeds.ac.uk) (<mailto:mm22kg@leeds.ac.uk>)
- Raghav Arora, [mm22ra@leeds.ac.uk](mailto:mm22ra@leeds.ac.uk) (<mailto:mm22ra@leeds.ac.uk>)
- Baha Abdul Vahab, [mm22bura@leeds.ac.uk](mailto:mm22bura@leeds.ac.uk) (<mailto:mm22bura@leeds.ac.uk>)
- Kartikeya Malimath, [sc22kpm@leeds.ac.uk](mailto:sc22kpm@leeds.ac.uk) (<mailto:sc22kpm@leeds.ac.uk>)

## Project Plan

### The Data

Anything which has a merit, its not necessary that there can't be any demerit of it. For eg: In case of Vehicles, Pedestrian Crossing, Road Safety seems like shouldn't have any problems because they have there existence to save one's time, provide services and safety. That is why we have selected UK Accidents and Casualties 2011 and 2012 data to analyse what problems are faced at the massive level and at what regions. We are analysing 2 datasets, First one is the main set which contains all the information about the location and number of accidents, casualties and the surrounding during the accidents, Second One is the mapping Excel which comprises of the keys and value for the later columns, For eg: for Accident Severity we have keys 1, 2 and 3 and in the mapping Excel we have values for it i.e. Fatal, Serious and Slight.

We have collected the data from [Kaggle \(https://www.kaggle.com/datasets/devansodariya/road-accident-united-kingdom-uk-dataset?select=UK\\_Accident.csv\)](https://www.kaggle.com/datasets/devansodariya/road-accident-united-kingdom-uk-dataset?select=UK_Accident.csv) and [data.gov.uk \(https://www.data.gov.uk/dataset/cb7ae6f0-4be6-4935-9277-47e5ce24a11f/road-safety-data\)](https://www.data.gov.uk/dataset/cb7ae6f0-4be6-4935-9277-47e5ce24a11f/road-safety-data) which is accumulated by the UK government from the traffic data. Unexpectedly, the data is almost clean and looks good to work upon directly.

The below table describes the key columns of our first dataset and their meaning:

Column Name	Description	Data Type
Accident_Index	unique identifier of an accident(combination of year of accident and accident reference)	object
Longitude	Longitude of the place where accident take place	float64
Latitude	Latitude of the place where accident take place	float64
Police_Force	The Police on whose shoulders the responsibility for investigating the accident is	int64(1-98)
Accident_Severity	The Severity of the accident (Fatal, Serious, Slight)	int64(1,2,3)
Number_of_Vehicles	Total Number of Vehicles involved in the accident	int64
Number_of_Casualties	Total Number of casualties due to the accident	int64
Date	The date of the accident	object
Day_of_Week	the day of the week(1-7{Monday-Sunday})	int64
Local_Authority_(District)	Under which district the accident take place	int64
Local_Authority_(Highway)	Under which Highway the accident take place	object

Column Name	Description	Data Type
Road_Type	Type of the road i.e. singleway or dualway etc	object
Pedestrian_Crossing-Human_Control	By whom the Pedestrian crossing was controlled or by none	object
Pedestrian_Crossing-Physical_Facilities	By which the Crossing was controlled such as has zebra crossing, had pedestrian junction at traffic signal etc	object
Light_Conditions	Was it dark or was during day light	object
Weather_Conditions	Conditions such as if there was fog, or was fine, was raining or snowing etc	object
Road_Surface_Conditions	Condition of the road such as dry, flooded, wet, Normal, Frost/Ice, Snow etc	object
Special_Conditions_at_Site	any Special conditions such as Mud, oil or diesel, Construction going on etc	object
Carriageway_Hazards	Anything or person that that can be a reason for the accident i.e. pedestrian or an animal etc	object
Urban_or_Rural_Area	Contains 1,2,3 - will be mapped with mapping Excel to show if it is Urban, Rural or Unclassified	object
Did_Police_Officer_Attend_Scene_of_Accident	Contains Yes or No	object
Year	Year in which accident happen, in our case it will be 2011 or 2012	int64

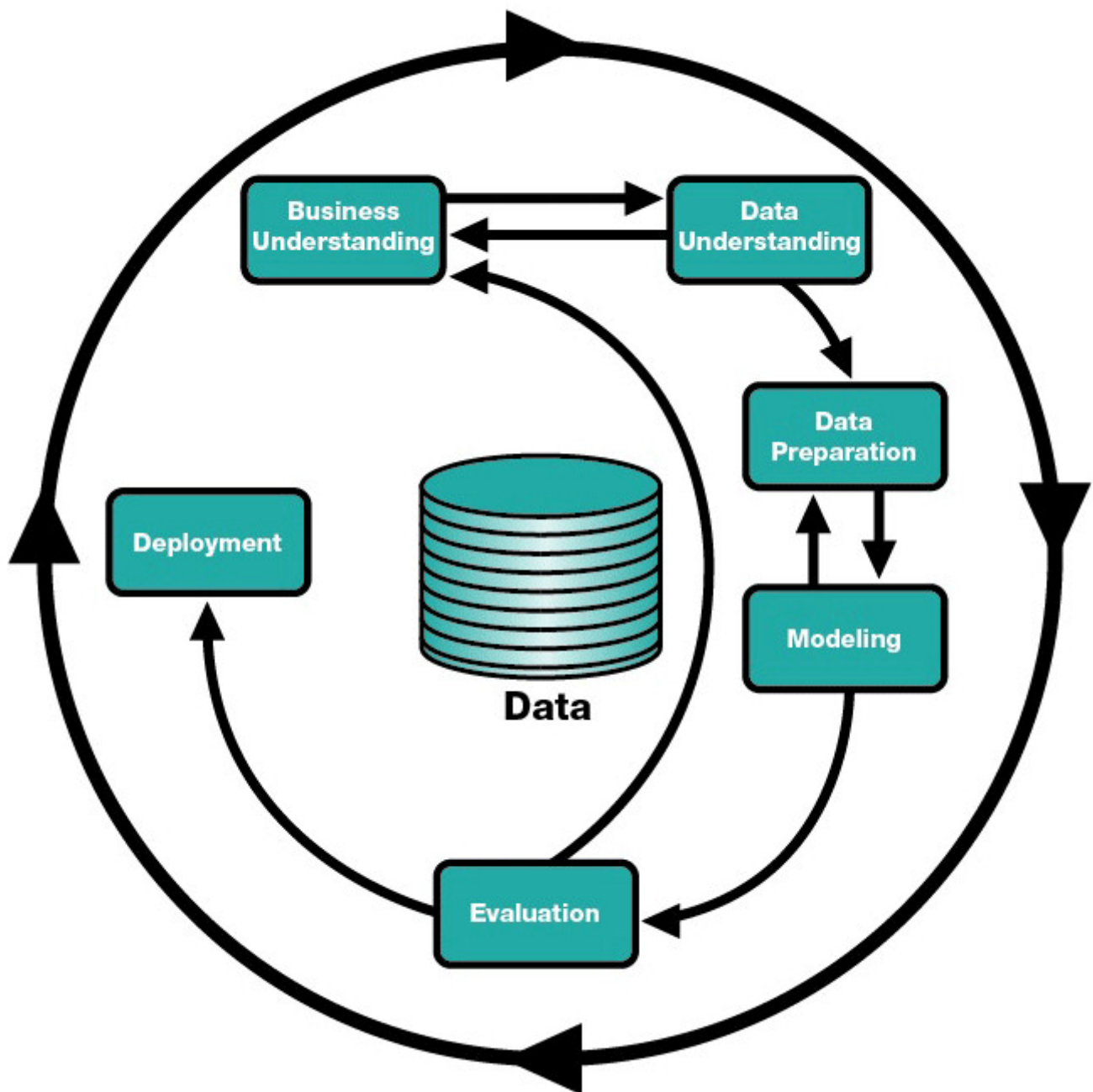
Below are the key columns of mapping Excel:

Column Name	Description	Format	Label
Field Name	Name of the column present in First Dataset	Key	Value

## Project Aim and Objectives

The major aim of our project is to identify what changes should be brought by an individual and by the govt so that the number of accidents can be decreased. The highways, carriageways, police, local authorities, Pedestrian Crossing are all made for the safety of the people and if in case of failure in 1 such thing there are accidents happening, those reasons should be addressed so that the proper measures can be taken. By this Road Accidents data we are trying to address the following:

- major reasons behind the fatal/severe accidents i.e. if they are due to human error or lack of physical facilities.
- the top cities which have the maximum number of accidents and the top cities that are safest to live i.e. which have lowest number of accidents happened.



- *Classification: separate data items into classes according to their characteristics (can be either a definite or a statistical kind of classification)*
- *Correlation: find correspondences between different attributes within a dataset*
- *Search: find solutions matching some criteria*
- *Visualisation: find informative ways to display the structure of a large and/or complex dataset*
- *Query Answering: create a system that enables one to retrieve information by evaluating some form of query representation*
- *Simulation: model the evolution of a complex process Here you should describe the general aim of your project in around 200-300 words.*

*This can be anything from classifying items according to their characteristic features (which mushrooms are poisonous?) to simulating an evolving process (will the rabbits eat all the carrots or get eaten by the foxes?)*

*Here some ideas of general types of processing functionality that you could implement:*

- *Classification: separate data items into classes according to their characteristics (can be either a definite or a statistical kind of classification)*
- *Correlation: find correspondences between different attributes within a dataset*
- *Search: find solutions matching some criteria*

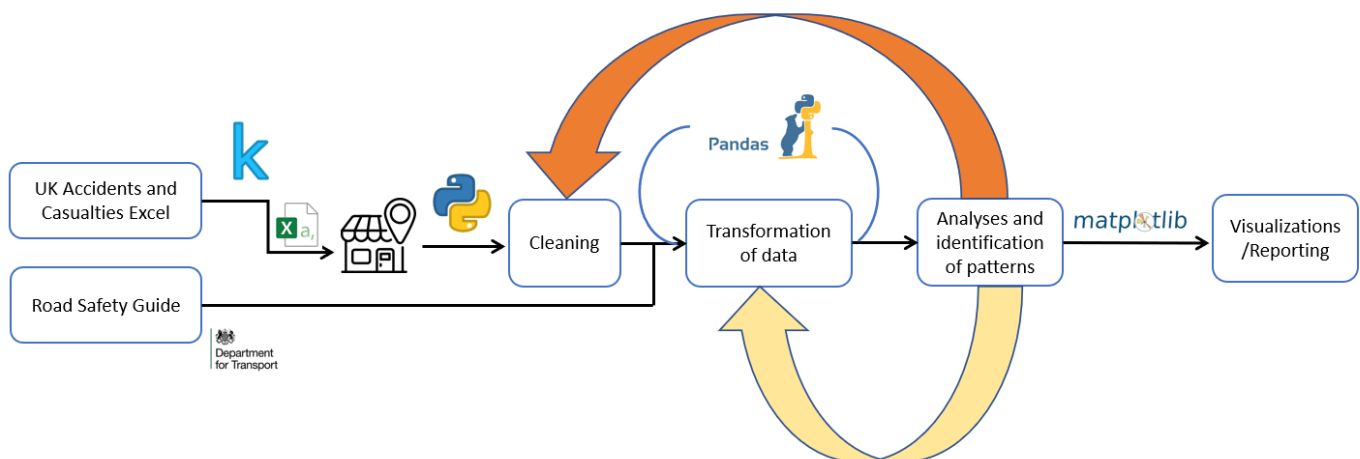
- *Visualisation*: find informative ways to display the structure of a large and/or complex dataset
- *Query Answering*: create a system that enables one to retrieve information by evaluating some form of query representation
- *Simulation*: model the evolution of a complex process

## Specific Objective(s)

- **Objective 1**: To identify the highest number of casualties in each year and their potential reasons and how the police responds at the scene of accidents?
- **Objective 2**: What are the top 10 cities that had the maximum number of casualties?
- **Objective 3**: To identify the relationship between the reasons that are affecting the accidents and the number of casualties over the year and their trend. Additionally we are trying to understand the police response as per the reason
- **Objective 4**: What is the percentage change of the number of accidents corresponding to the cities and the major physical reasons such as Road Type, urban or rural?

## System Design

### Architecture



At first comes the datasets which are required for analysis i.e. the UK Accidents/Casualties and Road Safety Data Guide(Mapping Excel) which then are combined on the basis of the FieldName and columns to populate the true value of the keys provided. Once the data is merged using left join from Accidents to Mapping Excel the data is cleaned i.e. the unimportant columns are dropped. Transformation of data is done by updating the datatypes or, filling the missing values or, filtering the 2 years of data. After this Analysis is done to identify the correlation between the features and if there is any interesting patterns among them. As we move forward with the analysis there are times when we have to transform the data again or even perform the cleaning task again. Visualizations are created once we feel the analysis seems appropriate.

## Processing Modules and Algorithms

- **Cleaning:**
  - Removed the following columns as were not required because there was lack of explanation.
    - Location\_Easting\_OSGR , Location\_Northing\_OSGR , 1st\_Road\_Class , 1st\_Road\_Number , Junction\_Control , 2nd\_Road\_Class , 2nd\_Road\_Number , LSOA\_of\_Accident\_Location , Time , Unnamed(0) ( This was coming because of the first row as a rownumber in Excel)

```
In [80]: data.head()
```

```
Out[80]:
```

	Unnamed: 0	Accident_Index	Location_Easting_OSGR	Location_Northing_OSGR	Longitude	Latitude	Police_Force	Accident_Severity	Number_of_Vehicles	N
0	0	200501BS00001	525680.0	178240.0	-0.191170	51.489096	1	2	1	
1	1	200501BS00002	524170.0	181650.0	-0.211708	51.520075	1	3	1	

- **Joining of Datasets:**

- We have only 1 main dataset i.e. UK Road Accidents but for some columns such as **Accident Severity, Urban\_or\_Rural\_Area, Local\_Authority\_(District)**, we don't have the mapping of the values for eg: what does 1 as severity means or 2 as Urban\_or\_Rural Means. So we Downloaded a mapping Excel from data.gov.uk site and added the mapping for these columns.

- **Transformation:**

- Updated the data types of Date column from object to Datetime64

- **Analysis and Visualization:**

- Perform Analysis to identify the regions with massive number of casualties, what combinations of reasons are the main cause of accidents taking place, Accidents of what severity happened at what regions, Percentage of Accidents happened over the Fatal Accidents and their plot their Visualizations.

## Program Code

Importing Libraries which we will be using in this project

In [1]:

```
1  ## Libraries
2  import pandas as pd
3  import matplotlib.pyplot as plt
4  import datetime
5  import numpy as np
6  import seaborn as sns
7  import decimal
8  from matplotlib.patches import ConnectionPatch
9  from matplotlib.gridspec import GridSpec
10 import calendar
11 import matplotlib.patches as mpatches
```

Reading first dataset 'UK\_Accident.csv' in a DataFrame using pandas

In [2]:

```
1 uk_acc_all_data = pd.read_csv("UK_Accident.csv")
```

The following cell performs cleaning of the dataset by removing the unwanted columns

In [3]:

```
1 uk_acc_all_data.drop(columns=['LSOA_of_Accident_Location',
2                               'Location_Easting_OSGR',
3                               'Location_Northing_OSGR',
4                               '1st_Road_Class',
5                               '1st_Road_Number',
6                               'Junction_Control',
7                               '2nd_Road_Class',
8                               '2nd_Road_Number',
9                               'Time',
10                              'Unnamed: 0'], inplace = True)
```

Filtering out the data for 2011 and 2012 in which we want to do analysis

In [4]:

```
1 uk_filter_data = uk_acc_all_data[uk_acc_all_data["Year"].isin([2011,2012])]
```

For our Analysis and Visualizations we require 3 new columns and these will be added to the uk\_filter\_data DataFrame

- Month - Extracting from Date column and converting the argument to numeric type
- Reason - Combining the values of Road\_Surface\_Conditions , Light\_Conditions and Weather\_Conditions
- First\_day\_of\_month - Extracting the month/year and creating generating the first day of every month and year

In [5]:

```
1 uk_filter_data["month"] = pd.to_numeric(uk_filter_data["Date"].str[3:5])
2 uk_filter_data["reason"] = (uk_filter_data["Road_Surface_Conditions"]+"+"+uk_filter_data["Light_Conditions"]+"+"+uk_filter_data["Weather_Conditions"])
3 uk_filter_data["First_day_of_month"] = pd.to_datetime(uk_filter_data["month"].astype(str) + "01" + uk_filter_data["Year"].astype(str))
```

The below command shows the sum of all the null values if any

In [6]:

```
1 uk_filter_data.isnull().sum()
```

Out[6]:

```
Accident_Index          0
Longitude               0
Latitude               0
Police_Force            0
Accident_Severity       0
Number_of_Vehicles      0
Number_of_Casualties     0
Date                   0
Day_of_Week             0
Local_Authority_(District)  0
Local_Authority_(Highway)  0
Road_Type               0
Speed_limit             0
Pedestrian_Crossing-Human_Control  0
Pedestrian_Crossing-Physical_Facilities  0
Light_Conditions        0
Weather_Conditions      0
Road_Surface_Conditions  0
Special_Conditions_at_Site  0
Carriageway_Hazards     0
Urban_or_Rural_Area     0
Did_Police_Officer_Attend_Scene_of_Accident  0
Year                   0
month                  0
reason                 0
First_day_of_month     0
dtype: int64
```

From the above command's result we can see that there is no null value coming in data now.

Reading second dataset into a DataFrame using pandas read\_excel function

In [7]:

```
1 #importing a new dataset- dim table
2 data_guide = pd.read_excel("Road-Safety-Open-Dataset-Data-Guide.xlsx")
```

Creating 3 different Dataframes for 3 separate columns for whose key's, the values needs to be populated. The below cell is mapping values for local\_authority\_district column

In [8]:

```
1 df_lad_mapping = data_guide[data_guide["field name"] == 'local_authority_district'][["code/format", "label"]]
2 .rename(columns = {"code/format": "Local_Authority_(District)", "label": "label"})
```

The below cell is mapping values for urban\_or\_rural\_area column

In [9]:

```
1 df_uor_mapping = data_guide[data_guide["field name"] == 'urban_or_rural_area'][["code/format", "label"]].\n2   .rename(columns = {"code/format": "Urban_or_Rural_Area", "label": "Urban_or_Rural_Area_Label"})
```

The below cell is mapping values for casualty\_severity column

In [10]:

```
1 df_sev_mapping = data_guide[data_guide["field name"] == 'casualty_severity'][["code/format", "label"]].\n2   .rename(columns = {"code/format": "Accident_Severity", "label": "Accident_Severity_Label"})
```

Left Joining the above 3 DataFrame's to the uk\_filter\_data insert the values of the column in the DataFrame

In [11]:

```
1 left1 = pd.merge(uk_filter_data, df_lad_mapping, on = 'Local_Authority_(District)', how = 'left')\n2 left2 = pd.merge(left1, df_uor_mapping, on = 'Urban_or_Rural_Area', how = 'left')\n3 left3 = pd.merge(left2, df_sev_mapping, on = 'Accident_Severity', how = 'left')\n4 uk_filter_data = left3
```

Checking is the data types of all the columns are as expected.



In [12]:

```
1 uk_filter_data.info()
```

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 331189 entries, 0 to 331188
Data columns (total 29 columns):
#   Column                                                                 Non-Null Count  Dtype
---  -
0   Accident_Index                                                         331189 non-null  object
1   Longitude                                                             331189 non-null  float64
2   Latitude                                                              331189 non-null  float64
3   Police_Force                                                           331189 non-null  int64
4   Accident_Severity                                                      331189 non-null  object
5   Number_of_Vehicles                                                    331189 non-null  int64
6   Number_of_Casualties                                                  331189 non-null  int64
7   Date                                                                  331189 non-null  object
8   Day_of_Week                                                            331189 non-null  int64
9   Local_Authority_(District)                                           331189 non-null  object
10  Local_Authority_(Highway)                                             331189 non-null  object
11  Road_Type                                                             331189 non-null  object
12  Speed_limit                                                            331189 non-null  int64
13  Pedestrian_Crossing-Human_Control                                     331189 non-null  object
14  Pedestrian_Crossing-Physical_Facilities                             331189 non-null  object
15  Light_Conditions                                                       331189 non-null  object
16  Weather_Conditions                                                    331189 non-null  object
17  Road_Surface_Conditions                                               331189 non-null  object
18  Special_Conditions_at_Site                                           331189 non-null  object
19  Carriageway_Hazards                                                   331189 non-null  object
20  Urban_or_Rural_Area                                                   331189 non-null  object
21  Did_Police_Officer_Attend_Scene_of_Accident                       331189 non-null  object
22  Year                                                                  331189 non-null  int64
23  month                                                                331189 non-null  int64
24  reason                                                                331189 non-null  object
25  First_day_of_month                                                    331189 non-null  datetime64
4[ns]
26  Local_Authority_District_value                                         331189 non-null  object
27  Urban_or_Rural_Area_value                                              331189 non-null  object
28  Accident_Severity_value                                                331189 non-null  object
dtypes: datetime64[ns](1), float64(2), int64(7), object(19)
memory usage: 75.8+ MB
```

As can be seen that Date and Severity is of Object dtype. Updating the data type of Date to datetime and Accident Severity to int

In [13]:

```
1 uk_filter_data["Date"] = pd.to_datetime(uk_filter_data["Date"])
2 uk_filter_data["Accident_Severity"] = pd.to_numeric(uk_filter_data["Accident_Severity"])
```

...

## Objective 1:

Creating a DataFrame df, grouping up on the basis of the new column First\_day\_of\_month created and aggregating on Number\_of\_Casualties

In [15]:

```
1 df = uk_filter_data.groupby(["First_day_of_month"])["Number_of_Casualties"].agg(['sum',
```

```
1 Below cell creates another Dataframe defining the `First_day_of_month` along with  
the casualties corresponding to the `Reason`
```

In [16]:

```
1 uk_bar_graph_data_natural = uk_filter_data\  
2     .groupby(["First_day_of_month", "reason"])["Number_of_Casualties"]\  
3     .sum().reset_index()
```

Fetching the top 5 reasons resulting in highest number of casualties for 2011 and 2012 years separately

In [17]:

```
1 # taking tops reasons to plot the data in pie chart of 2011  
2 df_2011 = uk_bar_graph_data_natural[uk_bar_graph_data_natural["First_day_of_month"]=="2011"]\  
3     .sort_values("Number_of_Casualties", ascending=False).head(5)  
4  
5 # taking tops reasons to plot the data in pie chart of 2012  
6 df_2012 = uk_bar_graph_data_natural[uk_bar_graph_data_natural["First_day_of_month"]=="2012"]\  
7     .sort_values("Number_of_Casualties", ascending=False).head(5)
```

Below command defines a DataFrame df\_bar\_police that tell us did the police attended the accident or not.

In [18]:

```
1 df_bar_police = uk_filter_data.groupby(["First_day_of_month", "Did_Police_Officer_Attend_Accident"]
```

Below function converts the number to percentages

In [19]:

```
1 def func(pct, allvals):  
2     absolute = int(pct/100.*np.sum(allvals))  
3     return "{:.1f}%\n({:d})".format(pct, absolute)
```

The below cell defines the function which will be called for plotting visuals for first objective i.e. trend of accidents and casualties with the year and month using Gridspec. The trends will be shown with line chart and the percentage of reasons

In [20]:

```
1 def plot_obj1():
2     fig = plt.figure(figsize=(15,15))
3     #using gridspec to create the visual in the single space
4     gs = GridSpec(2, 2, figure=fig)
5     #plotting for the first visual starting
6     ax1 = fig.add_subplot(gs[0, :])
7     ax1.set_title('Accident & Causality VS yearmonth')
8     l1 = ax1.plot("First_day_of_month", "sum", data=df, marker='o', markersize = 4, label="#
9     ax_new = ax1.twinx()
10    l2 = ax_new.plot("First_day_of_month", "count", data=df, marker='o', markersize = 4, col
11    #adding legends and label for 1st graph
12    lns = l1+l2
13    labs = [l.get_label() for l in lns]
14    ax1.legend(lns, labs)
15    ax1.set_xlabel("YearMonth")
16    ax1.set_ylabel("Causality")
17    ax_new.set_ylabel("Accidents")
18    #second graph plot start
19    ax2 = fig.add_subplot(gs[1, :-1])
20    ax2.set_title('Top 5 Reasons for Oct-2011')
21    reason1 = df_2011["reason"].to_list()
22    data1 = df_2011["Number_of_Casualties"].to_list()
23    wedges, texts, autotexts = ax2.pie(data1, autopct=lambda pct: func(pct, data1))
24    ax2.legend(wedges, reason1,
25               title="Reasons",
26               loc="center left",
27               bbox_to_anchor=(0, -0.1, 1, 0)
28               )
29    #using connection patch to create the arrow from 1st graph to second
30    con = ConnectionPatch( xyA=(datetime.date(2011, 10, 1), 18479),
31                           coordsA="data",
32                           xyB=(0,1),
33                           coordsB=ax2.transData,
34                           axesA=ax1,
35                           axesB=ax2,
36                           arrowstyle="-|>", shrinkA=5, shrinkB=5,
37                           mutation_scale=20, fc="w")
38    #displaying the arrow
39    ax2.add_artist(con)
40    #3rd plot start in gridspec
41    ax5 = fig.add_subplot(gs[1:2,1])
42    ax5.set_title('Top 5 Reasons for Oct-2012')
43    reason4 = df_2012["reason"].to_list()
44    data4 = df_2012["Number_of_Casualties"].to_list()
45    wedges, texts, autotexts = ax5.pie(data4, autopct=lambda pct: func(pct, data4))
46    ax5.legend(wedges, reason4,
47               title="Reasons",
48               loc="center left",
49               bbox_to_anchor=(0, -0.1, 0, 0)
50               )
51    #using connection patch to get the arrow from 1st visual to 3rd visual
52    con = ConnectionPatch( xyA=(datetime.date(2012, 10, 1), 21472),
53                           coordsA="data",
54                           xyB=(0,1),
55                           coordsB=ax5.transData,
56                           axesA=ax1,
57                           axesB=ax5,
58                           arrowstyle="-|>", shrinkA=5, shrinkB=5,
59                           mutation_scale=20, fc="w")
```

```

60 #displaying arrow
61 ax5.add_artist(con)

```

In the below command we are plotting the data between police response over the year

In [79]:

```

1 def plot_police():
2     fig, ax = plt.subplots(figsize=(10,5))
3     index = np.arange(24)
4     labels = df_bar_police[df_bar_police["Did_Police_Officer_Attend_Scene_of_Accident"]]
5     bar_width = 0.4
6     ax.bar(index-0.2, df_bar_police[df_bar_police["Did_Police_Officer_Attend_Scene_of_Accident"]])
7     ax.bar(index+0.2, df_bar_police[df_bar_police["Did_Police_Officer_Attend_Scene_of_Accident"]])
8     ax.plot(index,df["count"],bar_width)
9     ax.legend(title = "police arrived on location",fontsize=7, title_fontsize=7)
10    ax.set_xticks(index, labels, rotation=45, ha='right');

```

## Objective 2:

Below command creates new Dataframe by grouping District, Reason, Severity and summing up by the number of casualties.

In [22]:

```

1 df_uk_region_reason = uk_filter_data.groupby(["Local_Authority_District_value", "reason"])

```

- The first dataframe contains the information of Number of Casualties per every region and is done by group by and sum python functions
- The second contains the top 10 regions which reported the highest number of casualties
- And the third contains the bottom 10 regions which reported the least number of casualties

In [23]:

```

1 #creating a region specific dataframe
2 df_uk_only_region = df_uk_region_reason.groupby(["Local_Authority_District_value"])["Number_of_Casualties"].sum()
3 #identifying top 10 and bottom 10 Locations in uk
4 df_uk_top10_region = df_uk_only_region.sort_values(by=["Number_of_Casualties"],ascending=False)
5 df_uk_bottom10_region = df_uk_only_region[df_uk_only_region["Local_Authority_District_value"]<df_uk_top10_region["Local_Authority_District_value"].min()]

```

The below command takes the above grouped Dataframe and shows only the reason with the number of Casualties which then is used to identify the top 10 reasons that are resulting in the maximum number of Casualties.

In [24]:

```
1 # creating reason specific data frame
2 df_uk_only_reason = df_uk_region_reason.groupby(["reason"])[ "Number_of_Casualties"].sum
3 #identifiying top 10 reasons for the accidents
4 df_uk_top10_reason = df_uk_only_reason.sort_values(by=["Number_of_Casualties"],ascending=True)
```

In the below command we are fetching those records in the dataframe where those regions had the maximum number of casualties and then using this dataframe we are simultaneously fetching the records with the top 10 reasons of the mishappenings.

In [25]:

```
1 # joining reason and region data uk main data to get the final result
2 top_10_final_df= pd.merge(df_uk_region_reason, df_uk_top10_region["Local_Authority_District"],on="Local_Authority_District",how="inner")
3 top_10_final_df = pd.merge(top_10_final_df,df_uk_top10_reason["reason"],on="reason",how="inner")
```

In the Above command we fetched the top regions and reasons and in the below we are getting the bottom 10 regions and reasons where the casualties were highest

In [26]:

```
1 df_uk_only_reason_bottom = df_uk_region_reason.groupby(["reason","Local_Authority_District"]).sum
2 df_uk_bottom_reason = pd.merge(df_uk_only_reason_bottom,df_uk_bottom10_region["Local_Authority_District"],on="Local_Authority_District",how="inner")
3 df_uk_bottom10_reason = df_uk_bottom_reason.groupby("reason",dropna=False)[ "Number_of_Casualties"].sum
```

In the below command we are joining our all data of region and reason with the top 10 regions we are doing inner join this will filter our data and only keep the top 10 region data with all their corresponding reason

In the second command in the same df, now we are joining the data with the top 10 reasons for accidents with the previously filtered data, now finally we will only have the top 10 regions and corresponding to those region top 10 reasons for each region in the dataframe

In [27]:

```
1 bottom_10_final_df= pd.merge(df_uk_region_reason, df_uk_bottom10_region["Local_Authority_District"],on="Local_Authority_District",how="inner")
2 bottom_10_final_df = pd.merge(bottom_10_final_df,df_uk_bottom10_reason["reason"],on="reason",how="inner")
```

In the command below we are grouping, summing the data with the region and severity and getting a dataframe that has all the severity with the region

In the next statment we are joining it with the top 10 region that we identified earlier to only keep those regions

In [28]:

```
1 ## top severity
2 df_uk_region_severity = df_uk_region_reason.groupby(["Local_Authority_District_value", '
3 #joining with top 10 regions to get the information
4 df_uk_region_severity_top10 = pd.merge(df_uk_region_severity,df_uk_top10_region["Local_
```

In the command below we are doing the same thing but with the bottom 10 of our list.

In [29]:

```
1 ## bottom severirty
2 df_uk_region_severity_bottom10 = pd.merge(df_uk_region_severity,df_uk_bottom10_region['
3 df_uk_region_severity_bottom10.loc[-1] = ["Shetland Islands", 1,"Fatal",0]
```

In the command below, we are using the seaborn library to create a heatmap. We are explicitly providing the label for Xaxis. then we are creating a legend table without any border that will give the full form of the abbreviation

In [30]:

```
1 def heatmap():
2     heat= top_10_final_df.groupby(["Local_Authority_District_value","reason"])[ "Number_
3     #creating label as the reason are too big to disply in the visual
4     labelx = ["R1","R2","R3","R4","R5","R6","R7","R8","R9","R10"]
5     # labely = heat.reset_index()["Local_Authority_District_value"].to_list()
6     sns.heatmap(heat,cmap="RdBu",square = True,linewidth =1,xticklabels=labelx)
7     cell_text= [[i] for i in top_10_final_df["reason"].drop_duplicates().to_list()]
8     the_table = plt.table(cellText =cell_text ,rowLabels=labelx,collabels=["Fullform"],
9     the_table.auto_set_font_size(False)
10    the_table.set_fontsize(8)
```

In this command we are creating the data set that will help us in creating the bar chart with table.

data - here we are creating the list of list that will be passed into the bar graph plot, and because we are creating a stack plot we need to creat a list of list.

columns - this is the xaxis/ the table header that will be displayed in the visual

rows - this the row heading for the table

In [77]:

```
1 data = [df_uk_region_severity_top10[df_uk_region_severity_top10["Accident_Severity"]==1  
2         df_uk_region_severity_top10[df_uk_region_severity_top10["Accident_Severity"]==2  
3         df_uk_region_severity_top10[df_uk_region_severity_top10["Accident_Severity"]==3  
4         ]  
5 cell_text = [df_uk_region_severity_top10[df_uk_region_severity_top10["Accident_Severity"]==1  
6             df_uk_region_severity_top10[df_uk_region_severity_top10["Accident_Severity"]==2  
7             df_uk_region_severity_top10[df_uk_region_severity_top10["Accident_Severity"]==3  
8             ]  
9 # col = df_uk_region_severity_top10[["Local_Authority_District_value"]].drop_duplicates  
10 columns=('Barnet', 'Birmingham', 'Bradford', 'Cheshire East', 'Cheshire West', 'Cornwall', 'C  
11 rows = df_sev_mapping["Accident_Severity_value"].to_list()
```

In the command below we are plotting a bar graph with the table at its bottom that will act as the legend for the x-axis of the bar graph and table header for the table graph

In [32]:

```
1 def CbyR():  
2     # Get some pastel shades for the colors -BuPu  
3     colors = plt.cm.Accent(np.linspace(0, 0.5, len(rows)))  
4     n_rows = len(data)  
5     index = np.arange(len(columns)) + 0.3  
6     bar_width = 0.4  
7     # Initialize the vertical-offset for the stacked bar chart.  
8     y_offset = np.zeros(len(columns))  
9     fig, ax = plt.subplots(figsize = (15,7))  
10    # Plot bars and create text labels for the table  
11    for row in range(n_rows):  
12        ax.bar(index, data[row], bar_width, bottom=y_offset, color=colors[row])  
13        y_offset = y_offset + data[row]  
14    # Reverse colors and text labels to display the last value at the top.  
15    colors = colors[::-1]  
16    cell_text.reverse()  
17    # Add a table at the bottom of the axes  
18    the_table = ax.table(cellText=cell_text,  
19                        rowLabels=rows,  
20                        rowColours=colors,  
21                        colLabels=columns,  
22                        loc='bottom')  
23    the_table.auto_set_font_size(False)  
24    the_table.set_fontsize(10)  
25    # Adjust layout to make room for the table:  
26    # plt.subplots_adjust(left=0.2, bottom=0.2)  
27    plt.ylabel("Number of Casualties")  
28    plt.xticks([])  
29    plt.title('Casualties by Region')  
30    plt.show()
```

same as heatmap1 but passing different data points

In [33]:

```
1 def heatmap2():
2     heat_bottom= bottom_10_final_df.groupby(["Local_Authority_District_value", "reason"])
3     labelx = ["R1", "R2", "R3", "R4", "R5", "R6", "R7", "R8", "R9", "R10"]
4     # labely = heat.reset_index()["Local_Authority_District_value"].to_list()
5     sns.heatmap(heat_bottom, cmap="seismic", square = True, linewidth = 1, xticklabels=labelx)
6
7     cell_text= [[i] for i in bottom_10_final_df["reason"].drop_duplicates().to_list()]
8     the_table = plt.table(cellText = cell_text, rowLabels=labelx, colLabels=["Fullform"],
9     the_table.auto_set_font_size(False)
10    the_table.set_fontsize(8)
```

same as the above explanation for bar chart with table data prepration

In [68]:

```
1 data1 = [df_uk_region_severity_bottom10[df_uk_region_severity_bottom10["Accident_Severity"] == "R1"],
2          df_uk_region_severity_bottom10[df_uk_region_severity_bottom10["Accident_Severity"] == "R2"],
3          df_uk_region_severity_bottom10[df_uk_region_severity_bottom10["Accident_Severity"] == "R3"],
4          ]
5
6 cell_text1= [df_uk_region_severity_bottom10[df_uk_region_severity_bottom10["Accident_Severity"] == "R1"],
7             df_uk_region_severity_bottom10[df_uk_region_severity_bottom10["Accident_Severity"] == "R2"],
8             df_uk_region_severity_bottom10[df_uk_region_severity_bottom10["Accident_Severity"] == "R3"],
9             ]
10
11 col = df_uk_region_severity_bottom10[["Local_Authority_District_value"]].drop_duplicates()
12 columns = tuple(col["Local_Authority_District_value"])
13 # ('Barnet', 'Birmingham', 'Bradford', 'Cheshire East', 'Cheshire West', 'Cornwall', 'County')
14 rows = df_sev_mapping["Accident_Severity_value"].to_list()
```

code for 2nd table with bar chart- explanation same as the first code



In [69]:

```
1 def CbyR_best():
2     # Get some pastel shades for the colors -BuPu
3     # colors = plt.cm.Pastel2(np.linspace(0, 0.5, len(rows)))
4     colors = plt.cm.tab20(np.linspace(0, 0.5, len(rows)))
5     n_rows = len(data1)
6     index = np.arange(len(columns)) + 0.3
7     bar_width = 0.4
8     # Initialize the vertical-offset for the stacked bar chart.
9     y_offset = np.zeros(len(columns))
10    fig,ax = plt.subplots(figsize = (15,7))
11    # Plot bars and create text labels for the table
12    for row in range(n_rows):
13        ax.bar(index, data1[row], bar_width, bottom=y_offset, color=colors[row])
14        y_offset = y_offset + data1[row]
15    # Reverse colors and text labels to display the last value at the top.
16    colors = colors[::-1]
17    cell_text1.reverse()
18    # Add a table at the bottom of the axes
19    the_table = ax.table(cellText=cell_text1,
20                        rowLabels=rows,
21                        rowColours=colors,
22                        colLabels=columns,
23                        loc='bottom')
24    the_table.auto_set_font_size(False)
25    the_table.set_fontsize(10)
26    # Adjust layout to make room for the table:
27    # plt.subplots_adjust(left=0.2, bottom=0.2)
28    plt.ylabel("Number of Casualties")
29    plt.xticks([])
30    plt.title('Casualties by Region')
31    plt.show()
```

### Objective 3:

In the below command we are fetching the sum and count of Severities corresponding to the first day of every month in a year

In [36]:

```
1 df_severity = uk_filter_data.groupby(["First_day_of_month","Accident_Severity"])[ "Number"
```

In the below command we are manipulating the data 1 - group by to get only the required data points 2 - pivoting the group by data to get in desired format that will be usable in the stackplot

In [37]:

```
1 #fetching only 3 columns
2 df_severity = df_severity[["First_day_of_month","Accident_Severity","sum"]]
3 #setting the index from default to first day of month
4 df_severity = df_severity.set_index("First_day_of_month")
5 #made the columns as severity
6 df_severity = df_severity.pivot(columns='Accident_Severity')
7 #reseting index
8 df_severity = df_severity["sum"].reset_index()
9 #renaming the columns
10 df_severity = df_severity.rename(columns = {1:"y1",2:"y2",3:"y3"})
```

In the below command we are creating the plot we are simply using the stackplot function provided by the matplotlib and passing the data points that we prepared above

In [38]:

```
1 def scvsYear():
2     fig,ax1 = plt.subplots(figsize=(10,6))
3
4     ax1.set_title('Severity & Causality VS yearmonth')
5     ax1.stackplot(
6         'First_day_of_month',
7         'y3',
8         'y2',
9         'y1',
10        data=df_severity,
11        colors=['pink', 'purple', 'red'],
12        alpha=0.5
13    )
14    #adding legends and label for 1st graph
15    labels = ["Slight", "Serious", "Fatal" ]
16    ax1.legend(labels)
17    ax1.set_xlabel("YearMonth")
18    ax1.set_ylabel("Causality")
19    ax1.set_ylim(0, 26000,5000)
20    plt.xticks(rotation = 45);
```

In the command below we are grouping the data on reason, accident severity to get only the required column. Then we are join with the top reasons dataframe created in objective 2 to only get the data for those reasons

In [39]:

```
1 df_uk_reason_sev = uk_filter_data.groupby(["reason","Accident_Severity"])["Number_of_Ca
2 df_uk_reason_sev = pd.merge(df_uk_reason_sev,df_uk_top10_reason["reason"],on="reason",f
```

In the command below we are pivoting the data to suffice our need to create a horizontal bar chart.from the below command, we will have 4 column reason, and remaining will be the sum of casualties by each severity in each column

In [40]:

```
1 df_uk_reason_sev = df_uk_reason_sev.set_index("reason").pivot(columns = ["Accident_Seve
```

Here we are converting the data in a dictionary and then passing that data into the visual so that we can run a loop over it and plot the stack horizontal bar chart and give the different colour for each severity

In [41]:

```
1 category_names = ["Slight", "Serious", "Fatal",]
2 results = {
3     'R1': df_uk_reason_sev["Number_of_Casualties"].iloc[0].to_list()[::-1],
4     'R2': df_uk_reason_sev["Number_of_Casualties"].iloc[1].to_list()[::-1],
5     'R3': df_uk_reason_sev["Number_of_Casualties"].iloc[2].to_list()[::-1],
6     'R4': df_uk_reason_sev["Number_of_Casualties"].iloc[3].to_list()[::-1],
7     'R5': df_uk_reason_sev["Number_of_Casualties"].iloc[4].to_list()[::-1],
8     'R6': df_uk_reason_sev["Number_of_Casualties"].iloc[5].to_list()[::-1],
9     'R7': df_uk_reason_sev["Number_of_Casualties"].iloc[6].to_list()[::-1],
10    'R8': df_uk_reason_sev["Number_of_Casualties"].iloc[7].to_list()[::-1],
11    'R9': df_uk_reason_sev["Number_of_Casualties"].iloc[8].to_list()[::-1],
12    'R10': df_uk_reason_sev["Number_of_Casualties"].iloc[9].to_list()[::-1]
13 }
14 def severity():
15     labels = list(results.keys())
16     data = np.array(list(results.values()))
17     data_cum = data.cumsum(axis=1)
18     category_colors = plt.colormaps['tab20b'](
19         np.linspace(0.15, 0.85, data.shape[1]))
20
21     fig, ax = plt.subplots(figsize=(12, 5))
22     ax.invert_yaxis()
23     ax.xaxis.set_visible(False)
24     ax.set_xlim(0, np.sum(data, axis=1).max())
25
26     for i, (colname, color) in enumerate(zip(category_names, category_colors)):
27         widths = data[:, i]
28         starts = data_cum[:, i] - widths
29         rects = ax.barh(labels, widths, left=starts, height=0.5,
30             label=colname, color=color)
31
32         r, g, b, _ = color
33         text_color = 'white' if r * g * b < 0.5 else 'darkgrey'
34     ax.legend(ncol=len(category_names), bbox_to_anchor=(0, 1),
35         loc='lower left', fontsize='small')
36     plt.show()
```

## Objective 4:

Taking the sum of casualties for 2011 & 2012 by grouping the year and conditions: road, light and weather

In [42]:

```
1 df_combined = uk_acc_all_data[uk_acc_all_data["Year"].isin([2011, 2012])]
2 uk_bar_graph_data_natural = df_combined.groupby(
3     ["Year", "Road_Surface_Conditions", "Light_Conditions", "Weather_Conditions"])[
4     "Number_of_Casualties"].reset_index()
```

Creating two dataframes for 2011 & 2012 which is sorted based on the no of casualties, top 15 records is taken. New column is created in both the dataframes using the lambda function by concating the columns (Road\_Surface\_Conditions, Light\_Conditions, Weather\_Conditions) for all the rows

In [43]:

```
1 df_line_plot_2011 = uk_bar_graph_data_natural[uk_bar_graph_data_natural["Year"] == 2011]
2 df_line_plot_2011['Reasons'] = df_line_plot_2011.apply(lambda row: (row['Road_Surface_Conditions'] + row['Light_Conditions'] + '+' + row['Weather_Conditions']), axis=1)
3
4 df_line_plot_2012 = uk_bar_graph_data_natural[uk_bar_graph_data_natural["Year"] == 2012]
5 df_line_plot_2012['Reasons'] = df_line_plot_2012.apply(lambda row: (row['Road_Surface_Conditions'] + row['Light_Conditions'] + '+' + row['Weather_Conditions']), axis=1)
6
```

Merging both the years dataframes based on Reasons column, the top 10 records are taken. Storing the no of casualties as a list for 2011 & 2012 in two variables. Reasons is stored in a variable as list

In [44]:

```
1 merged_plot_df = pd.merge(df_line_plot_2011, df_line_plot_2012, on = 'Reasons', suffixes=('_2011', '_2012'))
2 merged_plot_df = merged_plot_df.sort_values("Reasons", ascending = False)
3
4 no_of_casualties_2011_list = merged_plot_df['Number_of_Casualties_2011'].tolist()
5 no_of_casualties_2012_list = merged_plot_df['Number_of_Casualties_2012'].tolist()
6 reasons_list_line = merged_plot_df['Reasons'].tolist()
```

Plotting the line chart based on the Reasons column for 2011 & 2012. Used matplotlib to specify the figure size and plot the x-y axis with the labels. x-axis has the top 10 reasons(R1-R10) and y-axis has the no of casualties

In [45]:

```
1 def reason_casualties():
2     plt.rcParams["figure.figsize"] = [12, 5]
3     plt.rcParams["figure.autolayout"] = True
4
5     r1 = mpatches.Patch(color = 'white', label = 'R1 - ' + reasons_list_line[0])
6     r2 = mpatches.Patch(color = 'white', label = 'R2 - ' + reasons_list_line[1])
7     r3 = mpatches.Patch(color = 'white', label = 'R3 - ' + reasons_list_line[2])
8     r4 = mpatches.Patch(color = 'white', label = 'R4 - ' + reasons_list_line[3])
9     r5 = mpatches.Patch(color = 'white', label = 'R5 - ' + reasons_list_line[4])
10    r6 = mpatches.Patch(color = 'white', label = 'R6 - ' + reasons_list_line[5])
11    r7 = mpatches.Patch(color = 'white', label = 'R7 - ' + reasons_list_line[6])
12    r8 = mpatches.Patch(color = 'white', label = 'R8 - ' + reasons_list_line[7])
13    r9 = mpatches.Patch(color = 'white', label = 'R9 - ' + reasons_list_line[8])
14    r10 = mpatches.Patch(color = 'white', label = 'R10 - ' + reasons_list_line[9])
15
16    plt.legend(handles=[r1, r2, r3, r4, r5, r6, r7, r8, r9, r10])
17
18    axes = plt.gca()
19
20    merged_plot_df.plot(kind = 'line', x = 'Reasons', y = 'Number_of_Casualties_2011',)
21    merged_plot_df.plot(kind = 'line', x = 'Reasons', y = 'Number_of_Casualties_2012',)
22
23    axes.set_xticks([0,1,2,3,4,5,6,7,8,9])
24    axes.set_ylabel('No of Casualties')
25    axes.set_title('Casualties based on the Road, Light and Weather Conditions')
26    axes.set_xticklabels(('R1', 'R2', 'R3', 'R4', 'R5', 'R6', 'R7', 'R8', 'R9', 'R10'))
27
28    plt.show()
```

Taking the sum of casualties for 2011 & 2012 by grouping the column

Did\_Police\_Officer\_Attend\_Scene\_of\_Accident(has two values Yes or No) and conditions: road, light and weather. New column is created in both the dataframes using the lambda function by concating the columns (Road\_Surface\_Conditions, Light\_Conditions, Weather\_Conditions) for all the rows.

In [46]:

```
1 df_bar_plot = df_combined.groupby(  
2     ["Did_Police_Officer_Attend_Scene_of_Accident", "Road_Surface_Conditions", "Light_C  
3 df_bar_plot = df_bar_plot.reset_index()  
4 df_bar_plot = df_bar_plot.sort_values(['Number_of_Casualties'], ascending = [False])  
5 df_bar_plot['Reasons'] = df_bar_plot.apply(lambda row: (row['Road_Surface_Conditions']  
6     row['Light_Conditions'] + '+' + row['Weather_Conditions']
```

Creating two dataframes depending if the police officer attended the scene of accident, sorted based on the no of casualties, top 20 records is taken. Merging both the years dataframes based on Reasons column, the top 10 records are taken.

In [47]:

```
1 df_bar_plot_police_attended_yes = df_bar_plot[(df_bar_plot["Did_Police_Officer_Attend_S  
2 df_bar_plot_police_attended_no = df_bar_plot[(df_bar_plot["Did_Police_Officer_Attend_Sc  
3 merged_bar_plot_df = pd.merge(df_bar_plot_police_attended_yes, df_bar_plot_police_atten  
4 merged_bar_plot_df = merged_bar_plot_df.sort_values("Reasons", ascending = False)
```

Storing the no of casualties in two variables when the police attended the event or not. Reasons is stored in a variable as list.

In [48]:

```
1 no_of_casualties_police_attended_yes_list = merged_bar_plot_df['Number_of_Casualties_Ye  
2 no_of_casualties_police_attended_no_list = merged_bar_plot_df['Number_of_Casualties_No'  
3 reasons_list_bar = merged_bar_plot_df['Reasons'].tolist()
```

Plotting the bar chart based on the whether the police attended the scene or not. Used matplotlib to specify the figure size and plot the x-y axis with the labels. x-axis has the top 10 reasons(R1-R10) and y-axis has the no of casualties

In [49]:

```
1 def Police_severity():
2     plt.rc('figure', figsize=(12, 5))
3
4     x_axis = np.arange(len(merged_bar_plot_df))
5
6     plt.bar(x_axis, no_of_casualties_police_attended_yes_list, width = 0.2, color = 'blue')
7     plt.bar(x_axis + 0.20, no_of_casualties_police_attended_no_list, width = 0.2, color = 'red')
8
9     plt.xticks(x_axis, ('R1', 'R2', 'R3', 'R4', 'R5', 'R6', 'R7', 'R8', 'R9', 'R10'))
10
11     y1 = mpatches.Patch(color = 'blue', label = 'Police attended the scene of accident')
12     y2 = mpatches.Patch(color = 'red', label = 'Police did not attend the scene of accident')
13     r1 = mpatches.Patch(color = 'white', label = 'R1 - ' + reasons_list_bar[0])
14     r2 = mpatches.Patch(color = 'white', label = 'R2 - ' + reasons_list_bar[1])
15     r3 = mpatches.Patch(color = 'white', label = 'R3 - ' + reasons_list_bar[2])
16     r4 = mpatches.Patch(color = 'white', label = 'R4 - ' + reasons_list_bar[3])
17     r5 = mpatches.Patch(color = 'white', label = 'R5 - ' + reasons_list_bar[4])
18     r6 = mpatches.Patch(color = 'white', label = 'R6 - ' + reasons_list_bar[5])
19     r7 = mpatches.Patch(color = 'white', label = 'R7 - ' + reasons_list_bar[6])
20     r8 = mpatches.Patch(color = 'white', label = 'R8 - ' + reasons_list_bar[7])
21     r9 = mpatches.Patch(color = 'white', label = 'R9 - ' + reasons_list_bar[8])
22     r10 = mpatches.Patch(color = 'white', label = 'R10 - ' + reasons_list_bar[9])
23
24     plt.legend(handles=[y1, y2, r1, r2, r3, r4, r5, r6, r7, r8, r9, r10])
25
26     axes = plt.gca()
27
28     axes.set_xticks([0,1,2,3,4,5,6,7,8,9])
29     axes.set_xlabel('Accident Reasons')
30     axes.set_ylabel('No of Casualities')
31     axes.set_title('Casualties based on the Road, Lightning and Weather Conditions')
32
33     plt.show()
```

## Project Outcome (10 + 10 marks)

*This section should describe the outcome of the project by means of both explanation of the results and by graphical visualisation in the form of graphs, charts or or other kinds of diagram*

*The section should begin with a general overview of the results and then have a section for each of the project objectives. For each of these objectives an explanation of more specific results relating to that objective should be given, followed by a section presenting some visualisation of the results obtained. (In the case where the project had just one objective, you should still have a section describing the results from a general perspective followed by a section that focuses on the particular objective.)*

*The marks for this section will be divided into 10 marks for Explanation and 10 marks for Visualisation. These marks will be awarded for the Project Outcome section as a whole, not for each objective individually. Hence, you do not have to pay equal attention to each. However, you are expected to have a some explanation and visualisation for each. It is suggested you have 200-400 words explanation for each objective.*

## Overview of Results

*Give a general overview of the results (around 200 words).*

# Objective 1

## Explanation of Results

### Causality trends over months and their top 5 reasons

Here we are trying to identify how the causality is changing over the months and how the factors like -

- Light conditions - *Day or Night*,
- Road conditions - *Wet or Dry*
- Weather Conditions - *Winds or no Winds*

impact the accident and the number of casualties.

The data suggests that there is a high relationship between the time of the day and road conditions. Along with this, we can see there is an increasing trend in a particular **time period of the year** where the accidents increase relatively and decrease after that particular period.

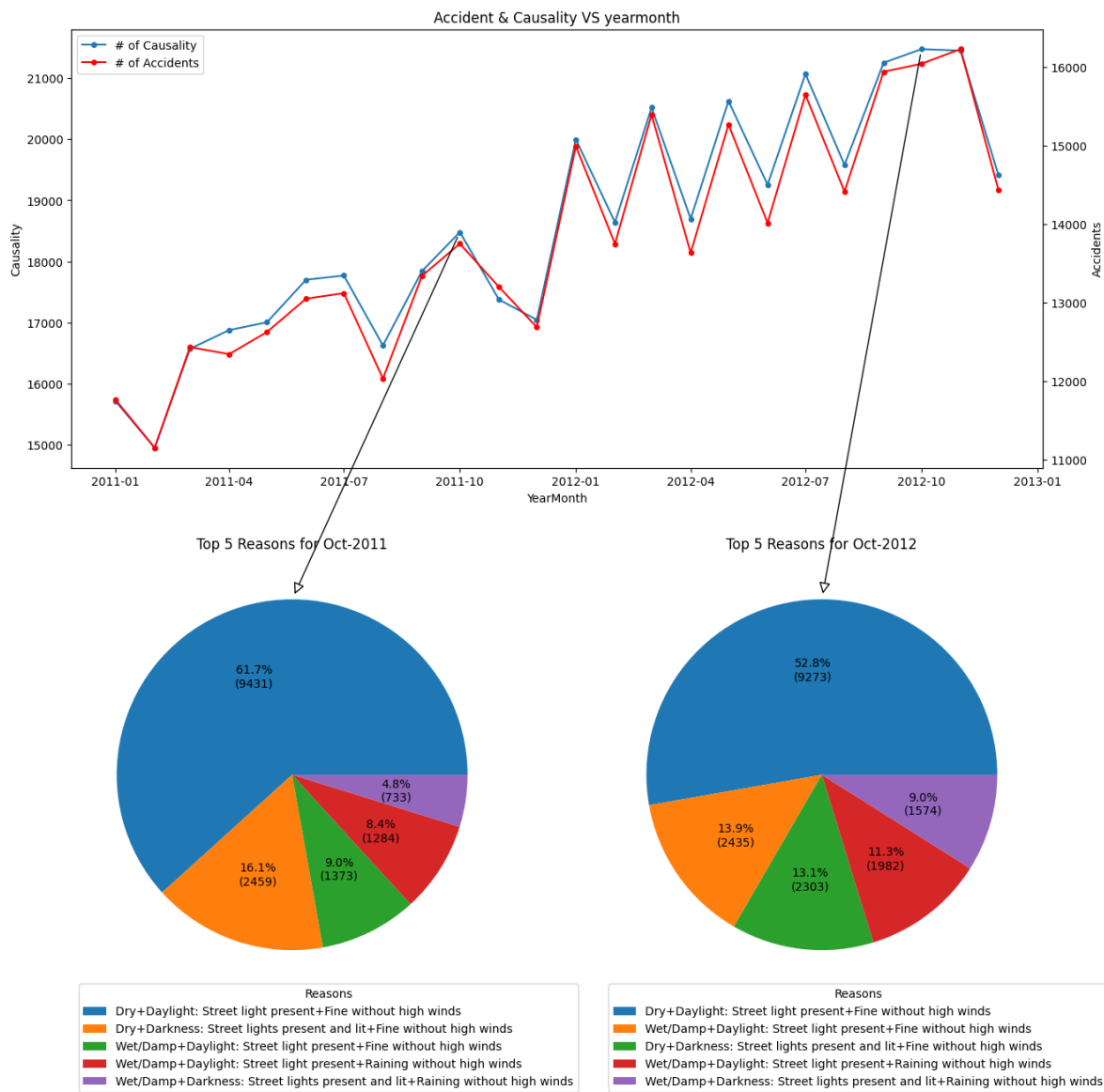
1. There are frequent peaks and bottoms in the values of *casualties and accidents* in our data i.e. there is a trough and crest after every month with an increase in number of both casualties and accidents.
2. There is a *strong* relationship between the year months and the number of accidents, we can see from the below graph, *title = "Accident & Causality VS year month"* that on the 10th month of each year the number of accidents is higher when in comparison to the rest of the months. For e.g. highest number of casualties in 2011, are reported in Oct(first arrow) i.e. 18479
3. We can see that during the first quarter of 2011, the line of casualties and accidents is *identical* which implies that for every accident there were quite less casualties as compared to others (the first 3 dots where the red and blue are overlapping). After that, the number of casualties started increasing as compared to the number of accidents and then again by the end of 2011 both started to overlap till the 1st month of 2012 and then again, the *same pattern* can be observed in 2012.
4. One more key observation is that during the troughs, casualties are much higher than the accidents as compared to the numbers during crest. This appeared an interesting discovery on which further analysis could be performed which after investigation came out that the top 5 reasons for accidents where there is a peak for each year is due to "dry road conditions where there is daylight and no high winds"

## Visualisation

For the 4th point the same can be observed in the 2 pie charts displayed below with the title "top 5 reasons for 2011", "top 5 reasons for 2012"

In [50]:

```
1 plot_obj1()
```

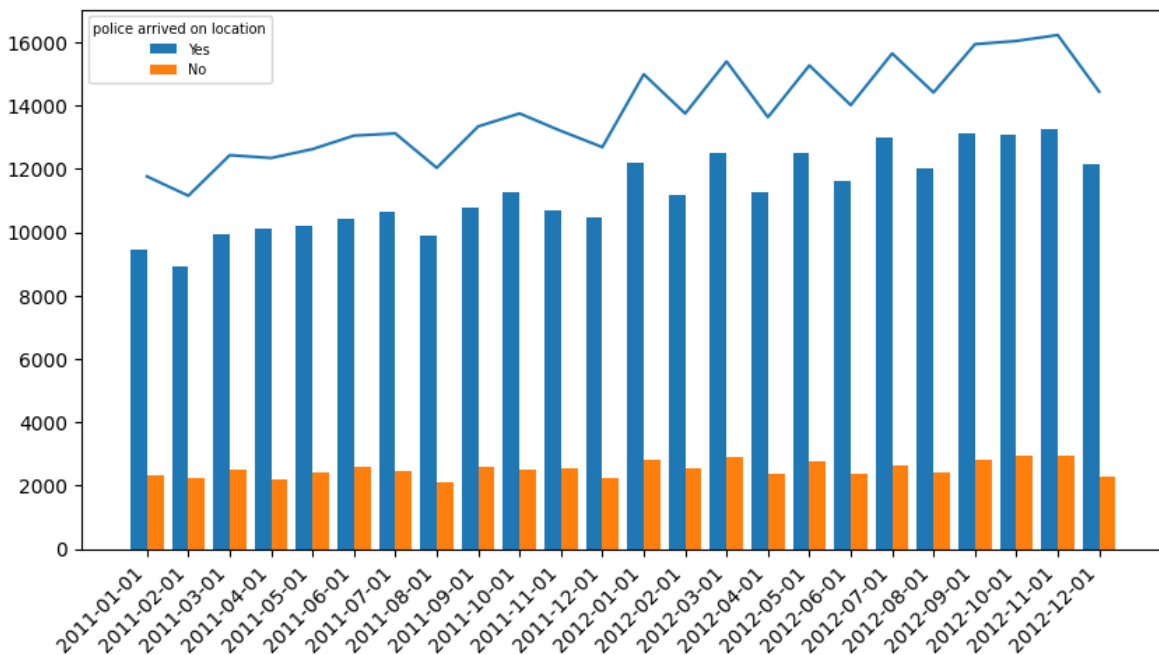


The below graph tells us that "police response" are respectively good when we are checking with the months



In [51]:

```
1 plot_police()
```



## Objective 2

### Explanation of Results

#### Analysis on Reason by region

We have identified the top 10 best and the worst regions to live in terms of casualties.

- Best region - **Orkney Island**
- Worst region - **Birmingham**

We have identified which reason has the highest impact on the worst places to live in terms of **# of casualty**, as could be seen from the **Worst Places to live in** heat map in the command below. For all the places **R4** (Dry+Daylight: Street light present+Fine without high winds) reason comes out to be dominant. In the same manner, we have identified the reason of accidents at the best places in the **best places to live in** heat map and there **R3** (Dry+Daylight: Street light present+Fine without high winds) comes out to be dominating.

**Note - Here you might be wondering why the different names R3 and R4 have the same reason in both scenarios. The reason behind that is, that for each scenario we are doing analysis, first identifying the top 10 reasons for that cluster. for e.g. we are identifying the top 10 reasons for the worst cities and those top 10 might be different from the top 10 reasons for the best cities hence the different abbreviations for each scenario.**

#### Analysis of Severity by region

We have identified the cities which have least and maximum number of casualties as per the severities during accidents in both best and worst places to live in.

If a person wants to live in the Worst region cluster and still wants to be safe then we can identify that as well.

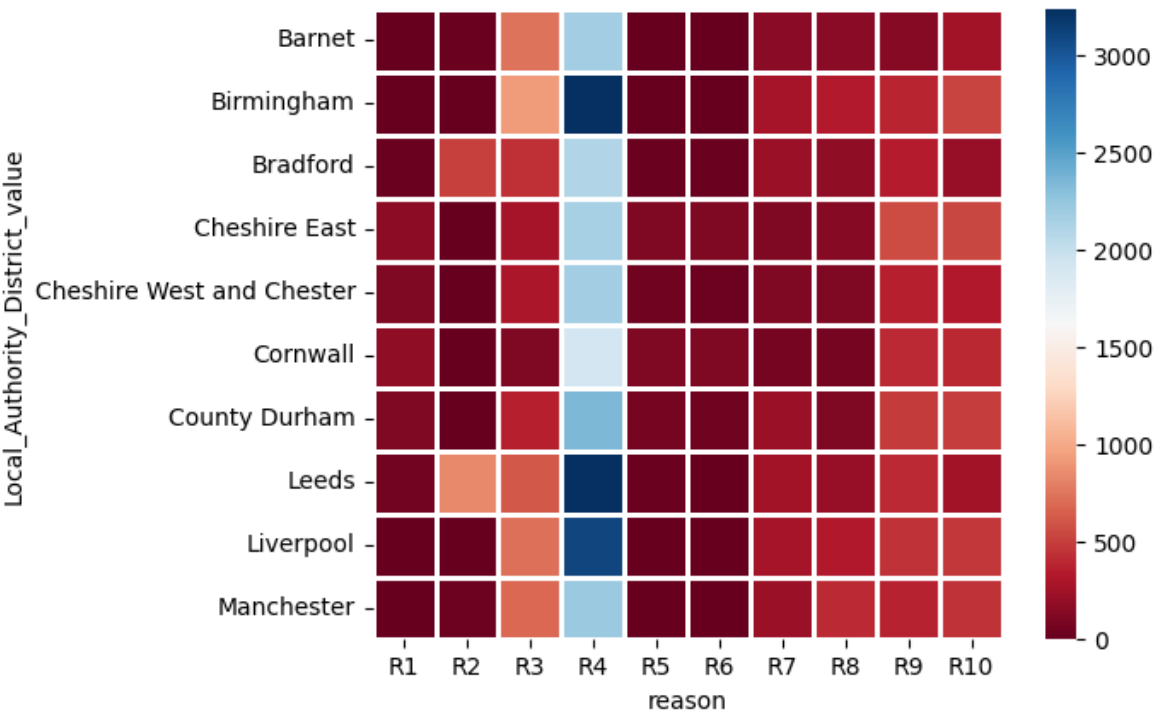
From the **causality by region and severity** graph we can see **Barnet** has the least number of fatal accidents but **County Durham** has the highest number of fatal accidents in that region. In the same way, if a person wants to be extra safe and wants to identify the safest place from the best region to live in the cluster that would be **Western isles** it has 0 fatal accidents and highest number of **fatal** accidents in this are in ***Rutland***

## Visualisation

### Worst places to live in

In [52]:

```
1 heatmap()
```

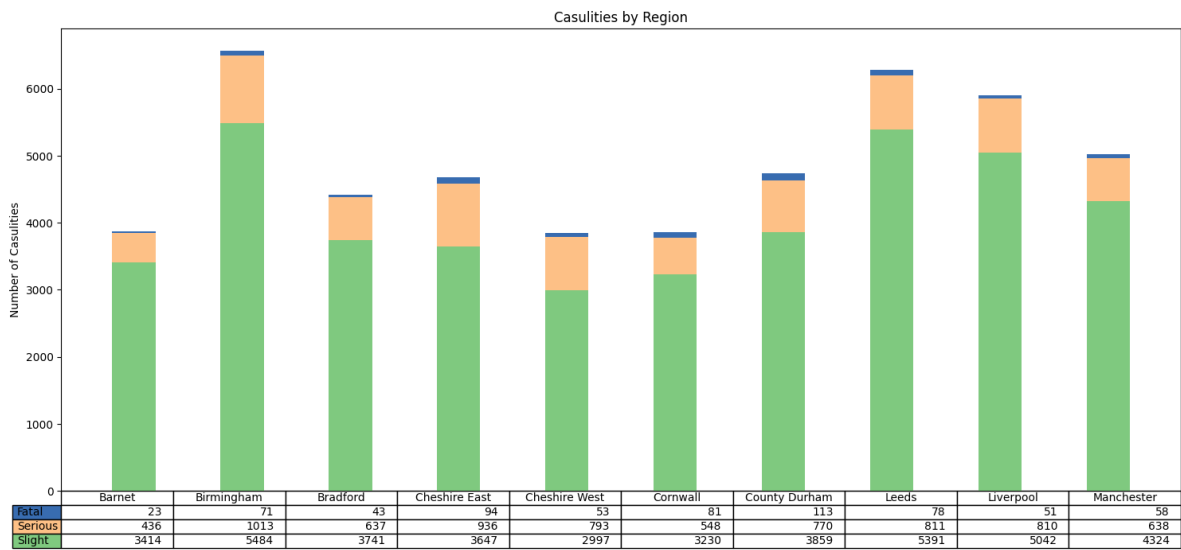


Fullform

- R1 Dry+Darkeness: No street lighting+Fine without high winds
- R2 Dry+Darkness: Street lighting unknown+Fine without high winds
- R3 Dry+Darkness: Street lights present and lit+Fine without high winds
- R4 Dry+Daylight: Street light present+Fine without high winds
- R5 Wet/Damp+Darkeness: No street lighting+Fine without high winds
- R6 Wet/Damp+Darkeness: No street lighting+Raining without high winds
- R7 Wet/Damp+Darkness: Street lights present and lit+Fine without high winds
- R8 Wet/Damp+Darkness: Street lights present and lit+Raining without high win
- R9 Wet/Damp+Daylight: Street light present+Fine without high winds
- R10 Wet/Damp+Daylight: Street light present+Raining without high winds

In [78]:

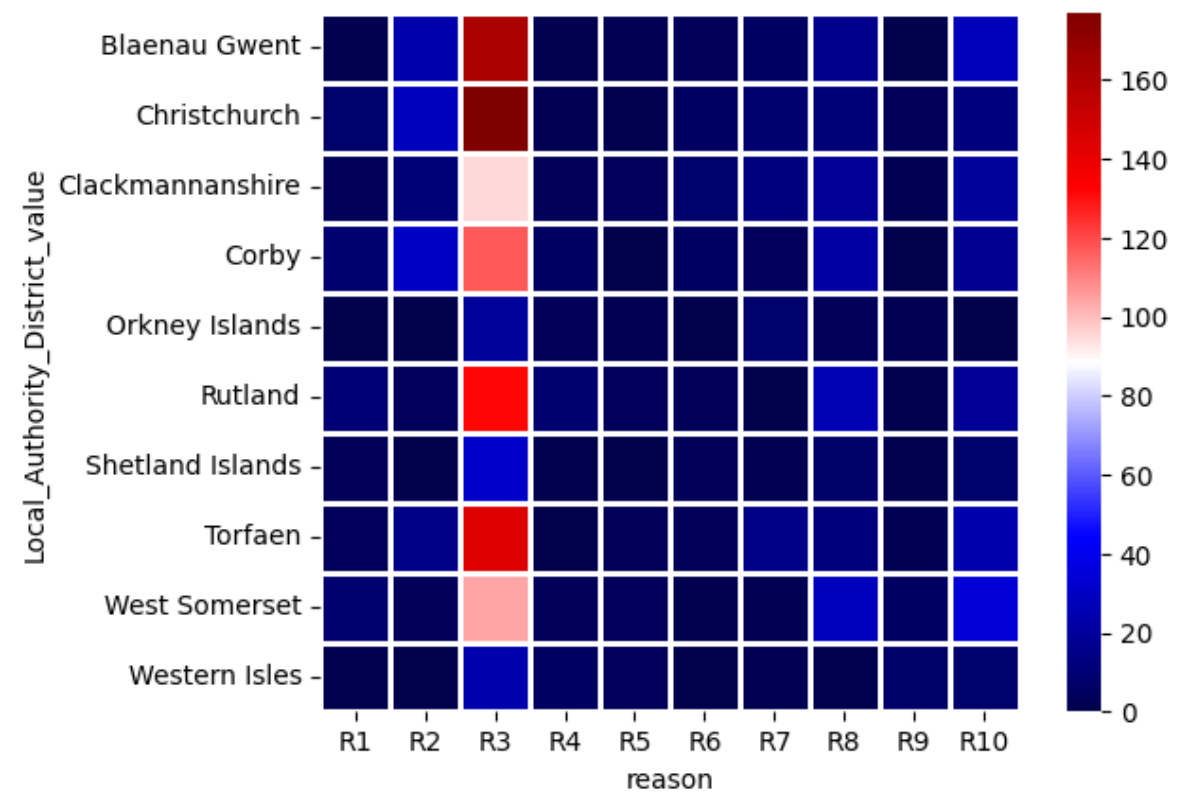
```
1 CbyR()
```



Best places to live in

In [54]:

```
1 heatmap2()
```

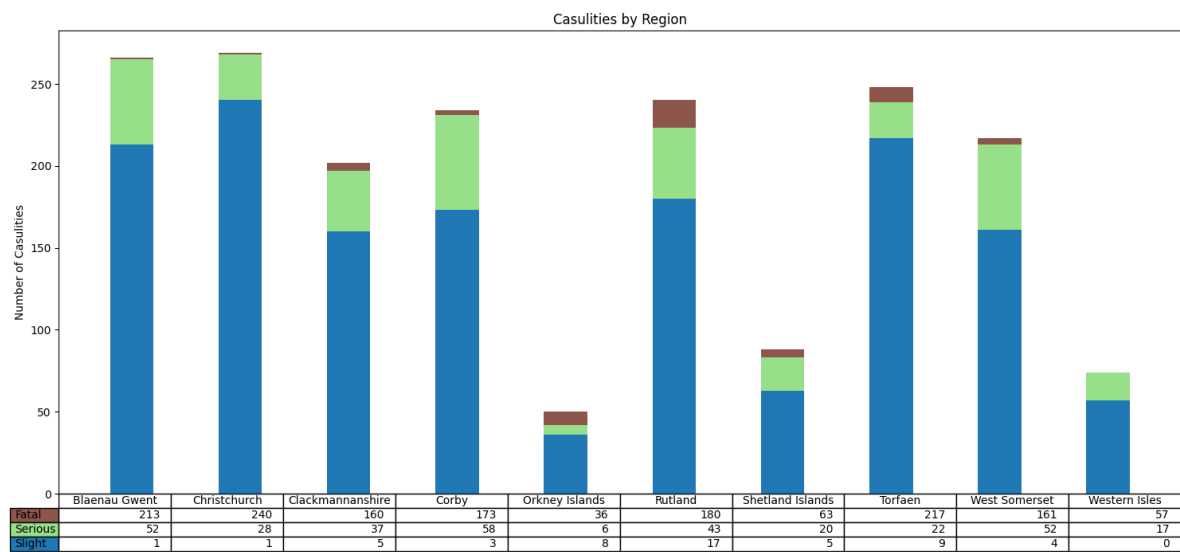


Fullform

- R1 Dry+Darkness: No street lighting+Fine without high winds
- R2 Dry+Darkness: Street lights present and lit+Fine without high winds
- R3 Dry+Daylight: Street light present+Fine without high winds
- R4 Wet/Damp+Darkness: No street lighting+Fine without high winds
- R5 Wet/Damp+Darkness: No street lighting+Raining without high winds
- R6 Wet/Damp+Darkness: Street lights present and lit+Fine without high winds
- R7 Wet/Damp+Darkness: Street lights present and lit+Raining without high wii
- R8 Wet/Damp+Daylight: Street light present+Fine without high winds
- R9 Wet/Damp+Daylight: Street light present+Raining without high winds
- R10 Wet/Damp+Daylight: Street light present+Raining with high winds

In [71]:

```
1 CbyR_best()
```



### Objective 3

Here we are identifying the month from the 2 years which has the highest number of ***fatal,serious*** and ***slight*** accidents. From the analysis we can conclude through the below table:

Highest number of accidents	Month
Fatal	Dec, 2012
Serious	May, 2012
Slight	Oct, 2012

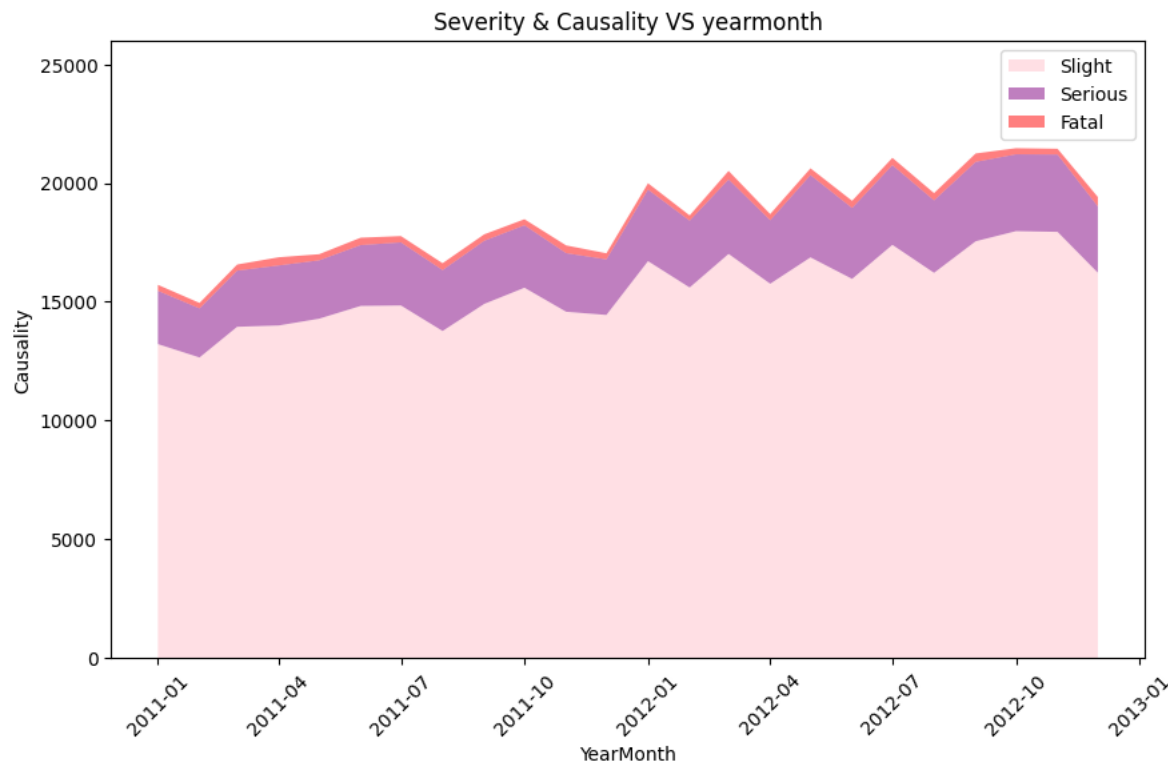
Secondly, we are identifying the top 10 reasons which have the highest number of casualties by severity, then we can see that **R4** (Dry+Daylight: Street light present+Fine without high winds) sores through the graph and has the maximum number of accidents across all the severities i.e. ***fatal,serious,slight***

### Visualisation

In [56]:

```
1 scvsYear()
```

C:\Users\kriti\AppData\Local\Temp\ipykernel\_23672\569235331.py:20: MatplotlibDeprecationWarning: Passing the emit parameter of set\_ylim() positionally is deprecated since Matplotlib 3.6; the parameter will become keyword-only two minor releases later.  
ax1.set\_ylim(0, 26000,5000)



In [57]:

```
1 severity()
```



# Objective 4

## Explanation of Results

The number of casualties for the reasons R1, R2 and R8 for 2011 and 2012 have some significant amount of difference. Whereas the reasons R5, R6, R7 and R10 have comparatively lesser differences among them in both years.

**R8** contributed to have the **most number of casualties** and on the other hand, **R7** has the **least no of casualties**. There can be below mentioned 2 explanations for the pattern observed in the reasons where we have **R8** (dry road conditions where there is daylight and no high winds dominating all the reasons:

1. These accidents are happening because people are driving for long periods, and they are dehydrated and hence resulting in accidents.
2. The dry weather is causing mirages and causing accidents.

Unfortunately, we don't have the data to make a conclusion on the same.

## Police performace over reasons

It can be clearly understood that the police made sure to attend the scene of the accidents no matter what the reason was. The number of times police did not attend was very less compared to the number of times they did attend the accident scene.

**Dry road conditions where there is daylight and no high winds** contributes to be the reason due to which the most number of casualties were there for both the cases where the police attend and do not attend the scene of the accident and on the other hand **R2** (Wet road conditions where there is Daylight and Raining with high winds) has the least no of casualties probably they are enjoying the weather from a safe distance.

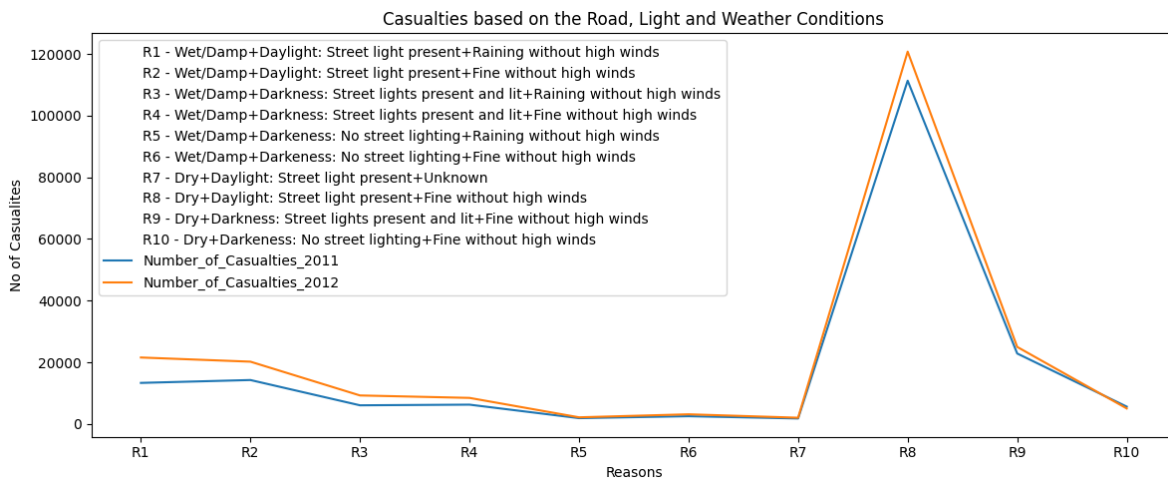
## Visualisation

Interpretation from the line graph: The Reasons(R1-R10) in the labels above are a combination of Road conditon + Lightning condition + Weather condition



In [58]:

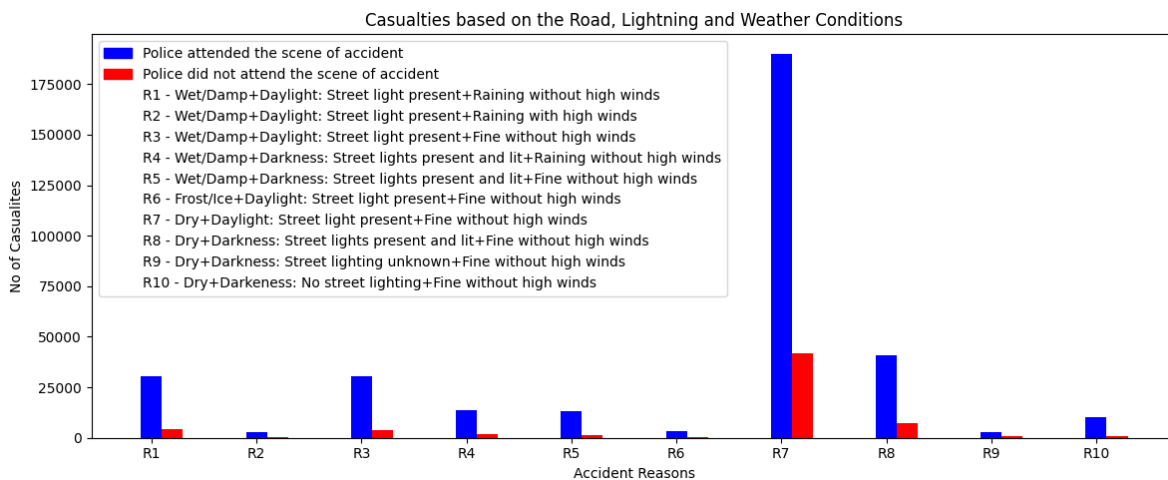
```
1 reason_casulities()
```



Interpretation from the bar graph: The Reasons(R1-R10) in the labels above are a combination of Road condition + Lightning condition + Weather condition

In [59]:

```
1 Police_severity()
```



## Conclusion (5 marks)

### Achievements

Unfortunately, the highest number of accidents were not during the conditions such as raining, snowing or at night, Instead, they have happened **during day-light, when the roads were clean** and there were no external forces that would act as catalyst. The **highest** number of accidents have taken place in the autumn month of both the **years(October)** which is considered to be one of the best months to visit UK. The **total number of accidents** for both years are **0.44 million**. Happily, the **fatal accidents are very low** i.e. in k's (thousands), to be precise 6.8k. Serious accidents are relatively high but still low in thousands i.e. 66K and when moving to **Slight accidents** they are approaching to **half a million** i.e. 0.37 million. One of the great things about UK is its Police Force Response, it doesn't matter what the condition is or the time of the year is, **Police response has**

**always been good.** Lastly we have seen that the highest number of accidents have taken place in **Birmingham followed by Leeds, Liverpool, and Manchester** which pushes us to investigate more on the reason behind these particular cities and any other external condition that can be an important factor.

## Limitations

Due to the limited set of factors such as surroundings, human control or physical facilities at the site of an accident but there is no information present if the **person was drunk driving or how was their mood or how is their physical and mental health or was they even eligible to ride the vehicle if they were riding one during accident.** We don't have the details of the casualty such as their age and sex itself and hence the precise reason for the high number of accidents in big cities such as Birmingham, Manchester, Leeds and Liverpool can't be commented on.

## Future Work

In Future work we would like to obtain more data related to the casualties as mentioned in the limitations as with that data we would be able to comment on the reason even more precisely and we hypothesise that there would be a **relationship between the drinking, licenced person with the accidents.**