



SSK4409-1: BIG DATA ANALYTICS

PROJECT

PROF. MADYA TS. DR. ISKANDAR BIN ISHAK

Prepared by:

No.	Name	Matric No.
1.	Anas Zulkifli bin Mohd Jeffry	206520
2.	Nik Muhammad Asyraf bin Nik Ismail	206630
3.	Amir Nurhakim bin Mohd Zaid	207092
4.	Muhammad Ikhwan Khuzairi bin Rozdi	206568

Part 1 - Big Data Platform	3
1. Apache Hadoop Installation	3
1.1. Installation details	3
1.2. Environment Information	4
1.3. Web Interfaces	4
2. MapReduce program on word counting on a text file	5
Part 2: Data Analytics	8
2.1 Data Analytics Tasks	8
2.1.1 Pre-processing	8
2.1.2 Data Cleaning	9
2.1.3 Modeling	9
2.1.4 Results	11
Part 3: Data Visualization	20
3.1 Public Tableau Installation	20
3.1.1 Installation Details	20
3.1.2 Creation of Public Tableau account in Public Tableau Site	21
3.2 Relationships in The Chosen Dataset with Tableau Analysis	21
3.2.1 Worksheet 1: Time Trend Analysis	21
3.2.2 Worksheet 2: Distribution in Urban and Rural Areas	22
3.2.3 Worksheet 3: Geographic Distribution	23
3.2.4 Worksheet 4: Severity on the Map	24
3.2.5 Worksheet 5: Factors by Severity	25
3.2.6 Worksheet 6: Light Conditions and Vehicle Types	26
3.3 Combination of Worksheets into a Single Story	27
3.3.1 Dashboard 1	27
3.3.2 Dashboard 2	28
3.3.3 Dashboard 3	29
3.4 Story Publication in Tableau Public Page	29

Part 1 - Big Data Platform

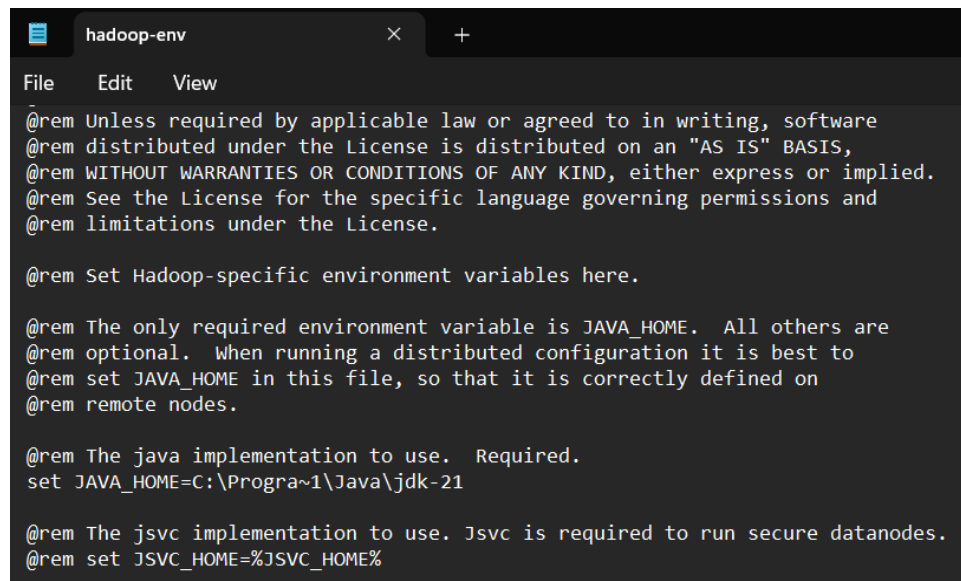
1. Apache Hadoop Installation

This part provides evidence of the successful installation of Apache Hadoop on our machine. Part 1 covers installation details, environment information, daemon status, web interfaces.

1.1. Installation details

Key configuration files, such as 'hadoop-env', 'core-site.xml', and 'hdfs-site.xml', were modified as per the installation requirements.

- hadoop-env

A screenshot of a text editor window titled 'hadoop-env'. The window shows the contents of the file, which includes a license notice, instructions on setting environment variables, and specific settings for JAVA_HOME and JSVC_HOME. The text is as follows:

```
File Edit View
@rem Unless required by applicable law or agreed to in writing, software
@rem distributed under the License is distributed on an "AS IS" BASIS,
@rem WITHOUT WARRANTIES OR CONDITIONS OF ANY KIND, either express or implied.
@rem See the License for the specific language governing permissions and
@rem limitations under the License.

@rem Set Hadoop-specific environment variables here.

@rem The only required environment variable is JAVA_HOME. All others are
@rem optional. When running a distributed configuration it is best to
@rem set JAVA_HOME in this file, so that it is correctly defined on
@rem remote nodes.

@rem The java implementation to use. Required.
set JAVA_HOME=C:\Progra~1\Java\jdk-21

@rem The jsvc implementation to use. Jsvc is required to run secure datanodes.
@rem set JSVC_HOME=%JSVC_HOME%
```

- core-site.xml

A screenshot of a text editor window titled 'core-site.xml'. The window shows the XML configuration for Hadoop, including the version, encoding, and a configuration block with a property for the default file system. The text is as follows:

```
C: > Hadoop > hadoop-2.9.2 > etc > hadoop > core-site.xml
1  <?xml version="1.0" encoding="UTF-8"?>
2  <?xml-stylesheet type="text/xsl" href="configuration.xsl"?>
3
4
5  <configuration>
6  <property>
7    <name>fs.defaultFS</name>
8    <value>hdfs://localhost:9000</value>
9  </property>
10 </configuration>
11
```

- hdfs-site.xml

```

hdfs-site.xml
C: > Hadoop > hadoop-2.9.2 > etc > hadoop > hdfs-site.xml
1  <?xml version="1.0" encoding="UTF-8"?>
2  <?xml-stylesheet type="text/xsl" href="configuration.xsl"?>
3
4
5  <configuration>
6    <property>
7      <name>dfs.replication</name>
8      <value>1</value>
9    </property>
10   <property>
11     <name>dfs.namenode.name.dir</name>
12     <value>C:\Hadoop\hadoop-2.9.2\data\namenode</value>
13     <final>true</final>
14   </property>
15   <property>
16     <name>dfs.datanode.data.dir</name>
17     <value>C:\Hadoop\hadoop-2.9.2\data\datanode</value>
18     <final>true</final>
19   </property>
20 </configuration>

```

1.2. Environment Information

The installed Apache Hadoop version is verified using the following command:

```

C:\Windows\System32\cmd.e  X  +  v
Microsoft Windows [Version 10.0.22621.2861]
(c) Microsoft Corporation. All rights reserved.

C:\Hadoop\hadoop-2.9.2\bin>hadoop version
Hadoop 2.9.2
Subversion https://git-wip-us.apache.org/repos/asf/hadoop.git -r 826afbeae31ca687bc2f8471dc841b66ed2c6704
Compiled by ajisaka on 2018-11-13T12:42Z
Compiled with protoc 2.5.0
From source with checksum 3a9939967262218aa556c684d107985
This command was run using /C:/Hadoop/hadoop-2.9.2/share/hadoop/common/hadoop-common-2.9.2.jar

```

1.3. Web Interfaces

Accessing Hadoop web interfaces to monitor system status

- HDFS NameNode

localhost:9870/dfshealth.html#tab-overview

Hadoop

Overview

Datanodes

Datanode Volume Failures

Snapshot

Startup Progress

Utilities

Overview '0.0.0.0:19000' (active)

Started:

Wed Jan 03 02:25:58 +0800 2024

Version:

3.2.1, rb3cbb467e22ea829b38084b7b01d07e0bf3842

Compiled:

Tue Sep 10 23:56:00 +0800 2019 by rohitsharmaks from branch-3.2.1

Cluster ID:

CID-a0bb6d7e-133e-4c4d-b391-bb880cb44f8

Block Pool ID:

BP-1641404058-10.106.7.193-1704208029490

Summary

Security is off.

Safemode is off.

1 files and directories, 0 blocks (0 replicated blocks, 0 erasure coded block groups) = 1 total filesystem object(s).

Heap Memory used 87.05 MB of 228.5 MB Heap Memory. Max Heap Memory is 889 MB.

Non Heap Memory used 49.35 MB of 50.88 MB Committed Non Heap Memory. Max Non Heap Memory is <unbounded>.

Configured Capacity:

475.83 GB

Configured Remote Capacity:

0 B

DFS Used:

322 B (0%)

- HDFS datanode

localhost:9864/datanode.html

Hadoop

Overview

Utilities

DataNode on 10.106.7.193:9866

Cluster ID:

CID-a0bb6d7e-133e-4c4d-b391-bb880cb44f8

Version:

3.2.1, rb3cbb467e22ea829b38084b7b01d07e0bf3842

Block Pools

Namenode Address	Block Pool ID	Actor State	Last Heartbeat	Last Block Report	Last Block Report Size (Max Size)
0.0.0.0:19000	BP-1641404058-10.106.7.193-1704208029490	RUNNING	2s	a few seconds	0 B (64 MB)

Volume Information

Directory	StorageType	Capacity Used	Capacity Left	Capacity Reserved	Reserved Space for Replicas	Blocks
C:\Hadoop\hadoop-3.2.1\data\dfs\data	DISK	322 B	124.43 GB	0 B	0 B	0

- YARN resource manager

localhost:8038/cluster

hadoop

All Applications

Cluster

About

Nodes

Node Labels

Applications

NEW

NEW SAVING

SUBMITTED

ACCEPTED

RUNNING

FINISHED

FAILED

KILLED

Scheduler

Tools

Cluster Metrics

Apps Submitted

Apps Pending

Apps Running

Apps Completed

Containers Running

Memory Used

Memory Total

Memory Reserved

VCores Used

VCores Total

Cluster Nodes Metrics

Active Nodes

Decommissioning Nodes

Decommissioned Nodes

Lost Nodes

Unhealthy Nodes

Rebooted Nodes

Scheduler Metrics

Scheduler Type

Scheduling Resource Type

Minimum Allocation

Maximum Allocation

Maximum Cluster Appli

Capacity Scheduler

[memory-mb (unit=M), vcores]

<memory:1024, vCores:1>

<memory:8192, vCores:4>

0

Show 20 entries

Search:

ID	User	Name	Application Type	Queue	Application Priority	StartTime	LaunchTime	FinishTime	State	FinalStatus	Running Containers	Allocated CPU Vcores	Allocated Memory MB	Reserved CPU Vcores	Reserved Memory MB	% of Queue	% of Cluster	Progress
No data available in table																		

Showing 0 to 0 of 0 entries

First

2. MapReduce program on word counting on a text file

- Open cmd in Administrator mode and move to "C:\Hadoop\hadoop-3.2.1\sbin" and start cluster

```
Administrator: Command Prompt
Microsoft Windows [Version 10.0.22621.2861]
(c) Microsoft Corporation. All rights reserved.

C:\Windows\System32>cd/

C:\>cd Hadoop

C:\Hadoop>cd hadoop-3.2.1

C:\Hadoop\hadoop-3.2.1>cd sbin

C:\Hadoop\hadoop-3.2.1\sbin>start-all.cmd
This script is Deprecated. Instead use start-dfs.cmd and start-yarn.cmd
starting yarn daemons

C:\Hadoop\hadoop-3.2.1\sbin>
```

Apache Hadoop Distribution - hadoop namenode

2022-03-03 12:04:24 Apache Hadoop Distribution - hadoop datanode

2022-03-03 12:04:24 Apache Hadoop Distribution - yarn resourcemanager

2022-03-03 12:04:24 Apache Hadoop Distribution - yarn nodemanager

2022-03-03 12:04:24 INFO: Registering org.apache.hadoop.yarn.server.nodemanager.webapp.JAXBContextResolver as a provider class

2022-03-03 12:04:24 INFO: Jan 03, 2024 3:23:45 AM com.sun.jersey.server.impl.application.WebApplicationImpl _initiate

2022-03-03 12:04:24 INFO: Initiating Jersey application, version 'Jersey: 1.19 02/11/2015 03:25 AM'

2022-03-03 12:04:24 INFO: Jan 03, 2024 3:23:46 AM com.sun.jersey.guice.spi.container.GuiceComponentProviderFactory getComponentProvider

2022-03-03 12:04:24 INFO: Binding org.apache.hadoop.yarn.server.nodemanager.webapp.JAXBContextResolver to GuiceManagedComponentProvider with the scope "Singleton"

2022-03-03 12:04:24 INFO: Jan 03, 2024 3:23:46 AM com.sun.jersey.guice.spi.container.GuiceComponentProviderFactory getComponentProvider

2022-03-03 12:04:24 INFO: Binding org.apache.hadoop.yarn.webapp.GenericExceptionHandler to GuiceManagedComponentProvider with the scope "Singleton"

2022-03-03 12:04:24 INFO: Jan 03, 2024 3:23:46 AM com.sun.jersey.guice.spi.container.GuiceComponentProviderFactory getComponentProvider

2022-03-03 12:04:24 INFO: Binding org.apache.hadoop.yarn.server.nodemanager.webapp.NMWebServices to GuiceManagedComponentProvider with the scope "Singleton"

- 2.2. Create an input directory in HDFS.
- 2.3. Copy the input text file named file1.txt in the input directory (input_dir) of HDFS.
- 2.4. Verify file1.txt available in HDFS input directory (input_dir).

```
Administrator: Command Prompt
C:\Hadoop\hadoop-3.2.1>cd sbin

C:\Hadoop\hadoop-3.2.1\sbin>start-all.cmd
This script is Deprecated. Instead use start-dfs.cmd and start-yarn.cmd
starting yarn daemons

C:\Hadoop\hadoop-3.2.1\sbin>cd/

C:\>hadoop dfsadmin -safemode leave
DEPRECATED: Use of this script to execute hdfs command is deprecated.
Instead use the hdfs command for it.
Safe mode is OFF

C:\>hadoop fs -mkdir /input_dir

C:\>hadoop fs -put C:/file1.txt /input_dir
2024-01-03 03:30:10,102 INFO sasl.SaslDataTransferClient: SASL encryption trust check: localhostTrusted = false, remoteHostTrusted = false

C:\>hadoop fs -ls /input_dir/
Found 1 items
-rw-r--r-- 1 anasz supergroup          76 2024-01-03 03:30 /input_dir/file1.txt

C:\>
```

- 2.5. Verify the content of the copied file.

```
C:\>hadoop dfs -cat /input_dir/file1.txt
DEPRECATED: Use of this script to execute hdfs command is deprecated.
Instead use the hdfs command for it.
2024-01-03 03:31:48,554 INFO sasl.SaslDataTransferClient: SASL encryption trust check: localhostTrusted = false, remoteH
ostTrusted = false
Install Hadoop
Run Hadoop Wordcount Mapreduce Example
I love Hadoop
Yarn
C:\>_
```

- ## 2.6. Run MapReduceClient.jar and also provide input and out directories.

```
Administrator: Command Prompt

Map output materialized bytes=119
Input split bytes=105
Combine input records=11
Combine output records=9
Reduce input groups=9
Reduce shuffle bytes=119
Reduce input records=9
Reduce output records=9
Spilled Records=18
Shuffled Maps =1
Failed Shuffles=0
Merged Map outputs=1
GC time elapsed (ms)=62
CPU time spent (ms)=0
Physical memory (bytes) snapshot=0
Virtual memory (bytes) snapshot=0
Total committed heap usage (bytes)=404750336

Shuffle Errors
BAD_ID=0
CONNECTION=0
IO_ERROR=0
WRONG_LENGTH=0
WRONG_MAP=0
WRONG_REDUCE=0


File Input Format Counters
  Bytes Read=76
File Output Format Counters
  Bytes Written=77

C:\>
```

2.7. Verify content for the generated output file.

```
C:\>hadoop dfs -cat /output_dir/*
DEPRECATED: Use of this script to execute hdfs command is deprecated.
Instead use the hdfs command for it.
2024-01-03 03:34:44,023 INFO sas1.SaslDataTransferClient: SASL encryption trust check: localhostTrusted = false, remoteH
ostTrusted = false
Example 1
Hadoop 3
I 1
Install 1
Mapreduce 1
Run 1
Wordcount 1
Yarn 1
love 1

C:\>
```



All Applications

Cluster

About

Nodes

Node Labels

Applications

NEW

NEW SAVING

SUBMITTED

ACCEPTED

RUNNING

FINISHED

FAILED

KILLED

Scheduler

Tools

Cluster Metrics

Apps Submitted	Apps Pending	Apps Running	Apps Completed	Containers Running	Memory Used	Memory Total	Memory Reserved	VCores Used
1	0	0	1	0	0 B	8 GB	0 B	0

Cluster Nodes Metrics

Active Nodes	Decommissioning Nodes	Decommissioned Nodes	Lost Nodes	Unhealthy Nodes	Rebooted Nodes
1	0	0	0	0	0

Scheduler Metrics

Scheduler Type	Scheduling Resource Type	Minimum Allocation	Maximum Allocation	Maximum
Capacity Scheduler	[memory-mb (unit=M), vcores]	<memory:1024, vCores:1>	<memory:8192, vCores:4>	0

Show: 20 entries

ID	User	Name	Application Type	Queue	Application Priority	StartTime	LaunchTime	FinishTime	State	FinalStatus	Running Containers	Allocated CPU VCoers	Allocated Memory MB	Reserved CPU VCoers	Reserved Memory MB	% of Queue
application_1704223424961_0001	anasz	word count	MAPREDUCE	default	0	Wed Jan 3 03:33:12 +0800 2024	Wed Jan 3 03:33:14 +0800 2024	Wed Jan 3 03:33:29 +0800 2024	FINISHED	SUCCEEDED	N/A	N/A	N/A	N/A	N/A	0.0

Showing 1 to 1 of 1 entries

Part 2: Data Analytics

2.1 Data Analytics Tasks

- Dataset link:
<https://www.kaggle.com/datasets/nezukokamaado/road-accident-casualties-dataset/data>

2.1.1 Pre-processing

- Importing essential libraries for data processing, visualization and statistics

```
# Packages
# Data Processing
import numpy as np
import pandas as pd
# Visualization
import matplotlib.pyplot as plt
plt.rcParams['figure.dpi'] = 200
import seaborn as sns
# Statistics
import math
from scipy import stats
from scipy.stats import norm
# # Deep Learning
# import tensorflow as tf
# File Path
import os
for dirname, _, filenames in os.walk('/kaggle/input'):
    for filename in filenames:
        print(os.path.join(dirname, filename))
```

- Defining root path and a random seed for reproducibility

```
# setting
path_root = "C:/Users/Asus/Desktop/UPM/sem 7/BDA/project/"
seed = 394
```

- Adjusting Pandas display settings

```
# pandas display setting
pd.set_option('display.max_columns', 200)
```

- Reading a CSV file from a specified path into a Pandas DataFrame

```
df_accident = pd.read_csv(path_root + "caraccident.csv")
```

- Standardize column names for easier referencing

```
# rename
df_accident.rename(columns = {'Accident Date': 'Accident_Date', 'District Area': 'District_Area'}, inplace = True)
```


- Categorical features processing

```
# categorical features
list_categorical_features = [
    'Accident_Severity', 'Light_Conditions', 'District_Area', 'Road_Surface_Conditions',
    'Road_Type', 'Urban_or_Rural_Area', 'Weather_Conditions', 'Vehicle_Type'
]
df_accident[list_categorical_features] = df_accident[list_categorical_features].astype('category')
df_accident.info()
```

2.1.2 Data Cleaning

- Refining the dataset by selecting relevant columns, converting a date column to DateTime, handling duplicates and missing values, removing specific rows, and verifying the resulting DataFrame's structure.

```
# subsetting
df_accident = df_accident[[
    # 'Index',
    'Accident_Severity', 'Accident_Date', 'Latitude',
    'Light_Conditions', 'District_Area', 'Longitude',
    'Number_of_Casualties', 'Number_of_Vehicles', 'Road_Surface_Conditions',
    'Road_Type', 'Urban_or_Rural_Area', 'Weather_Conditions',
    'Vehicle_Type'
]].copy()

# to_datetime
df_accident['Accident_Date'] = pd.to_datetime(df_accident['Accident_Date'], dayfirst = True)
df_accident.head()

# dropping duplicated rows and missing values
df_accident = df_accident.loc[~df_accident.duplicated()].reset_index(drop = True).copy()
df_accident = df_accident.dropna()

# remove "Unallocated" rows(only 3 cases)
df_accident = df_accident.loc[
    (df_accident['Urban_or_Rural_Area'] == "Urban") |
    (df_accident['Urban_or_Rural_Area'] == "Rural")
]

# check
df_accident.info()
```

2.1.3 Modeling

- Imports modules from the scikit-learn library for machine learning tasks.

```
from sklearn.model_selection import train_test_split
from sklearn.metrics import accuracy_score, classification_report, confusion_matrix

from sklearn.linear_model import LogisticRegression
from sklearn.tree import DecisionTreeClassifier
from sklearn.ensemble import RandomForestClassifier
from sklearn.neighbors import KNeighborsClassifier
```

- Prepares the feature set for modeling by removing unnecessary columns, creates a new DataFrame `df_X` with one-hot encoding for categorical variables to avoid multicollinearity, and sets up a new DataFrame `df_y` for the target variable "Accident_Severity".

```
df_accident.drop([
    "Accident_Severity", "Accident_Date", "District_Area", "Number_of_Casualties", "Number_of_Vehicles"
], axis = 1)

df_X = df_accident.drop([
    "Accident_Severity", "Accident_Date", "District_Area", "Number_of_Casualties", "Number_of_Vehicles"
], axis = 1)
df_X = pd.get_dummies(df_X, columns = [
    "Light_Conditions", "Road_Surface_Conditions", "Road_Type", "Urban_or_Rural_Area", "Weather_Conditions", "Vehicle_Type"
], drop_first = True)

df_y = df_accident["Accident_Severity"]

df_X.head()
```

- Splits the datasets into training and validation sets based on the target variable "Accident_Severity" to ensure a proportional representation of classes in both sets.

```
X_tr, X_val, y_tr, y_val = train_test_split(df_X, df_y, test_size = 0.3, random_state = seed, stratify = df_y)
```

- Initializes four classification models: Logistic Regression (`model_lr`), Decision Tree (`model_dt`), Random Forest (`model_rf`), and k-Nearest Neighbors (`model_knn`).

```
list_model = [model_lr, model_dt, model_rf, model_knn]
```

- Iterates through the list of models, fitting each model on the training data (`X_tr, y_tr`). Predicts the target variable on the validation data (`X_val`) and lastly prints the classification report, providing precision, recall, F1-score, and accuracy metrics for each model.

```

for model in list_model:

    print("")
    print(str(model))

    model.fit(X_tr, y_tr)
    y_pred = model.predict(X_val)

    report = classification_report(y_val, y_pred)
    print(report)

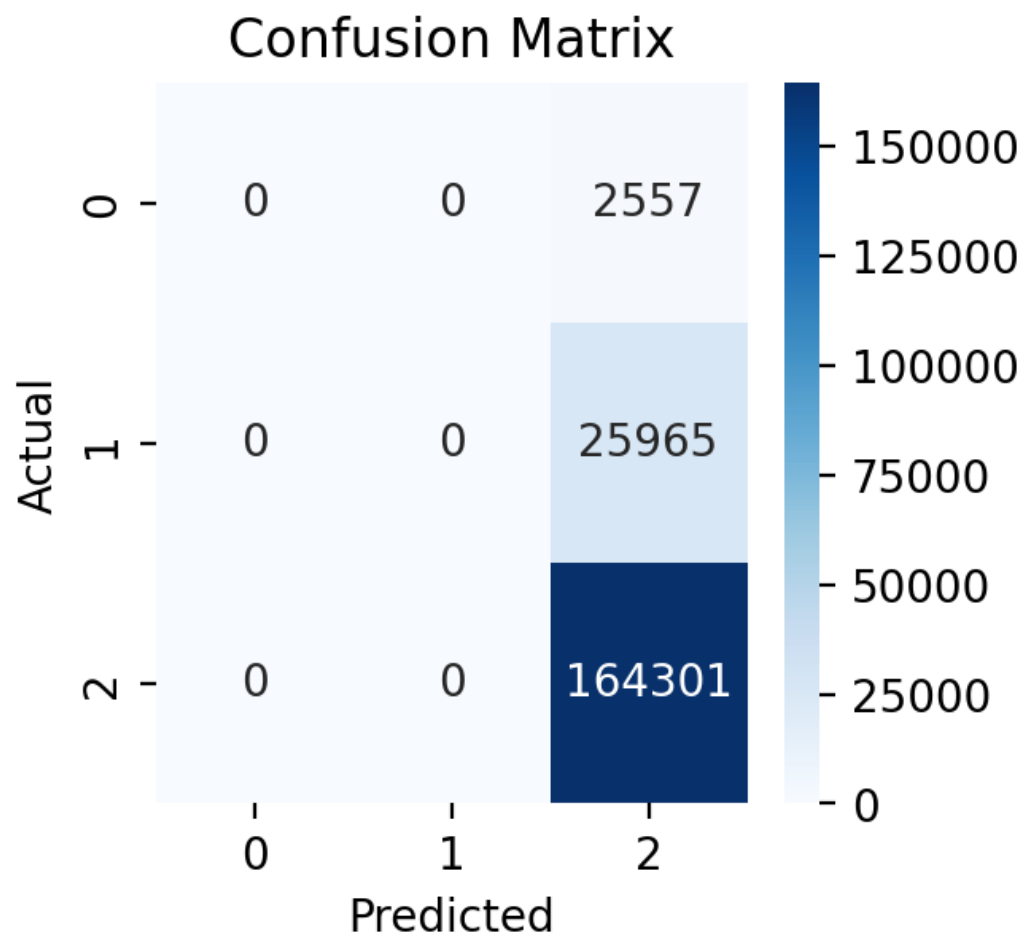
    # confusion matrix
    temp_confusion_matrix = confusion_matrix(y_val, y_pred)
    plt.figure(figsize = (3, 3), facecolor = "white")
    sns.heatmap(
        temp_confusion_matrix,
        annot = True, fmt = 'd', cmap = 'Blues'
    )
    plt.xlabel('Predicted')
    plt.ylabel('Actual')
    plt.title('Confusion Matrix')
    plt.show()

```

2.1.4 Results

- Confusion Matrix (Logistic Regression)

This model performed well overall with an accuracy of 85% but struggled when it came to predicting Fatal and Serious accidents as well as did not correctly identify any instances of these severe outcomes, resulting in 0% precision for both. However, the model excelled in correctly classifying less severe accidents, resulting in a 100% recall and a solid F1-score of 92% for the Slight class. While it may have limitations in dealing with more severe incidents, it proves reliable in recognizing less critical accidents.

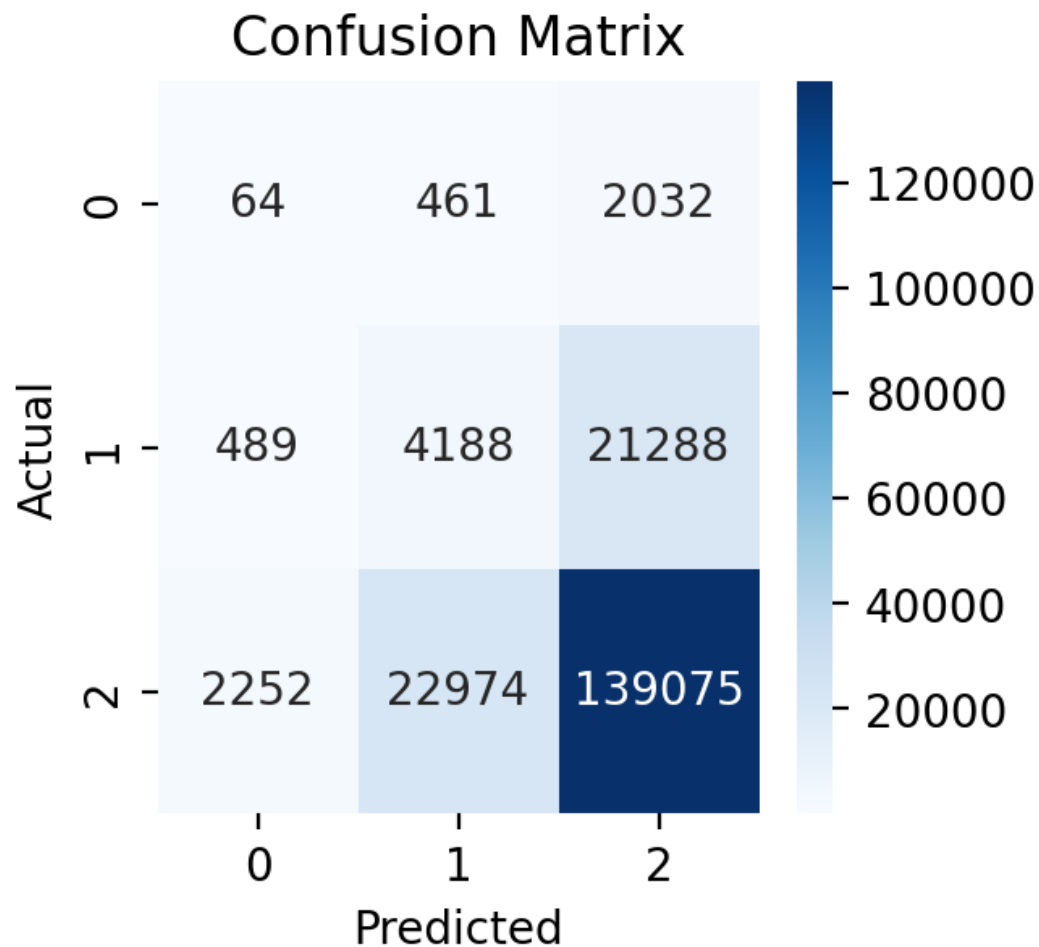


Output:

	precision	recall	f1-score	support
Fatal	0.00	0.00	0.00	2557
Serious	0.00	0.00	0.00	25965
Slight	0.85	1.00	0.92	164301
accuracy			0.85	192823
macro avg	0.28	0.33	0.31	192823
weighted avg	0.73	0.85	0.78	192823

- Confusion Matrix (Decision Tree)

With a 74% accuracy, the model improved in identifying Fatal and Serious accidents compared to Logistic Regression. However, it still had low precision for these severe cases (2% for Fatal and 15% for Serious), indicating the need for enhancements, especially in handling more severe accidents.

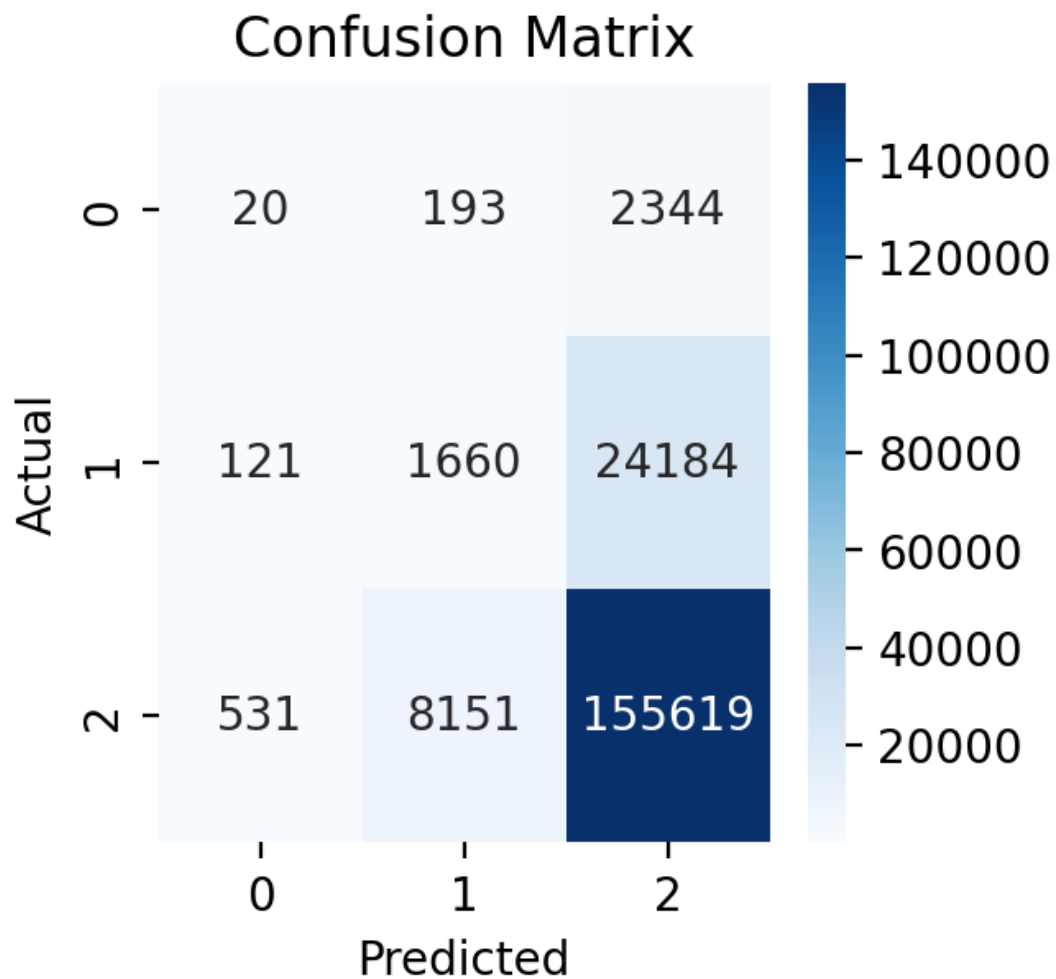


Output:

	precision	recall	f1-score	support
Fatal	0.02	0.03	0.02	2557
Serious	0.15	0.16	0.16	25965
Slight	0.86	0.85	0.85	164301
accuracy			0.74	192823
macro avg	0.34	0.34	0.34	192823
weighted avg	0.75	0.74	0.75	192823

- Confusion Matrix (Random Forest)

The model achieved an overall accuracy of 82%, showing better precision for Fatal (3%) and Serious (17%) accidents compared to the Decision Tree. However, it still has limitations in precision for these severe cases, and the model struggles to capture and classify instances of severe accidents accurately, as indicated by lower recall scores.

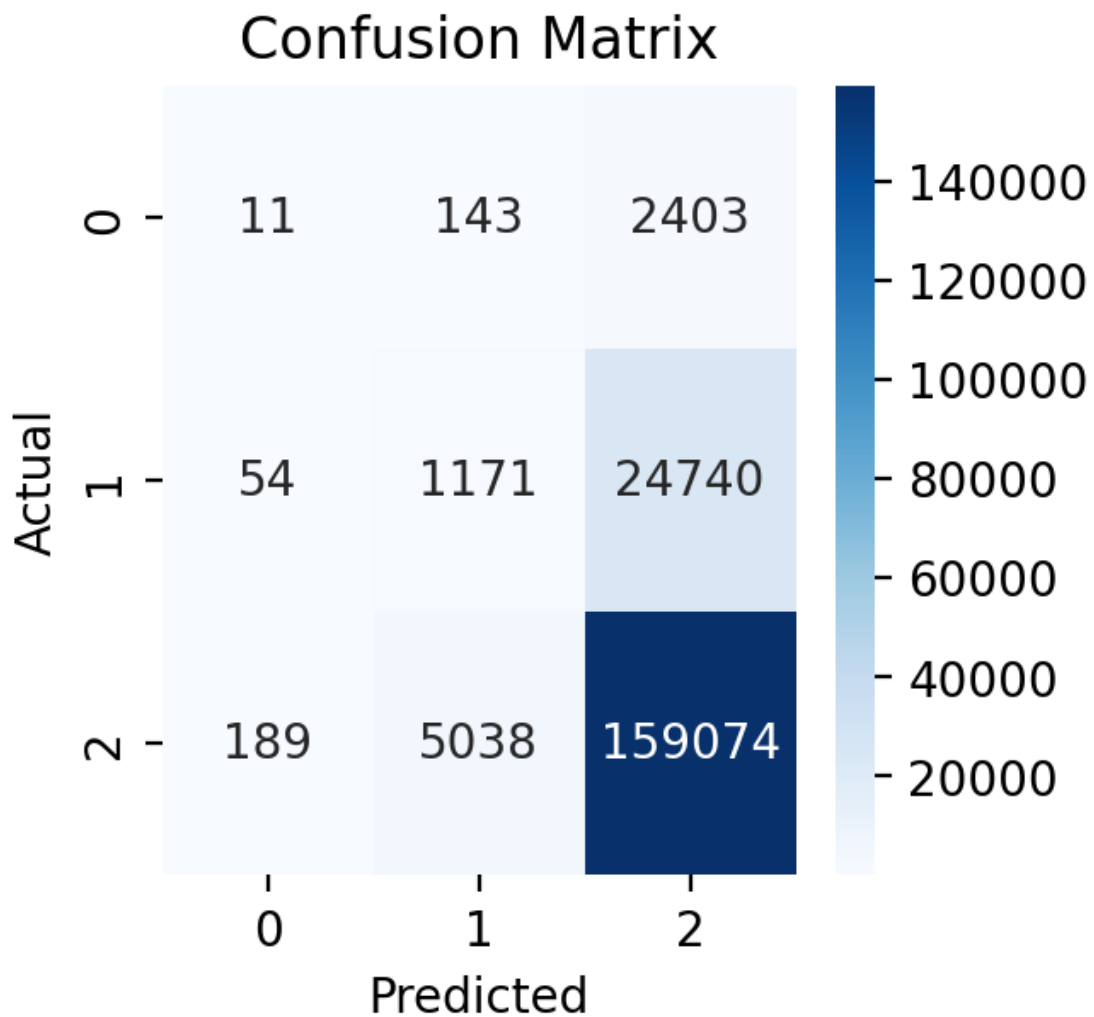


Output:

	precision	recall	f1-score	support
Fatal	0.03	0.01	0.01	2557
Serious	0.17	0.06	0.09	25965
Slight	0.85	0.95	0.90	164301
accuracy			0.82	192823
macro avg	0.35	0.34	0.34	192823
weighted avg	0.75	0.82	0.78	192823

- Confusion Matrix (K-Nearest Neighbor)

With an 83% accuracy, KNN showed relevant results in predicting all accident severity levels. It had balanced precision, recall, and F1-score for each class. Thus, KNN outperformed Logistic Regression, Decision Tree, and Random Forest in recall scores for Fatal and Serious accidents.

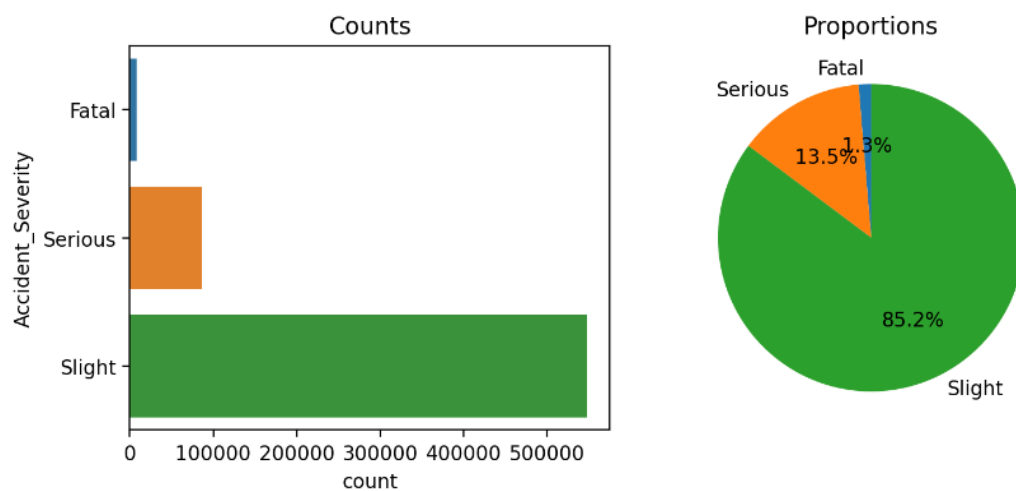


Output:

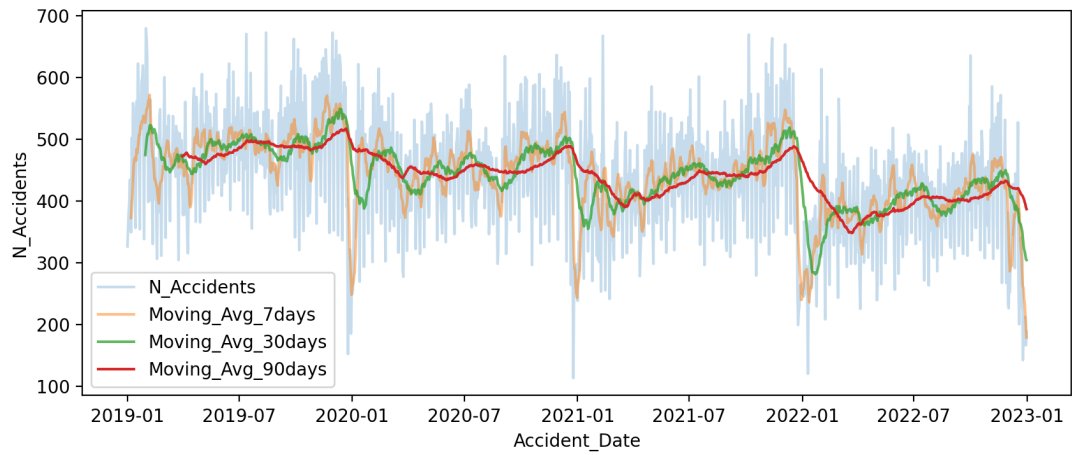
	precision	recall	f1-score	support
Fatal	0.04	0.00	0.01	2557
Serious	0.18	0.05	0.07	25965
Slight	0.85	0.97	0.91	164301
accuracy			0.83	192823
macro avg	0.36	0.34	0.33	192823
weighted avg	0.75	0.83	0.78	192823

- Result Visualization

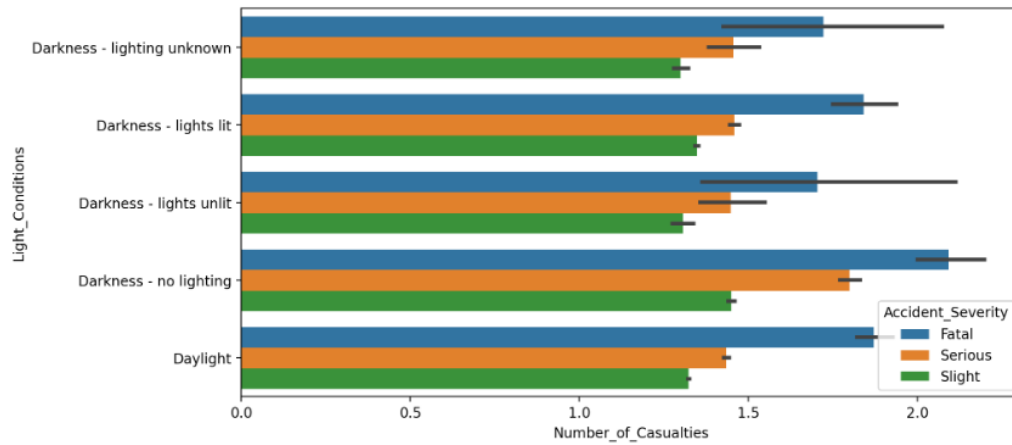
Distribution of: Accident_Severity

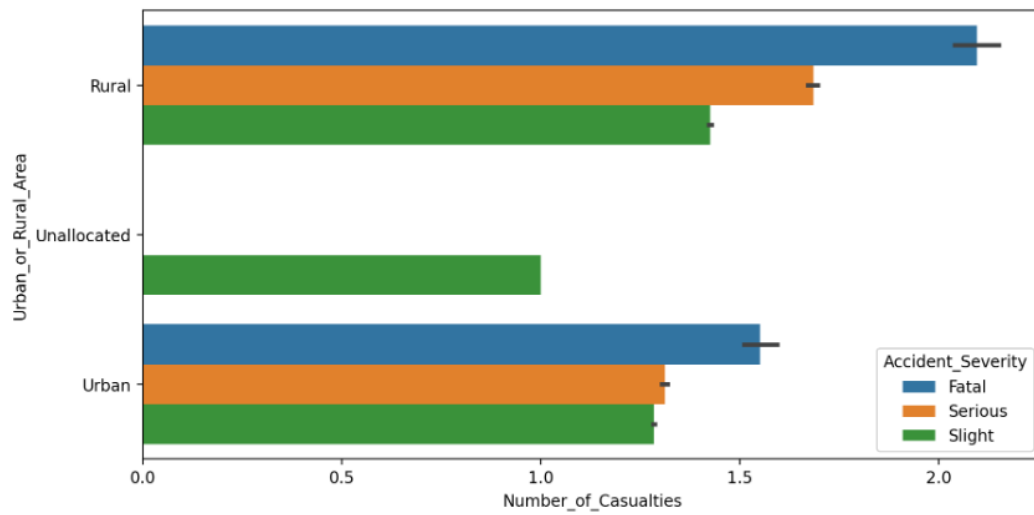
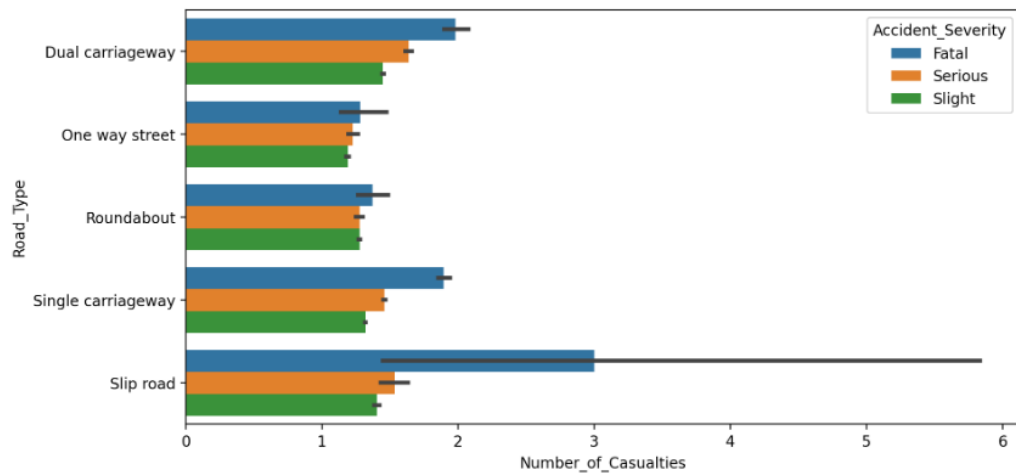
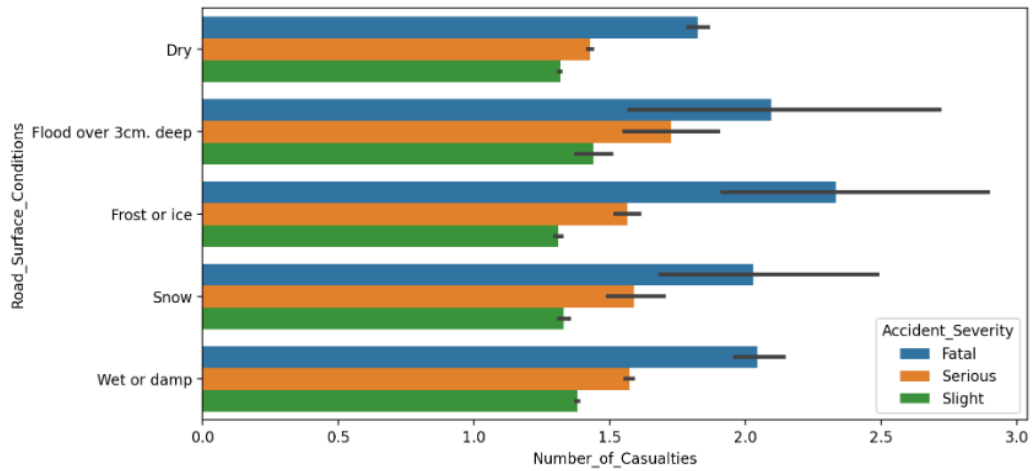


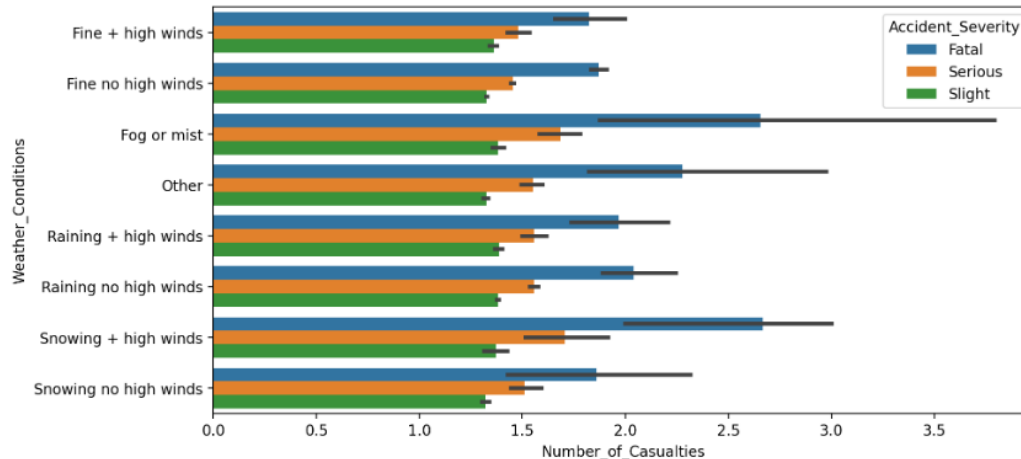
The majority of accidents, comprising 85.2%, are categorized as "Slight," denoting events with minor consequences. A lesser but notable proportion, accounting for 13.5%, falls under the "Serious" category, signifying incidents of greater impact. The least prevalent, at 1.3%, are "Fatal" accidents, emphasizing their infrequency and heightened gravity.



From the visualization of the trend in the number of car accidents occurring, there is a general decrease in the total incidents. It is however still noteworthy that specific periods throughout the year exhibit a distinct reduction in accident occurrences.







The boxplots showcase the impact of diverse variables—namely, light conditions, road surface conditions, road type, weather conditions, and urban or rural areas—on the severity of accidents. Within each distinct condition, these boxplots specify the statistical distribution of casualties across different severity levels, thereby offering a comprehensive visualization of accident occurrences. This results in visualization establishes a foundational basis for future research attempts focused on predicting the risk of car accidents by examining the interplay of these environmental conditions.

2.1.5 Summary

The results regarding critical aspects of road incidents offer valuable recommendations and insights to enhance road safety. Recognizing the severity of accidents is deemed crucial for providing effective road management strategies as well as the significance of understanding where and when accidents occur, guiding interventions in specific regions and timely safety measures. A comprehensive dataset is advocated, serving as a relevant foundation for ongoing research and informed policymaking. Weather and road conditions' impact on accident rates is highlighted, urging the development of weather-responsive safety protocols and awareness initiatives.

By identifying accident hotspots and associated risk factors, targeted preventative measures can be implemented to effectively allocate resources to high-risk areas. These results suggest utilizing data-driven techniques, such as predictive modelling, to proactively tackle road safety issues. Additionally, incorporating traffic collision analysis into urban planning can help create safer urban environments through

improved road design and infrastructure. Another key factor is understanding patterns of driver behavior and the importance of applicable educational campaigns and regulations to improve overall road safety. In conclusion, these results offer a comprehensive roadmap for evidence-based interventions, showcasing the value of a refined and targeted approach.

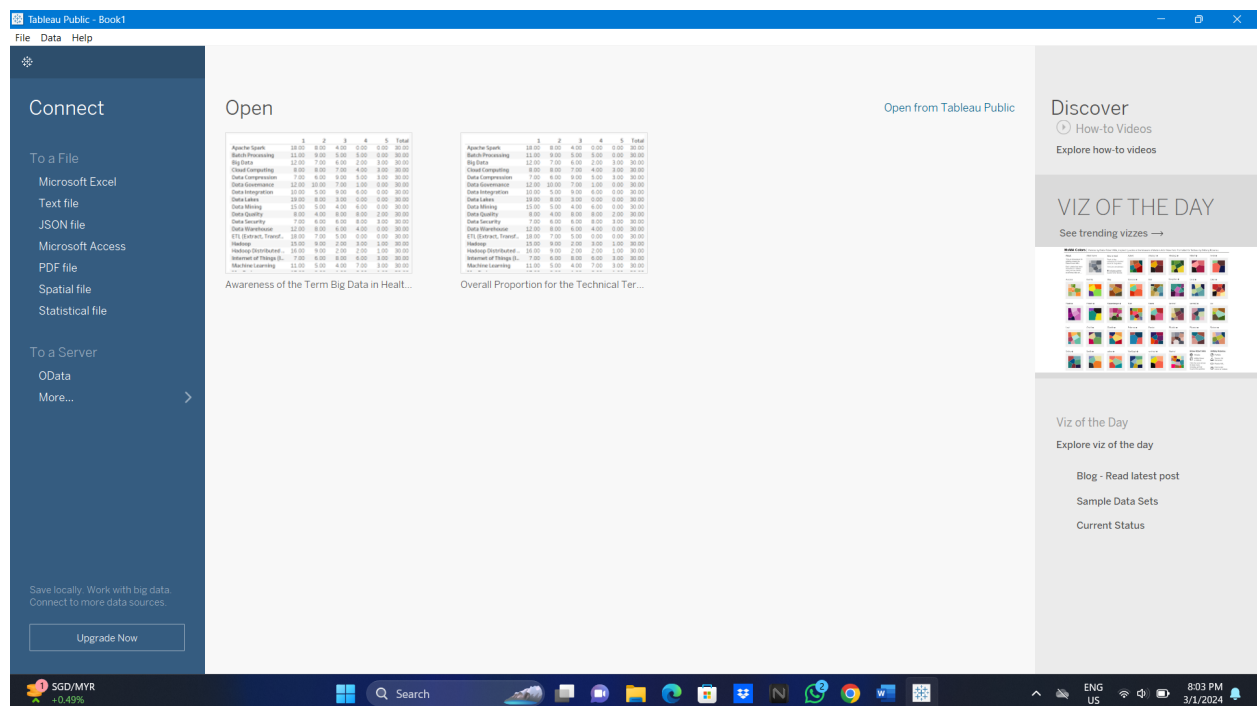
Part 3: Data Visualization

3.1 Public Tableau Installation

This part provides evidence of the successful installation of Public Tableau on our machine. Part 3 covers installation details, creation of a Tableau account on a Public Tableau Site, dataset analysis, data relationships analysis and a combination of Tableau worksheets within a single story.

3.1.1 Installation Details

- The start Screen of Tableau is shown:




3.1.2 Creation of Public Tableau account in Public Tableau Site

tableau public Create Learn

Learn all about Tableau Community Projects, exploring vizzes on Tableau Public, and more ways to connect with the DataFam. [Viz a little](#) →

Customize Banner




Nik Muhammad Asyraf Nik Ismail He/Him
Student at Universiti Putra Malaysia | Serdang, Selangor, Malaysia

[Edit Profile](#)

Vizzes 2 Favorites 0 Following 0 Followers 0 Stats [Create a Viz](#)

Awareness of the Term Big Data in Healthcare

Job Title

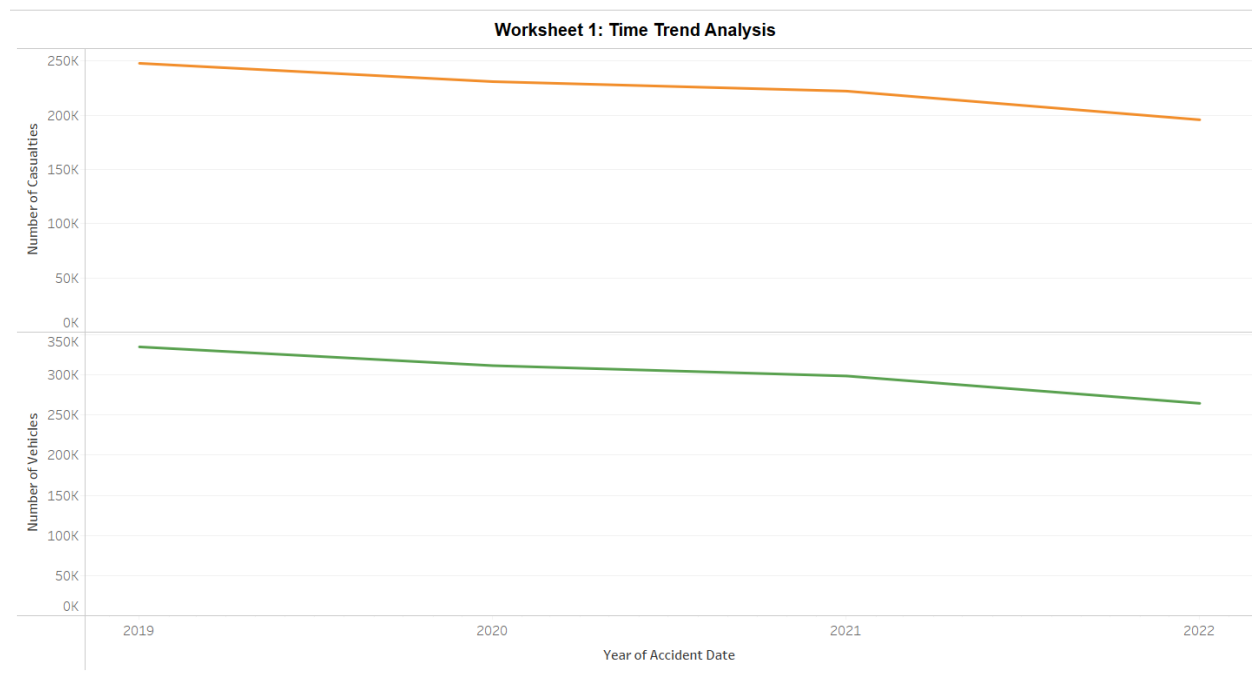


<Overall Proportion for the Technical Terms>

Technical Terms	1	2	3	4	5	Total
Apache Spark	10.00	0.00	4.00	0.00	0.00	10.00
Batch Processing	11.00	0.00	0.00	0.00	0.00	11.00
Hadoop	12.00	0.00	0.00	0.00	0.00	12.00
Cloud Computing	8.00	0.00	7.00	0.00	0.00	15.00
Data Compression	7.00	0.00	0.00	0.00	0.00	7.00
Data Governance	12.00	0.00	7.00	0.00	0.00	19.00
Data Integration	10.00	0.00	0.00	0.00	0.00	10.00

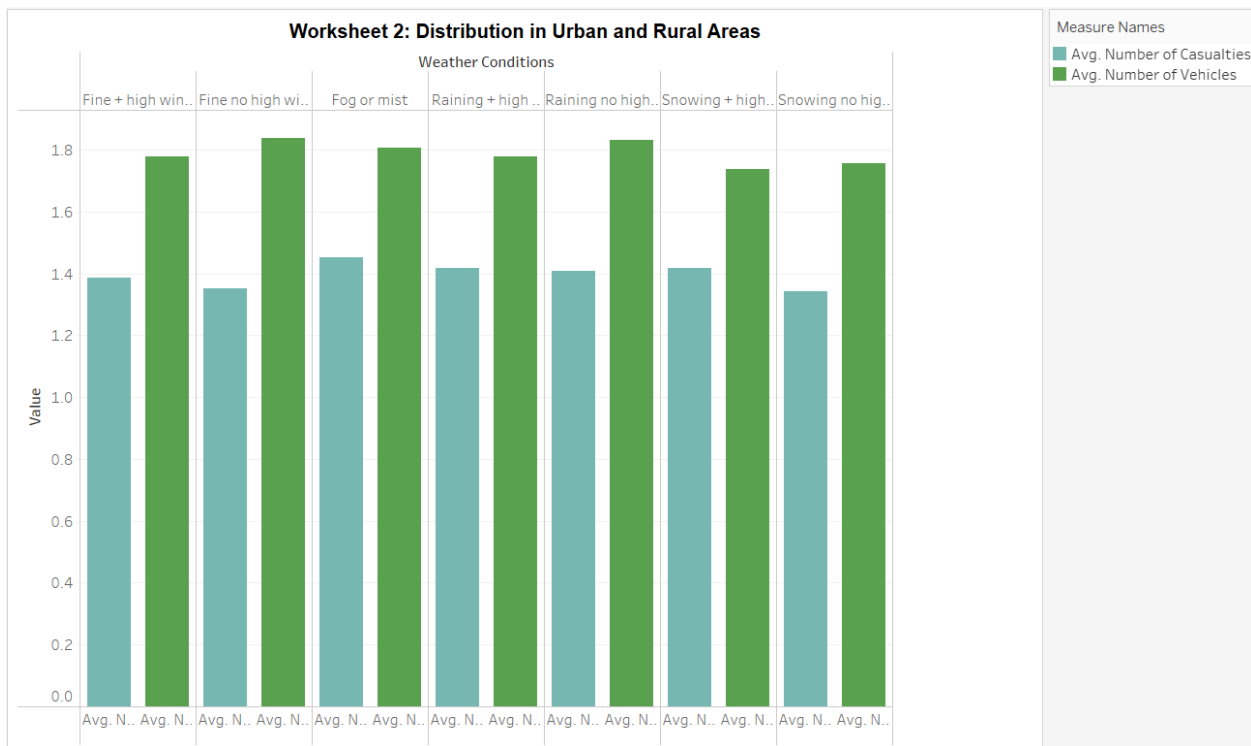
3.2 Relationships in The Chosen Dataset with Tableau Analysis

3.2.1 Worksheet 1: Time Trend Analysis



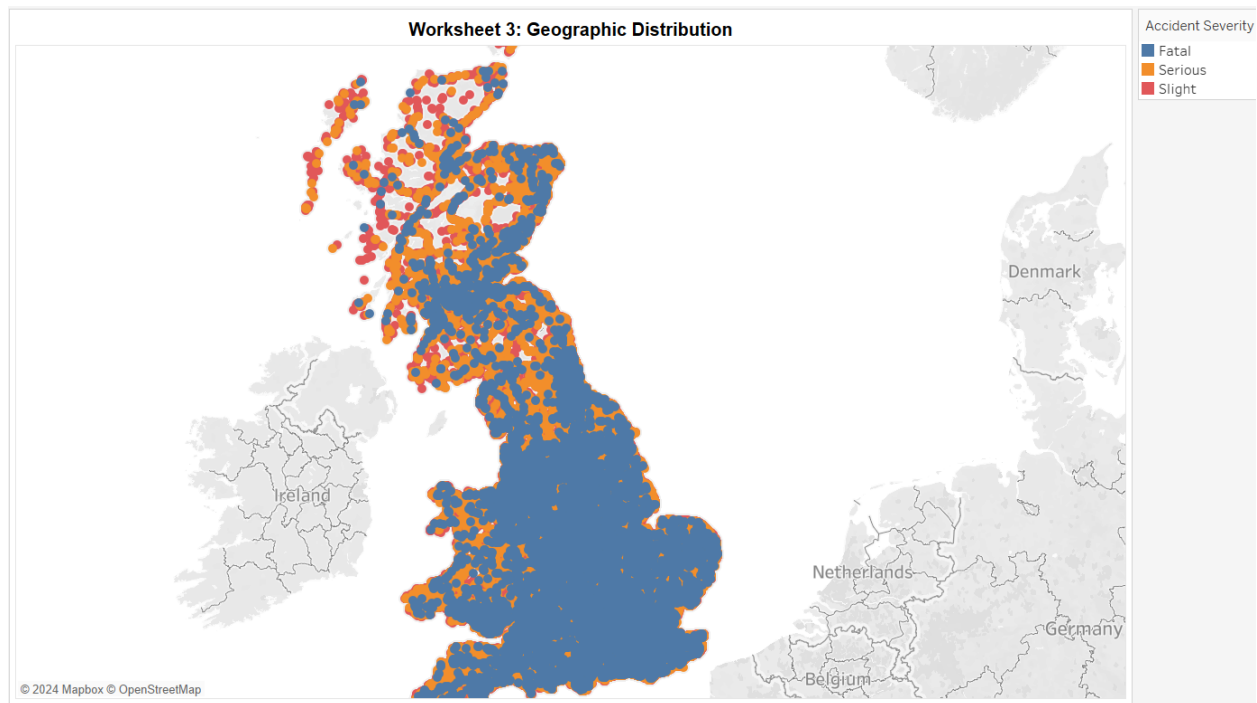
In Worksheet 1, the Time Trend Analysis graph presents a comprehensive overview of accident characteristics over the observed period. The x-axis represents the timeline, with years indicated for each point. The y-axes illustrate the total number of casualties and vehicles involved in accidents. Notably, both lines exhibit a gradual decrease over time, suggesting a positive trend in reducing the overall number of casualties and vehicles involved in accidents. This encouraging pattern may indicate successful interventions, improved safety measures, or changing traffic dynamics. Further analysis and correlation with external factors can provide deeper insights into the causes behind this positive trajectory, aiding in the formulation of effective accident prevention strategies.

3.2.2 Worksheet 2: Distribution in Urban and Rural Areas



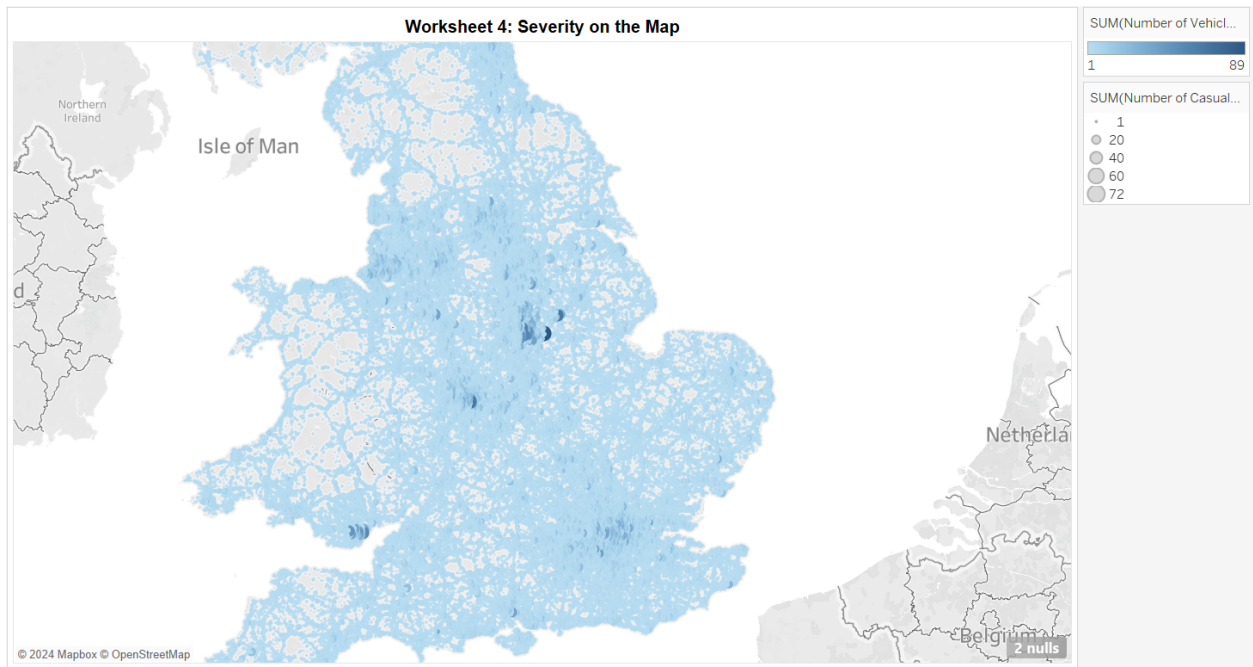
The Distribution in Urban and Rural Areas utilizes a clustered bar chart to highlight the impact of different weather conditions on accident characteristics. The x-axis displays various weather conditions, while the y-axis represents the average number of casualties and vehicles involved in accidents. The clustered bars showcase a clear distinction in the measure values, with an average of 1.4 for casualties and 1.8 for vehicles across different weather conditions. This chart allows for a quick comparison of the influence of weather on both casualties and vehicles, providing valuable insights for traffic management and safety measures.

3.2.3 Worksheet 3: Geographic Distribution



In Worksheet 3, the Geographic Distribution map employs longitude and latitude as columns and rows, respectively, creating a symbol map to visualize the spatial distribution of accidents across the United Kingdom. The color mark is used to represent the severity of accidents, with blue indicating fatal incidents, orange for serious, and red for slight. The map reveals a distinct pattern where fatal accidents are predominantly concentrated in the central to southern regions, extending towards the south, while the northern areas are coded with a higher prevalence of serious incidents.

3.2.4 Worksheet 4: Severity on the Map



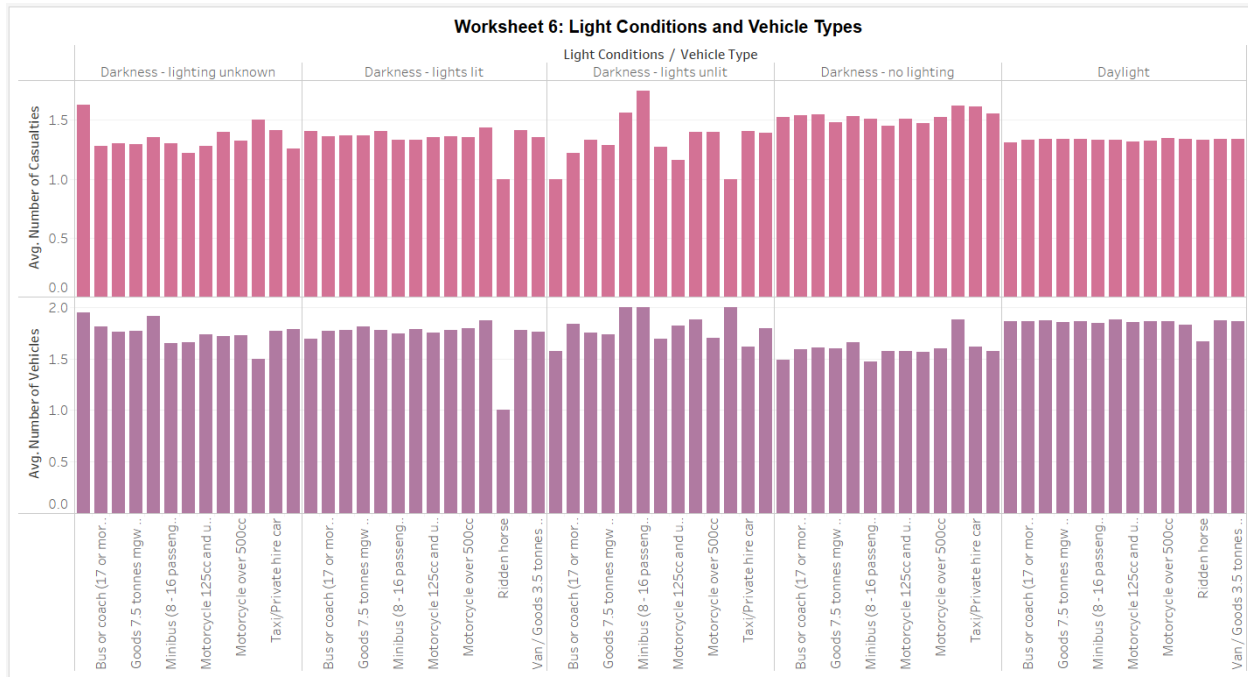
In Worksheet 4, the Severity on the Map graph utilizes latitude and longitude as rows and columns, respectively, creating a symbol map to visually represent the severity and impact of accidents. The filter allows users to explore accident severity levels. The color code, in blue, indicates the total number of vehicles involved in accidents, while the size mark represents the total number of casualties. Larger symbols denote a higher number of casualties, providing an immediate visual cue to the severity of accidents in specific geographic locations.

3.2.5 Worksheet 5: Factors by Severity



The Factors by Severity graph provides a comprehensive view of the distribution of accidents based on road surface conditions and road types. The graph employs stacked bars with AVG(Number of Casualties) and AVG(Number of Vehicles) as rows, and Road Surface Condition as columns. The filter allows users to explore the impact of different road types on accident severity levels. The color-coded bars vividly represent various road types, such as blue for dual carriageways, orange for one-way streets, red for roundabouts, light teal for single carriageways, and green for slip roads. The stacked structure of the bars allows for a clear comparison of the contribution of each road surface condition to different severity levels, offering insights into the factors influencing accident outcomes and facilitating targeted interventions for specific road types.

3.2.6 Worksheet 6: Light Conditions and Vehicle Types

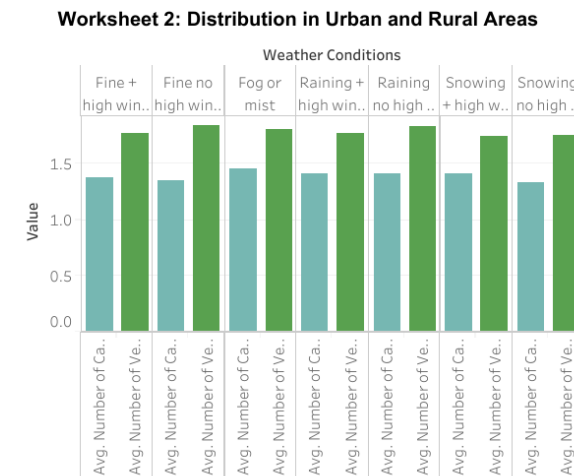
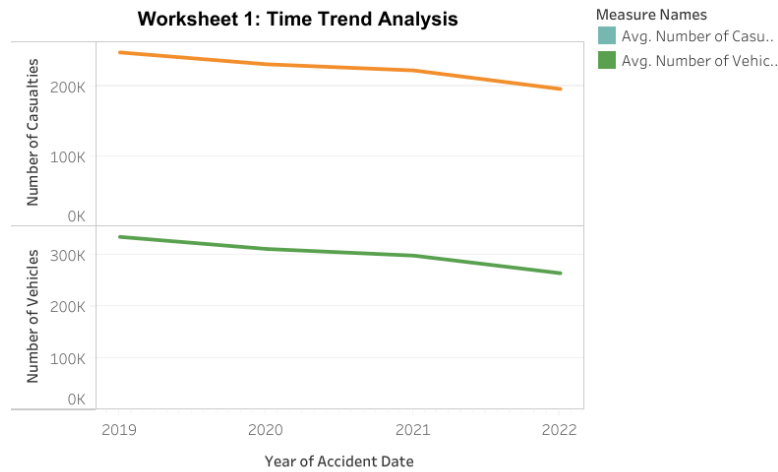
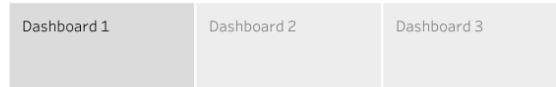


In Worksheet 6, the Light Conditions and Vehicle Types graph presents a detailed analysis of the distribution of accidents based on light conditions and vehicle types. The graph utilizes a side-by-side bar chart layout with AVG(Number of Casualties) and AVG(Number of Vehicles) as rows, Light Conditions as columns, and Vehicle Types as a filter. This visualization enables a nuanced exploration of how different light conditions impact accident severity and the involvement of various vehicle types. By segregating the data into distinct bars for each combination of light condition and vehicle type, the graph provides a comprehensive understanding of the relationships between these factors. Policymakers and safety experts can derive valuable insights from this visual representation to develop targeted strategies for improving road safety under specific conditions and vehicle scenarios.

3.3 Combination of Worksheets into a Single Story

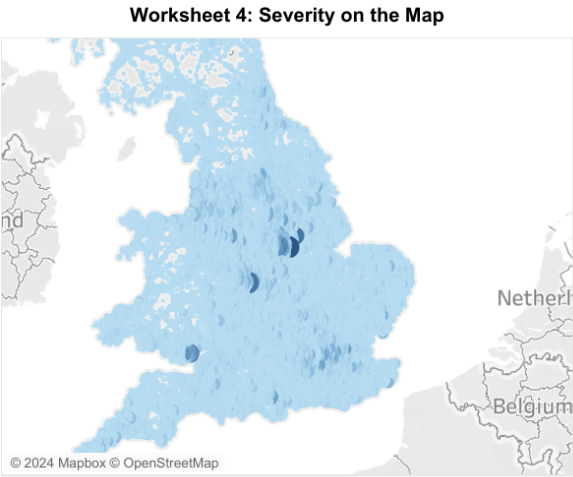
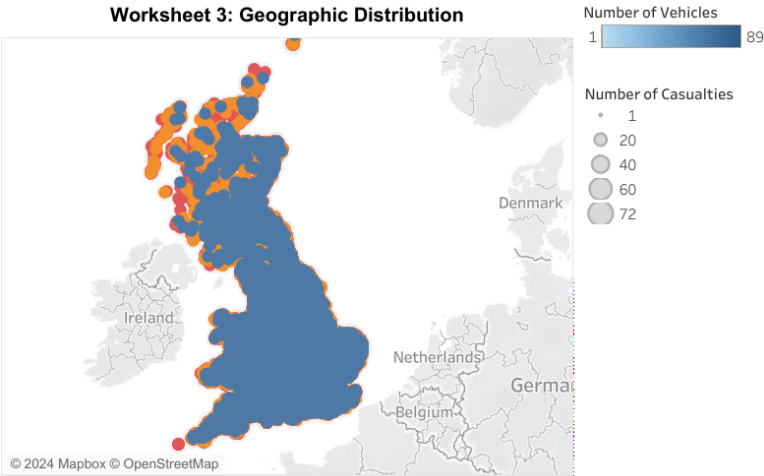
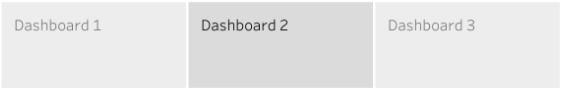
3.3.1 Dashboard 1

Story 1



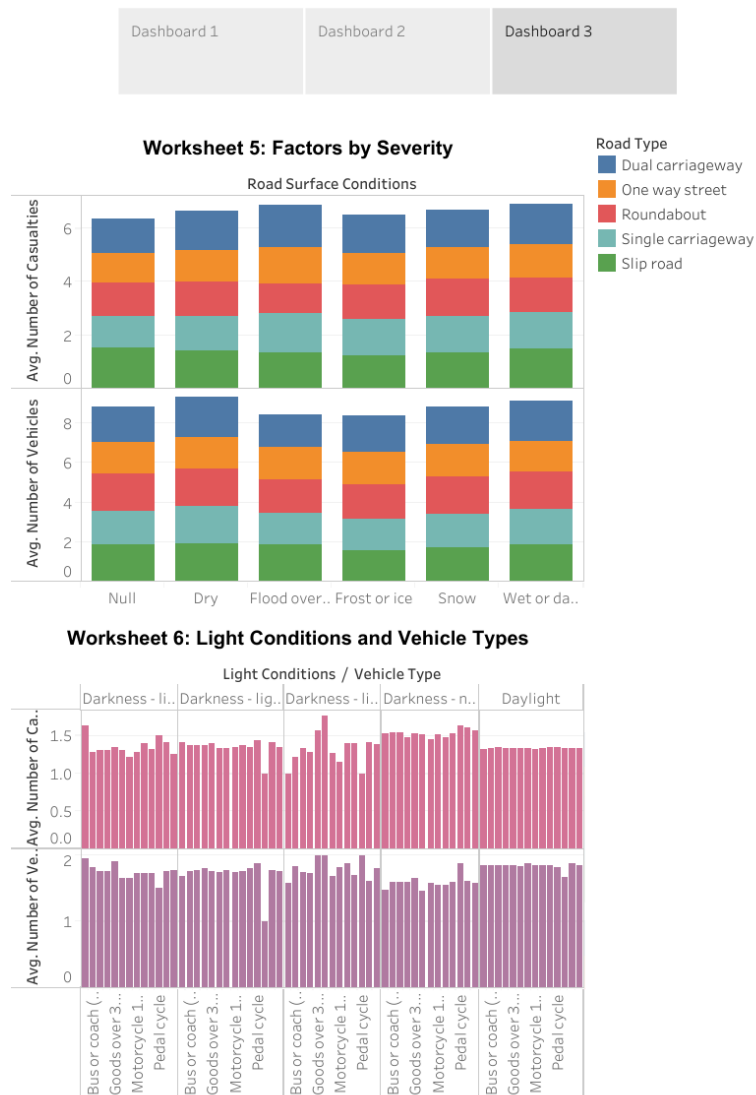
3.3.2 Dashboard 2

Story 1



3.3.3 Dashboard 3

Story 1



3.4 Story Publication in Tableau Public Page

- All Tableau worksheets were combined within a single Story, comprising 3 Dashboards. Each dashboard contains two worksheets. Link of Tableau publication:
https://public.tableau.com/views/GroupProjectPart3_17043908217930/Story1?:language=en-US&:display_count=n&:origin=viz_share_link