# Configuration Manual

MSc Research Project
Msc. in Data Analytics

## Melvin Akash Ambrose MohanDoss
Student ID: x22152601

School of Computing
National College of Ireland

Supervisor: Dr. Anu Sahni

## National College of Ireland

## MSc Project Submission Sheet

## School of Computing

| | |
|---|---|
| **Student Name:** | Melvin Akash AmbroseDoss |
| **Student ID:** | x22152601 |
| **Programme:** | Msc in Data Analytics                **Year:** 2023 -2024 |
| **Module:** | Msc Research Project |
| **Lecturer:** | Dr. Anu Sahni |
| **Submission Due Date:** | 14/12/2023 |
| **Project Title:** | Implementing Machine Learning Models for Predicting Road Accident Severity in Northern Ireland |
| **Word Count:** 730 | **Page Count:** 7 |

I hereby certify that the information contained in this (my submission) is information pertaining to research I conducted for this project.  All information other than my own contribution will be fully referenced and listed in the relevant bibliography section at the rear of the project.

<u>ALL</u> internet material must be referenced in the bibliography section.  Students are required to use the Referencing Standard specified in the report template.  To use other author's written or electronic work is illegal (plagiarism) and may result in disciplinary action.

| | |
|---|---|
| **Signature:** | Melvin Akash Ambrose MohanDoss |
| **Date:** | 14th December 2023 |

**PLEASE READ THE FOLLOWING INSTRUCTIONS AND CHECKLIST**

| | |
|---|---|
| Attach a completed copy of this sheet to each project (including multiple copies) | □ |
| **Attach a Moodle submission receipt of the online project submission,** to each project (including multiple copies). | □ |
| **You must ensure that you retain a HARD COPY of the project**, both for your own reference and in case a project is lost or mislaid.  It is not sufficient to keep a copy on computer. | □ |

Assignments that are submitted to the Programme Coordinator Office must be placed into the assignment box located outside the office.

| **Office Use Only** | |
|---|---|
| Signature: | |
| Date: | |
| Penalty Applied (if applicable): | |

# Configuration Manual

Melvin Akash Ambrose MohanDoss
Student ID: x22152601

# 1    Introduction

This Configuration Manual contains a list of all the requirements needed to replicate the study and its findings in a personal setting. All models constructed, data import and exploratory data analysis, data augmentation, and software and hardware requirements are covered.

# 2    System Specifications

This section covers Hardware and Software requirements.

## 2.1   Hardware Requirements



**Fig 1: Hardware requirement**

## 2.2   Software Requirements

- Jupyter Notebook (Version 6.5.2) or Google Colab
- Python (Version 3.10)
- MySql server (Version 8)
- Mysql Workbench

# 3    Data Collection

The Data is sourced from the UK government website.

Website link :
https://admin.opendatani.gov.uk/dataset?organization=police-service-of-northern-
ireland&tags=injury+collisions

# 4    Data Pre-processing

The total 9 datasets from the years 2020, 2021, 2022 are merged using Mysql and Python.

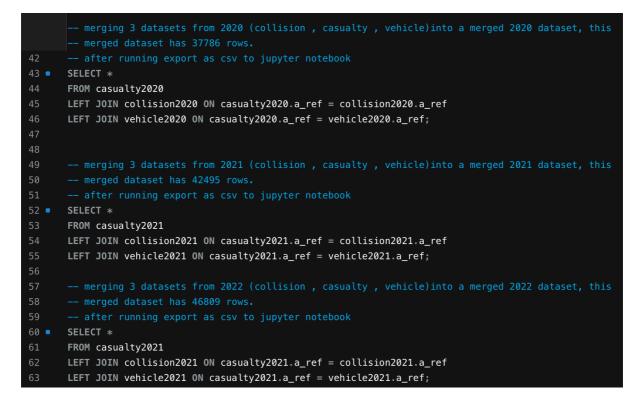First , the 9 datasets are loaded into MySql workbench and is merged year wise.

```sql
      -- merging 3 datasets from 2020 (collision , casualty , vehicle)into a merged 2020 dataset, this
      -- merged dataset has 37786 rows.
42    -- after running export as csv to jupyter notebook
43    SELECT *
44    FROM casualty2020
45    LEFT JOIN collision2020 ON casualty2020.a_ref = collision2020.a_ref
46    LEFT JOIN vehicle2020 ON casualty2020.a_ref = vehicle2020.a_ref;
47
48
49    -- merging 3 datasets from 2021 (collision , casualty , vehicle)into a merged 2021 dataset, this
50    -- merged dataset has 42495 rows.
51    -- after running export as csv to jupyter notebook
52    SELECT *
53    FROM casualty2021
54    LEFT JOIN collision2021 ON casualty2021.a_ref = collision2021.a_ref
55    LEFT JOIN vehicle2021 ON casualty2021.a_ref = vehicle2021.a_ref;
56
57    -- merging 3 datasets from 2022 (collision , casualty , vehicle)into a merged 2022 dataset, this
58    -- merged dataset has 46809 rows.
59    -- after running export as csv to jupyter notebook
60    SELECT *
61    FROM casualty2021
62    LEFT JOIN collision2021 ON casualty2021.a_ref = collision2021.a_ref
63    LEFT JOIN vehicle2021 ON casualty2021.a_ref = vehicle2021.a_ref;
```

**Fig 2: Merging the datasets year wise.**

The merged datasets from Mysql are then concatenated vertically in python.

```python
In [2]:  1  dataset_path_2020 = '/Users/melvinakash/Desktop/NCI/ric/datasets/merged_2020.csv'
         2  dataset_path_2021 = '/Users/melvinakash/Desktop/NCI/ric/datasets/merged_2021.csv'
         3  dataset_path_2022 = '/Users/melvinakash/Desktop/NCI/ric/datasets/merged_2022.csv'

In [3]:  1  pd.set_option('display.max_columns',None)
         2
         3  data_2020 = pd.read_csv(dataset_path_2020)
         4  data_2021 = pd.read_csv(dataset_path_2021)
         5  data_2022 = pd.read_csv(dataset_path_2022)

In [4]:  1  # Concatenate the DataFrames vertically (along rows)
         2  df = pd.concat([data_2020, data_2021 , data_2022], ignore_index=True)
```

**Fig 3: Concatenated vertically.**

# 5 Project Development

## 5.1 Importing Libraries

List of python libraries used :

```
In [29]:    1  import pandas as pd
            2  import numpy as np
            3  import missingno as msno
            4  from sklearn.model_selection import train_test_split
            5  from sklearn.preprocessing import StandardScaler
            6  from imblearn.over_sampling import SMOTE
            7  from sklearn.linear_model import LogisticRegression
            8  from sklearn.ensemble import RandomForestClassifier
            9  from sklearn.tree import DecisionTreeClassifier
           10  from sklearn.neighbors import KNeighborsClassifier
           11  from xgboost import XGBClassifier
           12  from sklearn.svm import SVC
           13  from sklearn.metrics import accuracy_score, classification_report
           14  from sklearn.neural_network import MLPClassifier
           15  from sklearn.metrics import confusion_matrix
           16  import matplotlib.pyplot as plt
           17  import seaborn as sns
           18  import plotly.express as px
```

**Fig 4: Libraries used**

Some of the main Libraries used in this project were Pandas , Numpy, Matplotlib and SMOTE.

## 5.2 Processing

- **Treatment of Missing Values:** The pattern is MCAR and 90% data were missing so the columns were deleted. .
- **Feature selection:** Generated histograms of all the data , and removed features which were unbalanced and could potentially lead to biasing.
- **Encoding :** One hot encoding and ordinal encoding are done to two variables a_District and a_wkday.
- The **final subset** of the filtered variables is shown in Figure 5.

```
In [10]:    1  #taking subset of necessary data after eval of columns
            2  data_subset = df[[ 'a_ref', 'a_District', 'a_type', 'a_veh', 'a_cas', 'a_wkday',
            3        'a_day', 'a_month', 'a_hour', 'a_min', 'a_gd1', 'a_gd2', 'a_ctype',
            4        'a_speed', 'c_class', 'c_sex', 'c_agegroup',
            5      'c_school', 'c_vtype', 'v_type', 'v_tow', 'v_man',
            6       'v_loc',  'v_impact',
            7       'v_sex', 'v_agegroup', 'v_hitr']]
```
**Fig 5: Subset of variables to ML models.**

.

- **Correlation** was checked with the target variable a_type.

3

The Data is split into train and test with a 80:20 split.

```
In [6]:    1  #over sampling using smote

In [7]:    1  pd.set_option('display.max_columns', None)
           2
           3  # Extract features and target variable
           4  X = data.drop('a_type', axis=1)  # Features
           5  y = data['a_type']  # Target variable
           6
           7  # Split the data into training and testing sets
           8  X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
           9
          10  # Apply SMOTE only to the training data
          11  smote = SMOTE(sampling_strategy='auto', random_state=42)
          12  X_train_synthetic, y_train_synthetic = smote.fit_resample(X_train, y_train)
          13
          14  # Combine the synthetic training data with the original training data
          15  X_train_combined = pd.concat([X_train, X_train_synthetic])
          16  y_train_combined = pd.concat([y_train, y_train_synthetic])
          17
          18
```

**Fig 5: SMOTE oversampler and test, train split.**

- The packages or libraries to perform the above tasks are shown in the Figure 4.

## 5.3 Modelling

- **Oversampling and test, train split:** The target variable is imbalanced , so we take samples of minority class and oversamples it .This is done using the SMOTE oversampler.

- The following code snippets contains implementation of four machine learning and one deep learning models.

- Each model is tuned with the best hyperparameters.

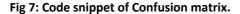### 5.3.1 Case Study 1 : Random Forest Classifier.

```
 6  # RandomForestClassifier
 7  model = RandomForestClassifier()
 8
 9  # Defining the hyperparameters and their possible values for grid search
10  param_grid = {
11      'n_estimators': [50, 100, 200],
12      'max_depth': [None, 10, 20, 30],
13      'min_samples_split': [2, 5, 10],
14      'min_samples_leaf': [1, 2, 4]
15  }
16
17  # Performing Grid Search
18  grid_search = GridSearchCV(model, param_grid, cv=5, scoring='accuracy')
19  grid_search.fit(X_train_combined, y_train_combined)
20
21  # Getting the best hyperparameters
22  best_params = grid_search.best_params_
23
24  # Training the model with the best hyperparameters
25  best_model = RandomForestClassifier(**best_params)
26  best_model.fit(X_train_combined, y_train_combined)
27
28  # Making predictions on the test set
29  y_pred = best_model.predict(X_test)
30
31  # Evaluating the model
32  accuracy = accuracy_score(y_test, y_pred)
33  print(f"Accuracy with the best hyperparameters: {accuracy:.4f}")
34
35  # Performing Randomized Search (alternative to Grid Search)
36  random_search = RandomizedSearchCV(model, param_distributions=param_grid, n_iter=10, cv=5,
37                                     scoring='accuracy', random_state=42)
38  random_search.fit(X_train_combined, y_train_combined)
39
40  # best hyperparameters from randomized search
41  best_params_random = random_search.best_params_
42
43  # Train the model with the best hyperparameters from randomized search
44  best_model_random = RandomForestClassifier(**best_params_random)
45  best_model_random.fit(X_train_combined, y_train_combined)
46
47  # Making predictions on the test set
48  y_pred_random = best_model_random.predict(X_test)
49
50  # Evaluating the model
51  accuracy_random = accuracy_score(y_test, y_pred_random)
52  report_random = classification_report(y_test, y_pred_random)
53
54
```

Accuracy with the best hyperparameters: 0.9849

**Fig 6: Code snippet of Random Forest Classifier.**

```
 1  #Confusion Matrix for RandomForestClassifier():
 2  confusion_mat = confusion_matrix(y_test, y_pred_random)
 3  print(f"Confusion Matrix for RandomForestClassifier():\n")
 4  # Visualize the confusion matrix as a heatmap
 5  plt.figure(figsize=(6, 4))
 6  sns.heatmap(confusion_mat, annot=True, fmt="d", cmap="Blues",
 7              xticklabels=['Class 1', 'Class 2', 'Class 3'],
 8              yticklabels=['Class 1', 'Class 2', 'Class 3'])
 9  plt.title(f'Confusion Matrix for RandomForestClassifier()')
10  plt.xlabel('Predicted')
11  plt.ylabel('Actual')
12  plt.show()
13
14  print('=' * 40)
```

**Fig 7: Code snippet of Confusion matrix.**

5

```
In [25]:   1 print(f"Best Hyperparameters for Decision Tree: {best_params}")

Best Hyperparameters for Decision Tree: {'max_depth': None, 'min_samples_leaf': 1, 'min_samples_split': 2, 'n_estimat
ors': 200}
```

**Fig 8: Code snippet of Best Hyperparameters.**

### 5.3.2   Case Study 2 : Decision Tree Classifier.

```
In [12]:    1 # Decision Tree
            2 dt_model = DecisionTreeClassifier()
            3
            4 # Defining hyperparameters for Decision Tree
            5 param_grid_dt = {
            6     'max_depth': [None, 5, 10, 15],
            7     'min_samples_split': [2, 5, 10],
            8     'min_samples_leaf': [1, 2, 4]
            9 }
           10
           11 # Performing Grid Search for Decision Tree
           12 grid_search_dt = GridSearchCV(dt_model, param_grid_dt, cv=5, scoring='accuracy')
           13 grid_search_dt.fit(X_train_combined, y_train_combined)
           14
           15 # Getting the best hyperparameters for Decision Tree
           16 best_params_dt = grid_search_dt.best_params_
           17
           18 # Training the Decision Tree model with the best hyperparameters
           19 best_dt_model = DecisionTreeClassifier(**best_params_dt)
           20 best_dt_model.fit(X_train_combined, y_train_combined)
           21
           22 # Making predictions on the test set using Decision Tree
           23 y_pred_dt = best_dt_model.predict(X_test)
           24
           25 # Evaluating the Decision Tree model
           26 accuracy_dt = accuracy_score(y_test, y_pred_dt)
           27 report_dt = classification_report(y_test, y_pred_dt)
           28
           29 #  Evaluating the model
           30 print(f"Accuracy for Decision Tree with hyperparameter tuning: {accuracy_dt:.4f}")
           31 print(f"Best Hyperparameters for Decision Tree: {best_params_dt}")
           32 print(f"Classification Report for Decision Tree with hyperparameter tuning:\n{report_dt}")
           33

Accuracy for Decision Tree with hyperparameter tuning: 0.9697
```

**Fig 9: Code snippet of Decision Tree Classifier.**

### 5.3.3   Case Study 3 : K-Nearest Neighbors

```
In [13]:  1   # K-Nearest Neighbors (KNN)
          2   knn_model = KNeighborsClassifier()
          3
          4   # Defining hyperparameters for KNN
          5   param_grid_knn = {
          6       'n_neighbors': [3, 5, 7],
          7       'weights': ['uniform', 'distance'],
          8       'p': [1, 2]
          9   }
         10
         11   # Performing Grid Search for KNN
         12   grid_search_knn = GridSearchCV(knn_model, param_grid_knn, cv=5, scoring='accuracy')
         13   grid_search_knn.fit(X_train_combined, y_train_combined)
         14
         15   # Get the best hyperparameters for KNN
         16   best_params_knn = grid_search_knn.best_params_
         17
         18   # Training the KNN model with the best hyperparameters
         19   best_knn_model = KNeighborsClassifier(**best_params_knn)
         20   best_knn_model.fit(X_train_combined, y_train_combined)
         21
         22   # Making predictions on the test set using KNN
         23   y_pred_knn = best_knn_model.predict(X_test)
         24
         25   # Evaluating the KNN model
         26   accuracy_knn = accuracy_score(y_test, y_pred_knn)
         27   report_knn = classification_report(y_test, y_pred_knn)
         28
         29   # Evaluating the model
         30   print(f"Accuracy for K-Nearest Neighbors with hyperparameter tuning: {accuracy_knn:.4f}")
         31   print(f"Best Hyperparameters for K-Nearest Neighbors: {best_params_knn}")
         32   print(f"Classification Report for K-Nearest Neighbors with hyperparameter tuning:\n{report_knn}")
         33

Accuracy for K-Nearest Neighbors with hyperparameter tuning: 0.9901
```

**Fig 10: Code snippet of K-Nearest Neighbors.**

### 5.3.4   Case Study 4 : Artificial Neural Network

```
In [14]:  1   # Artificial Neural Network (ANN)
          2   ann_model = MLPClassifier()
          3
          4   # Defining hyperparameters for ANN
          5   param_grid_ann = {
          6       'hidden_layer_sizes': [(50, 25), (100, 50), (150, 75)],
          7       'alpha': [0.0001, 0.001, 0.01]
          8   }
          9
         10   # Performing Grid Search for ANN
         11   grid_search_ann = GridSearchCV(ann_model, param_grid_ann, cv=5, scoring='accuracy')
         12   grid_search_ann.fit(X_train_combined, y_train_combined)
         13
         14   # Getting the best hyperparameters for ANN
         15   best_params_ann = grid_search_ann.best_params_
         16
         17   # Train the ANN model with the best hyperparameters
         18   best_ann_model = MLPClassifier(**best_params_ann)
         19   best_ann_model.fit(X_train_combined, y_train_combined)
         20
         21   # Making predictions on the test set using ANN
         22   y_pred_ann = best_ann_model.predict(X_test)
         23
         24   # Evaluating the ANN model
         25   accuracy_ann = accuracy_score(y_test, y_pred_ann)
         26   report_ann = classification_report(y_test, y_pred_ann)
         27
         28   # Evaluating the model
         29   print(f"Accuracy for Artificial Neural Network with hyperparameter tuning: {accuracy_ann:.4f}")
         30   print(f"Best Hyperparameters for Artificial Neural Network: {best_params_ann}")
         31   print(f"Classification Report for Artificial Neural Network with hyperparameter tuning:\n{report_ann}")

Accuracy for Artificial Neural Network with hyperparameter tuning: 0.8382
```

**Fig 11: Code snippet of ANN.**