

Predictive Analysis of Road Accidents in India: A Machine Learning Approach Configuration Manual

MSc Research Project Data Analytics

Harini Manjunatha Student ID: x22169288

School of Computing National College of Ireland

Supervisor: Shubham Subhnil

National College of Ireland Project Submission Sheet School of Computing



Student Name:	Harini Manjunatha
Student ID:	x22169288
Programme:	Data Analytics
Year:	2023
Module:	MSc Research Project
Supervisor:	Shubham Subhnil
Submission Due Date:	14/12/2023
Project Title:	Predictive Analysis of Road Accidents in India: A Machine
	Learning Approach
Word Count:	679
Page Count:	6

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:	Harini Manjunatha
Date:	31st January 2024

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):	

Predictive Analysis of Road Accidents in India: A Machine Learning Approach - Configuration Manual

Harini Manjunatha x22169288

1 Introduction

The configuration manual contains the detailed information on setting up system. The objective of this manual is to outline the process on conducting the study 'Predictive Analysis of Road Accidents in India: A Machine Learning Approach'. It specifies the necessary machine setup needed to create and execute the models. The instructions cover installing the required programmes and packages in addition to setting up the fundamental configuration that is necessary for the project to succeed.

2 System Specification

In this section, the hardware and software requirements used for the research is presented.

2.1 Hardware specification

Operating System	Mac OS
RAM	8 GB
System processor	Apple M2 (8-core)
Speed in Hz	3.49 G HZ
GPU	M2 10-Core GPU 13.5 W 1398 MHz

Figure 1: Hardware specification

2.2 Software specification

Programming Language	Python 3.9
Softwares	Anaconda and Jupyter notebook
Web browser	Safari

Figure 2: Software specification

3 Software Used

- Microsoft excel: Used for initial exploration of data.
- Jupyter Notebook ¹: Jupyter Notebook provides an interactive environment where we execute code in cells and visualize the results. It makes process easier to experiment with different machine learning models, algorithms, visualising the data and data preprocessing techniques.
- smartdraw ² : Online platform used for creating flowcharts.

4 Installation

The First step in the process involves launching the Anaconda application. The application provides an extensive set of software designed to perform diverse needs. Jupyter notebook is launched which offers an environment where we can execute the code in cells and visualise the results. Numerous exceptional Python libraries are available for comprehensive data analysis purposes. These libraries enhance the capability to explore, process, and derive insights from datasets efficiently. Figure 3 shows the interface of an Anaconda navigator.



Figure 3: Anaconda Interface

Figure 4 is the main page of the jupyter notebook. The new ipynb file is created to begin coding.

Figure 5 represents the code file of the project. The file is executed using the command restart and run all

¹Jupyter notebook: https://jupyter.org

²SmartDraw: https://cloud.smartdraw.com//?flags=129

💆 jupyter	Quit Logout
Files Running Clusters	
Select items to perform actions on them.	Upload New - 2
	Name 🕹 Last Modified File size
Desktop	2 hours ago
Documents	24 days ago
Downloads	11 minutes ago
Movies	a year ago
	19 days ago
D opt	a year ago
Pictures	9 months ago
	a year ago
VirtualBox VMs	a year ago
bankCusChurn_KNNPredictions_R.csv	7 months ago 4.69 kB
bankCusChurn_RFPredictions_R.csv	7 months ago 4.68 kB
🗆 🗋 d1.R	7 months ago 1.54 kB
□ 🗋 d2.R	7 months ago 1.56 kB

Figure 4: Jupyter Notebook

C Jupyter thesis Last Checkp	pint: 26/10/2023 (autosaved)	n Logout
File Edit View Insert C	ell Kernel Widgets Help	Trusted Python 3 (ipykernel) O
	► Ru Interrupt I, I + Restart 0, 0	
	Restart & Clear Output	
MSc Res	Reconnect Shutdo	notebook
Data Ana	Change kernel	
Predictiv Approact	Analysis of Road Accidents in	India: A Machine Learning
Harini Ma	injunatha - x22169288	
Required	Imports	
In [1]: import num import pan	y as np las as pd	
#Dataset in import csv	aport	
#Model bui from sklea from sklea from sklea from sklea from sklea from sklea from sklea	d m.sodel_selection import train_test m.ensemble import RandomForestRepres DecisionTreeRepressor m.edito DecisionTreeRepressor m.edito: import mean_auseidute_error m.edit_selection import foriSearch m.odel_selection import foriSearch m.edel_selection import foriSearch	split sor rr_z.core rmean_squared_error V
<pre># Analysis import sea import mat *matplotlii</pre>	oorn as sns olotlib.pyplot as plt o inline	

Figure 5: Code File on jupyter notebook

5 Project Development

5.1 Importing necessary library

Figure 6 shows the list of all necessary library imports used in the project. The numpy ³ and pandas ⁴ libraries enable efficient manipulation and analysis of data through functionalities like arrays, matrices, and data structures like DataFrames. Scikit-learn modules ⁵, regression models like RandomForestRegressor, DecisionTreeRegressor, LinearRegression and GridSearchCV are imported to construct, evaluate, and tune machine learning models.

³Numpy: https://numpy.org

⁴Pandas: https://pandas.pydata.org

⁵Scikit-learn: https://scikit-learn.org/stable/

	Required Imports
[1]:	import numpy as np import pandas as pd
	#Dataset import import csv
	<pre>#Model build from sklearn.model_selection import train_test_split from sklearn.ensemble import RandomForestRegressor from sklearn.tree import DecisionTreeRegressor from sklearn.neighbors import KNeighborsRegressor from sklearn.metrics import mean_squared_error, r2_score from sklearn.metrics import mean_absolute_error, mean_squared_error from sklearn.linear_model import GridSearchCV from sklearn.model_selection import GridSearchCV</pre>
	# Analysis import seaborn as sns import matplotlib.pyplot as plt *matplotlib inline import plotly.graph_objs as go import plotly.express as px sns.set() import plotly.express as px from tabulate import tabulate

Figure 6: Library Imports]

5.2 Import dataset

Ir

Figure 7 is the code snippet used for importing the dataset.



Figure 7: Dataset Import

5.3 Model Building

The dataset is spli into two sets. One for training and another for testing. Data from the year 2001 to 2020 is used for training whereas data of 2021 is used for testing. The code snippet which splits the data is shown in the Figure 8



Figure 8: Dataset Splitting

The code snippet in Figure 9 identify and eliminate outliers from the dataset using the z-score method.

Then code snippet in the Figure 10 implements a Random Forest regression model using scikit-learn's RandomForestRegressor ⁶. It initiates a grid search for the best hyperparameters using GridSearchCV ⁷ to optimize the model's performance. The code eval-

⁶Scikit-learn: https://scikit-learn.org/stable/

⁷Grid search: https://scikit-learn.org/stable/modules/generated/sklearn.model_selection.GridSearchCV.html

n []:	<pre># Outlier detection numeric_columns = df.select_dtypes(include=[np.number]).columns.tolist() numeric_columns.remove('Year')</pre>
	# Outlier detection using z-score method for each numeric column outlier_threshold = 3 outlier_indices = []
	<pre>for col in numeric_columns: mean = df[col].mean() std_dev = df[col].std() outliers = df[(df[col] = mean).abs() > outlier_threshold * std_dev] outliers_indices.extend(outliers.index.tolist())</pre>
	<pre>unique_outlier_indices = list(set(outlier_indices))</pre>
	df_outliers = df.drop(unique_outlier_indices)
	<pre>print("Indices of Detected Outliers:", unique_outlier_indices)</pre>

I

Figure 9: Outlier detection

uates the model's performance metrics using Mean Absolute Error (MAE), R-squared, Root Mean Squared Error (RMSE) for each target variable in test set.



Figure 10: Random Forest Regressor Model

Figure 11 visualizes the distribution of differences between the actual and predicted values from the regression model. Normal distribution plot provides an immediate understanding of the spread and density of prediction errors. The inclusion of vertical lines denoting the mean difference and standard deviations aids in interpreting the central tendency and variability of these prediction errors, offering valuable insights into the model's performance and the nature of discrepancies between actual and predicted values.

The Predicted output from the best performing model is visualised to understand the data pattern. The code snippet of the visualisation is given in Figure 12

```
# Calculate the differences between actual and predicted values
differences = test_labels.values - predictions
# Flatten the differences array for plotting
differences_flat = differences.flatten()
# Plot the normal distribution
plt.figure(figsize=(12, 6))
sns.histplot(differences_flat, kde=True, color='blue', stat='density')
plt.title('Difference between Actual and Predicted Values')
plt.xlabel('Difference')
plt.vlabel('Density')
# Add a vertical line at mean and standard deviations
mean_diff \= np.std(differences_flat)
std_dev_diff = np.std(differences_flat)
plt.axvline(mean_diff, color='red', linestyle='dashed', linewidth=1)
plt.axvline(mean_diff - std_dev_diff, color='green', linestyle='dashed', linewidth=1)
plt.axvline(mean_diff - std_dev_diff, color='green', linestyle='dashed', linewidth=1)
plt.axvline(mean_diff - std_dev_diff, color='green', linestyle='dashed', linewidth=1)
plt.legend(['Mean', 'Mean + 1 Std Dev', 'Mean - 1 Std Dev'])
plt.show()
```





Figure 12: Predicted result Visualisation