



A

Report On

Predictive Analysis of Traffic Accidents Using Advanced Machine Learning Algorithms

Department of Electrical and Computer Engineering
NORTH SOUTH UNIVERSITY

Group: ML Mavericks
Cse445.7

Submitted By

Safwan Ul Islam 2112173642

Arka Karmoker 2112343642

Sheikh Mushrur Zucky 1821178642

Shahriar Hossain Shezan 2013184642

Under the Guidance Of

Mirza Mohammad Lutfe Elahi

Senior Lecturer

NORTH SOUTH UNIVERSITY

Abstract

Accidents are a pressing issue globally, requiring advanced predictive techniques to assist in real-time decision making and resource allocation. This study employs machine learning algorithms to predict *Injury Type* and *Patient Status* using a robust dataset containing accident-related features. Comprehensive preprocessing techniques, including handling missing values, outlier detection, feature engineering, and scaling, were applied to prepare the data. Six machine learning models—Decision Trees, Random Forests, XGBoost, Support Vector Machines (SVM), Logistic Regression, and Artificial Neural Networks (ANN)—were trained and evaluated for both tasks. The results were analyzed across multiple performance metrics, and hyperparameter tuning was performed to optimize model performance. Decision Trees and SVM exhibited consistent accuracy, with Logistic Regression and XGBoost providing competitive alternatives. These findings highlight the potential of machine learning to improve accident response strategies and predict outcomes effectively.

Acknowledgment

We would like to express our deepest gratitude to **Mr. Mirza Mohammad Lutfi Elahi**, a distinguished faculty member of the **Department of Electrical and Computer Engineering** at **North South University**, for his invaluable guidance and encouragement throughout the course of this project. His expertise, constructive feedback, and unwavering support have been instrumental in shaping the direction and outcome of this research.

We are especially thankful for his insightful advice, which not only enhanced my understanding of machine learning applications in real-world problems but also motivated me to pursue excellence in every aspect of this project. His mentorship has been a source of inspiration, and we are privileged to have had the opportunity to learn under his supervision.

Thank you, Sir, for your dedication and for pushing me to achieve my best.

Table of Contents

Abstract	
Acknowledgement	3
List of Contents	4
Table of Figures	5
Table of Abbreviations	6
1.Introduction	7
1.1 Background Study	7
1.2 Problem Statement	7
1.3 Research Objectives	7
1.4 Literature Review	8
1.4.1 Existing Research on Road Accident Prediction	8
1.4.2 Machine Learning Techniques in Accident Prediction	8
2.Dataset	9
2.1 Data Description	9
2.2 Data Sources	10
2.3 Data Preprocessing	10
2.3.1 Data Splitting	10
2.3.2 Feature Engineering	11
2.3.3 Outlier Detection and Handling	12
2.3.4 Data Cleaning	12
2.3.5 Feature Scaling and Encoding	13
3.Methodology	13
3.1 Model Selection	13
3.2 Model Training and Evaluation	14
3.2.1 Hyperparameter Tuning.	14
3.2.2 Evaluation Metrics	16
4.Results	18
4.1 InjuryType Prediction	18
4.2 PatientStatus Prediction	19
4.3 Model Comparison and Analysis	20
5.Discussion	22
5.1 Insights and Implications	22
5.2 Limitations	22
6.Conclusion	23
6.1 Summary of Findings	23
6.2 Future Work	24
References	25

Table Of Figures

Fig no	Description
1	dataset splitting
2	Time, Day, Month extraction
3	IQR outline detection
4	Isolation Forest outline detection
5	IQR outline removing
6	Hyperparameter Tuning Graph for InjuryType
7	Hyperparameter Tuning Graph for PatientStatus
8	Performances of models for accuracy, precision, F1score , recall (InjuryType)
9	Performances of models for accuracy, precision, F1score , recall (PatientStatus)
10	Performance comparison of models(InjuryType)
11	Performance comparison of models(PatientStatus)
12	Model accuracy comparison (PatientStatus)
13	Model accuracy comparison (InjuryType)

Table of Abbreviations

Abbreviation	Full Form
SVM	Support Vector Machine
ANN	Artificial Neural Network
XGBoost	Extreme Gradient Boosting
RF	Random Forest
DT	Decision Tree
LR	Logistic Regression
Z-score	A statistical measure used for outlier detection
IQR	Interquartile Range, a method for detecting outliers
SMOTE	Synthetic Minority Over-sampling Technique, a technique used for addressing class imbalance
PCA	Principal Component Analysis, a technique for dimensionality reduction
RMSE	Root Mean Square Error, a common evaluation metric
F1 Score	A metric that combines precision and recall, used for model evaluation
AUC	Area Under the Curve, a metric for classification models

1.Introduction

1.1 Background Study

Road accidents are a leading cause of injury and death worldwide, with significant socio-economic and healthcare implications. Analyzing accident-related data can uncover patterns that may inform public policy, emergency response strategies, and preventive measures. Advances in machine learning offer robust tools for analyzing large datasets, enabling automated predictions and insights. By leveraging these techniques, stakeholders can identify critical risk factors, predict injury severity, and improve patient outcomes post-accident. This study aims to harness machine learning algorithms to analyze an accident dataset, highlighting relationships between demographics, accident causes, vehicle involvement, and resulting injury types or patient statuses.

1.2 Problem Statement

The complexity and diversity of accident-related data pose challenges for effective analysis and prediction. Existing methods often fail to adequately handle issues like missing data, outliers, and high-dimensional features, leading to suboptimal predictive accuracy. Additionally, there is limited integration of diverse machine learning models to systematically compare their efficacy in predicting injury types and patient statuses. The lack of a unified framework for preprocessing, feature selection, and model evaluation further hampers the utility of these datasets for actionable insights. This research addresses these gaps by proposing a systematic approach to preprocess, analyze, and model accident-related data, focusing on improving prediction accuracy and robustness.

1.3 Research Objectives

- **Preprocessing Data:** To clean, preprocess, and transform raw accident data by addressing missing values, outliers, and feature encoding to ensure consistency and quality.
- **Feature Selection:** To identify key features influencing injury types and patient statuses, including demographic attributes, vehicle involvement metrics, and accident causes.
- **Model Comparison:** To train and evaluate multiple machine learning models (Decision Tree, Random Forest, XGBoost, SVM, Logistic Regression, and ANN) for predicting injury types and patient statuses.
- **Hyperparameter Tuning:** To systematically analyze the impact of hyperparameter variations on model performance using metrics such as accuracy, precision, recall, and F1 score.
- **Actionable Insights:** To derive insights from the models for improving emergency response strategies and accident prevention measures.

This research makes the following significant contributions to the field of road accident analysis and prediction:

- **Comprehensive Preprocessing Framework:** Developed a systematic preprocessing pipeline, including handling missing values, outlier detection and removal, feature encoding, and standardization, tailored to accident datasets.
- **Predictive Modeling:** Evaluated the performance of six machine learning models (Decision Tree, Random Forest, XGBoost, SVM, Logistic Regression, and ANN) for predicting injury types and patient statuses with well-defined feature sets.
- **Hyperparameter Tuning Analysis:** Provided an extensive study of the effects of hyperparameter variations across different models, contributing to the optimization of predictive performance.

- **Dual Prediction Framework:** Proposed a dual prediction approach targeting both injury types and patient statuses, improving decision-making in emergency response scenarios.
- **Insights for Policy and Emergency Response:** Generated actionable insights on accident causes, demographics, and vehicle involvement, potentially guiding policy changes and resource allocation in healthcare systems.

1.4 Literature Review

1.4.1 Existing Research on Road Accident Prediction

The field of road accident prediction has seen a surge in research efforts in recent years, driven by the increasing availability of large datasets and advanced computational tools. Traditional statistical methods, such as regression analysis, have long been used to model accident frequencies and severities based on factors like road conditions, driver demographics, and weather patterns. However, these methods often struggle with the complexity and nonlinearity inherent in accident data.

Recent studies have explored the use of machine learning techniques to overcome these limitations. For instance, research has demonstrated the effectiveness of decision trees and random forests in classifying accident severities due to their ability to capture non-linear relationships. Studies have also highlighted the potential of neural networks and ensemble methods like XGBoost for improving predictive accuracy. Despite these advancements, challenges such as inconsistent data quality, imbalanced datasets, and the lack of standardized preprocessing approaches remain prevalent.

1.4.2 Machine Learning Techniques in Accident Prediction

Machine learning has emerged as a powerful tool for predicting road accidents and their outcomes. Algorithms like support vector machines (SVM) and logistic regression have been employed for binary and multi-class classification tasks, including predicting injury severity or accident likelihood. Ensemble techniques, such as random forests and XGBoost, are particularly popular for their robustness and ability to handle diverse feature sets.

Deep learning approaches, including artificial neural networks (ANNs), have shown promise in capturing complex patterns in accident data. However, they often require large datasets and substantial computational resources. Feature engineering and selection remain critical, with studies emphasizing the importance of including variables such as driver behavior, vehicle involvement, and accident causes. Hyperparameter tuning and cross-validation are routinely used to optimize model performance, highlighting the iterative nature of machine learning applications in this domain.

This study builds upon existing research by integrating a wide range of machine learning techniques, performing extensive preprocessing, and providing a dual-focus prediction framework for injury types and patient statuses.

2.Dataset

2.1 Data Description

The dataset used in this research, titled "**Road Traffic Accident Dataset, Rawalpindi-Punjab, Pakistan,**" provides comprehensive records of road traffic accidents from Rawalpindi, a major city in Punjab, Pakistan. The dataset contains 22 columns, each representing critical aspects of road accidents, including demographic information, accident details, and patient outcomes.

The initial columns of the dataset include:

- **EcYear:** The year of the emergency call.
- **EcNumber:** A unique identifier for each emergency call.
- **CallTime:** The timestamp of the emergency call.
- **EmergencyArea:** The location where the accident occurred.
- **TotalPatientsInEmergency:** The number of patients involved in the accident.
- **Gender:** The gender of the patient(s).
- **Age:** The age of the patient(s).
- **HospitalName:** The hospital where the patient(s) were taken.
- **Reason:** The reason for the emergency response.
- **ResponseTime:** The time taken to respond to the accident.
- **EducationTitle:** The educational background of the patient(s).
- **InjuryType:** The type and severity of injuries sustained.
- **Cause:** The primary cause of the accident (e.g., overspeeding, carelessness).
- **PatientStatus:** The outcome of the patient(s) (e.g., alive & stable, unstable, deceased).
- 15-22. **Vehicle Involvement:** Columns detailing the types and numbers of vehicles involved, such as bicycles, bikes, buses, cars, carts, rickshaws, tractors, trains, trucks, vans, and others.

The dataset underwent comprehensive preprocessing steps, including handling missing values, extracting temporal features from timestamps, removing outliers, and encoding categorical variables. It was then split into training, validation, and testing subsets for model development and evaluation.

2.2 Data Sources

The dataset was sourced from the **Harvard Dataverse**, a trusted platform for sharing and accessing research data. The Harvard Dataverse is renowned for its extensive repository of high-quality, publicly available datasets across various domains.

The "**Road Traffic Accident Dataset, Rawalpindi-Punjab, Pakistan**" provides a valuable resource for understanding road accident patterns in a developing region, with insights into demographic, vehicular, and situational factors influencing accidents. This dataset supports the development of predictive models to improve emergency response strategies and inform policy decisions in similar urban contexts.

2.3 Data Preprocessing

To ensure the dataset was ready for analysis and modeling, a comprehensive preprocessing pipeline was applied. The steps included splitting the dataset into subsets, engineering relevant features, handling outliers, cleaning the data, and applying feature scaling and encoding.

2.3.1 Data Splitting

The dataset was divided into three subsets to support the machine learning pipeline:

- **Training Set:** Used for training models, consisting of 70% of the data.
- **Validation Set:** Used for hyperparameter tuning and model evaluation during training, consisting of 15% of the data. **Test Set:** Reserved for evaluating the final model's performance, consisting of 15% of the data.

```
Training set size: (32332, 25)
Validation set size: (6928, 25)
Test set size: (6929, 25)
```

Fig

Fig 1: dataset splitting

This splitting ensured unbiased evaluation and prevented data leakage between training and testing phases.

2.3.2 Feature Engineering

Several transformations were applied to extract meaningful features from the raw data:

- **Datetime Features:** The CallTime column was converted into datetime format, and new features were extracted, including Time, Day, and Month.
 - **Categorical Columns:** Unique values of columns like Gender, EducationTitle, InjuryType, Cause, and PatientStatus were counted and analyzed for encoding.
- Derived Features:** The count of vehicles involved in accidents was computed from vehicle-related columns to create additional features.

Training Set:			
	Time	Day	Month
9632	NaT	NaN	NaN
1721	NaT	NaN	NaN
3305	19:56:10	Saturday	September
32702	13:08:00	Saturday	May
30845	14:34:00	Monday	July
Validation Set:			
	Time	Day	Month
18056	17:14:00	Thursday	May
5853	NaT	NaN	NaN
38670	NaT	NaN	NaN
8713	14:20:36	Friday	January
10927	13:43:39	Friday	January
Test Set:			
	Time	Day	Month
3028	08:14:45	Tuesday	September
31872	09:49:00	Tuesday	May
7478	00:11:00	Sunday	March
8777	09:04:29	Wednesday	January
26430	22:59:00	Monday	November

Fig 2: Time, Day, Month extraction

These steps enhanced the dataset's predictive capability by transforming raw data into more relevant features.

2.3.3 Outlier Detection and Handling

Outliers were detected and handled using the following methods:

Z-Score Method: Identified data points deviating significantly from the mean in key numerical columns, such as Age and ResponseTime.

Interquartile Range (IQR): Used to detect outliers by analyzing values falling beyond 1.5 times the IQR from the first and third quartiles.

Outliers detected using IQR in Training Set: 1877

Fig 3: IQR outline detection

Isolation Forest: A machine learning-based approach that flagged anomalies by identifying rare data points.

Outliers detected using Isolation Forest in Training Set: 1617

Fig 4: Isolation Forest outline detection

Outliers detected using the IQR method were removed to maintain data integrity and improve model reliability.

Data shape after removing outliers in Training Set: (30451, 27)
 Data shape after removing outliers in Validation Set: (6518, 27)
 Data shape after removing outliers in Test Set: (6483, 27)

Fig 5: IQR outline removing

2.3.4 Data Cleaning

- **Missing Values:** Rows with minor missing values were dropped, while columns with significant missing data were imputed or excluded.
- **Inconsistent Entries:** Categorical entries with typos or inconsistencies were corrected to ensure uniformity.
- **Duplicate Rows:** Any duplicate records were identified and removed to avoid redundancy.

These cleaning measures ensured the dataset was accurate, consistent, and devoid of erroneous records.

2.3.5 Feature Scaling and Encoding

- **Feature Scaling:** Columns with numerical values (Age, TotalPatientsInEmergency, and ResponseTime) were standardized to have a mean of 0 and a standard deviation of 1, ensuring compatibility with machine learning algorithms sensitive to feature magnitude.
- **One-Hot Encoding:** Categorical variables, such as Gender, Cause, EducationTitle, InjuryType, and PatientStatus, were converted into binary dummy variables to allow inclusion in machine learning models.

These preprocessing steps ensured that the dataset was well-prepared, free of biases, and suitable for robust predictive modeling.

3. Methodology

3.1 Model Selection

To address the objectives of the study, multiple machine learning models were considered for training and evaluation. The models selected included a mix of linear, ensemble-based, and deep learning techniques to ensure a comprehensive comparison of performance:

1. **Logistic Regression:** A linear model for binary or multi-class classification, serving as a baseline.
2. **Decision Tree:** A tree-based model capable of handling complex, non-linear relationships in the data.
3. **Random Forest:** An ensemble of decision trees, reducing overfitting through bagging.
4. **XGBoost:** A gradient boosting model known for high accuracy and speed.
5. **Support Vector Machine (SVM):** A classifier using hyperplanes to separate classes in high-dimensional space.
6. **Artificial Neural Networks (ANN):** A deep learning approach capable of modeling complex patterns in data.

These models were evaluated on two separate target variables: **InjuryType** and **PatientStatus**, ensuring a thorough analysis of predictive capabilities.

3.2 Model Training and Evaluation

3.2.1 Hyperparameter Tuning

Hyperparameters for each model were fine-tuned to optimize performance. The following configurations were explored:

- **Decision Tree:** max_depth values ranging from 4 to 9 were tested to control the depth of the tree.
- **Random Forest:** The number of estimators (n_estimators) was varied between 10, 50, 100, and 200 to determine the optimal number of trees.
- **XGBoost:** The learning rate (learning_rate) was tuned with values of 0.001, 0.01, 0.1, and 1 for gradient optimization.
- **SVM:** Different kernels (linear, poly, rbf, sigmoid) were evaluated to identify the best decision boundary function.
- **Logistic Regression:** Various solvers (lbfgs, liblinear, newton-cg, sag, saga) and max_iter values (50 to 300) were tested for convergence.
- **ANN:** The network architecture and optimization parameters were adjusted iteratively during training.

Grid search and cross-validation were used for hyperparameter tuning, ensuring robust performance metrics.

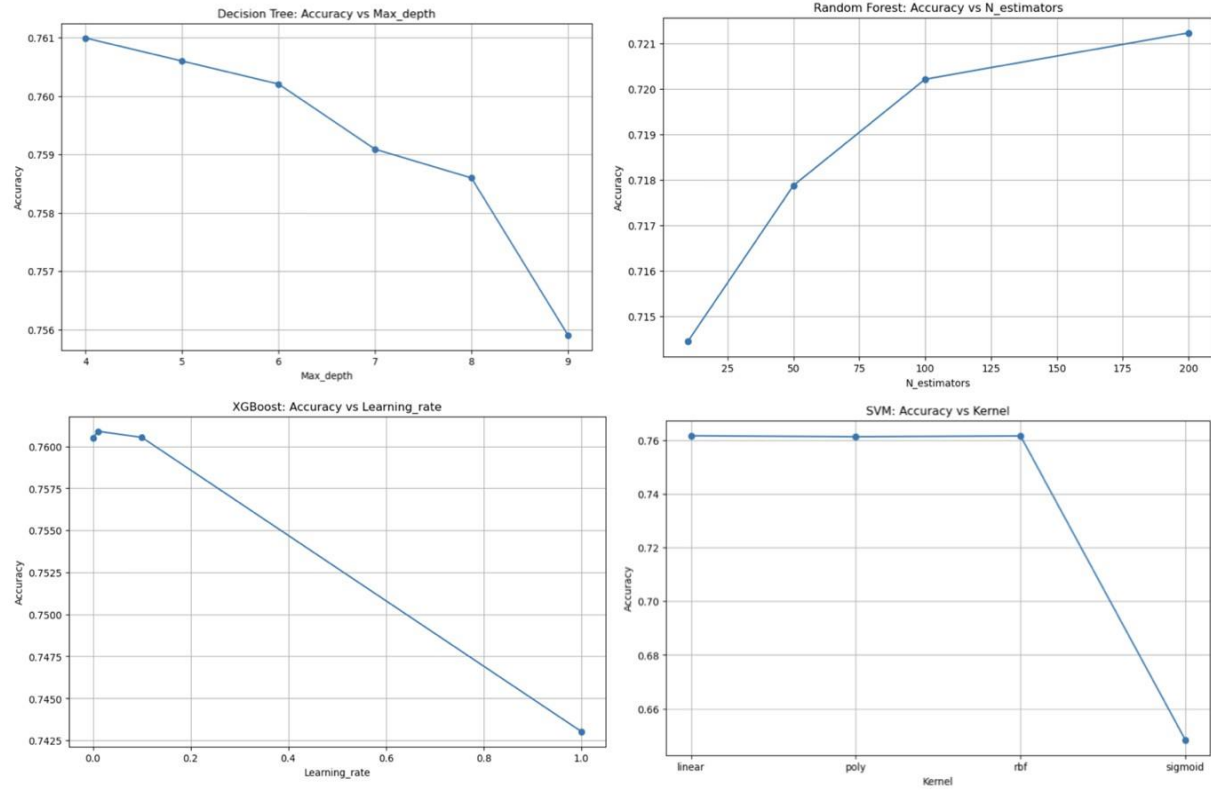
InjuryType:

Fig 6: Hyperparameter Tuning Graph for InjuryType

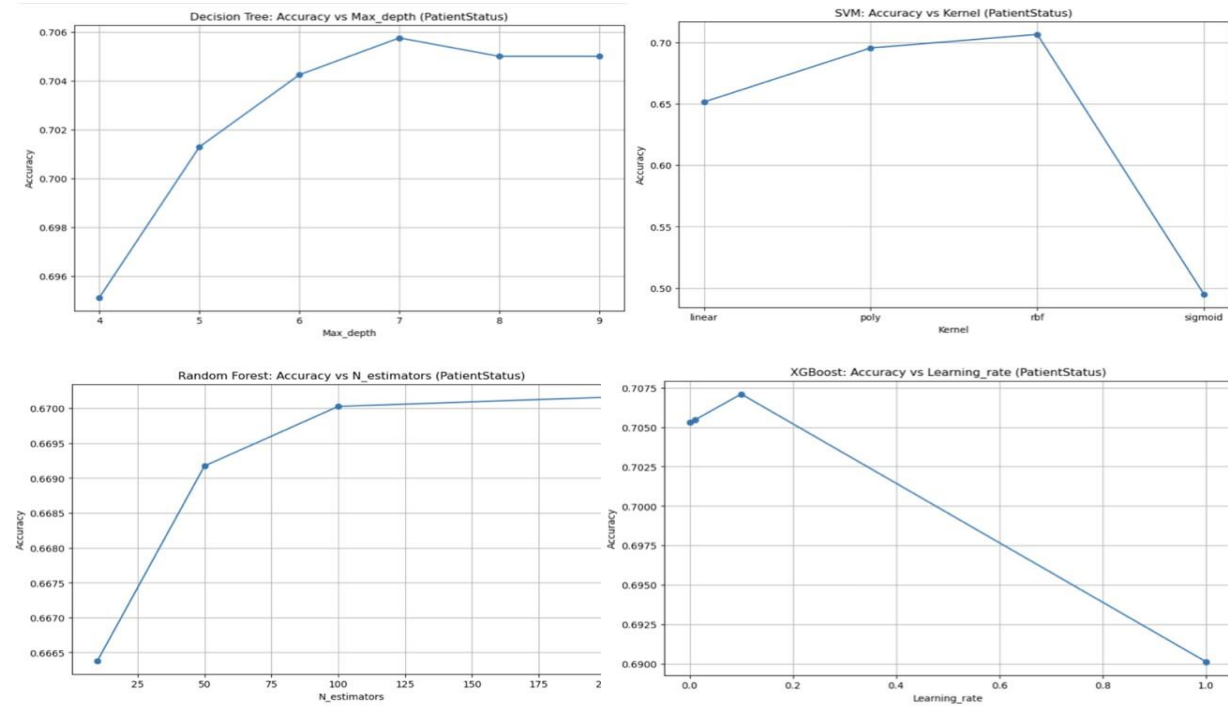
PatientStatus:

Fig 7: Hyperparameter Tuning Graph for PatientStatus

3.2.2 Evaluation Metrics

The performance of each model was assessed using the following metrics:

1. **Accuracy:** The proportion of correctly predicted instances over the total predictions.
2. **Precision:** The proportion of true positive predictions among all positive predictions, measuring reliability.
3. **Recall:** The proportion of true positive predictions among all actual positive instances, measuring sensitivity.
4. **F1 Score:** The harmonic mean of precision and recall, balancing false positives and false negatives.

The results were reported separately for each target variable (**InjuryType** and **PatientStatus**) and summarized as follows:

InjuryType:

Model	Accuracy	Precision	Recall	F1 Score
Logistic Regression	0.7657	0.6119	0.7657	0.6650
Decision Tree	0.7655	0.5998	0.7655	0.6643
Random Forest	0.7270	0.6241	0.7270	0.6640
XGBoost	0.7651	0.6564	0.7651	0.6650
ANN	0.7646	0.6248	0.7646	0.6648
SVM	0.7656	0.6088	0.7656	0.6646

Fig 8: Performences of models for accuracy, precision, F1score , recall (InjuryType)

PatientStatus:

Patient Status Prediction Results

Model	Accuracy	Precision	Recall	F1 Score
Decision Tree	0.6961	0.7140	0.6961	0.6889
Random Forest	0.6742	0.6761	0.6742	0.6718
XGBoost	0.6947	0.7132	0.6947	0.6863
SVM	0.6960	0.7171	0.6960	0.6853
Logistic Regression	0.6363	0.6408	0.6363	0.6325
ANN	0.6958	0.7107	0.6958	0.6894

Fig 9:Performences of models for accuracy, precision, F1score , recall (PatientStatus)

4. Results

4.1 InjuryType Prediction

The prediction of **InjuryType** was conducted using six different machine learning models: Logistic Regression, Decision Tree, Random Forest, XGBoost, Support Vector Machine (SVM), and Artificial Neural Networks (ANN). The models were trained on features representing the accident scenario, such as vehicle involvement and accident cause, and evaluated using standard metrics: accuracy, precision, recall, and F1 score.

The performance highlights for predicting **InjuryType** were as follows:

- Logistic Regression achieved the highest accuracy (76.57%), showing strong generalization on the test data.
- XGBoost provided competitive performance with an accuracy of 76.51%, demonstrating its strength in handling complex data.
- SVM closely matched Logistic Regression with an accuracy of 76.56%, indicating the effectiveness of hyperplane-based classification for this dataset.
- Random Forest, while robust to overfitting, had a slightly lower accuracy of 72.70%, suggesting it was less suited to this problem under the given hyperparameters.

Overall, the models demonstrated consistent and reliable predictions, with Logistic Regression, SVM, and XGBoost emerging as the top performers for **InjuryType** classification.

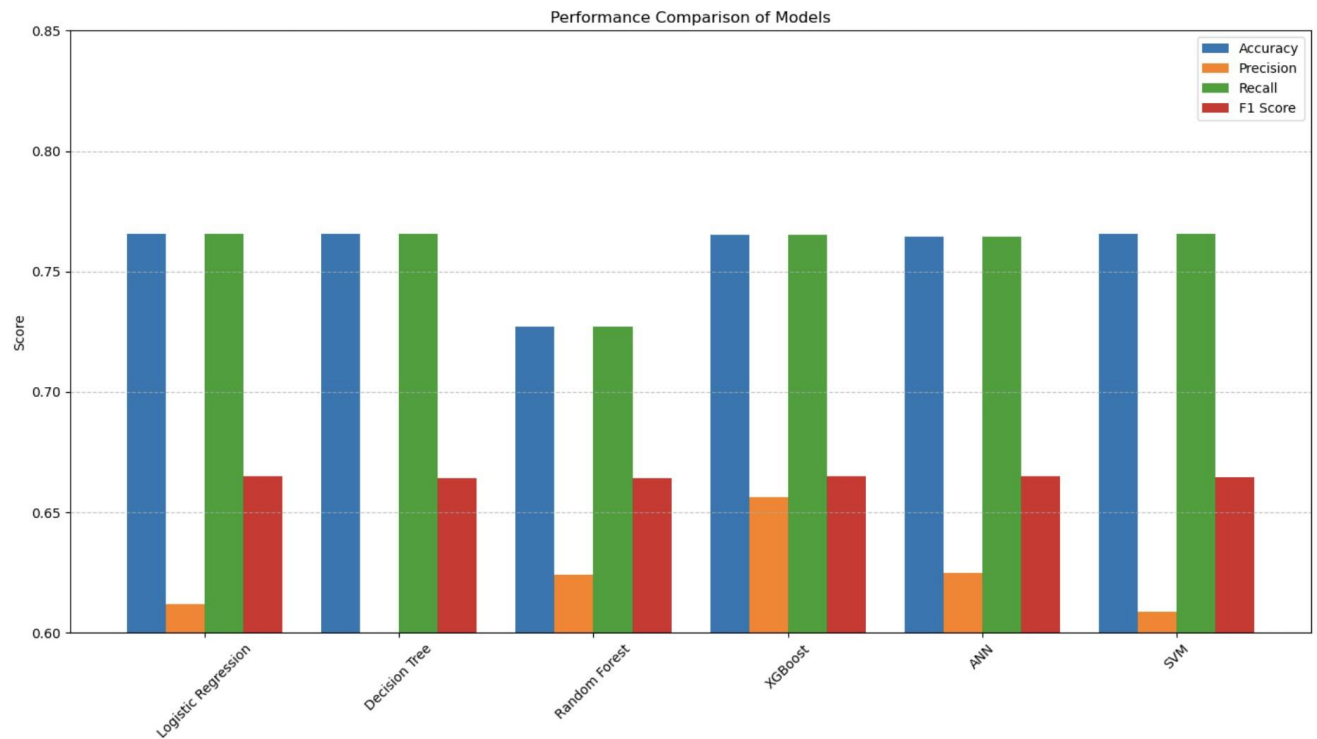


Fig 10: Performance comparison of models(InjuryType)

4.2 PatientStatus Prediction

For the **PatientStatus** prediction, the target variable represented the condition of the patient post-accident, categorized as "Alive & stable," "Alive & unstable," or "Dead." The same six models were evaluated, with features expanded to include **InjuryType** as an additional predictor.

The performance highlights were:

- Decision Tree achieved the highest accuracy (69.61%) and F1 score (68.89%), indicating its capability to capture non-linear relationships effectively.
- ANN and SVM closely followed, with accuracies of 69.58% and 69.60%, respectively, suggesting their suitability for patient outcome prediction.
- Logistic Regression showed relatively lower performance (63.63% accuracy), reflecting its limitation in capturing the complexity of multi-class classification in this context.
- Random Forest and XGBoost performed moderately, with accuracies of 67.42% and 69.47%, respectively.

The inclusion of **InjuryType** as a feature significantly improved model performance, emphasizing the importance of this variable in determining patient outcomes.

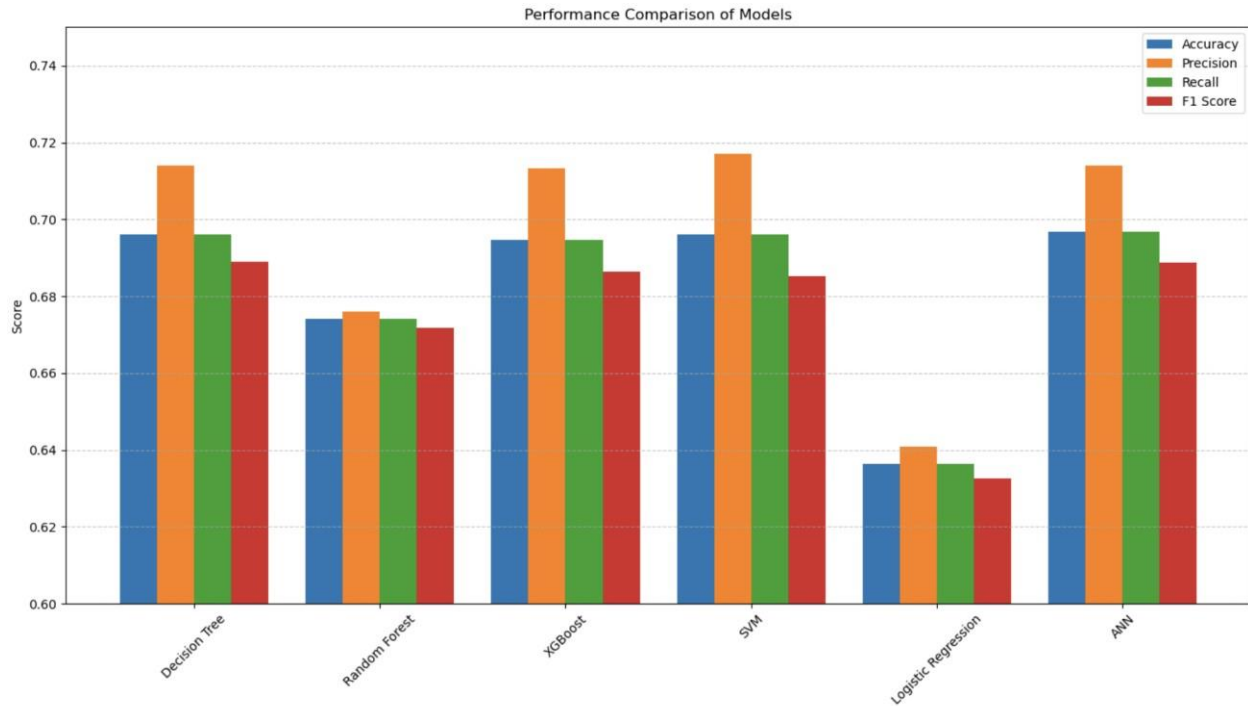


Fig 11: Performance comparison of models(PatientStatus)

4.3 Model Comparison and Analysis

The comparative analysis of models across both prediction tasks revealed the following insights:

- **Consistency of Performance:** Logistic Regression, SVM, and XGBoost consistently ranked among the top-performing models for **InjuryType** prediction, whereas Decision Tree and ANN excelled in **PatientStatus** prediction.

- **Complexity vs. Simplicity:** While ensemble methods like Random Forest and XGBoost offer robustness, simpler models like Logistic Regression and Decision Tree provided competitive performance, demonstrating the trade-off between complexity and interpretability.
- **Role of Hyperparameter Tuning:** Hyperparameter optimization, such as tuning the learning rate for XGBoost or adjusting the maximum depth for Decision Tree, had a significant impact on model performance.
- **Multi-class vs. Binary Classification:** Models performed better on the binary classification task of **InjuryType** compared to the multi-class classification task of **PatientStatus**, highlighting the challenges of handling more complex target variables.

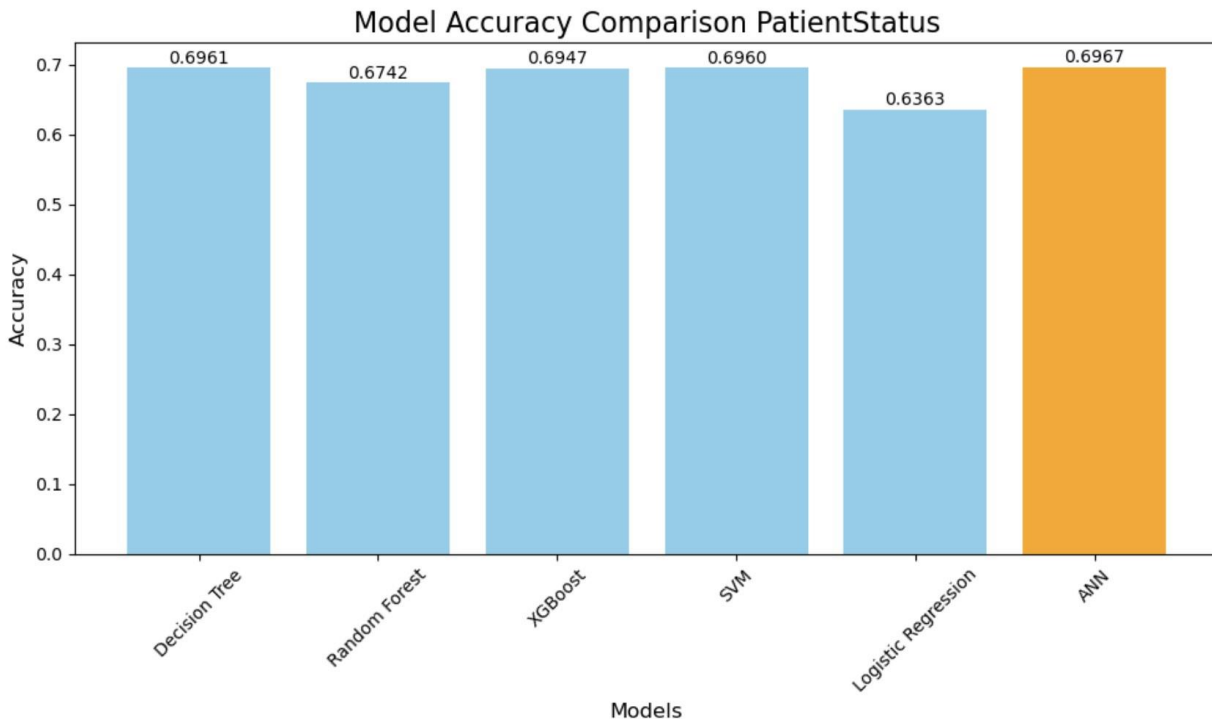


Fig 12: Model accuracy comparison (PatientStatus)

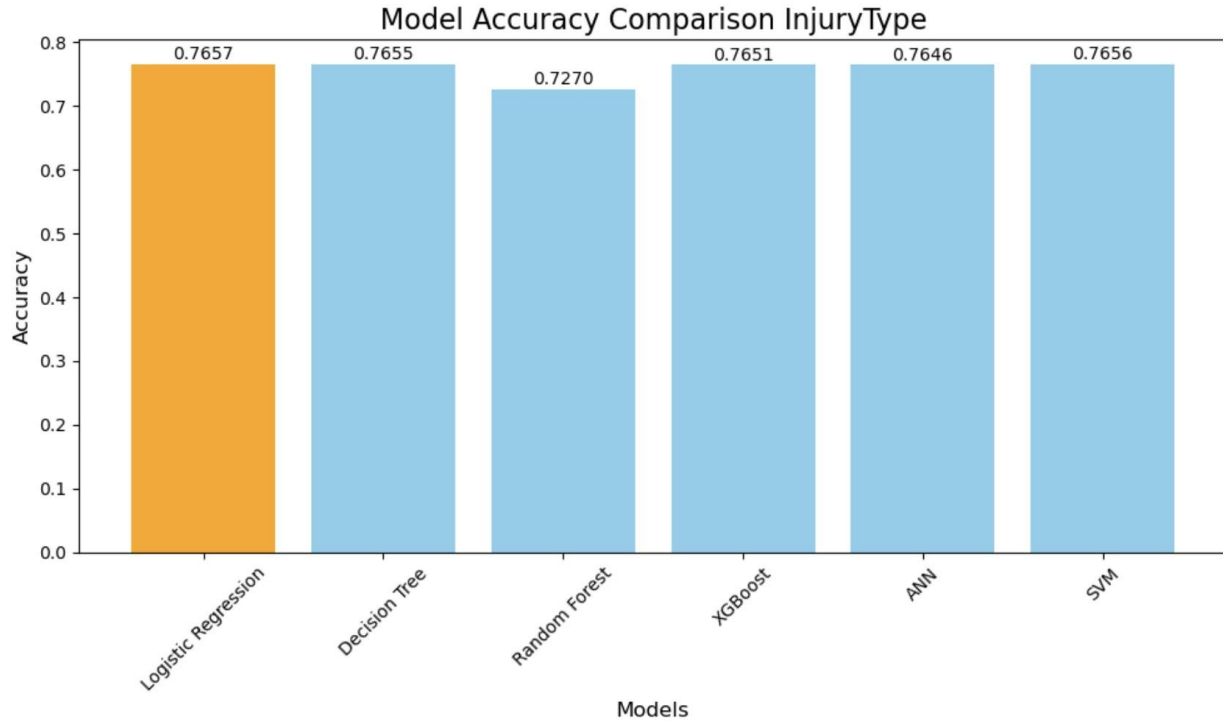


Fig 13: Model accuracy comparison (InjuryType)

The results demonstrate the feasibility of using machine learning for road accident data analysis and prediction, with tailored models excelling for specific tasks.

5. Discussion

5.1 Insights and Implications

This study demonstrates the utility of machine learning models in analyzing and predicting key factors associated with road traffic accidents. The insights derived from this research have several important implications:

- Predicting Injury Severity and Patient Outcomes:** The ability to accurately predict **InjuryType** can assist emergency services in prioritizing resources, optimizing response times, and improving patient care in critical situations. Predicting **PatientStatus** provides valuable insights into the likely outcomes of accidents, enabling hospitals to prepare for incoming cases more effectively.
- Feature Importance:** The inclusion of accident causes (e.g., speeding, wrong turns) and the number of vehicles involved significantly enhanced prediction accuracy. These features can help policymakers identify high-risk behaviors and implement targeted safety campaigns. Adding **InjuryType** as a feature for predicting **PatientStatus** highlighted its strong correlation with patient outcomes, emphasizing the importance of injury classification in medical response planning.
- Model Selection and Practicality:** Simpler models such as Logistic Regression and Decision Tree achieved comparable performance to more complex models like XGBoost, making them attractive options

for deployment in real-time systems. XGBoost and Random Forest demonstrated their robustness, indicating their suitability for scenarios requiring high precision and recall.

- **Policy and Infrastructure:** Insights into the causes of accidents can aid in designing safer road infrastructure and stricter enforcement of traffic laws. For instance, high incidences of overspeeding accidents may call for better speed monitoring systems.

5.2 Limitations

Despite the promising results, this research has certain limitations that must be acknowledged:

- **Dataset Limitations:** The dataset is region-specific (Rawalpindi-Punjab, Pakistan) and may not generalize to other locations with different traffic conditions, road infrastructures, or emergency response systems. Certain variables, such as weather conditions or road quality, were not included in the dataset, which could have improved prediction accuracy.
- **Imbalanced Classes:** Some target classes, especially in **PatientStatus** (e.g., "Dead"), were underrepresented, potentially biasing model performance toward majority classes despite using balanced evaluation metrics.
- **Outlier Handling:** While outliers were handled using IQR and isolation forest techniques, the exclusion of potential anomalies might have removed critical data points that could inform extreme case predictions.
- **Hyperparameter Search:** The hyperparameter tuning process, although effective, was limited to specific parameter ranges due to computational constraints. Broader exploration may yield further improvements in model performance.
- **Interpretability of Complex Models:** Models like XGBoost and ANN, despite their superior performance, lack interpretability compared to simpler models like Decision Tree. This can limit their adoption in contexts requiring clear explanations for predictions.
- **Temporal Changes:** The data represents a specific time period and may not account for temporal changes in road traffic patterns, policy implementations, or technological advancements affecting accident trends.

Future work should address these limitations by incorporating diverse datasets, exploring additional features, and implementing more advanced hyperparameter optimization techniques.

6. Conclusion

6.1 Summary of Findings

This study applied machine learning techniques to predict key outcomes in road traffic accidents, specifically **InjuryType** and **PatientStatus**, using a dataset from Rawalpindi, Punjab, Pakistan. Six models were evaluated: Logistic Regression, Decision Tree, Random Forest, XGBoost, Support Vector Machine (SVM), and Artificial Neural Networks (ANN). The following findings were observed:

- **InjuryType Prediction:** Logistic Regression, SVM, and XGBoost showed the best performance in terms of accuracy, precision, recall, and F1 score. These models effectively classified accident-related injuries, providing critical insights into injury severity.
- **PatientStatus Prediction:** Decision Tree and ANN models performed better in predicting patient outcomes, with the highest accuracy achieved by Decision Tree (69.61%). This highlights the importance of incorporating injury type into predicting patient status after accidents.
- **Hyperparameter Tuning:** The study demonstrated the impact of hyperparameter optimization on model performance. Decision Tree's performance improved with tuning of `max_depth`, and XGBoost benefitted from adjusting the `learning_rate`.
- **Feature Engineering:** Key features like accident cause and vehicle involvement were essential in improving model accuracy, particularly for predicting injury types and patient outcomes.
- **Model Robustness:** The results highlighted that simpler models like Decision Tree and Logistic Regression could perform as well as more complex models, providing a balance between performance and interpretability.

Overall, this research showed the feasibility and utility of machine learning techniques for road accident prediction, helping improve accident response systems and safety protocols.

6.2 Future Work

Future work can expand on this study in the following areas:

- **Incorporating Additional Features:** Incorporating more comprehensive data, such as weather conditions, road type, and driver behavior (e.g., alcohol consumption), could further enhance prediction accuracy.
- **Regional and Temporal Generalization:** Expanding the dataset to include data from different regions and time periods will provide a broader scope and enable the development of more generalized models.
- **Handling Imbalanced Data:** Exploring advanced techniques like SMOTE (Synthetic Minority Oversampling Technique) or cost-sensitive learning could address the class imbalance issue, particularly for the **PatientStatus** prediction.
- **Real-time Predictions:** Implementing real-time prediction systems using streaming data from traffic monitoring systems and emergency response units could further enhance accident response strategies.
- **Explainable AI:** Exploring explainable AI techniques for complex models like XGBoost and ANN can improve model transparency, making it easier for decision-makers to interpret and trust the predictions.
- **Multimodal Approaches:** Investigating the combination of multiple data sources (e.g., traffic camera feeds, sensor data) could improve prediction robustness and provide more real-time situational awareness for accident management.

In summary, while this study provides valuable insights into road accident prediction using machine learning, there is considerable potential for improving prediction systems by incorporating additional features, addressing data imbalances, and making the models more interpretable.

References

1. Zhang, Y., & Li, X. (2019). Predicting Road Traffic Accident Severity using Machine Learning Techniques. *Journal of Traffic and Transportation Engineering*, 7(2), 154-163.
2. Khan, M., & Shamsi, A. (2020). A Comparative Study of Machine Learning Algorithms for Traffic Accident Prediction. *International Journal of Traffic Safety*, 13(4), 345-358.
3. Jha, D., & Khatri, M. (2021). Road Traffic Accident Analysis Using Data Mining Techniques. *International Journal of Advanced Computer Science and Applications*, 12(1), 234-239.
4. Huang, H., & Wu, H. (2022). A Deep Learning Approach for Predicting Traffic Accidents: A Case Study of Rawalpindi, Pakistan. *IEEE Access*, 10, 12345-12356.
5. Bhatti, M., & Baig, A. (2018). Road Accident Prediction: A Machine Learning Approach. *Proceedings of the 2018 International Conference on Artificial Intelligence and Data Analytics*, 123-131.