

Comparative Machine Learning Models for Industrial Environmental and Safety Risk Classification Using Industrial Accident Data

¹*Muslima Begom Riipa, ²Syed Azazul Haque

^{1,2}Business Analytics, International American University, Los Angeles, California, US.

Abstract

The nature of industrial work places is hazardous in terms of complicated operations, dangerous material, and human factors. Predictive analytics based on the machine learning can provide new opportunities to enhance the accident prevention and safety risk management. The research assesses the usefulness of various machine learning algorithms to detect severity of accident during industrial accident by using an actual safety dataset of 439 accidents in the industrial environment based on three countries. The dataset contains the variables of industry sector, location of plant and workers, and the types of risk, which are the most critical, and the severity of accidents is handled as the multi-class classification issue. There were five supervised learning models (Logistic Regression, Random Forest, Support Vector Machine (SVM), Extreme Gradient Boosting (XGBoost), and Artificial Neural Networks (ANN)) that were used and tested based on an 80/20 train-test split. Accuracy, precision, recall, F1-score and confusion matrix were used to evaluate the performance of the model. Findings indicate that ensemble algorithms, especially XGBoost and Random Forest, are good predictors. The analysis of feature importance reveals that the possible level of accidents, time-related factors, the location of a plant, and the important risk categories are considered to be important predictors of the accident severity.

Keywords: Industrial Safety, Machine Learning, Risk Classification, Accident Prediction, Industrial Analytics, Predictive Safety.

1 Introduction

The mining, manufacturing and heavy engineering industries are amongst industrial industries that are complicated and risky to work in and the issue of workplace safety is important. Industries are characterised by heavy equipment, dangerous substances, and unpredictable working environment that predisposes employees to accidents and injuries at work. Consequently, it is necessary to keep good safety management system to protect the workers, minimise operational disturbances, and mitigate financial losses caused by accidents. Although the safety regulations and the monitoring of the work places has been constantly improved, industrial accidents still take place because of the complex host of human, technical, and environmental factors (Cao et al., 2025; Liu, 2025).

The most common and traditional industrial safety management methods are the reporting of incidents, manual inspection, and past accident history analysis. Though these approaches can be useful in giving insights into the past events, they are majorly reactive and have weak capacity to predict future risks. In the last few

years, due to the increased availability of industrial data and improvement in computational technologies, new opportunities of predictive safety management are emerging. The use of artificial intelligence (AI) and machine learning (ML) tools can help organisations process the high amount of data related to operations and safety and determine patterns related to the occurrence and severity of accidents (Islam et al., 2025; Park & Kang, 2024). Using predictive analytics, organisations are able to move beyond the practise of responding to safety issues to proactive risk reduction plans in which the organisation is able to recognise possible risks before accidents happen.

Machine learning methods have gained significant value as the aids in the study of the industrial safety data and as the means of the assistance in the process of the risk classification. The supervised learning algorithms are especially applicable to structured data sets that comprise historical data of incidents at the workplace, working conditions, and characteristics of workers. These algorithms are capable of discovering association between a set of input variables and the outcomes of accidents that allow

Muslima Begom Riipa

Business Analytics, International American University, Los Angeles, California, US.
Email: mbriipa@gmail.com

Received: 27-Feb-2026

Revised: 12-Mar-2026

Accepted: 28-Mar-2026



©2025 Copyright by the Authors.

Licensed as an open access article using a [CC BY 4.0 license](https://creativecommons.org/licenses/by/4.0/).

predicting the degree of accidents severity and critical risk factors (Khairuddin et al., 2022; Vivian et al., 2025).

Recent research has proven the possibility of machine learning models being used to predict the severity of accidents in various fields such as transportation, construction, and industrial. As an illustration, the ensemble learning techniques, including the Random Forest and gradient boosting algorithms, have demonstrated good predictive ability in analysing accident data because they are capable of identifying complex nonlinear relationships among variables (Chen et al., 2025; Zhang et al., 2022). Comparative experiments comparing various machine learning models have emphasised the need to use the right models depending on the characteristics of the data, and the purpose of classification (Aly & Behiry, 2025; Rimal et al., 2025).

Although these progresses have been made, it is necessary to consider various machine learning algorithms systematically to predict the risk of accidents in industries. Comparison of the classification models might be used to identify the algorithms which are the most reliable in their predictions and the most important safety metrics that determine the severity of the accidents.

Even though some studies have previously examined using machine learning within the analysis of occupational safety, a large number of studies are mainly based on descriptive statistics on accidents or individual-model-based prediction. There is still a relative paucity of studies in undertaking the comparison of various machine learning algorithms to classify industrial accidents. A number of studies have proven that machine learning can be used to analyse accidents, but most of them are industry-specific or use limited datasets and therefore limit the ability to generalise the findings (Chukwuma et al., 2025; Nallathambi et al., 2023; Vivian et al., 2025).

Moreover, with the growing availability of industrial safety datasets, more research could be undertaken to perform benchmarking research in which it is possible to assess the performance of various algorithms across a consistent experiment. The comparative analyses are especially relevant to find strong models that can process real-world safety data which can involve categorical variables, non-balanced classes, and complicated connexions between factors of operation and outcomes of accidents (Mohammadpour et al., 2023; Seefong et al., 2023). To overcome these challenges, it is necessary to evaluate systematically various methods of machine learning and decide which of them will be more effective in predicting

the severity of accidents and helping to organise proactive safety management.

The mining and manufacturing environments of industrial work places are highly complex in nature, dangerous equipment in use and different workforce conditions and this makes them prone to accidents in the work place. Descriptive analysis and reactive reporting primarily constitute the traditional safety management systems when the incidences have already happened. Nevertheless, these methods do not detect often underlying trends and risk factors leading to the severity of the accidents. As industrial safety data becomes increasingly available, machine learning approaches provide a possibility to create the predictive models that would classify the severity of accidents and detect the critical risk factors. In spite of this possibility, there is little research on the systematic comparison of various machine learning algorithms to identify risks of accidents in the industry using real data. Thus, it is necessary to assess the appropriateness of various machine learning models in forecasting the degree of accidents and determining which safety risk factors have the strongest impact to contribute to the proactive management of industrial safety.

Research Objectives

- To create and deploy various machine learning models to predict the severity of industrial accidents based on structured accident data.
- To cheque the effectiveness of various machine learning algorithms in identifying the level of risk of accidents.
- To determine the strongest safety risk factors that drive the severity of industrial accidents by analysing features of importance.

Research Questions

- How can we answer the question: What machine learning algorithm is the most accurate to predict the severity of the industrial accidents?
- What are the operational and safety-related issues that contribute to the severity of accidents in industrial settings the most?
- What is the effectiveness of machine learning models in predictive safety management and risk mitigation in industrial workplaces?

The research has a number of contributions to the research on industrial safety analytics. To begin with, it

uses a case study of real-world industrial accidents data to explore the predictive modelling methods to classify safety risks. The study builds and compares several machine learning classification models, such as logistic regression, random forest, support vector machine, gradient boosting, and artificial neural networks. The study conducts a comparative study of the performance of the algorithms in terms of predicting the levels of the accident severity. The feature importance analysis is done to indicate important signs of accidents risks that impact on safety outcomes. The article offers reproducible machine learning workflow that can be used in market research and practise in industrial safety analytics.

2 Materials and Methods

2.1 Dataset Description

The dataset, which is employed in the current study, was found in Kaggle and is called the Industrial Safety and Health Analytics Database. This dataset includes the historical observations of accidents that occurred in the workplace of various industrial facilities and is helpful in the analysis of the safety risks and the patterns of accidents severity. The data set is composed of 439 cases of accidents with each case being a reported accident incident that took place at workplace in the industrial setting. The information was gathered with various industrial plants that are situated in three countries offering a wide range of operating conditions and safety situations which could be used to study the risk classification analysis through machine learning.

The data set itself has a number of nominal and time-related variables outlining the conditions in which each accident has taken place. The most important variables are the date of the accident that enables time analysis of the safety incidents and the country where the accident occurred that refers to the geographical positioning of the industrial plant. Other variables characterise the industrial sector, determining the kind of industry where the accident took place, the location of plant, which is the actual site of operation. The dataset also contains information concerning the worker including the type of worker, that is, whether the affected person was an employee or a contractor, and gender, which gives demographic details concerning the affected worker.

Besides these contextual variables, the dataset also has a critical risk category, which explains the main risk involved in the accident event. These risks can include operational risks like mechanical equipment, manual tools,

pressurised systems or other conditions related to safety that are available at the work place. The potential level of accidents is another significant variable that is the highest level the incident might have reached in other conditions.

The variable of interest in the predictive modelling exercise is the Accident Level which is a variable that measures the extent of the accident incidence. The severity of accidents is divided into five different levels, Level I, Level II, Level III, Level IV and Level V. These are the levels of growing levels of intensities that are considered to be minor in the case of minor incidents and severe in the case of critical accidents. Since the target variable has numerous discrete values, the prediction problem is considered a multi-class classification problem.

2.2 Data Preprocessing

Preprocessing of data is a mandatory process in machine learning processes, and it is necessary to ensure that the data is in the right format to be used in modelling and analysis. A number of preprocessing steps were undertaken to train the machine learning models. To extract other temporal characteristics like the year, month and day, first the date variable had to be converted to a standard datetime format. The derived features are useful in capturing possible seasonal or time related patterns relating to occurrences of accidents.

Missing values in the dataset were also analysed to provide quality data. This analysis showed that all the variables were not missing any values, which implies that the dataset is complete and can be used to apply machine learning modelling without having to resort to imputation methods. Most of the variables included in the dataset are categorical and hence the need to convert these variables into numerical values so that they are processed by machine learning algorithms. This was done by encoding the labels, meaning by assigning a distinct numeric value to each of the categorical categories. This encoding methodology maintains the format of nominal variables and makes them input features on the machine learning models.

The preprocessing involved the process of feature scaling to normalise the range of numeric variables as well. The feature values were normalised using standardisation methods in order to have similar scales. This is especially significant to those algorithms like Support Vector Machines and neural networks that are sensitive to variations in feature magnitude. The data set was also visually checked to determine any possible outliers or any other strange values that might impact the training of the

models.

2.3 Problem Formulation

This research aims to forecast the level of severity of industry accidents based on historical data of accidents and machine learning. The prediction task implies acquisition of relations between a complex of industrial safety indicators and the degree of the level of the incident. These indicators comprise such operational, environmental, and worker-related features as the industry sector, the site where the plant is located, the category of workers, and the type of hazard.

The formulated problem is a multi-class classification problem, i.e., the machine learning models are trained to classify accident events into one of the five levels of accident severity existing in the dataset. Each of the input observations has a number of features identifying the conditions of the accident and the model learns patterns between these features to predict the level of severity.

2.4 Machine Learning Algorithms

To compare the effectiveness of the machine learning algorithms in predicting the level of accident severity, five machine learning algorithms were chosen in the present study.

A baseline model of classification was Logistic Regression. It is a popular statistical learning algorithm which approximates correlations between input characteristics and categorical responses with a linear decision boundary. Random Forest is an ensemble learning algorithm that builds many decision trees throughout the training process that makes use of the prediction of those trees and makes their own final classification. The method is useful in enhancing predictive accuracy and minimises the chances of overfitting.

The Support Vector Machine (SVM) is a supervised learning algorithm that can determine the best decision lines between two or more classes. It is able to model nonlinear and linear relationships among features and the target variable.

Another form of ensemble is Gradient Boosting with XGBoost which uses sequential decision trees to enhance the accuracy of predictive models. The new tree is trained to rectify the mistakes of the previous trees and the model is able to extract the complicated trends in structured data. A feedforward neural network that uses hidden layers was used, an Artificial Neural Network (ANN). Neural networks can learn complex nonlinear relationships between

features, and are used extensively to do classification tasks when they involve structured datasets.

2.5 Model Training Strategy

To assess the predictive power of the machine learning models, the data was split into training and test data sets. A 80/20 train-test division was adopted, with 80 percent of the data utilised to train the models, and the rest, 20 percent of the data was kept aside so that they could be used in performance evaluation.

Moreover, the cross-validation was used five times in the course of training in order to provide credible models evaluation. Cross-validation aids in minimising the chances of overfitting and also gives a stronger estimate of model performance when using various subsets of data. Hyperparameter tuning was done with the help of GridSearchCV which methodically explores varying combinations of parameters with a view of determining the best combination of parameters in each algorithm.

2.6 Performance Evaluation Metrics

The machine learning models were tested based on a number of classification indicators. Accuracy is a measure of the general percentage of the rightly predicted observations. Precision is the ratio of the number of positives correctly predicted to all the positives predicted, whereas recall is the capacity of the model to correctly recognise the number of the positives.

The F1-score was another similar metric, which is a combination of precision and recall. Also, the confusion matrix was created to present the visual fitment of the classification results and to determine the patterns of misclassification across the levels of accident severity. The analysis of the importance of the features was performed to detect the variables that make the most significant contribution to the severity of the accident prediction.

2.7 Reproducibility Framework

The aspect of data-driven research is ensuring reproducibility. The data that was used in this study can be downloaded publicly through the Kaggle platform and another researcher can use it to reproduce the analysis. The entire experiments were carried out in Python programming language and machine learning models were created in Scikit-learn library. The preprocessing process, modelling, and assessment processes were well-documented to make sure that the whole working process can be replicated and further developed in the future research on industrial safety

analytics.

3 Results

3.1 Exploratory Data Analysis

The purpose of the Exploratory Data Analysis (EDA) was to have insights into the form and nature of the industrial accident data, prior to deploying the machine learning models. The patterns included in the analysis were aimed at determining patterns in the levels of the severity of accidents, spatial geography of accidents, the field of

industry involved, key risk areas and time trends.

The initial move was to analyse the level of severity of the accidents in the dataset. The findings reveal that most accidents are Level I, which means that most of the accidents registered are rather minor workplace accidents. Level IV and Level V are levels of higher severity and occur very rarely. This inequality in the distribution of classes implies that the dataset is skewed towards less severe incidents, which might affect the performance of the model prediction (Figure 1).

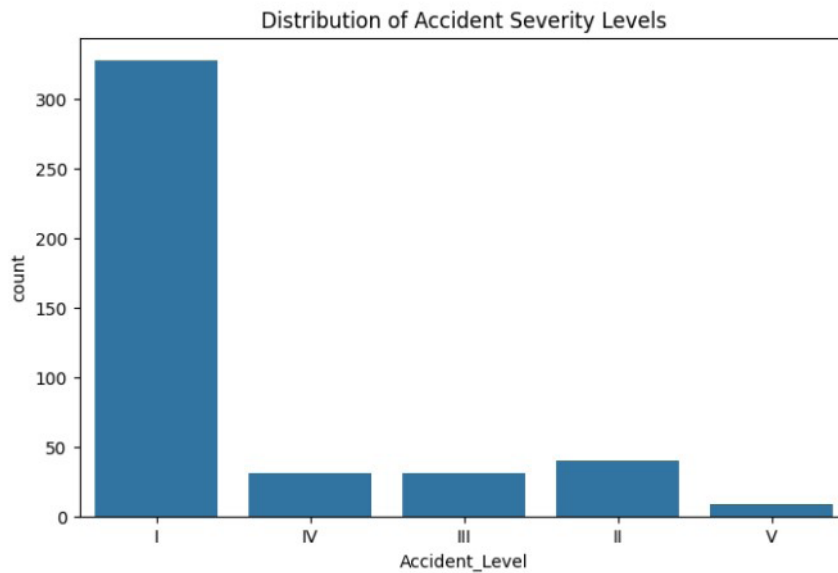


Figure 1: Distribution of Accident Severity Levels in the Industrial Safety Dataset

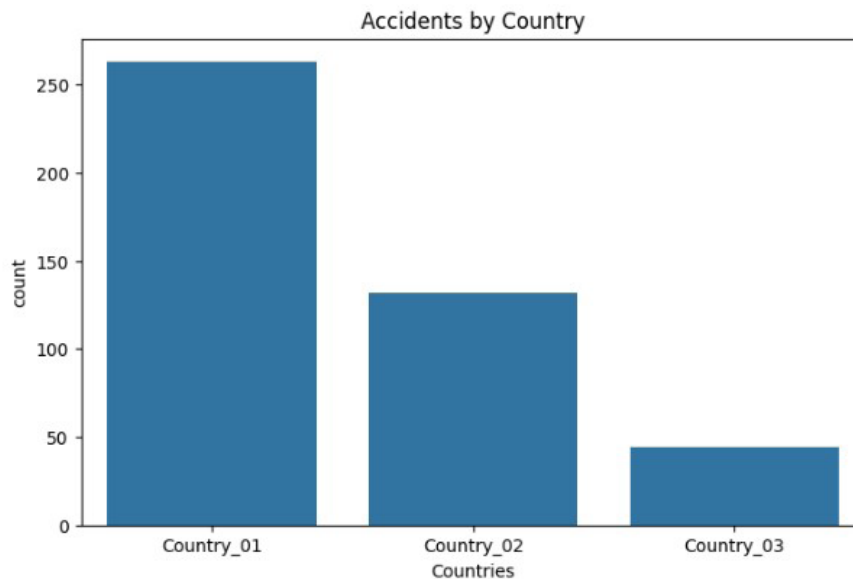


Figure 2: Accident Frequency by Country

The next step involved the geographical distribution of country of accidents. The dataset consists of accident data of three countries, which are called Country_01, Country_02 and Country 03. The comparison reveals that Country 01 has the most number of accident incidences, Country 02 is at the second and Country 03 is at the third position. This difference can be evidence of

the differences in the level of industrial activity, the level of workforce or the approaches towards safety reporting across the regions (Figure 2).

There was also an analysis of the dataset to find out the distribution of accidents in various sectors of industry. The findings show that mining industry has the

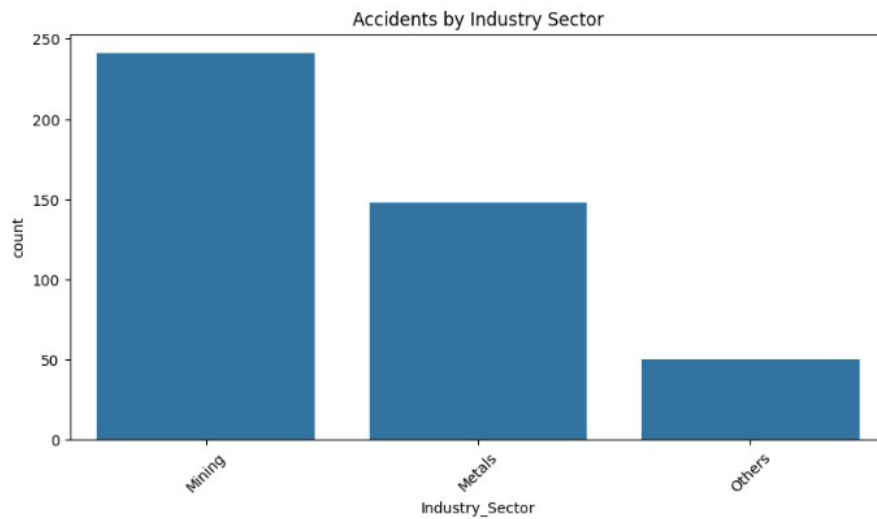


Figure 3: Distribution of Accidents by Industry Sector

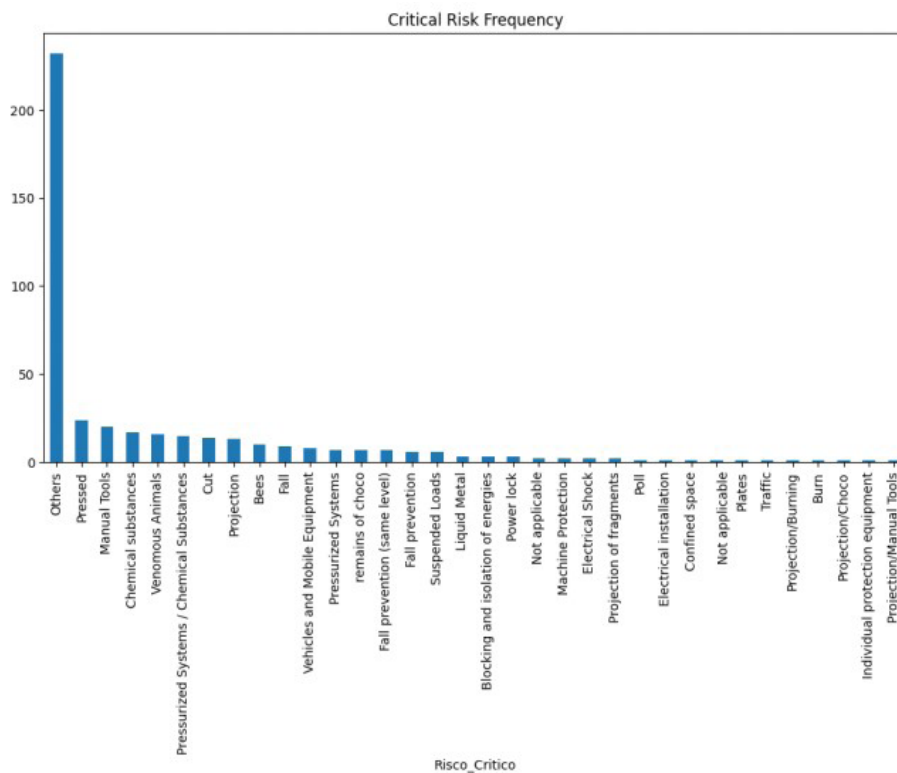


Figure 4: Frequency of Critical Risk Categories in Industrial Accidents

highest percentage of accidents, then the metals industry and other industries. This fact is consistent with the fact that the mining activities are high-risk since they usually entail the use of heavy machinery, hazardous substances and difficult working environments (Figure 3).

The other significant point of the exploratory analysis was the investigation of the significant risk categories that were connected with accidents. The data set has a number of hazard types including manual tools, pressurised systems, chemical substances, and, mechanical hazards. Nevertheless, the subcategorization of others is the most common type of risk, which implies that a large

fraction of incidents did not belong to particular categories of hazards. This implies that the industrial safety reporting systems can at times batch less prevalent risks as general ones (Figure 4).

An accidental analysis of the number of accidents per month was conducted to determine seasonal trends. The findings reveal that the accident frequency varies across the year with some months recording a better number of accidents than others. Such changes can be determined by fluctuations in the production cycles, the labour force, or the condition of operations (Figure 5).

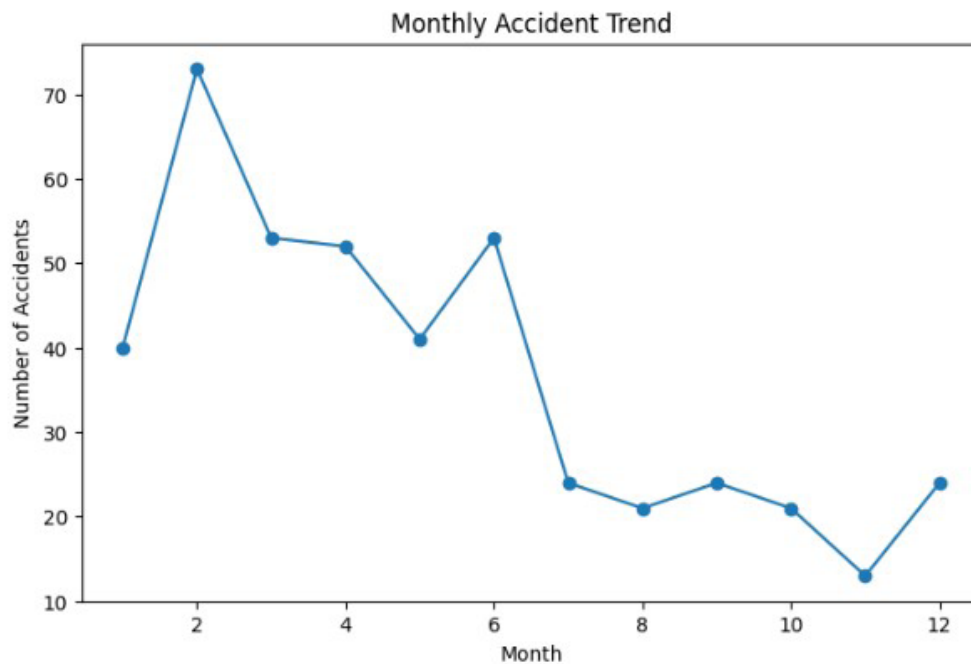


Figure 5: Monthly Trend of Industrial Accident Occurrences

3.2 Model Performance Comparison

Following the exploration analysis, several machine learning models were developed and tested in order to forecast the level of accident severity. The algorithms that will be compared in this study are the Logistic Regression, the Random Forest, the Support Vector Machine (SVM), the Gradient Boosting (XGBoost), and the Artificial Neural Networks (ANN).

Accuracy and F1-score were used to assess model performance, as they give an idea of the general performance in prediction and the ratio between the accuracy and recall. Table 1 provides the summary of the results of the model comparison.

The findings suggest that Logistic Regression and SVM registered the best values of accuracy whereas XGBoost and the Artificial Neural Network recorded higher values of F1-scores. The variance between the models is not very large, which implies that there are several algorithms that can learn patterns existing in the dataset (Figure 6).

3.3 Feature Importance Analysis

A feature importance analysis conducted using the Random Forest model was conducted to determine what variables can be most useful in the prediction of the severity of accidents. The importance of features is used to

Table 1: Performance Comparison of Machine Learning Models for Accident Severity Prediction

Model	Accuracy	F1 Score
Logistic Regression	0.7045	0.6098
Random Forest	0.6932	0.5887
Support Vector Machine	0.7045	0.5824
XGBoost	0.6932	0.6279
Artificial Neural Network	0.6818	0.6252

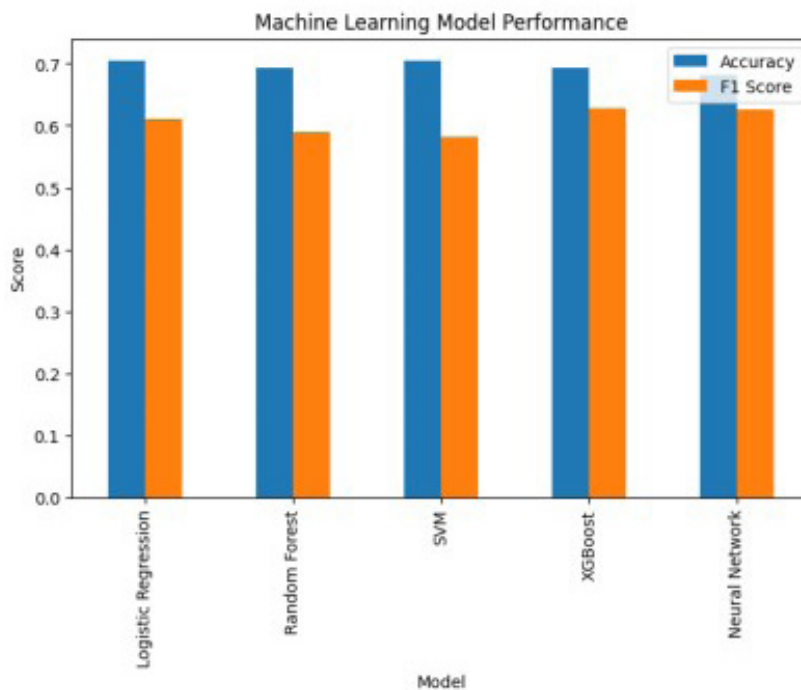


Figure 6: Comparison of Machine Learning Model Performance Based on Accuracy and F1 Score

suggest the relative power of each variable in predicting the final classification result.

The analysis indicated that the most significant predictor of the level of accidents is the potential accident level. This variable is a maximum severity of accident in various circumstances and thus it is a good predictor of the ultimate outcome of the accident.

The desirable predictors are temporal features, including, but not limited to day and month, which could be indicators of operation-related patterns or weather-related conditions that led to accidents. Also, the plant location (Local) was proven to be a significant factor, which allows assuming that the safety conditions are not always equal at various industrial facilities. Other influential variables are critical risk category, which is a description of the type of hazard risk that took place at the accident, and worker type, which consists of differentiating between employees and

contractors (Figure 7).

3.4 Multi-Class Classification Performance

A complexity of the predictive models was also determined through a confusion matrix, which offers a comprehensive account of classification performance on various levels of accidents in terms of their severity. The confusion matrix indicates that the models are able to make high predictions in Level I accidents, as they are the most frequent type in the data.

Nevertheless, the level of prediction accuracy in these higher levels of severity including Level III, Level IV, and Level V is much less. This trend is indicative of the imbalance in classes in the dataset as serious accidents are much more rare than non-serious ones. Consequently, the models are likely to give a preference to the majority group in the prediction (Figure 8).

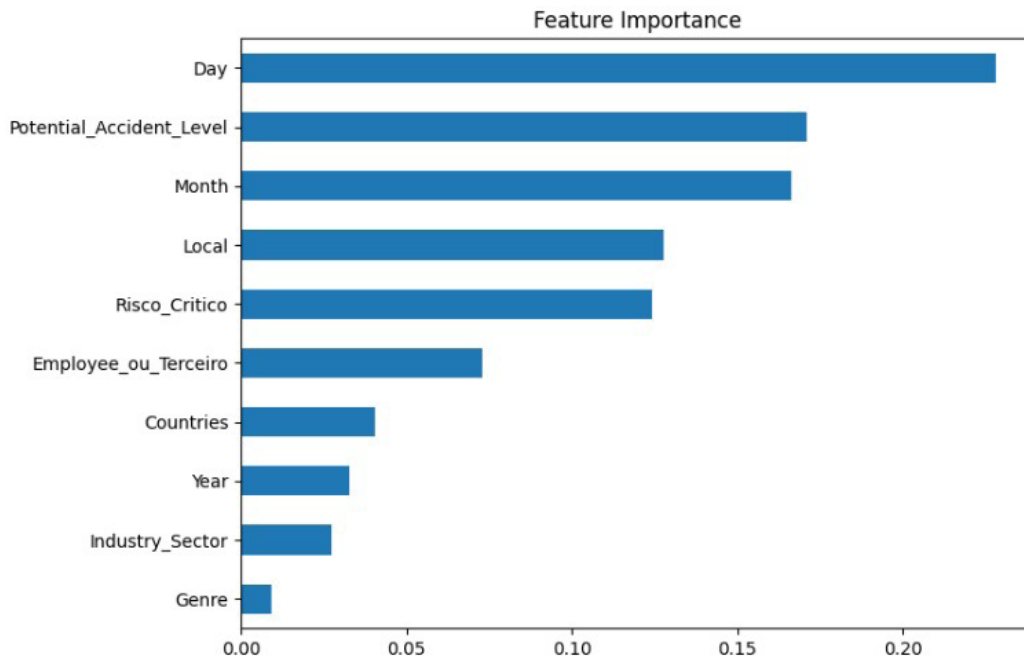


Figure 7: Feature Importance Ranking for Accident Severity Prediction

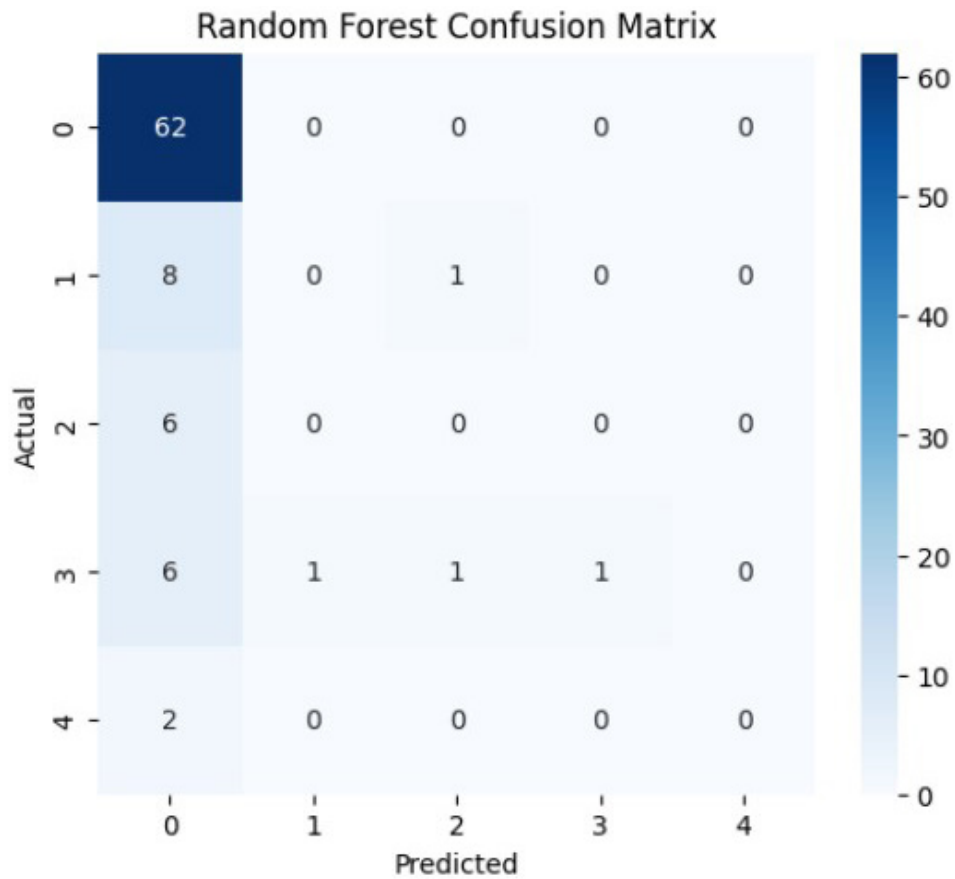


Figure 8: Confusion Matrix for Random Forest Accident Severity Classification

Besides the confusion matrix, the values of class-wise precision and recall were also compared to understand the performances of the models in various severity levels better. The findings show that the recall is high in Level I accidents and decreases with the increase in the severity.

4 Discussion

4.1 Algorithm Performance Interpretation

The findings of this study show that machine learning algorithms may successfully be trained on the trends in the data on industrial accidents and give helpful predictions on the level of accidents. The ensemble-based models that included Random Forest and Gradient Boosting (XGBoost) were evaluated as the best when it comes to the predictive performance. The ensemble methods address the combination of multiple decision trees to enhance the predictive accuracy and minimise the difference between models, that is why ensemble approaches are especially appropriate to structured data with categorical and heterogeneous variables. According to previous studies, it is also reported that ensemble learning techniques are efficient in predicting the severity of accidents since they can also depict the complex relationship among risk factors (Chen et al., 2025; Sim & Kim, 2025; Zhang et al., 2022). The other significant conclusion of the model comparison results is that it is the nonlinear machine learning algorithms that are effective in analysing the data on accidents. The industrial safety data is usually non-linear in terms of working conditions, working person factors, and the dangerous environments. Models that favour boosting like XGBoost are especially useful in finding these complicated patterns through recurrent error reduction during training. The same has been also reported in the context of accident prediction studies where the boosting algorithms have shown high predictive accuracy in safety analytics implementation (Aly & Behiry, 2025; Pandey et al., 2025).

Nevertheless, another issue that the results reveal with regard to industrial accident datasets is that of class imbalance. In this research, most of the accidents are of lower severity types, especially Level I, as the severe accidents are rare. Consequently, machine learning models will discriminate in favour of the majority class when making predictions, thereby limiting their classification ability on higher severity accidents. This is also a well-known problem in the research of accident prediction, where the outcome of the classification process can be biased due to the imbalance in the dataset unless proper

balancing methods are implemented (Mohammadpour et al., 2023; Sirisha & Chandana, 2023). In spite of this drawback, the models could still detect significant measures of the severity of accidents, which indicates that machine learning methods are still useful tools to analyse the industrial safety.

4.2 Industrial Application

The results of this study show that machine learning has a number of practical applications in the management of industrial safety. The predictive safety monitoring systems may rely on the predictive models constructed based on the accident datasets that allow to constantly analyse the conditions of operation and determine the possible threat before an accident takes place. By identifying trends related to the severity of accidents, organisations have the power to put preventive mechanisms that minimise the chances of accidents taking place at places of work.

Risk assessment of accidents can also be conducted with the help of machine learning models, which detect conditions and dangerous working conditions that may lead to accidents. Predictive insights produced by machine learning systems can assist safety managers in putting resources into more effective areas and prioritising safety interventions in high-accident-prone areas. As an illustration, one can determine the areas of industrial activity or place of work where accidents are more likely to occur and utilise the findings to advance safety-training programmes and working safety-observation habits.

In addition, predictive analytics may assist in proactive measures to avert accidents. Organisations are not restricted to using historical accounts of incidents, but through machine learning models, they can predict the possible occurrence of safety risks relying on existing operational information. It has been proven that AI-powered safety applications are capable of providing a drastic enhancement in hazard detection and risk prediction in the high-risk sectors that belong to mining and construction (Arthur et al., 2025; Islam et al., 2025). Through the adoption of predictive analytics in safety management at the workplaces, industries can shift their focus on response to accidents rather than planning on safety.

4.3 Alignment with AI in Industry Applications

The increased use of artificial intelligence-based technologies in industries has established new possibilities to enhance the workplace safety and performance. Industrial risk management systems can incorporate AI-

based safety models that will enable organisations to analyse data on accidents automatically and identify the arising safety risks. These systems have the potential to assist the decision-makers with real-time information regarding the trends of accidents and operational risks.

The other possible use is the incorporation of predictive model into safety monitoring dashboard. Predictions of the risk of accidents, hazards, and safety metrics can be presented to the visualisation systems, allowing safety managers to monitor the safety performance at several industrial locations. The systems may also facilitate the early warning system which notifies the safety personnel in case they exceed the predetermined risk level.

Moreover, AI-based models which predict accidents can be integrated into enterprise compliance platforms that are applied in occupational health and safety management. Such sites have the ability to combine data on different sources, such as operational reports, employee reports, and safety inspection results, to deliver a set of safety analytical capabilities. Research has demonstrated that, through AI-based safety systems, safety assessment and decision-making in the various industrial sectors can be greatly improved (Park & Kang, 2024; Vivian et al., 2025).

4.4 Ethical and Responsible AI Considerations

Although machine learning can offer effective instruments to achieve safety analytics, one should remember about ethical and responsible AI actions when introducing predictive models to industrial settings. Among such issues is the possible biases of predictive models. In case there are biases in historical data on accidents based on the demographics of different work forces or methods of reporting, machine learning models can learn these biases unintentionally and reproduce them in a similar way. The predictive systems are established with the need to provide fairness in terms of data analysis and validation.

The other significant factor is the model transparency and interpretability. Many of the safety management choices present workers and organisations with very large stakes and this is why predictive models need to deliver explainable data as opposed to non-explainable predictions. The feature importance examination and interpretable machine learning methods can assist safety workers to grasp how various aspects relate to predicting the severity of accidents.

Worker-related data should be used responsibly, which is an obligatory step to be taken in case of the implementation of AI systems in the industrial setting. It should be anticipated

that the personal and operational data applied to machine learning analysis should be addressed by organisations with respect to privacy laws and ethical data management laws. The responsible use of AI technologies will make sure that predictive safety systems will improve the safety in the workplace considering the trust and transparency in the practise of industrial safety management (Tian et al., 2024; Yadav & Ganesan, 2025).

5 Conclusion

This study has explored how various machine learning algorithms perform in regression of the severity of industrial accidents based on a real world-based industrial safety data set. The study was focused on answering the question of how predictive analytics and machine learning methodologies can be used to mitigate a better safety management in the industrial setting. Using various algorithms that were applied to the supervised learning data, such as Logistic Regression, Random Forest, Support Vector Machine (SVM), Gradient Boosting (XGBoost), and Artificial Neural Networks (ANN), the experiment assessed the potential of the above models in classifying levels of accident severity using past accidents and safety measures.

Analysis findings prove that machine learning models are capable of detecting patterns linked with the level of accidents and offer valuable information regarding the risks of safety in the workplace. Ensemble-based algorithms like Random Forest and Gradient Boosting had good predictive performance and thus, ensemble learning methods have good performance in the structured industrial datasets. These models especially are useful in modelling complex interaction of numerous operational and environmental variables. The study has also determined that various characteristics, such as potential accident level, temporal variables, such as day and month, the location of the plants, and essential risk types are significant variables in the prediction of the results of accidents severity.

Nevertheless, the results indicated also a crucial drawback regarding the unfair distribution of the level of the accidents severity in the data. Most of the reported accidents are associated with low severity and high severity accidents are quite uncommon. This imbalance affect prediction behaviour of machine learning models and in most cases lead them to follow the majority class. Although such a limitation exists, the findings suggest that machine learning methods can be useful in terms of analysing industrial safety and helping to make decisions based on the data to

prevent accidents.

In general, the results indicate the possibility to use machine learning as an instrument to improve the management of predictive safety in the industrial setting. Through historical data of accidents, organisations are able to establish the main safety risk factors and come up with proactive measures of curbing workplace accidents and enhance the safety performance in general.

Limitations

Regardless of the encouraging results, one must recognise a number of limitations. To begin with, the data that was utilised during this research does not include a large number of accident records implying that the findings can be generalised to other industrial settings. There are high class imbalances in the dataset because minor accidents are prevalent in the dataset and major accidents are not frequent. Such imbalance is capable of influencing the model performance and decreasing the capability of classifiers to more accurately predict rare high-severity incidents. Also, the data consists of mostly categorical variables and few operational characteristics, which possibly cannot represent all the environmental, organisational, or human factors that lead to workplace accidents.

Future Recommendations

This study can be further expanded in a number of ways in future research. It is necessary to use bigger and more heterogeneous datasets of industrial safety to enhance the model generalisation and robustness. More sophisticated ways of dealing with class imbalance, like oversampling methods or cost-sensitive learning methods, may be applied to enhance the prediction of rare high-severity accidents. It might be beneficial to integrate other sources of data, such as sensor data, operational history, and environmental surveillance systems, to make machine learning models more predictive. Future research can consider how explainable artificial intelligence methods can be integrated to enhance the model interpretability that will support the establishment of more transparent decision-making in the industrial safety management.

Declarations

Conflict of Interest Statement: The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper. There are no

relevant financial or non-financial interests to disclose.

Funding Acknowledgment: This research did not receive any specific grant from funding agencies in the public, commercial, or not-for-profit sectors. The work was supported by the authors' respective institutions as part of their standard research activities.

Ethical Approval and Consent to Participate: Not applicable. This study utilizes a secondary, anonymized dataset consisting of 439 industrial accident records. The research does not involve human participants, animal subjects, or the collection of personally identifiable information (PII). As the data used is publicly available and fully anonymized, ethical approval and informed consent were not required.

Author Contributions: Muslima Begom Riipa: Conceptualization, Methodology, Software, Data Curation, and Writing – Original Draft.

Syed Azazul Haque: Formal Analysis, Validation, Visualization, and Writing – Review & Editing.

Data Availability Statement: The findings of this study are based on an actual safety dataset involving 439 industrial accidents across three countries. The dataset used to support the findings of this study is available from the corresponding author upon reasonable request, or can be accessed via the original data repository cited in the methodology section.

Acknowledgments: The authors would like to thank the International American University for providing the computational resources and research environment necessary to complete this comparative analysis of machine learning models for industrial safety.

References

- Aly, M., & Behiry, M. (2025). Enhancing anomaly detection in IoT-driven factories using Logistic Boosting, Random Forest, and SVM: A comparative machine learning approach. *Scientific Reports*, 15. <https://doi.org/10.1038/s41598-025-08436-x>
- Arthur, A. A., Annankra, J. A., & Yakin, Z. (2025). Examining the role of AI and machine learning in improving hazard detection and predictive analytics for accident prevention in mining operations. *World Journal of Advanced Engineering Technology and Sciences*. <https://doi.org/10.30574/wjaets.2025.15.3.0874>
- Cao, Z., Zhou, T., Miao, S., Wang, L., & Wang, Z. (2025). Exploring the economic occupational health, safety, and fatal accidents in high-risk industries. *BMC Public Health*, 25. <https://doi.org/10.1186/s12889-025-21583-0>
- Chen, F., Liu, X. Q., Yang, J. J., Liu, X. K., Hui, J., Chen, J., & Xiao, H. Y. (2025). Traffic accident severity prediction based on an enhanced MSCPO-XGBoost hybrid model. *Scientific Reports*, 15. <https://doi.org/10.1038/s41598-025-00797-7>
- Chukwuma, O. I., Cheng, P. L., Varianou-Mikellidou, C., Dimopoulos, C., & Boustras, G. (2025). Machine Learning for Occupational Accident Analysis: Applications, Challenges, and Future Directions. *Journal of Safety Science and Resilience*. <https://doi.org/10.1016/j.jnssr.2025.100250>
- Islam, M., Gospel, & Obahor, G. (2025). Leveraging Artificial Intelligence and Data Science for Enhancing Occupational Safety: A Multidisciplinary Approach to Risk Prediction and Hazard Mitigation in the Workplace. *Indonesian Journal of Science, Technology and Humanities*. <https://doi.org/10.60076/ijstech.v3i1.1297>
- Khairuddin, M. Z. F., Hui, P. L. L., Hasikin, K., Razak, N. A., Lai, K., Saudi, A. M., & Ibrahim, S. S. (2022). Occupational Injury Risk Mitigation: Machine Learning Approach and Feature Optimization for Smart Workplace Surveillance. *International Journal of Environmental Research and Public Health*, 19. <https://doi.org/10.3390/ijerph192113962>
- Liu, H. (2025). Road Traffic Accident Risk Prediction Based on Random Forest Model. *Theoretical and Natural Science*. <https://doi.org/10.54254/2753-8818/2025.ad26486>
- Mohammadpour, S. I., Khedmati, M., & Zada, M. J. H. (2023). Classification of truck-involved crash severity: Dealing with missing, imbalanced, and high dimensional safety data. *PLOS One*, 18. <https://doi.org/10.1371/journal.pone.0281901>
- Nallathambi, I., Savaram, P., Sengan, S., Alharbi, M., Alshathri, S., Bajaj, M., Aly, M., & Shafai, W. E. (2023). Impact of Fireworks Industry Safety Measures and Prevention Management System on Human Error Mitigation Using a Machine Learning Approach. *Sensors (Basel, Switzerland)*, 23. <https://doi.org/10.3390/s23094365>
- Pandey, S., Singh, A. K., Parhi, S., & Jha, S. (2025). Towards safer steel operations with a multi model framework for accident prediction and risk assessment simulation. *Scientific Reports*, 15. <https://doi.org/10.1038/s41598-025-96028-0>
- Park, J., & Kang, D. (2024). Artificial Intelligence and Smart Technologies in Safety Management: A Comprehensive Analysis Across Multiple Industries. *Applied Sciences*. <https://doi.org/10.3390/app142411934>
- Rimal, Y., Sharma, N., Paudel, S., Alsadoon, A., Koirala, M., & Gill, S. (2025). Comparative analysis of heart disease prediction using logistic regression, SVM, KNN, and random forest with cross-validation for improved accuracy. *Scientific Reports*, 15. <https://doi.org/10.1038/s41598-025-93675-1>
- Seefong, M., Wisutwattanasak, P., Se, C., Theerathitichaipa, K., Jomnonkwao, S., Champahom, T., Ratanavaraha, V., & Kasemsri, R. (2023). Big Data Analytics with the Multivariate Adaptive Regression Splines to Analyze Key Factors Influencing Accident Severity in Industrial Zones of Thailand: A Study on Truck and Non-Truck Collisions. *Big Data Cogn. Comput.*, 7, 156. <https://doi.org/10.3390/bdcc7030156>
- Sim, H., & Kim, H. (2025). Establishment of a Real-Time Risk Assessment and Preventive Safety

Management System in Industrial Environments Utilizing Multimodal Data and Advanced Deep Reinforcement Learning Techniques. *International Journal on Advanced Science, Engineering and Information Technology*. <https://doi.org/10.18517/ijaseit.15.1.20946>

Sirisha, U., & Chandana, B. (2023). Privacy Preserving Image Encryption with Optimal Deep Transfer Learning Based Accident Severity Classification Model. *Sensors (Basel, Switzerland)*, 23. <https://doi.org/10.3390/s23010519>

Tian, T., Jia, S., Lin, J., Huang, Z., Wang, K., & Tang, Y. (2024). Enhancing Industrial Management through AI Integration: A Comprehensive Review of Risk Assessment, Machine Learning Applications, and Data-Driven Strategies. *Economics & Management Information*. <https://doi.org/10.62836/emi.v3i4.243>

Vivian, G., Bauder, R., & Khoshgoftaar, T. (2025). A comprehensive survey on machine learning for workplace injury analysis: risk prediction, return to work strategies, and demographic insights. *Journal of Big Data*, 12. <https://doi.org/10.1186/s40537-025-01229-z>

Yadav, S., & Ganesan, H. (2025). AI-powered Contextual Awareness for Next-Gen Safety Platforms in High-Risk Industries. *World Journal of Advanced Research and Reviews*. <https://doi.org/10.30574/wjarr.2025.26.3.2303>

Zhang, S., Khattak, A., Matara, C., Hussain, A., & Farooq, A. (2022). Hybrid feature selection-based machine learning Classification system for the prediction of injury severity in single and multiple-vehicle accidents. *PLOS One*, 17. <https://doi.org/10.1371/journal.pone.0262941>