



Predicting the probability of occupational accidents occurrence in a Portuguese retail company

Inês Sena ^a,* , Felipe G. Silva ^a, Ana Cristina Braga ^b, Florbela P. Fernandes ^a,
Clara B. Vaz ^a, Maria F. Pacheco ^a, Paulo Novais ^b, José Lima ^a, Ana I. Pereira ^{a,b}

^a Research Center in Digitalization and Intelligent Robotics (CeDRI) and Laboratório Associado para a Sustentabilidade e Tecnologia em Regiões de Montanha (SusTEC), Instituto Politécnico de Bragança, Campus de Santa Apolónia, Bragança, 5300-253, Bragança, Portugal

^b ALGORITMI Research Centre, LASI, University of Minho, Campus de Gualtar, Guimarães, 4710-057, Braga, Portugal

ARTICLE INFO

Keywords:

Data mining
Machine learning
Occupational accidents
Predictive analysis
Retail sector

ABSTRACT

Workplace accidents are a global problem impacting companies and society, as employee well-being and productivity/profit can be affected. Portugal ranks fifth among European Union countries despite efforts to reduce their frequency. Predictive solutions have demonstrated promising results in several economic sectors, but the retail sector, the country's third-largest in accident records, remains unexplored. This study proposes a predictive model based on the Multilayer Perceptron (MLP) algorithm to calculate the probability of risk situations occurring in a retail company. Ten databases provided by the company were analyzed and combined into a single dataset using impact scores. The predictive model was developed to predict risk situations in all the company's stores throughout two working days, the current and the next, and the four working shifts. The predictive model was implemented and tested in an integrated system for nine months and achieved 92% accuracy and a 29% precision rate in identifying risk situations. It is concluded that this approach provides a practical solution to assist companies and occupational health and safety teams prevent and minimize workplace accidents, contributing to increased safety and well-being.

1. Introduction

Industry 4.0 emerged with the promise of increasing productivity by integrating digital production systems with analyzing and communicating all data generated in an intelligent environment (Badri et al., 2018). Among the main advances, real-time communication, big data, human-machine cooperation, remote sensing, process monitoring and control, use of autonomous equipment, and interconnectivity stand out as vital and indispensable resources in the modern industry (Badri et al., 2018).

As practices to increase productivity evolve, safety and risk management practices must follow this evolution since occupational accidents considerably affect human capital and negatively impact the productivity and competitiveness of workers and companies (Khahro et al., 2020).

Consequently, in recent years, researchers have begun to explore innovative solutions such as predictive analytics to improve conditions and ensure safety in the workplace (Blanchard, 2021). Emerging approaches to prevent and minimize workplace accidents in several business sectors, such as construction (Sarkar et al., 2020), manufacturing (Sarkar et al., 2023), mining (Santibáñez et al., 2013), energy (Ajayi et al., 2020), have started to appear in the literature.

Despite these advancements, the retail sector remains underserved regarding publications and strategies to reduce workplace accidents, even though it involves tasks with significant risks for employees. In 2022, in Portugal, the retail sector ranked as the third economic activity with the highest workplace accidents (PORDATA, 2024).

Therefore, this study aims to develop a predictive model to predict risk situations in the retail sector over two consecutive days and for the four work shifts of the company in the case study, to anticipate such events. The case study will be Portugal's leading food retail company, with ten databases available.

To evaluate the predictive model performance, it will be deployed within the company through an integrated system.

This predictive model proposal is innovative in addressing a gap in the literature by developing an innovative solution for the retail sector, which remains unexplored. Furthermore, it differentiates itself by integrating several company information, providing a more complete and robust prediction of workplace accidents. Its ability to consider different work shifts reinforces attention to safety at multiple times, adjusting preventive measures to the specific needs of each period. With

* Corresponding author.

E-mail address: ines.sena@ipb.pt (I. Sena).

forecasts for two working days, the predictive model allows employees and key managers to proactively prepare for imminent risks, promoting more effective security management.

This paper is organized as follows. Section 2 provides a comprehensive literature review, which describes relevant studies on predicting workplace accidents and the main theoretical concepts necessary to better understand the study. The case study and the methodology for creating the dataset and the predictive model are presented in Section 3. The results for the prediction model in a real context are described in Section 4. Finally, Section 5 summarizes the study's key findings, discusses their implications, and offers insights for future research directions.

2. Literature review

The present study proposes developing an algorithm to estimate the probability of risk situations in the retail sector. Reviewing the research literature and theoretical concepts supporting the analysis is essential to validate and strengthen this approach.

2.1. Predict workplace accidents

Implementing intelligent solutions using Artificial Intelligence (AI) algorithms to predict and mitigate workplace accidents has emerged as a promising approach to improve workplace safety. However, the effectiveness of different algorithms and the definition of the most appropriate dataset to predict workplace accidents are still questionable. The literature review reveals that the majority of studies use Machine Learning (ML) and Deep Learning algorithms, highlighting Random Forest (RF) (Choi et al., 2020; Kang and Ryu, 2019; Poh et al., 2018), Artificial Neural Network (ANN) (Ayhan and Tokdemir, 2019; Koc et al., 2022; Sarkar et al., 2022), and Gradient Boosting Machine (GBM) (Ajayi et al., 2020).

The literature review is based on studies of construction and steel plants. These studies predominantly use historical data from accident records, considering predictors such as event description, accident cause, working conditions, employee characteristics (age, type of contract), and equipment status (Ayhan and Tokdemir, 2019; Choi et al., 2020; Kang and Ryu, 2019; Koc et al., 2022; Poh et al., 2018; Sarkar et al., 2020, 2023, 2022). RF-based approaches present the best accuracy values to predict work accidents. As in civil construction, RF has achieved up to 92% accuracy rate in predicting accidents (Choi et al., 2020; Kang and Ryu, 2019; Poh et al., 2018).

Although ML algorithms such as RF are widely used and present consistent results, there is a growing interest in exploring Deep Neural Networks (Koc et al., 2022; Sarkar et al., 2022).

Despite technological advances, most studies focus on specific sectors, such as construction and heavy industry, while areas, such as the retail sector, remain unexplored. Furthermore, many models limit themselves to using accident records as the main database, failing to incorporate contextual variables, such as task characteristics and company operational dynamics.

Based on this gap, this study proposes developing a predictive model based on the Multilayer Perceptron (MLP) algorithm. Although widely explored in other fields, its application to predicting workplace accidents remains largely unexplored (Mosavi et al., 2021; Ohadi et al., 2022; Qteat and Awad, 2021). The proposal stands out for integrating diverse data, including information about tasks, company population profile, and transaction volume, among others, to create a more robust and adaptable model for the retail sector.

2.2. Concepts

Predictive analysis uses techniques to predict future outcomes based on historical and current data (Liz-Domínguez et al., 2019) with applicability in several areas, particularly in predicting workplace accidents (Blanchard, 2021). Predictive models can be developed using Statistical Learning or ML (Liz-Domínguez et al., 2019).

As part of the literature review, it was observed that ML algorithms are generally explored to solve classification or regression problems, by predicting the severity, injuries, risks, and occurrence of workplace accidents (Oyedele et al., 2021). This study addresses a classification problem that aims to predict the probability of an accident occurring at a given time. Classification algorithms categorize new data based on patterns observed in historical data (Recal and Demirel, 2021).

The MLP was chosen among several ML algorithms. This is one of the most frequently used neural network algorithms, inspired by the functioning of the human brain. The MLP comprises several layers of interconnected artificial neurons capable of learning complex patterns in input data to perform accurate classifications. During training, synaptic weights are adjusted to optimize model performance using the back-propagation algorithm (Hastie et al., 2009).

The MLP model was implemented following standard configurations, with a parallel processing system consisting of input, hidden, and output layers (Ghodrati et al., 2018). In the input layer, each neuron represents a feature of the dataset. The hidden layer transforms the information received, and the number of layers and neurons is defined depending on the complexity of the problem. In the output layer, the number of neurons corresponds to the number of predicted dependent variables. Model training aims to reduce the error between the predictor and target variables (Ghodrati et al., 2018).

The confusion matrix was used to evaluate the model's performance, allowing the results to be interpreted. Matrix components include:

- True Positives (TP) – correctly predicted risk situations.
- True Negatives (TN) – correctly predicted non-risk situations.
- False Positives (FP) – non-risk situations data falsely predicted as risk situations.
- False Negatives (FN) – risk situations falsely predicted as non-risk situations.

Based on the matrix components, performance metrics were calculated, among them, *Accuracy* and *Precision* stood out (Grandini et al., 2020). *Accuracy* evaluates the model's ability to identify the class in the entire dataset correctly, defined in Eq. (1) where:

$$Accuracy = \frac{TP + TN}{TP + FP + TN + FN} \quad (1)$$

Precision, defined in Eq. (2), refers to the model's reliability in correctly identifying a specific class.

$$Precision = \frac{TP}{TP + FP} \quad (2)$$

In summary, this study develops a predictive model using the MLP algorithm to predict the probability of risk situations occurring at different times of the day. Model performance is evaluated based on the confusion matrix and the *Accuracy* and *Precision* metrics.

3. Methodology

This study's main objective is to develop a predictive model to minimize and prevent the risk of occupational accidents in the retail sector. A crucial factor in achieving this goal is the dataset to be used.

The case study, the datasets design approach, and the predictive model methodology, will be presented in this section.

Table 1
Description of the areas and sectors of the company's stores.

Area	Sections	Description of products sold or main task
Cashiers	Cashiers	Payment for the product.
Flows	Food	Milk, cereals, flour, biscuits, etc.
	Non-food	Drinks, pet food, cleaning products, etc.
Fresh Food	Bakery	Bread and cakes.
	Butchery	Fresh meat.
	Charcuterie	Cheese, ham, bacon, etc.
	Fishmonger	Fresh fish.
	Fruit and vegetables	Fresh fruit and vegetables.
Support	Decoration	Preparation and placement of promotional posters, etc.
	Food security	The team that analysis the food quality.
	Maintenance	Team to resolve machine and workplace problems.
	Management	Customer service, operations management, etc.
	Reception	Storing and receiving product stock., etc.
	Security	Elements that maintain local safety.
Textile	Textile	Clothing, shoes, interior and exterior home decoration.

3.1. Case study

The retail company in the case study has ten available databases with different information on more than 300 stores distributed throughout Portugal. It is important to emphasize that the data provided does not contain confidential information, such as health data or employees' identification. Each store has different work areas and sections, as shown in [Table 1](#).

Each store, area, and work section predominantly incorporates data from the ten databases provided by the company. Below is a brief description of each of the databases:

- **Accident Records** detailed information of each recorded incident. It is composed of 145 variables with information on the general characteristics of the injured workers (age, length of service, etc.), the condition of the accident (place, time, sector, task performed at the time of the accident, etc.), the damages (severity, type of injury, etc.) and the cause of the accident.
- **Demographic data** is information about the different characteristics of a population distributed across 17 variables. It includes information on the number of working hours, age, length of service, number of employees, type of contract (full-time or part-time), and the number of employees by education level.
- **Absenteeism** represents the absence of one or more employees during the workday, whether due to being late, partial-day absence, or missing multiple days. It comprises 149 variables, including the value of fees per year and month since 2020. The types of absenteeism rates include total sick leave, COVID-19, unjustified absences, parenting, accidents, and other causes.
- **Preventive Safety Actions** These are records of risk situations or unsafe conditions observed by members of the Occupational Health and Safety (OSH) team when visiting the store. It comprises 22 attributes associated with the verified situation, such as the number of risk situations observed, the number of unsafe conditions observed, the number of observations currently resolved, and the level of risk of the observation, among others.
- **Action plans** are prepared following the intervention of third parties or employees to correct and improve working conditions observed during a workplace audit. It comprises 23 variables that involve the action to be resolved, the typology (conforming or non-conforming), degree of criticality, occurrence process, cause, opening and closing date of resolution, among others.
- **Ergonomic Workplace Assessment (EWA)** consists of analysis values of the posture and movements adopted by employees when performing their duties at fourteen different points, by average

number of employees and security technicians, as well as total values through the calculation of several methods of Rapid Upper Limb Assessment (RULA), Rapid Entire Body Assessment (REBA), among others. It consists of 106 variables corresponding to the average values of employees and technicians for each point, the total values of each point, and the values calculated through each method for each process, micro, and macro task of each section and store.

- **Hazard Identification and Risk Assessments (HIRA)** reveals the risk levels associated with the details of micro tasks in each work section. It consists of 31 attributes that indicate the risk values for each process and the micro and macro tasks of each section of each store.
- **Transactions** daily transactions records since 2019. This dataset includes only three attributes: store designation, record date, and the number of tickets sold.
- **Audits** contains data from an investigation, often unannounced, carried out into operations to evaluate, among other things, the implementation of safety standards. These can be carried out by the company's occupational health and safety team members or by external evaluators, conducted by independent entities. The investigation is carried out by visualizing the fulfillment of several questions, resulting in a dataset of 22 attributes. These fall into the number of compliant issues, the number of non-compliant problems, the number of unobservable issues, the number of unfulfilled issues, the planned date of the visit, and the effective date, among others.
- **Training** corresponds to the essential instructions required to perform tasks safely and effectively. It includes 50 variables that reveal information such as the name of the training, type of standard, learning priority, and number of specific training carried out, general training carried out, specific training stipulated for each unit, and transversal training, among others.

The data must be analyzed and pre-processed to reduce the information according to the methodology outlined by [Silva et al. \(2023\)](#) to use the presented databases in the developed predictive model. In that study, databases were reduced using two statistical methods to find information related to the accident event for each database, resulting in a 48% reduction in information.

Furthermore, a column representing the company's four work shifts was added to each database so that the probability of risk situations occurring for each shift could be predicted. [Table 2](#) shows the different work shifts of the company.

Table 2
Designation of the company's work shifts.

Shift	Hours interval
Open	02:00 am to 07:59 am
Morning	08:00 am to 12:59 am
Afternoon	01:00 pm to 07:59 pm
Close	08:00 pm to 01:59 am

Table 2 provides an overview of the company's work shifts. The "Open" shift corresponds to the store's preparation period, typically starting around 6:00 am, though some stores begin earlier, at 4:00 am, to accommodate the production of bread and pastries. The "Morning" and "Afternoon" shifts represent the store's operational hours with customers. The "Close" shift includes the final operational hours and activities such as cleaning, organizing, and preparing the store for the next day. All times are listed assuming continuous operation, as insurers also classify accidents occurring on the commute to and from work as workplace accidents.

3.2. Dataset

Integrating data from all datasets into a unified dataset is difficult in these risk situations identification context. Through that, the developed approach is based on calculating impact scores for each database and linking them into a single dataset. Preliminary results on this approach can be seen in Sena et al. (2023), where an impact calculation was generated using the number of records in the **Preventive Safety Actions** and in the **Action Plans** databases, which achieved promising results in predicting work accidents and identifying the accident event. However, it did not contain the accident event as an influence on the impact score. Therefore, this initial strategy will be adapted and deepened considering ten databases.

In the **EWA** and **HIRA** databases, the impact score results from the normalization of the total values of the assessments of postures and tasks carried out in each store and section. The min-max normalization was applied, which converts numeric values to a specific range between 0 and 1. For the **EWA** database, the risk values associated with the postures and movements adopted by employees while performing their tasks in the store and section are considered. In the **HIRA** database, risk assessment values related to the performance of functions in each store and work sector are used. This process generates, in total, four impact scores.

For the remaining eight databases, it was necessary to establish an impact score calculation. The first step of the calculation is the same as presented in (Sena et al., 2023), in which the average number of records than occurred is given by Eq. (3).

$$\bar{r}_{uw} = \frac{\sum r_{uws}}{n_{uw}} \quad (3)$$

where the average of records per section (\bar{r}_{uw}) was obtained considering the records per section of each store (r_{uws}) and the number of sections (n_{uw}).

The number of accidents that occurred depends on the amount of information recorded on each database on the same registration date, which influences the Impact Score (I) of a given dataset, Equation Eq. (4):

$$\bar{I} = \frac{n_a - \bar{r}_{uw}}{n_a} \quad (4)$$

where n_a represents the number of accidents per work unit.

Finally, each score is normalized to a scale between 0 and 1 using the min-max approach, adapted for the present study in Eq. (5):

$$I = \frac{\bar{I} - I_{min}}{I_{max} - I_{min}} \quad (5)$$

The calculation of impact scores for the **Action Plans** and **Preventive Safety Actions** databases was carried out by relating the number of registered actions to the number of accidents that occurred in the same period, using Eq. (4) and, subsequently, normalizing the values based on Eq. (5). This process resulted in three I values: one referring to the **Action Plans** database and two associated with **Preventive Safety Actions**, considering the difference between the progress necessary for the plans to be completed and resolved.

Using the **Accident Records** database, eight I values related to location, time, age, and seniority were created. The location score is calculated for the store and the section by counting the number of accidents. Furthermore, the time score is calculated to estimate the number of accidents that occurred for each value of an hour, day of the week, month, and week of the year. Also, through this database, it is possible to obtain the impact score about the age and seniority of the employee who suffered the accident.

The **Audits** database stores information from periodic visits to its stores to verify the correct application of safety standards, recording in this database when deviations are detected. These issues are divided into four categories: compliant, non-compliant, non-observable, and non-applicable issues. In this context, problems are grouped by store to create audit notes, recording the number of issues for each category and causing four I values.

Employee absenteeism is recorded by store, year, and month and classified into nine categories: medical leave, parental leave, COVID-19, complaints, and unjustified absences, among others. To calculate the impact score based on this information, the **Absenteeism** database was integrated with the **Accident Records** database, and the resulting data was grouped by type of absenteeism, accounting for the number of accidents associated with each category. This process resulted in seven I values corresponding to each type of absenteeism identified in the company.

Similarly, the impact score was calculated for the **Demographic Data** database. In this case, the number of accidents related to different variables was determined, such as the workers' academic qualifications (elementary or secondary education), hours worked, number of employees, and workload category (full-time or part-time). This procedure generated six distinct I values corresponding to each variable.

The **Training** database contains information on training programs divided into specific and transversal for three risk levels (1 to 3). This database was merged with the **Accident Records** database to account for the number of accidents that occurred for each percentage of employees who completed this training, similar to the strategy used for absenteeism scores and demographics, resulting in seven impact scores, three for each type of program and one for total information.

Finally, for the **Transactions** database, it is necessary to use the algorithm described in Vaz et al. (2024), since daily store transaction data takes a few days to become available in the company's internal storage. Therefore, a transaction forecasting algorithm was necessary to support the predictive model for workplace accidents. After testing the study procedure (Vaz et al., 2024) for all of the company's stores, the Gradient Boosting Machine (GBM) achieves the best results for most stores. Thus, to obtain the I value for the **Transactions** database, the accident records are crossed with this information from all stores, in which the transaction data is grouped by value range every 500 transactions. From there, accident data is grouped by store and transaction range, counting the number of accidents.

To summarize the described information, Fig. 1 represents the methodology to develop the dataset based on the calculation of impact scores.

Fig. 1 reveals that this procedure generates 40 impact scores from 369 pre-processed information variables from 10 databases. The impact scores were combined into a single dataset using the common factors, store designation, and section ID.

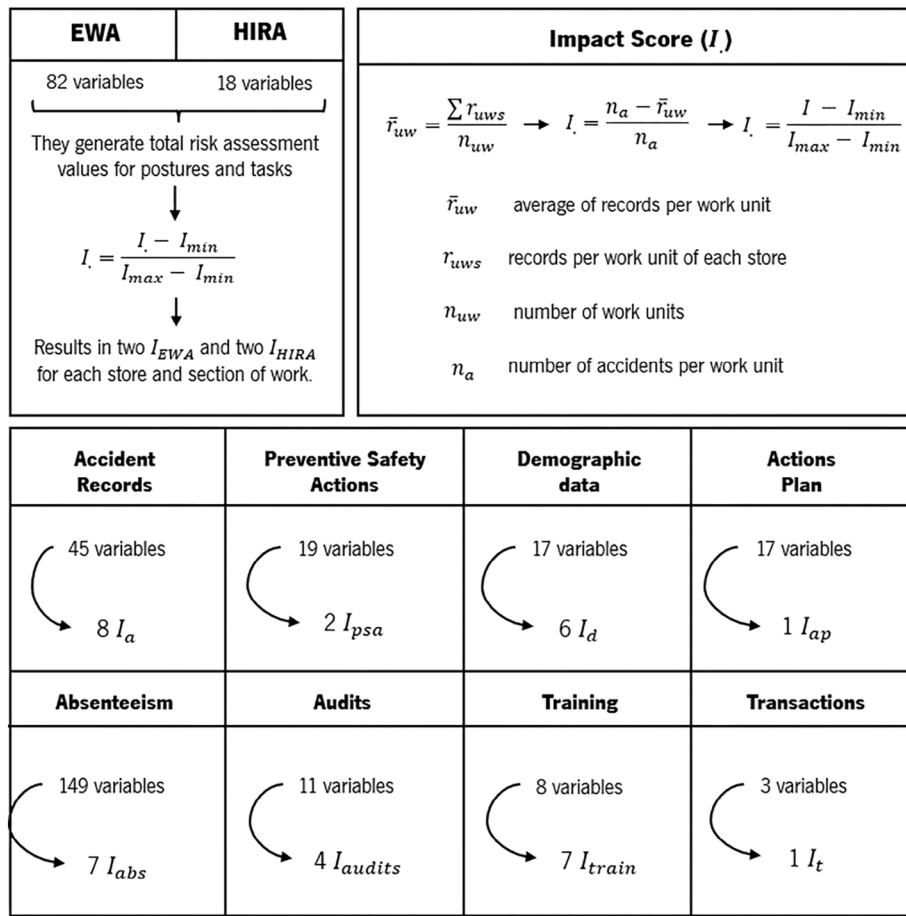


Fig. 1. Methodology used to develop the impact score dataset.

3.3. Predictive model

With the created dataset, the main structure of the predictive model was developed based on two fundamental components: the PodManager and the POD. PodManager acts as a central manager, allowing the integration of multiple PODs, which are individual modeling units. Each POD contains one or more Machine Learning (ML) models that contribute to the prediction values.

The impact score dataset was configured in a POD based on the Multilayer Perceptron (MLP) algorithm, tuned with specific parameters such as the maximum number of steps. After, the prediction is performed with the PodManager using the determined PODs from the dataset of impact scores to calculate the probability of risk situations for each combination of store, section, date, and shift analyzed. The predictive model structure can be seen in Fig. 2.

The proposed predictive model is a powerful tool for predicting risk situations. It integrates advanced ML techniques with a robust data processing and enrichment pipeline. This model enables a proactive approach to workplace risk management, identifying and mitigating potential issues before they occur.

4. Results and discussion

The proposed predictive model was implemented in an integrated system, which allows viewing and sending an alert message about the probability of risk situations occurring in a store/area/section of the company, to three types of users:

- **System administrators:** view and receive the stores' risk.

Table 3

Scale of risk level.

Risk level	Range for risk level (%)
Slight risk (r_s)	$0 \leq r_s < 50$
Moderate risk (r_m)	$50 \leq r_m < 75$
High risk (r_h)	$75 \leq r_h \leq 100$

- **Occupational Health and Safety (OSH) team technicians:** only observe and receive risk alerts from responsible stores.
- **Area/section managers and store directors:** have access to information about their store.

The predictive model was implemented in two steps. The first step, between July 2023 and January 2024, used the probability of risk accidents for the company's eight stores as a pilot study. From September 2023, the remaining 306 company stores were gradually included. Furthermore, since new impact scores were obtained, the algorithm was updated. The second step occurred between February and October 2024, considering all the impact scores and company stores.

The databases are updated daily through an Application Programming Interface (API) created by the company. Due to this information update, the algorithm is activated every day at midnight to process the necessary calculations, lasting approximately ten minutes.

For the first step, the predictive model output was according to a probability scale of risk situations occurrence between 0 and 100%, as shown in Table 3.

According to Table 3, when the predictive model indicated a probability greater than or equal to 75%, the platform sent a high-risk alert

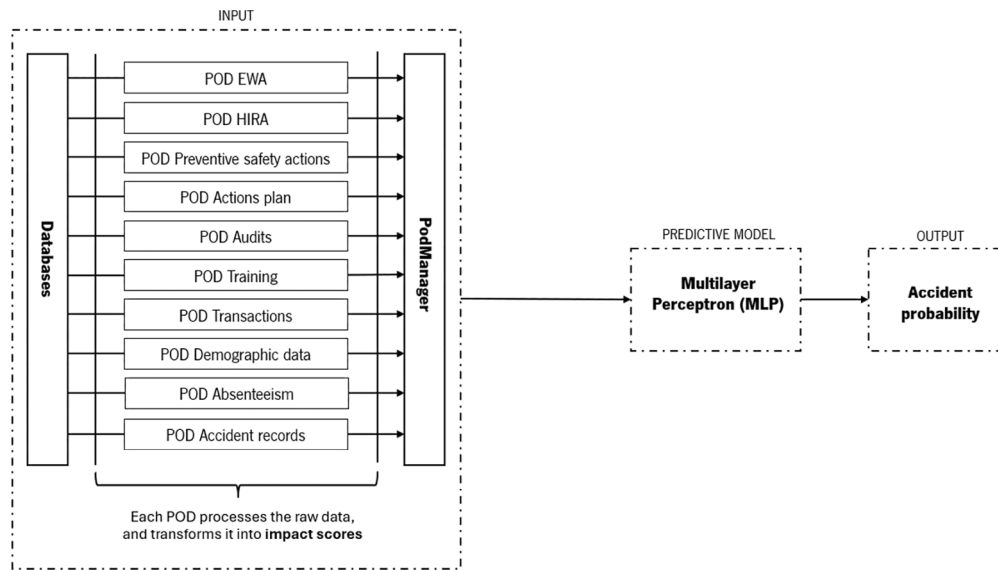


Fig. 2. Methodology for the operation of the predictive model.

Table 4

Metric results for the first algorithm implementation in the intelligent system.

TP	FP	FN	TN	Accuracy
76	15975	295	544086	97.10

Table 5

Scale for risk difference.

Risk level	Range for risk level (%)
Slight risk (dr_s)	$0 \leq dr_s < \epsilon_1$
Moderate risk (dr_m)	$\epsilon_1 \leq dr_m < \epsilon_2$
High risk (dr_h)	$\epsilon_2 \leq dr_h \leq 100$

(r_h) if a value between 50% and 75% was obtained, it sent an alert of moderate risk (r_m), to be aware of any situation, and in the case of values below 50%, no warnings were sent.

To evaluate the models' performance, metrics based on the confusion matrix and the Accuracy value achieved during this period were calculated; these values can be seen in Table 4.

During this period, 371 risk situations occurred, and the predictive model correctly predicted 20% of the events. The high Accuracy value achieved is shown in Table 4. However, the model is not very assertive in predicting risk situations. This can be justified by the imbalance between the number of records of non-risk situations and those of risk situations during this period.

Furthermore, a detailed data analysis revealed that certain stores rarely exhibited high probability values, yet risk situations still occurred. Data mining research identified that this discrepancy resulted from the larger volume of information available for some stores compared to others. Then, stabilizing a fixed alert risk scale is not an adequate solution.

Therefore, the second step considers the difference between the probability of risk situations occurring on the current day and the probability of risk situations occurring on the previous day, as shown in Table 5.

This risk calculation method was implemented for all of the company's stores between February and October 2024, where $\epsilon_1 = 2$, and $\epsilon_2 = 5$. Table 6 presents the metrics' values for performance analysis.

Table 6 shows a decrease in the final Accuracy of the predictive model to 92%, but there is an increase in assertiveness in identifying risk situations. The predictive model in 1255 risk situations manages to identify 29% correctly. The value of FP, corresponding to 8%,

Table 6

Metric results for the second algorithm implementation in the intelligent system.

TP	FP	FN	TN	Accuracy
358	453629	897	5011068	91.68

is notable. This happens due to the high record calculation for each combination of store, section, date, and work shift, which mainly results in non-risk situations and generates a high class imbalance.

The predictive model presents promising results, although there are still opportunities for improvement. Its main advantage is the easy adaptation to different organizational scenarios and the ability to generate different forecasts, allowing for more effective and specific preventive actions. The integration of multiple data sources contributes to greater forecast Accuracy. Additionally, automation of training and forecasting processes reduces the need for manual intervention, minimizing errors and increasing system efficiency.

In addition to the performance analysis, feedback was obtained from regular users of the application, including a store director, an Occupational Health and Safety (OSH) member, and a section manager. All participants agreed that unsafe behaviors are now analyzed more thoroughly and have become a daily topic of discussion between managers and employees. Overall, the application is considered intuitive, easy to use, and effective in highlighting critical areas within the store.

In summary, the proposed predictive model predicts the probability of risk situations using a Multilayer Perceptron (MLP) algorithm combined with a dataset of 40 impact scores. The intelligent system sends alerts to the main managers of each store with an Accuracy of 92% and a Precision of 29% in predicting risk situations through the predictive model, which presents promising results with opportunities for improvement. These results reveal that the platform provides an excellent solution to support companies and OSH teams in preventing and minimizing workplace accidents, consequently increasing safety and well-being.

From the literature review, the predictive models showing best results are RF, with an accuracy of 92%, GBM with 88%, and ANN with 84%, and using datasets based on accident records. In comparison, the predictive model herein presented uses MLP algorithm and a dataset of 40 impact values calculated using 10 databases. Furthermore, it stands out for predicting the probability of risk situations with an accuracy of 92%, covering two days and different working shifts for the company.

These results demonstrate that the proposed predictive model is more robust, and promising than the approaches presented in the literature, highlighting its potential for practical applications in the sector.

5. Conclusions

This study presents a predictive model to forecast risk situations in the retail sector for two consecutive days, the current and the next day, across four work shifts.

Based on the Multilayer Perceptron (MLP) algorithm, the predictive model was fed by an integrated dataset from ten databases provided by the company under study, using impact scores calculated based on the average number of information records and accidents. The model achieved 92% *Accuracy* and an 8% false positive (*FP*) rate. However, *Accuracy* in predicting risk situations was limited to 29%, reflecting the impact of data imbalance and the difficulty in predicting behavior events.

The results demonstrated the predictive model's potential as an innovative tool for minimizing and preventing occupational risks in organizations, highlighting its adaptability to different contexts.

The study faced significant challenges but provided promising insights into using Machine Learning (ML) algorithms to predict the probability of risk situations.

For future work, adjustments to the model are expected to be explored, including strategies to deal with data imbalance and increase the accuracy of predictions. Additional improvements may include updating the algorithm's output format, providing more detailed and interpretable information to users, identifying the most influential variables for increased risk, allowing for a more targeted prevention approach, and expanding the model to predict not only the risk situations but also different types of events, considering the most susceptible contexts. Furthermore, identifying data on employees' stress levels after and during their work shift is important to increase the identification of accidents and prevent such events due to human error.

These improvements could expand the model's practical applicability, strengthening its ability to support workplace safety management and contributing to safer and more efficient working environments.

CRedit authorship contribution statement

Inês Sena: Resources, Investigation, Conceptualization, Writing – original draft, Methodology, Data curation. **Felipe G. Silva:** Methodology, Data curation, Resources, Formal analysis. **Ana Cristina Braga:** Visualization, Supervision, Validation. **Florbela P. Fernandes:** Writing – review & editing, Validation, Visualization, Supervision. **Clara B. Vaz:** Writing – review & editing, Validation, Visualization, Supervision. **Maria F. Pacheco:** Writing – review & editing, Validation, Visualization, Supervision. **Paulo Novais:** Validation, Visualization, Supervision. **José Lima:** Visualization, Supervision, Writing – review & editing, Validation. **Ana I. Pereira:** Writing – review & editing, Validation, Project administration, Visualization, Supervision.

Declaration of competing interest

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.

Acknowledgments

The authors are grateful to the Foundation for Science and Technology (FCT, Portugal) for financial support through national funds FCT/MCTES (PIDDAC) to CeDRI UID/05757 (DOI: 10.54499/UIDB/05757/2020) and SusTEC LA/P/0007/2021 (DOI: 10.54499/LA/P/0007/2020). This work has been supported by NORTE-01-0247-FEDER-072598 iSafety: Intelligent system for occupational safety and well-being in the retail sector. Inês Sena was supported by FCT, Portugal PhD grant UI/BD/153348/2022. Also, thanks to the Mountains Research Collaborative Laboratory (MORE CoLAB) for

letting us test the algorithm in the intelligent system iSafety developed for them.

References

- Ajayi, A., Oyedele, L., Akinade, O., Bilal, M., Owolabi, H., Akanbi, L., Delgado, J.M.D., 2020. Optimised big data analytics for health and safety hazards prediction in power infrastructure operations. *Saf. Sci.* 125, 104656.
- Ayhan, B.U., Tokdemir, O.B., 2019. Predicting the outcome of construction incidents. *Saf. Sci.* 113, 91–104.
- Badri, A., Boudreau-Trudel, B., Souissi, A.S., 2018. Occupational health and safety in the industry 4.0 era: A cause for major concern? *Saf. Sci.* 109, 403–411.
- Blanchard, D., 2021. Safety technology. <https://www.ehstoday.com/safety-technology/article/21920103/>. (Accessed: 20 December 2021).
- Choi, J., Gu, B., Chin, S., Lee, J.S., 2020. Machine learning predictive model based on national data for fatal accidents of construction workers. *Autom. Constr.* 110, 102974.
- Ghodrati, N., Yiu, T.W., Wilkinson, S., Shahbapour, M., 2018. A new approach to predict safety outcomes in the construction industry. *Saf. Sci.* 109, 86–94.
- Grandini, M., Bagli, E., Visani, G., 2020. Metrics for multi-class classification: an overview. *arXiv preprint arXiv:2008.05756*.
- Hastie, T., Tibshirani, R., Friedman, J.H., Friedman, J.H., 2009. *The Elements of Statistical Learning: Data Mining, Inference, and Prediction*, vol. 2, Springer.
- Kang, K., Ryu, H., 2019. Predicting types of occupational accidents at construction sites in Korea using random forest model. *Saf. Sci.* 120, 226–236.
- Khahro, S.H., Ali, T.H., Memon, N.A., Memon, Z.A., 2020. Occupational accidents. *Current Sci.* 118 (2), 243–248.
- Koc, K., Ekmekcioğlu, Ö., Gurgun, A.P., 2022. Accident prediction in construction using hybrid wavelet-machine learning. *Autom. Constr.* 133, 103987.
- Liz-Domínguez, M., Caeiro-Rodríguez, M., Llamas-Nistal, M., Mikic-Fonte, F.A., 2019. Systematic literature review of predictive analysis tools in higher education. *Appl. Sci.* 9 (24), 5569.
- Mosavi, A., Samadianfard, S., Darbandi, S., Nabipour, N., Qasem, S.N., Salwana, E., Band, S.S., 2021. Predicting soil electrical conductivity using multi-layer perceptron integrated with grey wolf optimizer. *J. Geochem. Explor.* 220, 106639.
- Ohadi, S., Hashemi Monfared, S.A., Azhdary Moghaddam, M., Givchi, M., 2022. Feasibility of a novel predictive model based on multilayer perceptron optimized with harris hawk optimization for estimating of the longitudinal dispersion coefficient in rivers. *Neural Comput. Appl.* (ISSN: 0941-0643) 35 (9), 7081–7105. <http://dx.doi.org/10.1007/s00521-022-08074-8>.
- Oyedele, A., Ajayi, A., Oyedele, L.O., Delgado, J.M.D., Akanbi, L., Akinade, O., Owolabi, H., Bilal, M., 2021. Deep learning and boosted trees for injuries prediction in power infrastructure projects. *Appl. Soft Comput.* 110, 107587.
- Poh, C.Q., Ubeynarayana, C.U., Goh, Y.M., 2018. Safety leading indicators for construction sites: A machine learning approach. *Autom. Constr.* 93, 375–386.
- PORDATA, 2024. Acidentes de trabalho: total e por setor de atividade económica. <https://www.pordata.pt/portugal>. (Accessed: 23 March 2024).
- Qteat, H., Awad, M., 2021. Using hybrid model of particle swarm optimization and multi-layer perceptron neural networks for classification of diabetes. *Int. J. Intell. Eng. Syst.* 14 (3).
- Recal, F., Demirel, T., 2021. Comparison of machine learning methods in predicting binary and multi-class occupational accident severity. *J. Intell. Fuzzy Systems* 40 (6), 10981–10998.
- Santibáñez, F., Flores, C., Basso, F., Jiménez, A., Bravo, F., Núñez, F., Luco, H., Martínez, L., et al., 2013. Mining accident detection using machine learning methods. *IFAC Proc. Vol.* 46 (16), 31–33.
- Sarkar, S., Ejaz, N., Promod, C., Maiti, J., 2020. Pattern extraction using proactive and reactive data: A case study of contractors' safety in a steel plant. In: *Proceedings of ICETIT 2019: Emerging Trends in Information Technology*. Springer, pp. 731–742.
- Sarkar, S., Pramanik, A., Maiti, J., 2023. An integrated approach using rough set theory, ANFIS, and Z-number in occupational risk prediction. *Eng. Appl. Artif. Intell.* 117, 105515.
- Sarkar, S., Vinay, S., Djeddi, C., Maiti, J., 2022. Classification and pattern extraction of incidents: A deep learning-based approach. *Neural Comput. Appl.* 34 (17), 14253–14274.
- Sena, I., Braga, A.C., Novais, P., Fernandes, F.P., Pacheco, M.F., Vaz, C.B., Lima, J., Pereira, A.I., 2023. Exploring features to classify occupational accidents in the retail sector. In: *International Conference on Optimization, Learning Algorithms and Applications*. Springer, pp. 49–62.
- Silva, F.G., Sena, I., Lima, L.A., Fernandes, F.P., Pacheco, M.F., Vaz, C.B., Lima, J., Pereira, A.I., 2023. Data pruning approach in the retail sector. In: *Symposium of Applied Science for Young Researchers. SASYSR 2023*, <http://hdl.handle.net/10198/28842>.
- Vaz, C.B., Sena, I., Braga, A.C., Novais, P., Fernandes, F.P., Lima, J., Pereira, A.I., 2024. Predicting retail store transaction patterns: a comparison of ARIMA and machine learning models. In: *International Conference on Computational Science and Its Applications*, vol. 2280, Springer Nature Switzerland, pp. 268–283.