

RESEARCH ARTICLE

Data-Driven Bayesian Network Analysis of Railway Accident Risk

LEI SHI¹, YAZHI LIU¹, YOUPENG ZHANG^{1,2}, AND JUNYI LIANG¹¹School of Automation and Electrical Engineering, Lanzhou Jiaotong University, Lanzhou 730070, China²Rail Transit Electrical Automation Engineering Laboratory of Gansu Province, Lanzhou Jiaotong University, Lanzhou 730070, China

Corresponding author: Yazhi Liu (liuyazhiL@163.com)

This work was supported by the Key Program of Natural Science Foundation of Gansu Province, China, under Grant 21JR7RA292.

ABSTRACT Ensuring railway safety is a top priority, with a central focus on preventing accidents. By thoroughly analyzing data from railway accident investigations, we can pinpoint factors and patterns associated with different types of railway accidents. This proactive approach not only helps reduce the frequency of such incidents but also significantly boosts overall railway transportation safety. This paper investigates the impact of various risk factors on railway safety through the analysis of railway accidents by using data-driven Bayesian networks. First, key data representing the frequency of risk factors directly derived from railway accident reports are collected and analyzed. Then, the risk factors are incorporated into causal analysis for different types of railway accidents. Finally, a historical data-driven approach is utilized to model and gain new insights into the key risk factors causing different types of railway accidents. Meanwhile, a Tree-Augmented Naive Bayes (TAN) is employed to construct a model of interdependencies among risk factors, and the model is validated through sensitivity analysis and past accident records. The research findings demonstrate that the crucial risk factors for all types of accidents include undetected track damage, train operator skills, load, braking system conditions, train speed, traction system failures, level crossings, and bridge damage. Additionally, the research results highlight the differential impact of key factors on different types of accidents, providing a most probable explanation for observing the most likely configurations in the model for a specific scenario. This work contributes to accident prevention and safety decision-making.

INDEX TERMS Railway accidents, data-driven, Bayesian networks, risk analysis, tree-augmented naive Bayes.

I. INTRODUCTION

Most railway accidents, such as derailments, collisions, and fires, have characteristics of low probability and high consequences. Catastrophic railway accidents can cause great loss of life, social impact, and environmental damage [1]. The occurrence of railway accidents is attributed to both inevitable and accidental factors. To achieve long-term prevention of railway accidents, it is necessary to analyze and control various factors that lead to accidents and evaluate the probability of risks. Therefore, reducing the risk of railway accidents and guaranteeing safe, stable, and reliable train operations have always been a major concern.

The associate editor coordinating the review of this manuscript and approving it for publication was Dongxiao Yu¹.

The models and methods used for accident analysis can explain and even predict the mechanisms behind the occurrence of accidents, and implementation of effective countermeasures [2]. Risk analysis provides an effective approach to preventing railway accidents. In the research on risk analysis of railway accidents, historical data analysis methods are widely used [3], [4]. For instance, Liu et al. [5] applied knowledge graph theory to analyze railway accidents and investigated 214 railway accidents in the UK. Using railway accident data from the United States from 2000 to 2016, Zhang et al. [6] gained new insights into safety assessment and improvement through statistical risk analysis. Savage [7] analyzed the distribution of train-pedestrian fatal collision accidents in Chicago from 2004 to 2012 from a temporal and spatial aspect. Zhou and Lei [8] investigated the frequency

of 407 railway accidents/incidents in China based on accident types, and they revealed how actions and decisions at higher organizational and managerial levels in the railway system can cause accidents/incidents errors. The analysis of historical data is susceptible to sample selection bias, which may result in inaccurate estimates of accident occurrence frequency and patterns.

In the field of railways, quantitative safety risk analysis methods have been widely used, including Failure Mode and Effects Analysis (FMEA), Hazard and Operability Studies (HAZOP), Fault Tree Analysis (FTA), Event Tree Analysis (ETA), Bayesian Networks (BN), etc. For example, Xue and Yang [9] adopted a functional-based FMEA approach to investigate the impact of various interface failures between the railway signaling system and the platform door system on train operations. Based on this, clear safety requirements were formulated for the signaling system and safety analysis conclusions were provided for the platform door system. Bian and Li [10] identified risk factors in railway hazardous goods transportation processes and provided references and suggestions for managing railway hazardous goods transportation safety through HAZOP analysis. Liu et al. [11] utilized the FTA method to analyze operational accidents on the Yong-Wen railway in China. Ni and Tang [12] combined FTA with the Fuzzy Analytic Hierarchy Process (FAHP) to assess the risk of subway fires and then they proposed effective measures for preventing and controlling the occurrence and spread of subway fires. Li et al. [13] employed the event tree method to analyze the risk of trains accidentally entering work areas. Their study provided theoretical and technical support for optimizing on-site work protection management and preventing personal accidents. Baysari et al. [14] investigated 40 railway operational accidents in Australia using the Human Factors Analysis and Classification System (HFACS), focusing on human factors analysis and classification. Liang et al. [15] developed a general approach of Causal Statistic Risk Assessment based on hierarchical Causal Bayesian Networks (CSRA-CBN) to analyze the various impacting factors which may cause accidents, and identify the factors which contribute most to the accidents at Level Crossing (LX), thus allowing for risk quantification. Liang et al. [16] forward and reverse inferences based on the BN risk model are performed to predict LX accident occurrence and quantify the contribution degree of various impacting factors respectively, so as to identify the riskiest factors. Quantitative safety risk analysis methods typically rely on data, and inaccuracies or a lack of representative data may lead to inaccuracies in the analysis. Additionally, quantitative safety risk analysis methods may tend to focus on specific aspects, potentially overlooking the comprehensiveness of the entire system.

Weber et al. [17] pointed out that the number of publications on Bayesian networks in the field of risk analysis has been increasing every year because of its superiority in learning structure and inference algorithms. Compared to

other classical reliability analysis methods, BN maintains its advantage by establishing a multi-state variable model. For instance, BN has similar characteristics to Fault Trees (FT), and it is not only suitable for two-state variables but also can model multiple-state variables and multiple outputs. To overcome the limitations of FT in terms of static structure and uncertainty, FT can be mapped into BN [18], [19]. However, as the number of variables increases, the parameters and related functions will increase substantially, thus increasing system modeling complexity [17]. For example, Markov Chains (MC) utilize differences between variables to analyze the probabilities of failure events. MC can represent multiple-state variables, but as the number of variables increases, the system becomes more complex. Meanwhile, BN modeling requires relatively fewer parameters and smaller conditional probability tables. Besides, BN has become a popular method for railway risk modeling because of its ability to utilize expert knowledge or data-driven approaches. Huang et al. [20] presented a data-driven approach called Bayesian Network-K2 algorithm-Expectation Maximization (BN-K2-EM). They processed accident reports as a fault data matrix and employed the K2 algorithm and expectation-maximization algorithm to obtain the structure and parameters of the constructed Bayesian network. Then, the constructed Bayesian network was applied to predict and diagnose operational failures in high-speed trains. Li and Qi [21] analyzed the failure modes in the pantograph system using expert knowledge and historical data, and they constructed a fuzzy fault tree and transformed it into a BN model. When there is insufficient data from relevant accidents, expert knowledge remains important for railway accident modeling. Liang et al. [22] proposed a causal reasoning framework for risk analysis based on BN. This framework combines empirical knowledge with automatic learning methods by introducing causal structure constraints. The goal is to identify effective causal relationships while avoiding inappropriate structural connections. The researchers applied this framework to the risk analysis of the LX accidents in France.

The TAN model introduces a tree structure that allows for conditional dependencies among risk factors while maintaining a relatively low model complexity. This feature enables TAN to flexibly capture relationships between risk factors. The tree structure in the TAN model is utilized to represent conditional dependencies among risk factors, where each feature node, given a class node, is connected to other feature nodes. This structure provides a more accurate modeling of feature dependencies. In comparison to traditional naive bayes, TAN demonstrates better adaptability to the distribution of real-world data.

Compared to studies using expert knowledge in BN modeling, data-driven BN involves fewer subjective biases, and the analysis results have a certain degree of objectivity. However, data-driven BN models require collecting more empirical evidence before widespread practical applications. To bridge

this gap, this study employs newly collected raw data from railway accident reports to establish a data-driven BN model with the risk impact factor (RIF) structure. Based on this, this paper provides novel insights into distinguishing key risk factors for different types of railway accidents.

By comprehensively summarizing the current state of domestic and international research on railway accident analysis, it has been found that various data analysis approaches are used to identify relatively independent risk factors from a large volume of records. As for the modeling of accident processes, both domestic and international studies adopt specific accident analysis models to analyze the occurrence process of a particular incident or a specific class of accidents. The effectiveness of the models is mainly validated by accident statistics. However, there is a scarcity of research in railway accident analysis that focuses on modeling and exploring the relationships between different risk factors based on accident statistical data.

Therefore, this paper conducts a risk assessment of potential risk factors in railway accidents by establishing a TAN model. Innovations in this study include: (1) collecting and analyzing primary data representing the frequency of risk factors directly derived from railway accident reports; (2) incorporating each risk factor into causal analyses of different types of railway accidents; (3) employing a data-driven approach based on historical data to model and provide new insights into the key risk factors leading to different types of railway accidents. Addressing the issue of subjective judgment in risk analysis due to a lack of objective quantification, this paper adopts a data-driven approach to induce the Bayesian network structure from data, effectively mitigating subjective judgments of experts. While previous risk analyses often focused on the severity of accidents, this paper sets the target variable of the model as the accident type, predicting the potential types of accidents. To address the problem of local optima in optimizing Bayesian networks, this study utilizes the TAN method to optimize the naive Bayesian network, relaxing the assumption of attribute independence through a tree-like structure.

This paper explores how individual or combined risk factors influence different types of railway accidents. Based on the railway accident reports recorded by the RAIB in the United Kingdom from 2011 to 2020, a preliminary database was constructed. Then, by utilizing the accident data to construct an Augmented Naive Bayes model, a data-driven Bayesian network railway accident analysis approach is proposed. The modeling process of the TAN model is illustrated in Fig. 1.

The remainder of this article is structured as follows. Section II elaborates on the methods employed for the identification of RIFs, TAN-BN structure learning, and sensitivity analysis. In Section III, an analysis is conducted on the impact of various risk factors on different types of accidents, elucidating the combinations in which risk influencing factors manifest and providing reasoned interpretations for the observed outcomes. Finally, Section IV

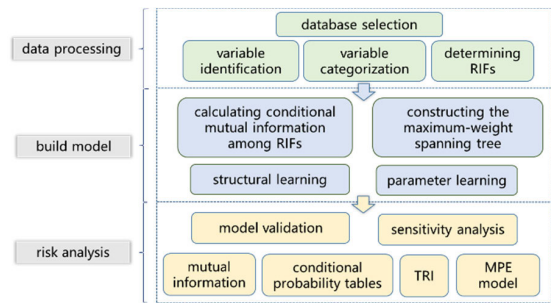


FIGURE 1. The process of establishing the TAN model.

provides a comprehensive summary and outlook for the entire paper.

II. METHODOLOGY

A. IDENTIFICATION OF RIFs

To analyze different types of railway accidents under various risk impact factors, it is necessary to identify and select the RIFs from accident reports. The data used in this study was obtained through a systematic analysis of publicly available railway accident cases from relevant transport organizations and the RAIB. The accident investigation reports released by the RAIB provide detailed investigative processes for significant accidents and incidents that have a great impact on the railway system. These reports provide comprehensive and clear material for researchers to reconstruct the accident scene and understand the process of accident occurrence.

The process of generating RIFs consists of four stages: (1) online database search, (2) reading and screening of accident reports, (3) extraction and analysis of accident report content, and (4) selection of RIFs. In this study, railway accident reports released by the Rail RAIB from January 2011 to January 2020 (www.gov.uk/raib) were read and screened, upon analyzing accident investigation reports, it was observed that certain accidents occurred due to violations of traffic rules by pedestrians, passengers, or car drivers, leading to collisions with trains. As accidents resulting from such factors do not align with the research objectives and may potentially compromise the accuracy of the analysis results, they were excluded from consideration. Additionally, in some accident reports, a single report documented multiple incidents. This paper treats such reports as multiple incidents for the purpose of analysis, and finally, 121 accident investigation reports were selected as the study sample. Further refinement and analysis of the reports were conducted, and through manual analysis of the original railway accident reports, 43 risk factors were first generated. Then, domain experts were invited to merge and transform risk factors with high similarity. Meanwhile, other factors exhibited certain interdependencies, and notable differences were observed, so they were retained as separate factors. Finally, a total of 21 risk impact factors that cause railway accidents were identified, as listed in Table 1.

Moreover, this study incorporates other external factors such as the location of railway accidents, occurrence time,

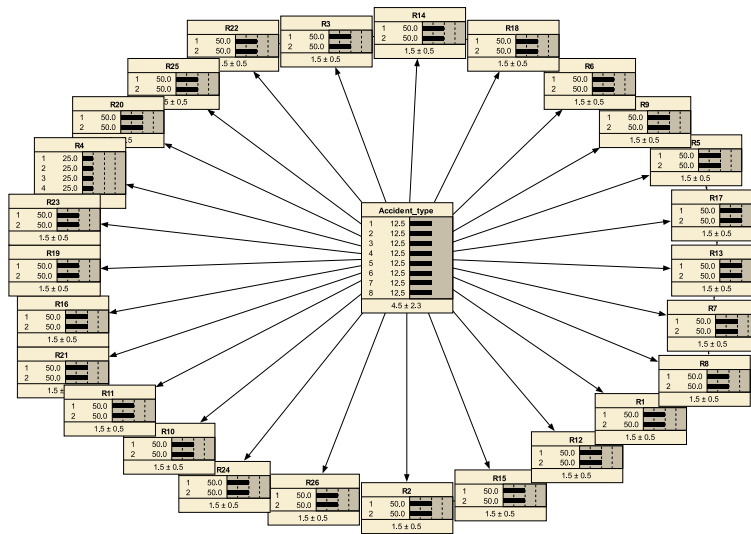


FIGURE 2. TAN-BN structure learning.

TABLE 1. RIFs contributing to railway accidents.

Number	Risk Factor	Frequency
1	Railway accidents occurring at level crossings	18.18%
2	Track misalignment (caused by long-term use of railway tracks, resulting in horizontal displacement and deviation of gauge beyond the safety limits)	1.67%
3	Turnout failure (infrastructure failure due to prolonged usage, causing improper functioning of switching devices)	0.83%
4	Tractive system failure (transformer malfunctions and auxiliary power equipment failures, etc.)	1.64%
5	Braking system failure (power supply device failure, control device failure, transmission device failure, brake device failure, etc.)	4.13%
6	Communication and signaling equipment failure (signal apparatus failure, station interlocking control system failure, etc.)	3.33%
7	Traction system failure (train wheel failure, wheel wear and tear)	1.65%
8	Failure of train-related components (detachment of train components, bogie failure, coupling component failure, train window damage, train door system failure, door detachment, etc.)	4.13%
9	Bridge damage (risk to railway train lines caused by bridge exceeding its service life or carrying capacity)	1.63%
10	Adverse weather conditions (heavy rain, strong winds, thick fog, icy and snowy environment, etc.)	16.53%
11	Geological hazards (floods, landslides, collapses, mudslides, ground subsidence, earthquakes, etc.)	3.27%
12	Inadequate staff competency management	17.35%
13	Dam damage	7.44%
14	Information communication failure (errors in ordering or transmitting information due to railway staff mistakes)	3.30%
15	Lack of job skills (inadequate training of railway staff's job skills)	24.79%
16	Fatigue driving by train operators (insufficient work capacity due to distractions, panic, etc., causing accidents)	8.26%
17	Inadequate regulations and systems (incomplete establishment of railway management system rules and regulations, inadequate staff training system, imperfect safety assessment system)	11.57%
18	Errors in train operation command (inappropriate personnel allocation, insufficient personnel reserve)	0.82%
19	Inadequate safety supervision (weak safety awareness, insufficient safety training)	22.31%
20	Undetected track damage	12.39%
21	Poor rail adhesion conditions	3.27%

train speed, and train load as objective risk impact factors. The definitions of variable states are mostly extracted from

the accident reports, with reference to and consolidation of variable state classifications in previous relevant studies. For

TABLE 2. 26 RIFs defined in railway accidents.

No	RIFs	Notation	Type	Value
1	Level Crossing	R1	Crossing, Non-crossing	1, 2
2	Occurrence Time	R2	6:00 to 18:00, 18:00 to the next day 6:00	1, 2
3	Train Speed	R3	Normal, Overspeed	1, 2
4	Load	R4	Normal, Full load, Empty load, Overload	1, 2, 3, 4
5	Weather Condition	R5	Normal, Rain, Snow, Fog, etc.	1, 2
6	Train Driver Skill	R6	Good skill, Insufficient skill	1, 2
7	Undetected Track Damages	R7	No track damages, Decayed wooden track	1, 2
8	Work Skills	R8	Good, Insufficient repair and signaling skills, Failure to identify and detect hazards timely	1, 2
9	Train Driver Fatigue	R9	Good, Distracted driving	1, 2
10	Brake System Condition	R10	Good, Brake mechanism failure	1, 2
11	Running Gear Condition	R11	Good, Wheel wear	1, 2
12	Inadequate Safety Supervision	R12	Good, Insufficient safety education for railway employees	1, 2
13	Dam Damages	R13	Good, Subsidence and slope failure	1, 2
14	Communication Failure	R14	Good, Improper decision-making	1, 2
15	Inadequate Staff Capability Management	R15	Good, Insufficient capability management for railway employees such as drivers, inspectors, maintenance workers, and machine operators	1, 2
16	Signal Equipment Failure	R16	Good, Signal equipment failure such as track circuit faults, line signal faults, alarm system faults	1, 2
17	Geological Disasters	R17	Good, Landslide	1, 2
18	Inadequate Regulations	R18	Good, Inadequate regulations	1, 2
19	Mistakes in Train Operation	R19	Good, Improper personnel deployment	1, 2
20	Failure of Vehicle Components	R20	Good, Bogie failure	1, 2
21	Traction System Failure	R21	Good, Traction equipment box failure	1, 2
22	Foreign Object Intrusion	R22	None, Tree leaves	1, 2
23	Track Deviation	R23	None, Lateral track deviation	1, 2
24	Bridge Damage	R24	None, Bridge collapse	1, 2
25	Rail Adhesion Conditions	R25	Good, Low rail adhesion conditions caused by fallen leaves, ice, rainwater, sand, etc.	1, 2
26	Switch Failure	R26	Good, Switch mechanical failure	1, 2

instance, the states of “accident location” in the accident reports are categorized as “level crossing” and “non-level crossing”. Finally, 26 risk impact factors and their categorized states for railway accidents were defined, as listed in Table 2.

B. TAN-BN STRUCTURE LEARNING

The concept of BN was proposed by Pearl in 1988 [23]. It is a graphical network based on probabilistic problems, composed of a directed acyclic graph (DAG) and conditional probability distributions (CPDs). The DAG visually represents the specific structure of the problem, and it consists of nodes and directed edges. The CPDs quantitatively represent the intrinsic relationships between the nodes. The directed edges between nodes represent the direct relationships among variables, so they determine the statistical correlations and conditional probability tables (CPTs) between nodes. By using a large amount of statistical information, this paper constructs CPTs to enable the BN model to have good capabilities for stochastic modeling and handling nonlinear relationships. Based on this, inference functionality is achieved under incomplete, imprecise, and uncertain information.

There are several data-driven approaches for BN modeling, such as the Naive Bayesian Network (NBN), the Augmented Naive Bayesian Network (ABN), and TAN. TAN is constructed based on the NBN by adding directed arcs between

attributes with strong dependency relationships, but there are some limitations on the connections between each attribute. This leads to a tree-like graphical model that represents the dependencies between attributes. Compared to the structure of an NBN, the TAN structure can fully utilize the dependency relationships among attribute variables. It improves upon the NBN by extending the structure, thereby preserving the learning capabilities of BN while reducing the complexity of BN [24], [25], [26]. Therefore, this study establishes a TAN-BN model to analyze the RIFs of various types of railway accidents.

A BN is used to encode the joint probability distribution of a set of random variables, denoted as a variable U . The BN is represented by DAG with annotations. Let $U = \{A_1, \dots, A_n, C\}$, where n represents the number of RIFs. Variable A_1, \dots, A_n corresponds to the RIFs, while variable C represents the class variable (accident type). In the TAN model, the structure considers the class variable as the root with no parent nodes, denoted as $\Pi C = \emptyset$ (where ΠC represents the parent set of C within variable U). Each RIF has a class variable as its unique parent node, denoted as variable $\Pi A_i = \{C\}$, $1 \leq i \leq n$. The BN defines the unique joint probability distribution over variable U as:

$$P(A_1, \dots, A_n, C) = P(C) \cdot \prod_{i=1}^n P(A_i|C) \quad (1)$$

For all instances of A_i , except for the class variable C serving as its parent node, at most one other attribute variable can serve as its parent node. Such a DAG is referred to as a tree

on $\{A_1, \dots, A_n\}$. the tree structure on variable A_1, \dots, A_n can be described using function π such that there exists a value i that precisely satisfies $\pi(i) = 0$, and there is no sequence of i_1, \dots, i_k that leads to $\pi(ij) = i_{j+1}, i \leq j \leq k$ and $\pi(i_k) = i_1$ (i.e., no cycle). This function is used to define the tree network, where $\pi(i) > 0, \Pi A_i = \{C, \dots, A_{\pi(i)}\}$; and $\pi(i) = 0, \Pi A_i = \{C\}$.

The essence of TAN learning is an optimization problem, where the TAN structure utilizes the conditional mutual information between attributes as proposed by Chow and Liu [24] in the learning and inference processes. This function is defined as:

$$I_P(A_i, A_j|C) = \sum_{a_{ii}, a_{ji}, c_i} P(a_{ii}, a_{ji}, c_i) \log \frac{P(a_{ii}, a_{ji}|c_i)}{P(a_{ii}|c_i)P(a_{ji}|c_i)}, i \neq j \quad (2)$$

where I_P denotes the conditional mutual information, a_{ii} represents the i th state of attribute variable A_i , a_{ji} represents the i th state of attribute variable A_j , and c_i represents the i th state of the class variable “accident type”. Essentially, the optimization problem of the TAN model involves searching for a tree structure definition function π on A_1, \dots, A_n that maximizes the log-likelihood and then using this function to construct the TAN model as the structural form of the target BN model.

The process of constructing a TAN model mainly consists of the following five steps:

Step 1: Compute the conditional mutual information $I_P(A_i, A_j|C)$, $i \neq j$ between each pair of attribute variables.

Step 2: Construct a complete undirected graph with vertices representing A_1, \dots, A_n , and vertex A_i is connected to A_j with edges weighted by $I_P(A_i, A_j|C)$.

Step 3: Construct a maximum-weight spanning tree.

Step 4: Choose a root variable from the attribute variables and set the direction of all edges to point away from the attribute variables, thereby transforming the obtained undirected tree into a directed tree.

Step 5: Add a class variable node and arcs between the class variable node and the attribute nodes to construct a TAN model.

In Step 3, the process of constructing the maximum-weight spanning tree is as follows: Firstly, sort the edges in descending order of their weights. Then, following the principle that the selected edges should not form cycles, choose the edges in descending order of their weights. In this way, the tree constructed from the selected edges will be the maximum-weight spanning tree.

C. SENSITIVITY ANALYSIS AND MODEL VALIDATION

1) MUTUAL INFORMATION

Mutual information is a measure provided by probability theory to quantify the degree of dependence between two random variables. In this paper, based on the theory of entropy, mutual information is adopted as an indicator of the uncertainty of a dataset. Since this paper aims to determine the relationship between RIFs and specific accident types,

in the calculation of mutual information, “accident type” is chosen as a fixed variable. The mutual information between “accident type” and RIFs is defined as:

$$I(s, \alpha_i) = - \sum_{s,i} P(s, \alpha_{ij}) \log_b \frac{P(s, \alpha_{ij})}{P(s)P(\alpha_{ij})} \quad (3)$$

where, s represents “accident type,” α_i denotes the i th RIF, α_{ij} represents the j th state of the i th RIF, and $I(s, \alpha_i)$ symbolizes the mutual information between “accident type” and the i th RIF. A larger value of mutual information implies a stronger correlation between α_i and “accident type.” Calculating the mutual information values in this approach allows for the selection of risk variables by identifying RIFs with minimal impact. The remaining RIFs, identified as important variables related to the selected accident type in the model, can be extracted to alleviate the subsequent computational workload.

2) SENSITIVITY ANALYSIS

Joint probability refers to the probability that multiple conditions occur simultaneously. In this study, the TAN-BN model is adopted to assign corresponding probability values to different states of the relevant RIFs. The probability distribution of different states of the class variable or target node can be calculated by fixing the states of other RIF variables. The sum of joint distribution probabilities for different states of RIFs is equal to 1. The specific calculation process is shown in Eq. (4), where s represents “accident type” and α_{ij} corresponds to the j th state of the i th RIF.

$$P(s, \alpha_{ij}) = P(s) \cdot P(\alpha_{ij}|s) \quad (4)$$

In risk management and analysis, given that multiple risk influencing factors may concurrently impact the occurrence of a risk, joint probability offers a method to quantify this relationship. If the joint probability of two variables is significantly higher than the probabilities of each being independent events, it can be inferred that there is a certain level of correlation between them.

After RIFs are filtered out through mutual information calculations, another form of sensitivity analysis, called scenario simulation, can be carried out. The traditional method of setting scenarios involves fixing all other nodes and then sequentially updating the states of the target node. This method is suitable for variables with two states but variables with more than two states. In such cases, RIFs with multiple states cannot be used in traditional scenario simulations.

To overcome the limitations of traditional scenario simulation, Chow and Liu [27] proposed a novel sensitivity verification method. It determines the influence of different RIFs in a combinatorial approach. By increasing the probability of the state that has the maximum impact on the target variable to 100%, a High-Risk Inference (HRI) can be calculated for a certain accident type. Similarly, by increasing the probability of the state that has the minimum impact on the target variable to 100%, a Low-Risk Inference (LRI) can be calculated for the same accident type. By calculating

the average of HRI and LRI, the True Risk Inference (TRI) under a specific accident type can be obtained. The specific calculation is defined as:

$$TRI = \frac{HRI + LRI}{2} \quad (5)$$

Therefore, to compare the impact of other relevant variable nodes on the “accident type,” the TRI is calculated for each selected RIF, and then the TRI values for all variables and all accidents are ranked. In this sensitivity analysis method, a higher TRI value indicates a greater influence of the corresponding node on the “accident type.”

3) MODEL VALIDATION

Sensitivity analysis needs to satisfy the following two criteria to demonstrate that the proposed research method is reasonable and logical [28], [29], [30].

Criterion 1: A slight increase or decrease in the prior probability of each parent node will invariably result in an increase or decrease in the posterior probability of the child node.

Criterion 2: The combined effect of probability changes from property x (evidence) on the value is always greater than the total impact of changes from the $x - y (y \in x)$ property set (sub-evidence).

Besides, the effectiveness of the proposed BN model is validated by simulating past railway accidents, and appropriate parameter settings are used to determine if the model can provide results that reflect real-world situations.

D. MOST PROBABLE EXPLANATION

To observe the connections between the nodes in a BN and identify the most likely states of the nodes, the TAN-BN model can be employed to perform re-reasoning for specific types of accidents and provide a reasonable explanation for the observed result. The maximum posteriori probability is a special case and is also known as the Most Probable Explanation (MPE). By setting the class variable state or target node state to the MPE mode, the BN model can observe the most likely occurrences of other nodes under specific types of accidents, known as the most probable RIF states, in this way, the causes of railway accidents can be predicted to some extent.

In the MPE mode, each node has at least a confidence level of 100% for its possible states, indicating the most likely situation. Other confidence levels represent lower probabilities, scaled accordingly. Some nodes may have multiple confidence levels of 100%, indicating that the states corresponding to these confidence levels have equal probabilities of occurrence under a certain type of accident.

III. RESULTS AND DISCUSSION

A. DESCRIPTION OF ACCIDENT TYPES

To generate RIFs in railway accidents, a case analysis was carried out following the procedure introduced in Section II. As listed in Table 3, in the quantitative analysis of BN

TABLE 3. Accident types.

Number	Accident Type
1	Derailment
2	Collision
3	Conflict
4	Explosion
5	Attempted Incident
6	Overrun Point
7	Loss of Control
8	Others

modeling, the accident types were defined as dependent variables, including derailment, collision, conflict, explosion, attempted incident, overrun point, loss of control, etc. Strictly speaking, “attempted incident” is not considered an accident type, and it was included as a category to enrich the accident cases because of its threat to railway safety operations [31].

B. DESCRIPTION OF ACCIDENT TYPES

To construct the BN model, 26 RIFs were selected to illustrate their relationship with the dependent variable (i.e., accident types). The Netica software package (<http://www.norsys.com>) was used to assist in the calculations, and the “Learn Network” function based on Eq. (2) was employed to construct the TAN network. The constructed BN structure is demonstrated in Fig. 2. After importing the data, domain experts thoroughly investigated the constructed TAN-BN network to guarantee that all connections between nodes were meaningful. In this study, since all the interconnections of the data reflected real-world scenarios, no adjustments were made in the fine-tuning process.

Based on the TAN model, the Netica software utilized the Counting Learning Algorithm (https://www.norsys.com/WebHelp/NETICA/XCountingLearning_Algorithm.htm) to perform parameter learning of the CPT in the case. Once the CPT was modeled, posterior probabilities for each variable could be obtained. By analyzing the probabilities of the variables, preliminary conclusions regarding railway accident safety warnings and accident prevention can be made. Based on this, the following insights are gained.

Fig. 3 represents the results obtained on the TAN model for all 26 retained RIFs. Among the accident types, derailment and collision are the most frequent, accounting for 32.6% and 31.0%, respectively. Compared to historical statistical data, the results obtained on the TAN model exhibit a high level of reliability. As shown in Table 4, the predicted probability for the “derailment” accident type differs from the historical data by only 1.28%, and the predicted probability for the “attempted incident” accident type only differs by 0.2%. These minor differences could be attributed to the “other” category, demonstrating the predictive accuracy of the constructed model.

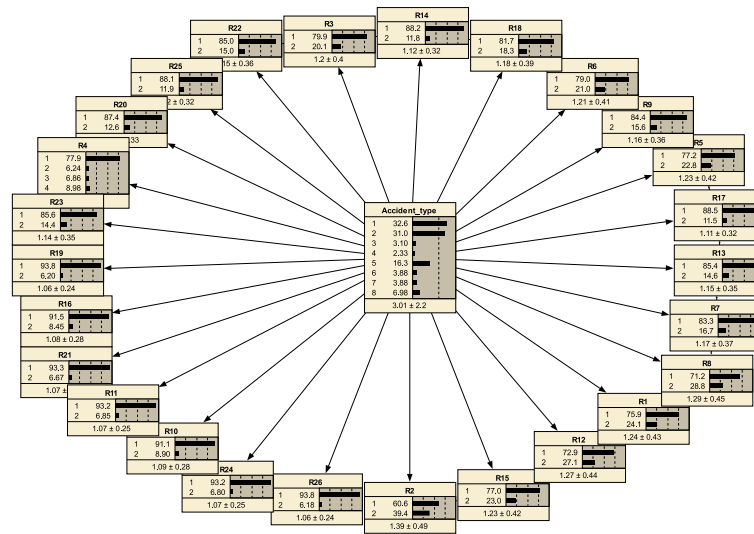


FIGURE 3. TAN-BN parameter learning.

TABLE 4. Comparison of historical data and TAN results.

Accident Type	Historical Data (%)	TAN Results (%)
Derailment	33.88	32.6
Collision	32.23	31.0
Conflict	2.48	3.10
Explosion	1.65	2.33
Attempted Incident	16.5	16.3
Overrun Point	3.30	3.88
Loss of Control	3.30	3.88
Other	6.61	6.98

C. SENSITIVITY ANALYSIS

1) MUTUAL INFORMATION

The mutual information table between “accident type” and RIFs is presented in Table 5. This table shows that risk factor variables with higher $I(s, \alpha_i)$ values are more indicative of the underlying impact on “accident type”. When “accident type” is the parent node, the risk factor variable “Undetected Track Damage” with a mutual information value of 0.10616 has the greatest influence on the accident type. Meanwhile, there are many RIFs with mutual information values smaller than 0.05. In this study, a threshold of 0.05 was chosen for further discussion and analysis. Based on this, eight variables were selected for further factor analysis, including “Undetected Track Damage,” “Train Driver Skill,” “Load,” “Braking System Condition,” “Train Speed,” “Traction System Failure,” “Level Crossing,” and “Bridge Damage”. However, this does not exclude the possibility of considering more factors with smaller mutual information values at an appropriate time. From a methodological perspective, prioritizing the influential individual RIFs using the mutual information-based ranking table is an effective approach.

2) SENSITIVITY ANALYSIS

After the important variables that affect the accident type are selected, the next step is to analyze how these variables (or the states of the variables) affect the target variable, i.e., accident type. Therefore, the joint probabilities of each variable with “accident type” are calculated, as shown in Table 6. Table 6 presents the states of each variable that have the highest and lowest impact on the accident type. For example, under the case of “track damage detected,” the likelihood of “derailment” is the highest (63.5%); under the case of “train driver skills insufficient,” the likelihood of “collision” is the highest (34.1%); when the train is in an “overspeed” state, the probability of “collision” is the highest (32.3%); when the train is in an “empty load” state, the likelihood of “derailment” is the highest (35.7%); when the train is in an “overload” state, the probability of “explosion” is the lowest (5.00%).

This demonstration demonstrates the impact of specific states of individual variables on the accident type. Meanwhile, it shows how different states of individual variables affect the probability of specific accident types. After TRI sensitivity analysis, Table 7 displays the TRI values of “train driver skills” for “collision,” and Table 8 displays the TRI values of risk variables for all accident types. By comparing the updated values of the target node, it demonstrates that the model complies with Criterion 1.

Specifically, in Table 7, the first row represents the baseline scenario, and the subsequent rows represent different scenarios when each state of the variable reaches 100%. To obtain the degree of impact of RIFs on the accident types, the TRIs are compared and ranked. The most important variables for “accident type” are given below:

Undetected track damage > Train overload > Traction system failure > Brake system condition > Level crossing > Train driver skills > Bridge damage > Train speed

TABLE 5. Mutual information related to “Accident Type”.

RIFs	Variance Reduction	Percentage (%)	Mutual Info	Percentage (%)	Variance of Belief
R7	0.01826	0.365	0.10616	4.44	0.0126630
R6	0.2327	4.65	0.07941	3.32	0.0023798
R4	0.1811	3.62	0.07630	3.19	0.0028541
R10	0.1845	3.69	0.06873	2.88	0.0014033
R3	0.1665	3.33	0.06865	2.87	0.0018828
R21	0.2098	4.2	0.06499	2.72	0.0014095
R1	0.0479	0.958	0.06122	2.56	0.0050607
R24	0.1699	3.4	0.05318	2.22	0.0013037
R18	0.2526	5.05	0.04855	2.03	0.0017683
R11	0.1349	2.7	0.04719	1.97	0.0010389
R17	0.08693	1.74	0.04147	1.74	0.0009358
R9	0.1777	3.55	0.04145	1.73	0.0013107
R26	0.1787	3.57	0.04139	1.73	0.0011039
R20	0.0101	0.202	0.04139	1.64	0.0015419
R19	0.1197	2.39	0.03448	1.44	0.0008366
R23	0.06698	1.34	0.03435	1.44	0.0007008
R14	0.1445	2.89	0.03271	1.37	0.0010618
R5	0.001429	0.0286	0.03164	1.32	0.0008809
R22	0.0474	0.948	0.02844	1.19	0.0007192
R25	0.02956	0.591	0.02533	1.06	0.0005831
R15	0.04514	0.903	0.0228	0.954	0.0008603
R16	0.08266	1.65	0.02278	0.953	0.0006581
R2	0.001058	0.0212	0.02187	0.915	0.0004305
R12	0.08212	1.64	0.01987	0.831	0.0012151
R13	0.0002215	0.00443	0.01952	0.817	0.0007745
R8	0.0003992	0.00798	0.01192	0.499	0.0005168

Undetected track damage may lead to instability during train operations, increasing the risk of accidents and causing train derailments, deceleration, or other safety issues. The skill level of the driver is directly related to the safety of operating the train; insufficient skills may result in operational errors and accidents. Insufficient driver skills can impact decision-making abilities during emergency situations, increasing the severity of accidents. Overloading may cause a decline in train performance, reduced braking system effectiveness, and an elevated risk of accidents. Accidents under overload conditions can jeopardize the safety of passengers and cargo. Brake system failures may lead to brake malfunction, increasing the risk of collisions. High speeds may make it more challenging to brake during emergencies, escalating the likelihood and severity of accidents. Accidents at high speeds can result in greater impact forces, causing more serious damage to personnel and equipment. Traction system failures may result in the train losing power, increasing the risk of railway accidents. Level crossing accidents may lead to collisions between vehicles and pedestrians or other traffic participants. The consequences of accidents may involve injuries to individuals and damage to vehicles. Bridge damage may cause accidents when trains pass through, especially on unstable or deteriorated bridges, leading to

train derailments or bridge collapses and causing significant accidents.

From another perspective, different accident types are associated with the priority levels of different variables. The priority list of the most important variables for different accident types is listed in Table 9. For instance, “Undetected track damage” is the most critical RIF for “derailment” and “collision,” while “Exceeding stopping point” is the least important RIF. Additionally, accidents such as “derailment,” “collision,” and “near miss” are more frequently caused by “Undetected track damage” compared to accidents like “conflict” and “loss of control.”

3) MODEL VALIDATION

To verify the effectiveness of the model, multiple RIFs were tested to investigate their combined impact on accident types. By considering the different states of the parent nodes, the change value for each state was obtained. “Undetected track damage” was selected as the first node. The state with the maximum change value in accident types (i.e., derailment) increased by 10%, while the state with the minimum change value in accident types decreased by 10%. This process is shown as “~10%” in Table 10. Then, the same method was

TABLE 6. Joint probability of TAN model.

R7								
	S1	S2	S3	S4	S5	S6	S7	S8
1	26.4	36.3	2.98	2.09	18.7	3.88	3.88	5.86
2	63.5	4.53	3.72	3.48	4.44	3.87	3.87	12.5
R6								
	S1	S2	S3	S4	S5	S6	S7	S8
1	38.3	30.2	2.26	2.01	16.5	1.76	2.44	6.59
2	10.5	34.1	6.33	3.53	15.5	12.0	9.42	8.48
R4								
	S1	S2	S3	S4	S5	S6	S7	S8
1	32.6	35.7	1.98	1.28	17.3	2.75	2.75	5.63
2	21.9	14.9	8.52	7.30	14.4	9.46	9.46	14.0
3	21.9	14.9	8.52	7.30	14.4	9.46	9.46	14.0
4	46.5	10.2	5.84	5.00	9.87	6.48	6.48	9.6
R10								
	S1	S2	S3	S4	S5	S6	S7	S8
1	33.4	33.6	2.58	1.80	17.1	2.79	2.18	6.57
2	26.2	10.1	7.26	6.55	9.70	12.6	17.5	10.2
R3								
	S1	S2	S3	S4	S5	S6	S7	S8
1	35.9	30.7	2.23	1.98	17.1	1.74	3.74	6.60
2	18.7	32.3	6.70	3.74	12.8	12.7	4.43	8.53
R21								
	S1	S2	S3	S4	S5	S6	S7	S8
1	34.2	32.8	2.52	1.28	16.7	3.31	2.72	6.42
2	16.0	12.4	8.96	12.9	12.0	9.58	15.6	12.6
R1								
	S1	S2	S3	S4	S5	S6	S7	S8
1	38.9	26.6	3.01	2.09	15.3	3.95	2.54	7.65
2	12.2	45.3	3.41	3.07	19.5	3.64	8.19	4.79
R24								
	S1	S2	S3	S4	S5	S6	S7	S8
1	38.3	30.2	2.26	2.01	16.5	1.76	2.44	6.59
2	10.6	34.1	6.33	3.53	15.5	12.0	9.42	8.48

TABLE 7. TRI of collision risk variable (Train Driver Skills).

Train driver skills						
1	2	collision	HRI	LRI	TRI	
/	/	31.0				
100%	0	30.2	3.1	0.8	1.95	
0	100%	34.1				

applied to the next RIF, and the cumulative change value was obtained and updated. This updating process continued until all RIF nodes were included. Similarly, the same updating process was applied to states 2, 3... 8 of the “accident type” until all states were included.

The first column in Table 10 shows the original values in the TAN model, while the other columns show the updated

change values of the results. Each state of the “accident type” is calculated separately, i.e., each row is computed based on the state changes of the RIFs within each accident type. It can be observed from Table 10 that the update values of the target node gradually increase or decrease as the RIFs change, thus validating Criterion 2.

The impact of small variations in variables on railway accidents may be minimal in certain situations, but in other cases, it could lead to significant outcomes. The occurrence of accidents is typically the result of complex interactions among multiple factors, and small changes may trigger chain reactions in these factors. Particularly in scenarios involving human operations or decision-making, minor errors or variations can have substantial impacts on the entire process. As shown in Table 10, when there is a slight variation in the human factor “train driver skill,” the probabilities of accidents classified as “conflict,” “attempted,” and “loss of control” all increase.

Moreover, railway accidents that were not previously included in the database were also simulated on the model to validate its effectiveness. For example, in this study, the RAIB report R142020 was taken for a case analysis. This event occurred at 22:44 on March 23, 2020, when a locomotive collided with the buffer stops at the end of a siding approximately 700 meters south of the Bromsgrove station. Less than a minute later, a northbound passenger train collided with the corner of the locomotive. All parameter settings of the proposed BN model were derived from the accident report as follows:

- (1) The accident occurred on the main line, not at a level crossing.
- (2) The accident occurred in very dark conditions. The weather was clear and dry, with recorded temperatures ranging from 2°C to 4°C. The local area had minimal lighting conditions.
- (3) The locomotive entered the siding, and the driver was distracted in the driving process, causing the locomotive to collide with the buffer stops and veer to the left.
- (4) The passenger train was traveling towards the Bromsgrove station at a speed of about 85 miles per hour (136 kilometers per hour). Even if the train driver had applied the brakes upon seeing the headlight in the darkness, he would not have been able to stop the train before the accident occurred.
- (5) After the collision with the buffer stops, the first action taken by the locomotive driver was to exit the locomotive and check for any derailment. This process took about 20 seconds, so the driver did not have time to issue any warnings to the signaler.

In the accident report, for situations where certain information was not recorded or updated, evidence was not collected from the accident report, and the other nodes maintained their generic original probabilities. As illustrated in Fig. 4, based on the above parameter settings, the probability of a collision occurring between the trains was up to 72.8%. This further verifies the effectiveness of the proposed model.

TABLE 8. The TRI of risk variables for all accident types.

Node	TRI								Average
	S1	S2	S3	S4	S5	S6	S7	S8	
R7	18.55	16.035	0.37	0.695	7.13	0.005	0.005	3.32	5.764
R6	13.85	1.95	2.035	0.76	0.5	5.12	3.49	0.945	3.581
R4	12.3	12.75	3.27	3.01	3.715	3.355	3.355	4.185	5.743
R10	3.6	11.75	2.34	4.485	3.7	4.905	7.66	3.425	5.233
R3	8.6	0.8	2.235	0.88	2.15	5.48	0.345	0.965	2.682
R21	9.1	10.2	3.22	5.81	2.35	3.135	6.44	3.09	5.418
R1	13.35	9.35	0.155	0.49	2.1	0.155	2.825	1.43	3.732
R24	13.84	1.92	2.03	0.77	0.4	5.15	3.50	0.94	3.569

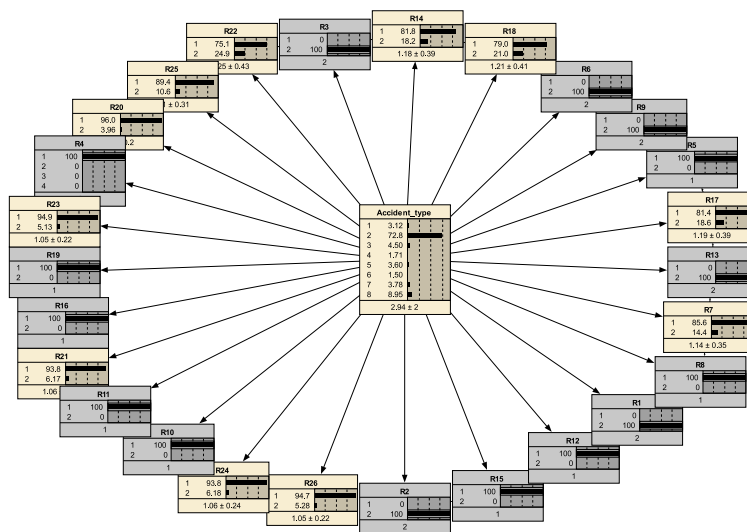


FIGURE 4. Verification of the model based on the previous railway accident.

TABLE 9. The priority of important variables.

Accident type	R7	R6	R4	R10	R3	R21	R1	R24
S1	1	2	5	8	7	6	4	3
S2	1	6	2	3	8	4	5	7
S3	7	6	1	3	4	2	8	5
S4	7	6	3	2	4	1	8	5
S5	1	7	2	3	5	4	6	8
S6	8	3	5	4	1	6	7	2
S7	7	4	5	1	6	2	6	3
S8	3	7	1	2	6	4	5	8

4) IMPLICATIONS

To observe the connections between related nodes in a BN and identify the most likely states within the nodes, the TAN-BN model can provide the MPE based on a determined

accident type. By exploiting the MPE mode of the BN, the most probable risk factors under the current accident type can be observed. This method provides a more comprehensive and reliable solution for analyzing railway accidents, predicting the causes of accident occurrences, and aiding in accident prevention.

As demonstrated in Fig. 5, under the MPE mode, “derailment” is the most likely accident type, while other RIFs exhibit their most likely states. That is, train accidents involving “derailment” usually occur under the following conditions.

(1) The accident occurs not at a level crossing, between 6:00 and 18:00, with a normal train speed and under good weather conditions.

(2) The train driver is fatigued, the track timber is rotten and deteriorated, and maintenance personnel lack sufficient skills to identify and detect relevant hazards in time.

TABLE 10. Accident probability with minor variations in variables.

R7	/	~10%	~10%	~10%	~10%	~10%	~10%	~10%	~10%
R6	/	/	~10%	~10%	~10%	~10%	~10%	~10%	~10%
R4	/	/	/	~10%	~10%	~10%	~10%	~10%	~10%
R10	/	/	/	/	~10%	~10%	~10%	~10%	~10%
R3	/	/	/	/	/	~10%	~10%	~10%	~10%
R21	/	/	/	/	/	/	~10%	~10%	~10%
R1	/	/	/	/	/	/	/	~10%	~10%
R24	/	/	/	/	/	/	/	/	~10%
S1	32.6	33.9	32.5	34.3	33.7	32.7	32.5	32.3	33.7
S2	31.0	31.3	31.1	30.0	30.7	31.2	31.6	31.4	30.5
S3	3.10	2.00	2.50	2.80	3.40	3.08	3.60	3.70	2.70
S4	2.33	2.80	2.61	2.60	2.62	2.34	2.30	2.64	2.80
S5	16.3	15.72	15.8	16.31	15.7	16.2	16.4	14.9	16.7
S6	3.88	4.30	3.90	3.60	3.62	3.89	3.81	4.00	3.80
S7	3.88	3.50	3.70	3.71	3.69	3.88	2.80	4.30	3.72
S8	6.98	6.50	7.90	6.60	6.59	6.99	6.90	7.20	6.10

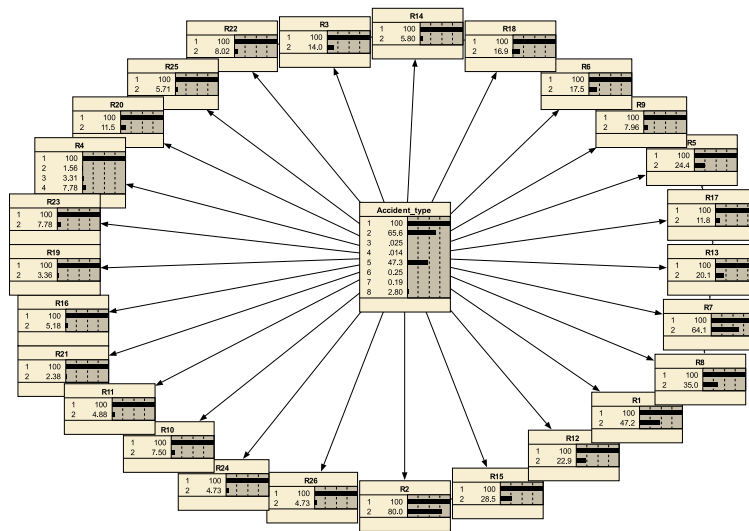


FIGURE 5. The MPE of the BN model.

This result implies that the railway operating group should review its management assurance procedures related to operational safety, also, it should take measures to ensure effective monitoring, auditing, and management review of its safety arrangements. This includes but is not limited to competency management of operational personnel, traffic acceptance, and general operating instructions. Meanwhile, train drivers need

to undergo appropriate brake performance tests at a high speed to ensure effective braking. Besides, considering the relevant laws, guidance, and good practices applicable in other industries, a re-assessment of driver health standards should be conducted to confirm their ability to effectively control the risks of the driving task. Additionally, the railway department should address fatigue risks appropriately and

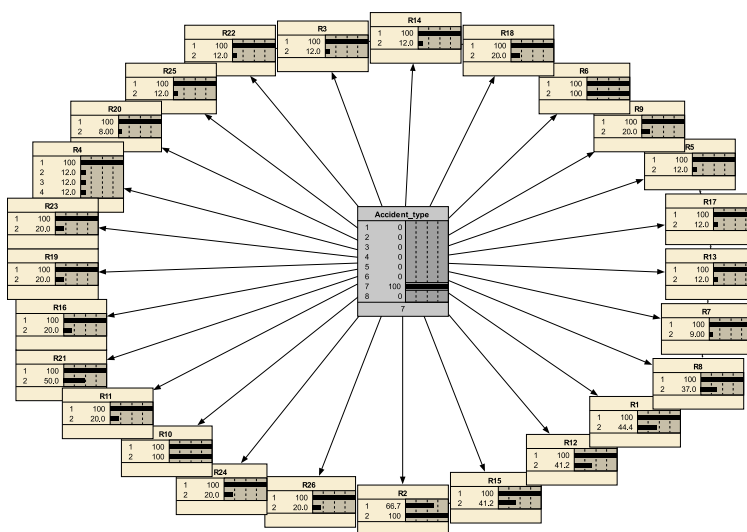


FIGURE 6. The MPE of the accident type “out of control”.

effectively by reviewing and improving its current fatigue risk management system for safety-critical personnel and ensuring compliance with relevant industry guidelines and best practices. For instance, in the case of accident reference number R092020, as suggested by RAIB, railway staff should be aware that stationary trains or vehicles may move without warning when they work or visit maintenance yards, sidings, and engineering works. Stationary trains or track vehicles may also conceal trains approaching another line, and staff must follow the requirements of their training and local procedures.

Studies have found that operational safety management has the greatest impact on most railway accidents, including dimensions of managing safety measures and safety training. Meanwhile, most railway accidents are caused by unsafe human behaviors, i.e., human factors. The supervision and control of train operations are based on human-environment communication, and improper human operations and delayed responses can cause serious railway accidents.

Similarly, when “accident type” is selected as state 7 (out of control), the MPE exhibits multiple 100% confidence bars, as demonstrated in Fig. 6. In this figure, multiple 100% confidence bars change when one state is selected. In the case of selecting the accident type in Fig. 5, the probability of a train experiencing an “out of control” situation is higher in the following conditions: accidents generally occurring between 18:00 and 6:00 in the next day, brake system failure, traction system failure or in good condition, insufficient skills of train drivers, and inadequate skills of personnel.

The establishment of a database based on railway accidents aims to comprehensively understand the safety aspects of railway systems and implement corresponding improvement measures. However, when utilizing such a database, it is crucial to acknowledge certain limitations and potential biases. The railway accident database may not be exhaustive, as some incidents may go unreported, undocumented, or excluded

from the database. This may lead to inaccurate estimations of accident frequency and types. Reports in the database may be influenced by preferences and standards of different organizations and individuals, introducing inconsistencies in the nature of accidents due to potential overestimation or underestimation. The timeliness of the database may experience delays, as the collection and compilation of accident data may take some time, hindering the timely analysis of recent changes and trends. Data quality may be impacted by errors, omissions, or inconsistencies during the collection and compilation processes, and inaccurate data can result in misleading conclusions. There may be selective bias in the database, recording only specific types or severities of accidents, potentially leading to an insufficiently comprehensive assessment of overall safety. Inconsistencies in variable selection and definitions within the database, or the absence of certain crucial variables, may affect a comprehensive evaluation of risk factors.

Despite the discussion of limitations associated with the database, each incident report sourced from the RAIB that we have chosen comprehensively documents detailed information leading to the occurrence of accidents. This includes an overview of the accident, geographical and natural environments, personnel conditions, machine and equipment status, operational conditions of the railway, consequences of the accident, and a detailed account of all relevant facts before, during, and after the occurrence of the accident. We have utilized the data records accumulated by the RAIB, a specialized institution responsible for railway accidents, over the past decade. The long-term accumulation of such data provides a comprehensive understanding of accident occurrences over an extensive time span, mitigating the impact of short-term random fluctuations on the results. We are cognizant of potential biases in the data arising from incomplete reporting or a tendency to more easily record specific types of accidents. To address this issue, we conducted a preliminary check of

data quality and approached potential biases cautiously in our analysis, ensuring the reliability of the study's conclusions. Through the integration of expert experience in establishing the TAN model, comparative analysis with historical data, model validation, and sensitivity analysis, the results of these analyses indicate that our research is founded on a reliable data basis.

IV. SUMMARY AND OUTLOOK

Compared to previous studies that focused on the causal factors related to the severity and probability of railway accidents, this study develops a new quantitative risk analysis method for railway accidents using a data-driven TAN approach. How different risk factors influence different types of railway accidents is investigated from an empirical and methodological perspective. First, to identify RIFs, railway accident reports released by the RAIB during the decade of 2011-2020 are selected to establish a railway accident database. Based on this, a risk-based TAN model is constructed to analyze RIFs in railway accidents. Finally, sensitivity analysis, scenario analysis, and MPE analysis are carried out to demonstrate the research findings.

By calculating mutual information, RIFs for different accident types are ranked. The results indicate that the key RIFs for railway accident types include "undetected track damage," "train driver skills," "loading," "braking system condition," "train speed," "traction system failure," "level crossing", and "bridge damage."

Scenario analysis provides reasonable explanations for the observed results and reveals the most probable scenarios for specific accident types. It helps to identify potential hazards and effectively assist railway authorities in formulating accident prevention measures.

While our study primarily focuses on the risk analysis of railway accidents in the United Kingdom, the methods employed in model establishment and performance validation are applicable to different contexts. Applying the findings of railway accident research to other geographical locations or populations may pose certain challenges and limitations. For instance, railway systems in different geographical locations may exhibit significant variations, including track design, technical standards, equipment configurations, and traffic regulations. This can necessitate consideration and adaptation to these differences when applying the results elsewhere. Populations in different regions possess distinct characteristics, encompassing cultural, socio-economic conditions, travel habits, among others. These factors may influence the occurrence and severity of accidents. However, these limitations serve to enhance the transparency of our study, providing a clear understanding of the complexities and uncertainties involved when applying research results to other geographical locations or populations. Consequently, decision-makers and practitioners can comprehensively grasp the applicability of research outcomes, fostering collaborative efforts across regions for a more judicious application of research findings and enhancing the broad applicability of the study.

Generally, the results obtained by the TAN model demonstrate the differences in important risk factors causing different accident types. This provides valuable insights for accident investigation and prevention. However, the MPE method has its limitations: its results may change with the introduction of irrelevant variables and even be deceptive under the most probable explanation.

Besides, there are limitations in data representation. In this study, 121 accident investigation reports were involved, and accident type 4 (explosion) accounted for only 1.65% of all accidents, i.e., two cases of explosion accidents. To obtain more representative results, continuous data collection is needed for model construction. In future work, more focus will be placed on evaluating variables that are difficult to measure in accident reports, i.e., studying the impact of human risk factors on railway accident risk analysis and investigating individual factors in railway accident risk analysis. Furthermore, employing a multidisciplinary approach that integrates knowledge from engineering, sociology, psychology, and other fields facilitates a comprehensive understanding of the responses, management, and impacts of railway accidents on diverse geographical locations and population groups. This approach contributes to a more nuanced comprehension of the complexity surrounding accident occurrences, thereby enhancing our ability to effectively address challenges posed by distinct geographical locations and populations. By implementing the above measures, it is possible to effectively address the limitations in data representation and the inadequacies in research direction identified in the study, thereby enhancing the accuracy and applicability of railway accident risk analysis.

REFERENCES

- [1] Y. Cao, Y. An, S. Su, G. Xie, and Y. Sun, "A statistical study of railway safety in China and Japan 1990–2020," *Accident Anal. Prevention*, vol. 175, Sep. 2022, Art. no. 106764, doi: [10.1016/j.aap.2022.106764](https://doi.org/10.1016/j.aap.2022.106764).
- [2] M. Ahmadi Rad, L. M. Lefsrud, and M. T. Hendry, "Application of systems thinking accident analysis methods: A review for railways," *Saf. Sci.*, vol. 160, Apr. 2023, Art. no. 106066, doi: [10.1016/j.ssci.2023.106066](https://doi.org/10.1016/j.ssci.2023.106066).
- [3] S. Kaeeni, M. Khalilian, and J. Mohammadzadeh, "Derailment accident risk assessment based on ensemble classification method," *Saf. Sci.*, vol. 110, pp. 3–10, Dec. 2018, doi: [10.1016/j.ssci.2017.11.006](https://doi.org/10.1016/j.ssci.2017.11.006).
- [4] V. N. M. G. I. Bargegol and M. Abolfazlzadeh, "Statistical analysis of the railway accidents causes in Iran," *Int. J. Eng.*, vol. 30, no. 12, pp. 1822–1830, Dec. 2017, doi: [10.5829/ije.2017.30.12c.02](https://doi.org/10.5829/ije.2017.30.12c.02).
- [5] J. Liu, F. Schmid, K. Li, and W. Zheng, "A knowledge graph-based approach for exploring railway operational accidents," *Rel. Eng. Syst. Saf.*, vol. 207, Mar. 2021, Art. no. 107352, doi: [10.1016/j.res.2020.107352](https://doi.org/10.1016/j.res.2020.107352).
- [6] Z. Zhang, T. Turla, and X. Liu, "Analysis of human-factor-caused freight train accidents in the United States," *J. Transp. Saf. Secur.*, vol. 13, no. 10, pp. 1157–1186, Dec. 2019, doi: [10.1080/19439962.2019.1697774](https://doi.org/10.1080/19439962.2019.1697774).
- [7] I. Savage, "Analysis of fatal train-pedestrian collisions in metropolitan Chicago 2004–2012," *Accident Anal. Prevention*, vol. 86, pp. 217–228, Jan. 2016, doi: [10.1016/j.aap.2015.11.005](https://doi.org/10.1016/j.aap.2015.11.005).
- [8] J.-L. Zhou and Y. Lei, "Paths between latent and active errors: Analysis of 407 railway accidents/incidents' causes in China," *Saf. Sci.*, vol. 110, pp. 47–58, Dec. 2018, doi: [10.1016/j.ssci.2017.12.027](https://doi.org/10.1016/j.ssci.2017.12.027).
- [9] Q. L. Xue and H. P. Yang, "Interface hazard analysis between urban rail signal system and platform screen door system based on functional FMEA," *Control Inform. Technol.*, no. 1, pp. 101–106, Feb. 2023, doi: [10.13889/j.issn.2096-5427.2023.01.016](https://doi.org/10.13889/j.issn.2096-5427.2023.01.016).

- [10] K. Bian and X. W. Li, "Analysis of risk factors in railway hazardous cargo transport based on HAZOP," *Log. Techn.*, vol. 34, no. 13, pp. 13–15, Jul. 2015, doi: [10.3969/j.issn.1005-152X.2015.07.004](https://doi.org/10.3969/j.issn.1005-152X.2015.07.004).
- [11] P. Liu, L. Yang, Z. Gao, S. Li, and Y. Gao, "Fault tree analysis combined with quantitative analysis for high-speed railway accidents," *Saf. Sci.*, vol. 79, pp. 344–357, Nov. 2015, doi: [10.1016/j.ssci.2015.06.017](https://doi.org/10.1016/j.ssci.2015.06.017).
- [12] P. Ni and Z. B. Tang, "Analysis of subway fire accident based on FTA and FAHP," *Bul. Sci. Techn.*, vol. 37, no. 9, pp. 109–112, Sep. 2021.
- [13] X. Y. Li, M. Li, Y. Gao, and Y. M. Sun, "Risk analysis of train breaking into operation site by accident based on event tree method," *Rail Qual. Contr.*, vol. 49, no. 7, pp. 35–38, Jul. 2021, doi: [10.3969/j.issn.1006-9178.2021.07.011](https://doi.org/10.3969/j.issn.1006-9178.2021.07.011).
- [14] M. T. Baysari, A. S. McIntosh, and J. R. Wilson, "Understanding the human factors contribution to railway accidents and incidents in Australia," *Accident Anal. Prevention*, vol. 40, no. 5, pp. 1750–1757, Sep. 2008, doi: [10.1016/j.aap.2008.06.013](https://doi.org/10.1016/j.aap.2008.06.013).
- [15] C. Liang, M. Ghazel, O. Cazier, and E.-M. El-Koursi, "Risk analysis on level crossings using a causal Bayesian network based approach," *Transp. Res. Proc.*, vol. 25, pp. 2167–2181, May 2017, doi: [10.1016/j.trpro.2017.05.418](https://doi.org/10.1016/j.trpro.2017.05.418).
- [16] C. Liang, M. Ghazel, and O. Cazier, "Using Bayesian networks for the purpose of risk analysis at railway level crossings," *IFAC-PapersOnLine*, vol. 51, no. 9, pp. 142–149, Jul. 2018, doi: [10.1016/j.ifacol.2018.07.024](https://doi.org/10.1016/j.ifacol.2018.07.024).
- [17] P. Weber, G. Medina-Oliva, C. Simon, and B. Iung, "Overview on Bayesian networks applications for dependability, risk analysis and maintenance areas," *Eng. Appl. Artif. Intell.*, vol. 25, no. 4, pp. 671–682, Jun. 2012, doi: [10.1016/j.engappai.2010.06.002](https://doi.org/10.1016/j.engappai.2010.06.002).
- [18] N. Khakzad, F. Khan, and P. Amyotte, "Safety analysis in process facilities: Comparison of fault tree and Bayesian network approaches," *Rel. Eng. Syst. Saf.*, vol. 96, no. 8, pp. 925–932, Aug. 2011, doi: [10.1016/j.res.2011.03.012](https://doi.org/10.1016/j.res.2011.03.012).
- [19] A. Bobbio, L. Portinale, M. Minichino, and E. Ciancamerla, "Improving the analysis of dependable systems by mapping fault trees into Bayesian networks," *Rel. Eng. Syst. Saf.*, vol. 71, no. 3, pp. 249–260, Mar. 2001, doi: [10.1016/s0951-8320\(00\)00077-6](https://doi.org/10.1016/s0951-8320(00)00077-6).
- [20] W. Huang, X. Kou, Y. Zhang, R. Mi, D. Yin, W. Xiao, and Z. Liu, "Operational failure analysis of high-speed electric multiple units: A Bayesian network-K2 algorithm-expectation maximization approach," *Rel. Eng. Syst. Saf.*, vol. 205, Jan. 2021, Art. no. 107250, doi: [10.1016/j.res.2020.107250](https://doi.org/10.1016/j.res.2020.107250).
- [21] X. Y. Li and J. P. Qi, "Reliability analysis of multi pantograph system based on fuzzy Bayesian network," *J. Rail. Sci. Eng.*, vol. 15, no. 6, pp. 1383–1390, Jun. 2018, doi: [10.19713/j.cnki.43-1423/u.2018.06.003](https://doi.org/10.19713/j.cnki.43-1423/u.2018.06.003).
- [22] C. Liang, M. Ghazel, O. Cazier, and L. Bouillaut, "Advanced model-based risk reasoning on automatic railway level crossings," *Saf. Sci.*, vol. 124, Apr. 2020, Art. no. 104592, doi: [10.1016/j.ssci.2019.104592](https://doi.org/10.1016/j.ssci.2019.104592).
- [23] N. Friedman, D. Geiger, and M. Goldszmidt, "Bayesian network classifiers," *Mach. Learn.*, vol. 29, no. 2, pp. 131–163, Nov. 1997, doi: [10.1023/a:1007465528199](https://doi.org/10.1023/a:1007465528199).
- [24] Z. Yang, Z. Yang, and J. Yin, "Realising advanced risk-based port state control inspection using data-driven Bayesian networks," *Transp. Res. Part A, Policy Pract.*, vol. 110, pp. 38–56, Apr. 2018, doi: [10.1016/j.tra.2018.01.033](https://doi.org/10.1016/j.tra.2018.01.033).
- [25] X. Liang, S. Fan, J. Lucy, and Z. Yang, "Risk analysis of cargo theft from freight supply chains using a data-driven Bayesian network," *Rel. Eng. Syst. Saf.*, vol. 226, Oct. 2022, Art. no. 108702, doi: [10.1016/j.res.2022.108702](https://doi.org/10.1016/j.res.2022.108702).
- [26] M. G. Madden, "On the classification performance of TAN and general Bayesian networks," *Knowl.-Based Syst.*, vol. 22, no. 7, pp. 489–495, Oct. 2009, doi: [10.1016/j.knosys.2008.10.006](https://doi.org/10.1016/j.knosys.2008.10.006).
- [27] C. Chow and C. Liu, "Approximating discrete probability distributions with dependence trees," *IEEE Trans. Inf. Theory*, vol. IT-14, no. 3, pp. 462–467, May 1968, doi: [10.1109/tit.1968.1054142](https://doi.org/10.1109/tit.1968.1054142).
- [28] H. Alyami, Z. Yang, R. Riahi, S. Bonsall, and J. Wang, "Advanced uncertainty modelling for container port risk analysis," *Accident Anal. Prevention*, vol. 123, pp. 411–421, Feb. 2019, doi: [10.1016/j.aap.2016.08.007](https://doi.org/10.1016/j.aap.2016.08.007).
- [29] Z. L. Yang, J. Wang, S. Bonsall, and Q. G. Fang, "Use of fuzzy evidential reasoning in maritime security assessment," *Risk Anal.*, vol. 29, no. 1, pp. 95–120, Jan. 2009, doi: [10.1111/j.1539-6924.2008.01158.x](https://doi.org/10.1111/j.1539-6924.2008.01158.x).
- [30] B. Jones, I. Jenkinson, Z. Yang, and J. Wang, "The use of Bayesian network modelling for maintenance planning in a manufacturing industry," *Rel. Eng. Syst. Saf.*, vol. 95, no. 3, pp. 267–277, Mar. 2010, doi: [10.1016/j.res.2009.10.007](https://doi.org/10.1016/j.res.2009.10.007).
- [31] J. Liu, F. Schmid, W. Zheng, and J. Zhu, "Understanding railway operational accidents using network theory," *Rel. Eng. Syst. Saf.*, vol. 189, pp. 218–231, Sep. 2019, doi: [10.1016/j.res.2019.04.030](https://doi.org/10.1016/j.res.2019.04.030).



LEI SHI received the B.S. degree from Northwest University for Nationalities, Lanzhou, China, in 2004, and the M.S. degree from Lanzhou University, Lanzhou, in 2014. He is currently pursuing the Ph.D. degree with the School of Automation and Electrical Engineering, Lanzhou Jiaotong University. His current research interests include railway safety and reliability analysis, and fault diagnosis.



YAZHI LIU received the B.S. degree in rail transportation signal and control from Lanzhou Jiaotong University, Lanzhou, China, in 2021, where she is currently pursuing the M.S. degree in transportation engineering with the School of Automation and Electrical Engineering. Her research interest includes reliability analysis.



YOUPENG ZHANG received the B.S. degree in electrical engineering from East China Jiaotong University, Nanchang, China, in 1992, and the M.S. degree in transportation information engineering and control from Lanzhou Jiaotong University, Lanzhou, China, in 1995. He is currently a Professor with the Rail Transit Electrical Automation Engineering Laboratory of Gansu Province. His research interests include traction power supply technology for electrified railway and power quality analysis.



JUNYI LIANG received the B.S. degree in rail transportation signal and control from Lanzhou Jiaotong University, Lanzhou, China, in 2022, where he is currently pursuing the M.S. degree in transportation engineering with the School of Automation and Electrical Engineering. His current research interests include high-speed train electromagnetic compatibility and electromagnetic exposure safety assessment.