CrossMark

**ARTICLE**

# Improved Bayesian Network-Based Risk Model and Its Application in Disaster Risk Assessment

Ming Li[1] · Mei Hong[1] · Ren Zhang[1,2]

**Abstract** The application of Bayesian network (BN) theory in risk assessment is an emerging trend. But in cases where data are incomplete and variables are mutually related, its application is restricted. To overcome these problems, an improved BN assessment model with parameter retrieval and decorrelation ability is proposed. First, multivariate nonlinear planning is applied to the feedback error learning of parameters. A genetic algorithm is used to learn the probability distribution of nodes that lack quantitative data. Then, based on an improved grey relational analysis that considers the correlation of variation rate, the optimal weight that characterizes the correlation is calculated and the weighted BN is improved for decorrelation. An improved risk assessment model based on the weighted BN then is built. An assessment of sea ice disaster shows that the model can be applied for risk assessment with incomplete data and variable correlation.

**Keywords** Bayesian network · Genetic algorithm · Grey relational analysis · Risk assessment

✉ Mei Hong
flowerrainhm@126.com

[1] Research Center of Ocean Environment Numerical Simulation, Institute of Meteorology and Oceanography, National University of Defense Technology, Nanjing 211101, China

[2] Collaborative Innovation Center on Forecast and Evaluation of Meteorological Disaster, Nanjing University of Information Science and Technology, Nanjing 210044, China

## 1 Introduction

Risk is the consequence of interactions between risk factors and risk-bearing objects (Grandell 1991) in a multidimensional and multilayered system. Risk assessment, the core of risk science, is a comprehensive evaluation and estimation of the occurrence of risks and losses (Zhang 2013), and constitutes an important research area in the field of management and decision making. Qualitative risk assessment is mainly based on expert knowledge, whereas quantitative risk assessment uses mathematical methods (Bühlmann 1996).

Considering risk cognition and risk information, both risk and its assessment are uncertain, or rather fuzzy and random (Paté-Cornell 1996). There are several culprits that create this uncertainty: (1) the randomness of attributes of risk such as time, frequency, and intensity; (2) the incompleteness and ambiguity of environmental information; and (3) dependency on subjective knowledge. Therefore, the expression and handling of fuzzy and random information is one of the key issues in modern risk assessment modeling (Yu 2017).

There are many risk assessment methods, including qualitative and quantitative ones. Classic methods, such as the analytic hierarchy process (Saaty 1980), fuzzy comprehensive assessment (Yang and Yang 1998), and grey system theory (Deng 1990), are used widely. The analytic hierarchy process (AHP) combines qualitative judgment with quantitative analysis to handle subjective preference in a quantitative way (Al-Harbi 2001). Fuzzy comprehensive assessment (FCA) can process ambiguous information with quantitative mathematical expressions (Ruan et al. 2005). Grey evaluation applies grey relational analysis (GRA) to assessment modeling with incomplete data (Zheng and Hu 2009). However, these classical methods

are mainly based on subjective experience and expert knowledge. Neither weight calculation in AHP nor affiliation determination in FCA takes advantage of objective data, indicating a relatively strong subjectivity and low credibility in assessment. Furthermore, these methods have defects in describing nonlinear interactions between risk factors.

In the past two decades, some emerging research techniques have been utilized in risk assessment. Prominent among them have been the neural network (NN) approaches discussed by Hagan and Beale (2002), cloud models (Li and Liu 2004), event tree analysis (ETA) applied by Ericson (2005) to system safety studies, and Petri net (PN) techniques that Girault and Valk (2003) introduced into systems engineering. NN in particular has capability in parallel computing, self-learning, and fault tolerance, which can achieve nonlinear modeling for complex systems. Cloud modeling combines fuzzy theory with probability theory for the expression of uncertain information. ETA and PN can describe the interinfluence of factors visually and achieve the rigorous inference of multisource information.

Although these emerging methods overcome many of the problems encountered in classic methods, they still have limitations in handling uncertainty. NN has a weak ability to express fuzzy and random knowledge; Cloud modeling cannot achieve integral reasoning of information; and ETA and PN are subjective in relationship analysis and information fusion. In general, the vital issues in risk assessment, that is, objective fusion and reasoning of multisource and uncertain information, have not been effectively tackled.

Enlightened by artificial intelligence (AI), which is well recognized in processing uncertainty, risk assessment can be propelled by means of AI algorithms. As one of the most promising technologies in AI, Bayesian Network modeling has drawn the attention of researchers. It has been preliminarily applied to risk assessment, such as catastrophic risk assessment (Li et al. 2010), health risk assessment (Liu et al. 2012), risk analysis of marine strategic aisles (Yang and Zhang 2014), and risk assessment for ship-bridge collisions (Yang 2015). BN modeling in recent studies can be summarized as comprising 3 steps: (1) the selection of indicators as network nodes; (2) the manual establishment of network structure based on casual analysis among nodes; and (3) the determination of network parameters according to expert knowledge. Structures and parameters, nevertheless, are mostly determined by experiential knowledge with strong subjectivity and little data.

In order to take full advantage of the potential of BN in uncertainty processing, several methodologies have been combined to promote more objective and quantitative modeling in risk assessment. Consequently, fuzzy mathematics theory (Zhang 2015), object-oriented analysis (Wang 2016), weight fusion (Liu 2016), and geographic information science (Grêt-Regamey and Straub 2006) were introduced. These fields of study extract quantitative data from the original information and determine the structure and parameters automatically by intelligent algorithms. Objective data mining, instead of subjective construction, can be the catalyst of BN application in the expression, fusion, and reasoning of uncertain information in risk assessment.

The BN application in risk assessment is not a doddle. The practical problems are usually full of qualitative description but also a lack of quantitative data, where expert knowledge is a feasible bridge that transfers data from the former to the latter. Nonetheless, the "quantitative data" transformed by expert knowledge is incomplete and subjective. There remains a barrier in constructing a BN objectively and scientifically with nonquantitative information and incomplete data. Moreover, the modeling process hitherto used ignores the conditional independence hypothesis of BN, that is, network nodes should be conditionally independent of each other. This assumption is hard to meet in practical application due to the strong correlation between assessment factors, even though the independence assumption should be satisfied as much as possible.

An improved BN model, with parameter retrieval and decorrelation ability, is proposed to deal with data incompleteness and factor correlation in quantitative risk assessment. The model makes an attempt to break through the restrictions of BN and promote its application to risk assessment.

## 2 Bayesian Network and Its Applicability

BN is an emerging AI algorithm and has been initially applied in risk assessment. We first make a brief introduction of basic concepts and mathematical principles, and then explain its applicability in risk assessment.

### 2.1 Bayesian Network Theory

BN, also known as Bayesian reliability network, is a combination of graph theory and probability theory (Shi 2012). It is not only a graphical description of the causal relationships between variables that provides a way to visualize knowledge, but it is also a probabilistic reasoning technique for uncertainty. BN is expressed as a complex and causal diagram intuitively, and can be denoted by a binary $B = \langle G, \theta \rangle$:

- $G = (V, E)$ denotes a directed acyclic graph, $V$ is a set of nodes denoting variables in problem domain, and $E$ is a set of arcs denoting the causal dependency between variables.
- $\theta$ is the network parameter including the prior probability and the conditional probability table (CPT) of nodes. It expresses the influence degree between nodes and reflects quantitative features in the knowledge domain.

Assume a set of variables $V = (V_1, \ldots, V_n)$. The mathematical basis of BN is the Bayesian formula.

$$P(V_i|V_j) = \frac{P(V_i, V_j)}{P(V_j)} = \frac{P(V_i) \cdot P(V_j|V_i)}{P(V_j)} \tag{1}$$

where $P(V_i)$ and $P(V_j)$ are prior probabilities, $P(V_j|V_i)$ is a conditional probability, and $P(V_i|V_j)$ is a posterior probability.

If the prior probability distribution of root nodes and the CPT of non-root nodes are given, the joint probability distribution can be reasoned by Eq. 2.

$$P(V_1, V_2, \ldots, V_n) = \prod_{i=1}^{n} P(V_i|\mathrm{Pa}(V_i)) \tag{2}$$

where $P(V_1, V_2, \ldots, V_n)$ is the joint probability distribution of variables, $\mathrm{Pa}(V_i)$ denotes the parent of $V_i$.

BN construction includes structural learning and parameter learning, or rather learning about the topology network and CPT. BN can be not only learned by intelligent algorithms with big data, but also can be constructed with relevant expert knowledge and experience (Wang 2010). Therefore, BN achieves an effective combination of qualitative analysis and quantitative calculation.

## 2.2 Applicability Analysis of Bayesian Network in Risk Assessment

Risk modeling with BN is a complex process that includes variable definition, node selection, data processing, structural learning, CPT learning, and probability reasoning. Assessment theory has also developed a complete system consisting of index selection, system establishment, weight calculation, and index fusion. Both processes are circular and revised constantly. The brief modeling steps are shown in Table 1.

As shown in Table 1, BN coordinates with risk assessment—the BN modeling process corresponds to the risk assessment process in each step. The probabilistic reasoning technique of BN can effectively achieve the fusion of uncertain information, which is vital for risk assessment. Therefore, BN is an effective model for dealing with uncertain risk assessment.

However, the problems of missing quantitative data and variable correlation mentioned in Sect. 1 limit the

application of BN in risk assessment. Our next step is to optimize the BN modeling based on a genetic algorithm (GA) and an improved grey relational analysis (GRA). The concrete technical flow is shown in Fig. 1.

## 3 Optimization Method of Bayesian Network

Concerning missing quantitative data, we introduce multivariate function planning and apply feedback error searching based on a GA. A CPT retrieval algorithm of qualitative nodes with small samples is proposed. As for the conditional independence hypothesis that each node should be conditionally independent, an improved GRA is adopted to get optimal weights, which is then blended into probability distributions to achieve decorrelation.

### 3.1 Retrieval of a Conditional Probability Table Based on a Genetic Algorithm

Parameter learning affects the accuracy of network reasoning directly. The existing algorithms are not applicable in practical problems due to the irregularity of data. We will analyze the data problem in CPT learning and propose a retrieval algorithm.
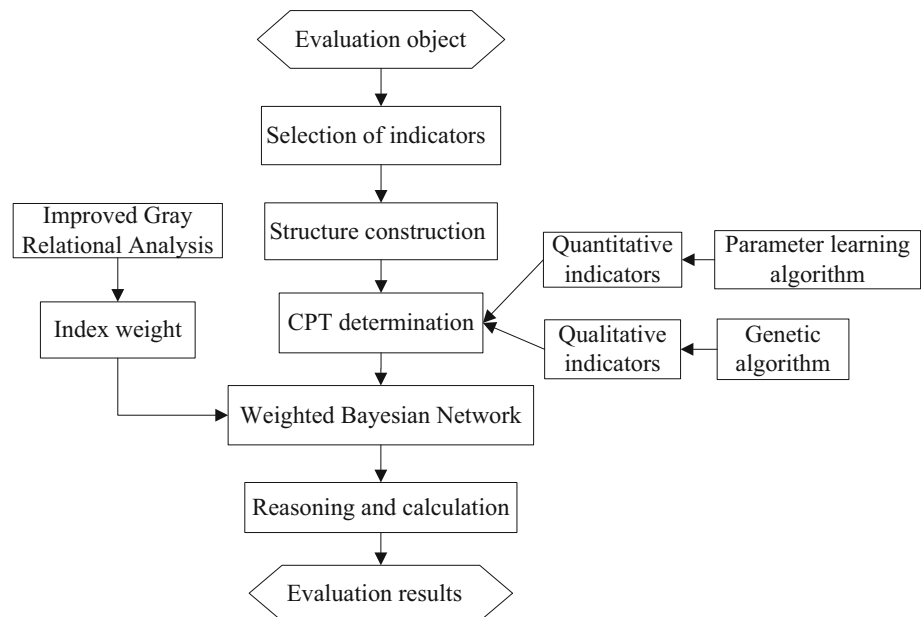
#### 3.1.1 Problem Analysis of Conditional Probability Table Learning

There are some parameter learning algorithms for complete data, including the maximum likelihood estimation method (Darwiche 2009), the Bayesian method, and the gradient descent method (Niculescu et al. 2006). And there are also some methods for dealing with incomplete data (Friedman et al. 1997; Basak et al. 2012), such as the expectation–maximization (EM) method and the Gibbs sample method. Although the existing algorithms could learn the parameters, they are no longer applicable in the following data conditions:

1. Non-quantitative data: a comprehensive description of a problem, especially of humanity and society, includes both quantitative data and qualitative language. When the raw data contains qualitative description, the above algorithms cannot be used directly for BN training. The traditional way to process qualitative data is to first quantify the qualitative information based on experience and knowledge, such as Delphi method (Okoli and Pawlowski 2005). There is no doubt that the process has a degree of subjectivity, which can easily cause a loss of objective information.

2. Incomplete data: The way of data collection and storage could result in missing data, but the above

**Table 1** Bayesian network modeling and assessment process

| BN modeling | Risk assessment |
| --- | --- |
| Step 1. Variable definition, node selection | Step 1. Indicator selection |
| Step 2. Structural learning | Step 2. System construction |
| Step 3. CPT learning | Step 3. Indicator weight calculation |
| Step 4. Reasoning and calculation | Step 4. Model construction |

**Fig. 1** The optimized modeling route of an ideal Bayesian network



algorithms require complete quantitative data. Although there are algorithms for processing missing data, which assume that the missing data is negligible, it is difficult to calculate CPT by the algorithms in the case of incomplete data, especially with excessively missing data. So how to calculate the CPT of qualitative and missing-data indicators objectively?

### 3.1.2 Retrieval Algorithm Design

Taking into account assessment results (such as disaster economic loss and casualty data) published by the authorities as a point of departure, we can get a true posterior probability distribution from an actual assessment and construct an error function. The feedback error function derived from searching for an optimal CPT is applied to reduce the learning errors. We propose the retrieval algorithm for parameter learning and use a genetic algorithm to search the optimal CPT of nodes.

Genetic algorithm (GA) is a randomized search method, which is derived from the evolution of biological circles. It was first proposed by Holland in 1992 (Holland 1992).

According to the GA process (Li 2002), we adapt the crossover and mutation operators, and fitness function to BN parameter learning. The CPT learning with GA (retrieval algorithm of CPT based on GA) is designed as follows.

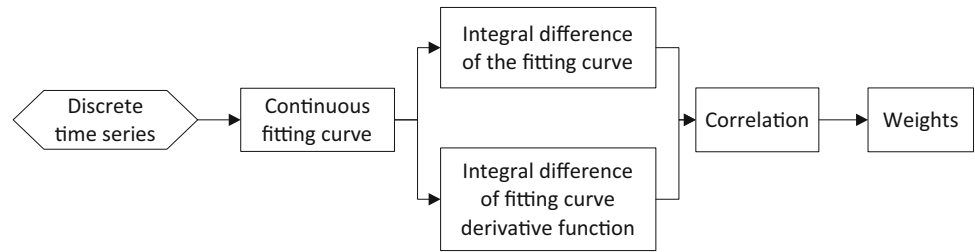| | |
| --- | --- |
| *Input* | CPT search space; Fitness function (error function) |
| *Output* | Optimal CPT |
| *Step* 1 | Creation of initial CPT population; |
| *Step* 2 | Statistical analysis of historical data to obtain posterior probability distribution of real assessment, to build the fitness function; |
| *Step* 3 | Perform crossover, mutation, and other genetic operations; |
| *Step* 4 | Selection according to fitness function; |
| *Step* 5 | Termination condition judgment, output optimal CPT |

The outstanding advantage is that the retrieval algorithm can achieve parameter learning with qualitative information and incomplete data. Based on objective data and error feedback, the CPT is retrieved by GA. In addition, learning efficiency is not related to the degree of missing data and the complexity of the network.

**Fig. 2** Improved thinking of grey relational analysis to correlate two dynamic sequences



## 3.2 Optimization of Weighted BN Based on Improved Grey Relational Analysis

BN is theoretically not suitable for modeling with related variables. We analyze the problem in the existing weighted BN, then apply an improved GRA to the weight calculation to improve the conditional independence hypothesis.

### 3.2.1 Problem Analysis of Variable Correlation

The conditional independence hypothesis simplifies the probability reasoning of a BN. On the basis of this assumption that non-associated nodes are conditionally independent, we use the great posteriori estimation (Liu 2016) to obtain the reasoning formula for posterior probability.

$$P(v|V_1, V_2 \ldots, V_n) = P(v) \prod_{i=1}^{n} P(V_i|v) \tag{3}$$

where under a given condition $v$, $V_1, V_2, \ldots, V_n$ are mutually conditional independent. $P(v|V_1, V_2 \ldots, V_n)$ is the posterior probability, $P(v)$ is the a priori probability, and $P(V_i|v)$ is the conditional probability.

Equation 3 is no longer applicable for related BN nodes in an actual problem. We need to improve the independence hypothesis and promote the application of BN in assessment. There are some ways to improve the hypothesis and give each node a different weight is an effective method. Liu (2016) adopted entropy weight to construct a weighted BN to assess flood disaster risk. Entropy weight is based completely on quantitative data, whose rationality requires further examination. Weight calculation is an essential part of the weighted BN. By considering the causal associations between BN nodes, the weight should be measured through the correlation of the nodes. In this way we improve the grey correlation method to analyze the correlation between variables and obtain the new weight to optimize the weighted BN.

### 3.2.2 Improved Grey Relational Analysis Design

Grey assessment is a method used to measure the correlation degree between factors based on the similarity or dissimilarity of the development trend (Gu et al. 2003). Grey relational analysis is used to judge the correlation between two sequences according to the geometry shapes of curves. The method makes up for the shortcomings of mathematical statistical methods, which require that samples should be subject to the typical probability distribution (Jiang et al. 2015). The specific process of GRA is as follows.

*Step* 1 Determine the reference sequence and comparison sequences;
*Step* 2 Data calculation and transformation (dimensionless, normalized);
*Step* 3 Calculate the difference sequence and the maximum, minimum difference;
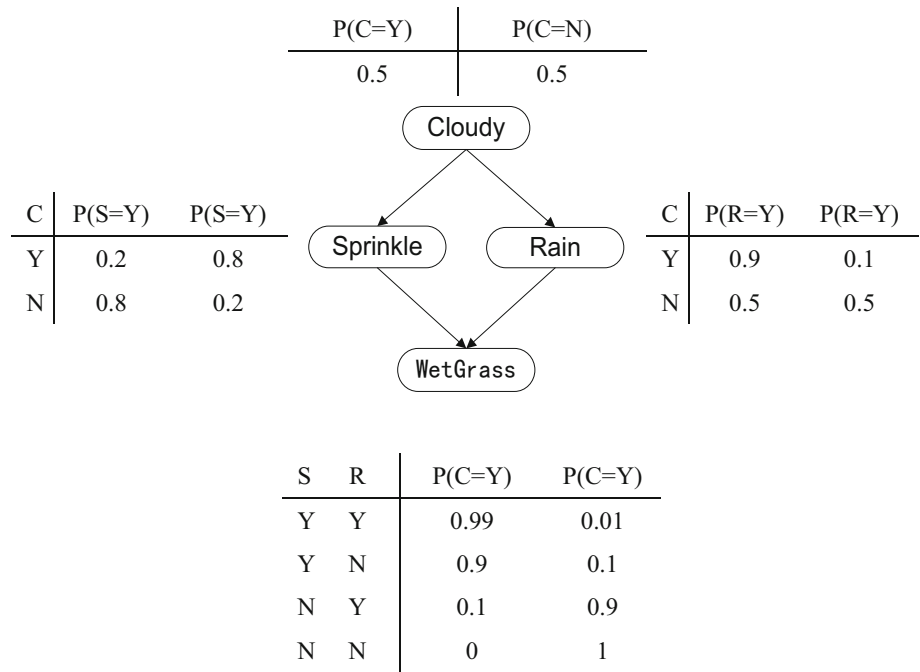*Step* 4 Calculate the correlation coefficient and correlation degree.

Grey relational analysis is a form of geometric processing that is a quantitative analysis of development trends in a dynamic process. The calculation of absolute difference at each moment ignores the variation between adjacent moments. Because it only takes into account the difference of each endpoint, the method cannot describe the correlation between two dynamic sequences.

To overcome this problem, we present the improved thinking as shown in Fig. 2. The sequence first is fitted nonlinearly, and the discrete function is replaced by a continuous function. Then the calculation of difference sequence is improved as shown in Eq. 4, and the changing trend is expressed in integral form. The correlation degree of the geometric shapes is described accurately.

$$\Delta_{0i} = \left| \int_{t}^{t+\Delta t} \{f_i(t) - f_0(t)\} dt \right| + \left| \int_{t}^{t+\Delta t} \{f_i'(t) - f_0'(t)\} dt \right|,$$
$$i = 1, 2, 3 \ldots m \tag{4}$$

**Fig. 3** "Wet Grass" Bayesian Network

| P(C=Y) | P(C=N) |
|--------|--------|
| 0.5 | 0.5 |

**Cloudy**

| C | P(S=Y) | P(S=Y) |
|---|--------|--------|
| Y | 0.2 | 0.8 |
| N | 0.8 | 0.2 |

**Sprinkle**     **Rain**

| C | P(R=Y) | P(R=Y) |
|---|--------|--------|
| Y | 0.9 | 0.1 |
| N | 0.5 | 0.5 |

**WetGrass**

| S | R | P(C=Y) | P(C=Y) |
|---|---|--------|--------|
| Y | Y | 0.99 | 0.01 |
| Y | N | 0.9 | 0.1 |
| N | Y | 0.1 | 0.9 |
| N | N | 0 | 1 |

where $\Delta_{0i}$ is the sequence difference, $f_i(t)$ is the fitting function of comparison sequence, $f_i'(t)$ is the derivative function of $f_i(t)$, $f_0(t)$ is the fitting function of reference sequence, and $f_0'(t)$ is the derivative function of $f_0(t)$.

The process of weight calculation based on correlation degree is as follows:

First calculate the sequence difference at each moment $\Delta_{0i}(k)$, the maximum and minimum difference $\Delta_{max}$, $\Delta_{min}$ according to Eq. 4. Then calculate the correlation coefficient.

$$\varepsilon_{0i}(k) = \frac{\rho\Delta_{max} + \Delta_{min}}{\Delta_{0i}(k) + \rho\Delta_{max}}, \quad k = 1, 2, 3, \ldots n \tag{5}$$

where resolution coefficient $\rho$ generally takes 0.5. Next calculate correlation degree.

$$\gamma_{0i} = \frac{1}{n}\sum_{k=1}^{n}\varepsilon_{0i}(k) \tag{6}$$

Finally, the weight of each indicator can be calculated.

$$w_i = \frac{\gamma_{0i}}{\sum_{i=1}^{m}\gamma_{0i}} \tag{7}$$

We integrate weights into the CPT to construct the weighted BN, and the reasoning formula 3 is improved as follows:

$$P(v|V_1, V_2\ldots, V_n) = P(v)\prod_{i=1}^{n}P(V_i|v)^{w_i} \tag{8}$$

where $w_i$ is the weight of each assessment indicator.

The improved GRA takes into full consideration both state and variation similarity, so it can measure the correlation of two sequences more accurately. Correspondingly, weights calculated by this method can better measure the interdependence between nodes and can improve the conditional independence assumption.

The new BN based on the above two algorithms could achieve objective assessment modeling in case of missing quantitative data and variable relation. In order to explain the improvement of modeling process, we make a comparison between the improved BN and classic BN (Table 2).

### 3.3 Algorithm Numerical Test

To verify the validity of the CPT retrieval algorithm, we undertake a comparative analysis with another method for parameter learning.

We use the classic BN "Wet Grass" in Fig. 3 for discussion. We first generate training data for this network, then calculate the CPT of node "Rain" by using an expectation–maximization algorithm and a retrieval algorithm. Comparison of the reasoning results is shown in Fig. 4 and Table 3. Considering the randomness of our search, we run a GA 100 times, each with a random initial population, and check whether the parameters converge to unique.

From Table 3, as the GA achieves the feedback of reasoning errors, the calculation error decreases by

**Fig. 4** Convergence curves of different algorithms (left for expectation–maximization (EM) algorithm, right for retrieval algorithm)
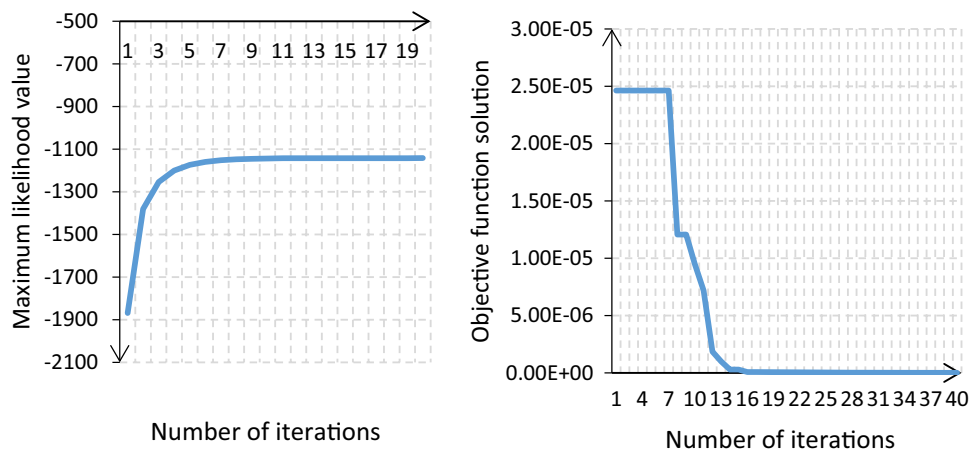


**Table 2** Comparison between the improved Bayesian network and classic Bayesian network models

| Modeling conditions | Quantitative data loss and strong variable relation | | |
| --- | --- | --- | --- |
| Modeling process | Structural learning | Parameter learning | Reasoning calculation |
| Classic BN | Manual construction based on knowledge | Expert scoring to quantify qualitative information | Simple reasoning without considering relation |
| Improved BN | Manual or automatic construction | Retrieval with GA | Weighted reasoning |

59.55%, and the convergence speed increases by 67.94%. In general, the retrieval algorithm not only has the ability to calculate the CPT with very limited quantitative data, but also it enhances accuracy.

## 4 Model Application

China is one of the countries that are seriously affected by sea ice in the world. Sea ice poses a major threat to offshore engineering facilities, resources, property, and personal safety, which has an adverse effect on the normal life of people in coastal areas. In this section, we select the Chinese Bohai Sea (37°N–41°N, 118°E–122°E) and four coastal provinces (Liaoning, Shandong, Hebei, and Tianjin) as the research setting in which to assess the potential for a sea ice disaster. The risk assessment is carried out in a *Matlab 2012a*[1] environment. The library functions come from the Bayesian Network Toolbox (BNT) written by Murphy.[2]

### 4.1 Node Selection and Structure Construction

Our purpose is to test the feasibility of our model. For the sake of comparison and verification, based on risk theory (Dilley et al. 2005) and expert knowledge (Sun and Shi 2012; Yuan et al. 2013), we select one representative indicator from each of three criteria in the risk system to illustrate the sea ice risk problem: danger, vulnerability, and precaution.

1. Danger: Maximum freezing range. This criterion can be represented by sea ice coverage, which is a quantitative indicator and a significant factor that contributes to sea ice disaster risk.
2. Vulnerability: Marine economic density. This factor is a ratio of total marine production to regional area. Marine economic production mainly includes marine fisheries, the offshore oil and gas, ocean engineering construction, and marine transportation industries, and so on.
3. Precaution: Social security level. This variable is a measure of societal preparedness, which is related to economy, medical care, transportation, and so on. The indicator contains all kinds of information, which is complex and non-quantitative.

Sea ice disaster risk also is measured by economic losses. An assessment system is shown in Table 4. The map of the corresponding BN structure is shown in Fig. 5.

**Table 3** Comparison of the conditional probability table calculation with different algorithms

|  | Number of iterations, time of convergence | CPT of node "rain" | | Posteriori probability | | Error |
|---|---|---|---|---|---|---|
| Expectation–maximization (EM) algorithm | 9 times 79.7752 s | 0.4758 | 0.5242 | 0.3352 | 0.6648 | 2.67% |
|  |  | 0.9143 | 0.0857 |  |  |  |
| Retrieval algorithm | 16 times 25.5728 s | 0.4960 | 0.5040 | 0.3458 | 0.6542 | 1.08% |
|  |  | 0.8997 | 0.1003 |  |  |  |
| True situation | – | 0.5 | 0.5 | 0.3529 | 0.6471 | – |
|  |  | 0.9 | 0.1 |  |  |  |

**Table 4** Assessment indicator system to measure potential sea ice hazard risk in the Bohai Sea

| Assessment target | Indicator | Attributes |
|---|---|---|
| $T$: Sea ice disaster | $d_1$: Maximum freezing range | Quantitative indicators |
| Economic losses | $d_2$: Marine economic density | Marine GDP/area; quantitative indicator |
|  | $d_3$: Social security level | Qualitative indicators |

## 4.2 Conditional Probability Distribution Calculation

After the construction of a Bayesian network structure, a CPT is calculated. We first need to process the raw data to get training samples. Then the retrieval algorithm with the GA developed in Sect. 3.1 is used to search for the probability distributions.

### 4.2.1 Data Processing

We collected indicator data from the provinces of Liaoning, Shandong, and Hebei and Tianjin Municipality between 1951 and 2015. The data of $T$, $d_1$, and $d_2$ are quantitative, continuous, and year-by-year statistics, whose sources are shown in Table 5. The information for $d_3$ reflects the qualitative description of experts. Because BN processes discrete data, the discretization of data was performed to produce samples and determine the number of states taken by nodes.

In risk assessment, the discrete states are determined depending on the risk level of the indicators. Based on the *Sea Ice Grade Standard* compiled by the State Oceanic Administration (2010) as shown in Table 6, we classified the disaster level into two states: high risk (including level I, II, III) and low risk (including level IV, V). That is, each network node has two states. Node value {1, 2} is on behalf of two states {high risk, low risk}. Discrete samples were generated as shown in Table 7. $T$, $d_1$, and $d_2$ were divided into two states, while $d_3$ lacks quantitative information. We used the first 50 years (1951–2010) of samples for BN training and the last 5 years' data (2011–2015) for a test.
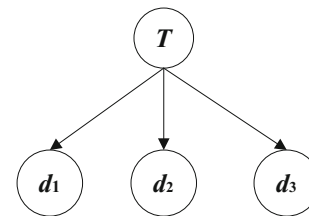


**Fig. 5** Bayesian network structure of sea ice disaster in the Bohai Sea region of China

### 4.2.2 Retrieval of Optimal Conditional Probability Table Based on Genetic Algorithm

For quantitative nodes $d_1$ and $d_2$, we calculated a priori probability distributions $P(d_1)$, $P(d_2)$ and CPT $P(d_1|T)$, $P(d_2|T)$ based on statistical analysis and a maximum likelihood estimation algorithm. For qualitative node $d_3$, we calculated $P(d_3)$ and $P(d_3|T)$ with the retrieval algorithm in Sect. 3.1. The specific steps are as follows.

*Step 1* Determine the initial species. To ensure the rationality of the search results, we first determine the search space based on expert knowledge as $P(d_3 = 1) \in [0, 0.35]$. Then we set the CPT as shown in Table 8.

*Step 2* Build the fitness function. According to Table 7, we analyzed the economic loss due to sea ice disaster in the Bohai Sea from 1951 to 2010 and determined the actual probability distribution of the sea ice disaster level (Table 9). Finally the objective function $f(x, y, z)$ representing assessment errors is constructed.

$$f(x, y, z) = |P(T = 1) - P(T_0 = 1)| + |P(T = 2) - P(T_0 = 2)| \tag{9}$$

**Table 5** Data sources for sea ice disaster and economic losses in Liaoning, Shandong, Hebei, and Tianjin, China, 1951–2015

| Assessment index | Data sources |
|---|---|
| Maximum freezing range | *China Marine Disaster Forty Years of Information Compilation* (Yang and Tian 1994), *Marine Disaster Bulletin*[a] |
| Marine GDP | *China Ocean Statistical Yearbook*,[b] |
| | *China Fishery Statistical Yearbook*[c] |
| Sea area | *Statistical Yearbook* compiled by provincial statistical offices[d] |
| Economic loss | *Marine Disaster Bulletin*, |
| | *Statistical Yearbook* compiled by provincial statistical offices |

[a]Marine Disaster Bulletin: http://www.soa.gov.cn/zwgk/hygb/zghyzhgb/

[b]China Ocean Statistical Yearbook: http://cyfd.cnki.com.cn/N2015050179.htm

[c]China Fishery Statistical Yearbook: https://www.douban.com/note/497925665/?type=like

[d]Statistical Yearbook: http://www.nianjianku.com/

**Table 6** Level criteria needed to determine sea ice disaster states (*Source* State Oceanic Administration 2010)

| Level<br>Name | I<br>Great disaster | II<br>Major disaster | III<br>Big disaster | IV<br>General disaster | V<br>Minor disaster |
|---|---|---|---|---|---|
| Maximum freezing range (n mile) | > 100 | 81–100 | 61–80 | 51–60 | ≤ 50 |
| Marine economic density (Billion/km$^2$) | > 0.5 | 0.41–0.5 | 0.31–0.4 | 0.21–0.3 | ≤ 0.2 |
| Economic loss (Billion) | > 60 | 41–60 | 21–40 | 1–20 | ≤ 1 |
| Clustering level | High risk level | | | Low risk level | |

**Table 7** Discrete samples for Bayesian network training and testing

| Indicator | Training Sample | | | | | | Test sample | | |
|---|---|---|---|---|---|---|---|---|---|
| | 1951 | 1952 | 1953 | 1954 | $\cdots$ | 2010 | 2011 | $\cdots$ | 2015 |
| $T$ | 2 | 2 | 1 | 2 | $\cdots$ | 1 | 2 | $\cdots$ | 1 |
| $d_1$ | 1 | 2 | 2 | 2 | $\cdots$ | 2 | 1 | $\cdots$ | 1 |
| $d_2$ | 1 | 1 | 1 | 1 | $\cdots$ | 2 | 2 | $\cdots$ | 2 |
| $d_3$ | Qualitative information (language description) | | | | | | | | |

**Table 8** Conditional probability distribution at high and low risk levels

| $d_3$ | P($d_3$) | P($d_3|T$) | |
|---|---|---|---|
| | | High risk level | Low risk level |
| High risk level | $z$ | $x$ | $1 - x$ |
| Low risk level | $1 - z$ | $y$ | $1 - y$ |

**Table 9** Actual probability distribution of sea ice disaster in the Bohai Sea

| $T$ | $P_0(T)$ |
|---|---|
| High risk level | 0.709 |
| Low risk level | 0.291 |

### 4.3 Weighted Bayesian Network Based on Improved Grey Relational Analysis

where P($T$) denotes the reasoning result, and P($T_0$) denotes the actual assessment result.

*Step 3* Search optimization. Based on the BNT (Jiang and Lin 2007), we design a mutation operator and a select operator, and search for the optimal CPT with minimizing inference errors. The results are shown in Fig. 6 and Table 10. The complete CPT is shown in Table 11.

We analyzed data from 1951 to 2010 to calculate weights and integrate them into the CPT. Because there is no quantitative data for $d_3$, the level of social security is evaluated by a 0–9 scale in order to determine its time series. The improved GRA was used to calculate the correlation between three indicators and the economic losses.

Then, we determined the weight according to Sect. 3.2.2. The results are shown in Table 12.

From Table 12, the improved correlation is more in accord with the correlation coefficient, so the dependency relationship between indicators is better explained with the weight. Finally, we integrated weights into the conditional probability and reasoned with Eq. 10. The network structure, CPT, weights, and reasoning mechanism of improved BN have been completed.
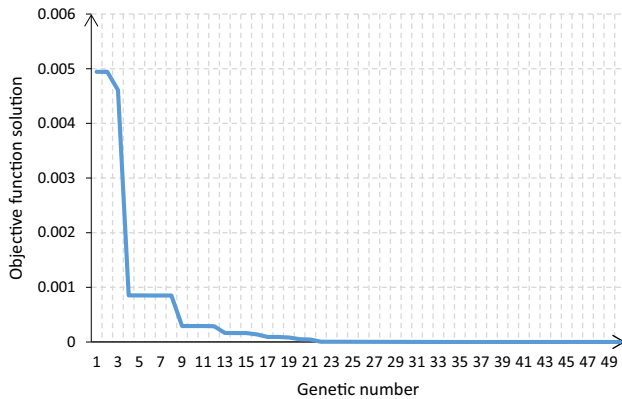


**Fig. 6** Convergence curve of the solution of the objective function

**Table 10** Optimal conditional probability table at high and low risk levels

| $d_3$ | $P(d_3)$ | $P(d_3|T)$ | |
|---|---|---|---|
| | | High risk level | Low risk level |
| High risk level | 0.312 | 0.719 | 0.281 |
| Low risk level | 0.688 | 0.296 | 0.704 |

**Table 11** Conditional probability table of each Bayesian network node

| $T$ | $d_i$ | | | | | |
|---|---|---|---|---|---|---|
| | $P(d_1|T)$ | | $P(d_2|T)$ | | $P(d_3|T)$ | |
| | 1 | 2 | 1 | 2 | 1 | 2 |
| 1 | 0.673 | 0.327 | 0.688 | 0.312 | 0.719 | 0.281 |
| 2 | 0.319 | 0.681 | 0.331 | 0.669 | 0.296 | 0.704 |

$$P(T|d_1, d_2, d_3) = P(T) \cdot P(d_1|T)^{w_1} \cdot P(d_2|T)^{w_2} \cdot P(d_3|T)^{w_3} \quad (10)$$

### 4.4 Reasoning Calculation and Model Discussion

Based on the improved BN, we input prior probability and reasoned according to Eq. 10 to evaluate sea ice disaster in the Bohai Sea from 2011 to 2015.

In order to test the effectiveness of the model, we compared it with the fuzzy comprehensive assessment (FCA) method and existing BN models. We respectively took the FCA method, the classic BN model, and the weighted BN model with the entropy method (Liu 2016) to carry out the same experiment and analyze the results.

1. The threat degree by the FCA is 0.616, while the probability distribution of sea ice risk is [0.697, 0.303]. Both results are high risk, which is relatively consistent. But FCA can only give a definite evaluation result, ignoring the uncertainty of risk. By contrast, the improved BN model can show all risk states and their probability. Moreover, the model contains less subjectivity than FCA.
2. Seen from Table 13, the assessment accuracy of the improved BN model can reach 97.89, 78.83% higher than the traditional BN and 60.41% higher than weighted BN with entropy method. The traditional BN model mainly fills missing data through expert scoring and does not consider correlation between variables, which makes the assessment subjective and inaccurate. The weighted BN model with entropy method considers correlation to some extent, but the weight calculation must be based on complete data filled with expert experience and the weight could not reflect the dependency between variables.

The improved BN model achieves objective risk assessment with BN in the case of missing data and variable correlation. In this article, we improve the weighted BN assessment model and apply GA to parameter learning. It not only realizes the CPT calculation with a lack of quantitative data but also achieves reasoning errors feedback. In addition, the conditional independence hypothesis of BN tends to be satisfied by integrating CPT with weights

**Table 12** Grey correlation analysis results for indicators of sea ice disaster

| | Maximum freezing range | Marine economic density | Social security level |
|---|---|---|---|
| Primitive correlation | 0.5831 | 0.5215 | 0.4941 |
| Improved correlation | 0.5746 | 0.5776 | 0.5571 |
| Correlation coefficient | 0.9488 | 0.9996 | 0.9326 |
| Weights | 0.3294 | 0.3469 | 0.3237 |

**Table 13** Reasoning results of sea ice risk by different Bayesian network models

|  | High risk level | Low risk level | Assessment error |
| --- | --- | --- | --- |
| Actual assessment | 0.712 | 0.288 | – |
| Improved BN | 0.697 | 0.303 | 2.11% |
| Weighted BN with entropy method | 0.674 | 0.326 | 5.33% |
| Traditional BN | 0.641 | 0.359 | 9.97% |

from GRA, which also makes the model more applicable in the actual situation.

## 5 Conclusion

To deal with uncertainties including randomness and fuzziness in risk assessment, BN is applied to construct an assessment model. Risk assessment with BN is carried out according to the following process: assessment system analysis, network structure establishment, node parameter learning, and reasoning. To address the CPT calculation of qualitative nodes and correlation of variables in BN application, we did the following work: (1) The CPT calculation for qualitative nodes is transformed into a multivariate function planning problem. The error function is constructed by combining the actual assessment result. Then GA is used for CPT retrieval with error feedback searching, which achieves CPT calculation with little quantitative data; (2) The weight calculation in weighted BN is improved with GRA. We improve the GRA by introducing variation rate into the calculation of sequence difference, which enhances the correlation degree between indicators. Thus, the optimized weight from correlation is obtained and the weighted BN is improved.

This research improves the application conditions of BN: big data and variable independence. In other words, the improved model can assess risk with BN in the case of incomplete data and correlated indicators. But there are also shortcomings in one respect—the optimal CPT obtained by GA must be derived from existing assessment truth, which to a certain extent limits use of the model.

## References

Al-Harbi, A.S. 2001. Application of the AHP in project management. *International Journal of Project Management* 19(1): 19–27.

Basak, A., S.V. Campus, I. Brinster, and O.J. Mengshoel. 2012. MapReduce for Bayesian network parameter learning using the EM algorithm. *Process of Big Learning Algorithms Systems and Tools* 15(1): 12–23.

Bühlmann, H. 1996. *Mathematical methods in risk theory*. Berlin: Springer.

Darwiche, A. 2009. *Modeling and reasoning with Bayesian networks: The maximum likelihood approach*. Cambridge: Cambridge University Press.

Deng, J.L. 1990. *Grey system theory tutorial*. Wuhan: Huazhong University of Science and Technology Press (in Chinese).

Dilley, M., R.S. Chen, U. Deichmann, A.L. Lerner-Lam, and M. Arnold, et al. 2005. *Natural disaster hotspots: A global risk analysis*. Disaster Risk Management Series No. 5. Washington, DC: World Bank, Hazards Management Unit.

Ericson, C.A. 2005. *Hazard analysis techniques for system safety*. Hoboken, NJ: Wiley.

Friedman, N., D. Geiger, and M. Goldszmidt. 1997. Bayesian network classifiers. *Machine Learning* 29(2): 131–163.

Girault, C., and R. Valk. 2003. *Petri nets for systems engineering*. Berlin: Springer.

Grandell, J. 1991. *Aspects of risk theory*. Berlin: Springer.

Grêt-Regamey, A., and D. Straub. 2006. Spatially explicit avalanche risk assessment linking Bayesian networks to a GIS. *Natural Hazards and Earth System Sciences* 6(6): 911–926.

Gu, C.D., H. Li, and S.H. Wu. 2003. Application of Grey System Theory to comprehensive assessment of new rice varieties. *Anhui Agricultural Sciences* 31(1): 98 (in Chinese).

Hagan, M.T., and M. Beale. 2002. *Neural network design*. Beijing: Machinery Industry Press.

Holland, J. 1992. Genetic algorithms. *Scientific American* 267(1): 66–72.

Jiang, W.D., and S.M. Lin. 2007. Bayesian learning and reasoning based on Bayesian network toolbox. *Information Technology* 6(2): 5–8 (in Chinese).

Jiang, S.Q., S.F. Liu, and Z.X. Liu. 2015. Grey relational decision model based on area. *Control and Decision* 4: 685–690 (in Chinese).

Li, M.Q. 2002. *Basic theory and application of genetic algorithm*. Beijing: Science Press (in Chinese).

Li, D.Y., and C.X. Liu. 2004. Universality of the normal cloud model. *China Engineering Science* 6(8): 28–34 (in Chinese).

Li, L., J. Wang, H. Leng, and C. Jiang. 2010. Assessment of catastrophic risk using Bayesian network constructed from domain knowledge and spatial data. *Risk Analysis* 30(7): 1157–1175.

Liu, F.R., C.F. Lu, C.W. Chen, and Y.S. Shen. 2012. Applying Bayesian belief networks to health risk assessment. *Stochastic Environmental Research and Risk Assessment* 26(3): 451–465.

Liu, R. 2016. *Research on risk assessment and modeling of flood disaster based on Bayesian network*. Shanghai: East China Normal University (in Chinese).

Niculescu, R.S., T.M. Mitchell, and R.B. Rao. 2006. Bayesian network learning with parameter constraints. *Journal of Machine Learning Research* 7(3): 1357–1383.

Okoli, C., and S.D. Pawlowski. 2005. The Delphi method as a research tool: An example, design considerations and applications. *Information and Management* 42(1): 15–29.

Paté-Cornell, M.E. 1996. *Uncertainties in risk analysis: Six levels of treatment. Reliability Engineering System Safety* 54(2): 95–111.

Ruan, B.Q., Y.P. Han, W. Hao, and R.F. Jiang. 2005. Fuzzy comprehensive assessment of water shortage risk. *Journal of Hydraulic Engineering* 36(8): 906–912 (in Chinese).

Saaty, T.L. 1980. *The analytic hierarchy process: Planning, priority setting, resource Allocation*. New York: McGraw-Hill press.

Shi, Z.F. 2012. *Bayesian network theory and its application in the military system*. Beijing: Defense Industry Press (in Chinese).

State Oceanic Administration. 2010. *Sea ice grade standard*. Beijing: National Marine Environment Forecast Center.

Sun, S., and P.J. Shi. 2012. Risk assessment of sea ice disaster in the Bohai Sea and the northern part of the Yellow Sea. *Journal of Natural Disasters* 4: 8–13 (in Chinese).

Wang, S.C. 2010. *Bayesian network learning, reasoning and application*. Shanghai: Lixin Accounting Publishing House (in Chinese).

Wang, W. 2016. *Object-oriented Bayesian network and its application in risk assessment*. Nanjing: Nanjing University of Aeronautics and Astronautics (in Chinese).

Yang, H.T., and S.Z. Tian. 1994. *Compilation of forty years of marine disasters in China*. Beijing: Ocean Publishing House (in Chinese).

Yang, L.Z., and R. Zhang. 2014. Security risk assessment of China's maritime energy strategy channel based on cloud model. *Military Operations and Systems Engineering* 28(1): 74–80 (in Chinese).

Yang, T., and X. Yang. 1998. Fuzzy comprehensive assessment, fuzzy clustering analysis and its application for urban traffic environment quality evaluation. *Transportation Research Part D Transport and Environment* 3(1): 51–57.

Yang, X.R. 2015. *Research on risk assessment of ship crash bridge based on Bayesian network*. Chongqing: Chongqing Traffic University (in Chinese).

Yu, L.Y. 2017. Discussion of uncertainty risk theory. *Modern Occupational Safety* 3: 75–77 (in Chinese).

Yuan, B.K., K.C. Guo, and X.Y. Wang. 2013. Preliminary study on single factor sea ice disaster index system and sea ice disaster classification method in China. *Ocean Forecast* 30(1): 65–70 (in Chinese).

Zhang, R. 2013. *Climate change and national ocean strategy: Impact and risk assessment*. Beijing: Meteorological Press (in Chinese).

Zhang, E.Y. 2015. *Risk assessment of port ship oil spill based on fuzzy Bayesian network*. Dalian: Dalian Maritime University (in Chinese).

Zheng, W., and Y. Hu. 2009. Grey evaluation method of knowledge management capability. In *Proceedings of 2009 second international workshop on knowledge discovery and data mining*, 23–25 January 2009, Moscow.