# Microblog Rumor Prediction Model and Analysis Based on Decision Tree

Yingqian Qiu

School of Economics and Management, Southwest Jiaotong University, Chengdu, 610031, China
`swjtu_qyq@163.com`

**Abstract.** The spread of rumors on Weibo seriously threatens social stability and personal interests, especially in public health events. Therefore, it is crucial to effectively identify and respond to rumors. This study uses machine learning methods to integrate and analyze Weibo content characteristics, dissemination characteristics, and user characteristics to construct a decision tree classification model. Since Weibo needs to increase its early warning capabilities and supervision of rumors, the optimization of this model pays special attention to the accurate prediction of rumors, that is, it pays more attention to the recall rate of rumors. After optimization, the decision tree model performs more accurately in rumor prediction, with the recall rate of rumors increasing to 0.88, and the precision of both categories has increased, with the rumor category increasing to 0.86 and the non-rumor category increasing to 0.82, showing higher recognition accuracy and reliability. Data analysis and model practice show that the consideration of comprehensive features helps to improve the accuracy of rumor identification. At the same time, the management suggestions put forward provide guidance for social platforms and regulatory authorities to cope with the challenges of rumor propagation.

**Keywords:** Rumors, decision trees, classification predictions, data analysis, management advice.

## 1    Introduction

Bakshy et al. found that user characteristics, especially the influence of users and the number of followers, contribute to improving prediction accuracy[5]. Enayet et al. focused on text features, manually extracting extensive text and user features from rumor data, and found that the LR model achieved the best results[6]. Shi Kaiwen et al. extracted text features and constructed a rumor identification model based on BiLSTM and deep neural networks from the perspective of content features[2]. Guo Miao et al. addressed the aggregation issue of related short texts, enabling them to be studied and discussed as longer texts[1]. Yin Chuanjuan designed a text clustering algorithm based on sudden topics, considering the characteristics of Weibo text data such as short length, poor language regularity, high noise, and large volume[4]. Lastly, Yang et al. found that prop-

agation features such as replies and reposts significantly affect the credibility assessment of Weibo[7]. Sun Ran et al. used a Transformer-CNN model to extract semantic features from Weibo text, integrating user, time, structure, and propagation dimensions to build a rumor identification model for sudden public health emergencies[3].

The Weibo platform is currently one of the most widely used platforms for public information retrieval and expression of opinions. This study defines rumors as false claims and employs Baidu API for sentiment classification of Weibo text content. Current research lacks multidimensionality in variable usage, and deep learning models require extensive parameters and computational resources, leading to high costs for training and inference. Therefore, focusing on a dataset of public health events on Weibo, this paper utilizes machine learning methods to construct a classification model that integrates analysis of content features, propagation characteristics, and user attributes. It aims to uncover patterns and features of rumor propagation, offering feasible solutions and decision support for rumor prevention and management.

## 2      Research Design and Algorithm Introduction

### 2.1     Study Design and Data Acquisition

**Research design.**

In the current development environment, this paper takes the rumor data in the field of public health on the Weibo platform as the starting point to construct a prediction model for Weibo rumors. First, the text content is processed by sentiment classification, and data preprocessing such as merging data sets, coding and standardization is performed; secondly, the data set is subjected to descriptive statistics, correlation analysis, variable impact analysis and other data visualization; finally, the model is constructed and optimized, and the model training and evaluation work such as validity test and factor importance verification is carried out. Finally, based on the conclusions, management suggestions are put forward from the three perspectives of supervision, model and user, and finally, the future prospects are put forward based on the shortcomings of this paper. The model improves the effectiveness and timeliness of Weibo supervision. The technical route studied in this paper is shown in Figure 1.
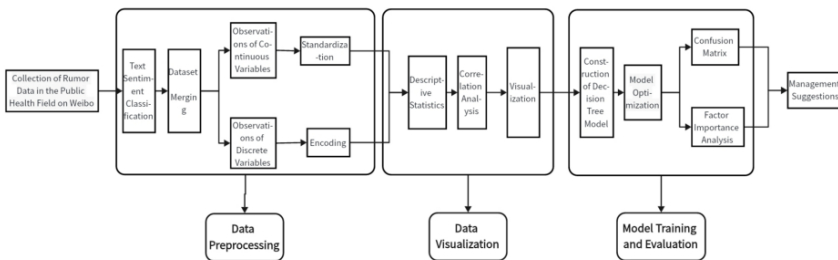


**Fig. 1.** Technology Roadmap

**Data collection.**

This report downloaded data related to rumors and their characteristics in public health from the Alibaba Cloud Tianchi website, and used the Baidu API for sentiment analysis of Weibo content, collecting a total of 3146 pieces of data. In terms of the features of the dataset, it mainly includes Weibo user features (ID, number of Weibo posts, number of followers, number of fans, personal introduction, authentication status), content features (special symbols, probability of positive emotions, probability of negative emotions, confidence, and emotional classification), and Weibo propagation features (number of likes, number of comments).

## 2.2    Decision Tree

Decision trees simulate the logical process of human decision-making and predict the category or value of new instances based on a series of rules. When classifying data, the criterion is the characteristics of the data, and the final decision is the category of the data. This paper uses the Gini index as an indicator of feature division. The lower the Gini index, the higher the purity of the data set and the smaller the uncertainty of the sample set. In the classification problem, assuming there are n classes and the probability that a sample point belongs to the kth class is pk, the Gini index of the probability distribution is defined as:

$$Gini(p) = \sum_{k=1}^{n} p_k(1 - p_k) = 1 - \sum_{k=1}^{n} p_k^2 \tag{1}$$

# 3      Result Analysis

## 3.1    Decision Tree Algorithm

The differences among users in terms of the number of Weibo posts, followers, fans, likes, and comments are substantial, indicating variability in sample users' influence and interaction capabilities. The average confidence level is approximately 0.69, suggesting a high confidence in judging sentiment probabilities. The average probabilities of positive sentiment and sentiment inclination are 0.55 and 1.14, respectively, indicating that the sentiment tendency of sample Weibo posts tends to be positive. The average value for whether the post contains special characters is about 0.56, indicating that most Weibo content includes special characters. The average values for personal verification and having a personal introduction are 0.49 and 0.84, respectively, showing that a significant number of users in the sample have personal introductions but are not verified.

## 3.2    Correlation Analysis

The number of comments and likes exhibit a very strong positive correlation, with a correlation coefficient as high as 0.97. The absolute correlation values between positive probability, negative probability, and sentiment inclination are as high as 0.94. There is a relatively high positive correlation between the number of followers and the number of Weibo posts, with a correlation coefficient of 0.53. User verification status shows a

positive correlation with the number of Weibo posts, with a correlation coefficient of 0.41. The feature with the highest correlation to the "abnormal" label is confidence, with a correlation coefficient of -0.26. Other features have minimal correlation with the label; therefore, in subsequent classification modeling, feature combination will be employed to enhance their relevance. The results are shown in Figure 2.

### 3.3    Analysis of the Impact of Weibo Posts, Followings, and Followers on Rumor Dissemination

As shown in Figure 3, the median and upper quartile of the total number of Weibo posts by bloggers identified as spreading rumors (R) are lower than those of bloggers who post non-rumors (NR), indicating that among users with a relatively large number of total Weibo posts, non-rumor spreaders are more prevalent, while rumor spreaders are more concentrated among those with a smaller number of total Weibo posts, indicating lower user activity. Similarly, rumors typically originate from users with fewer followers and followings, suggesting that rumor propagation is related to the social network structure and influence of users.
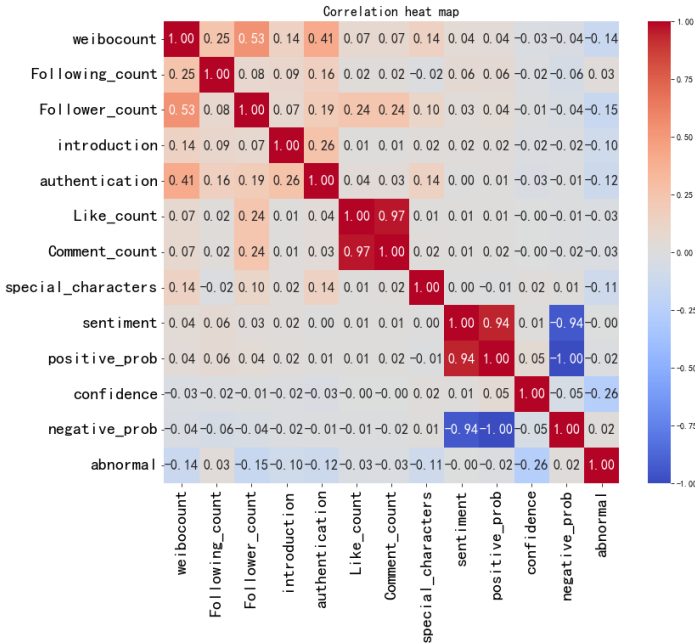


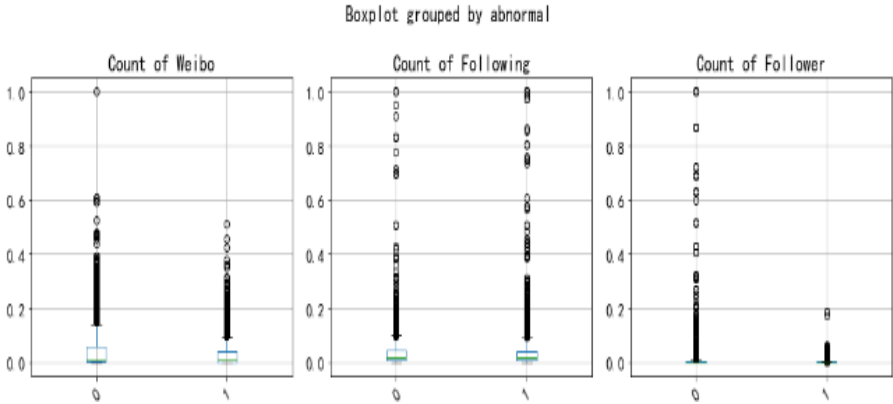**Fig. 2.** Correlation heat map visualization results

**Fig. 3.** Box plot visualization results

## 3.4    Model Validity Test

The spread of rumors can have serious negative consequences, making it crucial to ensure the accuracy and timeliness of regulatory efforts for maintaining a healthy order on Weibo. Therefore, this model focuses specifically on accurate rumor prediction, particularly emphasizing the recall of class 1. The study utilized cross-validation and grid search to optimize parameters of the decision tree model, particularly aiming to improve recall for rumor prediction, ensuring the model captures as many true positive samples (rumors) as possible. The optimal parameters identified through this process were: minimum samples required to split a node = 20, minimum samples required at each leaf node = 4, and maximum tree depth = 5.

Using these optimized parameters to maximize recall, a new decision tree model was constructed, and the resulting confusion matrix is depicted in Figure 4.
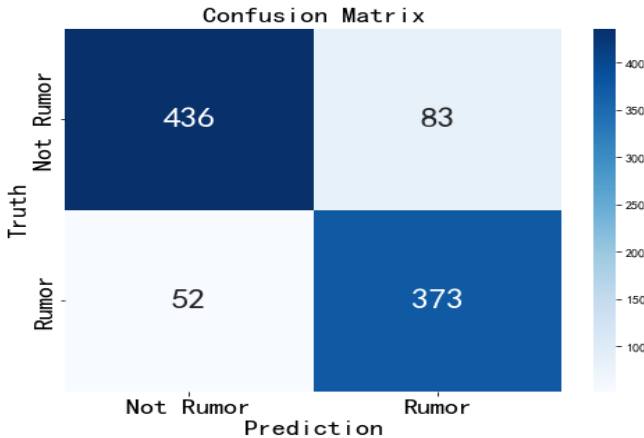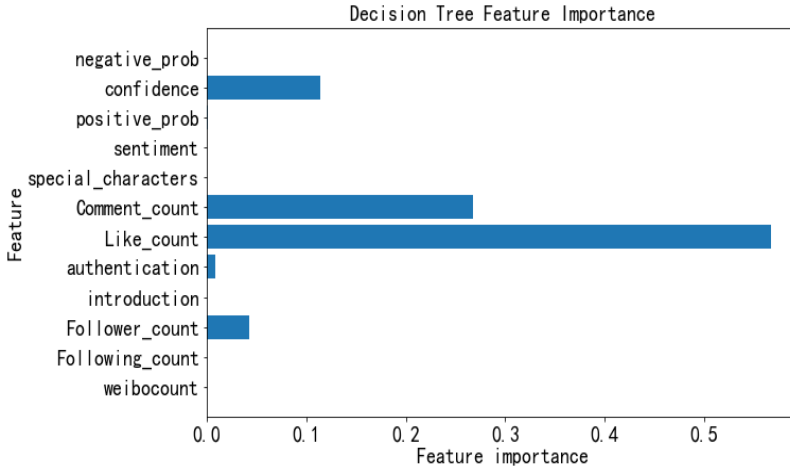


**Fig. 4.** Confusion Matrix

**Fig. 5.** Factor Importance

Visible in the confusion matrix are reduced false negatives and increased true negatives. Following optimization, the recall for rumors (class 1) improved to 0.88, indicating more rumors were correctly detected, reducing instances where rumors were mistakenly classified as non-rumors. Simultaneously, the average precision for both rumor and non-rumor categories improved to 0.86 and 0.82, respectively, suggesting enhanced model accuracy in classification and reduced misclassification rates.

### 3.5    Factor Importance

As shown in Figures 5, the importance of the number of likes is greater than 0.5, the number of comments is close to 0.3, the confidence is slightly greater than 0.1, the number of fans and personal certification are both less than 0.05, and the other indicators are all 0. Therefore, in the rumor prediction model constructed in this paper, the number of likes and comments is the most critical factor, followed by emotional confidence, while the number of fans and personal certification have relatively little influence.

## 4    Conclusions

In this paper, the rumor prediction model established identifies emotional credibility as the most significant factor among content features, highlighting the importance of Weibo content quality to users. Additionally, in terms of propagation characteristics, the number of likes and comments emerges as crucial factors, especially the role of likes in the dissemination of rumors. Furthermore, concerning user characteristics, the number of followers and personal verification are deemed important, underscoring their impact on rumor propagation and amplification.

Integrating the developed decision tree model into a real-time rumor monitoring system involves continuously monitoring content posted by bloggers, their features, and propagation characteristics for timely rumor identification. This integration aims to reduce the burden of manual review, promptly identify and manage potential misinformation. Moreover, establishing a model-based risk alert system enables the automatic detection of potential rumor risks in content posted by bloggers, issuing timely alerts and warnings to users or temporarily restricting content dissemination. Additionally, leveraging model outputs, educational efforts can target content potentially containing rumors to improve media literacy. Concurrently, governmental agencies can use monitoring results to promptly debunk misinformation, promote rational public information consumption, and maintain societal stability. Furthermore, businesses can utilize the model to avoid advertising within rumor-propagating content, safeguarding brand reputation and mitigating negative publicity.

## 5    Future Outlook

Future research hopes to expand the scale and diversity of datasets, collecting more extensive and varied data across different fields, regions, and social groups, to enhance the generalizability and adaptability of models. This will enable more accurate identification and response to various types of rumor propagation phenomena. Concurrently, leveraging computer vision and video analysis technologies to develop rumor detection algorithms tailored for image and video content will comprehensively cover multiple forms of dissemination.

## References

1. Guo M, Jiao Y. Predictive Analysis of Weibo Information Forwarding under the Background of Internet Public Opinion Communication and Evolution. Journal of Intelligence, 2016, 35(5): 7. DOI: 10.3969/j.issn.1002-1965.2016.05.009.
2. Shi K, Liu K. Identification of Weibo Rumors in Sudden Public Health Events. Library and Information Service, 2021, 65(13): 9. DOI: 10.13266/j.issn.0252-3116.2021.13.009.
3. Sun R, An L. Study on Rumor Identification in Sudden Public Health Events. Information and Documentation Work, 2021, 42(5): 8.
4. Yin C. Research on Detection of Weibo Hot Topics Based on Clustering Analysis [Dissertation]. Jiangsu: Nanjing University of Technology, 2012.
5. Bakshy E, Mason WA, Watts DJ. Everyone's an influencer: quantifying influence on Twitter [C]//Web Search and Data Mining. ACM, 2011. DOI: 10.1145/1935826.1935845.
6. Enayet O, El-Beltagy SR. NileTMRG at SemEval-2017 Task 8: Determining Rumour and Veracity Support for Rumours on Twitter [C]//Proceedings of the 11th International Workshop on Semantic Evaluation (SemEval-2017). 2017. DOI: 10.18653/v1/S17-2082.
7. Yang F, Yu X, Liu Y, et al. Automatic detection of rumor on Sina Weibo [C]//Knowledge Discovery and Data Mining. ACM, 2012. DOI: 10.1145/2350190.2350203.