http://www.sjst.psu.ac.th

*Original Article*

# A novel BNB-NO-BK method for detecting fraudulent crowdfunding projects

Qi Li, and Jian Qu*

*School of Engineering and Technology, Panyapiwat Institute of Management,
Pak Kret, Nonthaburi, 11120 Thailand*

## Abstract

Identifying fraudulent campaigns or messages remains a difficult task in the field of natural language processing. We proposed a hypothesis that if the well-known brands in the same category as the crowdfunding project can implement similar technology as crowdfunding projects, the project is considered to be more feasible. The opposite is considered more likely to be fraudulent. This research proposed a novel BNB-NO-BK method to detect fraudulent crowdfunding projects. A novel method called BNB, which was constructed by key-BERT, NLTK, and fine-tuned QA model for BERT, was proposed to extract the characteristics of crowdfunding projects. We proposed a novel NO (Nice Classification & Ontology) method for classifying the categories of projects, which constructed ontology trees based on the characteristics of the crowdfunding projects and our modified Nice Classification. Furthermore, we proposed a novel BK (Brand Knowledge) cross-checking method to extract the features of crowdfunding projects. Finally, we compared the performances of different machine learning methods for identifying fraudulent crowdfunding projects. Furthermore, to address the problem of possible bias caused by unbalanced data, we used data augmentation to process the dataset. Our proposed method achieved an accuracy of 95.71% in detecting fraudulent crowdfunding projects, which was superior to existing methods.

**Keywords**: crowdfunding projects, fine-tuneBert QA model, ontology, brand information retrieval, machine learning classifier

## 1. Introduction

Crowdfunding allows entrepreneurs to establish or boost their startups. A project with new ideas can be quickly realized through crowdfunding. However, crowdfunding platforms have no oversight on launched campaigns. Many entrepreneurs attempt to build fraudulent crowdfunding projects to gain personal profit. Fraudulent crowdfunding projects are divided into two types: the first is based on logically feasible and practical concepts. After the project sponsors raise funds, they use other reasons to announce the failure of the project in order to embezzle the project funds. This situation is related to the credibility of the sponsors. For example, The Doom That Came to Atlantic City was a game project that raised over $122,000. Unfortunately, the sponsors did not use the money raised to develop the game but to pay the rent, organize events, or do other things.

The second is also based on being logically feasible but impossibly to realize due to technical limitations. To obtain funds, the project sponsor used some fake pictures or videos to deceive the public. For example, artificial gill respirator –Triton. Jeabyun Yeon, the initiator of this project, claimed that this artificial gill respirator can separate oxygen from the water with micro‐battery power. Moreover, it allows divers to stay underwater for 45 minutes at a depth of 4.5 meters. To this end, the initiators also released a fake video of a diver diving underwater with this respirator. Therefore, once this artificial gill respirator was released, it received a lot of sponsorship and support from divers and gained more than $800,000 in just a few days.

However, marine biologist and diver Alistair Dove discussed this issue at Deep Sea News. Dr. Dove calculated that a person needs about 532.8 mg of oxygen per minute to breathe. To provide that amount of oxygen, the Triton is required to process approximately 90L of water per minute

---

*Corresponding author
Email address: jianqu@pim.ac.th

with 100% efficiency. This would require a larger pump than the Triton itself, but the Triton does not have any equipment to move water through the device. Not only that, storing and compressing oxygen also would require specialized equipment and energy that Triton could not realistically provide anyway. This article caused a sensation and the divers who supported and sponsored the project recognized that the project was a pipe dream. These diving enthusiasts or professional divers have basic knowledge of diving, yet they were still being scammed by this project. Identification of fraudulent information in crowdfunding projects is difficult even for humans and thus detection of fraudulent and fake information in crowdfunding projects is a very challenging task. Because of this, the existing methods for identifying fake information mainly focus on identifying fake information in the field of news on social media such as Facebook, Twitter, and YouTube (Shu, Sliva, Wang, Tang, & Liu, 2017).

The existing methods for identifying fake information in the news can be divided into three groups. These include fake information detection based on content, fake information detection based on social context, and fake information detection combined with external knowledge.

Content-based methods for detecting fake information include graphic content and text content of articles (Kochkina, Liakata, & Zubiaga, 2018; Liao et al., 2021; Qi, Cao, Yang, Guo, & Li, 2019; Vaibhav, Mandyam, & Hovy, 2019). In this case, image processing techniques are applied. Qi et al. (2019) proposed a fake image discriminator MVNN. The MVNN discriminator extracted the spatial and frequency domain features of an image, with the purpose to determine whether the image is manipulated by retouching software (Qi et al., 2019). However, for this research, there are two general types of image information in technological crowdfunding projects. One type is the design construction of the project, as design drawings are made using software anyway, so this method could not determine the authenticity of the images. The second type is object picture. The crowdfunding projects, even the fraudulent ones, did use authentically photographed pictures. The project sponsor might build a device that "looked" real, but there was no guarantee that the device would have functions and features promised by the project. Therefore, the method for detecting fake information in the news using graphic content, which identifies whether a crowdfunding project is fraudulent by detecting whether the image has been synthesized, has a relatively low accuracy rate. Methods based on text content can solve this problem.

Kochkina et al. (2018) worked on methods for detecting fake news based on text content. They proposed to use deep learning techniques to extract text features as an embedded representation of news texts. The embedding vectors were entered into the classifier for classification to obtain results of fake news detection. Some methods that did not apply directly to fake news or fraudulent information detection might be used for such tasks. For example, BERT. Devlin, Chang, Lee, and Toutanova (2019) proposed a BERT-based method which was highly generalizable and can handle multiple tasks. However, when applying BERT to this study, the extracted text features of crowdfunding projects were mostly long sentences that did not effectively represent the text information. BERT model was applied to our crowdfunding project and achieved 51.98% precision and 68.40% f-measure. Therefore, we proposed a novel method called BNB, which is constructed by BERT, NLTK, and fine-tuned QA model for BERT to extract characteristics of crowdfunding projects. Our method can extract the feature information of the crowdfunding project where the existing Bert model fails.

The existing method proposed by Perez, Machado, Andrews, and Kourtellis (2020) detected fake charitable crowdfunding projects based on text and image features. It achieved an accuracy of 90.14% in the classification of charitable crowdfunding projects. In the process of the experiment, the method identified the emotion of the portrait in the image and analyzed the emotion, named entities, and word importance in the text. Still, the difference between technological crowdfunding projects and charitable crowdfunding projects is that the images used in technological crowdfunding projects are design drawings or physical images. Therefore, we cannot identify any emotional features from the images to verify crowdfunding projects, and these features of the text have little effect in helping us detect fake crowdfunding projects. Hence, the method proposed by Perez et al. (2020) when applied to our research only achieved an accuracy rate of 21.32%.

The second group detected fake information using social media and user credibility. These methods required additional social information and actual diffusion behavior because they rely on the credit rating of users (Dou, Shu, Xia, Yu, & Sun, 2021; Jiang, Chen, Zhang, Chen, & Liu, 2019; Lu & Li, 2020). Jiang et al. (2019) proposed a method to detect fake news based on social media and user credibility. They modeled the network of news dissemination and the social network of users, and classified the node information by the classifier to achieve fake news detection. However, for crowdfunding projects, there is little direct information on social media between project initiators and crowd-funders. Furthermore, the actual diffusion degree of crowdfunding projects and the actual diffusion behavior of crowd-funders cannot be detected. Therefore, this method is limited to fake news detection with more interaction and dissemination behavior, thus resulting in a low recall.

The third group of methods is to detect fake information in the news by combining external knowledge. These methods require the construction of knowledge graphs made of objective facts of relevant news and news content to detect fake information (Bondielli & Marcelloni, 2019; Pickering, 2001; Qu, Nguyen & Shimazu, 2016). Bondielli et al. (2019) proposed an integrated model that combined pre-trained models and statistical features. They used the presence of attributes in news or tweets as statistical features to classify fake news. However, crowdfunding projects are generally novel and innovative projects, and the statistical features of such projects do not appear as often as those in the news. Therefore, the accuracy of the method, which detects fake crowdfunding projects by fusing statistical features, is not high. Pickering (2001) proposed a method to detect fake news by fusing the descriptive information of entities in common sense, the topic information in the news, and the news content. However, for some whimsical and fake crowdfunding projects, the knowledge graph that fuses the common sense information does not detect the fraudulent information in the projects. For example, the crowd-funders who participated in the Triton are some diving enthusiasts or even professional

divers. Before being able to dive at all, divers need to undergo professional training and be well-versed in diving skills and knowledge. Yet, they were still unable to sell Triton as a fake project. Thus, we proposed a novel BK cross-checking method to solve this problem.

As technology projects on crowdfunding platforms are generally innovative, such start-up crowdfunding projects contain a lot of technology, but whether the technology can be achieved is difficult to detect. Therefore, we proposed a hypothesis that if a well-known brand in the same category as the crowdfunding projects can implement similar technology as the crowdfunding project, then we consider the project to be more feasible and less likely to be a fraudulent project. Otherwise, it would be considered more likely to be a fraudulent project.

In order to apply our hypothesis, we proposed a novel hybrid method to detect fake information in crowdfunding projects. We proposed a novel method called BNB (BERT-NLTK-BERT's QA) to extract characteristics of the crowdfunding project, which is constructed by BERT, NLTK, and fine-tuned QA model for BERT. Secondly, we used the characteristics of projects to construct a novel ontology tree based on our modified Nice Classification. Then, a method for measures of association rule (QU *et al.*, 2012) from web-data was proposed to retrieve and count the brand knowledge (BK) for categories of crowdfunding projects. Finally, we used the cross-checking method to check the features of crowdfunding projects, which are web-data association results, to predict the feasibility of crowdfunding projects.

## 2. The Approach

In this section, a novel method to detect fake crowdfunding projects is discussed. The flowchart of our method is shown in Figure 1. The flowchart with example of our method is shown in Figure 2.

### 2.1 Information retrieval of crowdfunding projects

There are two sources to obtain information of crowdfunding projects. The first method is to obtain crowdfunding project information from web-pages. This method has low recall because there is significant noise in the information fragments returned from web-pages. The second method is to obtain crowdfunding project information from official website. We used the second approach to retrieve crowdfunding project information from the official website of that project, which had higher accuracy.

### 2.2 Characteristics extraction of crowdfunding projects

In the characteristics extraction process of crowdfunding projects, we extracted keywords from the text information of crowdfunding projects as the characteristics information. We applied the embedding model of BERT to extract keyword characteristics on our dataset. However, the experimental results show that the keyword characteristics extraction accuracy of the embedding model of BERT is low. The result of characteristics extraction by embedding model for BERT is shown in section 3.3. In our research, BERT-1 is

a semantic similarity model, and the model is called "distilbert-base-nli-mean-tokens". BERT-2 is multilingual knowledge distilled version 1 of the multilingual universal sentence encoder, and the model is called "distiluse-base-multilingual-cased-v1". BERT-3 is multilingual knowledge distilled version 2 of the multilingual universal sentence encoder, and the model is named "distiluse-base-multilingual-cased-v2". BERT-4 is paraphrase recognition model, and the model is named "paraphrase-distilroberta-base-v1". BERT-5 is multilingual version of "paraphrase-distilroberta-base-v1" and the model is called "paraphrase-xlm-r-multilingual-v1" (Reimers & Gurevych, 2019).

To improve the accuracy of characteristics extraction for crowdfunding projects, we proposed a novel method called BNB (BERT-NLTK-BERT's QA) that extracted characteristics from the text of crowdfunding projects. The characteristics of projects extracted by the first BERT model and NLTK helped us understand the meaning of the project information, which enabled us to construct sentence pairs of questions and answers. The question and answer sentence pairs were used as training sets to retrain the QA model for BERT. The fine-tuned QA model for BERT extracted characteristics of crowdfunding projects.

### 2.3 Classification of crowdfunding projects

This section introduces methods for the classification of crowdfunding projects. We proposed a method called NO, which constructed an ontology tree based on modified Nice Classification. We first obtained the tendency of categories for crowdfunding projects based on the extracted characteristics of projects. Since most technological crowdfunding projects are novel or even seminal, the original Nice Classification cannot find the category of the project we need. Therefore, we modified the Nice Classification in combination with the characteristics of crowdfunding projects. The original Nice Classification contained 45 first-level classes, 136 second-level classes, and 10,606 third-level classes. We modified Nice Classification into a form containing 35 first-level classes, 115 second-level classes, and 9,566 third-level classes (Qi & Qu, 2022). Lastly, the ontology tree was constructed by the categories of crowdfunding projects. The categories of some projects are displayed in the ontology tree, as shown in Figure 3.

### 2.4 Measures of association rule on web-data for crowdfunding projects

We proposed a BK (Brand Knowledge) method with measures of association rule from web-data, which associates brand official website, shopping site, characteristics of crowdfunding projects, and categories of crowdfunding projects to measure the characteristics of crowdfunding projects.

The principle of our method is that it checks the technology of the proposed Kickstarter project and Indiegogo project against the top-rated brand name company in the same product category. For example, CST-01 is an 0.80 mm thin flexible wristwatch. The micro-energy cell (MEC) that powers the watch can charge in 10 minutes and the watch is able to function for a month between charges. However, the famous brand names in the smartwatch industry, such as Apple,

Figure 1.　The flowchart of detecting fake crowdfunding projects

Figure 2.   The flowchart with an example of detecting fake crowdfunding projects

Samsung, Huawei, and Xiaomi, cannot achieve this technology. We believed that if the top-ranked brand can produce similar products in the same product category, then the crowdfunding project is more feasible and less likely to have fraudulent information, and vice versa. The top-ranked brand usually has a more capable team and more money than the crowdfunding initiator. If a technology that is difficult even for teams with large manpower, material resources and financial power, then it would be much more difficult for the crowdfunding initiator to implement.

Figure 3.   The ontology tree for the categories of crowdfunding projects

In addition, in our preliminary research, we found that if the characteristics of a crowdfunding project can be retrieved on a shopping site, the probability that the project may be implemented is greatly increased, and vice versa. Therefore, we used the results of retrieving project features in shopping sites as one of the features for detecting fraudulent projects.

## 2.5 Machine learning

We proposed a novel cross-check method to detect fraudulent crowdfunding projects. We obtained the features of the projects based on the results of association rule measures from web-data. The results of detection were taken as a feature. These features can be used to detect fraudulent crowdfunding projects with machine learning classifiers. Thus, the feasibility of fraudulent crowdfunding projects is predicted.

We used different algorithms to classify fraudulent crowdfunding projects, including Decision Tree (DT) (Boonchuay, Sinapiromsaran, & Lursinsap, 2017), k-Nearest Neighbor (KNN), Artificial Neural Network (ANN), Random Tree (RF), Support Vector Machine (SVM), Logistic Regression (SVM) (LR(SVM)), and Logistic Regression (LR). Furthermore, we also optimized the model parameters with evaluation, we used Backward Elimination to select features, and used AdaBoost to boost the classifier.

## 3. Results and Discussion

### 3.1 Dataset

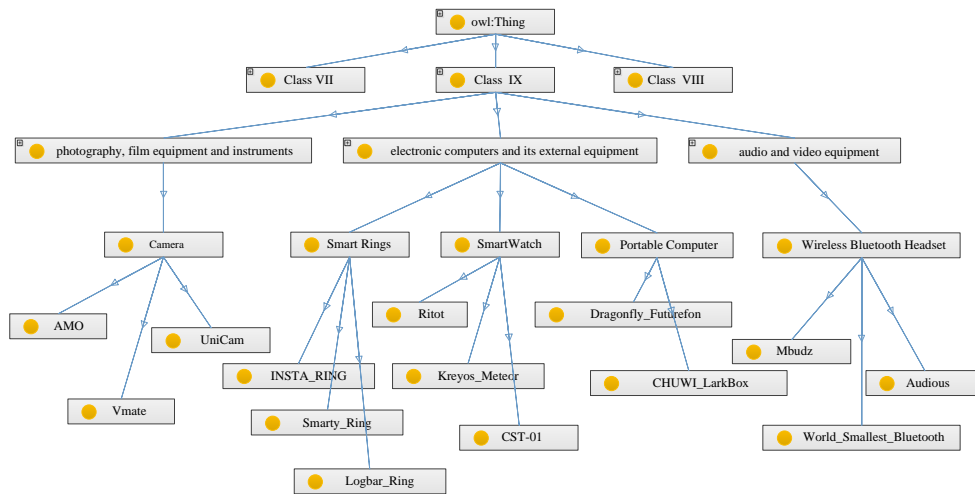The dataset I of this research is based on two crowdfunding platforms, Kickstarter and Indiegogo. Dataset I has a total of 70 crowdfunding projects in the technological category with 246,517 parameters, which contains 20 fake crowdfunding projects and 50 genuine crowdfunding projects. A flowchart describing the amount of data for the research is shown in Figure 4.

In our dataset, fraudulent technology crowdfunding projects and genuine crowdfunding projects cannot be balanced due to the high level of invisibility of fraudulent technology crowdfunding projects. Therefore, we used data augmentation methods to process the dataset. We built Dataset II to reduce the possible bias caused by the unbalanced dataset. We processed our dataset using synonym replacement methods in data augmentation. The synonym replacement approach is based on TF-IDF replacement (Xie, Dai, Hovy, Luong, & Le, 2020). In our study, information with low TF-IDF probability is unimportant and we can perform a replacement on it. Therefore, we replaced words with low TF-IDF scores without affecting the underlying meaning of the sentences. Our experiment utilized TF-IDF to calculate the word frequency of the 20 fake crowdfunding projects in our dataset and get the scores. Moreover, we generated a new dataset by replacing words in the original dataset based on synonyms in WordNet. The newly generated dataset II contains 50 fake crowdfunding projects and 50 genuine crowdfunding projects.

We also preprocessed the data before applying it to the model. During the preprocessing of the dataset, we performed a data cleansing operation on the data. Meta data were removed using regular expressions. The process and an example are shown in Figure 5.

### 3.2 Baseline comparison to existing method

In this section, we explained and compared the baseline, result of our method, and the result from an existing paper. In this research, we used the accuracy of human identification as a baseline. Since many fake information of crowdfunding projects are difficult or even impossible to identify artificially, the accuracy of this baseline was 75%. The accuracy of the method proposed by this research to detect fake crowdfunding projects achieved 95.71%. The method proposed by Perez et al. (2020) when applied to our research only achieved an accuracy of 21.32%.

### 3.3 Project text feature extraction results

In this section, we compared the accuracy of the fine-tuned QA model for BERT and the BERT embedding

Figure 4.   A flowchart describing the amount of data for the research

model for extracting characteristics. The accuracy of characteristics extraction with fine-tuned QA model is evaluated by F1-score and EM-score. The F1-score is 45.30%, and the EM-score is 32.25%. The accuracy that the text feature extraction by fine-tuned QA model achieved was 92.38%. The accuracy of characteristics extraction with the BERT embedding model is evaluated by precision, recall, and F1-score, as shown in Figure 6.

## 3.4 Results of Nice Classification

According to the original Nice Classification standard, we can only identify the categories of 59 projects. There were 11 projects whose categories cannot be identified (Qi *et al*., 2022). From our modified Nice Classification, we can identify categories of all projects in our dataset.

## 3.5 Results of machine learning training features

In this section, we compare the results of different algorithms for classification, including DT, RF, KNN, ANN, SVM, LR, and LR(SVM) algorithms. We compared classifiers with Backward Elimination and AdaBoost. The results of feature selection by Backward Elimination are shown in Table 1. The description of Feature-ID in Table 2 is shown in Table 1. In addition, the results of the meta-level classifier and the classifier with evaluation optimized parameters are also compared. The accuracies of classification for the different algorithms are shown in Table 3. The confusion matrix is shown in Table 4.

Figure 5. A flowchart for preprocessing of data and example



Figure 6. The results of characteristic extraction by embedding model for BERT

## 4. Discussion

In the original Nice Classification, there were 11 projects whose categories cannot be determined. Some of the categories of crowdfunding projects that could not be retrieved in the original Nice Classification were innovative, such as "bottomless trash cans" and "solar highways". The original Nice Classification was wider (containing only the general category of such items). However, most Kickstarter project requires a more refined range of categories. For

example, of "Wireless Mechanical Keyboard", only "keyboard" is displayed in the original Nice Classification.

There are five projects that are not retrievable because there is no refined classification. In addition, it will not be retrieved because the synonyms of the word are retrievable, such as "underwater respirator" in the Nice Classification is "diving respirator". There are four projects that were not retrieved because the synonyms were not retrievable.

For fake crowdfunding project classification, we compared the effectiveness of some classification algorithms in an existing classification task for fake crowdfunding projects. We found that DT, RF, KNN, and ANN algorithms generally work much better than SVM, LR, and LR(SVM). The DT, RF, KNN, and ANN classification algorithms can correctly predict all the fake crowdfunding projects. However, they predict some real projects as fake ones, these incorrect predictions being projects ID40, ID42, ID47, ID54, and ID56. We found that this was because the product features of these five projects had neutral feature scores in the process of network knowledge feedback. For example, feature 1 can prove that the product is real with a score of 3, but feature 2, feature 3, and feature 4 can prove that the product is fake with a score of 1 each. The classifier will generate an error in this case.

Table 1. Comparison of results for Backward Elimination selected features

| Feature ID | ANN | | | | RF | | | | DT | | | | KNN | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | B | BA | OB | OBA | B | BA | OB | OBA | B | BA | OB | OBA | B | BA | OB | OBA |
| 1 | √ | √ | √ | √ | ✗ | √ | √ | √ | √ | √ | √ | √ | √ | √ | √ | √ |
| 2 | √ | √ | √ | √ | √ | ✗ | ✗ | √ | ✗ | √ | ✗ | √ | √ | √ | √ | √ |
| 3 | ✗ | ✗ | √ | √ | √ | √ | √ | √ | √ | √ | ✗ | √ | √ | √ | √ | √ |
| 4 | √ | √ | √ | √ | √ | √ | √ | √ | √ | √ | ✗ | √ | √ | √ | √ | √ |
| 5 | √ | √ | √ | √ | √ | √ | √ | √ | √ | √ | ✗ | √ | √ | √ | √ | √ |
| 6 | √ | √ | √ | √ | √ | √ | ✗ | √ | √ | √ | √ | √ | ✗ | √ | √ | √ |
| 7 | ✗ | √ | √ | ✗ | √ | √ | ✗ | √ | √ | ✗ | √ | ✗ | √ | √ | √ | √ |
| 8 | √ | √ | √ | √ | √ | √ | √ | √ | √ | √ | √ | √ | √ | √ | √ | √ |
| 9 | √ | √ | √ | √ | √ | √ | √ | √ | ✗ | √ | √ | √ | √ | √ | √ | √ |
| 10 | √ | √ | √ | √ | √ | √ | √ | √ | √ | √ | ✗ | √ | √ | √ | √ | √ |
| 11 | √ | √ | √ | √ | √ | ✗ | ✗ | √ | √ | √ | √ | √ | √ | ✗ | √ | √ |
| 12 | √ | √ | ✗ | √ | √ | √ | √ | √ | √ | √ | √ | √ | √ | ✗ | √ | √ |
| 13 | √ | √ | √ | √ | √ | √ | √ | √ | √ | √ | √ | √ | √ | √ | √ | √ |
| 14 | √ | √ | √ | √ | √ | √ | √ | ✗ | √ | √ | √ | √ | √ | √ | √ | √ |

Table 1.　Continued.

| Feature ID | ANN | | | | RF | | | | DT | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | B | BA | OB | OBA | B | BA | OB | OBA | B | BA | OB | OBA |
| 1 | √ | √ | √ | √ | √ | √ | √ | √ | √ | √ | √ | √ |
| 2 | √ | √ | √ | √ | √ | √ | √ | √ | √ | √ | √ | √ |
| 3 | × | √ | × | √ | √ | √ | √ | √ | √ | √ | √ | √ |
| 4 | √ | × | √ | × | √ | √ | × | × | √ | √ | √ | √ |
| 5 | √ | √ | × | √ | √ | √ | √ | √ | × | × | × | × |
| 6 | × | √ | × | √ | √ | √ | √ | √ | √ | √ | × | √ |
| 7 | √ | √ | × | √ | × | × | √ | × | √ | √ | √ | √ |
| 8 | √ | √ | √ | √ | × | √ | × | √ | × | √ | √ | √ |
| 9 | √ | √ | √ | √ | √ | √ | √ | √ | × | √ | √ | √ |
| 10 | √ | √ | √ | √ | √ | √ | √ | √ | √ | √ | × | √ |
| 11 | √ | √ | √ | √ | √ | √ | √ | √ | √ | √ | √ | √ |
| 12 | √ | √ | × | √ | √ | √ | √ | √ | √ | √ | √ | √ |
| 13 | √ | √ | √ | √ | √ | √ | √ | √ | √ | √ | √ | √ |
| 14 | √ | √ | √ | √ | √ | √ | √ | √ | √ | √ | √ | √ |

Note: B: Backward Elimination; BA: Backward Elimination + AdaBoost; OB: Optimize Parameters (Evolutionary)+ Backward Elimination; OBA: Optimize Parameters (Evolutionary)+ Backward Elimination + AdaBoost

Table 2.　The description of feature ID

| Feature ID | Input form | Feature name |
|---|---|---|
| 1 | K AND "C" site: official websites of the brand | SUM1 |
| 2 | "K" AND C site: official websites of the brand | SUM2 |
| 3 | "K" site: official websites of the brand | SUM3 |
| 4 | K site: official websites of the brand | SUM4 |
| 5 | K AND "C" site: shopping URLs | sum1 |
| 6 | "K" AND C site: shopping URLs | sum2 |
| 7 | "K" site: shopping URLs | sum3 |
| 8 | K site: shopping URLs | sum4 |
| 9 | "K" | sum5 |
| 10 | "K" AND "C" site: shopping URLs | sum6 |
| 11 | K AND "C" site: shopping URLs & brand URLs | s1 |
| 12 | "K" AND C site: shopping URLs & brand URLs | s2 |
| 13 | "K" site: shopping URLs & brand URLs | s3 |
| 14 | K site: shopping URLs & brand URLs | s4 |

Table 3.　The accuracy of classification for different algorithms

| | | Meta-level | B | A | BA | OB | OA |
|---|---|---|---|---|---|---|---|
| DT | I | 85.71% | 78.57% | 85.71% | 92.86% | 81.43% | 92.86% |
| | II | 70.00% | 80.00% | 92.00% | 90.00% | 89.00% | 93.00% |
| RF | I | 85.71% | 81.43% | 87.14% | 90.00% | 82.86% | 91.43% |
| | II | 74.00% | 79.00% | 87.00% | 90.00% | 87.00% | 94.00% |
| KNN | I | 81.43% | 82.86% | 81.43% | 85.71% | 82.86% | 87.14% |
| | II | 74.00% | 75.00% | 86.00% | 91.00% | 78.00% | 92.00% |
| ANN | I | 85.71% | 82.86% | 85.71% | 85.71% | 85.71% | 88.57% |
| | II | 73.00% | 76.00% | 81.00% | 85.00% | 79.00% | 90.00% |
| SVM | I | 74.29% | 75.71% | 82.86% | 82.86% | 75.71% | 82.86% |
| | II | 68.00% | 70.00% | 72.00% | 82.00% | 75.00% | 72.00% |
| LR | I | 88.57% | 80.00% | 88.57% | 88.57% | 84.29% | 88.57% |
| | II | 69.00% | 73.00% | 74.00% | 77.00% | 76.00% | 74.00% |
| LR(SVM) | I | 82.86% | 78.57% | 82.86% | 84.29% | 82.86% | 82.86% |
| | II | 65.00% | 68.00% | 66.00% | 70.00% | 70.00% | 66.00% |

Note: I: Dataset I; II: Dataset II; B: Backward Elimination; A: AdaBoost; BA: Backward Elimination + AdaBoost; OB: Optimize Parameters (Evolutionary)+ Backward Elimination; OA: Optimize Parameters (Evolutionary)+AdaBoost; OBA: Optimize Parameters (Evolutionary)+ Backward Elimination + AdaBoost

Table 4. The confusion matrix of classification

| Method | | Meta-level | | A | | B | | BA | | OA | | OB | | OBA | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | True0 | True1 | True0 | True1 | True0 | True1 | True0 | True1 | True0 | True1 | True0 | True1 | True0 | True1 |
| **DT** | | | | | | | | | | | | | | | |
| I | pred.0 | 14 | 4 | 14 | 4 | 10 | 5 | 15 | 0 | 15 | 0 | 10 | 3 | 15 | 0 |
| | pred.1 | 6 | 46 | 6 | 46 | 10 | 45 | 5 | 50 | 5 | 50 | 10 | 47 | 5 | 50 |
| II | pred.0 | 34 | 14 | 38 | 8 | 50 | 8 | 50 | 8 | 47 | 8 | 48 | 5 | 50 | 6 |
| | pred.1 | 16 | 36 | 12 | 42 | 0 | 42 | 0 | 42 | 3 | 42 | 2 | 45 | 0 | 44 |
| **RF** | | | | | | | | | | | | | | | |
| I | pred.0 | 11 | 1 | 12 | 1 | 12 | 5 | 14 | 1 | 15 | 1 | 10 | 2 | 17 | 0 |
| | pred.1 | 9 | 49 | 8 | 49 | 8 | 45 | 6 | 49 | 5 | 49 | 10 | 48 | 3 | 50 |
| II | pred.0 | 34 | 10 | 36 | 7 | 41 | 4 | 45 | 5 | 46 | 9 | 50 | 6 | 50 | 5 |
| | pred.1 | 16 | 40 | 14 | 43 | 9 | 46 | 5 | 45 | 4 | 41 | 0 | 44 | 0 | 45 |
| **KNN** | | | | | | | | | | | | | | | |
| I | pred.0 | 10 | 3 | 10 | 3 | 11 | 3 | 11 | 1 | 12 | 1 | 10 | 2 | 13 | 1 |
| | pred.1 | 10 | 47 | 10 | 47 | 9 | 47 | 9 | 49 | 8 | 49 | 10 | 48 | 7 | 49 |
| II | pred.0 | 41 | 17 | 39 | 14 | 50 | 14 | 50 | 9 | 41 | 13 | 50 | 8 | 50 | 7 |
| | pred.1 | 9 | 33 | 11 | 36 | 0 | 36 | 0 | 41 | 9 | 37 | 0 | 42 | 0 | 43 |
| **ANN** | | | | | | | | | | | | | | | |
| I | pred.0 | 11 | 1 | 11 | 1 | 11 | 3 | 11 | 1 | 14 | 2 | 12 | 2 | 15 | 1 |
| | pred.1 | 9 | 49 | 9 | 49 | 9 | 47 | 9 | 49 | 6 | 48 | 8 | 48 | 5 | 49 |
| II | pred.0 | 30 | 7 | 33 | 7 | 36 | 5 | 38 | 3 | 36 | 7 | 48 | 8 | 50 | 8 |
| | pred.1 | 20 | 43 | 17 | 43 | 14 | 45 | 12 | 47 | 14 | 43 | 2 | 42 | 0 | 42 |
| **SVM** | | | | | | | | | | | | | | | |
| I | pred.0 | 3 | 1 | 10 | 2 | 4 | 1 | 14 | 6 | 10 | 2 | 5 | 2 | 13 | 4 |
| | pred.1 | 17 | 49 | 10 | 48 | 16 | 49 | 6 | 44 | 10 | 48 | 15 | 48 | 7 | 46 |
| II | pred.0 | 32 | 14 | 33 | 13 | 34 | 12 | 38 | 6 | 35 | 10 | 34 | 12 | 38 | 6 |
| | pred.1 | 18 | 36 | 17 | 37 | 16 | 38 | 12 | 44 | 15 | 40 | 16 | 38 | 12 | 44 |
| **LR** | | | | | | | | | | | | | | | |
| I | pred.0 | 13 | 1 | 13 | 1 | 11 | 5 | 13 | 1 | 13 | 1 | 12 | 3 | 13 | 1 |
| | pred.1 | 7 | 49 | 7 | 49 | 9 | 45 | 7 | 49 | 7 | 49 | 8 | 47 | 7 | 49 |
| II | pred.0 | 36 | 17 | 36 | 13 | 36 | 12 | 38 | 11 | 36 | 10 | 36 | 12 | 38 | 11 |
| | pred.1 | 14 | 33 | 14 | 37 | 14 | 38 | 12 | 39 | 14 | 40 | 14 | 38 | 12 | 39 |
| **LR (SVM)** | | | | | | | | | | | | | | | |
| I | pred.0 | 14 | 6 | 14 | 6 | 12 | 7 | 14 | 5 | 14 | 6 | 13 | 5 | 14 | 5 |
| | pred.1 | 6 | 44 | 6 | 44 | 8 | 43 | 6 | 45 | 6 | 44 | 7 | 45 | 6 | 45 |
| II | pred.0 | 43 | 28 | 43 | 25 | 43 | 27 | 43 | 23 | 43 | 23 | 43 | 27 | 43 | 23 |
| | pred.1 | 7 | 22 | 7 | 25 | 7 | 23 | 7 | 27 | 7 | 27 | 7 | 33 | 7 | 27 |

Note: DT-I, RF-I, KNN-I, ANN-I, SVM-I, LR-I, LR(SVM)-I: The method acts on dataset I; DT-II, RF-II, KNN-II, ANN-II, SVM-II, LR-II, LR(SVM)-II: The method acts on dataset II; B: Backward Elimination; A: AdaBoost; BA: Backward Elimination + AdaBoost; OB: Optimize Parameters (Evolutionary)+ Backward Elimination; OA: Optimize Parameters (Evolutionary)+AdaBoost; OBA: Optimize Parameters (Evolutionary)+ Backward Elimination + AdaBoost;

## 5. Conclusions

In this research, we proposed a novel method called BNB-NO-BK for detecting fraudulent crowdfunding projects. Firstly, the novel BNB approach constructed a fine-tuned Bert QA model to extract characteristics of crowdfunding projects using the auxiliary information of crowdfunding projects provided by BERT and NLTK. In the task of extracting characteristics of crowdfunding projects, the accuracy of the BNB method is 92.38%. Secondly, we used a NO method that contains a modified Nice Classification and Ontology to classify the categories of crowdfunding projects. We can classify all categories of crowdfunding projects by our modified Nice Classification. Then, we proposed a novel BK (Brand Knowledge) method to generate features for machine learning classifiers. The method used measures of association rules in web‑data to cross‑check techniques for crowdfunding projects. We generated 14 features for machine learning classifiers with the BK method. Finally, the machine learning classifier we applied was used to predict and classify fraudulent crowdfunding projects. Furthermore, to validate the effectiveness of our proposed method, we conducted experiments on a dataset of technology‑based crowdfunding projects (TCPD). We applied our method to detect TCPD and calculated the accuracy of our method. In addition, we used the method in the existing study to detect TCPD. The experimental results showed that the detection by the existing study was only 21.32% in accuracy of detecting fake crowdfunding projects, while our research method was 95.71% accurate in detecting fake crowdfunding projects.

## Acknowledgements

## References

Bondielli, A., & Marcelloni, F. (2019). A survey on fake news and rumour detection techniques. *Information Sciences, 497*, 38-55. doi:10.1016/j.ins.2019.05.035

Boonchuay, K., Sinapiromsaran, K., & Lursinsap, C. (2017). Boundary expansion algorithm of a decision tree induction for an imbalanced dataset. *Songklanakarin Journal of Science Technology, 39*(5), 665-673. doi: 10.14456/sjst-psu.2017.82

Devlin, J., Chang, M. W., Lee, K., & Toutanova, K. (2018). Bert: Pre-training of deep bidirectional transformers for language understanding. *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics*, 4171–4186. doi:10.18653/v1/N19-1423

Dou, Y., Shu, K., Xia, C., Yu, P. S., & Sun, L. (2021). User preference-aware fake news detection. *Proceedings of the 44ᵗʰ International ACM SIGIR Conference on Research and Development in Information Retrieval,* 2051–2055. doi:10.1145/3404835.3462990

Jiang, S., Chen, X., Zhang, L., Chen, S., & Liu, H. (2019). User-characteristic enhanced model for fake news detection in social media. *CCF International Conference on Natural Language Processing and Chinese Computing Cha,* 634-646. doi:10.1007/978-3-030-32233-5_49

Kochkina, E., Liakata, M., & Zubiaga, A. (2018). All-in-one: Multi-task learning for rumour verification. *Proceedings of the 27ᵗʰ International Conference on Computational Linguistics,* 3402–3413. Retrieved from https://aclanthology.org/C18-1288

Liao, Q., Chai, H., Han, H., Zhang, X., Wang, X., Xia, W., & Ding, Y. (2021). An integrated multi-task model for fake news detection. *IEEE Transactions on Knowledge and Data Engineering*, 1-1. doi:10.1109/TKDE.2021.3054993

Lu, Y.-J., & Li, C.-T. (2020). GCAN: Graph-aware co-attention networks for explainable fake news detection on social media. *Proceedings of the 58ᵗʰ Annual Meeting of the Association for Computational Linguistics,* 505-514. doi:10.18653/v1/2020.acl-main.48

Perez, B., Machado, S. R., Andrews, J. T., & Kourtellis, N. (2020). I call BS: Fraud detection in crowdfunding campaigns (Preprint arXiv:16849). doi:10.48550/arXiv.2006.16849

Pickering, S. (2001). Common sense and original deviancy: News discourses and asylum seekers in Australia. *Journal of Refugee Studies, 14*(2), 169-186. doi:10.1093/jrs/14.2.169

Qi, P., Cao, J., Yang, T., Guo, J., & Li, J. (2019). Exploiting multi-domain visual information for fake news detection. *IEEE International Conference on Data Mining,* 518-527. Retrieved from https://arxiv.org/pdf/1908.04472.pdf

Qi, L & Qu, J. (2022). Product ontology construction for crowdfunding projects. *7ᵗʰ International Conference on Business and Industrial Research*, 568-572. doi:10.1109/ICBIR54589.2022.9786391

Qu, J., Nguyen, L. M., & Shimazu, A. (2016). Cross-language information extraction and auto evaluation for OOV term translations. *China Communications, 13*(12), 277-296. doi:10.1109/CC.2016.7897550

Qu, J., Theeramunkong, T., Le Ming, N., Shimazu, A., Nattee, C., & Aimmanee, P. (2012). A flexible rule-based approach to learn medical English-Chinese OOV term translations from the web. *International Journal of Computer Processing of Languages, 24*(2), 207-236. doi:10.1142/S1793840612400132

Reimers, N., & Gurevych, I. (2019). Sentence-bert: Sentence embeddings using siamese bert-networks (Preprint arXiv:1908.10084). Retrieved from file:///C:/Users/user/Downloads/Nikolaev_N_Sentence_BERT_Sentence_Embeddings_using_Siamese_BERT.pdf

Shu, K., Sliva, A., Wang, S., Tang, J., & Liu, H. (2017). Fake news detection on social media: A data mining perspective. *SIGKDD Explorations Newsletter, 19*(1), 22–36. doi:10.1145/3137597.3137600

Vaibhav, V., Mandyam, R., & Hovy, E. (2019). Do sentence interactions matter? Leveraging sentence level representations for fake news classification. *Proceedings of the Thirteenth Workshop on Graph-Based Methods for Natural Language Processing, Hong Kong,* 134-139. doi:10.18653/v1/D19-5316

Xie, Q., Dai, Z., Hovy, E., Luong, T., & Le, Q. (2020). Unsupervised data augmentation for consistency training. *Advances in Neural Information Processing Systems, 33*, 6256-6268. Retrieved from https://dl.acm.org/doi/abs/10.5555/3495724.3496249