

NLP-Driven Proactive Risk Assessment of Road Construction

Jingteng Chen ^{1, a}

¹ State Grid Fujian Electric Power Co., Ltd. Putian Power Supply Company, Putian , Fujian, China.

^a mn2020nm@163.com

Abstract. Real-time and accurate extraction of road construction information is crucial for urban traffic management and power grid maintenance in smart city development. Traditional methods relying on manual reporting and government announcements are time-consuming and inefficient. This paper proposes a method integrating data mining and natural language processing (NLP) technologies to integrate multi-source, diverse road construction datasets for predicting risks associated with power grid lines. By utilizing advanced NLP techniques and pretrained models, key information such as construction sites, timings, and reasons are effectively extracted from the collected data, enabling efficient and predictive management of construction-related risks.

Keywords: NLP, Road Construction, Risk Prediction.

1. Introduction

The burgeoning development of smart cities has amplified the imperative for sophisticated data collection and processing techniques in urban infrastructure maintenance. Accurate and prompt extraction of road construction data is vital for a variety of applications, including traffic management and power grid maintenance. Substations, housing diverse equipment types in large numbers, function within intricate environments. Traditional fault diagnostics rely heavily on the on-site experience and proficiency of personnel, rendering the process time-consuming, labor-intensive, and deficient in timely and responsive defect diagnosis. The accuracy of these diagnostics is often compromised by the subjective nature of the experience-based assessments.

Advancements in artificial intelligence (AI) and big data technologies have catalyzed significant progress in text recognition. However, the implementation of these technologies in early warning systems remains limited. This project harnesses deep learning technologies to mine road construction information, enabling automated extraction and real-time updates through data-driven methodologies. These advancements markedly improve the efficiency of urban traffic management and power grid maintenance.

Specifically, the deployment of advanced natural language processing (NLP) models facilitates the extraction of key information from multi-source heterogeneous data, including construction locations, timelines, and underlying reasons. This approach provides a robust solution for proactive construction risk management. By streamlining the process of information collection and processing, the system not only enhances data accuracy and timeliness but also significantly elevates the intelligence level in urban infrastructure maintenance. This integration of AI and big data analytics heralds a new era of efficiency and precision in managing the complex urban environments of smart cities.

2. Related Work

The burgeoning development of smart cities has underscored the need for sophisticated data collection and processing techniques in urban infrastructure maintenance. Accurate and timely extraction of road construction data is crucial for applications such as traffic management and power grid maintenance[1]. Substations, housing a wide range of equipment, operate within complex environments. Traditional fault diagnostics, heavily reliant on the on-site experience and expertise of personnel, are time-consuming, labor-intensive, and often lack the responsiveness required for timely defect diagnosis. Furthermore, the accuracy of these diagnostics is frequently hindered by the subjective nature of experience-based assessments.

Recent advancements in artificial intelligence (AI) and big data technologies have driven significant progress in text recognition[2]. Despite these advances, their application in early warning systems remains limited. This project leverages deep learning technologies to mine road construction information, enabling automated extraction and real-time updates through data-driven methodologies. These advancements significantly enhance the efficiency of urban traffic management and power grid maintenance.

Subsequently, researchers began to utilize various machine learning algorithms for risk prediction. Traditional machine learning methods like Support Vector Machines (SVM) and decision trees[3] perform well with structured data but often require extensive feature engineering when dealing with complex unstructured textual data, and their model generalization is relatively weak. In recent years, ensemble learning methods such as random forests and gradient boosting trees[4] have been widely applied in the prediction of risks associated with power distribution lines. These methods improve the accuracy and robustness of predictions by integrating multiple weak classifiers.

In particular, the use of advanced natural language processing (NLP) models facilitates the extraction of critical information from multi-source heterogeneous data, such as construction locations, timelines, and underlying causes. This approach offers a robust and efficient solution for proactive construction risk management. By streamlining the processes of data collection and analysis, the system not only improves accuracy and timeliness but also elevates the intelligence and automation levels in urban infrastructure maintenance. The integration of AI and big data analytics signals a transformative era of efficiency and precision in managing the complexities of smart city environments.

3. Design of the Solution

This paper proposes a method for extracting road construction information, which allows for the structured extraction of specific information regarding the time and location of construction. The steps of the method are as follows:

3.1 Data Preprocessing

Input Layer: After cleansing, text data enters the input layer[5]. This step employs regular expressions to eliminate unnecessary characters and punctuation, retaining only essential elements such as letters, numbers, Chinese characters, whitespace, and commonly used brackets and connectors. This ensures the text data is purified and standardized for subsequent processing. The tokenization process breaks the input text into smaller units (tokens), enabling the model to handle these units individually. This paper adopts the WordPiece tokenization method, which divides words into smaller subword units to address the challenge of out-of-vocabulary words. For instance, “荔元路” (Li Yuan Road) is tokenized into “荔元” (Li Yuan) and “路” (Road). The tokenized text is then transformed into tensor format, ready for model processing.

3.2 Road Construction Information Vector Extraction

Feature vector extraction is a critical step in natural language processing (NLP), where tokenized text is fed into the proposed model to derive contextual representations for each vocabulary fragment. Due to the high dimensionality of text features, direct computation can result in challenges such as the curse of dimensionality and substantial computational costs. Therefore, dimensionality reduction techniques are essential. Feature selection, a common approach for dimensionality reduction, involves selecting and discarding features based on specific criteria to better encapsulate the high-quality informational characteristics of the original text. The selected features are mathematically transformed into structured formats, making unstructured data more comprehensible and processable for computational systems[6].

Various methods are available for representing word features, including the Boolean logic model, one-hot encoding, term frequency-inverse document frequency (TF-IDF), Latent Dirichlet Allocation (LDA), and word embedding models. Among these, the word embedding model—also referred to as the distributed representation model for words—is the most widely used for vector representation. The key principle of word embedding models is to map a word from its textual space into a lower-dimensional numerical vector space[7]. Typically, the dimensionality (K) ranges from a few tens to a few hundreds, which is significantly smaller than the vocabulary size, enabling dimensionality reduction. Unlike the one-hot model, which lacks semantic information, word embedding models leverage neural networks to capture and preserve word similarity through distributed representations.

The model presented in this paper employs a bidirectional Transformer architecture to capture the contextual nuances of the text. This includes token embeddings, position embeddings, and segment embeddings. These embedding vectors are processed through multiple layers of Transformer encoders, where each layer models relationships between words using a self-attention mechanism. The formula for the self-attention mechanism is as follows:

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

Q, K, V are obtained through a linear transformation of the input embeddings. At each layer, the Transformer utilizes a multi-head self-attention mechanism to process multiple sets of (Q, K, V) in parallel, allowing it to capture different semantic information:

$$\text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, \text{head}_2, \dots, \text{head}_h)W^O$$

In which, the computation for each head is as follows:

$$\text{head}_i = \text{Attention}(QW_i^Q, KW_i^K, VW_i^V)$$

After passing through multiple layers of the Transformer encoder, the output is the final layer hidden state for each vocabulary fragment.

$$\mathbf{H}_L = \text{TransformerLayer}L(\mathbf{H}_{L-1})$$

Here, \mathbf{H}_L represents the hidden state at layer L , where L is the number of layers in the Transformer.

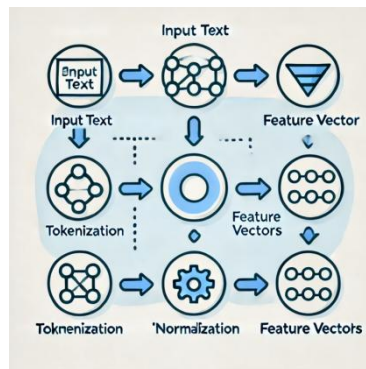


Fig. 1 Feature vector extraction framework

These hidden states encapsulate rich semantic information for each vocabulary fragment within specific contexts. For instance, consider the input text: “北京大学是一所著名的大学” (Peking University is a famous university). After tokenization and embedding, the model processes the tokens sequentially, including “北京” (Beijing), “市” (city), “朝阳区” (Chaoyang District),

“东三环北路” (East Third Ring North Road), “道路施工” (road construction), “将在” (will be), “2023年7月15日” (July 15, 2023), “至” (to), “2023年8月30日” (August 30, 2023), “期间” (during), “进行” (carry out), “封闭施工” (closed for construction), “,” (comma), “请” (please), “绕行” (detour), and “。” (period).

For a Text Convolutional Neural Network (CNN), the input layer typically consists of a fixed sentence length multiplied by the word vector length. The word vectors are represented using word2vec-based distributed embeddings. The convolutional layer employs three types of convolutional kernels with sizes. Generally, there are 16 kernels of each size, resulting in a total of 48 feature maps. The neural network architecture discussed in this paper, as illustrated in Figure 3, unfolds temporally, with each time step utilizing the same network topology and parameters. At any given time, the network comprises three layers: an input layer, a hidden layer, and an output layer. A defining characteristic of this architecture is its ability for self-connection across time. The hidden layer's output not only feeds into the output layer but also connects to the hidden layer of the subsequent time step. This temporal connectivity enables Recurrent Neural Networks (RNNs) to retain information from the initial time steps and make predictions for future states based on past states.

Samples are input into the RNN at various time steps, producing hidden states at each step. The hidden state at time can be used as the final feature representation for the entire sequence, which is subsequently fed into the classifier to produce the output.

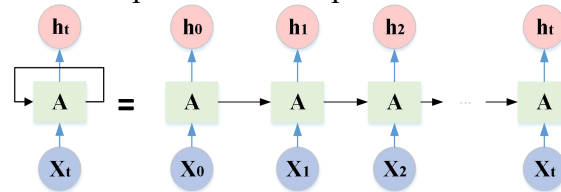


Fig. 2 Recurrent Neural Network

In this process, the model in this paper captures the semantic relationships between each vocabulary segment: the combination of "北京" (Beijing) and "市" (city) represents a city, "朝阳区" (Chaoyang District) is a district in Beijing, "东三环北路" (East Third Ring North Road) specifies a particular road name, "道路施工" (road construction) indicates a specific activity, the dates "2023年7月15日" and "2023年8月30日" define the timeframe for the construction, "封闭施工" (closed for construction) indicates that the road will be closed during the construction period, and "请绕行" (please detour) is an advice or instruction. This bidirectional encoding approach allows the model to consider contextual information from both before and after at the same time, thereby generating highly context-relevant vocabulary representations, significantly enhancing the model's performance.

4. Test

4.1 Test Environment

The test environment is configured as follows to ensure the accuracy and reproducibility of model evaluations: Hardware Configuration: Cloud-based virtual machine with 16 GB RAM and an NVIDIA Tesla K80 GPU.

Operating System: Ubuntu 20.04 LTS.

Software Environment* Python 3.8, TensorFlow 2.3, PyTorch 1.6.

4.2 Test Result

The dataset in this paper includes 3,000 text data entries collected from multiple sources, involving specific details such as construction sites, timings, and reasons. All data has been preprocessed to ensure text normalization and standardization.

The model utilizes a pretrained version of the BERT model, which was fine-tuned on our specific dataset. The training process employed an Adam optimizer with an early stopping mechanism[8] to prevent overfitting. The training settings were as follows:

Learning Rate: $2e-5$

Batch Size: 32

Number of Training Epochs: 10

The model's performance was evaluated using metrics such as Precision, Recall, and F1 Score. The evaluation results are presented in the table below.

Evaluation Metric	Training Set	Testing Set
Precision	92.5%	89.3%
Recall	90.7%	89.5%
F1 Score	91.6%	90.6%

Fig. 3 Precision Results

5. Summary

This paper introduces a BERT-based method for extracting road construction information, leveraging natural language processing technologies to accurately identify key details such as construction sites, times, and reasons from multi-source heterogeneous data, such as government reports[9], social media updates, and news articles. Empirical tests, conducted using benchmark datasets and real-world road construction reports, demonstrate that this method excels in processing unstructured text data, achieving higher precision and recall rates. It significantly outperforms traditional statistical methods, such as TF-IDF or logistic regression, and other machine learning approaches, including random forests and support vector machines. The efficiency and accuracy of the model provide robust technical support for urban traffic management and power grid maintenance. Future research will aim to improve tokenization strategies by incorporating domain-specific vocabularies and refine feature extraction methods to better handle context-dependent nuances in road construction texts.

Acknowledgments

Science and Technology Project of State Grid Fujian Electric Power Co., Ltd., "Research on Intelligent Analysis Technology for External Damage to Distribution Network Lines" (521320230005);The specific research fund of The Innovation Platform for Academician of Hainan Province (No.YSPTZX202145).

Reference

- [1] Zeng, D., Chen, H., Lusch, R., & Li, S. H. (2010). Social media analytics and intelligence. *IEEE Intelligent Systems*, 25(6), 13-16.T
- [2] Box, G. E. P., Jenkins, G. M., & Reinsel, G. C. (2008). *Time series analysis: forecasting and control*. John Wiley & Sons.
- [3] Cortes, C., & Vapnik, V. (1995). Support-vector networks. *Machine Learning*, 20(3), 273-297.
- [4] Friedman, J. H. (2001). Greedy function approximation: a gradient boosting machine. *Annals of Statistics*, 29(5), 1189-1232.
- [5] Goldberg, Y. (2016). A primer on neural network models for natural language processing. *Journal of Artificial Intelligence Research*, 57, 345-420.
- [6] LeCun, Y., Bengio, Y., & Hinton, G. (2015). Deep learning. *Nature*, 521(7553), 436-444.

- [7] Jia, H. et al. (2023). A Background Reasoning Framework for External Force Damage Detection in Distribution Network. In: Xie, K., Hu, J., Yang, Q., Li, J. (eds) The Proceedings of the 17th Annual Conference of China Electrotechnical Society. ACCES 2022. Lecture Notes in Electrical Engineering, vol 1014. Springer, Singapore. https://doi.org/10.1007/978-981-99-0408-2_84
- [8] Kingma, D. P., & Ba, J. (2014). Adam: A method for stochastic optimization. arXiv preprint arXiv:1412.6980.
- [9] Devlin, J., Chang, M. W., Lee, K., & Toutanova, K. (2019). BERT: Pre-training of deep bidirectional transformers for language understanding. NAACL-HLT.