

# Legal Privacy Protection Machine Learning Method Based on Word2Vec Algorithm

Rongrong Wang

 <https://orcid.org/0009-0009-8000-379X>

*Zhe Jiang J.R.C. Law Firm, China*

## ABSTRACT

This study uses Word2Vec's word vector representation technology to finely capture the semantic relationships of vocabulary in legal texts through the Skip-gram model. By introducing Hierarchical Softmax optimization, a legal privacy protection model based on Word2Vec algorithm is ultimately designed. The results showed that the model performed better than other comparative algorithms in both the macro classification performance (F1\_macro) and the micro classification performance (F1\_micro). In practical legal sensitive word recognition tasks, the accuracy, recall rate, and F1 score of the model reached 92.56%, 88.78%, and 90.62%, respectively. Therefore, the proposed model effectively improved the accuracy of identifying sensitive legal privacy words and providing new methods for the personal information security protection system.

## KEYWORDS

Word2Vec, LSTM-B, Privacy Protection, Machine Learning, Text Recognition

## INTRODUCTION

In consequence of the accelerated pace of technological advancement in the sphere of internet technology, the number of internet users in China has surpassed one billion, with the penetration rate of the internet exhibiting a marked increase (Amin et al., 2023). However, in this context, information security issues have become increasingly prominent, with frequent leaks of personal legal information. In particular, the irregular collection and processing of personal information by internet enterprises has become the focus of social attention (Nsugbe, 2023). As an important resource for internet enterprises, the collection and analysis of user information provides enterprises with the possibility of personalized services and precision marketing but also exposes users to potential threats of privacy disclosure (Hebbi & Mamatha, 2023; Wang et al., 2021a). The formulation and implementation of privacy policies are often controlled by enterprises, which can easily lead to the abuse of users' personal information. Although relevant laws and regulations such as the Personal Information Protection Law have been implemented, in practical operation, it is difficult for infringed users to obtain effective compensation, and the role of legal protection is limited (Wang et al., 2021b). With the increasing maturity of machine learning in recent years, applying it to legal privacy protection (LPP) has gradually become an effective choice (Chen et al., 2022). Given this, the study uses word2vec's word vector representation technology to finely capture the semantic relationships of vocabulary in legal texts through the skip-gram model. The introduction of hierarchical softmax (HS) optimizes the efficiency of traditional softmax, and combines word vector distance and TextRank algorithm to identify sensitive content in text. Finally, the long short-term memory network classifier (LSTM-B)

DOI: 10.4018/IJISP.365911

This article published as an Open Access article distributed under the terms of the Creative Commons Attribution License (<http://creativecommons.org/licenses/by/4.0/>) which permits unrestricted use, distribution, and production in any medium, provided the author of the original work and original publication source are properly credited.

is introduced, and an LPP machine learning model based on the word2vec algorithm (LPPM-W2V) is designed. This study aims to explore the LPPM-W2V method and improve the personal information security protection system.

## LITERATURE REVIEW

The progress of the internet has made more people begin to pay attention to the LPP problem on the internet. Wang proposed a new word2vec-based unsupervised anomaly detection method (LogUAD) to address the challenges faced by large-scale distributed system log analysis. Compared with log cluster, LogUAD's F1 score has increased by 67.25% (Wang et al., 2022). Ren et al. (2022) proposed a new k-mer context free alignment method (kmer2vec) to address the problem of high computational complexity and difficulty in effectively capturing sequence context structure in traditional multiple sequence alignment methods when dealing with large numbers of sequences. The kmer2vec performed well in constructing phylogenetic trees and species clustering, with a much faster running speed than the multiple sequence alignment method. Gao et al. (2022) proposed a method that combined the dynamic topic model and the word2vec model to address the issue of quantifying semantic distribution and its changing characteristics in topic evolution research. The themes in library and information science mostly corresponded to multiple semantic concepts, and there were three evolutionary patterns: convergence, diffusion, and stability, while the popularity of themes was independent of their evolutionary dynamics.

Lakshmanan and Anandha (2024) proposed a method that combines blockchain technology and a new privacy protection model to address the security vulnerabilities and privacy issues that medical data may face when transmitted through open communication channels. This method improved the privacy, reliability, and practicality of electronic medical models through two stages of data cleaning and recovery, and had significant advantages in comparison. Schulze et al. (2023) proposed a deep learning-based internal discrimination algorithm and an outer discrimination algorithm to handle the issue of privacy leakage in laparoscopic surgery videos. This algorithm was used to automatically identify and mask non-abdominal areas in videos to protect patient privacy. The outer discrimination algorithm performed well in classifying non-abdominal frames, with an average F1 score of  $0.96 \pm 0.01$  (binary classification) and  $0.97 \pm 0.01$  (fifth classification). Cao et al. (2022) proposed a blockchain privacy protection data mining algorithm based on decision tree classification to address the problems of low mining accuracy and high data noise in privacy protection data mining methods in blockchain. The mining accuracy of this algorithm was always above 90%, and the data noise was stable below 0.6 dB.

Although the above-mentioned studies have proposed innovative privacy protection or data processing methods in specific fields and achieved certain results, there are still some common or individual shortcomings. Firstly, some studies lack broad applicability validation for datasets of different types or sizes, and the results may not be universally applicable. Secondly, some methods may sacrifice computational efficiency or resource consumption while pursuing performance improvement, which may pose challenges in practical applications (Dhinakaran & Prathap, 2022). Therefore, based on word2vec's word vector representation, this study introduces HS to optimize the efficiency of traditional softmax, and ultimately constructs the LPPM-W2V model.

## LPPM-W2V MODEL

### Improvement of Word2vec

In LPP, identifying and processing sensitive information is an important task for Dynamic Topic Model. Word2vec is utilized to generate word vectors. The model is a shallow and double-layer neural network, which is defined by words and requires guessing the input of adjacent positions (Balfagih et

17 more pages are available in the full version of this document, which may be purchased using the "Add to Cart" button on the publisher's webpage: [www.igi-global.com/article/legal-privacy-protection-machine-learning-method-based-on-word2vec-algorithm/365911](http://www.igi-global.com/article/legal-privacy-protection-machine-learning-method-based-on-word2vec-algorithm/365911)

## Related Content

---

### Trust and Reliability Management in Large-Scale Cloud Computing Environments

Punit Gupta (2021). *Large-Scale Data Streaming, Processing, and Blockchain Security* (pp. 66-89).

[www.irma-international.org/chapter/trust-and-reliability-management-in-large-scale-cloud-computing-environments/259465](http://www.irma-international.org/chapter/trust-and-reliability-management-in-large-scale-cloud-computing-environments/259465)

### A Dynamic Subspace Anomaly Detection Method Using Generic Algorithm for Streaming Network Data

Ji Zhang (2015). *Handbook of Research on Emerging Developments in Data Privacy* (pp. 403-425).

[www.irma-international.org/chapter/dynamic-subspace-anomaly-detection-method/123543](http://www.irma-international.org/chapter/dynamic-subspace-anomaly-detection-method/123543)

### Strategies for Mitigating Security Concerns in IoT-Enabled Smart Cities

Ravikumar, Shilpa Singhal, Santushti Betgeriand Sushil Kumar Singh (2024). *Secure and Intelligent IoT-Enabled Smart Cities* (pp. 239-273).

[www.irma-international.org/chapter/strategies-for-mitigating-security-concerns-in-iot-enabled-smart-cities/343453](http://www.irma-international.org/chapter/strategies-for-mitigating-security-concerns-in-iot-enabled-smart-cities/343453)

### Formal Reliability Analysis of Engineering Systems

Naeem Abbasi, Osman Hasanand Sofiène Tahar (2014). *Network Security Technologies: Design and Applications* (pp. 224-238).

[www.irma-international.org/chapter/formal-reliability-analysis-of-engineering-systems/105810](http://www.irma-international.org/chapter/formal-reliability-analysis-of-engineering-systems/105810)

### Computer Security Practices and Perceptions of the Next Generation of Corporate Computer Use

S.E. Kruckand Faye P. Teer (2008). *International Journal of Information Security and Privacy* (pp. 80-90).

[www.irma-international.org/article/computer-security-practices-perceptions-next/2477](http://www.irma-international.org/article/computer-security-practices-perceptions-next/2477)