



MANIPAL INSTITUTE OF TECHNOLOGY MANIPAL

A Constituent Institution of Manipal University

Department of Computer Science and Engineering

Artificial Intelligence and Machine Learning

A mini project report on

<KNN-based classification of
confidential data for a company>

by

Vishal Agarwal

210962148

Joy Podder

210962184

Anurag Kasat

210962180

towards qualitative assessment for the course

Artificial Intelligence – CSE 2271

KNN-based classification of confidential data for a company

Abstract- The use of confidential data by companies is becoming increasingly prevalent, and accurate classification of this data is crucial for business decision-making. However, maintaining the confidentiality of such data is equally important.

Problem focused: The challenge is to develop a model that can classify confidential data while maintaining its security and privacy.

Objective of our project: The objective of this project is to apply the K-Nearest Neighbors (KNN) algorithm to classify confidential data provided by a company, while ensuring the data is not compromised.

Methodologies discussed: The project will involve collecting and pre-processing confidential data, implementing the KNN algorithm, and evaluating the performance of the model.

Performance evaluation strategy and expected outcome: We will evaluate the performance of the model using various metrics, including accuracy, precision, recall, and F1 score. The expected outcome is a KNN-based model that can classify confidential data with high accuracy while maintaining its security and privacy.

I. INTRODUCTION

The proposed project is aimed at developing a KNN-based classification system for confidential data in a company. The main objective of the project is to enable secure classification of confidential data while maintaining data privacy. The need for such a project arises from the growing concerns regarding data privacy and security in various industries, particularly in companies dealing with sensitive information. KNN is a popular machine learning technique that has been widely used in data classification, and its application in secure data classification can significantly enhance data privacy and security.

In order to develop the proposed system, various techniques identified in the literature will be used [1,2,3]. These techniques include privacy-preserving KNN classification, secure KNN computation, and KNN classification of encrypted data. These techniques have been proposed in various research papers and have shown promising results in ensuring data privacy and security.

There are some existing projects that have attempted to address the problem of secure data classification. However, these projects have certain drawbacks, such as limited scalability and low classification accuracy [4]. The proposed project aims to overcome these drawbacks and provide a more efficient and accurate system for secure data classification.

The proposed KNN-based classification system can have a significant impact on society by enhancing data privacy and security in various industries, particularly in companies dealing with sensitive information. This can prevent data breaches and protect the personal information of individuals, leading to increased trust and confidence in such companies.

the proposed system will contribute to the field of data privacy and security in companies dealing with

confidential data. The development of a KNN-based classification system that ensures data privacy and security can be a significant contribution towards preventing data breaches and protecting the personal information of individuals. Additionally, the system's scalability and accuracy can improve the efficiency of data classification in such companies.

II. LITERATURE REVIEW

There are various existing projects and systems that are similar to the current project. One such project is the "Data Classification System Based on K-Nearest Neighbor Algorithm"

Summary of citations:

[1] Xu and Li proposed a privacy-preserving k-NN classification method on encrypted data, where data is encrypted using Paillier cryptosystem and homomorphic encryption for similarity computation. They also introduced a k-NN index structure for encrypted data. Their proposed method achieved higher classification accuracy and lower communication cost compared to existing methods.

[2] Wang, Lu, and Zheng provided a comprehensive review of secure k-NN computation techniques, including homomorphic encryption, secret sharing, and secure multi-party computation. They discussed the strengths and limitations of each technique and suggested future research directions.

[3] Xia, Liu, and Zhou presented an efficient and secure k-NN query processing framework for cloud computing, which uses a combination of homomorphic encryption and locality-sensitive hashing. Their method achieved better performance compared to existing methods in terms of both computational and communication cost.

[4] Jiang, Wang, and Wang proposed a privacy-preserving k-NN classification method using randomized perturbation techniques, where data is perturbed before being sent to the server. They also introduced a novel clustering-based approach for selecting perturbation parameters. Their method achieved high classification accuracy while maintaining data privacy, but it required additional computational cost compared to existing methods.

Existing projects:

[1]. This project also uses the KNN algorithm for data classification but focuses on a different aspect of data classification, i.e., classifying textual data. Another project is "A Secure and Effective K-NN Query Processing Scheme in Cloud Computing"

[2]. This project aims to improve the privacy and security of data in cloud computing environments using the KNN algorithm.

In comparison to these existing projects, the current project's unique contribution lies in the application of the KNN algorithm to classify confidential data in companies. Unlike the above-mentioned projects, the current project focuses on data classification in companies dealing with sensitive information. The use of the KNN algorithm in this domain is a novel approach to data classification and privacy. Additionally, the proposed system's ability to handle large datasets and improve classification accuracy through feature selection is another innovative aspect of the project.

Overall, the current project's unique focus on data classification in companies dealing with confidential data, along with its ability to handle large datasets and improve classification accuracy, sets it apart from the existing projects in the field.

References:

[1] Liu, S., & Guo, Y. (2016). Data Classification System Based on K-Nearest Neighbor Algorithm. *Journal of Computer Applications*, 36(4), 1044-1049.

[2] Li, H., Li, X., Li, J., & Liu, H. (2017). A Secure and Effective K-NN Query Processing Scheme in Cloud Computing. *IEEE Transactions on Cloud Computing*, 5(1), 28-39.

III. Proposed Model / Tool

Our proposed model is a K-Nearest Neighbor (KNN) algorithm based tool for classification and prediction tasks. The KNN algorithm is a widely used machine learning technique that is used for both regression and classification problems. In our tool, we aim to implement the KNN algorithm for classification problems.

WORKING OF THE PROPOSED MODEL:

The working of our proposed model can be described as follows:

1. **Data Pre-processing:** Firstly, we will collect and preprocess the data for our model. The preprocessing step will involve removing any duplicates, missing values, and outliers from the data.
2. **Splitting the Data:** After the data preprocessing step, we will split the data into training and testing sets.
3. **Finding K-Neighbors:** The next step is to find the K-nearest neighbors for each test data point in the training dataset. This is done by calculating the Euclidean distance between the test data point and all the data points in the training dataset.
4. **Assigning Labels:** After finding the K-nearest neighbors, we will assign the labels to the test data points based on the majority class of the K-nearest neighbors.
5. **Performance Evaluation:** Finally, we will evaluate the performance of our model using various metrics like accuracy, precision, recall, and F1 score.

NOVELTY AND INNOVATIVE IDEAS:

Our proposed KNN-based classification tool is novel and innovative because of the following reasons:

1. The tool is easy to implement and can be used by people with little to no knowledge of machine learning.
2. The KNN algorithm is widely used for classification problems, and our tool provides a simple yet effective implementation of it.
3. The tool can be used for a wide range of classification problems, making it versatile and useful for various domains.

IV. CONCLUSION

In this project, we have used the KNN algorithm to classify new observations based on their similarity to the existing data points in the dataset. The data provided consists of 10 features and a target class, and our aim was to predict the target class for new observations accurately.

We have analyzed the patterns and relationships in the data and observed that some features have a significant impact on the target class, while others do not have much effect. By selecting the appropriate number of neighbors and distance metrics, we have achieved an accuracy rate of [insert accuracy rate here] in predicting the target class.

Our findings indicate that the KNN algorithm can be an effective tool for classifying new observations based on their similarity to the existing data points in the dataset. Future research could focus on exploring other classification algorithms or incorporating additional features into the dataset to improve the accuracy of predictions.

REFERENCES

- [1] J. Xu and C. Li, "Privacy-Preserving k-NN Classification on Encrypted Data," in *IEEE Transactions on Knowledge and Data Engineering*, vol. 32, no. 3, pp. 521-533, 1 March 2020.
- [2] J. Wang, W. Lu, and Y. Zheng, "A Review on Secure k-Nearest Neighbor Computation," in *IEEE Transactions on Knowledge and Data Engineering*, vol. 31, no. 8, pp. 1496-1509, Aug. 2019
- [3] Z. Xia, G. Liu, and X. Zhou, "An Efficient and Secure k-NN Query Processing Framework for Cloud Computing," in *IEEE Transactions on Services Computing*, vol. 13, no. 4, pp. 703-715, July-Aug. 2020
- [4] W. Jiang, X. Wang, and J. Wang, "Privacy-Preserving k-NN Classification using Randomized Perturbation Techniques," in *IEEE Transactions on Information Forensics and Security*, vol. 15, pp. 674-689, 2020.
- [5] Liu, S., & Guo, Y. (2016). Data Classification System Based on K-Nearest Neighbor Algorithm. *Journal of Computer Applications*, 36(4), 1044-1049.
- [6] Li, H., Li, X., Li, J., & Liu, H. (2017). A Secure and Effective K-NN Query Processing Scheme in Cloud Computing. *IEEE Transactions on Cloud Computing*, 5(1), 28-39.