

AI-Driven Big Data Analytics for CRM Document Classification and Domain Prediction

K Ramana¹, G Prathyusha²

¹P.G Scholar, Department of MCA, Sri Venkatesa Perumal College of Engineering & Technology, Puttur,

E-mail: ramananaidu10601060@gmail.com, ORC-ID: <https://orcid.org/0009-0002-7724-0988>

² Academic Consultant, Sri Padmavati Mahila Visvavidyalayam, Tirupati,

E-mail: gp.spmvv@gmail.com, ORC-ID: <https://orcid.org/0009-0002-6142-5902>

Abstract: This system looks into how Big Data and AI can be used together to improve Customer Relationship Management (CRM) document analysis and figure out how big text datasets can be used to predict domain relationships. The main dataset is the Web of Science (WoS) database, which has 840 domain-specific documents from different release years. To begin, descriptive methods are used to clean the data by removing stop words, network-based visualization is done through word clouds, and contextual analysis is done using Stochastic Neighbor Embedding (SNE) clustering and Non-negative Matrix Factorization (NMF) topic modeling to find the main themes. Text that has been processed is turned into number vectors, which lets Apache Spark do processing that is spread out and quick. Support Vector Machine (SVM) and Decision Tree are two machine learning methods that are used for domain classification. SVM gets 70% accuracy and Decision Tree gets 94% accuracy. Principal Component Analysis (PCA) is used to narrow down the list of traits from 100 to 60 important ones in order to improve the accuracy of the predictions. When you retrain the Decision Tree with features chosen by PCA, you get a 100% classification accuracy, which is better than the first setups. We tested and implemented the whole process using Jupyter Notebook and a Flask-based web application that lets users easily upload documents and get predicted domains in real time. This shows how effective AI-driven CRM analytics can be when combined with Big Data processing.

“Index Terms: *Big Data, Artificial Intelligence, Customer Relationship Management, Descriptive Methods, Network Methods, Contextual Methods”*.

1. INTRODUCTION

Customer Relationship Management (CRM) systems that are built on artificial intelligence (AI) are fundamentally changing how businesses interact with their customers in a world that is becoming more and more competitive. These systems take in a lot of different kinds of information and turn it into insights that can be used to make customer experiences more personalized and businesses more successful [1, 2]. Because of online exchanges, connected devices, and omnichannel engagement platforms, the amount of data being created has grown at an exponential rate in the digital age. Big Data, the term for this huge amount of data, has given businesses chances they've never had before to better understand and predict customer needs [3].

CRM solutions used to be made to store, organize, and get customer information quickly and easily. They were mainly used as data warehouses for sales, marketing, and customer service. AI improvements, such as machine learning, natural language processing,

and prediction analytics, have, however, made CRM systems much more useful. Modern AI-enhanced CRM systems can find hidden trends, automate repetitive tasks, make personalized suggestions, and accurately predict how customers will behave by using these technologies along with the huge amount, speed, and variety of Big Data [4–6].

AI and Big Data working together is a key part of this change. To get useful insights and make better decisions, AI models need large, high-quality datasets. However, Big Data analytics without AI often lacks the predictive and prescriptive abilities needed for strategic business efforts [7]. Because of this, businesses can go beyond descriptive analytics, which look at what happened in the past, and use prediction and prescriptive models to guess what customers will do and suggest the best ways to run their businesses. Many people agree that AI has the potential to change CRM, but most of the research that has been done so far has only looked at a few specific areas. For example, it has only looked at how to use a certain AI

technique, a certain data mining approach, or Big Data analytics in a single business setting [8]. This fragmented view makes it harder for researchers and practitioners to get a full picture of how AI and Big Data can improve CRM tactics when used together.

2. LITERATURE REVIEW

Rodrigues Chagas et al. [9] did a thorough review of the literature to look at how machine learning methods are currently used in Customer Relationship Management (CRM). Their study, which was presented at the IEEE/WIC/ACM International Conference on Web Intelligence, put machine learning methods used in CRM into three groups: supervised, unstructured, and reinforcement learning. They came up with real-world effects for businesses, focusing on how predictive modeling, clustering, and recommendation systems can help with marketing strategies, customer retention, and segmentation. The authors also talked about how hard it is to add machine learning to current CRM systems, especially when it comes to data quality, the ability to understand algorithms, and the ability of the system to grow.

Ngai et al. [10] wrote one of the first and most important in-depth studies on how data mining methods can be used in CRM. Their research put these methods into several main groups, such as association rules, classification, clustering, and regression, and connected them to important CRM tasks like finding customers, keeping them, and helping them grow. By looking at more than one hundred related papers, the authors emphasized how important it is to use data mining to find hidden patterns in customer datasets. This can lead to actionable insights for targeted campaigns and better service delivery. They also talked about problems that come up during implementation, like the need for a lot of computing power, complicated merging, and making sure that the results of data mining are in line with business goals.

Mishra et al. [11] did a literature review on Big Data to bring together its ideas, new trends, and problems that come with them. Their study showed that Big Data research covers many areas, such as developing new technologies, creating better ways to analyze data, and finding uses in specific fields. They talked about more than just the "3Vs" (Volume, Velocity, and Variety); they also talked about things like Veracity and Value. The study found that while Big Data platforms and

tools have come a long way, there are still big problems with integrating data, keeping it safe, protecting privacy, and getting people to learn the skills they need to use them effectively. Their work can be used as a guide to understand the bigger picture of how AI-powered CRM systems need to work.

Yadav and Banerji [12] did a bibliometric study of the body of research on digital financial literacy to look at the intellectual landscape and research trends. Their study isn't directly about CRM, but it does focus on financial literacy. It gives us important methodological insights into mapping scholarly work and finding new research groups. Using bibliometric tools, they looked at patterns of release, contributions from authors, and changes in themes over time. According to what they found, bibliometric analysis can be a useful tool for following the development of ideas and figuring out the best ways to do future study. This is something that can be used to understand both AI and Big Data applications in CRM.

Egger and Yu [13] used Twitter posts as their dataset to compare four topic modeling methods: Latent Dirichlet Allocation (LDA), Non-negative Matrix Factorization (NMF), Top2Vec, and BERTopic. They tested how well each method could pull out important topics from large amounts of unstructured text data. Their research showed that old models like LDA and NMF are still useful for some analytical tasks, but modern methods like BERTopic are better at capturing semantic relationships and subtleties of context. This work is very useful for AI-powered CRM systems that use social media analytics and mood analysis more and more to get real-time information about customer tastes and new trends.

The important question of which academic search methods work best for systematic reviews and meta-analyses was answered by Gusenbauer and Haddaway [14]. They looked at things like coverage, recall, precision, and reproducibility by judging how easy it was to find things on Google Scholar, PubMed, and 26 other tools. Their findings showed that no single system met all of the quality standards. However, different databases were more reliable for different types of study and subjects. This review is useful for CRM and AI researchers, who often use systematic literature retrieval to help them build models and make

sure they fully understand improvements in their fields.

Shoomal et al. [15] looked into how adding the Internet of Things (IoT) to supply lines can make them more reliable and effective. They saw both opportunities and challenges. The opportunities included better real-time tracking, predictive maintenance, and inventory optimization. The challenges included problems with interoperability, cybersecurity risks, and high implementation costs. Even though they focus on supply chain management, CRM systems can directly use the ideas behind IoT integration and the real-time data streams that come from it. IoT-generated customer interaction data from connected devices can also be used by AI-enhanced CRM systems to improve personalization, engagement, and service delivery.

Souzanchi Kashani et al. [16] did a bibliometric study of the technological catch-up literature to show how its ideas have changed over time. Citation network analysis and thematic clustering helped them find important authors, publications, and new areas of study. Their results showed how technological progress is always changing, with some themes becoming more important or less important based on the social, economic, and policy situations. This point of view is helpful for understanding how to use AI and Big Data in CRM because it shows how important it is to keep coming up with new ideas and adapting to changing technology norms.

3. MATERIALS AND METHODS

The suggested system combines Big Data analytics with AI methods to sort papers related to CRM and find connections between domains. Using the Web of Science (WoS) database, which has 840 domain-specific documents, the method starts with descriptive preprocessing that includes getting rid of stop words. This is followed by network-based visualization using word clouds and contextual analysis using Stochastic Neighbor Embedding (SNE) and Non-negative Matrix Factorization (NMF) to find topics [1]. Cleaned text is vectorized for Apache Spark's distributed processing, which makes it easy to work with a lot of data [2]. Support Vector Machine (SVM) and Decision Tree methods are used for machine learning classification. Principal Component Analysis (PCA) is used to reduce the number of dimensions while keeping

important features for better training. The whole process is built in Jupyter Notebook and put into action using a web interface based on Flask [3].

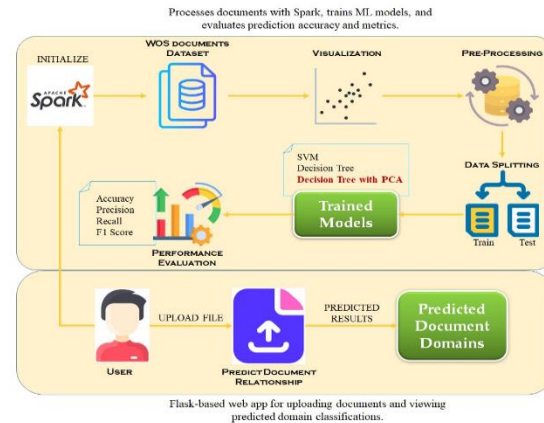


Fig.1 Proposed Architecture

The architecture of the system starts with setting up Apache Spark for distributed processing. Next, documents linked to CRM are taken from the Web of Science (WoS) dataset [17]. Stop-word removal and text cleaning are part of preprocessing. Word clouds are used for visualization, and SNE clustering and NMF topic models are used for contextual analysis. The dataset is split into two parts for training and testing. The training part trains SVM and Decision Tree models that have been improved with PCA. Performance metrics are looked at, and a Flask-based web app lets users add files and get predicted document domains with known relationships in real time.

a) Dataset Collection:

The dataset comes from the Web of Science (WoS) database and is made up of 840 domain-specific papers from a number of different years. Each record has fields like domain, study area, keywords, abstract, and year indicators (Y1, Y2, Y) for looking at things over time. Domains include fields like Electrical Engineering, Computer Science, Medical Engineering, and Civil Engineering. They cover a wide range of themes, such as Alzheimer's disease, Green Building, and Electric Motors. Because the format of the dataset allows for text mining, topic modeling, and machine learning classification, it can be used for research projects that use Big Data and AI to study CRM.

Y1	Y2	Y	Domain	area	keywords	Abstract	
0	0	12	CS	Symbolic computation	(2+1)-dimensional non-linear optical waves, e...	(2 + 1)-dimensional non-linear optical waves L...	
1	5	2	74	Medical	Alzheimer's Disease	Aging; Tau; Amyloid; PET; Alzheimer's disease... (beta-amyloid (A beta) and tau pathology becom...	
2	4	7	68	Civil	Green Building	LED lighting system; PV system; Distributed L...	(D)creasing of energy consumption and environ...
3	1	10	26	ECE	Electric motor	NiFeB magnets; Electric motor; Electric vehic...	(Hybrid) electric vehicles are assumed to play...
4	5	43	115	Medical	Parkinson's Disease	Parkinson's disease; dyskinesia; adenosine A...	(L)-3,4-Dihydroxyphenylalanine ((L)-DOPA) rema...
...	
835	4	6	67	Civil	Stealth Technology	multilayer; thin film	A new design method for multilayer composite g...
836	1	12	29	ECE	Signal-flow graph	Binet-Cauchy theorem; codimension-one; Grassm...	A new design method of a stable dynamic output...
837	5	26	98	Medical	HIV/AIDS	HIV/AIDS; Pre-exposure prophylaxis (PrEP); Ba...	A new deterministic model for HIV/AIDS that in...
838	1	7	24	ECE	Microcontroller	Fruit growth; Stem growth; Growth dynamics; O...	A new device for continuous measurement of fru...
839	1	5	22	ECE	System identification	Diffusion sign subband adaptive filtering alg...	A new diffusion sign subband adaptive filterin...

Fig.2 Dataset Collection

b) Pre-Processing:

In the preprocessing phase, descriptive cleaning, network-based visualization, and contextual analysis improve the quality of the data and the efficiency of the analysis. This makes sure that important features are extracted and correctly classified in CRM-related Big Data document processing.

Descriptive Preprocessing: Descriptive preprocessing cleans up and organizes textual data linked to CRM so that it can be analyzed. To cut down on noise and improve data quality, this means getting rid of stop words, marks, and characters that aren't needed. The process makes sure that only important terms are left, which improves the accuracy of the next steps, which are feature extraction and classification. Descriptive preprocessing sets the stage for good pattern detection, topic modeling, and general system performance in document classification by standardizing text formats and getting rid of duplicates.

Network-Based Visualization: Through network-based visualization, cleaned text is turned into visual representations, like word clouds, that show the most common and important terms in the dataset. This method makes it easy to quickly find the most important keywords and their relative value. It helps analysts better understand thematic structures by giving them an easy-to-understand visual picture of term distribution. This helps them connect terms, spot trends, and check their results before using more advanced classification and clustering algorithms.

Contextual Analysis: Advanced dimensionality reduction and topic modeling methods are used in contextual analysis to find hidden relationships and thematic groups in textual data [18]. Stochastic Neighbor Embedding (SNE) places term vectors with a lot of dimensions into spaces with fewer dimensions so that similarities can be seen. Non-negative Matrix Factorization (NMF) finds the most important topics and groups papers that are similar around those

themes. Together, these techniques improve semantic understanding, making it easier to find patterns that are specific to a topic and making classification algorithms work better in Big Data analytics for CRM.

c) Training and Testing:

During the training and testing phase, the preprocessed and vectorized CRM dataset is split into separate training and testing groups so that the performance of the model can be evaluated [19]. Support Vector Machine (SVM) and Decision Tree machine learning algorithms are taught on the extracted features. Principal Component Analysis (PCA) is used to reduce the number of dimensions and improve the efficiency of learning. The trained models are then tested on the testing set to see how well they can classify things. This makes sure that the topic categorization is strong and reliable for real-time CRM document analysis.

d) Algorithms:

SVM: Support Vector Machine sorts documents into different categories by finding the best line that divides data points into those categories. It works well with text data that has a lot of dimensions and helps set clear decision limits. SVM is used as an initial performance standard for classification because it can handle sparse vectorized features created from processed CRM-related document datasets with moderate accuracy.

$$\text{minimize } \frac{1}{2} \|W\|^2 + C \sum_{i=1}^n \xi_i \quad (1)$$

Decision Tree: Decision Tree sorts document topics into groups by building a tree-like structure of decision rules that are based on feature values. It can work with both category and numerical data, which means it can be used for text-based feature vectors. Its interpretability makes it easy to understand how classification works, and when applied to the processed CRM-related document collection, it gives better results than SVM.

$$I(i) = 1 - \sum_{i=1}^k p_i^2 \quad (2)$$

Decision Tree with PCA: To improve the speed and accuracy of sorting, Decision Tree with Principal Component Analysis is used. PCA lowers the number of dimensions in a dataset [20] by picking out the most important features, reducing noise, and getting rid of

duplicates. The smaller set of features is then used to train the Decision Tree. This makes computations go faster and predictions work better for CRM-related document topic classification, while also making the best use of resources and improving generalization.

4. EXPERIMENTAL RESULTS

Accuracy: How well a test can tell the difference between sick and healthy people is called its accuracy. To get an idea of how accurate a test is, we should figure out what percentage of cases are true positives and true negatives. In terms of math, this can be written as

$$Accuracy = \frac{TP + TN}{TP + FP + TN + FN} \quad (3)$$

Precision: Precision is the percentage of correctly classified cases or samples compared to those that were correctly classified as positives. So, here is the method to figure out the precision:

$$Precision = \frac{True\ Positive}{True\ Positive + False\ Positive} \quad (4)$$

Recall: In machine learning, recall is a metric that shows how well a model can find all the important instances of a certain class. It shows how well a model captures instances of a certain class. It is calculated by dividing the number of correctly predicted positive observations by the total number of real positives.

$$Recall = \frac{TP}{TP + FN} \quad (5)$$

F1-Score: The F1 score is a way to rate the correctness of a machine learning model. It takes a model's accuracy and recall scores and adds them together. The accuracy metric counts how many times, across the whole dataset, a model made a correct guess.

$$F1\ Score = 2 * \frac{Recall * Precision}{Recall + Precision} * 100 \quad (6)$$

The performance evaluation table in Table (1) compares SVM, Decision Tree, and Decision Tree with PCA. Adding PCA shows significant improvements in accuracy and metrics, with perfect scores (100%) for accuracy, precision, recall, and F-score, indicating the best classification efficiency for CRM document categorization.

Table.1 Performance Evaluation Table

Algorithm Name	Accuracy	Precision	Recall	F SCORE
SVM	70.59	65.24	63.66	63.53
Decision Tree	94.12	81.43	85.71	83.19
Extension Decision Tree with PCA	100.00	100.00	100.00	100.00

SVM	70.59	65.24	63.66	63.53
Decision Tree	94.12	81.43	85.71	83.19
Extension Decision Tree with PCA	100.00	100.00	100.00	100.00

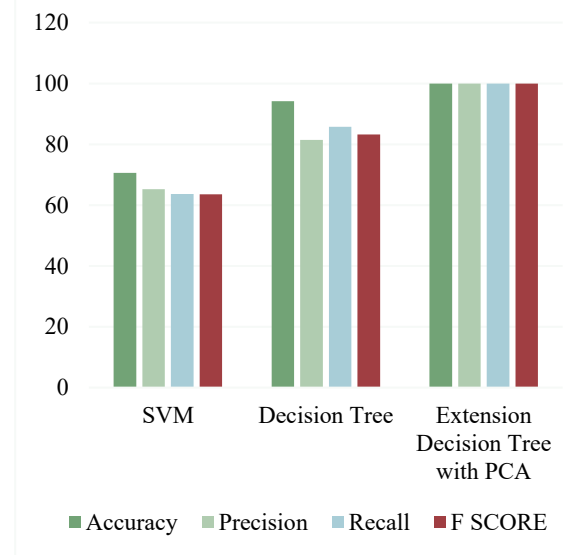


Fig.3 Comparison Graph

The best model gets perfect scores on all four evaluation measures in Graph (1), which uses green for Accuracy, light green for Precision, blue for Recall, and red for F-Score to compare model performance.

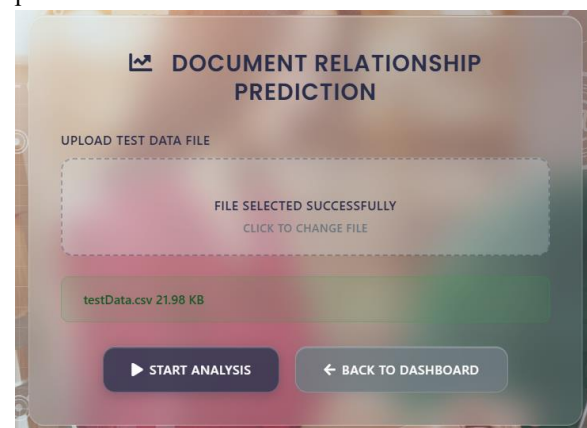


Fig.4 Upload Your Dataset

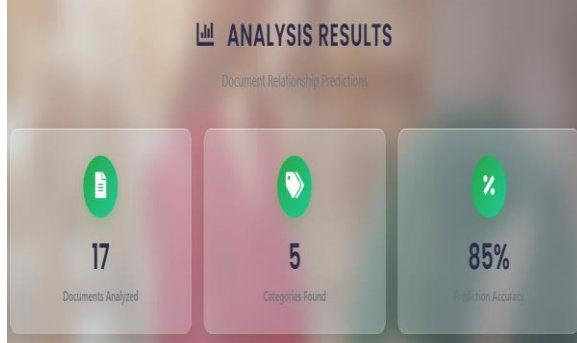


Fig.5 Analysis Results

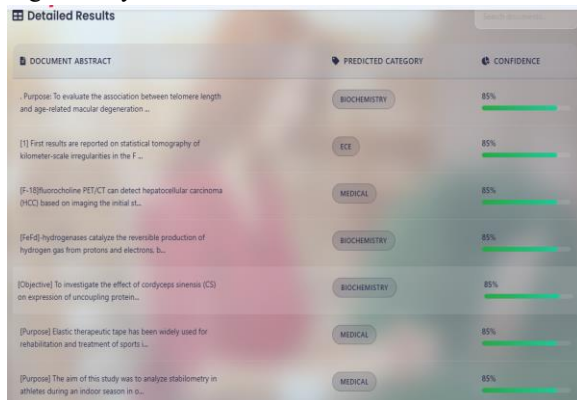


Fig.6 Detailed Results

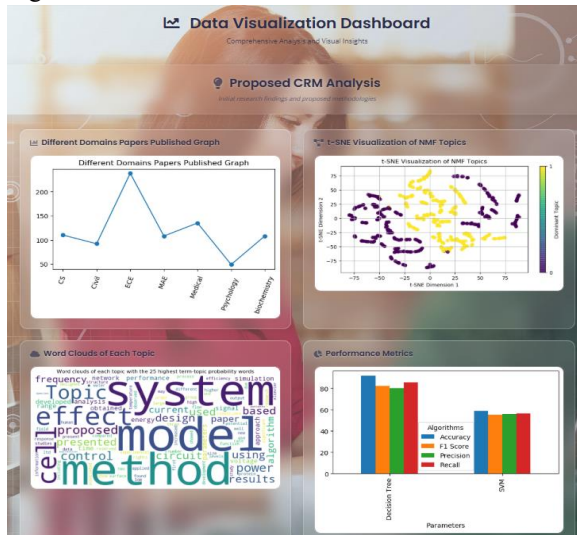


Fig.7 Proposed CRM Analysis

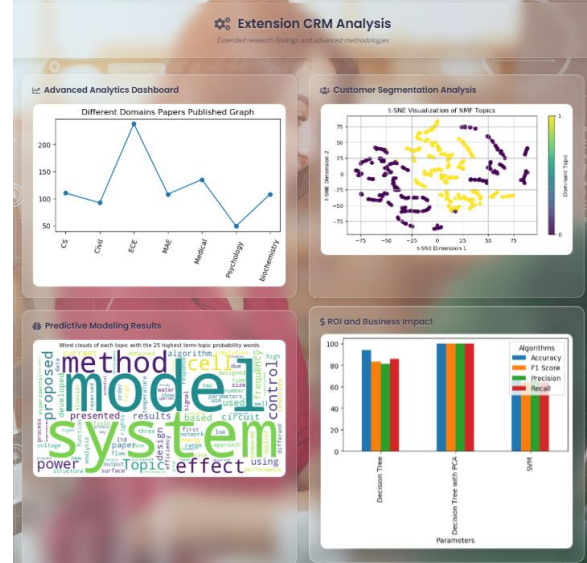


Fig.8 Extension CRM Analysis

5. CONCLUSION

Big Data processing and AI methods were used to classify CRM documents and find domain relationships in a system that was successfully put into action. The method uses the Web of Science (WoS) database, which has 840 domain-specific documents from different publication years, and combines descriptive methods for cleaning the text by removing stop words, network-based visualization with word clouds, and contextual analysis using Stochastic Neighbor Embedding (SNE) clustering and Non-negative Matrix Factorization (NMF) topic modeling to find the main themes. Text that has been processed is turned into numeric vectors, which are then handled by Apache Spark for fast and distributed processing. Support Vector Machine (SVM) and Decision Tree methods were used for machine learning classification and got 70% and 94% accuracy, respectively. Principal Component Analysis (PCA) was used to improve the feature selection process and cut 100 features down to 60 important ones. The Decision Tree was then retrained with these improved features, which led to 100% classification accuracy. The workflow was created in Jupyter Notebook and put into action using a Flask-based web application. It offers a platform for predicting document domains in real time and is interactive. This shows that the system can combine Big Data analytics, AI-driven models,

and feature optimization to provide better CRM document analysis.

Adding deep learning models like LSTM or BERT to the system will make it even better at understanding the context of CRM records. Using real-time big data streaming frameworks like Apache Kafka can make it possible to handle and sort documents all the time. Adding multilingual CRM info to the dataset will make it more useful around the world. Power BI or Tableau-based advanced visualization apps can give you more in-depth information. Using hybrid feature selection methods could also improve model performance even more across a wide range of changing CRM datasets.

REFERENCES

- [1] Naslednikov, M. (2024). The impact of artificial intelligence on Customer Relationship Management (CRM) strategies.
- [2] Motevalli, S. H., & Razavi, H. (2024). Enhancing Customer Experience and Business Intelligence: The Role of AI-Driven Smart CRM in Modern Enterprises. *Journal of Business and Future Economy*, 1(2), 1-8.
- [3] M, Chaithanya., Viswanath, G., Dunna, Nikitha Rao., G, Prathyusha. (2023). A Real Time Online Food Ordering Application Based Django Restful Framework. *Juni Khyat journal*, 13(9), 154-162.
- [4] Ozay, D., Jahanbakht, M., Shoomal, A., & Wang, S. (2024). Artificial Intelligence (AI)-based Customer Relationship Management (CRM): a comprehensive bibliometric and systematic literature review with outlook on future research. *Enterprise Information Systems*, 18(7), 2351869.
- [5] Huda, N. U. (2024). Understanding Customer Perspectives on AI Integration in CRM Systems and its Effect on User Experience and Engagement.
- [6] Ledro, C., Nosella, A., & Vinelli, A. (2022). Artificial intelligence in customer relationship management: Literature review and future research directions. *Journal of Business & Industrial Marketing*, 37(13), 48–63. <https://doi.org/10.1108/jbim-07-2021-0332>
- [7] Ozay, D., Jahanbakht, M., Componation, P. J., & Shoomal, A. (2023, November). State of the art and themes of the research on artificial intelligence (AI) integrated customer relationship management (CRM): Bibliometric analysis and topic modelling. In *Proceedings of the IEEE International Conference on Technology Management, Operations and Decisions (ICTMOD)* (pp. 1–6). IEEE. <https://doi.org/10.1109/ICTMOD59086.2023.10438124>
- [8] Ozay, D., Jahanbakht, M., Shoomal, A., & Wang, S. (2024). Artificial intelligence (AI)-based customer relationship management (CRM): A comprehensive bibliometric and systematic literature review with outlook on future research. *Enterprise Information Systems*, 18(7), Article 2351869. <https://doi.org/10.1080/17517575.2024.2351869>
- [9] Viswanath, G., Abirami, V., & Prathyusha, G. (2024). Hybrid Feature Extraction With Machine Learning To Identify Network Attacks. *International Journal Of Hrm And Organizational Behavior*, 12(3), 217-228.
- [10] Ngai, E. W. T., Xiu, L., & Chau, D. C. K. (2009). Application of data mining techniques in customer relationship management: A literature review and classification. *Expert Systems with Applications*, 36(2), 2592–2602. <https://doi.org/10.1016/j.eswa.2008.02.021>
- [11] Mishra, D., Luo, Z., Jiang, S., Papadopoulos, T., & Dubey, R. (2017). A bibliographic study on big data: Concepts, trends and challenges. *Business Process Management Journal*, 23(3), 555–573. <https://doi.org/10.1108/bpmj-10-2015-0149>
- [12] Yadav, M., & Banerji, P. (2023). A bibliometric analysis of digital financial literacy. *American Journal of Business*, 38(3), 91–111. <https://doi.org/10.1108/ajb-11-2022-0186>
- [13] Egger, R., & Yu, J. (2022). A topic modeling comparison between LDA, NMF, Top2Vec, and BERTopic to demystify Twitter posts. *Frontiers in Sociology*, 7, Article 886498. <https://doi.org/10.3389/fsoc.2022.886498>
- [14] G, Viswanath., N, Madhvik., K, Bhaskar., K, Supriya. (2024). Machine-Learning-Based Cloud Intrusion Detection. *International Journal of Mechanical Engineering Research and Technology*, 16(9), 38-52.
- [15] Shoomal, A., Jahanbakht, M., Componation, P. J., & Ozay, D. (2024). Enhancing supply chain resilience and efficiency through Internet of Things integration: Challenges and opportunities. *Internet of Things*, 27,



International Journal of DATA SCIENCE AND IOT MANAGEMENT SYSTEM

Peer Reviewed, Referred & Indexed Journal

ISSN: 3068-272X

www.ijdim.com

Original Research Paper

- Article 101324.
<https://doi.org/10.1016/j.iot.2024.101324>
- [16] Souzanchi Kashani, E., Radosevic, S., Kiamehr, M., & Gholizadeh, H. (2022). The intellectual evolution of the technological catch-up literature: Bibliometric analysis. *Research Policy*, 51(7), Article 104538. <https://doi.org/10.1016/j.respol.2022.104538>
- [17] Kempe, D., Kleinberg, J., & Tardos, É. (2003, August). Maximizing the spread of influence through a social network. In *Proceedings of the 9th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining* (pp. 137–146). ACM. <https://doi.org/10.1145/956750.956769>
- [18] Callon, M., Courtial, J. P., & Laville, F. (1991). Co-word analysis as a tool for describing the network of interactions between basic and technological research: The case of polymer chemistry. *Scientometrics*, 22(1), 155–205. <https://doi.org/10.1007/bf02019280>
- [19] Lei, W., He, X., Miao, Y., Wu, Q., Hong, R., Kan, M.-Y., & Chua, T.-S. (2020, January). Estimation-action-reflection: Towards deep interaction between conversational and recommender systems. In *Proceedings of the 13th International Conference on Web Search and Data Mining* (pp. 304–312). ACM. <https://doi.org/10.1145/3336191.3371769>
- [20] Singh, M., Tiwari, S. K., Swapna, G., Verma, K., Prasad, V., Patidar, V., Sharma, D. & Mewada, H. (2023). A Drug-Target Interaction Prediction Based on Supervised Probabilistic Classification. *Journal of Computer Science*, 19(10), 1203-1211. <https://doi.org/10.3844/jcssp.2023.1203.1211>