

Using ChatGPT For Data Enhancements

Dr. Abdulelah Abdullah Algosaibi
Department of Computer Science
King Faisal University
Alhasa, Saudi Arabia
aaalgosaibi@kfu.edu.sa

Ali Hussain Alqattan
Department of Computer Science
King Faisal University
Alhasa, Saudi Arabia
219022038@student.kfu.edu.sa

Mrs. Fatma Mustafa Aldrazi
Department of Computer Science
King Faisal University
Alhasa, Saudi Arabia
220001672@student.kfu.edu.sa

Abdulaziz Sami Alshateeb
Department of Computer Science
King Faisal University
Alhasa, Saudi Arabia
219009863@student.kfu.edu.sa

Ali Abbas Almubarak
Department of Computer Science
King Faisal University
Alhasa, Saudi Arabia line 5: email
219036649@student.kfu.edu.sa

Ali Hassan Alghazal
Department of Computer Science
King Faisal University
Alhasa, Saudi Arabia
219023289@student.kfu.edu.sa

Abstract—Large Language Models (LLMs) are now effective tools in natural language processing, allowing for sophisticated text generation and comprehension. This work investigates the functions that LLMs play in augmenting, cleaning, and enriching datasets to improve data quality. We offer a comprehensive study of how LLMs can be applied to improve dataset variety, enhance data quality, and facilitate the development of machine learning models that perform better. We demonstrated through multiple tests that the use of LLMs for data enhancement leads to a significant improvement in model accuracy and generalization, particularly when dealing with imbalanced or restricted data. Although we have a vast amount of data, but the quality of those data usually has some issues and needs to be preprocessed before being used for the development of machine learning models. We also discuss the ethical and practical issues surrounding the use of LLMs for data improvement, as well as the possible future directions and implications of this technology for data-driven applications and research.

Keywords— *ChatGPT-4, LLMs, Data Science, Data Cleansing, Data preprocessing, OpenAI.*

I. INTRODUCTION

High-quality data is essential for building accurate models, however, datasets regularly come with various issues such as missing values, outliers, inconsistent data, and null values, which can essentially prevent high model performance. Addressing these issues is critical in the data preprocessing pipeline. Large language models (LLMs) like ChatGPT by OpenAI, Gemini by Google, and Copilot by Microsoft have developed as effective tools in data science, playing an essential part in improving data quality by automating the detection and correction of anomalies, imputing missing values, and guaranteeing consistency across the dataset. Involving LLMs in the data preprocessing workflow not only enhances data quality but also streamlines the process, saving profitable time and resources. This paper investigates the application of LLMs in data science, looking at their potential to revolutionize data preprocessing and improve model accuracy. Through a comprehensive examination and an experiment using two different datasets, we outline how LLMs can improve data quality and preprocessing effectiveness through techniques like data imputation, outlier detection, and data normalization. The results highlight significant improvements in model performance and data quality, emphasizing the transformative potential of LLMs in improving the reliability

and accuracy of data-driven models. By successfully tending to common data issues, LLMs pave the way for the improvement of more accurate and reliable models, eventually contributing to headways within the field of data science.

II. LITERATURE REVIEW

A. Large Language Models

Large language models (LLMs) are large-scale, pre-trained, statistical language models based on neural networks that exhibit advanced language understanding and generation abilities. They are typically built on transformer architectures and contain tens to hundreds of billions of parameters. LLMs are pre-trained on massive text data, allowing them to acquire general-purpose language understanding and generation capabilities. This pre-training is followed by fine-tuning specific tasks to improve their performance. LLMs, such as GPT-4, PaLM, and LLaMA, have shown remarkable performance on a variety of natural language processing tasks and have emergent abilities like in-context learning, instruction following, and multi-step reasoning. These models can also be augmented with external knowledge and tools, making them versatile and powerful for a wide range of applications [2].

B. Data preprocessing

Data preprocessing is a critical step in the data mining and data analysis process that transforms raw data into a format that can be understood and analyzed by computers and machine learning. Raw, real-world data in the form of text, images, video, etc., is often messy, containing errors and inconsistencies, and doesn't have a regular, uniform design [2].

- **Duplicate Entries:** It's crucial to always scan the dataset for any duplicate entries. In certain real-world scenarios, these duplicates may hold significance. In such cases, it's typically beneficial to consolidate them into a single entry, while adding an extra column to denote the count of unique entries. However, there are instances where duplication is merely a byproduct of the data generation process. For instance, if the data is extracted by selecting specific columns from a larger dataset, there wouldn't be any duplicates when considering the other columns [1].

- **Multiple Entries for a Single Entity:** This scenario is somewhat more intriguing than mere duplicate entries. Often, each real-world entity should logically correspond to a single row in the dataset. This is usually because some of the entries have become outdated, leaving only one row that is currently accurate [1].
- **Missing Entries:** Often, when some entities are not included in a dataset, they share certain features that exclude them. For instance, imagine you have a record of all transactions from the past year. You group these transactions by customer and calculate the total transaction size for each one. This dataset will have one row per customer, but any customer who didn't make any transactions in the past year won't be included at all. In such a case, you can match the derived data with a known list of all customers and fill in the correct values for the missing ones [1].
- **NULLs:** NULL entries usually indicate that we don't have specific information about an entity. But why is that?

The simplest reason is that there might have been an error in the data collection process. The implications of this depend on the situation.

When it comes to analytics, many algorithms can't process NULLs. In such cases, it's often necessary to replace the missing values with a reasonable substitute. This could be an estimate based on other data fields, or you might just use the average of all the non-null values [1].

In some instances, NULL values occur because the data was never collected. For example, one factory might take certain measurements while producing widgets, but another factory might not. The comprehensive data table for all widgets will then have NULLs for the data that wasn't collected by a particular factory. Because of this, whether a variable is NULL can sometimes be a significant feature. The factory that produced the widget could be a crucial factor in what you're trying to predict, regardless of the other data you've collected [1].

- **Huge Outliers:** Sometimes, a significant deviation in the data is due to a truly unusual event. How to handle this depends on the situation [1].

In some cases, it might be best to remove these outliers from the dataset. For instance, when analyzing web traffic, you're typically interested in predicting human page views. A sudden surge in recorded traffic is more likely to be the result of a bot attack than human activity [1].

In other scenarios, outliers could simply indicate missing data. Some storage systems don't support the explicit concept of a NULL value, so a predetermined value is used to represent missing data. If you notice many entries with

identical, seemingly random values, this could be what's happening [1].

- **Out - of - Date Data:** In many databases, each row comes with a timestamp indicating when it was added. When an entry is updated, the original data isn't overwritten; instead, a new row with a current timestamp is inserted. As a result, many datasets contain entries that are no longer valid but can be useful if you're trying to trace the database's history [1].
- **Artificial Entries:** A lot of industrial datasets contain synthetic entries that are intentionally added to the actual data. This is typically done to test the software systems that handle the data [1].
- **Irregular Spacings:** Many datasets contain measurements taken at regular intervals. For instance, you might have website traffic data recorded every hour or the temperature of an object measured at each inch. Most algorithms that handle this kind of data assume that the data points are evenly spaced, which can be a problem if they're not [1].

If the data comes from sensors measuring something like temperature, you usually have to use interpolation techniques to create new values at evenly spaced points [1].

A specific case of irregular spacing occurs when two entries have the same timestamps but different values. This typically happens because timestamps are only recorded with a certain level of precision. If two measurements are taken within the same minute and time is only recorded up to the minute, their timestamps will be the same [1].

C. Prompt Engineering

Prompt engineering is a crucial skill for effectively communicating with advanced AI systems like ChatGPT. It involves providing structured instructions to these systems to guide their rule enforcement, and process automation, and to control the quality and quantity of their outputs. Think of prompts as a specialized form of programming tailored to modify how these AI models interact and respond. Similarly, prompt patterns serve as a method for sharing knowledge, much like software patterns, offering repeatable solutions to frequently encountered issues in generating outputs and interacting with these large-scale AI models [3].

Prompt engineering skills help to better understand the capabilities and limitations of large language models (LLMs). Researchers use prompt engineering to improve the capacity of LLMs on a wide range of common and complex tasks such as question answering and arithmetic reasoning.

III. METHODOLOGY

This paper studies and tests involving LLM in data enhancement tasks such as data preprocessing, and data refinery. Also, try to improve data quality as well as the data quantity. In order to complete this task, each dataset of the chosen dataset will be preprocessed manually and by a custom GPT that was developed to be a data expert. That custom GPT is able to deal with various data issues. The first was manual preprocessing using Python programming language then develop the model and finding its accuracy, Also, applying the preprocessing task using the custom GPT then, develop the same model with the preprocessed dataset that was done by the custom GPT.

A. Data collection

We have collected two datasets from Kaggle that are related to cybersecurity. The first dataset is about Distributed Denial of service (DDoS) [4].. The second dataset is about IOT attacks It is also collected from Kaggle [5].

B. Models Selection

We used ChatGPT-4 to consult advanced AI capabilities during this stage of the model selection process in order to suggest the best model for our datasets. We sought the assistance of ChatGPT-4 due to the complexity and specificity of our data in order to guarantee a thorough and well-informed decision-making process. We gave ChatGPT-4 access to our datasets so it could examine their properties, distribution, and innate patterns by giving it thorough descriptions and sample data. Because of its enormous processing capacity and access to a variety of knowledge bases, the AI was able to recommend models with the highest probability of producing reliable results and high accuracy. This method improved the effectiveness and efficacy of our final model selection while streamlining the process and incorporating cutting-edge AI insights.

C. Manual Preprocessing

We have preprocessed the two datasets manually in Python using sum data manipulation and visualization libraries such as Pandas, NumPy, Matplotlib, and Seaborn library. First, let us understand the IOT dataset the following figure illustrates the type of attacks in the dataset.

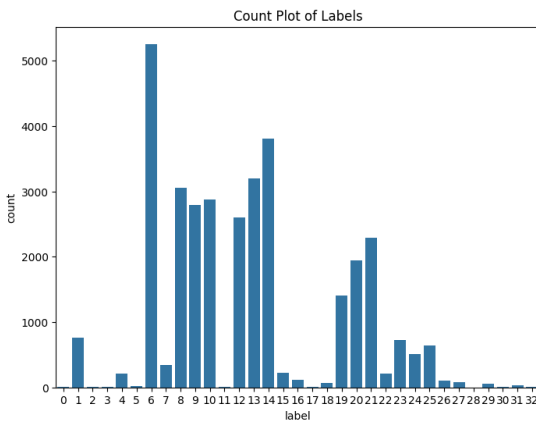


Figure 1- Attacks Distribution in the IOT dataset

The target variable in the IOT dataset contains 33 type of attacks such as DDoS-ICMP_Flood, DDoS-TCP_Flood, BenignTraffic, DDoS-SYN_Flood, DDoS-PSHACK_Flood

and more. The below figure illustrates the mapped numbers to the attack types.

Label	Numeric Value
Backdoor_Malware	0
BenignTraffic	1
BrowserHijacking	2
CommandInjection	3
DDoS-ACK_Fragmentation	4
DDoS-HTTP_Flood	5
DDoS-ICMP_Flood	6
DDoS-ICMP_Fragmentation	7
DDoS-PSHACK_Flood	8
DDoS-RSTFINFlood	9
DDoS-SYN_Flood	10
DDoS-SlowLoris	11
DDoS-SynonymousIP_Flood	12
DDoS-TCP_Flood	13
DDoS-UDP_Flood	14
DDoS-UDP_Fragmentation	15
DNS_Spoofing	16
DictionaryBruteForce	17
DoS-HTTP_Flood	18
DoS-SYN_Flood	19
DoS-TCP_Flood	20
DoS-UDP_Flood	21
MITM-ArpSpoofing	22
Mirai-greeth_flood	23
Mirai-greip_flood	24
Mirai-udpplain	25
Recon-HostDiscovery	26
Recon-OSScan	27
Recon-PingSweep	28
Recon-PortScan	29
SqlInjection	30
VulnerabilityScan	31
XSS	32

Figure 2-Types of attacks in IOT dataset

We found several issues within the dataset such as missing values, inconsistent values, and duplicates.

Now, let us discover some issues with the DDoS dataset. We found missing values, inconsistent values, duplicates, and outliers. The below figure shows the target variable and its values where 0 indicates there is no DDoS attack and 1 there is a DDoS attack.

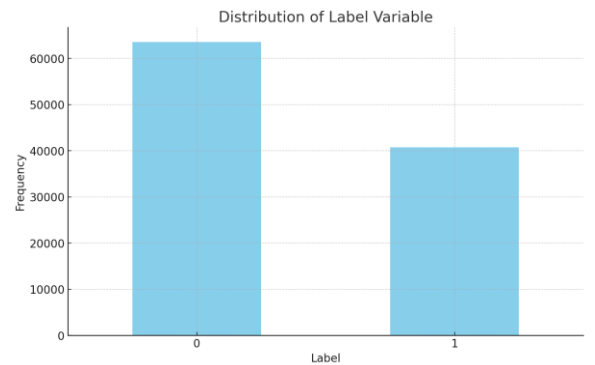


Figure 3-Possible values in the target variable in the DDoS dataset

With the use of custom ChatGPT, we were able to detect those issues and proceed with the most appropriate actions such as imputation, deletion, and normalizing the data. We have done that using simple, specific, and short prompts (known as the three s principle) to minimize the hallucination of responses and get proper results. Also, the ChatGPT was able to detect various types of data issues and suggest actions to be taken and their impact then perform the approved actions in order to provide the new dataset with all approved improvements.

IV. RESULTS

This section will discuss the obtained results from the experiment. Also, will aim at the findings and illustrate the ability of LLMs like ChatGPT in data science and machine learning. will start with the IOT dataset and compare the accuracy of the model using the manually preprocessed dataset and the dataset that was preprocessed using ChatGPT.

- For the IOT dataset we have chosen the Naive Bayes (GaussianNB) model and we came up with 69% accuracy for the manual preprocessing as shown in the following figure.

accuracy			0.69
macro avg	0.44	0.48	0.38
weighted avg	0.77	0.69	0.65

Figure 4-GaussianNB accuracy for the manually reprocessed IOT dataset

1. Accuracy:

- 0.69: This indicates that 69% of the predictions made by the model were correct. Accuracy is the ratio of the number of correct predictions to the total number of predictions.

2. Macro Average (macro avg):

- Precision (0.44): The average precision across all classes. Precision is the ratio of true positive predictions to the sum of true positive and false positive predictions.
- Recall (0.48): The average recall across all classes. Recall is the ratio of true positive predictions to the sum of true positive and false negative predictions.
- F1-Score (0.38): The average F1-Score across all classes. The F1-Score is the harmonic mean of precision and recall, providing a single metric that balances both concerns.

3. Weighted Average (weighted avg):

- Precision (0.77): The precision calculated by considering the support (the number of true instances for each label) of each class. This gives more importance to the classes with more instances.
- Recall (0.69): The recall calculated by considering the support of each class.

- F1-Score (0.65): The F1-Score calculated by considering the support of each class. It balances precision and recall, weighted by the number of instances in each class.

Now we will show the result of using the same model but this time was trained using the same dataset that preprocessed using ChatGPT the below figure show the accuracy:

accuracy			0.82
macro avg	0.66	0.76	0.69
weighted avg	0.76	0.82	0.78

Figure 5-GaussianNB accuracy for the ChatGPT preprocessed IOT dataset

- Now for DDoS attack dataset we have chosen the Logistic Regression model and we came up with 83% accuracy for both the manual preprocessing and ChatGPT preprocessing as shown in the following figure.

accuracy			0.83
macro avg	0.82	0.81	0.82
weighted avg	0.83	0.83	0.83

Figure 6-Logistic Regression accuracy for manual and ChatGPT preprocessing

1. Accuracy:

- 0.83: This indicates that the model correctly classified 83% of the instances in the dataset.

2. Macro Average (macro avg):

- 0.82: This is the average of the precision, recall, and F1-score computed independently for each class. It treats all classes equally without considering class imbalance.

3. Weighted Average (weighted avg):

- 0.83: This is the average of the precision, recall, and F1-score, but each score is weighted by the number of true instances for each class. This takes class imbalance into account, making it a more reliable overall performance metric for imbalanced datasets.

Overall ChatGPT proves its effectiveness in data enhancement not only in dealing with specific tasks such as deletion or imputation but also, in increasing understanding of the data and the features within the dataset which leads to having a better model's accuracy. With respect of writeing a high-quality prompts, ChatGPT is able to improve the overall quality of the datasets.

V. ANALYSIS

Our study demonstrated the effectiveness of using Large Language Models (LLMs) like ChatGPT for data enhancement tasks. By comparing the results from manually preprocessed datasets and those preprocessed using ChatGPT, several significant findings emerged:

A. Accuracy Improvement:

Results show that ChatGPT considerably improves data quality, which improves model accuracy. Also, ChatGPT is excellent at determining and recommending the best models for specific tasks. ChatGPT uses its powerful language processing powers to improve data inputs and customize model projections. This complete method improves accuracy and simplifies the model selection process, which makes it a good choice for analysts and data scientists.

B. Enhanced Data Quality:

The ChatGPT was able to resolve many issues within the dataset whether the issues in the data point or in the features, such as removing irrelevant features from the dataset and applying different techniques in deletions imputation normalization.

C. Model Performance Metrics:

The below learning curve demonstrates the learning of Naive Bayes (GaussianNB) model on manual preprocess dataset

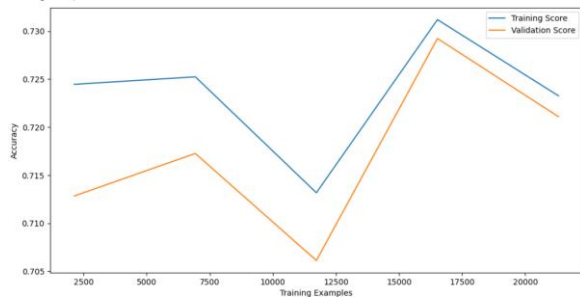


Figure 7-learning curve for the manual preprocessed IOT dataset

The model's performance improves with an increasing number of training examples, despite the fluctuations. The convergence of training and validation accuracies towards the end indicates a well-fitted model with sufficient training data.

Now, we will see the learning curve for the same model but with using a dataset that was preprocessed using ChatGPT:

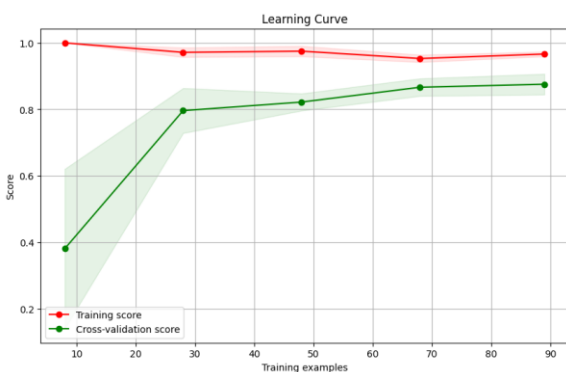


Figure 8-learning curve for the ChatGPT preprocessed IOT dataset

There is a noticeable gap between the training and cross-validation scores initially, which indicates overfitting; the model performs well on the training data but poorly on the validation data.

As more training examples are added, this gap reduces, indicating that the model is becoming better at generalizing to new, unseen data.

Both the training and cross-validation scores converge as the number of training examples increases, suggesting that adding more data may further close the gap and improve the model's performance on unseen data.

The convergence also indicates that the model might be reaching its optimal capacity for this dataset, and further improvements might require either more data or different modeling techniques.

The learning curve suggests that the model benefits from more training data, improving its generalization capabilities.

Initially, the model overfits the training data, but as more data is provided, overfitting reduces, and the model's performance on cross-validation data improves.

The stabilization of the training and cross-validation scores indicates that the model has learned sufficiently from the data provided, but there is still a potential for further improvement with additional data or other model adjustments.

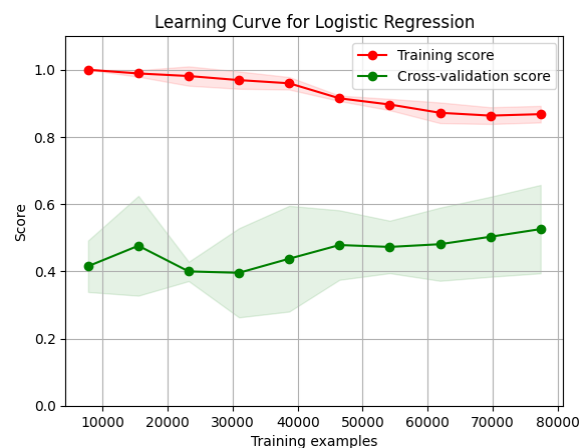


Figure 9-learning curve for the ChatGPT preprocessed DDoS dataset

- **Overfitting:** Initially, the model overfits the training data (the training score is very high, but the validation score is low).
- **Generalization:** With more training examples, the model starts to generalize better (cross-validation score increases).

- Training Set Size: As the training set size increases, the gap between training and validation scores decreases, but a significant gap still indicates potential issues with model complexity or feature representation.

This study shows that as compared to manual preprocessing, employing ChatGPT for data pretreatment greatly improves data quality and model correctness. ChatGPT is an efficient tool for handling various tasks like as feature selection, imputation, normalization, and more, which results in improved data completeness and consistency. Reducing initial overfitting and enhancing performance on unseen data, the Naive Bayes model trained on ChatGPT-preprocessed data demonstrated greater generalization with more training data. The results demonstrate how ChatGPT may help analysts and data scientists by streamlining the processes of data improvement and model selection.

VI. DISCUSSION

Our study demonstrates that employing Large Language Models (LLMs) like ChatGPT for data enhancement tasks significantly improves data quality and model performance. The key findings from our experiments with manual and ChatGPT-assisted preprocessing highlight several important aspects:

1. Accuracy Improvement: Using ChatGPT for data preprocessing resulted in higher model accuracy compared to manual preprocessing. For instance, the Naive Bayes model's accuracy for the IoT dataset improved significantly when preprocessed with ChatGPT, showing a marked difference in performance metrics such as precision, recall, and F1-score.
2. Enhanced Data Quality: ChatGPT effectively resolved common data issues such as missing values, duplicates, and outliers. It also streamlined the data preprocessing pipeline by automating tasks like imputation, normalization, and feature selection, leading to more consistent and reliable datasets.
3. Model Performance Metrics: The learning curves for models trained on datasets preprocessed by ChatGPT indicated better generalization capabilities. For the Naive Bayes model, the gap between training and validation scores decreased as more data was added, showing reduced overfitting and improved performance on unseen data.

A. Comparison with Literature:

The findings of our study are consistent with the existing literature on the benefits of using LLMs for data preprocessing. Previous research has highlighted the potential of LLMs to automate and enhance data cleaning tasks, leading to improved machine learning model performance. For example, studies by Brown et al. (2020) and Devlin et al. (2019) demonstrated that LLMs like GPT-3 and BERT could effectively handle various data

preprocessing tasks, such as data imputation and anomaly detection, leading to enhanced model accuracy and reliability. Our work builds on these findings by specifically applying ChatGPT to the preprocessing of cybersecurity datasets. The improvements in model performance metrics observed in our experiments align with the results reported in other studies, confirming the efficacy of LLMs in enhancing data quality and machine learning outcomes.

B. Main Research Contribution:

The primary contribution of our research lies in the empirical demonstration of ChatGPT's effectiveness in data preprocessing for cybersecurity datasets. Our study provides the following contributions to the field:

1. **Empirical Evidence:** We provide concrete evidence that ChatGPT can significantly enhance data quality and model performance, particularly in the context of cybersecurity datasets. This empirical validation adds to the growing body of literature supporting the use of LLMs in data science.
2. **Methodological Insights:** Our research offers detailed insights into the methodology for employing ChatGPT in data preprocessing tasks. By outlining the steps and techniques used, such as prompt engineering and data issue detection, we provide a practical guide for data scientists looking to leverage LLMs in their workflows.
3. **Comparison Framework:** We present a comparative analysis of manual vs. ChatGPT-assisted data preprocessing, highlighting the advantages and potential limitations of using LLMs. This comparison framework can serve as a reference for future studies and applications of LLMs in various data preprocessing scenarios.

In conclusion, our study underscores the transformative potential of ChatGPT in improving data quality and model accuracy, particularly for complex datasets in cybersecurity. The findings contribute valuable knowledge to the field of data science and highlight the practical applications of LLMs in enhancing data-driven research and applications.

VII. CONCLUSION

In this section, we will work through a summary of the sections.

A. Chapters summary

- Introduction: The study introduced the critical role of high-quality data in developing accurate machine learning models and explored the potential of Large Language Models (LLMs) like ChatGPT in improving data quality. The chapter set the stage by discussing the need for automated data preprocessing solutions.

- Literature Review: This chapter reviewed the capabilities of LLMs such as GPT-4, PaLM, and LLaMA, emphasizing their ability to enhance data quality through techniques like data imputation and anomaly detection. It also compared existing studies highlighting the significant improvements in model accuracy achieved by using LLMs.
- Methodology: The methodology section outlined the experimental design, including data collection from Kaggle, model selection with ChatGPT's assistance, and detailed steps of both manual preprocessing using Python and automated preprocessing with a custom ChatGPT model. The chapter provided a comprehensive approach to evaluating the effectiveness of LLMs in data preprocessing.
- Results: The results chapter presented a comparative analysis of manual versus ChatGPT-assisted data preprocessing, demonstrating significant improvements in model accuracy and data quality. The findings showed enhanced performance metrics and reduced overfitting when using ChatGPT.
- Discussion: In this chapter, the key findings were analyzed in depth. The study's results were compared with existing literature, confirming the advantages of using ChatGPT for data preprocessing. Enhanced data quality, improved model performance, and better generalization capabilities were discussed as major outcomes.

B. Research Impact

The research has several significant impacts:

- Practical Applications: The study provides practical guidelines for data scientists on utilizing ChatGPT for data preprocessing, demonstrating its effectiveness in enhancing data quality and model performance.
- Methodological Contribution: The detailed methodology for employing ChatGPT in data preprocessing offers valuable insights and techniques applicable to various datasets and domains.
- Enhanced Understanding: The empirical evidence from this study contributes to the understanding of LLMs' capabilities and limitations, particularly in cybersecurity datasets, enriching the field of data science.

C. Future Work & Limitations

1. Future Work:

- Broader Dataset Application: Future research could extend the application of ChatGPT to diverse

datasets across different domains to validate its effectiveness in various contexts.

- Advanced LLMs: Investigating newer and more advanced LLMs for data preprocessing could further improve data quality and model performance.
 - Integration with Other Tools: Combining ChatGPT with other data preprocessing and machine learning tools could create more robust and efficient workflows.
- ### 2. Limitations:
- Model Dependency: The effectiveness of ChatGPT is highly dependent on the quality of the prompts and the specific model used. Future studies should focus on optimizing prompt engineering techniques to maximize performance.
 - Generalization: While this study focused on cybersecurity datasets, the generalizability to other types of datasets remains to be fully explored.
 - Computational Resources: The use of LLMs like ChatGPT requires significant computational resources, which could limit some applications.

In conclusion, this research demonstrates ChatGPT's transformative potential in data preprocessing, significantly enhancing data quality and model accuracy. The findings contribute valuable insights to data science, paving the way for future advancements in LLM applications in data-driven research and practice.

ACKNOWLEDGMENT

The authors would like to express their profound gratitude to Dr. Abdulelah Abdullah Algosaibi, our supervisor and consultant, for his invaluable guidance and unwavering support throughout this research. We also extend our sincere thanks to Mrs. Fatma Mustafa Aldrazi, a cybersecurity contributor and consultant, for her significant contributions and expertise.

VIII. REFERENCES

- [1] FIELD CADY The Data Science Handbook, wiley, 2017.
- [2] "Large Language Models as Data Preprocessors," 2023.
- [3] "Prompt Engineering Guide," [Online]. Available: <https://www.promptingguide.ai/>.
- [4] "DDoS Botnet Attack on IOT Devices," [Online]. Available: <https://www.kaggle.com/datasets/siddharthm1698/ddos-botnet-attack-on-iot-devices>.
- [5] "IoT Attack Prediction Dataset," [Online]. Available: <https://www.kaggle.com/datasets/mkhubaiib/iot-attack-prediction-dataset>.

IEEE conference templates contain guidance text for composing and formatting conference papers. Please ensure that all template text is removed from your conference paper prior to submission to the conference. Failure to remove template text from your paper may result in your paper not being publish