

1. Table des matières

4.	Chapter 4.....	59
4.1	Introduction.....	59
4.2	Evaluation Metrics for Assessing Model Performance	59
4.2.1	Confusion Matrix.....	59
4.2.2	Accuracy.....	60
4.2.3	Precision	60
4.2.4	Recall	61
4.2.5	F1 Score.....	61
4.3	Model Performance Analysis.....	61
4.4	Comparative Analysis with Other Approaches.....	63
4.5	Model Performance Evaluation on New Data	64
4.6	Conclusion	64
5.	Bibliographie	65

4. Chapter 4

Test and Evaluation

4.1 Introduction

The detection of SQL injections using AI techniques, such as machine learning and deep learning algorithms, has gained the interest of many researchers in this field. These techniques have been shown to be effective at identifying SQL injection attacks with high accuracy.

In this chapter, we will discuss the test and evaluation of our model for detecting SQL injection attacks. We will use a variety of metrics, including accuracy, precision, recall, and F1 score. We will also compare the performance of our model to other machine learning algorithms and related works.

4.2 Evaluation metrics for assessing model performance

In the field of machine learning and classification tasks, evaluation metrics play a crucial role in assessing the performance of models. These metrics provide quantitative measures that help us understand the accuracy, effectiveness, and reliability of model predictions. When evaluating the performance of classification models, it is essential to examine the appropriate evaluation metrics that provide insights into their strengths and weaknesses. In this section, we will explore some of the most commonly used evaluation metrics that provide valuable insights into the performance of classification models.

4.2.1 Confusion matrix

The confusion matrix provides a tabular representation of the model's predictions against the actual labels. It allows us to visualize the distribution of true positives, true negatives, false positives, and false negatives, providing valuable insights into the model's performance [26].

	Positive Prediction	Negative Prediction
Positive Class	True Positive (TP)	False Negative (FN)
Negative Class	False Positive (FP)	True Negative (TN)

TP (True Positives): Correctly predicted positive instances.

TN (True Negatives): Correctly predicted negative instances.

FP (False Positives): Incorrectly predicted positive instances.

FN (False Negatives): Incorrectly predicted negative instances.

4.2.2 Accuracy

Accuracy is a commonly used evaluation metric that measures the overall correctness of model predictions. It calculates the ratio of correct predictions to the total number of instances. Accuracy provides a general overview of the model's performance across all classes [26].

$$\text{Accuracy} = \text{Correct Predictions} / \text{Total Predictions}$$

4.2.3 Precision

Precision focuses on the proportion of correctly identified positive predictions (true positives) out of the total positive predictions made by the model. It helps assess the model's ability to minimize false positives [27].

$$\text{Precision} = \text{TruePositive} / (\text{TruePositive} + \text{FalsePositive})$$

4.2.4 Recall

Recall, also known as sensitivity or true positive rate, measures the proportion of true positive predictions captured by our model out of the total actual positive instances. It reflects the model's ability to minimize false negatives [27].

$$\text{Recall} = \text{TruePositive} / (\text{TruePositive} + \text{FalseNegative})$$

4.2.5 F1 Score

The F1 score is a combined metric that balances precision and recall. It provides a harmonic mean of these two measures and offers a comprehensive evaluation of the model's performance [28].

$$\text{F-Measure} = (2 * \text{Precision} * \text{Recall}) / (\text{Precision} + \text{Recall})$$

By analyzing these evaluation metrics, we can gain a deeper understanding of how our classification model performs and identify areas for improvement. These metrics provide valuable insights into the model's strengths and weaknesses, allowing us to make informed decisions to enhance its performance.

4.3 Model performance analysis

In this section, we evaluate the performance of our SQL injection detection model based on the BERT architecture. As already discussed the model was trained on a dataset containing 22,599 samples, consisting of both SQL injections instances and normal SQL queries.

We assess the model's performance using a variety of evaluation metrics, including accuracy, precision, recall, F1 score, and the confusion matrix. These metrics provide valuable insights into the model's ability to correctly classify SQL injection instances and non-malicious queries.

The performance results obtained are as follows:

Accuracy: The model achieved an accuracy of **99.98%** during training and validation, indicating its high level of accuracy in classifying the queries.

F1 Score: The F1 score, which considers both precision and recall, also reached an impressive value of **99.98%**. This indicates a balance between correctly identifying SQL injection attacks (precision) and capturing all actual SQL injection instances (recall).

Precision: The precision of the model, which measures the proportion of true positive predictions out of all positive predictions, was **99.96%**. This indicates a very low rate of false positives, meaning that the model maintains a high level of confidence in its SQL injection detection predictions.

Recall: The recall of the model, which measures the proportion of true positive predictions out of all actual positive instances, was **100%**. This indicates the model's ability to capture nearly all instances of SQL injection attacks, resulting in a low rate of false negatives.

The confusion matrix provides a more detailed breakdown of the model's performance:

	Positive Prediction	Negative Prediction
Positive Class	2284	0
Negative Class	1	2235

The confusion matrix reveals that out of the 2284 positive instances (normal queries), the model correctly identified 2284 instances as positive (true positives). It also correctly classified 2235 out of 2236 negative instances as negative (SQL injections) while misclassifying one instance as positive (false positive).

These results are consistent with the findings of Srishti Lodha and Atharva Gundawar from the Department of Computer Science and Engineering at Vellore Institute of Technology [29], who made a similar study using the BERT architecture for SQL injection detection. Their research demonstrated comparable performance and highlighted the effectiveness of the BERT model in accurately identifying SQL injection attacks.

Overall, the performance results demonstrate the high accuracy, precision, recall, and F1 score of our SQL injection detection model. These results, in alignment with the work of Srishti Lodha and Atharva Gundawar, underscore the effectiveness of the BERT model for accurately detecting SQL injection attacks while minimizing false positives and false negatives.

4.4 Comparative analysis with other approaches

In this section, we present a comparative analysis of our BERT-based SQL injection detection model with other commonly used approaches available in the literature. While we didn't evaluate the performance of the alternative approaches ourselves, we have collected and compiled information from various sources on their reported performance. By comparing our model with these approaches, we aim to provide insights into the effectiveness of our BERT-based model for SQL injection detection. The comparison is based on commonly used evaluation metrics such as accuracy, precision, recall, and F1 score. The findings of this analysis are summarized in the following table.

Model name	Training Accuracy	Validation Accuracy	Precision	Recall	F1
KNN	100%	99.12%	98.85%	99.52%	99.18%
SVM	92.78%	92.44%	90.21%	96.36%	93.19%
BERT	99.99%	99.98%	99.96%	100%	99.98%

The compared approaches were trained on a dataset of approximately 42,000 data points, while our dataset consisted of 22,599 samples. This disparity arises from the fact that we obtained the original dataset from the same resource, but we had to perform data cleaning and preprocessing to ensure its quality.

The comparison shows that BERT performed the best among the compared models in terms of key evaluation metrics, including accuracy, precision, recall, and F1 score. These results

affirm the effectiveness of BERT in detecting SQL injection attacks, highlighting its superior performance in our study.

4.5 Model performance evaluation on new data

In this section, we assess the performance of our deep learning model, trained on a specific dataset, on a new dataset obtained from Kaggle called "sqli". This dataset consists of 1100 samples and serves as a valuable benchmark to evaluate our model's capabilities on unseen data.

We started by examining the overlap between the training data and the new dataset. Through careful analysis, we identified 90 common sentences shared between the two datasets. We also quantified the similarity between the training data and the new dataset, revealing a remarkable similarity score of 9.48%. This score highlights the percentage of sentences present in the training dataset that also appear in the test dataset, underscoring the model's proficiency in handling comparable contexts and retaining its predictive power across different datasets.

Ultimately, the model's performance on the new data is outstanding, boasting an accuracy of 99.73%. This accuracy metric signifies the percentage of correct predictions made by the model when tested on the new dataset. The exceptionally high accuracy demonstrates the robustness and reliability of our deep learning model in identifying SQL injections, even when confronted with previously unseen data.

4.6 Conclusion

In conclusion, our model for detecting SQL injection attacks has shown impressive effectiveness. Through testing and evaluation using various metrics, we have demonstrated its accuracy in identifying SQL injection attacks also the comparison with other models further validates its superior performance.

5. Bibliographie

- [26] M. & L. G. Sokolova, A systematic analysis of performance measures for classification tasks, 2009.
- [27] D. M. Powers, "Evaluation: From precision, recall and F-measure to ROC, informedness, markedness and correlation," *Journal of Machine Learning Technologies*, 2011.
- [28] C. J. Van Rijsbergen, Information Retrieval (2nd ed), 1979.
- [29] S. L. a. A. Gundawar, SQL Injection and Its Detection Using Machine Learning Algorithms, 2023.