

COMPARISON OF SUPPORT VECTOR MACHINE WITH NAIVE BAYES CLASSIFIER METHOD FOR CLASSIFICATION OF LECTURER PERFORMANCE

REMIANUS TUNTI

Department of Computer Science
Dili Institute of Technology, Dili, Timor-Leste

ARÃO TERNORIO ALVES SANTOS

Department of Computer Science
Dili Institute of Technology, Dili, Timor-Leste

JACOB SOARES

Department of Computer Science
Dili Institute of Technology, Dili, Timor-Leste

ANGELINA SOARES

Department of Computer Science
Dili Institute of Technology, Dili, Timor-Leste

REGINA XIMENES VIEGAS

Department of Computer Science
Dili Institute of Technology, Dili, Timor-Leste

ABSTRACT

Classification is a supervised learning technique that extracts models from training data to identify predefined categories for new test data, and high accuracy is crucial for accurate predictions. This research aims to compare the Support Vector Machine (SVM) and Naïve Bayes Classifier (NBC) methods for classifying lecturer performance. This comparative study utilizes datasets from three units, namely data from research scorecard evaluations, teaching scorecards, and lecturer evaluations from students. The comparison results employ a confusion matrix table to determine the accuracy, precision, and recall of the dataset. The results of comparing 150 Research Scorecard datasets by dividing 70% of the data for training and 30% for testing using the NBC and SVM methods yielded identical accuracy, precision, and recall results, all at 100%. Meanwhile, the comparison of 150 lecturer evaluations from student datasets using the NBC method resulted in higher accuracy, precision, and recall, all at 100%, compared to using the SVM method, which achieved an accuracy of 93.18%, precision of 72.04%, and recall of 75%.

KEYWORDS: classification, SVM, NBC, lecturer performance.

INTRODUCTION

Classification is a supervised learning technique that extracts models from training data to identify predefined categories for new test data (Morán-Fernández, 2021). It involves three phases: formation, validation, and testing (Tharwat, 2021). High accuracy is crucial for accurate predictions (Zhou, 2014), but performance bias poses a challenge (Soleymani, 2019)

and imbalanced datasets are commonly encountered in real-life applications (Jedrzejowicz, 2021).

Several studies on the performance evaluation of classifiers in data mining such as analysis of teacher performance using multiple classifiers (Ahmad & Rashid, 2016; Kumar Pal & Pal, 2013), classifications on instructor performance (Agaoglu, 2016), to predict instructor performance (Ahmed, 2016), sentiment analysis to classify student-lecturer comments (Rakhmanov, 2020), sentiment analysis and opinion mining on educational data (Shaik, 2023) and experimental comparison of multilabel methods (García-Pedrajas, 2024).

The Naive Bayes Classifier (NBC) and Support Vector Machine (SVM) are widely used classification techniques. The NBC algorithm is a popular choice for big data analysis due to its efficient structure, while Naive Bayes is a probabilistic classifier that uses the Bayesian theorem with a strong assumption of independence (Chen et al, 2021; Perez, 2021). Bayes theorem is crucial for inferential statistics and advanced machine learning models as it updates hypothesis probability based on new evidence (Berrar, 2018) and this model achieves higher classification accuracy with less complexity (Salmi and Rustam, 2019). Whereas SVMs are highly effective and reliable algorithms for regression and classification across various fields (Cervantes et al, 2020) and effectively tackles the challenges of large data categorization, especially in multidomain applications in large data environments (Suthaharan, 2016). For problems involving regression and classification, SVMs may provide both linear and nonlinear decision boundaries (Somvanshi et al, 2016). Traditional SVMs are based on identifying a hyperplane in a higher-dimensional space that effectively divides various data classes (Amaya-Tejera et al, 2024).

To increase classification performance, numerous approaches are commonly used, such as study conducted by Rahman (2018) on the application of feature selection with information gain for document classification, selection of accurate and significant features in attack detection system alerts on computer networks (Alhaj, 2016), improve classification performance in the credit scoring problem (Jadhav, 2018) and the study about the significance of feature selection in achieving classification accuracy in bank marketing datasets (Prasetyo et al., 2021).

Thus, the aim of this research is to compare the SVM and NBC methods for classifying lecturer performance at the Dili Institute of Technology. This comparative study utilizes datasets from three units, namely data from research scorecard evaluations, teaching scorecards, and lecturer evaluations from students. The comparison results employ a confusion matrix table to determine the accuracy, precision, and recall of the dataset.

RELATED WORK

Classification is a supervised machine learning model with the objective of predicting categorical class labels for new instances based on previous observations (Sadiq et al., 2020). The classification process consists of two main phases: model development for training and model evaluation using testing data (Jalota & Agrawal, 2019b). In the principal three phases of the classification process, namely formation phase, validation phase, and testing phase, various steps are involved (Tharwat, 2021).

Several studies related to classification model include performance analysis of lecturers with Multiple Classifiers at Kurdistan-Iraq University (Ahmad & Rashid, 2016), evaluation of teaching quality with the Flipped Classroom model in colleges and universities (Fu & Li, 2022), application of multiclassification models to evaluate teaching quality in the art department (Hua et al., 2022), and evaluation of English teaching quality using online analytical processing by combining classification algorithms at the college and university levels in China (Zhang et al., 2022).

Some commonly used methods in supervised machine learning include Support Vector Machine (SVM) and Naive Bayes Classifier (NBC). SVM is a powerful machine

learning technique that minimizes structural risk (Roy & Chakraborty, 2023). It serves as an algorithm for classification and regression tasks, offering advanced capabilities and parameter optimization (Cervantes et al., 2020). SVM is widely used in data mining due to its efficiency, generalizability, and ability to find optimal solutions (Gaye et al., 2021). Its applications include face detection, handwriting recognition, and various real-world scenarios (Ghosh et al., 2019).

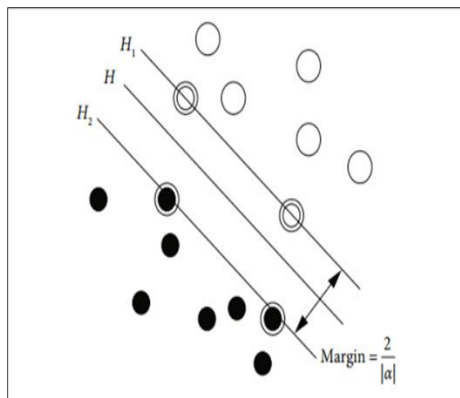


Figure 2.1. SVM Schematic Diagram (Fu & Li, 2022)

Meanwhile, the NBC) method are used to calculate the highest probability value as a classification process (Atmadja et al., 2020). NBC is potentially good at serving as a classification model due to its simplicity and accuracy (Dangi et al, 2014). Naive Bayes algorithm is one of the most effective methods in the field of text classification, but only in the large training sample set can it get a more accurate result (Y. Huang and L. Li, 2011). Naive Bayes is one of the most well-known data mining algorithms for classification. Naive Bayes is a simple and effective learning theory that does not need various parameters (Ramadhani et al, 2021).

$$P(H|X) = \frac{P(H|X)P(H)}{P(X)}$$

Where X data with undetermined classes, H hypothesis that data X belongs to a specific class, P(H|X) probability of hypothesis H given condition X (posterior probability), P(H) probability of hypothesis H (prior probability), P(X|H) probability of X given condition H and P(X) probability of X

One way to improve model accuracy is through feature selection techniques. One commonly used technique is through information gain. Information gain is a method of feature evaluation that is widely used in the field of machine learning (Lei, 2012). It is a technique for feature selection that can reduce the size of a given feature by optimizing each attribute's value and providing a relative increase for that particular feature (Zareapoor & Seeja, 2015). Information is typically used in a variety of applications and is based on entropy metrics. The previously discussed beta can be used to determine relevance and reduce growth (Cherrington et al., 2019).

RESEARCH METHODOLOGY

The total dataset used in this research consists of 450 datasets, comprising 150 datasets of research scorecard with 69 attributes each, 150 datasets of teaching scorecard with 92 attributes each, and 150 datasets of lecturer evaluations from students with 35 attributes each. These datasets were obtained from the CARPS-CS, CEQA, and Academic Department of DIT.

These attributes will be assigned their codes, standardized as per standard format, and any missing values will be completed. The total attributes after these processes amount to 101 attributes ready to undergo the feature selection process, with 38 attributes from Research Score Card, 37 from Teaching Score Card, and 26 from Student Evaluation.

attributes, in the teaching scorecard dataset is 37 attributes, and in the evaluation of lecturers by student dataset is 26 attributes.

Data splitting as a general approximation used for model validation, the dataset will be split into two parts: training data and testing data. This model will be trained using the training data and validated using the testing data (Joseph, 2022). The total estimation of datasets used for this research will be divided into 70% for data training and 30% for data testing. From the total of 150 datasets for each unit, 104 will be allocated for data training and 46 will be allocated for data testing. This scaling utilization is also applied in research conducted by (Abbi Nizar Muhammad et al., 2019), which combines the NBC method with SVM and demonstrates its superior accuracy level and strong performance.

In this phase, we will measure the performance of the Naïve Bayes Classifier and Support Vector Machine methods in predicting accuracy and relevancy for the Research Score Card, Teaching Score Card, and Student Evaluation datasets. To conduct this performance evaluation, we will use the Confusion Matrix as a method for accurate calculation, based on the concept of data mining. This formula calculates various outputs such as Accuracy, Precision, and Recall. Related research on this accuracy test has already been conducted by researchers (Shahi et al., 2018) to test the accuracy of classifying Nepali news using the Naïve Bayes Classifier, Support Vector Machine, and Neural Networks methods. Additionally, research conducted by (Ma et al., 2020) has tested Precision and Recall for the classification of spam emails using the Naïve Bayes Classifier and Support Vector Machine methods.

Table 1
Confusion Matrix for Multiclass

| | | Predicted Number | | | |
|---------------|---------|------------------|----------|-----|----------|
| | | Class 1 | Class 2 | ... | Class n |
| Actual Number | Class 1 | X_{11} | X_{12} | ... | X_{1n} |
| | Class 2 | X_{21} | X_{22} | ... | X_{2n} |
| | . | . | . | . | . |
| | . | . | . | . | . |
| | Class n | X_{n1} | X_{n2} | ... | X_{nn} |

RESULT

To perform performance testing on both methods using a dataset of 150 instances, the dataset will be divided into 70% for training and 30% for testing. The dataset will be divided into each unit as follows: For the Research Score Card, 104 instances will be used for training and the remaining 46 for testing. For the Evaluation of Lecturers by Students, 106 instances will be used for training and the remaining 44 for testing, and for the Teaching Score Card, 105 instances will be used for training and the remaining 45 for testing.

The three datasets will undergo testing using the Rapid Miner Studio 102 platform to analyze the comparison results of the Naïve Bayes Classifier and Support Vector Machine methods. The aim of this testing is to determine the levels of accuracy, precision, and recall of both methods. Below are the comparison results of Accuracy, Precision, and Recall for the NBC and SVM methods for the Research Score Card, Evaluation of Lecturers by Students, and Teaching Score Card.

Comparison Results for NBC and SVM Methods for the Research Score Card Dataset

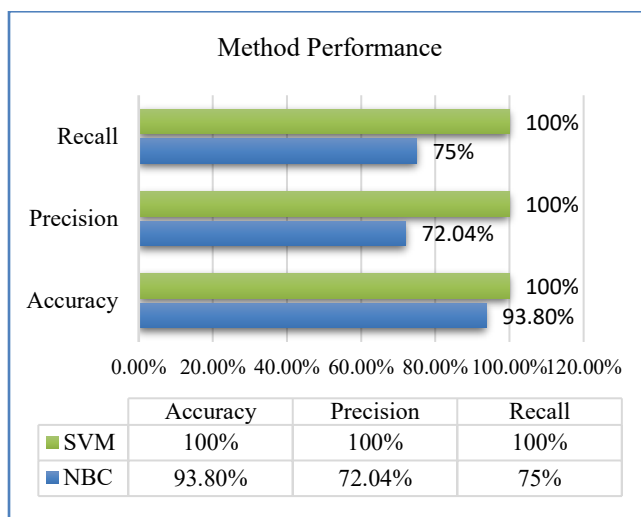


Figure 2
Comparison Results for NBC and SVM Methods for the Research Score Card Dataset

Comparison of the Accuracy, Precision, and Recall results from the NBC Model for the Research Score Card dataset shows that the prediction for 46 test instances achieves a true accuracy rate of 100%, indicating that the accuracy level of this model is excellent. The precision testing results also indicate a value of 100%, meaning that the predictions for this test data align with the actual data. Similarly, the recall testing results indicate a value of 100%, signifying that the information obtained from the prediction results of this model is highly accurate. Meanwhile, the testing results of the SVM model also demonstrate a true accuracy rate of 100% for the Research Score Card dataset, indicating that the accuracy level of this model is excellent. The precision testing results also show a value of 100%, indicating that the predictions for this test data match the actual data. Likewise, the recall testing results show a value of 100%, indicating that the information obtained from the prediction results of this model is highly accurate.

Comparison Results for NBC and SVM Methods for the Evaluation of Lecturers by Students Dataset

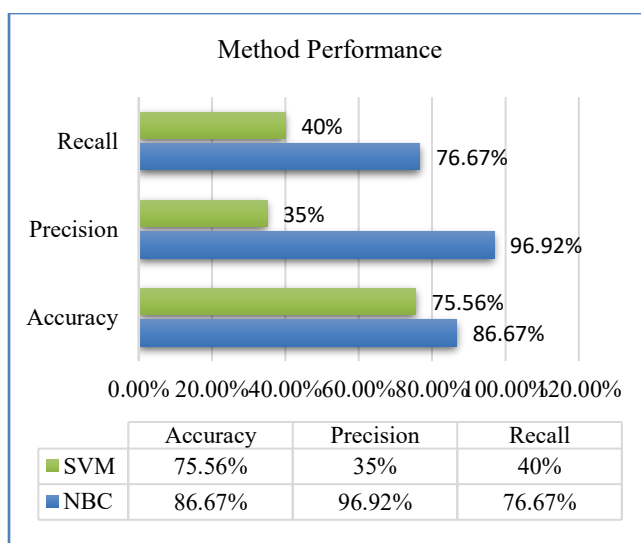


Figure 3
Comparison Results for NBC and SVM Methods for the Evaluation of Lecturers by Students Dataset

Comparison of the Accuracy, Precision, and Recall results from the NBC model for the Student Evaluation Dataset show that the prediction for 44 test instances of Student Evaluation achieves an accuracy rate of 93.18%. The precision testing results indicate a value of 72.04%, meaning that the predictions for this test data have 27.96% correct predictions within the actual class. Meanwhile, the recall testing results show a value of 75%, indicating that 25% of the information obtained from the prediction results does not match this model. On the other hand, the testing results of the SVM model show an accuracy rate of 100%, indicating that the accuracy level of this model is excellent. The precision testing results also indicate a value of 100%, meaning that the predictions for this test data match the actual data. Similarly, the recall testing results show a value of 100%, indicating that the information obtained from the prediction results of this model is highly accurate.

Comparison Results for NBC and SVM Methods for the Teaching Score Card Dataset

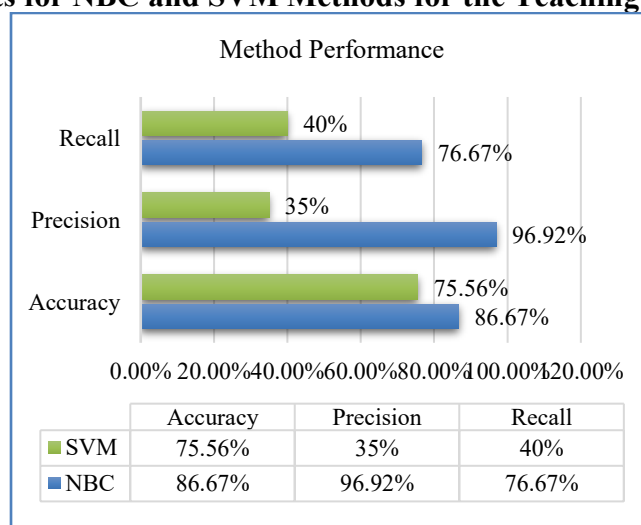


Figure 4

Comparison Results for NBC and SVM Methods for the Teaching Score Card Dataset

Comparison of the Accuracy, Precision, and Recall results from the NBC model for the Teaching Score Card dataset show that the prediction for 45 test instances achieves an accuracy rate of 86.67%. The precision testing results indicate a value of 96.92%, meaning that the predictions for this test data have 3.08% correct predictions within the actual class. Meanwhile, the recall testing results show a value of 76.67%, indicating that 23.33% of the information obtained from the prediction results does not match this model. Meanwhile, the testing results of the SVM model show an accuracy rate of 75.56%. The precision testing results indicate a value of 35%, meaning that the predictions for this test data have 65% correct predictions within the actual class. Similarly, the recall testing results show a value of 40%, indicating that 60% of the information obtained from the prediction results does not match this model.

The Comparison Results of Performance Between the NBC and SVM methods

Comparison results for the test data using Naive Bayes and SVM methods to evaluate faculty performance indicate variations in prediction across three scenarios: For the first scenario with the Research Score Card test data consisting of 46 instances, both methods demonstrate excellent prediction. Naive Bayes achieves an accuracy, precision, and recall of 100%, as does SVM. In the second scenario with Evaluation of Lecturers by Students test data of 44 instances, SVM outperforms Naive Bayes with perfect accuracy, precision, and recall scores of 100%, compared to Naive Bayes' accuracy of 93.18%, precision of 72.04%, and recall of 75%. In the third scenario, assessing Teaching Score Card test data with 45 instances, Naive Bayes yields superior predictions with an accuracy of 86.67%, precision of 96.92%, and recall of 76.67%, compared to Naive Bayes' accuracy of 75.56%, precision of

35%, and recall of 40%. These comparisons reveal that the choice of method can significantly impact prediction performance across different evaluation scenarios.

DISCUSSION

The use of feature selection with information gain in this research did not have a significant effect on increasing accuracy, precision and recall as in research conducted by (Omuya et al, 2021), where the results of applying feature selection with information gain using the NBC method had an effect on increasing accuracy from 94.89% rose to 97.81%, precision from 95% rose to 97.80% and recall from 94.90% rose to 97.80%. Likewise, feature selection with information gain using the SVM method, where there was an increase in accuracy from 67.77% to 100%, precision from 50% to 100% and recall percentage from 63% to 100%. Similar research was conducted by (Vijayashree & Sultana, 2018) where the use of feature selection with information gain using the NBC method had a significant effect on increasing the accuracy percentage from 79.35% to 82.65%. However, using the SVM method experienced a decrease in accuracy from 75.23% down to 74.12%.

Some of the researchers' findings regarding factors that influence increasing accuracy, precision and recall, apart from feature selection, are the complexity and pattern of the dataset being tested. This is proven by increasing the number of datasets, both teaching scorecards and evaluation datasets from students, and research scorecards. The results of the Teaching Scorecard dataset test using the NBC method show that there is a significant effect on increasing the number of datasets from 150 to 1000 datasets with the data pattern in the form of a series of numbers from 0 to 4 (0,1,2,3,4). Using the number of data sets with these data patterns, the comparison results obtained in sequence, namely accuracy, precision and recall, were previously 86.67%, 96.92%, 76.67%, increasing to 98.00%, 98.10% and 98.00%. Likewise, the test results using the SVM method, obtained sequential comparison results, namely accuracy, precision and recall were previously 75.56%, 35.00%, 40.00%, increasing to 92.00%, 92.71% and 92.00%. Further evidence is also found in the evaluation dataset of students using the NBC method, showing that there is a significant effect on increasing the number of datasets from 150 to 1000 datasets with the data pattern in the form of a series of numbers from 1 to 5 (1,2,3,4,5). With the number of datasets and data patterns, the comparison results obtained sequentially, namely accuracy, precision and recall, previously were 93.18%, 72.04%, 75.00%, increasing to 100%, 100% and 100%. Test results using the SVM method, obtained sequential comparison results, namely accuracy, precision and recall were previously 100%, 100%, 100%, decreasing to 96.67%, 97.14% and 96.67%.

Specifically, for the Research Scorecard dataset, the factor that influences accuracy, precision and recall using both the NBC and SVM methods is increasing the number and pattern of the dataset, where the data pattern used previously was 0, 2, 3, 5, 7, 8, 10, 15, 20. In the tests carried out by researchers using the NBC method, significant results were obtained with the first scenario with a dataset of 150 and the data pattern used a series of numbers 0 to 4 (0,1,2,3,4) and obtained accuracy results, precision and recall respectively are 93.33%, 95.00% and 93.33%. The second scenario involves increasing 1000 datasets with better accuracy, precision and recall results respectively, namely 98.00%, 98.10%, 98.88%. In testing using the SVM method on data patterns with rows of numbers 0 to 4 with a dataset of 150, the accuracy, precision and recall test results were obtained sequentially, namely 73.33%, 76.00% and 73.33%. By increasing the number of datasets to 1000 datasets, test results obtained were accuracy of 92.33%, precision of 92.97% and recall of 92.33%.

CONCLUSION AND IMPLICATION

The results of comparing 150 Research Scorecard datasets by dividing 70% of the data for training and 30% for testing using the NBC and SVM methods yielded identical accuracy, precision, and recall results, all at 100%. Meanwhile, the comparison of 150 lecturer evaluations from student datasets using the NBC method resulted in higher accuracy,

precision, and recall, all at 100%, compared to using the SVM method, which achieved an accuracy of 93.18%, precision of 72.04%, and recall of 75%. Subsequently, the comparison results for 150 Teaching Scorecard datasets using the NBC method showed an accuracy of 86.67%, precision of 96.92%, and recall of 76.67%, whereas using the SVM method resulted in an accuracy of 75.56%, precision of 35%, and recall of 40%.

Based on the test results and findings related to this research, it is indicated that both methods do not exhibit a significant influence whether undergoing feature selection process or not towards the improvement of accuracy, precision, and recall. The findings observed in the Teaching Scorecard dataset and lecturer evaluations from students are influenced by the magnitude or insignificance of the dataset size, while in the research scorecard dataset, it is influenced by two factors: the number of datasets and the data pattern. To achieve satisfactory accuracy, precision, and recall results in dataset testing, it is necessary to establish a standard number of datasets with an appropriate data pattern, such as the following data patterns: 0 to 4 (0, 1, 2, 3, 4) or 1 to 5 (1, 2, 3, 4, 5) with a dataset count of 1000.

ACKNOWLEDGEMENTS

The author expresses gratitude to the Dili Institute of Technology (DIT) for the research grant and extends special thanks to the Center for Applied Research and Policies Studies and Community Services (CARPS-CS), the Center for Evaluation and Quality Assurance (CEQA), and the Academic Department for facilitating the entire research process.

REFERENCES

1. Abbi Nizar Muhammad, Saiful Bukhori, & Priza Pandunata. (2019). Proceedings, 2019 International Conference on Computer Science, Information Technology, and Electrical Engineering (ICOMITEE 2019): October 16th-17th 2019, Jember, Indonesia.
2. Amaya-Tejera, N., Gamarra, M., Vélez, J. I., & Zurek, E. (2024). A distance-based kernel for classification via Support Vector Machines. *Frontiers in Artificial Intelligence*, 7. <https://doi.org/10.3389/frai.2024.1287875>
3. Berrar, D. (2019). Bayes' Theorem and Naive Bayes Classifier. In *Encyclopedia of Bioinformatics and Computational Biology* (Vols. 1–3, pp. 403–412). Elsevier. <https://doi.org/10.1016/B978-0-12-809633-8.20473-1>
4. Cherrington, M., Airehrour, D., Lu, J., Thabtah, F., Xu, Q., & Madanian, S. (2019, October). Particle swarm optimization for feature selection: A review of filter-based classification to identify challenges and opportunities. In *2019 IEEE 10th Annual Information Technology, Electronics and Mobile Communication Conference (IEMCON)* (pp. 0523-0529). IEEE.
5. Firmahsyah, F., & Gantini, T. (2016). Penerapan metode content-based filtering pada sistem rekomendasi kegiatan ekstrakurikuler (Studi Kasus di Sekolah ABC). *Jurnal Teknik Informatika dan Sistem Informasi*, 2(3).
6. Jalota, C., & Agrawal, R. (2019b). Analysis of Educational Data Mining using Classification. *International Conference on Machine Learning, Big Data, Cloud and Parallel Computing (Com-IT-Con)*, India, 14th -16 thFeb 2019.
7. Joseph, V. R. (2022). Optimal ratio for data splitting. *Statistical Analysis and Data Mining*, 15(4), 531–538. <https://doi.org/10.1002/sam.11583>
8. Joseph, V. R. (2022). Optimal ratio for data splitting. *Statistical Analysis and Data Mining*, 15(4), 531–538. <https://doi.org/10.1002/sam.11583>
9. Lei, S. (2012). A feature selection method based on information gain and genetic algorithm. *Proceedings - 2012 International Conference on Computer Science and Electronics Engineering, ICCSEE 2012*, 2, 355–358. <https://doi.org/10.1109/ICCSEE.2012.97>

10. Ma, T. M., Yamamori, K., & Thida, A. (2020). A Comparative Approach to Naïve Bayes Classifier and Support Vector Machine for Email Spam Classification. 2020 IEEE 9th Global Conference on Consumer Electronics, GCCE 2020, 324–326. <https://doi.org/10.1109/GCCE50665.2020.9291921>
11. Mofizur Rahman, C., Afroze, L., Sultana Refath, N., & Shawon, N. (n.d.). Iterative Feature Selection Using Information Gain & Naïve Bayes for Document Classification.
12. Morán-Fernández, L., Bólon-Canedo, V., & Alonso-Betanzos, A. (2022). How important is data quality? Best classifiers vs best features. *Neurocomputing*, 470, 365–375. <https://doi.org/10.1016/j.neucom.2021.05.107>
13. Muhammad, A. N., Bukhori, S., & Pandunata, P. (2019, October). Sentiment analysis of positive and negative of youtube comments using naïve bayes–support vector machine (nbsvm) classifier. In 2019 International Conference on Computer Science, Information Technology, and Electrical Engineering (ICOMITEE) (pp. 199-205). IEEE.
14. Odhiambo Omuya, E., Onyango Okeyo, G., & Waema Kimwele, M. (2021). Feature Selection for Classification using Principal Component Analysis and Information Gain. *Expert Systems with Applications*, 174. <https://doi.org/10.1016/j.eswa.2021.114765>
15. Omuya, E. O., Okeyo, G. O., & Kimwele, M. W. (2021). Feature selection for classification using principal component analysis and information gain. *Expert Systems with Applications*, 174, 114765.
16. Rakhmanov, O. (2020). A Comparative Study on Vectorization and Classification Techniques in Sentiment Analysis to Classify Student-Lecturer Comments. *Procedia Computer Science*, 178, 194–204. <https://doi.org/10.1016/j.procs.2020.11.021>
17. Ramanda Hasibuan, M. (2019). Pemilihan Fitur dengan Information Gain untuk Klasifikasi Penyakit Gagal Ginjal menggunakan Metode Modified K-Nearest Neighbor (MKNN) (Vol. 3, Issue 11). <http://j-ptiik.ub.ac.id>
18. Rashid, T. A., & Ahmad, H. A. (2016). Lecturer performance system using neural network with Particle Swarm Optimization. *Computer Applications in Engineering Education*, 24(4), 629–638. <https://doi.org/10.1002/cae.21737>
19. Shahi, T. B., & Pant, A. K. (2018, February). Nepali news classification using Naive Bayes, support vector machines and neural networks. In 2018 international conference on communication information and computing technology (iccict) (pp. 1-5). IEEE.
20. Varghese, S. M., & Sushmitha, M. N. (2014). Efficient Feature Subset Selection Techniques for High Dimensional Data. *International Journal of Innovative Research in Computer and Communication Engineering*, 2(3).
21. Vijayashree, J., & Sultana, H. P. (2018). A Machine Learning Framework for Feature Selection in Heart Disease Classification Using Improved Particle Swarm Optimization with Support Vector Machine Classifier. *Programming and Computer Software*, 44(6), 388–397. <https://doi.org/10.1134/S0361768818060129>
22. Zareapoor, M., & K. R, S. (2015). Feature Extraction or Feature Selection for Text Classification: A Case Study on Phishing Email Detection. *International Journal of Information Engineering and Electronic Business*, 7(2), 60–65. <https://doi.org/10.5815/ijieeb.2015.02.08>