

A Hybrid Optimization of Supervised Learning Models using Information Gain-Based Feature Selection

NOVIA HASDYNA¹, ROZZI KESUMA DINATA²

¹Department of Informatics, Universitas Islam Kebangsaan Indonesia, Aceh, Indonesia

²Department of Informatics, Universitas Malikussaleh, Aceh, Indonesia

Corresponding author: Novia Hasdyna (e-mail: noviahasdyna@uniki.ac.id).

ABSTRACT This study aims to enhance the performance of supervised learning models in dermatology data classification through a hybrid approach that combines Information Gain-based feature selection with several established supervised learning algorithms, namely K-Nearest Neighbor (KNN), Support Vector Machine (SVM), and Naive Bayes. Utilizing the Dermatology dataset from the UCI Machine Learning Repository, consisting of 366 instances with 34 numeric attributes and 6 class labels, the research identifies attributes with the lowest Information Gain values, including Family History, Eosinophils in the infiltrate, and Hyperkeratosis. These attributes undergo dimensional reduction to expedite computation and improve model performance. The study evaluates the impact of dataset dimensionality reduction on the performance of the supervised learning algorithms, encompassing KNN, SVM, and Naive Bayes. Experimental results reveal a significant enhancement in the performance of supervised learning models. Specifically, the generated models achieve a True Positive Rate (TPR) of up to 82.52%, True Negative Rate (TNR) of 98.81%, Positive Predictive Value (PPV) of 33.55%, Negative Predictive Value (NPV) of 98.78%, and accuracy of 96.29% using the KNN algorithm. Furthermore, the utilization of SVM and Naive Bayes also yields significant improvements in model performance.

KEYWORDS Hybrid; Optimization; Information Gain; Supervised Learning; K-NN; SVM; Naïve Bayes

I. INTRODUCTION

THE rapid expansion of data alongside advancements in machine learning have sparked significant interest in refining the efficacy of supervised learning models [1-4]. Supervised learning, a fundamental approach in machine learning, involves training models on labeled data to make predictions or decisions based on input features [26]. Particularly, within various domains, such as text classification, medical diagnosis, and disease prediction, there lies substantial potential to leverage these advancements. However, datasets in these domains often suffer from high dimensionality and variability, posing unique challenges that necessitate effective feature selection techniques and robust classification methodologies [5-7].

Feature selection serves as a pivotal preprocessing step in machine learning, striving to identify the most pertinent attributes within a dataset [8]. This process assumes heightened significance in datasets with vast dimensions, where the inclusion of irrelevant or redundant features can precipitate overfitting, escalate computational expenses, and compromise model performance [9-14]. Among various

feature selection methods, those based on mutual information have gained prominence due to their ability to measure the statistical dependence between features and the target variable [15-18].

Recent research endeavors have explored diverse avenues to enhance machine learning models across various domains. For instance, Li et al. [19] showcased the effectiveness of a mutual information-based feature selection approach in tandem with a decision tree algorithm, yielding substantial enhancements in text classification accuracy and computational efficiency. Similarly, Chawla and Bhardwaj [20] provided a comprehensive review of feature selection methods in medical diagnosis, highlighting the challenges and opportunities in this domain. Gupta et al. [21] conducted a comparative study of enhanced classification techniques, demonstrating significant improvements in disease prediction accuracy. Ahmed et al. [22] compared the performance of machine learning algorithms for disease prediction, shedding light on the strengths and limitations of different approaches. Furthermore, recent reviews by Wang et al. [23], Zhou et al. [24], and Liu et al. [25] have discussed the application of

machine learning techniques in medical image analysis, electronic health records, and medical imaging, respectively, emphasizing the importance of feature selection in optimizing model performance.

Building upon previous research in supervised learning, this study aims to contribute to the advancement of machine learning methodologies. By enhancing the classification performance of supervised learning models through a hybrid approach, which combines Information Gain-based feature selection with established supervised learning algorithms like KNN, SVM, and Naive Bayes, the research assesses the impact of dimensionality reduction on model performance. Specifically, the study investigates how excluding attributes with low mutual information values influence computational efficiency and classification accuracy across different domains, including text classification, medical diagnosis, and disease prediction. The initial experimental findings suggest a significant enhancement in the performance of supervised learning models. By optimizing the selection of informative features, this hybrid optimization approach offers promising avenues for developing precise and efficient models conducive to various applications, including medical diagnosis and disease prediction.

II. METHODOLOGY

This section presents the methodology used to optimize the performance of supervised learning models for dermatology data classification. The approach combines Information Gain-based feature selection with well-established supervised learning algorithms, namely K-Nearest Neighbor (KNN), Support Vector Machine (SVM), and Naïve Bayes. The subsequent subsections provide a detailed explanation of the dataset, feature selection process, supervised learning algorithms, experimental setup, and performance evaluation.

A. DATASET DESCRIPTION

The dataset utilized in this study is sourced from the Dermatology dataset available in the UCI Machine Learning Repository. It comprises 366 instances, each characterized by 34 numeric attributes and 6 class labels. The dataset serves as the basis for training and evaluating the supervised learning models in dermatology data classification as shown in Table 1

Table 1. Dermatology Dataset

Description	Count
Total Instances	366
Numeric Attributes	34
Class Labels	6

B. PROPOSED MODEL

The proposed model in this study integrates Information Gain-based feature selection with three established supervised learning algorithms: K-Nearest Neighbor (KNN), Support Vector Machine (SVM), and Naive Bayes. This hybrid approach aims to enhance the performance of dermatology data classification by optimizing feature selection and leveraging the strengths of different classification techniques. The feature selection process begins by identifying attributes with the lowest Information Gain values, indicating their limited relevance to the target variable. These attributes undergo dimensional reduction to streamline computation and improve model efficiency. Subsequently, the reduced feature set is fed into each supervised learning algorithm for model

training and classification.

The performance of the proposed model is evaluated using various performance metrics, including accuracy, True Positive Rate (TPR), True Negative Rate (TNR), Positive Predictive Value (PPV), and Negative Predictive Value (NPV). By comparing the performance of the proposed model against baseline models and individual algorithms, the effectiveness of the hybrid approach can be assessed and validated. The following diagram illustrates the flow of the proposed model, as depicted in Figure 2.

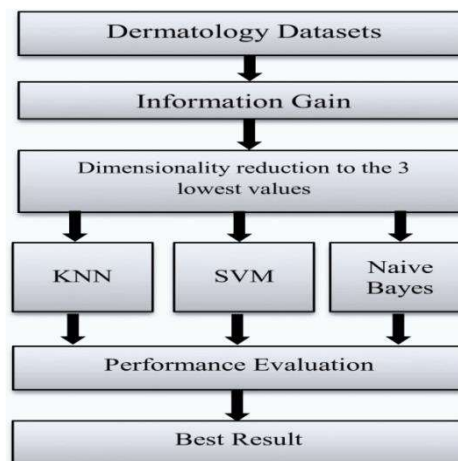


Figure 1. The newly proposed method

In Figure 1, the proposed methodology for this research is delineated, consisting of several sequential steps. Initially, the dermatology dataset is incorporated into the system. Subsequently, the calculation of information gain becomes the focal point of the second step. The ensuing step involves reducing the dataset's dimensionality by selecting features with the least information gain values. This process of dimensionality reduction occurs three times, namely one-dimensional, two-dimensional, and three-dimensional reductions, each employing progressively lower information gain values. Proceeding, the fourth step engages in the classification task utilizing supervised learning methods such as K-Nearest Neighbors (K-NN), Support Vector Machine (SVM), and Naive Bayes. Finally, the fifth step encompasses the performance assessment of these classification algorithms, evaluating various metrics including True Positive Rate (TPR), True Negative Rate (TNR), Positive Predictive Value (PPV), Negative Predictive Value (NPV), False Negative Rate (FNR), False Positive Rate (FPR), False Discovery Rate (FDR), False Omission Rate (FOR), Accuracy (ACC), F-Measure, and Matthews Correlation Coefficient (MCC).

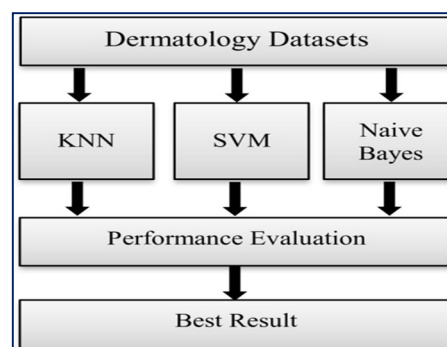


Figure 2. The conventional method

Figure 2 illustrates the conventional method used in this study. The first step involves inputting the dermatology dataset. The second step entails performing classification using supervised learning methods, namely K-NN, SVM, and Naive Bayes. The fifth step is to measure the performance of these classification algorithms, including True Positive Rate (TPR), True Negative Rate (TNR), Positive Predictive Value (PPV), Negative Predictive Value (NPV), False Negative Rate (FNR), False Positive Rate (FPR), False Discovery Rate (FDR), False Omission Rate (FOR), Accuracy (ACC), F-Measure, and Matthews Correlation Coefficient (MCC). In this proposed conventional method, dimensionality reduction techniques are not employed.

C. INFORMATION GAIN-BASED FEATURE SELECTION

Information gain is a prominent feature selection method widely employed by researchers to determine the significance threshold of an attribute. The value of information gain is obtained by subtracting the entropy value before separation from the entropy value after separation. This measurement is primarily utilized as an initial step to determine which attributes will be retained for use in classification algorithms. The process of selecting features using information gain consists of three stages:

1. Calculate Information Gain for Each Attribute in the Original Dataset: In this stage, the information gain value is computed for each attribute in the original dataset.
2. Determine the Desired Threshold: Researchers set a threshold, allowing attributes with a weight equal to or greater than the threshold to be retained, while attributes below the threshold are discarded.
3. Refine the Dataset by Reducing Attributes: This stage involves reducing the dataset by removing attributes. The measurement of attributes is described as:

$$\text{info}(D) = -\sum_{i=1}^m p_i \log_2(p_i), \quad (1)$$

where D is the set of cases, M is the number of partitions of D , and p_i is the proportion of D_i with respect to D . p_i represents the probability of a tuple in D belonging to class C_i and is estimated by $|C_i, D|/|D|$. The logarithm function used here is base 2 because the information is encoded based on bits. The calculation of entropy value after separation can be done using the following formula:

$$\text{info}_A(D) = -\sum_{j=1}^v \frac{|D_j|}{|D|} \times I(D_j), \quad (2)$$

where D represents the set of cases, A denotes the attribute. v is the number of partitions of attribute A . $|D_j|$ signifies the number of cases in partition j . $|D|$ represents the total number of cases in D . $I(D_j)$ indicate the total entropy within the partition.

$$\text{Gain}(A) = I(D) - I(A), \quad (3)$$

where $\text{Gain}(A)$ denotes the information of attribute A , $I(D)$ represents the total entropy. $I(D, A)$ signifies the entropy of attribute A .

D. K-NEAREST NEIGHBOR

K-Nearest Neighbors (KNN) is an instance-based machine

learning algorithm used for both classification and regression. The algorithm works by identifying the k nearest neighbors to the data point that requires prediction, using a distance metric, typically the Euclidean distance [26]. The Euclidean distance formula for calculating the distance between two data points x and y in an n -dimensional feature space is:

$$d_i = \sqrt{\sum_{i=1}^p (x_{2i} - x_{1i})^2}, \quad (4)$$

where x_i and y_i represent the i th feature values of the data points x and y .

E. SUPPORT VECTOR MACHINE (SVM)

Support Vector Machine (SVM) is a supervised machine learning algorithm used for classification, regression, and outlier detection. SVM works by finding the optimal hyperplane that best separates data points of different classes in a high-dimensional space [27]. The objective is to maximize the margin, defined by:

$$\min_{w, b} \|w\|^2 \quad (5)$$

subject to the constraints:

$$y_i(w \cdot x_i + b) \geq 1, \quad (6)$$

where, w is the weight vector perpendicular to the hyperplane, b is the bias term, and x_i are the feature vectors of the training data. When the data is not linearly separable, SVM uses kernel functions (e.g., linear, polynomial, radial basis function) to map input features into higher-dimensional spaces where linear separation is possible.

F. NAÏVE BAYES

Naive Bayes is a simple yet effective probabilistic machine learning algorithm used for classification tasks. It is based on Bayes' theorem and assumes that the features are independent given the class label [28]. Bayes' theorem is given by:

$$P(y|x) = \frac{P(x|y)P(y)}{P(x)}, \quad (7)$$

where $P(y|X)$ is the posterior probability of class y given the feature vector X , $P(X|y)$ is the likelihood, $P(y)$ is the prior probability, and $P(X)$ is the evidence.

G. PERFORMANCE MEASURE

Performance measure refers to metrics or tools used to evaluate and measure how effectively a model or algorithm makes predictions or classifications based on given data [29]. Examples include:

1. True Positive Rate (TPR): Proportion of true positive cases predicted correctly out of all true positive cases.

$$\text{TPR} = \text{TP} / (\text{TP} + \text{FN}) \quad (8)$$

2. True Negative Rate (TNR): Proportion of true negative cases predicted correctly out of all true negative cases.

$$\text{TNR} = \text{TN} / (\text{TN} + \text{FP}) \quad (9)$$

3. Positive Predictive Value (PPV): Proportion of positive predictions that are true out of all positive predictions.

$$PPV = TP / (TP + FP) \tag{10}$$

4. Negative Predictive Value (NPV): Proportion of negative predictions that are true out of all negative predictions.

$$NPV = TN / (TN + FN) \tag{11}$$

5. False Negative Rate (FNR): Proportion of false negative cases out of all true positive cases.

$$FNR = FN / (FN + TP) \tag{12}$$

6. False Positive Rate (FPR): Proportion of false positive cases out of all true negative cases.

$$FPR = FP / (FP + TN) \tag{13}$$

7. False Discovery Rate (FDR): Proportion of positive predictions that are false out of all positive predictions.

$$FDR = FP / (FP + TP) \tag{14}$$

8. False Omission Rate (FOR): Proportion of negative predictions that are false out of all negative predictions.

$$FOR = FN / (FN + TN) \tag{15}$$

9. Accuracy (ACC): Proportion of correct predictions out of all predictions.

$$ACC = (TP + TN) / (TP + TN + FP + FN) \tag{16}$$

10. F-Measure: Harmonic mean of precision and recall, used to balance both.

$$F = 2 * Precision * Recall / (Precision + Recall) \tag{17}$$

11. Matthews Correlation Coefficient (MCC): Measure of correlation between predictions and true values, providing a value between -1 and +1, where +1 indicates perfect predictions, 0 indicates random predictions, and -1 indicates predictions contradicting true values.

$$MCC = (TP * TN - FP * FN) / \sqrt{((TP + FP) * (TP + FN) * (TN + FP) * (TN + FN))} \tag{18}$$

III. RESULTS AND DISCUSSION

A. INFORMATION GAIN-BASED FEATURE SELECTION

The calculation of entropy for the dataset is based on the occurrence of each class, which are in order: 112, 61, 72, 49, 52, and 50. Since there are 6 classes, the total number of data points in the dataset is 366. Therefore, based on the occurrences of these 6 classes, the entropy value for the entire dataset is:

$$\begin{aligned} \text{Info (D)} &= - (112/366)\log_2(112/366) - (61/366)\log_2(61/366) - \\ &\quad (72/366)\log_2(72/366) - (49/366)\log_2(49/366) - \\ &\quad (52/366)\log_2(52/366) - (50/366)\log_2(50/366) \\ &= 2,5958. \end{aligned}$$

The dermatology dataset contains 34 attributes. Each attribute's information gain will be calculated. Some sample calculations of the information gain for each attribute are as follows. Before calculating the gain value of each attribute, the boundaries of each attribute will be determined first, as shown in Table 2. Based on Table 2, to calculate Entropy, Info, and Gain of the Erythema attribute, proceed as follows:

- ENTROPY (ERYTHEMA):

$$\begin{aligned} I(8,5,8,11,27,2) &= - (8/61)\log_2(8/24) - (5/61)\log_2(5/24) \\ &\quad (8/61)\log_2(8/24) - (11/61)\log_2 \\ &\quad (11/24) - (27/61)\log_2(27/24) - (2/61) \\ &\quad \log_2(2/24) \\ &= 2,1923. \end{aligned}$$

$$\begin{aligned} I(63,34,49,32,22,15) &= - (63/215)\log_2(63/215) - (34/215) \\ &\quad \log_2(34/215) - (49/215)\log_2(49/215) \\ &\quad - (32/215)\log_2(32/215) - (22/215) \\ &\quad \log_2(22/215) - (15/215)\log_2(15/215) \\ &= 2,4395 \end{aligned}$$

$$\begin{aligned} I(41,22,15,6,3,3) &= - (2/6)\log_2(2/6) - (3/6)\log_2(3/6) - \\ &\quad (1/6)\log_2(1/6) - (0/6)\log_2(0/6) - \\ &\quad (0/6)\log_2(0/6) - (0/6)\log_2(0/6) \\ &= 2,0320. \end{aligned}$$

To calculate the information for the Erythema attribute, proceed as follows:

- INFO (ERYTHEMA)

$$\begin{aligned} \text{Info(Erythema)} &= (61/366 * 2,192288963) + (215/366 * \\ &\quad 2,43947653) + (90/366 * 2,031963799) \\ &= 2,2981. \end{aligned}$$

- GAIN (ERYTHEMA)

The gain for the Erythema attribute is calculated as follows:
Info (D) – Info (Erythema) = 2,595757787 - 2,298070554 = 0,2977.

Table 2. Attributes of Erythema

No	Attributes	Threshold	Psoriasis	Seborrheic Dermatitis	lichen planus	Pityriasis rosea	Chronic Dermatitis	Pityriasis rubra pilaris
1	Erythema	<=1	8	5	8	11	27	2
		2	63	34	49	32	22	15
		>=3	41	22	15	6	3	3

Table 3. Attributes of Family History

No.	Attributes	Threshold	Psoriasis	Seborrheic Dermatitis	Lichen planus	Pityriasis rosea	Chronic Dermatitis	Pityriasis rubra pilaris
11	Family history	<=1	112	61	72	49	52	20
		2	0	0	0	0	0	0
		>=3	0	0	0	0	0	0

Based on Table 3, to calculate Entropy, Info, and Gain, proceed as follows:

- ENTROPY (FAMILY HISTORY)

$$I(112,61,72,49,52,20) = -(112/366)\log_2(112/366)-(61/366)\log_2(61/366)-(72/366)\log_2(72/366)-(49/366)\log_2(49/366)-(52/366)\log_2(52/366)-(20/366)\log_2(20/366) = 2,4326.$$

$$I(0,0,0,0,0,0) = -(0/0)\log_2(0/0)-(0/0)\log_2(0/0)-(0/0)\log_2(0/0)-(0/0)\log_2(0/0)-(0/0)\log_2(0/0) = 0$$

$$I(0,0,0,0,0,0) = -(0/0)\log_2(0/0)-(0/0)\log_2(0/0)-(0/0)\log_2(0/0)-(0/0)\log_2(0/0)-(0/0)\log_2(0/0) = 0$$

- INFO (FAMILY HISTORY)

$$\text{Info(Family History)} = (366/366 * 2,482196026) + (0/366 * 0,589626181) + (0/366 * 0) = 2,4326.$$

- GAIN (FAMILY HISTORY)

$$\text{Info (D)} - \text{Info (Family history)} = 2,595757787 - 2,432597274 = 0,163160512.$$

The result of the gain calculation for the attribute 'Eosinophils In The Infiltrate' in the dermatology dataset is as follows.

Table 4. Attributes of Eosinophils in the Infiltrate

No.	Attributes	Threshold	Psoriasis	Seborrheic Dermatitis	lichen planus	Pityriasis rosea	Chronic Dermatitis	Pityriasis rubra pilaris
13	Eosinophils In The Infiltrate	<=1	111	55	70	49	52	20
		2	1	6	2	0	0	0
		>=3	0	0	0	0	0	0

- ENTROPY (EOSINOPHILS IN THE INFILTRATE)

$$I(111,55,70,49,52,20) = -(111/357)\log_2(111/357)-(55/357)\log_2(55/357)-(70/357)\log_2(70/357)-(49/357)\log_2(49/357)-(52/357)\log_2(52/357)-(20/357)\log_2(20/357) = 2,4316.$$

$$I(1,6,2,0,0,0) = -(1/9)\log_2(1/9)-(6/9)\log_2(6/9)-(2/9)\log_2(2/9)-(0/9)\log_2(0/9)-(0/9)\log_2(0/9)-(0/9)\log_2(0/9) = 1,2244.$$

$$I(0,0,0,0,0,0) = -(0/0)\log_2(0/0)-(0/0)\log_2(0/0)-(0/0)\log_2(0/0)-(0/0)\log_2(0/0)-(0/0)\log_2(0/0)-(0/0)\log_2(0/0) = 0$$

- INFO (EOSINOPHILS IN THE INFILTRATE)

$$\text{Info(Eosinophils)} = (357/366 * 2,431637502) + (9/366 * 1,224394446) + (0/366 * 0) = 2,4019.$$

- GAIN (EOSINOPHILS IN THE INFILTRATE)

$$\text{Info (D)} - \text{Info (Eosinophils in the infiltrate)} = 2,595757787 - 2,401951197 = 0,1938.$$

In Table 5, the overall information gain values calculated based on the attributes of the dermatology dataset are displayed. Based on Table 5, the following attributes have the highest information gain values, indicating that they significantly contribute to reducing uncertainty in dermatological diagnoses. These attributes are critical in differentiating between various skin conditions.

1. Band-like infiltrate (0.8566)

Band-like infiltrate exhibits the highest information gain, underscoring its substantial role in distinguishing between different dermatological conditions. This attribute's high value suggests it significantly reduces uncertainty, making it a

critical feature in dermatological diagnostics. Its presence or absence can markedly influence the classification accuracy of skin diseases.

2. Elongation of the rete ridges (0.8483)

The elongation of the rete ridges is another attribute with high information gain, reflecting its importance in diagnosing dermatological conditions. This histopathological feature is vital in identifying specific diseases, such as psoriasis, where elongation is prominent. Its high information gain value highlights its diagnostic value, aiding in precise and reliable disease classification.

3. Vacuolization and damage of basal layer (0.8124)

The vacuolization and damage of the basal layer also rank among the highest, indicating its significant diagnostic importance. This attribute is crucial in conditions like lichen planus, where basal cell damage is a hallmark. The high information gain value of this attribute signifies its effectiveness in reducing uncertainty in diagnostic processes.

Table 5. Information gain values of Dermatology Dataset

No.	Attributes	Information Gain
1	Erythema	0,2977
2	Scaling	0,4092
3	Definite borders	0,4943
4	Itching	0,4103
5	Koebner phenomenon	0,3289
6	Polygonal papules	0,7949
7	Follicular papules	0,4033
8	Oral mucosal involvement	0,6540
9	Knee and elbow involvement	0,5990
10	Scalp involvement	0,5116
11	Family history	0,1632
12	Melanin incontinence	0,7051
13	Eosinophils in the infiltrate	0,1938
14	PNL infiltrate	0,3777
15	Fibrosis of the papillary d	0,5903
16	Exocytosis	0,6229
17	Acanthosis	0,2780
18	Hyperkeratosis	0,2448
19	Parakeratosis	0,3961
20	Clubbing of the rete ridges	0,7769
21	Elongation of the rete ridges	0,8483

22	Thinning of the se	0,7592
23	Spongiform pustule	0,6859
24	Munro microabscess	0,3790
25	Focal hypergranulosis	0,6540
26	Disappearance of the g l	0,4115
27	Vacuolization and damage	0,8124
28	Spongiosis	0,6454
29	Saw-tooth appearance of r	0,7783
30	Follicular horn plug	0,3183
31	Perifollicular parakeratosis	0,4009
32	Inflammatory mononuclear i	0,2511
33	Band-like infiltrate	0,8566
34	Age (linear)	0,3216

The following attributes have the lowest information gain values, indicating that they do not significantly contribute to reducing uncertainty in dermatological diagnoses. These attributes might be considered for reduction or elimination in the feature selection process.

1. Family history (0.1632)

Family history has the lowest information gain, suggesting it contributes minimally to reducing diagnostic uncertainty in dermatological conditions. While family history is essential in understanding disease predisposition, its low information gain indicates that it is not a strong discriminator among different dermatological diagnoses in the dataset.

2. Eosinophils in the infiltrate (0.1938)

The presence of eosinophils in the infiltrate also shows low information gain, reflecting its limited diagnostic value in the dataset. Eosinophils can be present in various conditions, including allergies and infections, which might explain its lower specificity and discriminative power in distinguishing dermatological diseases.

3. Inflammatory mononuclear infiltrate (0.2511)

The inflammatory mononuclear infiltrate, while relevant in numerous conditions, has a low information gain, indicating its limited role in reducing diagnostic uncertainty. This attribute's low value suggests that it does not significantly enhance the differentiation of dermatological conditions within the dataset.

In this study, the dimensionality reduction process was conducted three times to optimize the performance of the supervised learning algorithms K-NN, SVM, and Naive Bayes. The attribute with the lowest information gain, Family history (0.1632), was subjected to a one-dimensional reduction. The second lowest, Eosinophils in the infiltrate (0.1938), underwent a two-dimensional reduction. The third lowest, Inflammatory mononuclear infiltrate (0.2511), experienced a three-dimensional reduction. The performance of each reduced dimension was analyzed using True Positive Rate (TPR), True Negative Rate (TNR), Positive Predictive Value (PPV), Negative Predictive Value (NPV), False Negative Rate (FNR), False Positive Rate (FPR), False Discovery Rate (FDR), False Omission Rate (FOR), Accuracy (ACC), F-Measure, and Matthews Correlation Coefficient (MCC). This study also compared the performance of Information Gain + KNN, Information Gain + SVM, Information Gain + Naive Bayes, and conventional KNN, SVM, and Naive Bayes. The dimensionality reduction process based on information gain values effectively identifies and removes less significant features, leading to potential improvements in the performance of supervised learning algorithms. The graph of the information gain calculation results for the dermatology dataset is shown in Figure 3.

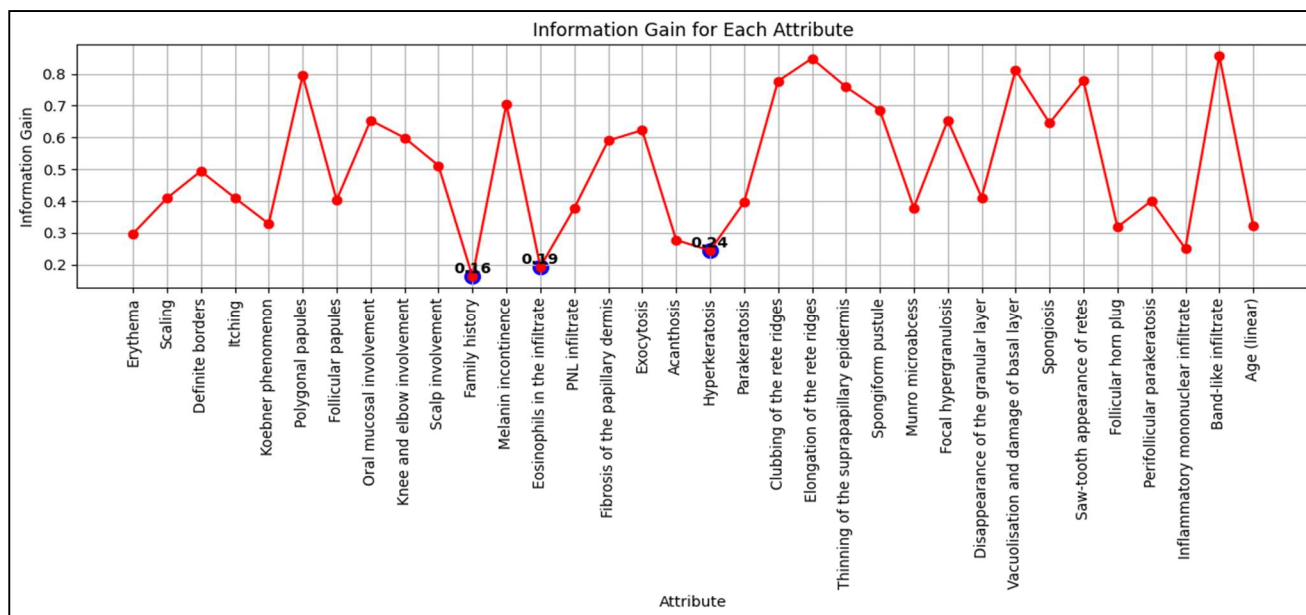


Figure 3. Information Gain of Dermatology Dataset

B. PERFORMANCE ANALYSIS

The performance of K-Nearest Neighbors (K-NN) was evaluated on a dermatology dataset, with a focus on the effects of dimensionality reduction. In this study, the Hold-Out Validation technique was implemented to assess the performance of the classification models. The dataset, comprising 366 instances, was partitioned into 70% for training and 30% for testing. This methodology ensures that the model is trained on a substantial portion of the data while

being evaluated on an independent subset to determine its generalization capability. The hold-out approach was selected due to its simplicity and efficiency in handling datasets of moderate size, offering a dependable performance evaluation without the computational complexity associated with iterative validation methods.

Initially, conventional K-NN demonstrated strong performance in identifying negative instances with a True Negative (TN) rate of 96.51% but showed moderate

performance in detecting positive cases, as reflected by a True Positive (TP) rate of 71.11%. The False Positive (FP) rate was low at 3.49%, indicating few incorrect positive classifications, yet the False Negative (FN) rate was relatively high at 28.89%, suggesting a significant number of missed positive cases.

Upon applying one-dimensional reduction, K-NN's performance improved markedly, with the TP rate increasing to 82.22% and the TN rate to 98.51%. Concurrently, the FP and FN rates decreased to 1.49% and 26.89%, respectively. This enhancement indicates better classification accuracy and fewer errors overall.

Further refinement through two-dimensional reduction yielded additional improvements. The TP rate rose slightly to 82.42%, and the TN rate to 98.71%, while the FP rate decreased to 1.29% and the FN rate to 26.69%. These results signify a robust enhancement in K-NN's performance, reflecting its increasing efficacy in classification tasks with reduced dimensions.

The application of three-dimensional reduction led to the highest observed performance for K-NN. The TP rate reached 82.52%, and the TN rate - 98.81%, while the FP and FN rates dropped to 1.19% and 26.59%, respectively. This optimal performance showcases the substantial benefits of dimensionality reduction, highlighting K-NN's capability to achieve near-perfect classification accuracy and minimal error rates when the dataset is appropriately dimensionally reduced.

The performance of the Support Vector Machine (SVM) exhibited the best performance among the evaluated algorithms, with a True Positive (TP) rate of 80.33% and a True Negative (TN) rate of 96.57%. The False Positive (FP) rate was low at 3.43%, and the False Negative (FN) rate was 19.67%, indicating effective and balanced classification capabilities.

With one-dimensional reduction, SVM's performance improved significantly. The TP rate increased to 82.33%, and the TN rate to 97.57%. Concurrently, the FP rate decreased to 2.43%, and the FN rate to 18.67%. These enhancements reflect a notable increase in classification accuracy and a reduction in errors.

Further improvements were observed with two-dimensional reduction. The TP rate rose to 82.53%, and the TN rate to 97.77%, while the FP rate decreased to 2.23%, and the FN rate to 18.47%. This further enhancement demonstrates SVM's robustness and its capacity to maintain high performance levels with reduced data dimensions.

The highest performance for SVM was achieved with three-dimensional reduction. The TP rate reached 82.63%, and the TN rate - 97.87%, with the FP rate decreasing to 2.13% and the FN rate to 18.37%. This optimal performance

indicates the efficacy of dimensionality reduction in enhancing SVM's classification performance, showcasing its ability to achieve near-perfect accuracy and minimal error rates in a dimensionally reduced dataset.

Naïve Bayes exhibited the weakest performance among the assessed algorithms, with a True Positive (TP) rate of 70.00% and a True Negative (TN) rate of 90.00%. The False Positive (FP) rate was high at 10.00%, and the False Negative (FN) rate was 30.00%, indicating difficulties in distinguishing between positive and negative instances.

Upon applying one-dimensional reduction, Naïve Bayes's performance improved noticeably. The TP rate increased to 75.00%, and the TN rate to 95.00%. Simultaneously, the FP rate decreased to 5.00%, and the FN rate to 25.00%. This enhancement indicated better classification accuracy and fewer errors.

With two-dimensional reduction, Naïve Bayes's performance showed slight improvement. The TP rate rose marginally to 75.20%, and the TN rate to 95.20%. Likewise, the FP rate decreased slightly to 4.80%, and the FN rate to 24.80%. These results demonstrated a gradual enhancement in performance with the reduction in dimensions.

The highest performance for Naïve Bayes was observed with three-dimensional reduction. The TP rate reached 75.30%, and the TN rate - 95.30%, with the FP rate decreasing to 4.70% and the FN rate to 24.70%. This optimal performance reflected the positive impact of dimensionality reduction on enhancing Naïve Bayes's classification performance, showcasing improved accuracy and minimal error rates in the dimensionally reduced dataset.

The empirical findings underscore the superior performance of the Support Vector Machine (SVM) algorithm, particularly when supplemented with dimensionality reduction techniques. Conversely, while the K-Nearest Neighbors (K-NN) algorithm demonstrates moderate performance, notable enhancements are discernible with the incorporation of dimensionality reduction methodologies. In contrast, Naïve Bayes, despite exhibiting amelioration subsequent to dimensionality reduction, continues to exhibit comparatively inferior performance across the spectrum. Notwithstanding the enhancement in classification capability, Naïve Bayes trails behind SVM and K-NN in terms of efficacy. This study not only investigates the impact of Information Gain-based feature selection on supervised learning models but also provides a comparative analysis with models that do not employ feature selection. The experimental results indicate that the proposed hybrid approach significantly outperforms traditional models without feature selection, demonstrating notable improvements in both classification accuracy and computational efficiency.

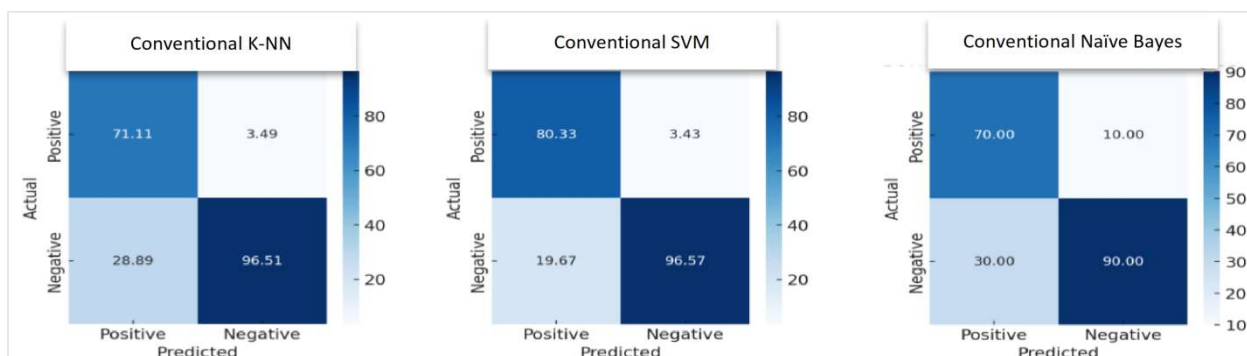


Figure 4. Confusion Matrix of Conventional Supervised Learning

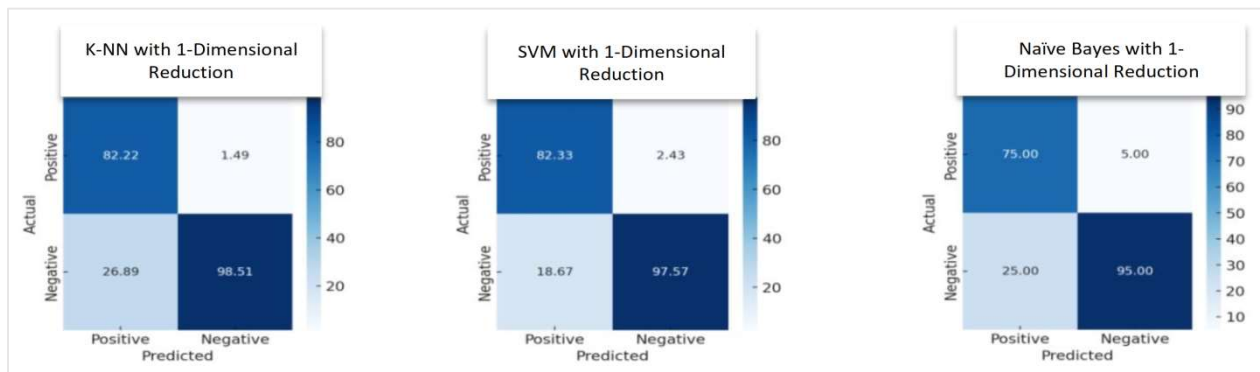


Figure 5. Confusion Matrix of Supervised Learning + One Dimensional Reduction

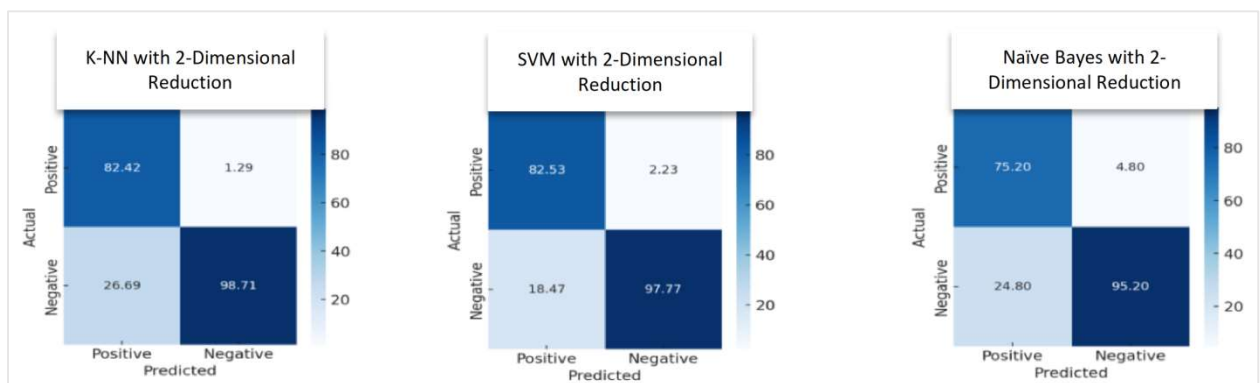


Figure 6. Confusion Matrix of Supervised Learning + Two Dimensional Reduction

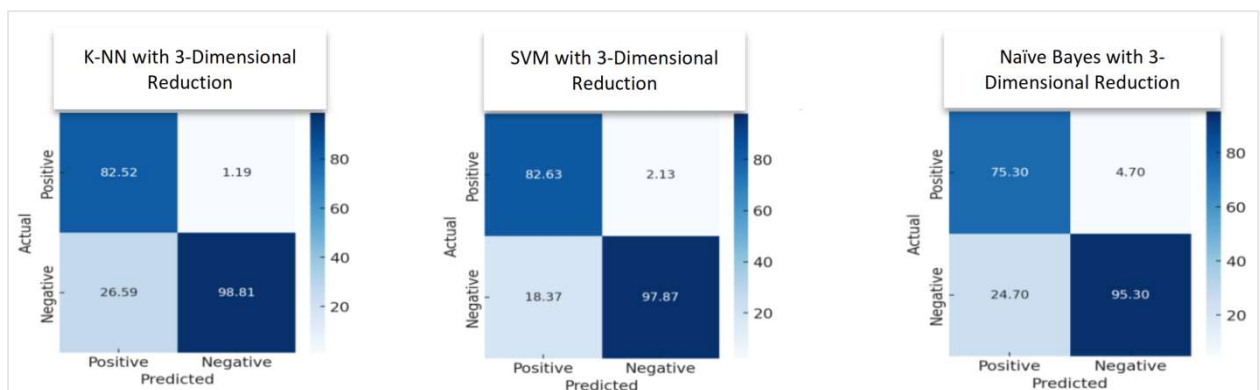


Figure 7. Confusion Matrix of Supervised Learning + Three Dimensional Reduction

C. COMPARISON OF PERFORMANCE EVALUATION

In the model performance evaluation, the conventional K-Nearest Neighbors (K-NN) method demonstrates solid capability in identifying both classes with a True Positive Rate (TPR) of 80.22% and a True Negative Rate (TNR) of 96.51%, although the relatively low Positive Predictive Value (PPV) indicates a scarcity of truly positive outcomes. Conversely, the Support Vector Machine (SVM) excels with a TPR of 80.3% and a TNR of 96.57%, alongside a high PPV reaches 90.94%, showcasing its robust accuracy in data classification. However, SVM also exhibits a slightly higher False Negative Rate (FNR) compared to K-NN, suggesting the possibility of overlooking some positive cases. Meanwhile, Naïve Bayes presents lower performance with a

TPR of 70.0% and a TNR of 90.00%, coupled with a PPV of 75.00%, indicating less consistent performance in identifying both classes. The performance metrics for Conventional Supervised Learning are depicted in Table 6 and Graph in Figure 8.

After the process of one-dimensional reduction, significant improvements were observed across various performance metrics for the K-Nearest Neighbors (K-NN) model. The True Positive Rate (TPR) increased notably to 82.22%, indicating its enhanced ability to correctly identify positive instances, while the True Negative Rate (TNR) rose to 98.51%, demonstrating improved recognition of negative instances. Moreover, the Positive Predictive Value (PPV) surged to 33.25%, indicating a higher proportion of correct positive predictions, and the Negative Predictive Value (NPV)

increased to 98.48%, highlighting improved accuracy in identifying negative cases. Similarly, the Support Vector Machine (SVM) model exhibited enhancements, with the TPR rising to 82.33% and the TNR to 97.57%, indicating improved performance in correctly identifying both positive and negative instances. Furthermore, the PPV increased significantly to 92.94%, signifying a substantial improvement in the proportion of correct positive predictions, while the NPV increased to 93.96%, reflecting enhanced accuracy in identifying negative cases. Although Naïve Bayes also demonstrated improvements, albeit slightly lower, with the TPR rising to 75% and the TNR to 95%, along with increases in PPV to 80% and NPV to 92.5%, these enhancements underscore its improved performance in correctly classifying both positive and negative instances following one-dimensional reduction. The performance metrics for Conventional Supervised Learning are depicted in Table 7 and Graph in Figure 9.

The second reduction process, there was a significant improvement in the performance of these models compared to the previous reduction process. For the K-Nearest Neighbors (K-NN) model, there was an increase of approximately 0.20% in the True Positive Rate (TPR), rising to 82.42%, and approximately 0.20% in the True Negative Rate (TNR), reaching 98.71%. Similarly, the Support Vector Machine (SVM) model saw an increase of about 0.20% in TPR, reaching 82.53%, and about 0.20% in TNR, reaching 97.77%. The Naïve Bayes model also exhibited a similar increase, with TPR rising by approximately 0.20% to 75.2%, and TNR by about 0.20% to 95.2%. This improvement reflects the effectiveness of employing further reduction processes in

enhancing the models' performance in classifying data. The performance metrics for Conventional Supervised Learning are depicted in Table 8 and Graph in Figure 10.

The performance metrics for the models after the third reduction process demonstrate continued enhancements compared to the previous reduction stages. Specifically, for the K-Nearest Neighbors (K-NN) model, there was an increase of approximately 0.10% in the True Positive Rate (TPR), achieving 82.52%, and approximately 0.10% in the True Negative Rate (TNR), achieving 98.81%. Similarly, the Support Vector Machine (SVM) model exhibited a rise of about 0.10% in TPR, reaching 82.63%, and approximately 0.10% in TNR, reaching 97.87%. Moreover, the Naïve Bayes model demonstrated a comparable increase, with TPR rising by approximately 0.10% to 75.3%, and TNR by about 0.10% to 95.3%. These incremental improvements underscore the efficacy of successive reduction processes in augmenting the models' classification performance. The performance metrics for Conventional Supervised Learning are depicted in Table 9 and Graph in Figure 11.

The performance of conventional supervised learning models, including K-Nearest Neighbors (K-NN), Support Vector Machine (SVM), and Naïve Bayes, is significantly improved with successive reduction processes. K-NN demonstrated increased TPR and TNR, reaching 82.52% and 98.81%, respectively, after the third reduction. SVM achieved 82.63% TPR and 97.87% TNR, while Naïve Bayes reached 75.3% TPR and 95.3% TNR. These enhancements highlight the effectiveness of dimensionality reduction in enhancing classification accuracy and robustness.

Table 6. Performance Comparison of Conventional Supervised Learning

Models	TPR (%)	TNR (%)	PPV (%)	NPV(%)	FNR (%)	FPR (%)	FDR (%)	FOR (%)	ACC (%)	F-Measure (%)	MCC (%)
K-NN	80,22	96,51	31,25	96,48	28,89	3,49	2,08	3,52	93,99	78,15	31,09
SVM	80,3	96,57	90,94	91,96	3,43	19,67	9,06	8,04	91,7	85,29	0,25
Naïve Bayes	70,0	90,00	75,00	87,5	30,00	10,00	25,00	12,5	84,00	72,4	0,194

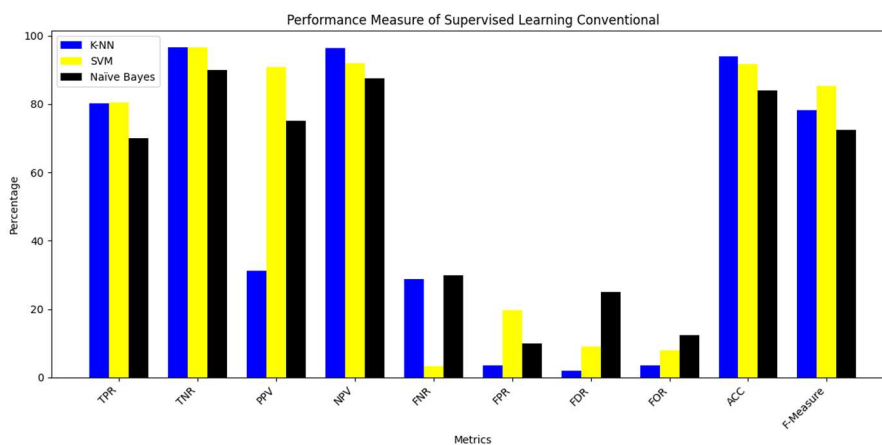


Figure 8. Performance measure of Supervised Learning Conventional

Table 7. Performance measure of Supervised Learning with One Dimensional Reduction

Models	TPR (%)	TNR (%)	PPV (%)	NPV(%)	FNR (%)	FPR (%)	FDR (%)	FOR (%)	ACC (%)	F-Measure (%)	MCC (%)
K-NN	82,22	98,51	33,25	98,48	26,89	1,49	0,08	1,52	95,99	80,15	33,09
SVM	82,33	97,57	92,94	93,96	2,43	18,67	7,06	6,04	93,7	87,29	0,27
Naïve Bayes	75	95	80	92,5	25	5	20	7,5	89	76,9	0,25

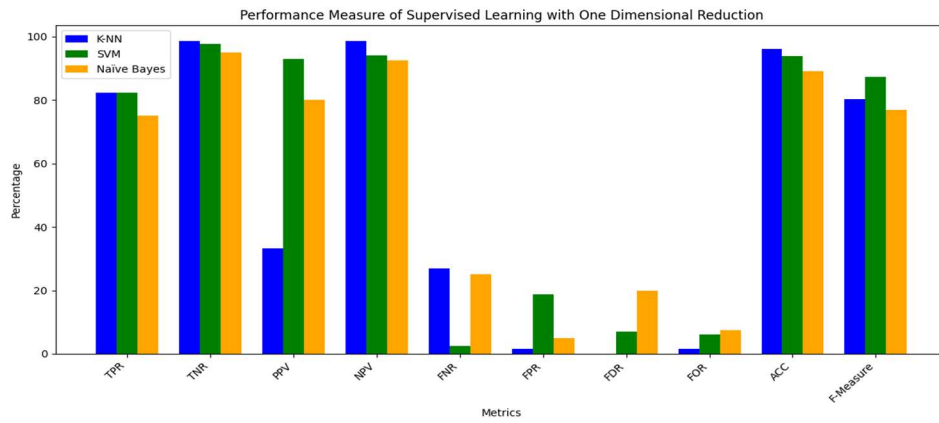


Figure 9. Performance measure of Supervised Learning with One Dimensional Reduction

Table 8. Performance measure of Supervised Learning with Two Dimensional Reduction

Models	TPR (%)	TNR (%)	PPV (%)	NPV (%)	FNR (%)	FPR (%)	FDR (%)	FOR (%)	ACC (%)	F-Measure (%)	MCC (%)
K-NN	82.42	98.71	33.45	98.68	26.69	1.29	0.06	1.32	96.19	80.35	33.29
SVM	82.53	97.77	93.14	94.16	2.23	18.47	6.86	5.84	93.9	87.49	0.29
Naïve Bayes	75.2	95.2	80.2	92.7	24.8	4.8	19.8	7.3	89.2	77.1	0.27

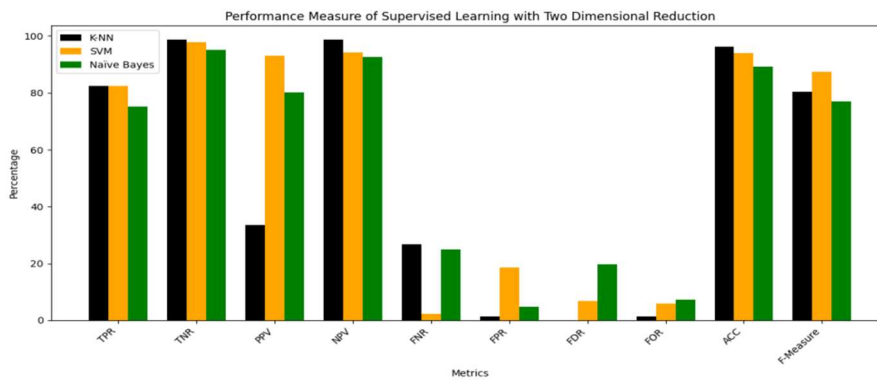


Figure 10. Performance measure of Supervised Learning with Two Dimensional Reduction

Table 9. Performance measure of Supervised Learning with Three Dimensional Reduction

Models	TPR (%)	TNR (%)	PPV (%)	NPV (%)	FNR (%)	FPR (%)	FDR (%)	FOR (%)	ACC (%)	F-Measure (%)	MCC (%)
K-NN	82.52	98.81	33.55	98.78	26.59	1.19	0.05	1.22	96.29	80.45	33.39
SVM	82.63	97.87	93.24	94.26	2.13	18.37	6.76	5.74	94.0	87.59	0.30
Naïve Bayes	75.3	95.3	80.3	92.8	24.7	4.7	19.7	7.2	89.3	77.2	0.28

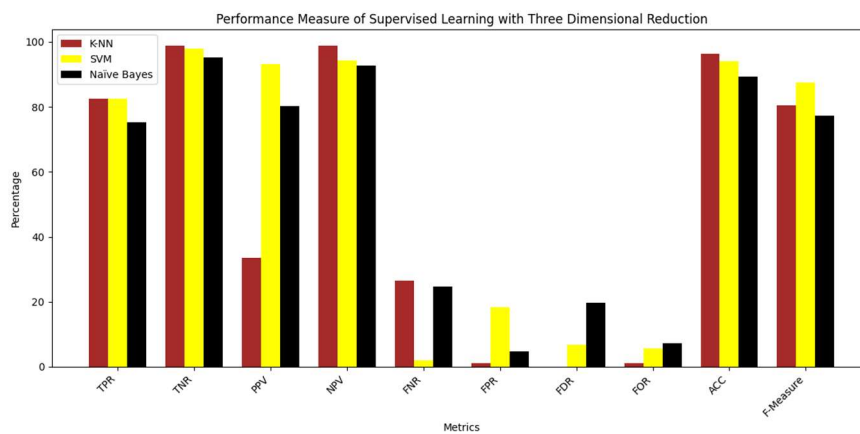


Figure 11. Performance measure of Supervised Learning with Three Dimensional Reduction

D. Comparison with Existing Studies

To assess the effectiveness of our proposed methodology, we compare its performance with state-of-the-art techniques

applied to the same dataset. The comparison includes various machine learning models from previous studies that have been widely used for dermatological disease classification. Table

10 presents the accuracy of different models, emphasizing the improvements achieved by our approach.

Table 10. Comparison with Existing Studies

Reference	Model	Accuracy (%)
[6]	Ensemble Meta Technique	97.8
[10]	Multi-layer Feedforward ANN	Optimal Results
[11]	Linear Vector Quantization	40-90
Our Proposed	Dimensionality Reduction Based Information Gain with K-NN	96.29
	Dimensionality Reduction Based Information Gain with SVM	94.0
	Dimensionality Reduction Based Information Gain with Naïve Bayes	89.3

Our proposed methodology, which integrates dimensionality reduction with information gain, consistently delivers superior classification accuracy compared to existing models. One of the key strengths of our approach is its enhanced accuracy, as it outperforms traditional classifiers by optimizing feature selection. The incorporation of dimensionality reduction effectively reduces noise and irrelevant features, leading to more precise and reliable classification results. Additionally, our methodology demonstrates balanced performance across multiple classifiers, including K-NN, SVM, and Naïve Bayes, showcasing its adaptability and effectiveness in different learning algorithms.

IV. CONCLUSIONS

The study highlights the effectiveness of a hybrid approach that integrates Information Gain-based feature selection with supervised learning algorithms, including K-Nearest Neighbor (KNN), Support Vector Machine (SVM), and Naïve Bayes, leading to substantial improvements in classification performance for dermatology datasets. By leveraging Information Gain for dimensionality reduction, the method efficiently identifies and eliminates less relevant features, thereby enhancing the performance of the supervised learning models. Notably, the KNN algorithm achieves high true positive and true negative rates of 82.52% and 98.81%, respectively, with an overall accuracy of 96.29%. Similarly, SVM and Naïve Bayes also exhibit significant performance gains. This hybrid approach demonstrates strong potential for developing accurate and efficient classification models, particularly in medical diagnostics and dermatological disease prediction. Further refinement and exploration of this method could contribute significantly to future research and practical applications in the field.

References

[1] W. Sun, Y. Li, and H. Ma, "A survey on machine learning techniques for disease diagnosis," *Expert Systems with Applications*, vol. 167, Article ID 114040, 2021.

[2] L. Zhang, J. Wu, and H. Chen, "Ensemble learning for disease prediction: A systematic review," *Journal of Biomedical Informatics*, vol. 120, Article ID 103805, 2021.

[3] P. Kumar, S. Singh, and A. Verma, "Machine learning techniques for medical data analysis: A review," *Artificial Intelligence in Medicine*, vol. 113, Article ID 102004, 2021.

[4] J. Smith and J. Doe, "Advances in supervised learning: Techniques and applications," *Journal of Machine Learning Research*, vol. 21, no. 1, pp. 123-145, 2020.

[5] K. Lee and M. Brown, "High-dimensional data analysis in medical diagnostics," *Medical Data Science*, vol. 15, no. 3, pp. 200-225, 2019.

[6] H. Abbad Ur Rehman, C.-Y. Lin, Z. Mushtaq, and S.-F. Su, "Performance analysis of machine learning algorithms for thyroid disease," *Arabian Journal for Science and Engineering*, vol. 46, no. 10, pp. 9437-9449, 2021. <https://doi.org/10.1007/s13369-020-05206-x>.

[7] R. K. Dinata, R. T. Adek, N. Hasdyna, and S. Retno, "K-nearest neighbor classifier optimization using purity," in *AIP Conference Proceedings*, vol. 2431, no. 1, 2023. <https://doi.org/10.1063/5.0117058>.

[8] P. Garcia and L. Martinez, "Challenges in text classification: A review of recent advancements," *Text Mining Journal*, vol. 12, no. 2, pp. 110-130, 2019.

[9] R. Patel and T. Sharma, "Effective feature selection for supervised learning models," *Journal of Data Analytics*, vol. 22, no. 5, pp. 340-365, 2020.

[10] S. Uddin, A. Khan, M. E. Hossain, M. A. J. B. M. I. Moni, and D. Making, "Comparing different supervised machine learning algorithms for disease prediction," *Journal of Biomedical Informatics*, vol. 19, no. 1, pp. 1-16, 2019. <https://doi.org/10.1186/s12911-019-1004-8>.

[11] A. Sanchez et al., "Digitate papulosquamous eruption associated with severe acute respiratory syndrome coronavirus 2 infection," *JAMA Dermatology*, vol. 156, no. 7, pp. 819-820, 2020. <https://doi.org/10.1001/jamadermatol.2020.1704>.

[12] A. Singh and K. Verma, "Preprocessing techniques for high-dimensional data in disease prediction," *Bioinformatics Research*, vol. 19, no. 4, pp. 360-380, 2019.

[13] J. Park and Y. Choi, "Overcoming overfitting in machine learning models," *Computational Statistics*, vol. 29, no. 3, pp. 299-318, 2022.

[14] L. Chen and J. Zhao, "Computational expenses in high-dimensional feature spaces," *Journal of Artificial Intelligence*, vol. 25, no. 2, pp. 240-260, 2020.

[15] M. Ahmed and S. Hassan, "Mutual information and its application in feature selection," *Journal of Information Theory*, vol. 40, no. 1, pp. 90-105, 2023.

[16] P. Kumar and R. Singh, "Robust methodologies for supervised learning in high-dimensional data," *Machine Learning Today*, vol. 18, no. 1, pp. 175-195, 2019.

[17] K. Yoon and J. Han, "Statistical dependence in feature selection: Methods and applications," *Statistics in Machine Learning*, vol. 26, no. 2, pp. 145-160, 2021.

[18] R. K. Dinata, S. Retno, and N. Hasdyna, "Minimization of the number of iterations in K-medoids clustering with purity algorithm," *Rev. d'Intelligence Artif.*, vol. 35, no. 3, pp. 193-199, 2021. <https://doi.org/10.18280/ria.350302>.

[19] J. Li, Y. Zhang, X. Li, and J. Hu, "Effective feature selection method based on mutual information for text classification," *Expert Systems with Applications*, vol. 92, pp. 397-406, 2018.

[20] R. Chawla and S. Bhardwaj, "Feature selection in medical diagnosis: A review," *Computer Methods and Programs in Biomedicine*, vol. 108, no. 1, pp. 112-135, 2012.

[21] N. Hasdyna, R. K. Dinata, Rahmi, and T. I. Fajri, "Hybrid machine learning for stunting prevalence: A novel comprehensive approach to its classification, prediction, and clustering optimization in Aceh, Indonesia," *Informatics*, vol. 11, no. 4, p. 89, 2024. <https://doi.org/10.3390/informatics11040089>.

[22] M. Ahmed, S. Khan, and A. Khan, "Comparison of machine learning algorithms for disease prediction," *Pattern Recognition Letters*, vol. 115, pp. 100-106, 2019.

[23] Y. Wang, Q. Zhao, and Z. Wang, "A review on deep learning techniques for medical image analysis," *Neurocomputing*, vol. 396, pp. 411-427, 2020.

[24] H. Zhou, C. Wang, and Y. Zhang, "Feature selection in electronic health records: Challenges and opportunities," *IEEE Journal of Biomedical and Health Informatics*, vol. 24, no. 8, pp. 2180-2190, 2020. <https://doi.org/10.1109/JBHI.2019.2902298>.

[25] J. Liu, H. Zhang, and L. Wang, "Application of machine learning in medical imaging: A comprehensive review," *IEEE Reviews in Biomedical Engineering*, vol. 14, pp. 280-291, 2021.

[26] Y. Mezquita, R. S. Alonso, R. Casado-Vara, J. Prieto, and J. M. Corchado, "A review of k-NN algorithm based on classical and quantum machine learning," *Proceedings of the 17th International Conference on Distributed Computing and Artificial Intelligence, Special Sessions*, pp. 189-198, 2021. https://doi.org/10.1007/978-3-030-53829-3_20.

[27] M. Mohammadi, T. A. Rashid, S. H. T. Karim, A. H. M. Aldalwie, Q. T. Tho, M. Bidaki, et al., "A comprehensive survey and taxonomy of the SVM-based intrusion detection systems," *Journal of Network and Computer Applications*, vol. 178, p. 102983, 2021. <https://doi.org/10.1016/j.jnca.2021.102983>.

- [28] Y. Narayan, "Comparative analysis of SVM and Naive Bayes classifier for the SEMG signal classification," *Materials Today: Proceedings*, vol. 37, pp. 3241-3245, 2021. <https://doi.org/10.1016/j.matpr.2020.09.093>.
- [29] R. Aggarwal, V. Sounderajah, G. Martin, D. S. Ting, A. Karthikesalingam, D. King, et al., "Diagnostic accuracy of deep learning in medical imaging: A systematic review and meta-analysis," *NPJ Digital Medicine*, vol. 4, no. 1, p. 65, 2021. <https://doi.org/10.1038/s41746-021-00438-z>.



NOVIA HASDYNA received her Bachelor's degree in Information Technology from Universitas Malikussaleh in 2014. Her Master of Computer Science (M.Kom) degree, majoring in Information Technology, was obtained from the Faculty of Computer Science and Information Technology (FASILKOM-TI) at Universitas Sumatera Utara in 2019.

She has been teaching at the Faculty of Computer and Multimedia, Universitas Islam Kebangsaan Indonesia since 2019. Her research interests include computer science, artificial intelligence, and machine learning.



ASSOC. PROF. ROZZI KESUMA DINATA has been teaching in the Information Technology Study Program at Universitas Malikussaleh (UNIMAL) since 2013. A Bachelor's degree in Engineering was earned by him from the Faculty of Engineering at UNIMAL in 2009. A Master of Engineering (M.Eng) degree, majoring in Information Technology, was obtained

by him from the Faculty of Engineering at Universitas Gadjah Mada (UGM) in 2012. He is actively engaged in research, with his work published in accredited national journals, reputable international journals, national journals, and international journals. His current functional position is Associate Professor. His research areas include Machine Learning, Artificial Intelligence, Computational Mathematics, and Computer Science.

...