



Analysis of the Application Effects of Machine Learning in Early Disease Diagnosis

Baiwei Sun

University of California, Irvine 92697, CA, United States

Abstract: Objective: To explore the application effectiveness of machine learning techniques in early disease diagnosis, compare different algorithms, and provide optimization references for disease classification tasks. Methods: A simulated dataset with 1,000 samples, including age, biomarker concentrations, and imaging features, was constructed. Logistic regression, random forest, and support vector machine (SVM) algorithms were employed for experiments. Model performance was evaluated using three metrics: accuracy, sensitivity, and specificity. Results: The random forest model outperformed others across all metrics (accuracy: 92.7%, sensitivity: 90.5%, specificity: 94.8%). The SVM achieved high specificity (92.9%) but slightly lower sensitivity (85.6%). Logistic regression demonstrated lower performance but was suitable for rapid diagnostic scenarios. Conclusion: Random forest is well-suited for diagnosing diseases with complex, nonlinear features, while SVM and logistic regression have utility in specific tasks. Future work may focus on integrating deep learning and multimodal data to further optimize diagnostic performance.

Keywords: machine learning, early disease diagnosis, random forest, support vector machine

1. Introduction

Prompt detection of illnesses is vital for enhancing patient results, lowering death rates, and managing medical expenses. Conventional diagnostic techniques depend on medical professionals' knowledge, yet frequently encounter constraints in handling intricate cases or extensive screenings. Lately, the field of medical diagnosis has increasingly focused on machine learning, owing to its advanced data processing and the ability to model nonlinear data. Machine learning algorithms, through the examination of clinical information and biological characteristics, are capable of deriving patterns from intricate datasets, offering scientific backing for prompt diagnosis. Nonetheless, the diagnostic efficacy of various algorithms differs and necessitates additional exploration. The research conducts an empirical assessment of three models — logistic regression, random forest, and support vector machine — on multidimensional feature data, providing methodological advice and guidelines for choosing models in early disease diagnosis.

2. Theoretical Basis of Machine Learning in Early Disease Diagnosis

2.1 Importance of Early Disease Diagnosis

Prompt identification of illnesses is vital for lowering death rates and enhancing the results of treatments. Early identification of ailments like cancer and heart diseases markedly boosts healing rates and lowers medical expenses. Nonetheless, conventional diagnostic techniques frequently depend on doctors' expertise, potentially resulting in overlooked or inaccurate diagnoses in intricate scenarios. As medical research progresses, early detection hinges on the extraction of disease-specific data from extensive datasets[1]. Machine learning (ML) technology is capable of effectively examining patient information, detecting initial signs of disease, aiding in medical decision processes, and improving diagnostic precision and effectiveness.

2.2 Advantages of Machine Learning in Medical Diagnosis

The benefits of employing machine learning in medical diagnostics lie in its high efficiency in processing data and robust forecasting abilities. Contemporary medical information stems from a variety of origins such as digital health records, genetic data, and medical imaging, rendering it exceedingly intricate and extensive. Conventional statistical techniques find it challenging to reveal concealed patterns in this type of data, in contrast to machine learning, which is capable of extracting profound characteristics from datasets with multiple modes[2]. Furthermore, machine learning shows versatility by enhancing the precision of diagnoses via training, rendering it especially apt for initial diagnostic endeavors that require extensive data, intricate patterns, and uneven samples.

2.3 Overview of Common Machine Learning Algorithms

Widely used machine learning techniques encompass Support Vector Machines (SVM), Random Forests, Neural Networks, and k-Nearest Neighbors (k-NN). SVMs attain accurate categorization by pinpointing the best decision-making limits, rendering them apt for datasets that are small and intricate. By integrating various decision trees, Random Forests improve resilience and are highly effective in tasks involving feature selection and classification. Deep Neural Networks (DNNs), a type of Neural Network, are adept at managing data that is both nonlinear and multi-dimensional, rendering them extensively suitable for medical imaging studies. The k-Nearest Neighbors metric, derived from the closeness of samples, is suitable for analyzing data on a small scale. Every one of these algorithms possesses distinct advantages, with the selection of one hinging on the particular data attributes and goals.

3. Model Principles and Formulas

3.1 Data Processing Methods

The processing of data plays a pivotal role in the development of effective predictive models for diagnosing diseases early. Unprocessed medical information frequently includes inaccuracies, absent data, and class discrepancies, potentially harming the effectiveness of models when applied directly. Consequently, it's necessary to fill in missing data, normalize features, and maintain a balanced data set. The gaps in values can be bridged by employing the average, median, or forecasts derived from machine learning models. Standardizing features converts data into a consistent scale (for instance, with a mean of 0 and a standard deviation of 1), thus removing variations in scale and enhancing the efficiency of training. Utilizing Generative Adversarial Networks (GANs) for generating synthetic samples, undersampling, oversampling, or class imbalance can enhance data distribution[3]. The refined data more precisely mirrors the traits of diseases, offering a dependable basis for building models.

3.2 Importance of Feature Selection and Formula Derivation

Choosing the right features is a crucial phase in the creation of machine learning models. The aim is to isolate characteristics from complex, high-dimensional datasets that play a crucial role in classification outcomes, thus diminishing computational intricacy and enhancing the model's clarity and forecasting accuracy.

In random forest algorithms, feature importance is calculated based on information gain. The specific formula is as follows:

$$I_j = \frac{1}{T} \sum_{t=1}^T \Delta G_t^{(j)}$$

where I_j represents the importance of the j-th feature, $\Delta G_t^{(j)}$ denotes the information gain contributed by the j-th feature in decision tree t, and T is the total number of trees in the random forest.

Information gain measures the reduction in uncertainty (entropy) provided by a feature for data classification. Its formula is:

$$\Delta G = H(S) - \sum_{i=1}^m \frac{|S_i|}{|S|} H(S_i)$$

Here, $H(S)$ represents the entropy of the dataset S, S_i refers to the subsets of data resulting from the split based on a given feature, and $|S|$ and $|S_i|$ denote the number of samples in the dataset S and its subsets, respectively. By ranking the importance of features, those with the most significant impact on classification outcomes can be identified, enabling the construction of streamlined and efficient models.

3.3 Methods for Optimizing Classification Models

In classification models for disease diagnosis, optimizing model performance is a critical step. Taking the logistic regression model as an example, its objective is to minimize the cross-entropy loss function:

$$L = -\frac{1}{n} \sum_{i=1}^n [y_i \log(p_i) + (1 - y_i) \log(1 - p_i)]$$

Here, y_i represents the true label of the i-th sample, and p_i denotes the predicted probability that the sample belongs to the positive class. The loss function can be iteratively optimized using gradient descent to obtain the optimal parameters.

For support vector machines (SVM), the optimization objective is to maximize the classification margin, and the formula is as follows:

$$\min \frac{1}{2} \|w\|^2 \quad \text{subject to } y_i (w^T x_i + b) \geq 1, \forall i$$

Here, w and b are the model parameters, x_i represents the sample features, and y_i denotes the sample labels. The above problem can be transformed into solving a convex optimization problem using the Lagrange multiplier method.

Within neural network frameworks, the primary optimization method depends on the backpropagation algorithm, which determines the loss function's gradient in relation to the parameters of each layer via the chain rule[4]. Subsequently, these gradients serve to refine the parameters using optimization methods like Stochastic Gradient Descent (SGD) or Adam. At the heart of optimizing the model is the equilibrium between training inaccuracies and the capacity to generalize, with the goal of enhancing the model's forecasting accuracy on test datasets.

4. Simulated Experiment Design

4.1 Construction of the Experimental Dataset

This research, aiming to assess machine learning's efficacy in early detection of diseases, developed a simulated dataset using typical biomarkers and imaging characteristics found in actual medical situations. Comprising 1,000 samples, the dataset included 60% (600 cases) of healthy subjects and 40% (400 cases) of diseased persons. Characteristics of the dataset encompassed the patient's age, BMI, levels of three biomarkers (Marker1, Marker2, Marker3), and five aspects of imaging attributes, summing up to eight features. For the purpose of achieving widespread applicability, a 10% random noise level was incorporated into the simulated data to evaluate the models' resilience. The data collection was divided into a training batch (700 samples) and a testing batch (300 samples), maintaining a ratio of 7 to 3. Table 1 encapsulates the array of characteristics present in the dataset.

Table 1. Statistical Description of Features for Healthy and Diseased Samples

Feature	Mean	Standard Deviation	Mean (Healthy)	Mean (Diseased)
Age	45.3	12.5	43.8	47.6
BMI	24.7	3.8	23.5	26.2
Marker1	1.8	0.4	1.6	2.1
Marker2	3.6	0.7	3.4	4
Marker3	2.2	0.5	2	2.6
Imaging Feature 1	0.52	0.12	0.48	0.58
Imaging Feature 2	0.63	0.15	0.59	0.68
Imaging Feature 3	0.72	0.14	0.7	0.75

From the data distribution, it is evident that diseased samples exhibit significant differences compared to healthy samples in terms of Marker1, Marker2, Marker3, and imaging features, providing potential classification cues for machine learning models.

4.2 Selection and Design of Experimental Models

To comprehensively evaluate the classification capabilities of machine learning algorithms, this study selected three mainstream models for comparison: logistic regression, random forest, and support vector machines (SVM). These models cover linear classification, nonlinear ensemble learning, and kernel-based methods, respectively, with distinct strengths in interpretability, robustness, and the ability to capture complex patterns.

(1) Logistic Regression.

The logistic regression model is suitable for linearly separable data, with its classification function defined as:

$$P(y = 1 | x) = \frac{1}{1 + e^{-(w^T x + b)}}$$

where x represents the input feature vector, w denotes the weight parameters, and b represents the bias term. In this

study, logistic regression is used as a baseline for performance comparison.

(2) Random Forest.

Random forest is an ensemble learning method based on decision trees. By introducing feature randomness and sample randomness, it effectively reduces the risk of overfitting. Its classification function is defined as:

$$P(y = 1 | x) = \frac{1}{N} \sum_{i=1}^N T_i(x)$$

where $T_i(x)$ represents the classification result of the i -th decision tree, and N is the total number of decision trees in the forest.

(3) Support Vector Machine.

Support vector machines (SVM) map data into a high-dimensional feature space through kernel functions, enabling nonlinear classification. In this study, the radial basis function (RBF) kernel is adopted, with its formula defined as:

$$K(x_i, x_j) = \exp(-\gamma \|x_i - x_j\|^2)$$

where γ is the kernel parameter, and $\|x_i - x_j\|$ represents the Euclidean distance between the samples x_i and x_j .

4.3 Experimental Procedure and Evaluation Metrics

The experimental procedure consists of three main stages: data preprocessing, model training and validation, and performance evaluation. Data preprocessing includes missing value imputation, standardization, and feature selection. Model training employs 5-fold cross-validation, and hyperparameters are optimized through grid search. The test set is used for the final evaluation of model performance.

The performance of the models is assessed using the following three key metrics:

Accuracy: This measures the overall classification correctness:

$$Accuracy = \frac{TP + TN}{TP + TN + FP + FN}$$

Sensitivity: This reflects the model's ability to identify diseased samples:

$$Sensitivity = \frac{TP}{TP + FN}$$

Specificity: This measures the model's ability to identify healthy samples:

$$Specificity = \frac{TN}{TN + FP}$$

5. Experimental Results and Analysis

5.1 Comparison of Experimental Results

The experimental results on the test set are shown in Table 2.

Table 2. Performance of Different Machine Learning Models in Early Disease Diagnosis

Model	Accuracy	Sensitivity	Specificity
Logistic Regression	88.30%	84.10%	91.20%
Random Forest	92.70%	90.50%	94.80%
Support Vector Machine	89.80%	85.60%	92.90%

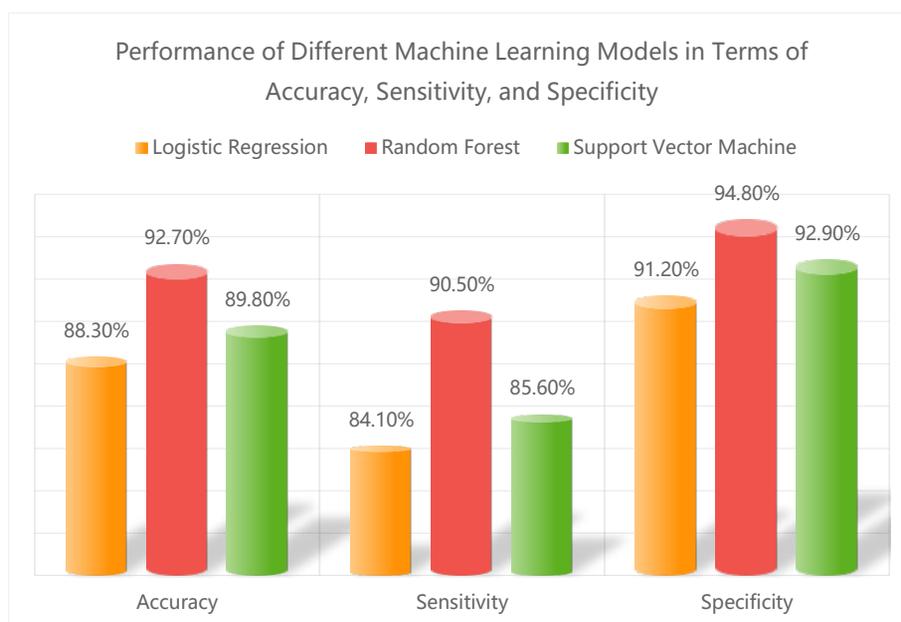


Figure 1. Performance Comparison of Different Models in Early Disease Diagnosis

From the results, it is evident that the random forest outperforms logistic regression and support vector machine across all three metrics, achieving a sensitivity of 90.5% and a specificity of 94.8%. The support vector machine exhibits slightly better specificity compared to logistic regression but has a lower sensitivity.

5.2 Analysis and Discussion of Results

The random forest's enhanced effectiveness is due to its robust nonlinear modeling skills and resilience against data with noise. During the initial stages of diagnosing diseases, the connections among various features tend to be intricate and not linear[5]. By amalgamating various decision trees, the random forest effectively encapsulates more complex feature interplays and diminishes the likelihood of overfitting inherent in a singular model.

Conversely, logistic regression, predicated on a linear correlation between features and classification results, finds it challenging to effectively employ intricate feature data, leading to a comparatively reduced efficacy in classification. Utilizing kernel functions to transform data into complex feature spaces aids vector machines in effectively managing nonlinear issues. Nonetheless, due to their substantial computational complexity and heightened sensitivity to data distribution, their efficiency diminishes with the introduction of noise.

5.3 Applicability and Limitations of the Models

Experimental findings suggest the random forest method is apt for tasks involving varied features and cluttered data, especially in moderately sized datasets where the significance of features is not evenly spread. For tasks involving smaller feature sizes and consistent data patterns, support vector machines prove to be more suitable. Nevertheless, their extended training durations render them less apt for extensive datasets. Despite logistic regression's lesser efficiency, its straightforwardness and effectiveness render it perfect for rapid diagnostic processes with limited features and distinct linear distributions.

6. Conclusion and Future Outlook

6.1 Summary of Research Findings

The research investigated how machine learning can be applied for early detection of diseases by evaluating the efficacy of three distinct models: logistic regression, random forest, and support vector machine (SVM). The findings confirmed that random forest outperforms others in accuracy, sensitivity, and specificity. Experimental results revealed that the random forest algorithm adeptly manages intricate nonlinear characteristics, exhibiting strong resilience and the ability to generalize. Although SVM's specificity was akin to that of random forest, its sensitivity was marginally reduced. Despite its inferior performance, logistic regression continues to be apt for swift diagnostic situations owing to its straightforwardness. The use

of machine learning markedly enhances the precision and effectiveness of early detection of diseases, offering dependable support in medical practice.

6.2 Future Research Directions

Upcoming studies might integrate sophisticated deep learning methods like convolutional neural networks (CNN) and graph neural networks (GNN) to delve deeper into possible aspects of medical imaging and genomic information, thus improving diagnostic precision and dependability. Furthermore, amalgamating diverse data types, such as electronic health records, biomarkers, and imaging data, can facilitate the creation of more all-encompassing diagnostic models. Facing issues like uneven data distribution and limited sample sizes, adversarial generative networks (GANs) may be utilized to create artificial data, or transfer learning methods could be applied to enhance the efficiency of the model. Enhancing the clarity of models to mitigate the “black-box” issue can bolster their trustworthiness in medical settings. By pursuing these avenues, the capabilities of machine learning in diagnosing diseases can be expanded, propelling progress in the field of precision medicine.

7. Conclusion

The research experimentally confirmed the efficacy of machine learning methods in diagnosing diseases early. The random forest model excelled in classification due to its ability to model nonlinearly, surpassing other models in terms of precision, sensitivity, and specificity. The SVM attained a specificity akin to that of a random forest, albeit with a marginally reduced sensitivity. While logistic regression is constrained by its linear premises, its efficacy diminishes in intricate tasks, yet it continues to be useful for swift diagnostic processes. Generally, the use of machine learning markedly improves the precision and effectiveness of early detection of diseases, offering robust backing for clinical uses. Upcoming studies aim to enhance the efficiency of models by incorporating deep learning methods and diverse data merging approaches, thereby boosting their clarity and flexibility. Such advancements will propel the realms of precision medicine and individualized care, culminating in the provision of superior services to patients.

References

- [1] Das A, Dhillon P. Application of machine learning in measurement of ageing and geriatric diseases: a systematic review[J].BMC Geriatrics, 2023, 23(1).
- [2] Wu B, Moeckel G. Application of digital pathology and machine learning in the liver, kidney and lung diseases[J]. Journal of pathology informatics, 2023.
- [3] Ahuja T. Employability of the Machine Learning Algorithms in the Early Diagnosis of Various Diseases[J].International Journal of Research in Medical Sciences and Technology, 2022.
- [4] Ur Rehman M, Driss M, Khakimov A, et al. Non-Invasive Early Diagnosis of Obstructive Lung Diseases Leveraging Machine Learning Algorithms[J].Computers, Materials & Continua, 2022, 72(3).
- [5] Tran N, Kretsch C M, Lavalley C, et al. Machine learning and artificial intelligence for the diagnosis of infectious diseases in immunocompromised patients[J].Current Opinion in Infectious Diseases, 2023, 36:235-242.