

Classification of Dengue Fever Disease Using a Machine Learning turnitin (1).docx

by 1 1

Submission date: 22-Jul-2024 11:56AM (UTC+0800)

Submission ID: 2420536261

File name: Classification_of_Dengue_Fever_Disease_Using_a_Machine_Learning_turnitin_1_.docx (47.07K)

Word count: 2368

Character count: 13979

Classification of Dengue Fever Disease Using a Machine Learning-Based Random Forest Algorithm

Arif Fitra Setyawan¹, Amelia Devi Putri Ariyanto², Fari Katul Fikriah³

¹Program Studi Sistem dan Teknologi Informasi, Fakultas Keperawatan Bisnis dan Teknologi, Universitas Widyadarmas Husada Indonesia, Sp. (ETS) Sp. (ETS)
²Program Studi Sistem dan Teknologi Informasi, Fakultas Keperawatan Bisnis dan Teknologi, Universitas Widyadarmas Husada Indonesia, Sp. (ETS) Sp. (ETS) Proper No Sps (ETS)
³Program Studi Sistem dan Teknologi Informasi, Fakultas Keperawatan Bisnis dan Teknologi, Universitas Widyadarmas Husada Indonesia, Sp. (ETS) Sp. (ETS) Proper No Sps (ETS)

*Email: lariffitrasetyawan@gmail.com, ameliadev26@gmail.com, farichatulfikriyah45@gmail.com

Abstract (10pt)

Dengue Hemorrhagic Fever (DHF) is a tropical disease that often results in high morbidity and mortality rates. Early diagnosis of DHF is crucial to mitigate its adverse effects. However, manual diagnostic processes are often inefficient and prone to errors. This study aims to develop a DHF classification model using the Random Forest algorithm, which is expected to assist in the early diagnosis of this disease. The methodology used in this research is CRISP-DM (Cross-Industry Standard Process for Data Mining), which includes the stages of Business Understanding, Data Understanding, Data Preparation, Modeling, Evaluation, and Deployment. Data was obtained from kaggle.com, and during the Data Preparation stage, missing values were removed, categorical features were encoded, data was normalized, and split into training and testing sets. The research results show that the Random Forest model has an accuracy of 88.5%, precision of 88.2%, recall of 65.2%, F1-score of 74.9%, and ROC AUC of 0.810. Feature importance analysis revealed that the Gender_Male and Body_Pain features have the largest contributions in DHF classification. Although the model demonstrated high accuracy and precision, the lower recall value indicates that some positive cases were missed, requiring further improvements. The Random Forest can be used as a tool for early DHF diagnosis, but further adjustments are necessary to enhance its performance. This research provides insights into the contributing factors for DHF diagnosis and the practical application potential of this model in medical decision support systems.

Keywords: CRISP-DM, Dengue Hemorrhagic Fever, Classification, Machine Learning, Random Forest

INTRODUCTION

Diseases related to the environment remain a public health issue to this day. One of the diseases caused by poor environmental sanitation conditions is dengue fever (Dengue Hemorrhagic Fever or DHF) [1]. Dengue fever is an infectious disease caused by the dengue virus. Dengue Hemorrhagic Fever (DHF) is an epidemic that affects various countries worldwide, with over 500,000 cases reported annually. [2]. This has become a significant global health issue, especially in tropical and subtropical regions. The disease spreads rapidly and is often fatal because many patients die due to delayed treatment. [3]. According to the World Health Organization (WHO), Dengue Hemorrhagic Fever (DHF) or Dengue Fever (DF) is a disease caused by the bite of an Aedes mosquito infected with one of the

four types of dengue virus, presenting with clinical manifestations such as fever, muscle and/or joint pain, along with leukopenia, rash, lymphadenopathy, thrombocytopenia, and hemorrhagic diathesis [4]. Early detection and management of this disease are crucial to reducing its impact. The diagnosis of dengue fever often relies on laboratory test results and clinical symptoms. This can be a challenging task that requires deep medical knowledge. Additionally, the early symptoms of dengue fever often resemble those of other illnesses, such as the flu or other infectious diseases, making it difficult to make an accurate diagnosis.

Information technology and machine learning have unlocked significant potential for supporting disease diagnosis. Data mining is the process of using statistical, mathematical, and

artificial intelligence techniques to extract and identify information and patterns from large datasets. Its primary goal is information extraction, particularly in the form of classification. [5]. Although various terms such as knowledge mining or knowledge discovery are used to refer to data mining, this concept has become a popular method in society. Data mining is also known as knowledge extraction, pattern analysis, data archaeology, information harvesting, and pattern discovery. By uncovering the facts stored within data, data mining generates valuable knowledge through deep analysis processes.

Classification in data mining is a technique that groups data into predefined categories. The goal is to predict the class value of unknown objects by identifying patterns that represent data classes. This process involves validating the model using training data for model creation and testing the model's accuracy with test data. Classification enables the prediction of the name or value of data objects by utilizing identified patterns [6].

To detect early symptoms of diabetes promptly, it is necessary to develop a model and classification using data mining. This research will utilize the Cross-Industry Standard Process for Data Mining (CRISP-DM) methodology and the Random Forest algorithm. It is hoped that this approach can serve as a tool for the public to independently detect diabetes. [7]. Random Forest is a machine learning algorithm used to develop decision trees. It can be considered as a combination of multiple decision trees [8]. Predictions from the Random Forest algorithm are obtained through the majority vote of each individual decision tree (voting process for classification and averaging for regression) [9].

In previous research on Random Forest, it was proven effective in classifying and detecting stroke symptoms, with the initial data divided into 80% for training and 20% for testing. Validation was performed using cross-validation, achieving a training score of 96%, an accuracy of 95%, and an AUC value of 0.80, indicating strong model performance [10]. Research using the Random Forest algorithm has also been conducted to predict diabetes, with results indicating that the algorithm can predict diabetes with high performance. This algorithm has proven to be highly reliable as a reference for developing predictive models for similar cases, as reflected by an AUC value reaching 100%. [11]. Research comparing the Random Forest algorithm and

SVM showed that Random Forest outperformed with an accuracy of 88.2% on the test data and 98.8% on the training data, which is better than SVM. After hypertuning the SVM algorithm, its accuracy improved to 81%, approaching the accuracy level of Random Forest [12]. The Random Forest machine learning algorithm has proven effective in classifying data and predicting outcomes based on patterns identified within the data.

The main objective of this research is to develop a classification model with a high accuracy rate in recognizing dengue fever cases. This model is expected to accurately differentiate between dengue and non-dengue cases based on patient data. Through this study, a reliable classification model is hoped to be produced that can be implemented in healthcare facilities, supporting better clinical decision-making and improving overall public health. This research also aims to identify the most significant clinical features in dengue classification, which can assist in medical decision-making. Thus, it is anticipated that this study will contribute significantly to enhancing the quality of dengue diagnoses and assist healthcare professionals in providing faster and more accurate treatment.

RESEARCH METHOD

This research applies the CRISP-DM (Cross-Industry Standard Process for Data Mining) methodology, which includes six main phases: Business Understanding, Data Understanding, Data Preparation, Modeling, Evaluation, and Deployment [13]. Here are the details of each research phase:

2.1. Business Understanding

This initial step focuses on interpreting business objectives and requirements, which are then translated into insights to define the main problem that can be addressed through data mining [14]. In the Data Understanding phase, clinical data related to Dengue Fever (DF) is sourced from the dataset available on Kaggle.com. This dataset includes essential information such as symptoms, vital signs, laboratory results, and final diagnosis. The initial data exploration process involves understanding the data distribution, identifying missing values, and detecting outliers that may affect the analysis results. Descriptive analysis is performed to gain initial insights into the data characteristics, such as the number of cases, the demographic distribution of patients, and common symptom

patterns. A deep understanding of the data structure and quality is crucial to ensure that the data is ready for the preprocessing phase and the development of machine learning models.

2.2. Data Understanding

In this phase, data visualization is performed to understand the data and clean it by addressing missing data or removing problematic features, aiming to produce a better and more generalizable machine learning model [15]. Clinical data related to Dengue Fever (DF) is sourced from a dataset available on Kaggle.com. This dataset includes crucial information such as symptoms, vital signs, laboratory results, and final diagnoses. The initial data exploration process involves understanding data distribution, identifying missing values, and detecting outliers that may impact the analysis results. Descriptive analysis is performed to gain preliminary insights into data characteristics, including the number of cases, patient demographic distribution, and common symptom patterns. A thorough understanding of the data's structure and quality is essential to ensure the data is ready for the preprocessing stage and the development of machine learning models.

2.3. Data Preparation

Data preparation is a crucial stage in the CRISP-DM methodology aimed at ensuring the data for machine learning models is clean, relevant, and in the appropriate format. In this study, the steps taken include removing missing values to prevent disruption to the model, and encoding categorical features like gender using one-hot encoding to ensure compatibility with the Random Forest algorithm. Irrelevant features are also removed to reduce complexity and improve prediction accuracy. Additionally, data normalization is performed if necessary, to ensure numerical features fall within a uniform range, aiding faster convergence of the algorithm. The dataset is then split into training (80%) and testing (20%) data to prevent overfitting and ensure good model generalization. Finally, balancing techniques such as oversampling or undersampling are applied if class imbalance is present, to prevent bias towards the majority class. These steps ensure optimal data for training the Random Forest model, contributing to

better performance in classifying Dengue Fever.

2.4. Modeling

In the Modeling phase, the Random Forest algorithm is used to develop a classification model for Dengue Fever (DF). Initially, the dataset is divided into training and testing sets using k-fold cross-validation to prevent overfitting. The model is trained on the training data, with hyperparameter tuning performed via grid search or random search to optimize its performance. After training, the model is evaluated using the testing data with metrics such as accuracy, precision, recall, F1-score, and AUC-ROC to assess its effectiveness. Feature importance analysis is also conducted to identify the most significant features for DF classification.

2.5. Evaluation

In the Evaluation phase, the classification model for Dengue Fever (DF) is assessed to measure its performance in predicting DF cases using relevant metrics. After training and optimizing the Random Forest model through hyperparameter tuning, we test the model on a separate testing dataset. The evaluation involves using several metrics, including accuracy, precision, recall, F1-score, and Area Under the Receiver Operating Characteristic Curve (AUC-ROC).

RESULT AND DISCUSSION

This study employs the Random Forest algorithm to classify Dengue Fever (DF) using the CRISP-DM (Cross-Industry Standard Process for Data Mining) methodology. The process begins with the Business Understanding phase, which establishes the goal of developing a predictive model to assist in the early diagnosis of DF. In the Data Understanding phase, data from kaggle.com is analyzed to understand distribution, identify missing values, and explore feature relationships.

In the Data Preparation phase, the data is cleaned by removing missing values, encoding categorical features, and splitting it into training and testing datasets. Data Preparation involves clearing missing values, encoding categorical features, and dividing the data into 80% training and 20% testing sets.

During the Modeling phase, the Random Forest algorithm is utilized to build the classification model, with hyperparameter tuning conducted via Grid Search to find the optimal parameter combination. The model

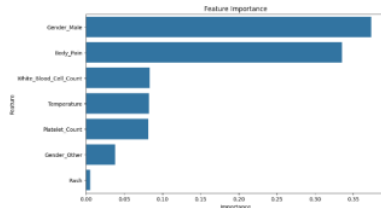
is trained on the training data and tested on the testing data to evaluate its performance.

Model evaluation is carried out using metrics such as accuracy, precision, recall, F1-score, and ROC AUC. The Random Forest model achieves an accuracy of 88.5%, indicating that it correctly classifies about 88.5% of the tested data. This high accuracy suggests the model's potential for DF classification. A precision of 88.2% means that 88.2% of the model's positive predictions are truly positive. The recall of 65.2% shows that the model detects 65.2% of all actual positive cases. Although this recall value is lower than precision, it is acceptable depending on the clinical context and application priorities. Improving recall is crucial to minimize missed positive cases. The F1-score of 0.749 balances precision and recall, reflecting a reasonable performance in balancing true positive detection and minimizing false negatives. The ROC AUC value of 0.810 indicates that the model performs well in distinguishing between positive and negative classes, with a value above 0.8 indicating strong performance. The results from the confusion matrix are shown in the following table:

Table 1. confusion matrix

		Actual values	
		Positive (1)	Negative (0)
Predicted values	Positive (1)	713	23
	Negative (0)	92	172

The model produces 713 true negatives (TN) and 172 true positives (TP). However, there are 23 false positives (FP) and 92 false negatives (FN). The higher number of FN compared to FP suggests that the model tends to overlook positive cases more frequently, indicating a potential area for further improvement. The Feature Importance results are displayed in the following graph; here is the Feature Importance chart generated:



Picture 1. Feature Importance chart

Feature importance analysis reveals that Gender_Male and Body_Pain are the most significant contributors to the model, with importances of 37.4% and 33.5%, respectively. They are followed by White_Blood_Cell_Count (8.3%) and Temperature (8.2%). The feature Rash has the lowest contribution at 0.6%. This analysis helps identify which features most impact the model's predictions and can provide further insights for medical research.

CONCLUSION

This study developed a classification model for Dengue Fever (DF) using the Random Forest algorithm and the CRISP-DM methodology. The evaluation results indicate that the model performs well with an accuracy of 88.5%, precision of 88.2%, recall of 65.2%, F1-score of 74.9%, and ROC AUC of 0.810. The Data Preparation phase involved removing missing values, encoding categorical features, eliminating irrelevant features, normalizing data, and splitting the data into training and testing sets. Feature importance analysis revealed that Gender_Male and Body_Pain are the most significant contributors. Despite the high accuracy and precision, the lower recall indicates that some positive cases were missed, suggesting the need for further improvement. This study provides insights into the factors contributing to DF diagnosis and demonstrates the potential practical application of the model in medical decision support systems.

Classification of Dengue Fever Disease Using a Machine Learning turnitin (1).docx

ORIGINALITY REPORT

18%

SIMILARITY INDEX

14%

INTERNET SOURCES

12%

PUBLICATIONS

9%

STUDENT PAPERS

PRIMARY SOURCES

- 1** Fari Katul Fikriah, Amelia Devi Putri Ariyanto, Arif Fitra Setyawan. "KLASIFIKASI HASIL MRI TUMOR OTAK DENGAN EKTRAKSI FITUR GRAY LEVEL CO-OCCURANCE MATRIX (GLCM)", Rabit : Jurnal Teknologi dan Sistem Informasi Univrab, 2024
Publication 1%
- 2** fastercapital.com
Internet Source 1%
- 3** repositorio.iscte-iul.pt
Internet Source 1%
- 4** www.ijnrd.org
Internet Source 1%
- 5** www.mdpi.com
Internet Source 1%
- 6** Submitted to University of Bradford
Student Paper 1%
- 7** Submitted to INTI Universal Holdings SDM BHD 1%

8	library.unisel.edu.my Internet Source	1 %
9	Submitted to Queen's University of Belfast Student Paper	1 %
10	Yuejiang Chen, Yingjing He, Jiang-Wen Xiao, Yan-Wu Wang, Yuanzheng Li. "Hybrid model based on similar power extraction and improved temporal convolutional network for probabilistic wind power forecasting", Energy, 2024 Publication	1 %
11	Ali Hakan Işık, Özlem Özmen, Ömer Can Eskicioğlu, Nimet Işık, Sadettin Melenli. "Classification and Diagnosis of Mammary Tumors in Dogs Using Deep Learning Techniques", Traitement du Signal, 2023 Publication	1 %
12	Hung Viet Nguyen, Haewon Byeon. "Prediction of ECOG Performance Status of Lung Cancer Patients Using LIME-Based Machine Learning", Mathematics, 2023 Publication	1 %
13	Submitted to British University in Egypt Student Paper	1 %
14	Sandhya Rani Bansal, Savita Wadhawan, Rajeev Goel. "mRMR-PSO: A Hybrid Feature	1 %

Selection Technique with a Multiobjective Approach for Sign Language Recognition", Arabian Journal for Science and Engineering, 2022

Publication

15

Submitted to University of North Texas

Student Paper

1 %

16

Submitted to Louisiana State University

Student Paper

1 %

17

"Intrusion detection by machine learning = Behatolás detektálás gépi tanulás által", Corvinus University of Budapest, 2020

Publication

<1 %

18

dokumen.pub

Internet Source

<1 %

19

healthdocbox.com

Internet Source

<1 %

20

Stella Roussou, Eva Michelaraki, Christos Katrakazas, Amir Pooyan Afghari et al.

"Unfolding the dynamics of driving behavior: a machine learning analysis from Germany and Belgium", European Transport Research Review, 2024

Publication

<1 %

21

essuir.sumdu.edu.ua

Internet Source

<1 %

22	neuroquantology.com Internet Source	<1 %
23	jurnal.kdi.or.id Internet Source	<1 %
24	tutorsonspot.com Internet Source	<1 %
25	Diego José Gallardo Romero, Orly Enrique Apolo-Apolo, Manuel Pérez Ruíz. "Estimating Optimal Harvest Time and Yield in Tomatoes Using Deep Learning Techniques: A Preliminary Study", 2023 IEEE International Workshop on Metrology for Agriculture and Forestry (MetroAgriFor), 2023 Publication	<1 %
26	ijlbpr.com Internet Source	<1 %
27	www.ijritcc.org Internet Source	<1 %
28	www.tgc.com Internet Source	<1 %

Exclude quotes Off

Exclude matches Off

Exclude bibliography Off

Classification of Dengue Fever Disease Using a Machine Learning

turnitin (1).docx

PAGE 1



Sp. This word is misspelled. Use a dictionary or spellchecker when you proofread your work.



Sp. This word is misspelled. Use a dictionary or spellchecker when you proofread your work.



Sp. This word is misspelled. Use a dictionary or spellchecker when you proofread your work.



Sp. This word is misspelled. Use a dictionary or spellchecker when you proofread your work.



Sp. This word is misspelled. Use a dictionary or spellchecker when you proofread your work.



Sp. This word is misspelled. Use a dictionary or spellchecker when you proofread your work.



Proper Nouns You may need to use a capital letter for this proper noun.



Sp. This word is misspelled. Use a dictionary or spellchecker when you proofread your work.



Sp. This word is misspelled. Use a dictionary or spellchecker when you proofread your work.



Sp. This word is misspelled. Use a dictionary or spellchecker when you proofread your work.



Sp. This word is misspelled. Use a dictionary or spellchecker when you proofread your work.



Sp. This word is misspelled. Use a dictionary or spellchecker when you proofread your work.



Proper Nouns You may need to use a capital letter for this proper noun.



Sp. This word is misspelled. Use a dictionary or spellchecker when you proofread your work.



Sp. This word is misspelled. Use a dictionary or spellchecker when you proofread your work.



Sp. This word is misspelled. Use a dictionary or spellchecker when you proofread your work.



Sp. This word is misspelled. Use a dictionary or spellchecker when you proofread your work.



Sp. This word is misspelled. Use a dictionary or spellchecker when you proofread your work.



Proper Nouns You may need to use a capital letter for this proper noun.



Sp. This word is misspelled. Use a dictionary or spellchecker when you proofread your work.



Sp. This word is misspelled. Use a dictionary or spellchecker when you proofread your work.



Sp. This word is misspelled. Use a dictionary or spellchecker when you proofread your work.



Sp. This word is misspelled. Use a dictionary or spellchecker when you proofread your work.



Sp. This word is misspelled. Use a dictionary or spellchecker when you proofread your work.



Proper Nouns You may need to use a capital letter for this proper noun.



Sp. This word is misspelled. Use a dictionary or spellchecker when you proofread your work.



Sp. This word is misspelled. Use a dictionary or spellchecker when you proofread your work.



Sp. This word is misspelled. Use a dictionary or spellchecker when you proofread your work.



Sp. This word is misspelled. Use a dictionary or spellchecker when you proofread your work.



Sp. This word is misspelled. Use a dictionary or spellchecker when you proofread your work.



Proper Nouns You may need to use a capital letter for this proper noun.



Sp. This word is misspelled. Use a dictionary or spellchecker when you proofread your work.



Sp. This word is misspelled. Use a dictionary or spellchecker when you proofread your work.



Sp. This word is misspelled. Use a dictionary or spellchecker when you proofread your work.



Sp. This word is misspelled. Use a dictionary or spellchecker when you proofread your work.



Sp. This word is misspelled. Use a dictionary or spellchecker when you proofread your work.



Proper Nouns You may need to use a capital letter for this proper noun.



Sp. This word is misspelled. Use a dictionary or spellchecker when you proofread your work.



Sp. This word is misspelled. Use a dictionary or spellchecker when you proofread your work.



Sp. This word is misspelled. Use a dictionary or spellchecker when you proofread your work.



Sp. This word is misspelled. Use a dictionary or spellchecker when you proofread your work.



Article Error You may need to use an article before this word.



Article Error You may need to remove this article.



S/V This subject and verb may not agree. Proofread the sentence to make sure the subject agrees with the verb.



Article Error You may need to use an article before this word.



Dup. Did you mean to repeat this word?

PAGE 2



Sp. This word is misspelled. Use a dictionary or spellchecker when you proofread your work.



Proofread This part of the sentence contains an error or misspelling that makes your meaning unclear.



Missing ", " Review the rules for using punctuation marks.



Article Error You may need to use an article before this word.



Article Error You may need to use an article before this word.

PAGE 3



P/V You have used the passive voice in this sentence. You may want to revise it using the active voice.



Sp. This word is misspelled. Use a dictionary or spellchecker when you proofread your work.



Sp. This word is misspelled. Use a dictionary or spellchecker when you proofread your work.



Article Error You may need to remove this article.



Confused You have used either an imprecise word or an incorrect word.



Article Error You may need to use an article before this word.



P/V You have used the passive voice in this sentence. You may want to revise it using the active voice.



Sp. This word is misspelled. Use a dictionary or spellchecker when you proofread your work.



Sp. This word is misspelled. Use a dictionary or spellchecker when you proofread your work.



Sp. This word is misspelled. Use a dictionary or spellchecker when you proofread your work.

PAGE 4



Article Error You may need to remove this article.



Possessive Review the rules for possessive nouns.



P/V You have used the passive voice in this sentence. You may want to revise it using the active voice.