Determining Important Features for Dengue Diagnosis using Feature Selection Methods

Yulianti Paula Bria^{1,*}⁽⁰⁾, Paskalis Andrianus Nani^{2,}⁽⁰⁾, Yovinia Carmeneja Hoar Siki^{3,}⁽⁰⁾, Natalia Magdalena Rafu Mamulak^{4,}⁽⁰⁾, Emiliana Metan Meolbatak^{5,}⁽⁰⁾, Robertus Dole Guntur^{6,}⁽⁰⁾

1.2.3.4.5 Universitas Katolik Widya Mandira, Jl. San Juan No. 1 Penfui Timur, Kabupaten Kupang, Nusa Tenggara Timur 85361, Indonesia

⁶Universitas Nusa Cendana, Jl. Adisucipto Penfui, Kupang, Nusa Tenggara Timur 85001, Indonesia

(Received: August 25, 2024; Revised: October 6, 2024; Accepted: November 21, 2024; Available online: December 27, 2024)

Abstract

This research aims to determine the important features including symptoms and risk factors for dengue diagnosis. This study's dataset was obtained from medical records collected from two hospitals in Indonesia from patients with dengue and nondengue diseases. Four feature selection methods including feature importance, recursive feature elimination, correlation matrix and KBest were leveraged to determine significant features. Feature importance employed a tree-based classifier to derive the importance scores of the features. Recursive feature elimination employed a machine learning classifier to choose the most important features from the given dataset. Correlation matrix was employed to select the best features because it has the ability to use the correlation between each feature with the target. Univariate feature selection – Kbest has the ability to choose the best features based on univariate statistical tests. Important features were also gathered from fifteen Indonesian medical doctors to confirm the results. We used six machine learning techniques for dengue prediction. The random forest classifier yields the highest accuracy for the best combination of features with the accuracy of 0.93 (LR: 0.90 (0.04), KNN: 0.89 (0.04), XGBoost: 0.91 (0.03), RF: 0.93 (0.04), NB: 0.88 (0.09), SVM: 0.89 (0.04)) and precision of 0.90 (LR: 0.86 (0.22), KNN: 0.67 (0.14), XGBoost: 0.77 (0.13), RF: 0.90 (0.13), NB: 0.66 (0.20), SVM: 0.66 (0.18)). This study shows the significant features for dengue diagnosis including fever, fever duration, headache, muscle and joint pain, nausea, vomiting, abdominal pain, shivering, malaise, loss of appetite, shortness of breath, rash, bleeding nose, bitter mouth, temperature and age. This knowledge is pivotal to educate society to seek medical advice when dengue symptoms appear to avoid severe conditions. Arthralgia/joint pain and myalgia/muscle pain are the most significant features for the dengue prediction. This knowledge is important for medical doctors as a starting point for clinic

Keywords: Dengue Fever, Dengue Diagnosis, Feature Selection, Machine Learning

1. Introduction

Dengue infection is a life-threatening disease spread by female mosquitos, Aedes aegypti. This disease is one of the most prevalent diseases in many countries including Indonesia. Over 3.9 billion people across more than 132 countries are at risk of being infected with dengue, with an estimated 40,000 death per year [1]. This number of cases is much bigger than malaria cases worldwide, which is accounted for 249 million cases all over the world [1]. In 2022, Indonesia contributed to 143,266 dengue cases with the mortality rate in the same year with 1,237 people [2]. Based on the report from Ministry of Health of Indonesia in the week-19 2023, Indonesia had 31,380 dengue cases which claimed 246 people [2]. This indicates that dengue eradication must be prioritized by the government and society without ignoring other priority health problems such as tuberculosis, malaria, stunting, etc.

Early-stage dengue diagnosis is challenging since dengue shares similar symptoms to other diseases including malaria, typhoid fever, and even COVID-19. Malaria, for example, shares the same symptoms with dengue fever such as fever, nausea, vomiting and headache [3]. Some countries have their own identified symptoms for dengue fever. Australia, for example, defines the combination of fever, headache, arthralgia, myalgia, rash, nausea and vomiting as dengue symptoms [4]. Whereas Singapore uses the combination of fever, headache, backache, myalgia, rash, abdominal

^{*}Corresponding author: Yulianti Paula Bria (yulianti.bria@unwira.ac.id)

[©]DOI: https://doi.org/10.47738/jads.v5i4.445

This is an open access article under the CC-BY license (https://creativecommons.org/licenses/by/4.0/). © Authors retain all copyrights

discomfort and thrombocytopenia for dengue symptoms [4]. In Indonesia, the guideline for dengue diagnosis and treatment is issued by Ministry of Health of Indonesia, which is used as a reference for medical personnel [5]. This is adopted from WHO dengue case classification [6].

The number of deaths in East Nusa Tenggara (NTT) province from dengue cases in 2022 was 29 out of 3,376 cases [7]. These cases spread all over NTT's districts. Most of the death cases were because of the severe conditions. People often visit the nearest medical center when they identify rash or severe conditions because of the lack of knowledge of dengue symptoms and risk factors [8]. Understanding important features of dengue is beneficial to avoid the progression to severe condition, which can avoid death. This information is helpful to seek medical advice as soon as dengue symptoms appear. The important features are pivotal to develop early-stage dengue detection tools to assist in dengue diagnosis from other dengue-like symptoms diseases such as malaria, typhoid fever and even COVID-19.

Even though there are some guidelines used to diagnose and treat dengue [6], different countries have different symptoms [4]. Therefore, it is essential to identify first significant symptoms that contribute most for dengue prediction in Indonesia especially in East Nusa Tenggara Province, which will be done in this study. This study will also use the combination of symptoms and dengue risk factors that contribute most for dengue diagnosis.

To obtain significant features from datasets, we use feature selection methods. Feature selection methods are often used to minimize the number of input variables that are considered to be the most significant to a machine learning model to improve the model performance [9], [10]. In recent years, numerous publications focus on the implementation of feature selection methods for disease prediction [10], [11]. In the classification stage, most researchers use machine learning techniques such as BayesNet [9], [10], support vector machine [9], [11] and tree-based classifiers [9], [10]. This study aims to provide information about the most stable classifiers for dengue prediction.

The use of tree-based machine learning techniques including decision tree and random forest for dengue prediction has been conducted by Sarma et al. [12]. However, this study only compared the use of these two classifiers in dengue prediction with fair average accuracy results. The use of feature selection for dengue fever has been implemented successfully by Ramasami et al. [13]. They focus on applying feature selection process and relative analysis to enhance the performance of dengue prediction models. In this study, apart from selecting the best model/classifier for dengue prediction, we also apply feature selection methods for the combination of symptoms and risk factors in Indonesian setting to obtain the most significant features for dengue diagnosis. We also include the interview results with the fifteen Indonesian medical doctors to affirm the most significant features determined. We aim to raise the awareness of society regarding the important symptoms and risk factors for dengue diagnosis to avoid the late detection of dengue, which can lead to death.

In Indonesia, dengue prediction research has been focused on the use of machine learning techniques for predicting the dengue outbreak [14], predicting number of dengue incidents [15], forecasting model for dengue fever [16], and focusing spatial modelling for dengue fever [17]. To the best of our knowledge, this study is the first study to elaborate some feature selection methods to determine significant features for dengue diagnosis based on medical records collected. The results will be compared with the knowledge gathered from the fifteen Indonesian medical doctors' knowledge to confirm the results. This study also aims to provide information on important symptoms and factors for dengue diagnosis in Indonesian context and to point out the symptoms that should not be ignored or should be prioritized by medical doctors when diagnosing potential dengue patients.

2. The Proposed Method

Figure 1 shows the approach to obtain significant features for dengue diagnosis. 561 medical records from dengue and nondengue diseases patients including malaria, COVID-19 and typhoid were collected from two Indonesian hospitals in East Nusa Tenggara Province. These medical records which are called dengue dataset consist of 36 symptoms and two risk factors that can be seen in Appendix 1. To find the most significant features for dengue prediction, we employed four commonly used feature selection methods including recursive feature elimination (RFE), feature importance (FI), correlation matrix from Pearson's correlation coefficient (PCC) and KBest with their own thresholds for obtaining the significant features. Each feature selection method generated one feature set. We also conducted interviews with fifteen Indonesian medical doctors, which then generated important symptoms and risk factors for

clinical dengue diagnosis. This knowledge from the fifteen medical doctors was then formed as extra feature sets. The feature sets from the four feature selection methods and from the medical doctors were used as the dataset. We leveraged six commonly used machine learning techniques including logistic regression (LR), k-nearest neighbour (KNN), eXtreme gradient boosting (XGBoost), random forest (RF), Naïve Bayes (NB) and support vector machine (SVM) to show the performance comparison of the feature sets using the performance metrics of accuracy and precision. The most accurate and precise feature set which then determined as the most significant combination of features for dengue diagnosis.



Figure 1. The approach for determining important features for dengue diagnosis

2.1. Data Collection – Medical Records Collection

To obtain the dengue dataset, we conducted the data collection in two hospitals in Kewapante Hospital, Maumere in Sikka District and Soe Hospital in South Central Timor District of NTT Province. Medical records were collected in the department of medical records of each hospital after obtaining the data collection approvals from the hospital directors in each hospital. Medical records of patients diagnosed with dengue fever or other dengue-like symptoms diseases, such as malaria, typhoid fever, COVID-19, dyspepsia, pneumonia, and gastritis, were collected for the years 2017-2023. These two hospitals' medical records were paper-based, requiring manual recording using an Excel spreadsheet. The features recorded from the medical records collected are age, gender, temperature, all recorded symptoms, duration of fever, working diagnosis, laboratory test results and final diagnosis.

The characteristics of collected medical records from the two Indonesian hospitals can be seen in Appendix 1. The total medical records collected (n) is 561 records. The medical records consist of 473 nondengue cases and 88 dengue cases. Features in the form of symptoms are indicated using S and features in the form of risk factors are indicated using F. The collected dataset will then be named as a dengue dataset, which has 36 symptoms (S1 – S36) and two risk factors (F1 – F2). Most of the symptoms are binary in the form of 1 for Yes or Female and 0 for No or Male. The duration of

fever (S2), temperature (S25) and age (F1) are in the form of number. The target in the dataset is Diagnosis, which is the form of the binary value (1 for dengue and 0 for nondengue diseases). % was used to show the percentage of binary values based on n values. Mean was used to show the average of the numerical values. Whereas standard deviation (SD) was used to show how dispersed the set of data is for the numerical values.

2.2. Interview Results with Fifteen Indonesian Medical Records

To confirm the results from the significant features obtained from the feature selection process, we interviewed 15 Indonesian medical doctors about important symptoms and risk factors for clinical dengue diagnosis. The 15 medical doctors work in hospitals and medical centers in East Nusa Tenggara Province. Before the interview process, the 15 medical doctors agreed to sign the consent forms for the interview. These 15 medical doctors were provided with the structured interview questions regarding important symptoms and risk factors for clinical dengue diagnosis. We used the list of the symptoms and risk factors in Appendix 1 as the baseline to determine the important features for clinical dengue diagnosis. We then asked them further questions regarding other symptoms and risk factors outside the given list that are considered important for clinical diagnosis of dengue. The questions were in Bahasa Indonesia. Therefore, we translated it in English.

The fifteen medical doctors indicated by D were asked using the same questions regarding important symptoms and risk factors for clinical diagnosis of dengue fever. The answers from the medical doctors were indicated using Y for Yes and – for No. As shown in table 1, all the fifteen medical doctors agreed that fever is an the most considered symptom for dengue prediction, followed by rash and bleeding nose (14 doctors), fever duration and abdominal pain (13 doctors). Some symptoms were not considered important by medical doctors including chest pain (S11), sneezing (S15), coughing (S16), sore throat (S21), blurry vision (S22), diarrhea (S24), sweating (S26), swallowing pain (S27), pale (S28), jaundice (S29), anemia (S30), black water (S31), constipation (S32), flatulence (S33), feeling anxious (S34) and bleeding coughing (S35). The interview results showed that two extra features arose in the interviews that were considered important by some medical doctors including orbital pain and whether the patients live in endemic areas of dengue or not. However, since we did not have this information in the medical records collected, we did not consider these two symptoms in this study.

Symptom and risk factor	D1	D2	D3	D4	D5	D6	D7	D8	D9	D10	D11	D12	D13	D14	D15	Total Y
Fever (S1)	Y	Y	Y	Y	Y	Y	Y	Y	Y	Y	Y	Y	Y	Y	Y	15
Fever duration (S2)	Y	Y	Y	Y	Y	Y	Y	Y	Y	Y	Y	-	-	Y	Y	13
Headache (S3)	Y	Y	-	-	Y	Y	Y	Y	Y	Y	-	-	-	-	Y	9
Arthralgia/joint pain and Myalgia/muscle pain (S4)	Y	Y	-	-	Y	Y	Y	Y	Y	Y	-	-	Y	-	-	9
Nausea (S5)	Y	-	-	-	-	-	-	-	-	-	Y	-	-	-	-	2
Vomiting (S6)	Y	-	-	Y	-	-	-	-	-	-	Y	-	-	-	-	3
Abdominal pain (S7)	Y	Y	Y	Y	Y	Y	Y	Y	Y	Y	Y	Y	Y	-	-	13
Shivering (S8)	Y	Y	-	-	Y	Y	Y	Y	Y	Y	-	-	-	-	-	8
Body pain (S9)	Y	Y	-	-	Y	Y	Y	Y	Y	Y	-	-	-	-	-	8
Heartburn (S10)	-	-	-	-	-	-	-	-	-	-	-	-	Y	-	-	1
Dizziness (S12)	Y	Y	-	-	Y	Y	Y	Y	Y	Y	-	-	Y	-	-	8
Malaise (S13)	Y	Y	-	-	Y	Y	Y	Y	Y	Y		Y			Y	10
Loss of appetite (S14)	Y	Y	-	-	Y	Y	Y	Y	Y	Y	-	-	-	-	-	8
Shortness of breath (S17)	-	-	-	Y	-	-	-	-	-	-	-	-	-	-	-	1
Rash (S18)	Y	Y	Y	Y	Y	Y	Y	Y	Y	Y	Y	-	Y	Y	Y	14
Bleeding nose (S19)	Y	Y	Y	Y	Y	Y	Y	Y	Y	Y	-	Y	Y	Y	Y	14
Bitter mouth (S20)	Y	Y	-	-	Y	Y	Y	Y	Y	Y	-	-	-	-	-	8
Seizure (S23)	Y	-	-	-	-	-	-	-	-	-	-	-	-	-	-	1
High Temperature (S25)	Y	Y	Y	-	Y	Y	Y	Y	Y	Y	-	-	Y	-	Y	11

Table 1. The summarized symptoms and risk factors from the fifteen Indonesian medical doctors

Orbital pain	-	Y	-	Y	Y	Y	Y	Y	Y	Y	-	-	-	-	-	8
Loss of consciousness (S36)	Y	Y	-	Y	Y	Y	Y	Y	Y	Y	-	Y	Y	-	-	11
Age (F1)	-	Y	-	-	Y	Y	Y	Y	Y	Y	Y	-	-	Y	-	9
Gender (F2)	-	Y	-	-	Y	Y	Y	Y	Y	Y	-	-	-	Y	-	8
Endemic area	-	Y	-	Y	Y	Y	Y	Y	Y	Y	Y	-	Y	-	Y	11

D: medical doctor; Y: Yes

3. Methodology

3.1. Machine Learning Techniques Used

In this study, we employ commonly used machine learning techniques in dengue prediction including SVM [18], [19], [20], RF [19], [20], XGBoost [18], LR [18], [19], KNN [21] to develop dengue classifiers that can accurately distinguishing dengue from nondengue diseases.

3.2. Performance Metrics Used

To evaluate the performance of the classifiers, we use two performance metrics including accuracy and precision. The formula for the two-performance metrics can be seen in equations (1) and (2).

$$Accuracy = \frac{(TN+TP)}{(TN+TP+FN+FP)}$$
(1)

$$Precision = \frac{TP}{(TP+FP)}$$
(2)

Note: TP is the number of dengue records that are correctly classified; TN is the number of nondengue records that are correctly classified; FP is the number of nondengue records classified as dengue; FN is the number of dengue records classified as nondengue.

3.3. Feature Selection Methods

Feature selection filters redundant or irrelevant features [22]. By reducing the number of features, it will minimize the computational cost of the prediction and increase the performance of the machine learning classifier. The feature selection methods assess the relationship between each feature and the target feature and choose the input features that have the strongest correlation with the target feature [23]. The higher the score, the more the feature is related to the target feature. In this study, feature selection methods used to determine important features are feature importance [24], RFE [25], correlation matrix from Pearson's correlation coefficient (PCC) [3], [23] and KBest [22].

3.3.1. Recursive Feature Elimination (RFE)

RFE is one of the feature selection methods that employs a machine learning classifier to choose the most important features from the given dataset. This method is widely used for feature selection in healthcare [25], [26]. This method works by evaluating the importance of the features in the dataset in the form of importance scores, which later will be ranked. The least important features will then be removed. This step will be iterated until the optimal number of features that gives the best model performance is satisfied [27]. RFE normally is paired with Support Vector Machine [26], [28], [29]. However, in this study, other machine learning classifiers that yields high performance will be leveraged to obtain reliable features. In RFE, the most significant features are selected as number one.

3.3.2. Feature Importance (FI) – Random Forest (RF)

Tree-based feature importance classifiers are reliable for future selection. The tree-based classifiers such as random forest can deal with missing data, numerical and categorical data. It also has the functionality to derive the importance scores of the features without additional cost to the training process [30] and handle a higher number of features [31]. In this study, we use RF, which is founded by Breiman [32] employs simple probability to choose the significant features [33]. It uses subsets of sample data and maps the random sample of feature subspaces by creating multiple (k) decision trees to train and predict samples [34].

3.3.3. Correlation Matrix - Pearson's Correlation Coefficient

The correlation matrix has the ability to use the correlation between each feature with the target to select the best features [23], [35], [36]. The correlation matrix has been widely used in healthcare for disease prediction [3], [23]. In the Pearson's correlation matrix, the calculation of the linear correlation strength between features x and y are made to have a pairwise comparison of all n features. Equation 3 shows the Pearson's correlation coefficient r [36].

$$r = \frac{n(\Sigma xy) - (\Sigma x)(\Sigma y)}{\sqrt{[n\Sigma x^2 - (\Sigma x)^2]}\sqrt{[n\Sigma x^2 - (\Sigma x)^2]}}$$
(3)

3.3.4. Univariate Feature Selection – KBest

Univariate feature selection selects the best features based on univariate statistical tests [37]. In this study, we employ the scikit.learn, which allows feature selection routines using SelectKBest function [37]. The SelectKBest function determines the feature scores and the correlation between each feature and the target feature [23]. The higher score indicates is highly correlated to the target feature.

3.4.Ethical Statement

This research was approved by Human Ethics Committee of Widya Mandira Catholic University (reference number: 001/WM.H9/LPPM/SKKEP/X/2023). We obtained the permission to collect the data from the Department of Permission Affair in Sikka, South Central Timor, and in East Nusa Tenggara Province and the directors of the two hospitals including Kewapante Hospital in Sikka and Soe Hospital in South Central Timor. To ensure the confidentiality of medical records collected, we did not include the patients name as part of the dengue dataset.

4. Results and Discussion

4.1. Results

Table 2 shows the four feature selection scores from FI, RFE, CM and KBest for features that meet the threshold. The threshold value for each feature selection method is used to obtain the most important features from FI (\geq =0.030), RFE (1), PCC (\geq =0.100) and KBest (>1.000). From this first process of filtering, some features are eliminated.

Footure (notation)	FI (>0.030):	RFE (1):	PCC (>=0.100):	KBest (>=1.000):	Number of
reature (notation)	(FSFI)	(FSRFE)	(FSPCC)	(FSKBest)	occurrences
age (F1)	0.163a	3	0.230 c	0.316	2
gender (F2)	0.023	1 b	0.070	2.429 d	2
fever (S1)	0.054 a	1 b	0.330 c	0.705	3
fever_duration (S2)	0.141 a	3	0.160 c	0.153	2
headache (S3)	0.037 a	1 b	0.080	3.478 d	3
muscle_joint_pain (S4)	0.038 a	1 b	0.110 c	7.296 d	4
nausea (S5)	0.032 a	1 b	0.140 c	0.116	3
vomiting (S6)	0.023	3	0.080	4.056 d	1
abdominal_pain (S7)	0.022	1 b	0.130 c	8.929 d	3
shivering (S8)	0.013	2	0.120 c	8.349 d	2
body_pain (S9)	0.009	1 b	0.010	0.053	1
heartburn (S10)	0.028	2	0.060	1.826 d	1
chest_pain (S11)	0.007	2	0.090	4.500 d	1
dizziness (S12)	0.016	1 b	0.080	3.906 d	2
malaise (S13)	0.069 a	1 b	0.020	0.350	2
loss_of_appetite (S14)	0.037 a	1 b	0.240 c	0.347	3
sneezing (S15)	0.029	2	0.150 c	0.133	1

Table 2. The number of occurrences of features in the four feature selection methods with their selection results

Journal of Applied Data Science	es
Vol. 6, No. 1, January 2025, pp.	47-59

coughing (S16)	0.062 a	1 b	0.200 c	0.242	3
shortness_of_breath (S17)	0.038 a	1 b	0.230 c	0.301	3
rash (S18)	0.013	1 b	0.200 c	0.236	2
bleeding_nose (S19)	0.078 a	1 b	0.380 c	0.955	3
bitter_mouth (S20)	0.070 a	1 b	0.060	2.149 d	3
sore_throat (S21)	0.000	1 b	0.060	2.088 d	2
blurry_vision (S22)	0.000	1 b	0.030	0.372	1
seizure (S23)	0.000	1 b	0.050	1.317 d	2
diarrhea (S24)	0.011	1 b	0.070	2.501 d	2
temperature (S25)	0.086 a	3	0.090	4.181 d	2
Total selected features	13	19	13	14	
blurry_vision (S22) seizure (S23) diarrhea (S24) temperature (S25) Total selected features	0.000 0.000 0.011 0.086 a 13	1 b 1 b 1 b 3 19	0.030 0.050 0.070 0.090 13	0.372 1.317 d 2.501 d 4.181 d 14	1 2 2 2

a: selected feature for FI; b: selected feature for RFE; c: selected feature for PCC; and d: selected feature for KBest

Table 2 also shows the total number of significant features for each feature selection method. Feature importance from RF has 13 significant features. RFE selects 19 significant features. There are 13 significant features for PCC and 14 significant features for KBest respectively. These results show that each feature selection method has its own combination of significant features. In Colum 6 of table 2, we total number of occurrences for each feature based on the given thresholds from the four feature selection methods. The higher the number of occurrences, the more significant the feature. There are some features that are significant for three or four feature selection methods. Muscle_joint_pain (S4), for example, is the only feature choosed by the four feature selection methods. This indicates that this feature is the most significant features >=2 and FS3 from selected features >=1. FS4 consists of FS1 and selected features = 1. FS5 consists of FSPCC and the selected symptoms from 15 medical doctors.

In order to choose the significant features for dengue diagnosis based on various combination of features, we compare feature sets (FSs) generated. FSFI consists of F1, S1, S2, S3, S4, S5, S13, S14, S16, S17, S19, S20, S24, S25. FSRFE consists of F2, S1, S3, S4, S5, S7, S9, S12, S13, S14, S16, S17, S18, S19, S20, S21, S22, S23, S24. FSPCC consists of F1, S1, S2, S4, S5, S7, S8, S14, S15, S16, S17, S18, S19. FSKBest consists of F2, S3, S4, S6, S7, S8, S10, S11, S12, S20, S21, S23, S24, S25. FS1 consists of S1, S3, S4, S5, S7, S14, S16, S17, S19, S20. FS2 consists of FS1, S2, S13, S18, S8, S12, S21, S23, S24, S25, F1, F2. FS3 consists of FS2, S6, S10, S9, S11, S15, S22. FS4 consists of FS1, S6, S10, S9, S11, S15, S22. FS5 consists of S1, S2, S3, S4, S5, S6, S7, S8, S13, S14, S15, S16, S17, S18, S19, S20, S25, F1. Table 3 shows the performance comparison from various features sets generated. As shown in table 3, the most stable performance for almost all machine learning classifiers is FS5. Therefore, the most significant features for dengue prediction are the combination of features of FS5. The random forest classifier yields the highest accuracy for FS5 with the accuracy of 0.93 and precision of 0.90.

Feature set (FS)	LR		KNN		XGE	Boost	R	F	Ν	В	SVM	
	Acc	Pre	Acc	Pre	Acc	Pre	Acc	Pre	Acc	Pre	Acc	Pre
	(SD)	(SD)	(SD)	(SD)	(SD)	(SD)	(SD)	(SD)	(SD)	(SD)	(SD)	(SD)
FSFI	0.91	0.89	0.90	0.70	0.90	0.71	0.91	0.81	0.82	0.50	0.89	0.66
	(0.04)	(0.17)	(0.04)	(0.17)	(0.03)	(0.14)	(0.04)	(0.21)	(0.10)	(0.15)	(0.05)	(0.23)
ECDEE	0.89	0.78	0.86	0.57	0.83	0.47	0.83	0.47	0.42	0.21	0.87	0.80
FSKFE	(0.06)	(0.32)	(0.06)	(0.25)	(0.07)	(0.23)	(0.06)	(0.23)	(0.07)	(0.05)	(0.05)	(0.40)
ESDCC	0.91	0.84	0.90	0.69	0.90	0.70	0.92	0.83	0.89	0.70	0.90	0.70
rsree	(0.04)	(0.19)	(0.03)	(0.14)	(0.04)	(0.20)	(0.04)	(0.16)	(0.06)	(0.18)	(0.04)	(0.17)
ESKBost	0.85	0.10	0.82	0.15	0.79	0.23	0.80	0.22	0.30	0.18	0.84	0.00
rskbest	(0.04)	(0.30)	(0.03)	(0.32)	(0.05)	(0.21)	(0.03)	(0.30)	(0.04)	(0.04)	(0.04)	(0.00)
FS1	0.88	0.73	0.88	0.69	0.89	0.76	0.89	0.71	0.82	0.54	0.88	0.87

Table 3. Performance comparison of features sets generated with the standard deviation values

	(0.05)	(0.33)	(0.06)	(0.24)	(0.07)	(0.24)	(0.07)	(0.23)	(0.16)	(0.19)	(0.05)	(0.30)
FS2	0.90	0.72	0.90	0.70	0.89	0.70	0.92	0.88	0.66	0.31	0.88	0.59
	(0.04)	(0.30)	(0.04)	(0.17)	(0.01)	(0.16)	(0.03)	(0.13)	(0.10)	(0.09)	(0.05)	(0.24)
EC2	0.90	0.74	0.89	0.66	0.92	0.77	0.93	0.89	0.56	0.25	0.88	0.61
1.35	(0.05)	(0.33)	(0.03)	(0.14)	(0.02)	(0.18)	(0.03)	(0.13)	(0.08)	(0.06)	(0.04)	(0.22)
ES/	0.90	0.77	0.88	0.65	0.87	0.59	0.88	0.58	0.57	0.25	0.88	0.87
г54	(0.05)	(0.35)	(0.05)	(0.22)	(0.05)	(0.17)	(0.04)	(0.18)	(0.09)	(0.07)	(0.05)	(0.31)
FS5	0.90	0.86	0.89	0.67	0.91	0.77	0.93	0.90	0.88	0.66	0.89	0.66
	(0.04)	(0.22)	(0.04)	(0.14)	(0.03)	(0.13)	(0.04)	(0.13)	(0.09)	(0.20)	(0.04)	(0.18)

Acc: Accuracy; Pre: Precision; SD: Standard deviation

4.2. Discussion

4.2.1. Principal Results

Table 3 shows that significant features for dengue prediction are fever (S1), fever duration (S2), headache (S3), muscle joint pain (S4), nausea (S5), vomiting (S6), abdominal pain (S7), shivering (S8), malaise (S13), loss of appetite (S14), sneezing (S15), coughing (S16), shortness of breath (S17), rash (S18), bleeding nose (S19), bitter mouth (S20), temperature (S25) and age (F1). However, not all these features are dengue symptoms. It is important to note that the dataset consists of dengue records and other medical records including malaria, COVID-19, dyspepsia, gastritis, typhoid fever and pneumonia. We will discuss which symptoms and risk factors that are important for dengue predictions or dengue diagnosis with the confirmation of medical doctors' knowledge.

Fever, fever duration and high temperature are three important dengue symptoms. For fever, three out of four feature selection methods select this symptom as an important feature. All fifteen medical doctors interviewed also agree that one of the most important dengue features is fever. Even though only two feature selection methods including FI and PCC chose fever duration as an important feature, fever normally starts 4-10 days after infection and last for 2-7 days [38]. Based on the medical doctors interviewed and medical records collected, temperature also has a significant contribution for dengue prediction that can reach 39-40°C. Eleven medical doctors interviewed agree that high temperature of fever is important to distinguish dengue from other diseases such as malaria and typhoid fever. Therefore, it is important to include fever, fever duration and high temperature of fever as three important features for dengue diagnosis.

Arthralgia/joint pain and myalgia/muscle pain are two symptoms that are considered as the most significant features for the dengue prediction and dengue diagnosis [39]. All the four feature selection methods indicate these two symptoms are important for distinguishing dengue from other diseases including malaria, typhoid fever, COVID-19, dyspepsia and pneumonia. Nine medical doctors also consider these two symptoms as significant symptoms for dengue diagnosis.

Headache is one of the most important symptoms in diagnosing and predicting dengue [39]. The three feature selection methods other than PCC consider this symptom essential for dengue diagnosis. In addition, it is also confirmed by nine medical doctors interviewed.

Nausea is considered as one of the most significant symptoms for dengue diagnosis [39]. That also applies for vomiting [40]. However, if persistent vomiting occurs then the individual might progress to the severe state [39]. Two medical doctors agree that nausea is part of dengue symptoms whereas three medical doctors agree that vomiting is an important symptom for dengue diagnosis. In the prediction perspective, nausea is more considered significant because it is selected by three feature selection methods. Whereas vomiting is least significant as only KBest selects this symptom. However, these two symptoms are highly correlated, thus it is important to consider both symptoms as dengue symptoms.

Loss of appetite is considered a symptom that can indicate individuals suffer from dengue. Eight medical doctors interviewed confirm that this symptom is also considered as a dengue symptom. This symptom is also selected by three feature selection methods other than KBest.

Even though shivering is associated with malaria [3], [39], shivering is also important for dengue diagnosis and prediction. Eight medical doctors interviewed also agree that shivering is also a dengue symptom. In the dengue prediction perspective, shivering is also an important feature for dengue prediction as it is selected by two feature selection methods including PCC and KBest as part of significant features.

Malaise is an important symptom for dengue diagnosis, and it normally happens when individuals are in severe condition [39]. In addition, ten medical doctors also confirm that this symptom is essential in dengue diagnosis. It is also selected by two feature selection methods including FI and RFE.

Bleeding nose is one of the most important symptoms in dengue diagnosis as part of bleeding manifestations [39], [40]. This symptom with other bleeding manifestations indicate that individuals progress is in severe condition. Fourteen medical doctors interviewed agree that to determine an individual suffers from dengue is to check the presence of the bleeding nose. Moreover, three feature selection methods selected this symptom as a significant feature for dengue prediction.

Similar to the bleeding nose, the presence of rashes in skin is also pivotal in distinguishing dengue from other similar diseases such as malaria and typhoid fever [39]. Fourteen medical doctors confirm that a rash in an individual's body is a distinguishing symptom that led their initial diagnosis to dengue. This symptom is also selected in two feature selection methods including RFE and PCC.

Abdominal pain is considered as one of the dengue symptoms especially when someone in the severe state [39], [40]. Thirteen medical doctors also confirm that this symptom is essential to determine dengue from other diseases. This symptom is also selected by three feature selection methods other than FI as a significant symptom for dengue prediction.

Shortness of breath or fast breathing is one of dengue symptoms that indicates the severe state of dengue [39]. This is also confirmed by one medical doctor interviewed. This symptom is also selected by three feature selection methods other than KBest as the important feature for dengue prediction.

Age can be considered as one of the important risk factors for dengue diagnosis [41]. Even though six medical doctors do not consider this factor as an important feature for dengue diagnosis, nine medical doctors include this factor as feature that should not be overlooked when diagnosing potential dengue patients. Two feature selection methods including FI and PCC also consider this factor important for dengue prediction. Normally, individuals younger than 15 years old are prone to dengue infection [42].

Bitter mouth is associated with malaria as this symptom is considered as one of malaria symptoms [43]. However, interestingly eight medical doctors interviewed agree that this symptom also can be found in individuals who suffer from dengue. This symptom also appears in three feature selection methods other than PCC. Thus, this symptom should not be ignored when diagnosing potential dengue patients.

From the dengue prediction perspective, sneezing and coughing are important features. Three feature selection methods select this symptom as significant features for dengue prediction. However, no medical doctors confirm that sneezing and coughing are part of dengue symptoms. Sneezing and coughing might be the distinguished symptom to determine COVID-19 from dengue. It is important to know that the dataset consists of medical records from COVID-19 patients. Besides, sneezing and coughing are known as COVID-19 [44]. Therefore, sneezing and coughing are important for dengue prediction but not necessarily are dengue symptoms.

Based on the discussion above, we conclude that there are 16 features that are significant for dengue prediction including fever, fever duration, headache, muscle and joint pain, nausea, vomiting, abdominal pain, shivering, malaise, loss of appetite, shortness of breath, rash, bleeding nose, bitter mouth, temperature and one risk factor feature including age.

4.2.2. The Implications of Significant Features for Clinical Practice

The fifteen significant features excluding bitter mouth symptom are commonly used in the clinical diagnosis of dengue. An Indonesian digital health, ayosehat.kemkes.go.id [45], shows that fever, fever duration, headache, muscle and joint pain, nausea, vomiting, malaise, loss of appetite, shortness of breath, rash, bleeding nose, temperature, abdominal pain,

shivering and age have been considered as symptoms and the risk factor for dengue diagnosis. The finding of this study show that the medical doctors should not ignore bitter mouth symptom when diagnosing potential dengue patients. Arthralgia/joint pain and myalgia/muscle pain are considered as the most significant features for the dengue prediction. Therefore, in diagnosing potential dengue patients, medical doctors should consider this symptom as the starting point.

4.2.3. Limitations

This study does not include other features such as orbital pain, history of previous suffering from dengue and history of visiting endemic dengue areas. In this study, all this information were not found in the medical records collected. We also realized that to draw more comprehensive conclusion, the number of medical records used as the dataset should be increased. The interview results with medical records did not involve the inter-rater reliability as the interview results were used for the affirmation of the significant features generated by machine learning techniques based on the given dataset.

4.2.4. Possible Future Work

The significant features as results from this study can be used to develop reliable and powerful machine learning techniques, which later can be used to develop early-stage dengue prediction tools. We can also further extend the study to rank the significant features.

5. Conclusion

In conclusion, there are four findings of this study. First, there are 17 symptom features including fever, fever duration, headache, muscle and joint pain, nausea, vomiting, abdominal pain, shivering, malaise, loss of appetite, sneezing, coughing, shortness of breath, rash, bleeding nose, bitter mouth, temperature and one risk factor feature including age that are important for dengue prediction. However, sneezing and coughing are not necessarily important for dengue diagnosis. Second, arthralgia/joint pain and myalgia/muscle pain are the most significant features for the dengue prediction. Third, even though a bitter mouth symptom is highly related to malaria diagnosis, this study suggests that the medical doctors should not ignore the bitter mouth symptom in diagnosing dengue as this symptom is also important for dengue prediction. Fourth, random forest classifier yields the most stable performance for dengue prediction. Knowledge of these features are essential to educate society about significant symptoms and risk factors for dengue to avoid progression to severe conditions, which can lead to death. The findings of this study can also be used as a reference for medical doctors in differentiating dengue from nondengue diseases including malaria, COVID-19 and typhoid fever.

6. Declarations

6.1. Author Contributions

Conceptualization: Y.P.B., P.A.N., Y.C.H.S., N.M.R.M., E.M.M., and R.D.G.; Methodology: Y.P.B., P.A.N., Y.C.H.S., N.M.R.M., E.M.M., and R.D.G.; Software: Y.P.B.; Validation: Y.P.B., E.M.M., and R.D.G.; Formal Analysis: Y.P.B., E.M.M., and R.D.G.; Investigation: Y.P.B.; Resources: E.M.M.; Data Curation: E.M.M.; Writing Original Draft Preparation: Y.P.B., E.M.M., and R.D.G.; Writing Review and Editing: Y.P.B., P.A.N., Y.C.H.S., N.M.R.M., E.M.M., and R.D.G.; Visualization: Y.P.B.; All authors have read and agreed to the published version of the manuscript.

6.2. Data Availability Statement

The data presented in this study are not publicly available to maintain the confidentiality of the medical records collected.

6.3. Funding

This study was funded by Widya Mandira Catholic University Kupang East Nusa Tenggara Province Indonesia (044/WM.H9/SKP/IX/2023)

6.4. Institutional Review Board Statement

This research was reviewed and approved by Human Ethics Committee of Widya Mandira Catholic University (reference number: 001/WM.H9/LPPM/SKKEP/X/2023).

6.5. Informed Consent Statement

We obtained written consents from the 15 medical records interviewed.

6.6. Declaration of Competing Interest

The authors declare that we have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

References

- [1] World Health Organization, "Vector-borne diseases," Sep. 2024, accessed Oct. 9, 2024.
- [2] Ministry of Health of Indonesia, "Information on DHF Cases in 2023 Week 19," accessed Aug. 19, 2023. (In Indonesian)
- [3] Y. P. Bria, C. H. Yeh, and S. Bedingfield, "Significant symptoms and nonsymptom-related factors for malaria diagnosis in endemic regions of Indonesia," *Int. J. Infect. Dis.*, vol. 103, no. 1, pp. 194–200, 2021.
- [4] World Health Organization, "Update on the dengue situation in the Western Pacific Region," vol. 481, 2015, accessed Mar. 20, 2024.
- [5] Ministry of Health of Indonesia, "National guidelines for medical services for the management of dengue infection in children and adolescents," pp. 1–67, 2021, accessed Jan. 10, 2024. (In Indonesian)
- [6] World Health Organization, "Dengue guidelines for diagnosis, treatment, prevention and control," 2009, accessed Feb. 15, 2024.
- [7] Central Bureau of Statistics East Nusa Tenggara Province, "Number of disease cases according to regency/city and type of disease (Inhabitant) in 2022," 2022, accessed Jan. 10, 2024. (In Indonesian)
- [8] A. N. Rakhmani and L. Zuhriyah, "Knowledge, attitudes, and practices regarding dengue prevention among health volunteers in an urban area Malang, Indonesia," *J. Prev. Med. Public Health*, vol. 57, no. 2, pp. 176–184, 2024.
- [9] Z. Noroozi, A. Orooji, and L. Erfannia, "Analyzing the impact of feature selection methods on machine learning algorithms for heart disease prediction," *Sci. Rep.*, vol. 13, no. 1, pp. 1–15, 2023.
- [10] J. D. Álvarez, J. A. Matias-Guiu, M. N. Cabrera-Martín, J. L. Risco-Martín, and J. L. Ayala, "An application of machine learning with feature selection to improve diagnosis and classification of neurodegenerative disorders," *BMC Bioinformatics*, vol. 20, no. 1, pp. 1–12, 2019.
- [11] J. Song, Z. Li, G. Yao, S. Wei, L. Li, and H. Wu, "Framework for feature selection of predicting the diagnosis and prognosis of necrotizing enterocolitis," *PLoS One*, vol. 17, no. 8, pp. 1-22, Aug. 2022.
- [12] D. Sarma, S. Hossain, T. Mittra, M. A. M. Bhuiya, I. Saha, and R. Chakma, "Dengue prediction using machine learning algorithms," in 2020 IEEE 8th R10 Humanitarian Technology Conference (R10-HTC), Kuching, Malaysia, 2020, vol. 12. No. 1, pp. 1–6.
- [13] V. Ramasamy, S. Vadivel, S. Kothandapani, J. Mahilraj, P. Sivaram, and B. Sharma, "An optimal feature selection with neural network-based classification model for dengue fever prediction," in 6th Int. Conf. Inf. Syst. Comput. Networks, 2023, vol. 6, no. 1, pp. 1–5.
- [14] A. L. Ramadona, Y. Tozan, J. Wallin, L. Lazuardi, A. Utarini, and J. Rocklöv, "Predicting the dengue cluster outbreak dynamics in Yogyakarta, Indonesia: A modelling study," *Lancet Reg. Heal. Southeast Asia*, vol. 15, no. 1, pp. 1-8, 2023.
- [15] I. N. Tanawi, V. Vito, D. Sarwinda, H. Tasman, and G. F. Hertono, "Support vector regression for predicting the number of dengue incidents in DKI Jakarta," *Proceedia Comput. Sci.*, vol. 179, no. 1, pp. 747–753, 2021.
- [16] N. A. Lestari, R. Tyasnurita, R. A. Vinarti, and W. Anggraeni, "Long short-term memory forecasting model for dengue fever cases in Malang regency, Indonesia," *Procedia Comput. Sci.*, vol. 197, no. 1, pp. 180–188, 2022.
- [17] S. A. Thamrin, Aswi, Ansariadi, A. K. Jaya, and K. Mengersen, "Bayesian spatial survival modelling for dengue fever in Makassar, Indonesia," *Gac. Sanit.*, vol. 35, no. 1, pp. S59–S63, 2021.

- [18] A. Joshi and C. Miller, "Review of machine learning techniques for mosquito control in urban environments," *Ecol. Inform.*, vol. 61, no. 101241, pp. 1-14, 2021.
- [19] W. Hoyos, J. Aguilar, and M. Toro, "Dengue models based on machine learning techniques: A systematic literature review," *Artif. Intell. Med.*, vol. 119, no. 102157, pp. 1-16, Aug. 2021.
- [20] M. S. G. Shaikh, D. B. SureshKumar, and D. G. Narang, "Development of optimized ensemble classifier for dengue fever prediction and recommendation system," *Biomed. Signal Process. Control*, vol. 85, no. 104809, pp. 1-13, Mar. 2023.
- [21] Y. P. Bria, C. H. Yeh, and S. Bedingfield, "Machine learning classifiers for symptom-based malaria prediction," in *Proc. Int. Jt. Conf. Neural Networks*, vol. 2022, no. 1, pp. 1-6, Jul. 2022.
- [22] P. Qiu and Z. Niu, "TCIC_FS: Total correlation information coefficient-based feature selection method for high-dimensional data," *Knowledge-Based Syst.*, vol. 231, no. 107418, pp. 1-12, 2021.
- [23] E. M. Senan, I. Abunadi, M. E. Jadhav, and S. M. Fati., "Score and correlation coefficient-based feature selection for predicting heart failure diagnosis by using machine learning algorithms," *Comput. Math. Methods Med.*, vol. 2021, no. 1, pp. 1-16, 2021.
- [24] R. C. Chen, C. Dewi, S. W. Huang, and R. E. Caraka, "Selecting critical features for data classification based on machine learning methods," J. Big Data, vol. 7, no. 52, pp. 1-26, 2020.
- [25] T. E. Mathew, "A logistic regression with recursive feature elimination model for breast cancer diagnosis," *Int. J. Emerg. Technol.*, vol. 10, no. 3, pp. 55–63, 2019.
- [26] H. M. Alshanbari, T. Mehmood, W. Sami, W. Alturaiki, M. A. Hamza, and B. Alosaimi, "Prediction and classification of COVID-19 admissions to intensive care units (ICU) using weighted radial kernel SVM coupled with recursive feature elimination (RFE)," *Life*, vol. 12, no. 7, pp. 1-10, 2022.
- [27] P. Misra and A. S. Yadav, "Improving the classification accuracy using recursive feature elimination with cross-validation," *Int. J. Emerg. Technol.*, vol. 1, no. 3, pp. 659–665, 2020.
- [28] X. Huang, L. Zhang, B. Wang, F. Li, and Z. Zhang, "Feature clustering based support vector machine recursive feature elimination for gene selection," *Appl. Intell.*, vol. 48, no. 3, pp. 594–607, 2018.
- [29] S. Zhou, T. Li, and Y. Li, "Recursive feature elimination based feature selection in modulation classification for MIMO systems," *Chinese J. Electron.*, vol. 32, no. 4, pp. 785–792, 2023.
- [30] A. Alsahaf, N. Petkov, V. Shenoy, and G. Azzopardi, "A framework for feature selection through boosting," *Expert Syst. Appl.*, vol. 187, no. 115895, pp. 1-10, Feb. 2021.
- [31] R. C. Chen, C. Dewi, S. W. Huang, and R. E. Caraka, "Selecting critical features for data classification based on machine learning methods," *J. Big Data*, vol. 7, no. 52, pp. 1-26, 2020.
- [32] L. Breiman, "Random forests," Mach. Learn., vol. 45, no. 8, pp. 5–32, 2001.
- [33] S. Gündoğdu, "Efficient prediction of early-stage diabetes using XGBoost classifier with random forest feature selection technique," *Multimed. Tools Appl.*, vol. 82, no. 22, pp. 34163–34181, 2023.
- [34] K. Mei, M. Tan, Z. Yang, and S. Shi, "Modeling of feature selection based on random forest algorithm and Pearson correlation coefficient," J. Phys. Conf. Ser., vol. 2219, no. 012046, pp. 1-9, 2022.
- [35] I. M. Nasir, M. A. Khan, M. Yasmin, J. H. Shah, M. Gabryel, R. Scherer, and R. Damaševičius, "Pearson correlation-based feature selection for document classification using balanced training," *Sensors*, vol. 20, no. 23, pp. 1-18, 2020.
- [36] C. A. Rickert, M. Henkel, and O. Lieleg, "An efficiency-driven, correlation-based feature elimination strategy for small datasets," APL Mach. Learn., vol. 1, no. 016105, pp. 1-14, 2023.
- [37] D. P. M. Abellana, and D. M. Lao, "A new univariate feature selection algorithm based on the best-worst multi-attribute decision-making method," *Decis. Anal. J.*, vol. 7, no. 100240, pp. 1-13, 2023.
- [38] G. Gupta, S. Khan, V. Guleria, A. Almjally, B. I. Alabduallah, T. Siddiqui, B. M. Albahlal, S. A. Alajlan, and M. AL-subaie, "DDPM: A dengue disease prediction and diagnosis model using sentiment analysis and machine learning algorithms," *Diagnostics*, vol. 13, no. 6, pp. 1-15, 2023.
- [39] World Health Organization, "Dengue and severe dengue," 2023, accessed Aug. 6, 2023.
- [40] S. N. N. Tatura, D. Denis, M. S. Santoso, R. F. Hayati, B. J. Kepel, B. Yohan, R. T. Sasmono, "Outbreak of severe dengue associated with DENV-3 in the city of Manado, North Sulawesi, Indonesia," *Int. J. Infect. Dis.*, vol. 106, no. 1, pp. 185–196, 2021.

- [41] M. S. Santoso, B. Yohan, D. Denis, R. F. Hayati, S. Haryanto, L. Trianty, R. Noviyanti, M. L. Hibberd, and R. T. Sasmono, "Diagnostic accuracy of 5 different brands of dengue virus non-structural protein 1 (NS1) antigen rapid diagnostic tests (RDT) in Indonesia," *Diagn. Microbiol. Infect. Dis.*, vol. 98, no. 2, pp. 1-7, 2020.
- [42] Ministry of Health of Indonesia, "Dengue hemorrhagic fever," 2022, accessed Feb. 5, 2024. (In Indonesian)
- [43] E. Nwokolo, C. Ujuju, J. Anyanti, C. Isiguzo, I. Udoye, E. Bongos-Ikwue, O. Ezire, M. Raji, and W. A. Oyibo, "Misuse of artemisinin combination therapies by clients of medicine retailers suspected to have malaria without prior parasitological confirmation in Nigeria," *Int. J. Heal. Policy Manag.*, vol. 7, no. 6, pp. 542–548, 2018.
- [44] N. S. Romero-Castro, I. C. Hernández, M. E. G. Reyes, M. H. Hernández, A. G. Verónica, S. Paredes-Solis, and S. R. Fernández, "Clinical signs and symptoms associated with COVID-19: A cross-sectional study," *Int. J. Odontostomatol.*, vol. 16, no. 1, pp. 112–119, 2022.
- [45] Ministry of Health of Indonesia, "Demam berdarah dengue," 2022, accessed Oct. 9, 2024. (In Indonesian)