

Comparative Analysis of Machine Learning Classifiers for the Prediction of Malaria Incidence Attributed to Climatic Factors

Pallavi Mohapatra (✉ mohapatra.pallavi@gmail.com)

Asian Institute of Technology <https://orcid.org/0000-0003-3491-7469>

N.K. Tripathi

Asian Institute of Technology

Indrajit Pal

Asian Institute of Technology

Sangam Shrestha

Asian Institute of Technology

Research

Keywords: Machine Learning, Malaria prediction, J48 Decision Tree, WEKA, Multilayer Perceptron

Posted Date: August 7th, 2020

DOI: <https://doi.org/10.21203/rs.3.rs-52162/v1>

License:  This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

1 Research Article

2 **Comparative Analysis of Machine Learning Classifiers for the Prediction of Malaria**

3 **Incidence Attributed to Climatic Factors**

4

5 Pallavi Mohapatra*

6 Remote Sensing and Geographic Information System, Asian Institute of Technology, Pathum

7 Thani, Thailand,

8 Email: mohapatra.pallavi@gmail.com

9

10 N. K. Tripathi

11 Remote Sensing and Geographic Information System, Asian Institute of Technology, Pathum

12 Thani, Thailand,

13 Email: nitinkt@ait.ac.th

14

15 Indrajit Pal

16 Disaster Preparedness Mitigation and Management, Asian Institute of Technology, Pathum

17 Thani, Thailand.

18 Email: indrajit-pal@ait.ac.th

19

20 Sangam Shrestha

21 Water Engineering and Management, Asian Institute of Technology, Pathum Thani, Thailand;

22 Email: sangam@ait.ac.th

23 * Author for correspondence (e-mail: mohapatra.pallavi@gmail.com; Tel: +66-873369790)

24 **Abstract**

25 **Background:** India has a rising rate of malaria as well as a high mortality rate despite awareness
26 and efforts being focused on the issue. Some regions are profoundly affected than others, such as
27 in Odisha, where the prevalence of malaria is nearly a third of the whole country. This study
28 investigated the influence of climate factors on the incidence of malaria in the Sundargarh
29 district in the state of Odisha, India.

30 **Methods:** Block-wise observed station rainfall data was sourced from the Special Relief
31 Commissioners' (SRC) web portal. Gridded surface maximum temperature and relative humidity
32 data were accessed from the European Center for Medium-range Weather Forecast (ECMWF)
33 reanalysis data archive. Malaria incident data were collected from the Directorate of Public
34 Health, Government of Odisha. WEKA machine learning tool with two classifier techniques,
35 Multi-Layer Perceptron (MLP) and J48 with 10-fold cross-validation, percentile split (66%), and
36 supplied test options, were used for the Malaria prediction. A comparative analysis was carried
37 out on both techniques to ascertain the superior model amongst the two, concerning the
38 prediction accuracy of malaria in the context of a varying climate. Classifier accuracy, Root
39 Mean Square Error (RMSE), Kappa, and ROC scores were the indicators used for the analysis.

40 **Results:** The results suggested that J48 had exhibited a better skill to MLP and illustrated less
41 error with a positive kappa. In particular, the 10-fold cross-validation method had better
42 performance over the percentile Spilt (66%) and supplied test options. J48 demonstrated less
43 error (RMSE = 0.6), better kappa = 0.63, and higher accuracy = 0.71), suggesting it as most
44 suitable model. Further, seasonal temperature and humidity variation had shown a better
45 association with malaria incidents in comparison to rainfall.

46 **Conclusion:** The performance of the machine learning methods for Sundargarh was particularly
47 better during the monsoon and post-monsoon when the events are at the peak. The results were
48 encouraging for the utilization of climate forecast for prediction of malaria incidences. It is thus
49 recommended that the J48 classifier machine learning technique could be adopted for the
50 development of malaria early warning system.

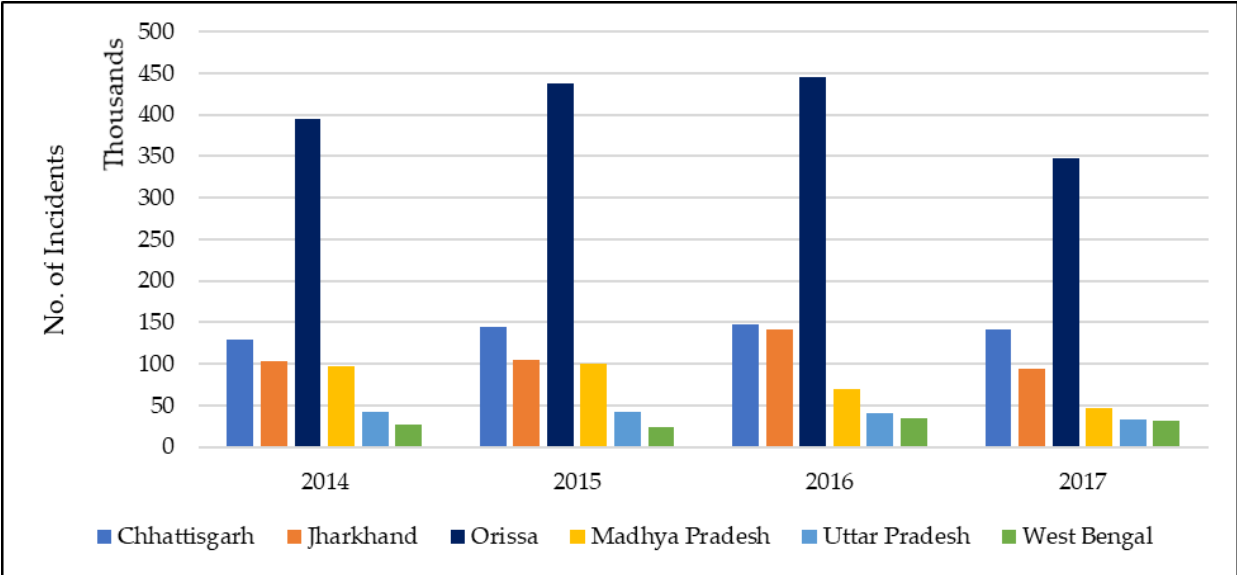
51 **Keywords:** Machine Learning, Malaria prediction, J48 Decision Tree, WEKA, Multilayer
52 Perceptron

53

54 **Background**

55 Malaria remains one of the perennial public health concerns in many parts of the world, even
56 with the efforts put in place by the World Health Organization (WHO) and other international
57 and national bodies to curb it [1]. According to Kovats et al. [2], malaria is characterized by
58 seasonal transmission and distribution of vectors and is influenced by seasonal climatic
59 variations [3]. This is because both vectors and parasites tend to be sensitive to changes in
60 atmospheric temperature and moisture [4]. The distribution of malaria is limited by the climate
61 tolerance of the mosquito vectors, and the biological restrictions that limit the incubation and the
62 survival of the infective agent in the vector population [5]. The examination of how climate
63 conditions could affect the spreading of malaria can be approached by closely monitoring
64 various aspects that change in the climate and the surrounding environment. Van et al. [6]
65 examined the spatio-temporal effects of climate change on malaria. They established that
66 significant changes in the temperature and rainfall patterns could lead to an increase in the
67 spreading of malaria.

68 Malaria prevalence depends on the parasite *Plasmodium* and population dynamics of *Anopheles*
 69 mosquito [7]. The development, as well as the survival rates of both *Plasmodium* parasites and
 70 the *Anopheles* mosquitoes, is dependent on weather. As Kakmeni et al. [8] explain, more
 71 specifically, the temperature is key to the persistence of these parasites. Current evidence
 72 suggests that inter-decadal and inter-annual variability of the climate have a direct effect on the
 73 epidemiology of some of the critical vector-borne diseases [2]. Odisha, an eastern coastal state in
 74 India, has the maximum number of incidents of malaria and casualties since 2014, compared to
 75 other states (provinces) of India, as per the statistics (as shown in Figure 1) provided by the
 76 National Vector Borne Disease Control Programme [9]. The figure shows only the states with
 77 the highest number of incidents. The geographical positioning of the state made it susceptible to
 78 climate extremes and adversely affecting human health.

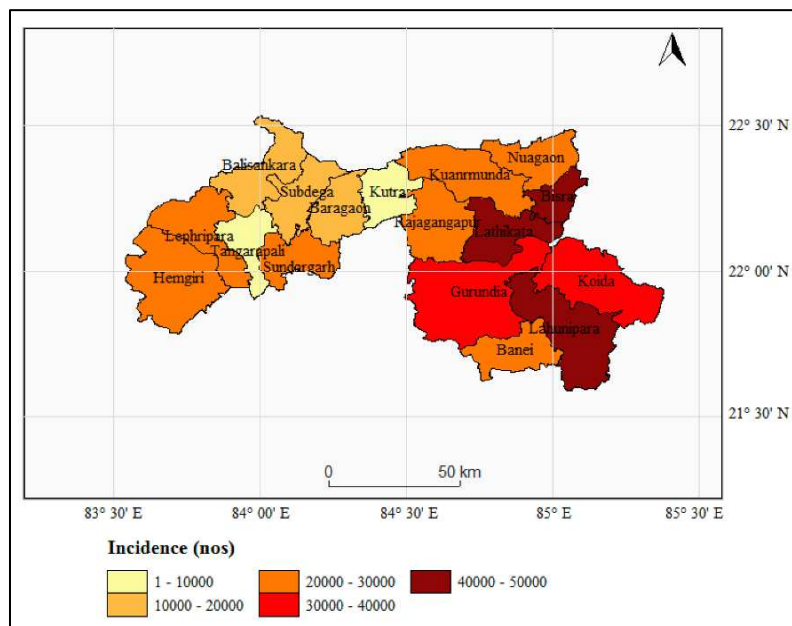


79 **Figure 1.** Malaria incidents across six different states of India for the period of 2014-2017.
 80

81 (Source: NVBDCP Malaria situation reports)

82 The current research performs an analysis using advanced machine learning approaches to
 83 determine how different climate conditions are related to the transmissions of malaria and the

84 possibility of accurately predicting malaria incidence for the Sundargarh district in Odisha, India.
 85 Sundargarh district has the second-highest malaria incidents in the state, followed by the
 86 Rayagada district. Block-wise cumulative incidence map for the district is presented in Figure 2,
 87 which suggest the eastern region is severely affected by Malaria, in particular. The study has
 88 evaluated the efficiency of the machine learning algorithms by determining the accuracy level at
 89 which they were able to predict the malaria incidents. Further, the findings would also encourage
 90 better utilization of climate forecasts to predict potential malaria risk.



91
 92 **Figure 2.** Cumulative malaria incidence map for the period 2002 to 2017 (Data Source:
 93 Directorate of Public Health, Odisha)

94 **Malaria as a Public Health Concern**

95 There were 219 million cases of malaria globally in 2017 [1], and an estimated 228 million cases
 96 of malaria occurred worldwide in 2018 [10]. The burden was most substantial in the African
 97 region, where an estimated 93% of all malaria deaths occurred, and in children aged under five
 98 years, who accounted for 61% of all deaths [1]. Almost 85% of all malaria cases globally were in

99 19 countries, including India and 18 African countries. In India, seven states accounted for 90%
100 of the estimated cases in 2018, counting to 5.7 million cases [10]. Malaria is prevalent in eastern,
101 central, and north-eastern states, especially in ethnic groups, usually, dominate these tribal areas.
102 Inequality and poverty in this area play a crucial role in the spreading as well [11,12]. The most
103 vulnerable community to malaria is the tribal populations in India habitually reside in remote
104 areas with complex topography and dense forest with limited to no access to basic facilities [13].

105 **Influence of Climate on Malaria**

106 Statistical methods were used by several researchers to investigate the association of climatic
107 factors and malaria incidents which included the multiple polynomial regression to model
108 malaria incidents in India [14], semi-parametric Poisson distribution methods to model the
109 influence of temperature and rainfall on malaria incidence in Zambia [15], distributed non-linear
110 lag model to associate malaria to meteorological factors in China [16], hierarchical Bayesian
111 framework to model effects of weather and climate on malaria distributions in West Africa [17]
112 and the time series regression models [18,19]. All the models have shown reasonable skills over
113 the respective regions. Neter et al. [20] recommend the use of Multiple Linear Regression (MLR)
114 in the analysis of data because this model can determine the relative influence of one or more
115 predictor variables. Other advanced computational models include Artificial Neural Network
116 (ANN) models, which are relatively simple to interpret [21] and, as such, require less formal
117 training. They also can, implicitly, detect complex non-linear relationships between the set of
118 variables being investigated. Yao et al. [22] demonstrated that using neural network for data
119 analysis could detect all possible interactions among the predictors.

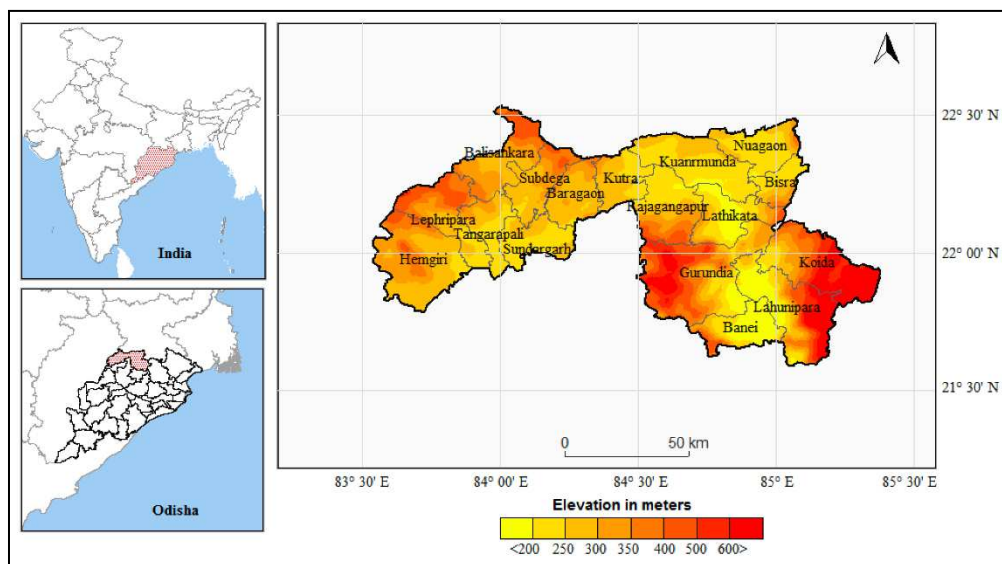
120 This research acknowledges the existing problem of the malaria epidemic in the region and the
121 influence of increased frequency of extreme climate events such as floods, heatwaves, drought,

122 which contribute further to the escalation of malaria spreading. Malaria is no doubt a significant
123 threat to human life, and climate variability plays a key role in the survival and abundance of the
124 disease vectors. The researchers analyzed the data using the machine learning methods and
125 quantify the accuracy and skill level for predicting the incidence.

126 **Methods**

127 **Study Area**

128 The statistics from the National Vector Borne Disease Control Programme revealed that Odisha,
129 a coastal province in India, records the maximum number of casualties due to vector-borne
130 diseases and especially from malaria [9]. This research targeted the Sundargarh district, which
131 records one of the highest numbers of cases of malaria incidents in the state. The geographical
132 location of Odisha made it susceptible to increased occurrences of climate extremes and
133 adversely affecting human health due to the environmental changes. Sundargarh district forms
134 the north-western part of Odisha state and is the second-largest district in the state, accounting
135 for 6.23% of the total area. The geographical area of the district is 9712 square km. The district
136 spreads from 21°36'N to 22°32'N and 83°32'E to 85°22'E [23].



137

138 **Figure 3.** Study area with elevation for Sundargarh district in Odisha, India

139

140 **Topography**

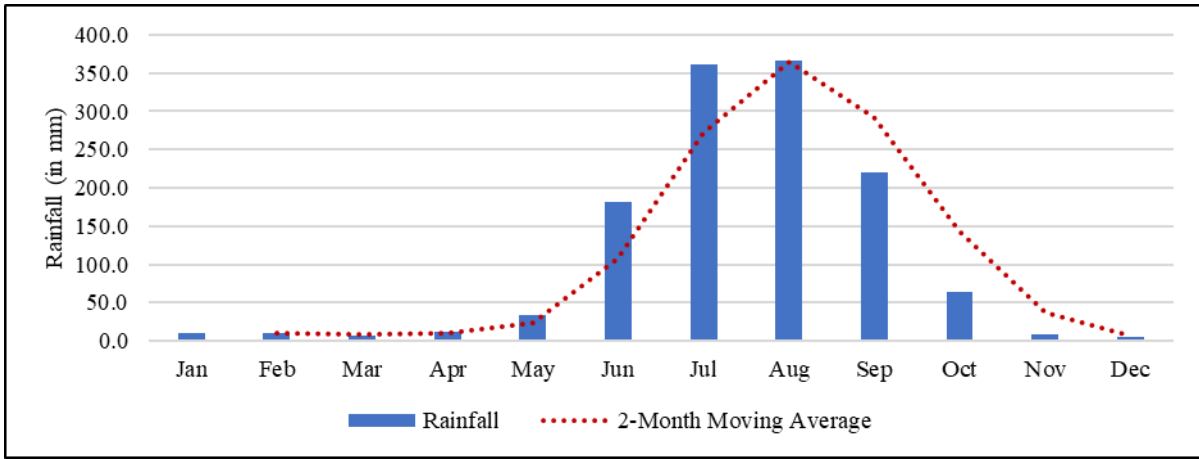
141 The district exhibits an ideal ecological condition for the malaria transmission topographically
142 with its undulating uplands intersected by forested hills and widely diversified tracts of
143 mountains. The areas covered by western blocks are long undulating tracts of about 700 ft. (213
144 mt.) above the sea level, dotted with hill-ranges and isolated peaks of considerable height. At the
145 same time, the far eastern and southern-central blocks are mostly an isolated hilly tract with an
146 average elevation of about 800 ft. (244 mt.) above sea level.

147 **Climate of Sundargarh**

148 From the southwest monsoon, the area received rainfall between June and September and
149 characterized as a tropical humid climate region (shown in Figure 4 (a) and (b)). The average
150 annual temperature ranges between 22°C and 27°C and the average annual rainfall ranges
151 between 1600 and 2000 mm. The weather seasons are hot and dry summer from April to mid-
152 June, monsoon from mid-June to September, autumn from October to November, winter from
153 December to January, and spring from February to March. The maximum temperature during
154 summer rises to 40–45°C and the minimum temperature during winter falls to 5-10°C.

155

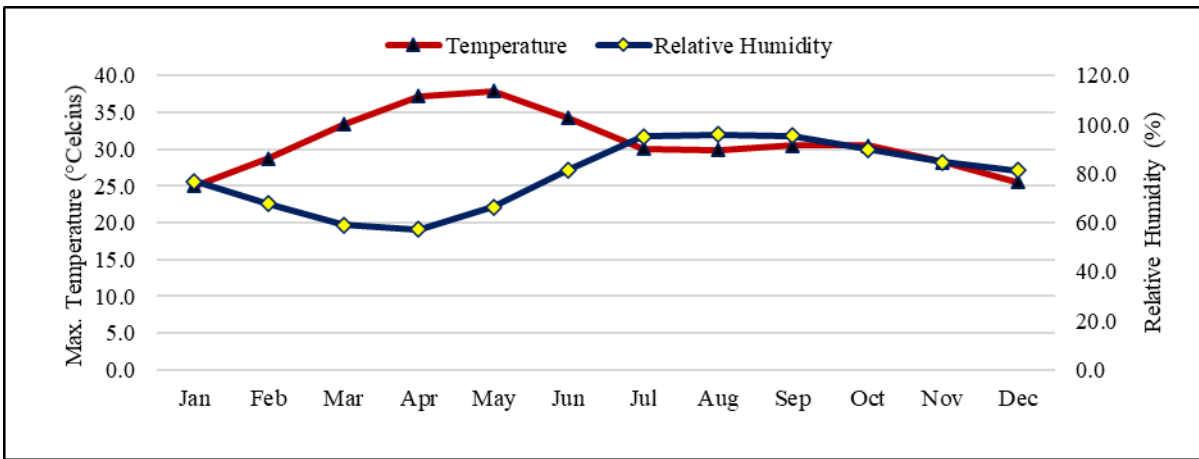
156



157

158

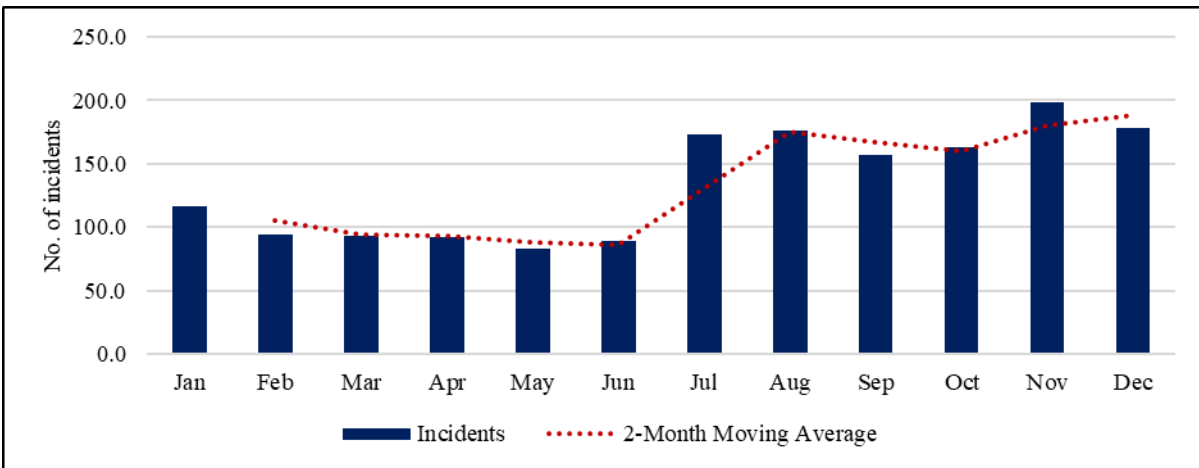
(a)



159

160

(b)

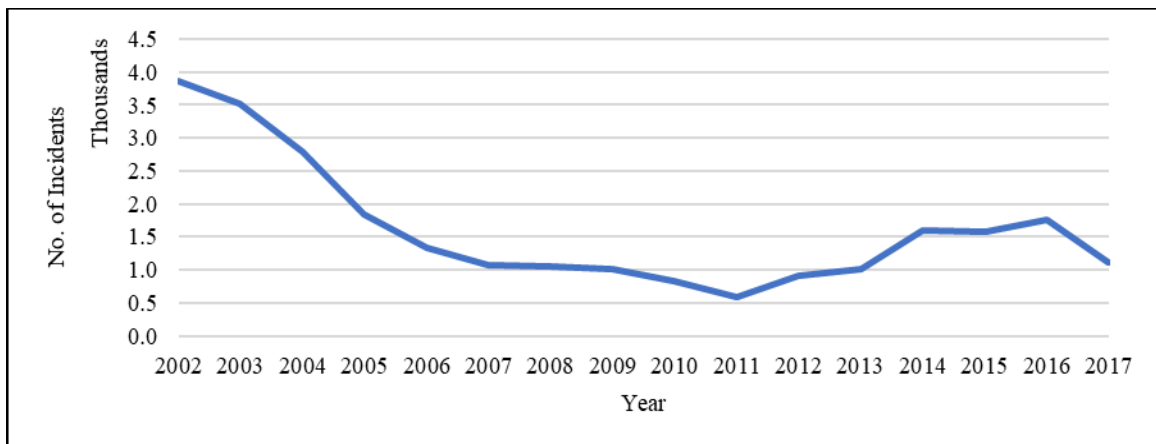


161

162

(c)

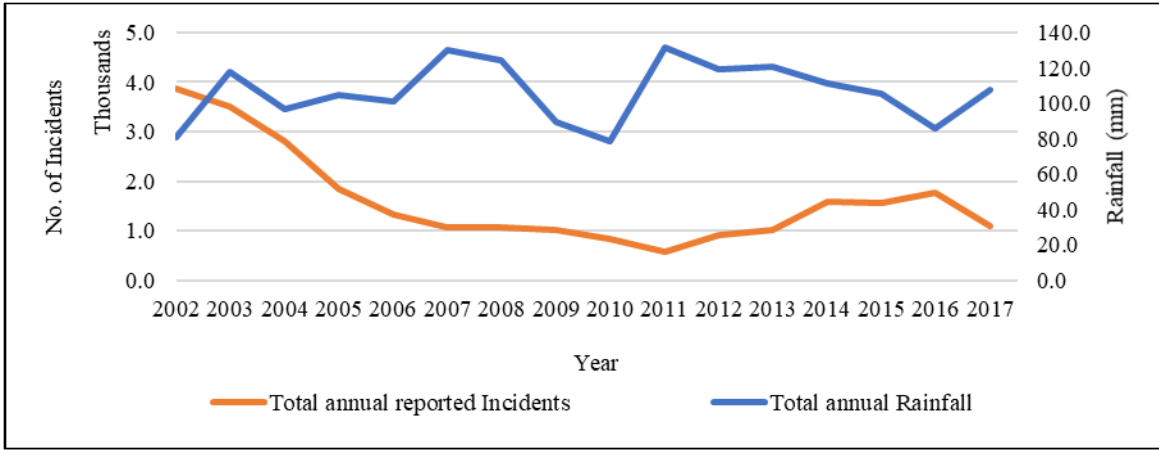
163 **Figure 4.** Average monthly rainfall (in mm) (a), maximum surface temperature (in °C), relative
164 humidity (in %) (b), and incidents reported (c) for the period of 2002 to 2017.
165 Seasonal variation in temperature plays an important role as far as vector diseases, and their
166 transmission is concerned. According to Servadio et al. [24], the most significant effect of
167 climate change on the transmission of malaria would be felt at the extreme temperature ranges.
168 This reiterates the importance of studying climate variability and determining how its changes
169 can affect the transmission of malaria, not just in places that are already affected by the disease,
170 but also in fresher areas. Figure 4 provides the seasonal influence of the incidents to the rainfall,
171 maximum temperature, and relative humidity, and Figure 5 shows the trend of incidents for 2002
172 to 2017 and its relationship with the three climate factors under consideration in this study.



173

174

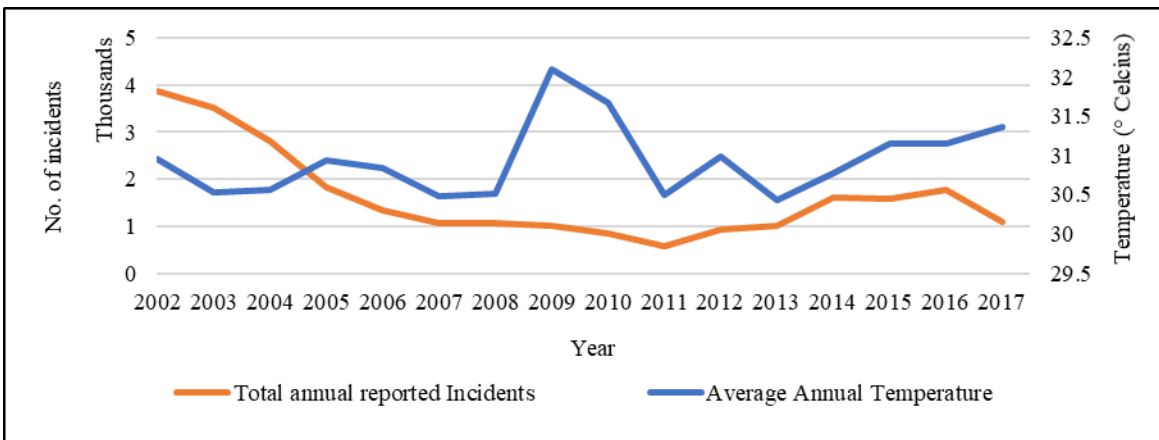
(a)



175

176

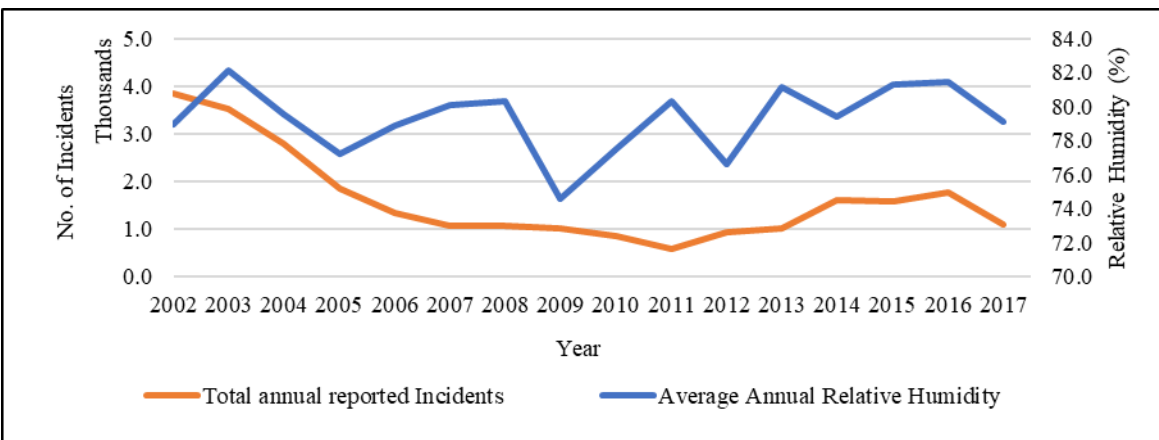
(b)



177

178

(c)



179

180

(d)

181 **Figure 5.** The trend of annual incidents (a); and comparison with average annual rainfall (b);
182 temperature (c); and relative humidity (d) for the period of 2002 to 2017

183 **Data Used**

184 The variables used in this research were climate parameters like rainfall (RF), relative humidity
185 (RH), and surface (2-meter height from ground) maximum temperature ($T_{2_{max}}$) and malaria
186 incidents. While the climate parameters are treated as independent variables, the malaria incident
187 records are considered the dependent variable. For this analysis, historical meteorological rainfall
188 data from 17 blocks for the district are accessed from the Odisha Government portal of special
189 relief commissioner (http://srcodisha.nic.in/rain_fall.php), which is a publicly accessible portal.
190 Surface maximum temperature and relative humidity data with a horizontal resolution of 0.1-
191 degree was obtained from the Copernicus Climate data store (CDS) of the European Center for
192 Medium-Range Weather Forecast (ECMWF). The most recent ECMWF Reanalysis (ERA5-
193 Land) is a reanalysis of the global atmosphere covering the data-rich period since 1981 and
194 continuing in real-time. More details about the dataset can be found from the Copernicus climate
195 data store
196 (<https://cds.climate.copernicus.eu/cdsapp#!/dataset/10.24381/cds.e2161bac?tab=overview>). The
197 data was taken at daily temporal time scales for both climate parameters. Many studies have
198 demonstrated use of gridded reanalysis data in the absence of actual ground observation [25, 26,
199 27, 28] as it is the closest possible representation of the actual observations. Monthly malaria
200 incident datasets at block level were collected from the Directorate of Public Health Services,
201 Government of Odisha. For consistency in the analysis, all data were collected for the identical
202 period of 2002-2017. Table 1 summarizes the data used in the study.

203 **Table 1.** Data Source and the Attributes

Type of Data	Data Source	Period	Spatial Scale	Temporal scale
Malaria Incidents	The Directorate of Public Health, Odisha	2002 -2017	Block	Monthly
Rainfall	Special Relief Commissioner, Odisha	2002-2017	Block/Station	Daily
Surface Max. Temperature	ECMWF Reanalysis land data (ERA5-Land)	2002-2017	Gridded. 0.1°x0.1°; Native resolution is 9 km	Daily
Relative Humidity	ECMWF Reanalysis land Data (ERA5-Land)	2002-2017	Gridded. 0.1°x0.1°; Native resolution is 9 km	Daily

204 Data Preparation

205 The data received were at different temporal scales and required to be brought to a standard
206 spatial and temporal scale. As the malaria incidents are the parameter to be predicted and were
207 available at a monthly scale, the climate data (which are at daily time scales) were statistically
208 averaged over the month. The ERA5-Land data for maximum temperature and relative humidity
209 were extrapolated to produce average spatial data for the blocks. The percentile (p=25, p=50,
210 p=75, and p=95) were computed for each sample to define the spread of the variables. The next
211 step in the analysis was changing the historical data into range values, and turning it from
212 numerical to nominal (Low, Medium, High, and Very High). A sample conversion is shown in
213 Table 2.

214 **Table 2.** Conversion of Numeric data to Nominal data

Year	RF (mm)	T _{2max} (°C)	RH (%)	Incidents	RF.	T _{2max}	R.H.	Incidents
Numeric Data					Nominal Data			
2002	24	24.7	82.6	111	L	L	L	L

2002	15	32.8	66.1	129	L	M	L	M
2002	43	36.3	72.0	175	M	H	L	M
2002	308	28.8	98.1	250	H	L	VH	H
2002	169.5	29.5	96.9	195	H	L	H	H
2002	84.4	30.5	92.9	177	M	M	M	M

215 L=Low, M=medium, H=High, VH=Very High

216 **Weka Machine Learning Tool**

217 Waikato Environment for Knowledge Analysis (WEKA) is a collection of machine learning
218 algorithms that accurately perform data mining tasks [29]. WEKA contains tools that facilitate
219 data preparation, regression, classification, association rules mining, clustering, and
220 visualization. WEKA, through its machine learning platform, enables the algorithm to learn
221 about data as samples and with or without the interference of any other explicit programs
222 [30,31]. More detail about the tool is available at <https://www.cs.waikato.ac.nz/~ml/weka/>.
223 Multilayer Perceptron (MPL) and J48 classifier techniques in the Weka tool recently being used
224 in successfully predicting malaria incidents [32-34]. Researchers around the globe also used it
225 for prediction of dengue [35-38] and other public health issues such as Cholera [39], diabetes
226 [40-42], heart diseases [43, 44].

227 **Multiple Layer Perceptron**

228 A Multilayer Perceptron (MLP) is a class of feed-forward artificial neural networks [45]. It
229 constitutes at least three layers of nodes, a hidden layer, an input layer, and an output layer. Each
230 of these nodes, except the input nodes, is a neuron that uses a non-linear activation function.
231 ANN has, for a long time, been a robust perceptive classifier for tasks not just in medical
232 diagnosis, but also for early detection of diseases [46]. MLP uses a supervised learning technique
233 that is referred to as propagation for training the network [45]; it is a modification of the standard

234 linear perceptron. As such, it can distinguish data that is not separable. A perceptron produces a
235 single output based on several real-valued inputs by forming a linear combination using its input
236 weights [46]. Which can be represented in the following form:

$$y = \varphi (\sum_{i=1}^n w_i x_i + b) = \varphi (w^T x + b), \quad (1)$$

237 Where w denotes the vector of weights, x is the vector of inputs, b is the bias, and φ is the non-
238 linear activation function.

239 MLP is composed of an input layer to receive the signal, an output layer that decides or predict
240 the input, and in between those two, an arbitrary number of hidden layers that are the actual
241 computational engine of the MLP. They train on a set of input-output pairs and learn to model
242 the Correlation (or dependencies) between those inputs and outputs [46]. Training involves
243 adjusting the parameters, or the weights and biases, of the model to minimize error. Back-
244 propagation is used to make those weight and bias adjustments relative to the error, and the error
245 itself can be measured in a variety of ways.

246 **J48 Classifier Model – A Decision Tree Based Method**

247 J48 in WEKA is the implementation of the C4.5 decision tree [45, 47]. Dangare & Apte [44]
248 defined J48 classification as building models of classes from records that contain class labels. A
249 decision tree algorithm is used to find how the attribute-vector is likely to behave for an array of
250 instances. The algorithm generates rules that would be used for the prediction of the targeted
251 variables and accounts for any missing values present in the model and the output. Some
252 algorithms perform classification recursively until each leaf has been deemed pure [45]. In other
253 words, the classification of data would be as perfect as possible. The objective of the J48
254 classification is to reduce the impurity or uncertainty in data as much as possible. A subset of
255 data is pure if all instances belong to the same class. The heuristic is to choose the attribute with

256 the maximum Information Gain or Gain Ratio based on information theory. Entropy is a measure
 257 of the uncertainty associated with a random variable. We choose the attribute with the highest
 258 gain to split the current tree. Assuming the attributes are categorical, a tree is constructed in a
 259 top-down recursive manner. At the start, all the training samples are at the root, and samples are
 260 partitioned recursively based on selected attributes. Attributes are selected based on an impurity
 261 function (e.g., information gain). This process uses the "Entropy," i.e., a measure of the disorder
 262 of the data [45, 47, 48]. The Entropy of \vec{E} is calculated as:

$$Entropy(\vec{E}) = - \sum_{j=1}^n \frac{|E_j|}{|\vec{E}|} \log \frac{|E_j|}{|\vec{E}|} \quad (2)$$

263 iterating over all possible values of \vec{E} . The conditional Entropy is

$$Entropy(j | \vec{E}) = \frac{|E_j|}{|\vec{E}|} \log \frac{|E_j|}{|\vec{E}|} \quad (3)$$

264 and finally, the gain is

$$gain(\vec{E}, j) = entropy(\vec{E}) - Entropy(j|\vec{E}) \quad (4)$$

265 The aim is to maximize the gain, dividing by overall Entropy due to split argument \vec{E} by value j.

266 **Predictive Modeling using MLP and J48**

267 The steps followed for the predictive modelling using MLP and J48 are presented in Figure 6.

268 The climate datasets collected from the respective sources were reprocessed to monthly scale as

269 the prediction of malaria incidents was expected to be carried out at a monthly time scale. The

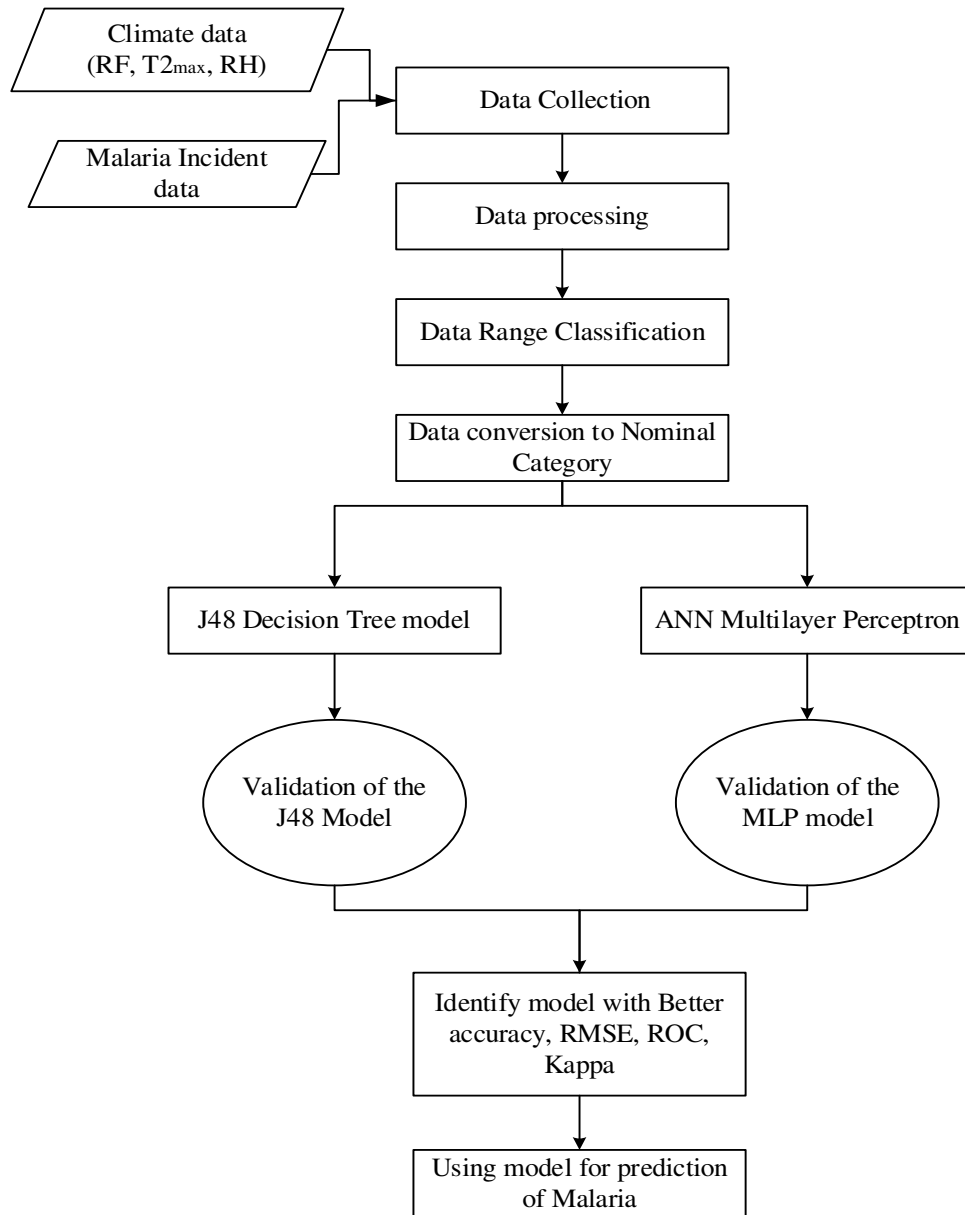
270 numerical monthly malaria incident and climate data were then, transformed to nominal range,

271 before they were fed to the MLP and J48 classifier models. The WEKA tool [30] was used as a

272 base platform for all the analyses. The datasets were split into two sets; first set for training of

273 the model and the second set for testing or the prediction. Different test options used include; (a)

274 10-fold cross-validation method, in which all samples were divided as ten equal sets, from which
275 1 set is used for testing, and the rest nine sets for the training of the model; (b) percentage split
276 method, in which the data is distributed as a percent of the total number of samples with 34%
277 data are used for testing and rest 66% data are for training; and (c) a supplied test set, which
278 enables users to decide on the distribution of the samples for training and prediction. Various
279 indicators including the RMSE, Kappa, ROC, accuracy was used to evaluate the performance of
280 the techniques used for the prediction. From the investigation, the better performing technique
281 with the most appropriate test option was identified. Further, the technique and test methods
282 would be used as a malaria prediction engine for the prediction of malaria though a Malaria early
283 warning system.



284

285

Figure 6. Detailed Methodology for Predictive Modeling of Malaria

286

Performance Indicators

287

With the classifiers, we investigate how good both the models are, and this is done by examining

288

the number of correctly classified instances to the number of incorrectly classified cases from the

289

supplied datasets. The performance of the machine learning analysis methods was evaluated

290

through different indicators which were inbuilt in the tool. These include the Root Mean Square

291 Error (RMSE), the accuracy, the kappa, and the Receiver Operating Characteristics (ROC)
 292 values. This section provides a brief about each of these indicators and its significance. The
 293 confusion matrix (Table 3 and Table 4) provides a simplified structure of the representation of
 294 the observed and expected samples to segregate the classifications into four classes: True
 295 Positives (*a*), False Positives (*b*), False Negatives (*c*) and True Negatives (*d*).

296 **Table 3.** Confusion Matrix for observed agreement

		Observed		
Expected		Positive	Negative	Total
	Positive	<i>a</i> (TP)	<i>b</i> (FP)	<i>a+b</i>
	Negative	<i>c</i> (FN)	<i>d</i> (TN)	<i>c+d</i>
Total		<i>a+c</i>	<i>b+d</i>	<i>N</i>

297 TP=True positives; FP = False Positives; FN=False Negatives; TN=True Negatives

298 The observed agreement is the frequency with which the two variants (observed and expected)
 299 agreed. From the confusion matrix, the observed agreement can be determined as:

$$Observed\ agreement = \frac{a + d}{N} \quad (5)$$

300 **Table 4.** Confusion Matrix for expected agreement

		Observed		
Expected		Positive	Negative	Total
	Positive	$\frac{(a+b)(a+c)}{N}$	$\frac{(a+b)(b+d)}{N}$	<i>a+b</i>
	Negative	$\frac{(a+c)(c+d)}{N}$	$\frac{(c+d)(b+d)}{N}$	<i>c+d</i>
Total		<i>a+c</i>	<i>b+d</i>	<i>N</i>

$$expected\ agreement = \frac{expected\ (a) + expected\ (d)}{N} \quad (6)$$

302 Where, expected (*a*) = $\frac{(a + b)(a + c)}{N}$ and expected (*d*) = $\frac{(c + d)(b + d)}{N}$

303 **Accuracy**

304 The percentage of correctly classified instances is often called accuracy. The basic formula for
305 calculation of prediction accuracy can be described as (referring to the confusion matrix for
306 observed agreement):

$$accuracy = \frac{a + d}{N} \quad (7)$$

307 Where $a = True\ Positives$ and $d = Correct\ Negatives$.

308 **Kappa Coefficient**

309 Kappa is the measurement of the inter-rater reliability, which represents the extent to which the
310 data collected in the study are correct representations of the variables measured [49]. The
311 formula for kappa is:

$$Kappa = \frac{P_o - P_e}{1 - P_e} \quad (8)$$

312 Where, $P_o = Observed\ agreement$; $P_e = Expected\ agreement$

313 Kappa coefficients are interpreted using the guidelines outlined by [50], where the strength of the
314 kappa coefficients is interpreted in the following manner: 0.01-0.20 slight; 0.21-0.40 fair; 0.41-
315 0.60 moderate; 0.61-0.80 substantial; 0.81-1.00 almost perfect. A negative kappa would indicate
316 agreement worse than that expected by chance.

317 **Root Mean Square Error (RMSE)**

318 RMSE is used to measure the difference between the expected and the observed values from the
319 environment that is being modeled [51]. The RMSE values can be used to distinguish model
320 performance in a training period with that of a validation period as well as to compare the
321 individual model performance to that of other predictive models. The RMSE of a model prediction
322 for the estimated variable X_{pred} is defined as the square root of the mean squared error:

$$RMSE = \sqrt{\frac{\sum_{i=1}^n (X_{obs,i} - X_{pred,i})^2}{n}} \quad (9)$$

323 Where X_{obs} = observed values

324 X_{pred} = modelled values at time/place i .

325 n = total number of sample datasets

326 Receiver Operating Characteristics (ROC)

327 ROC is a curve that characterizes the randomly chosen probability of positive instance over
328 negative instances [51]. It is a measure of the skill of different classifiers with the true positives

329 (TP) to the false-positive rates (FPR). Setting $P_{1,j}$ as the prediction probability for the j^{th}
330 observed event, and $P_{0,i}$ as the prediction probability of an event for the i^{th} non-event, the ROC

331 score, A , can be

$$A = \frac{1}{n_0 n_1} \sum_{i=1}^{n_0} \sum_{j=1}^{n_1} I(P_{0,i}, P_{1,j}) \quad (10)$$

332 where n_0 is the number of non-events and n_1 the number of events and the scoring rule $I(P_{0,i},$

333 $P_{1,j})$ is defined as;

$$I(P_{0,i}, P_{1,j}) = \begin{cases} 0.0 & \text{if } P_{1,j} < P_{0,i} \\ 0.5 & \text{if } P_{1,j} = P_{0,i} \\ 1.0 & \text{if } P_{1,j} > P_{0,i} \end{cases} \quad (11)$$

334 In the ROC score, a hit is the selected observations are events. The proportion of all events thus

335 selected is calculated, and is known as the hit rate (HR):

$$HR = \frac{\text{No. of TP}}{\text{No. of Event}} \quad (12)$$

336 Some non-events may have been selected incorrectly; these are known as false positives. The

337 proportion of non-events incorrectly chosen [the false-positive rate (FPR)] is:

$$FPR = \frac{\text{No. of False positives}}{\text{No. of non Event}} \quad (13)$$

338 The ROC classifications are excellent, good, fair, poor, fail having range of [0.90-1], [0.80-0.90],
 339 [0.70-0.80], [0.60-0.70], [0.50–0.60], respectively [51].

340 **Taylor Diagram**

341 Taylor Diagram [52], provides a concise statistical summary and illustrates the matching patterns
 342 of data, for their Coefficients of Correlation, the root-mean-square, and their variances and ratio.
 343 Plots are patterned that precisely indicate measures and scores. The statistical display of the three
 344 factors on Taylor Diagram mathematically represented using the following formula:

$$E'^2 = \sigma_r^2 + \sigma_t^2 - 2\sigma_r\sigma_t\rho \quad (14)$$

345 Where; ρ = Correlation Coefficient

346 E' = Centered RMS difference between the observation and the prediction

347 σ_r, σ_t = Variances of the observation and the prediction respectively

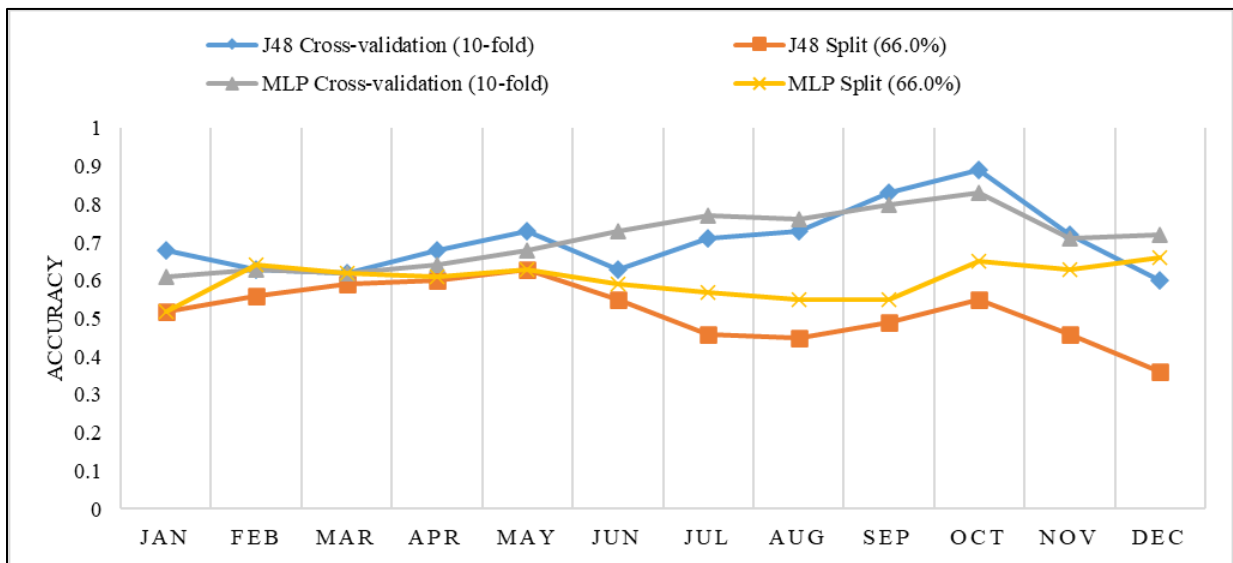
348 **Results and Discussion**

349 This section presents the analyzed data, and it constitutes the examination of how different
 350 climatic conditions influence the incidence of malaria. The comparisons between the two
 351 classifier methods were made for each of the 17 blocks in the district, and performance accuracy
 352 was evaluated month wise. The focus of this study, however, is to find the climatic influence of
 353 these incidents using some advanced machine learning techniques and ways to provide early
 354 warning about the possible future outbreaks.

355 **Comparison of MLP and J48 Results**

356 All three test options were used to assess the performance of both MLP and J48 methods, a) 10-
 357 fold cross-validation, and b) percent split (66%) and the user-supplied test sets. The following
 358 section discussed the outcome of the model prediction. An initial evaluation was completed over

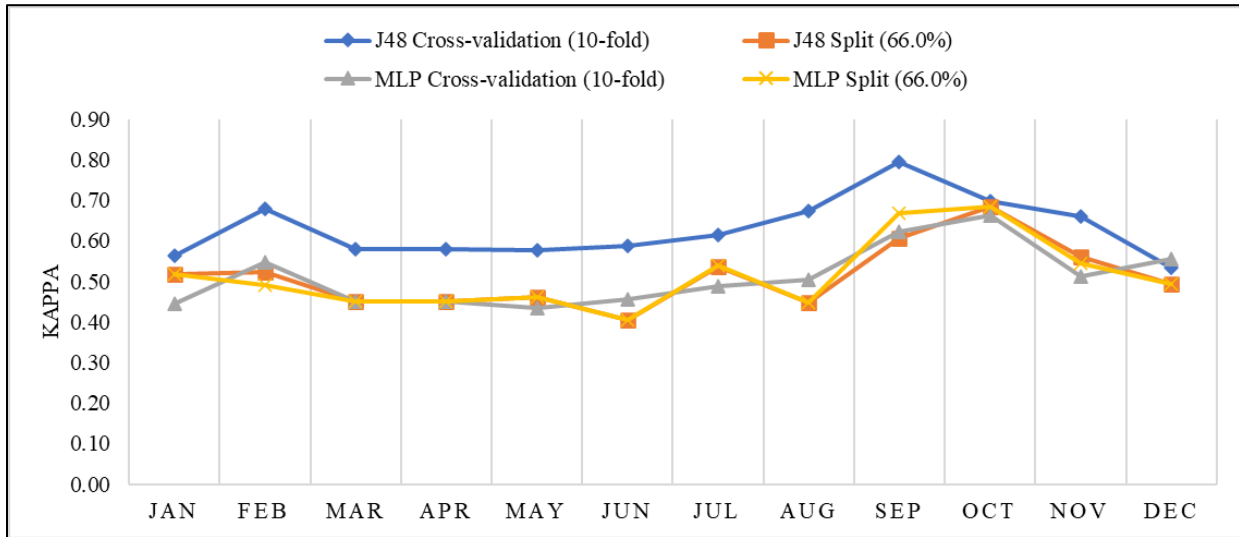
359 the Sundargarh district to investigate which of the machine learning method performed better
 360 compared to the others, before going for block-wise performance. All evaluations were done for
 361 the different months of the year to understand the prediction skill based on the monsoonal
 362 rainfall, temperature variation. While the results for prediction accuracy shows that the
 363 performance of both MLP and J48 is not very significantly different, but J48 shows better
 364 performance to MLP. In a similar study conducted by Gupta, Kumar & Sharma [53], where more
 365 attributes were analyzed and larger volumes of data used, the prediction using J48 has also
 366 turned out to be better. Besides, for both the classifiers, the 10-fold cross-validation classification
 367 testing option outperforms the percentage split (66%) method for the whole district, as shown in
 368 Figure 6 (Accuracy), Figure 7 (Kappa) and Figure 8 (RMSE), respectively.



369
 370 **Figure 7.** Month-wise Comparison of the Accuracy of the J48 (Cross-validation & Percent Split
 371 model) to the Multi-Layer Perceptron model for Accuracy for Sundargarh District

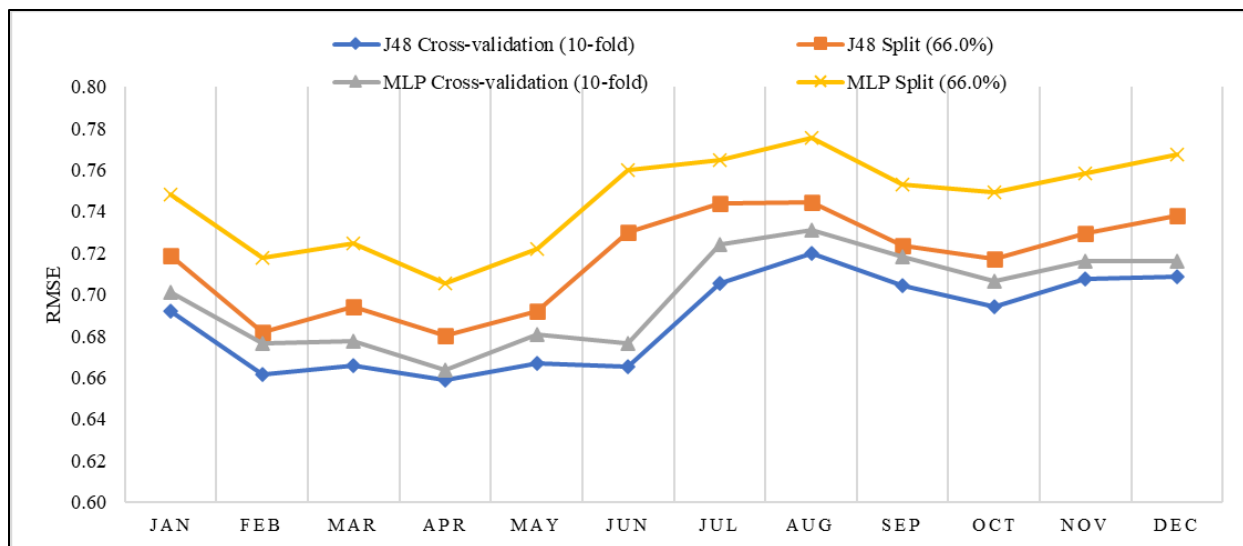
372 Figure 7 suggests that while the prediction accuracy for both the cross-validation method (J48
 373 and MLP) has improved during the mid-monsoon (July-August period) to late monsoon

374 (September-October) at the same time the prediction accuracy for the split method has declined.
 375 For all models, the dry season has less accuracy.



376
 377 **Figure 8.** Month-wise Comparison of Kappa of the J48 (Cross-validation & Percent Split model)
 378 to the Multi-Layer Perceptron model for Accuracy for Sundargarh District

379 If we consider kappa, the J48 cross-validation method has significantly better kappa in
 380 comparison to the other three methods, visibly after the monsoon onset, it shows a better
 381 agreement. While almost all methods have shown poor agreement during the drier summer
 382 period, results suggest the superiority of the J48 with cross-validation over the MLP percent
 383 split, MLP Cross-validation, and J48 Percent split (Figure 8).



384
 385 **Figure 9.** Month-wise Comparison of RMSE of the J48 (Cross-validation & Percent Split model)
 386 to the Multi-Layer Perceptron model for Accuracy for Sundargarh District

387 The prediction error (Figure 9) depicts that the cross-validation method has much less error
 388 compared to the percentage split method. The MLP percent split method consistently depicts the
 389 largest errors across seasons. Table 5 shows that the MLP has better accuracy, especially during
 390 the wet season, while less agreement to the observed condition as depicts comparatively lower
 391 kappa. For this, J48 has better performance during the wet season. Errors are more substantial for
 392 both the models during the wet period, and RMSE is low during the dry period. Since the data
 393 was analyzed monthly, J48 can be considered a more reliable predictor of malaria for the weather
 394 variables. During the monsoon and post-monsoon seasons, it has comparable RMSE and higher
 395 kappa (with highest values in September = 0.79 and October = 0.70) indicated that it performed
 396 better compared to MLP with Kappa (September = 0.62 and October = 0.66).

397 **Table 5.** Month-wise Performance metrics for RMSE, KAPPA, and Accuracy for the cross-
 398 validation Classifier

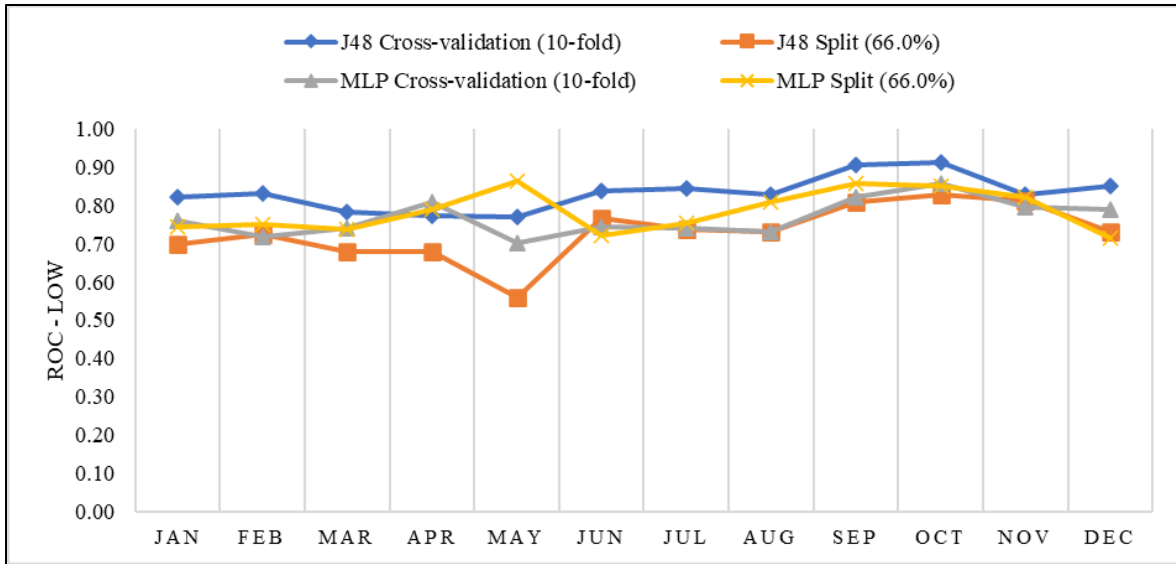
Month	Accuracy	Kappa	RMSE
-------	----------	-------	------

	J48	MLP	J48	MLP	J48	MLP
Jan	0.68	0.61	0.56	0.44	0.69	0.70
Feb	0.63	0.63	0.68	0.55	0.66	0.68
Mar	0.62	0.62	0.58	0.45	0.67	0.68
Apr	0.68	0.64	0.58	0.45	0.66	0.66
May	0.73	0.68	0.58	0.44	0.67	0.68
Jun	0.63	0.73	0.59	0.45	0.67	0.68
Jul	0.71	0.77	0.61	0.49	0.71	0.72
Aug	0.73	0.76	0.67	0.50	0.72	0.73
Sep	0.83	0.80	0.79	0.62	0.70	0.72
Oct	0.89	0.83	0.70	0.66	0.69	0.71
Nov	0.72	0.71	0.66	0.51	0.71	0.72
Dec	0.60	0.72	0.54	0.56	0.71	0.72

399 Highlighted values are for Accuracy ≥ 0.70 ; Kappa ≥ 0.60 ; RMSE ≤ 0.70

400 ROC score, as explained earlier, is generally a measure of the skill of the classifier. Evaluation
401 with ROC requires the grouping of the prediction models into three distinct prediction categories,
402 e.g., 1) High, 2) medium, and 3) Low. The evaluation shows how well the three categories of
403 events can be predicted.

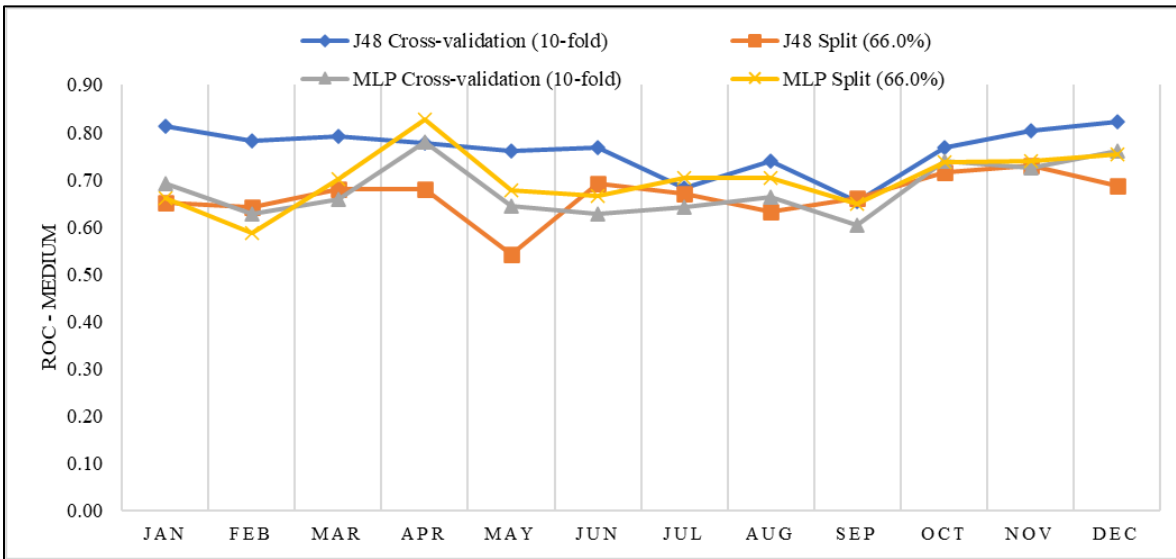
404



405

406

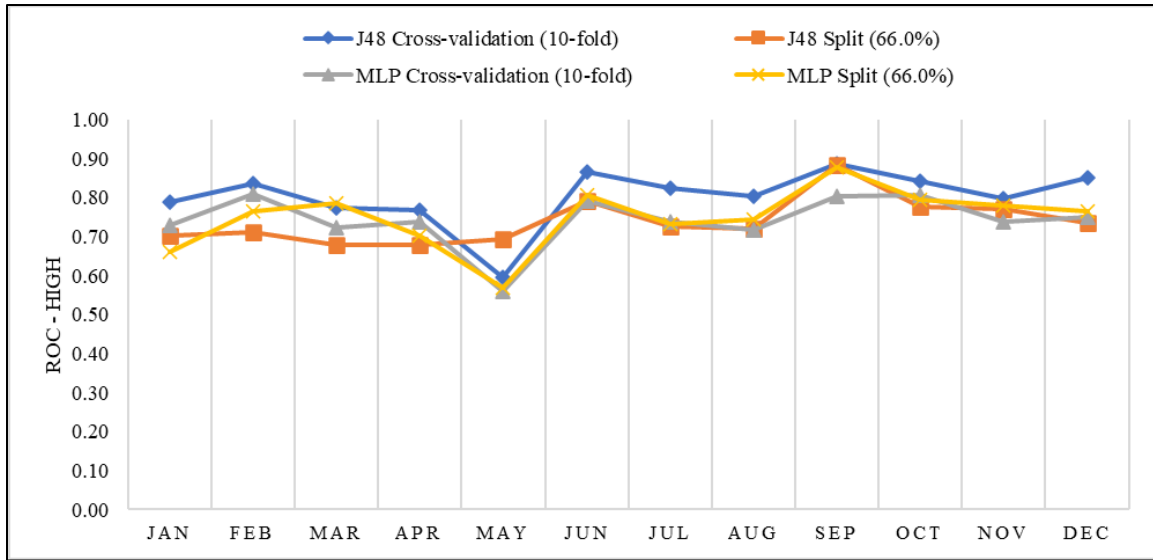
(a)



407

408

(b)



(c)

409

410

411 **Figure 10.** Month-wise Comparison of ROC for ((a) Low, (b) medium, and (c) High prediction

412 category) of the J48 (Cross-validation & Percent Split model) to the Multi-Layer Perceptron

413 model for Accuracy for Sundargarh District

414 Figure 10 shows the performance of the three different event categories, and it is evident that the

415 skill of the J48 method is comparatively better than MLP. Table 6 lists the ROC scores of all the

416 four classifiers and provides a comparative analysis of the three-event categories.

417 **Table 6.** Month-wise ROC prediction skill scores for all four classifiers for both J48 and MLP

Month	J48						MLP					
	Cross-validation			Split (66.0%)			Cross-validation			Split (66.0%)		
	High	Med	Low	High	Med	Low	High	Med	Low	High	Med	Low
Jan	0.79	0.81	0.82	0.70	0.65	0.70	0.73	0.69	0.76	0.66	0.66	0.74
Feb	0.84	0.78	0.83	0.71	0.64	0.73	0.81	0.63	0.72	0.77	0.59	0.75
Mar	0.77	0.79	0.78	0.68	0.68	0.68	0.72	0.66	0.74	0.79	0.70	0.74
Apr	0.77	0.78	0.77	0.68	0.68	0.68	0.74	0.78	0.81	0.70	0.83	0.79
May	0.60	0.76	0.77	0.69	0.54	0.56	0.56	0.65	0.70	0.57	0.68	0.86

Jun	0.87	0.77	0.84	0.79	0.69	0.77	0.79	0.63	0.74	0.81	0.67	0.72
Jul	0.82	0.68	0.84	0.73	0.67	0.74	0.74	0.64	0.74	0.73	0.70	0.75
Aug	0.80	0.74	0.83	0.72	0.63	0.73	0.72	0.66	0.73	0.74	0.70	0.81
Sep	0.89	0.65	0.91	0.88	0.66	0.81	0.80	0.60	0.82	0.88	0.65	0.86
Oct	0.84	0.77	0.91	0.78	0.72	0.83	0.81	0.74	0.86	0.80	0.74	0.85
Nov	0.80	0.80	0.83	0.77	0.73	0.81	0.74	0.72	0.80	0.78	0.74	0.82
Dec	0.85	0.82	0.85	0.74	0.69	0.73	0.75	0.76	0.79	0.76	0.75	0.72

418 Highlighted values are ROC scores ≥ 0.75

419 The J48 cross-validation method has better performance in terms of predicting the high and low
420 events across the year. During the post-monsoon season prediction skill for high events
421 (Sep=0.89, Oct=0.84, Nov=0.80, Dec=0.85) and (Sep=0.91, Oct=0.91, Nov=0.83, Dec=0.85) for
422 skill for predicting “low” events. At the same time, the percent split method has comparatively
423 less skill. This probably could be because of the fewer sample datasets used in the Percent split
424 methods for the training of the model, which might not be adequate. MLP has even poorer results
425 depicting skill for high events with ROC=0.56 for May and consistently poor throughout the
426 year. This shows that the models are generally poor during the early to the mid-monsoon period
427 (May-August), irrespective of the model classifier technique used. At the same time, the
428 prediction of the medium category event is challenging for both models as well. For J48 cross-
429 validation, the lowest value (ROC=0.65) in September and for percent split method ROC=0.54
430 in May. While in the MLP cross-validation ROC=0.60 in September and ROC=0.59 in February
431 month, respectively.

432 **Performance Evaluation of the Models at smaller Administrative Units**

433 The final exercise for the prediction model was to supply the classifiers with a user-defined set of
434 datasets for training and isolate a specific year or years for prediction at block level. So, for this

435 test, data from 2002 till 2015 was used for training and 2016 and 2017 for prediction. The results
 436 presented in Table 7 suggest the method has comparable performance to J48 and MLP. It has
 437 less RMSE and better accuracy and higher kappa values. Further investigating the performance
 438 of the Supplied test set, it was found that its accuracy of prediction is better compared to the
 439 cross-validation method, especially for central and western blocks, including Sundargarh,
 440 (Accuracy =1.0, Kappa=1.0, RMSE=0.19), Tangarpali (Accuracy =1.0, Kappa=1.0,
 441 RMSE=0.17), and Kutra (Accuracy =1.0, Kappa=1.0, RMSE=0.16). The block-wise comparison
 442 was presented in Table 7.

443 **Table 7.** Comparisons of Accuracy, RMSE, and Kappa for all blocks for J48 cross-validation
 444 and supplied set classifiers

Blocks	J48 Cross-Validation			J48 Supplied test set		
	Accuracy	Kappa	RMSE	Accuracy	Kappa	RMSE
Hemgiri	0.57	0.40	0.89	0.71	0.16	0.37
Lephripara	0.69	0.33	0.88	0.83	0.00	0.34
Tangarpali	0.94	0.55	0.76	1.00	1.00	0.17
Sundargarh	0.88	0.51	0.80	1.00	1.00	0.19
Subdega	0.71	0.60	0.86	0.71	0.26	0.36
Baragaon	0.81	0.63	0.81	0.79	0.17	0.32
Balisankara	0.55	0.31	0.89	0.79	0.00	0.40
Kutra	0.98	0.52	0.72	1.00	1.00	0.16
Rajgangpur	0.69	0.51	0.86	0.63	0.00	0.39
Kuanrmunda	0.63	0.57	0.87	0.63	0.30	0.41
Nuagaon	0.58	0.55	0.90	0.29	0.08	0.53
Bisra	0.53	0.38	0.89	0.58	0.39	0.49
Lathikata	0.72	0.50	0.86	0.42	0.13	0.50

Bonai	0.57	0.50	0.91	0.46	0.18	0.48
Lahunipara	0.70	0.60	0.86	0.54	0.47	0.48
Gurundia	0.53	0.45	0.90	0.17	0.07	0.58
Koira	0.51	0.39	0.90	0.17	0.29	0.58

445 Highlighted values are for Accuracy ≥ 0.70 ; Kappa ≥ 0.55 ; RMSE ≤ 0.50

446 ROC scores for both the classifiers of the J48 model were compared, and the results are
447 presented in Table 8. The results suggest that the model performance is satisfactory for the
448 central blocks like Kutra (High=0.99, Med=0.80, low=0.90), Subdega (High=0.91, Med=0.76,
449 Low=0.86), Rajgampur (High=0.82, Med=76, Low=84). While blocks such as Bonai
450 (High=0.76, Med=0.63, Low=0.78), Koira (High=0.77, Med=0.75, Low=0.72), Gurundia
451 (High=0.89, Med=0.71, Low=0.76), and Balisankara (High=0.73, Med=0.69, Low=0.74),
452 depicts considerably lower accuracy and prediction skills. While blocks with plain land and
453 forest cover performed much better compared to the highly elevated regions.

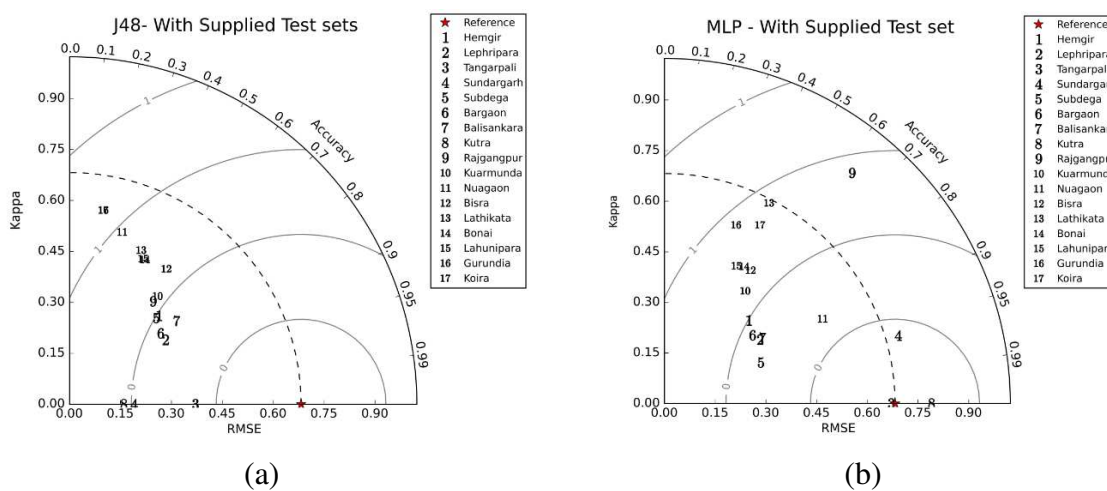
454 **Table 8.** Comparative analysis of ROC scores for all blocks in the District for J48 cross-
455 validation and supplied set classifiers

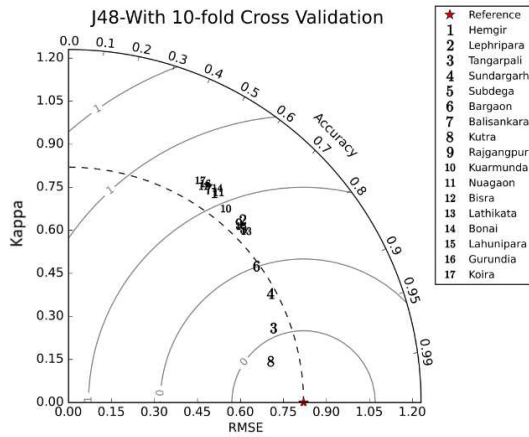
Blocks	Cross-Validation			Supplied Test		
	High	Med	Low	High	Med	Low
Hemgiri	0.83	0.73	0.72	0.76	0.74	0.94
Lephripara	0.75	0.76	0.72	0.79	0.61	0.68
Tangarpali	0.69	0.75	0.75	0.81	0.79	0.85
Sundargarh	0.62	0.74	0.70	0.45	0.69	0.81
Subdega	0.91	0.76	0.86	0.98	0.84	0.93
Baragaon	0.93	0.77	0.90	0.73	0.76	0.60
Balisankara	0.73	0.69	0.74	0.12	0.34	0.64
Kutra	0.99	0.80	0.90	0.86	0.94	0.88

Rajgangpur	0.82	0.76	0.84	0.92	0.89	0.77
Kuanrunda	0.97	0.71	0.84	0.93	0.65	0.83
Nuagaon	0.81	0.67	0.81	0.60	0.71	0.53
Bisra	0.88	0.70	0.84	0.49	0.79	0.73
Lathikata	0.88	0.89	0.66	0.65	0.87	0.43
Bonai	0.76	0.63	0.78	0.46	0.82	0.59
Lahunipara	0.90	0.76	0.94	0.31	0.39	0.39
Gurundia	0.89	0.71	0.76	0.53	0.20	0.33
Koira	0.77	0.75	0.72	0.25	0.21	0.34

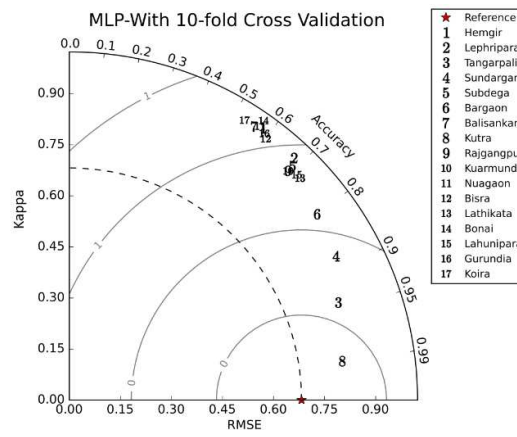
456 Highlighted values are ROC scores ≥ 0.75

457 Figure 11 demonstrates the comparison of the three test options used in both classifier techniques
 458 J48 and MLP using the Taylor diagram. All three indicators (Accuracy, Kappa, and RMSE) were
 459 represented in three different axes. RMSE in X-axis, Kappa in Y-axis, and accuracy in the arc,
 460 respectively.

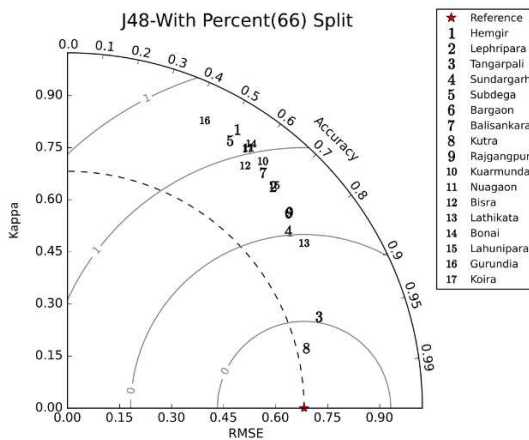




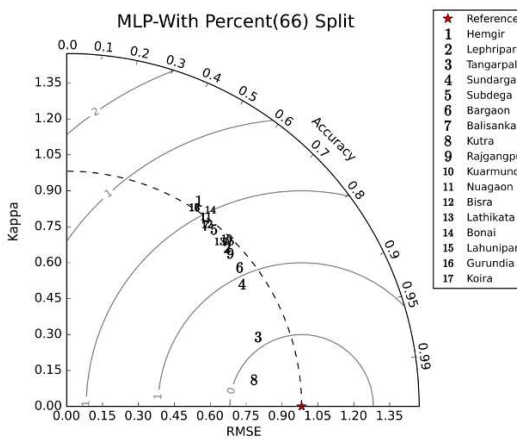
(c)



(d)



(e)



(f)

461 **Figure 11.** Block-wise accuracy (arc), RMSE (x-axis) and Kappa (Y-axis) for the J48 and the
 462 Multi-Layer Perceptron model with a Supplied test set, 10-fold Cross-validation & Percent Split
 463 classifiers

464 While the numbers 1 to 17 represent each block of the district, its position in the plotting space
 465 determines its corresponding Accuracy, RMSE, and Kappa. Block-wise analysis with all the
 466 results suggests that with the 10-fold cross-validation and the supplied test set option has yielded
 467 promising results in comparison to the percent split and supplied test options. Especially blocks
 468 (8=Kutra, 3=Tangarpali, 4=Sundargarh) from the central to western plain have better

469 performance to the blocks with varying topography (17=Koirā, 1=Hemgiri, 14=Bonai). The
 470 supplied test options depicted smaller RMSEs but also have inconsistency with Accuracy and
 471 Kappa (either Accuracy=1.0 or very low), making it unreliable for use in predictions. So, it is
 472 evident that flat terrains with lower variability in the rainfall, temperature, and humidity provides
 473 reliable performance than to the regions with higher variability.

474 **Month-wise Nominal Relationship**

475 Establishing a conventional relationship between the monthly and seasonal variation of the
 476 climatic parameters to the incidents using the nominal range was beneficial to the evaluation
 477 process. The nominal range derived from the continuous numeric data range [54] is represented
 478 in Table 9. Based on the spectrum, Table 10 provides a summary of the relationship between the
 479 nominal rainfall, temperature, and humidity range to that of the malaria incidents.

480 **Table 9.** Climate and Incident data range Evaluation

Range	RF	T2max	RH	Incidents
Low	0 - 23	28.3 - 30.3	68.3 - 82.6	35 - 78
Medium	23.1 - 178	30.4 - 33.6	82.7 - 93	78.1 - 173
High	178.1 - 445	33.7 - 38.3	93.1 - 96.7	173.1 - 460
Very High	> 445	> 38.3	> 96.7	> 460

481

482 **Table 10.** Month-wise nominal relationship between incident data and climate data

Month	RF	T _{2max}	RH	Incidents
January	Low	Low	Low	Low
February	Low	Low	Low	Low
March	Low	Medium	Low	Medium
April	Low	Very High	Low	Medium

May	Medium	Very High	Medium	Medium
June	High	High	Medium	Medium
July	High	Medium	High	High
August	Very High	Medium	High	High
September	Very High	Low	High	Medium
October	Medium	Low	Medium	Medium
November	Low	Low	Medium	High
December	Low	Low	Low	High

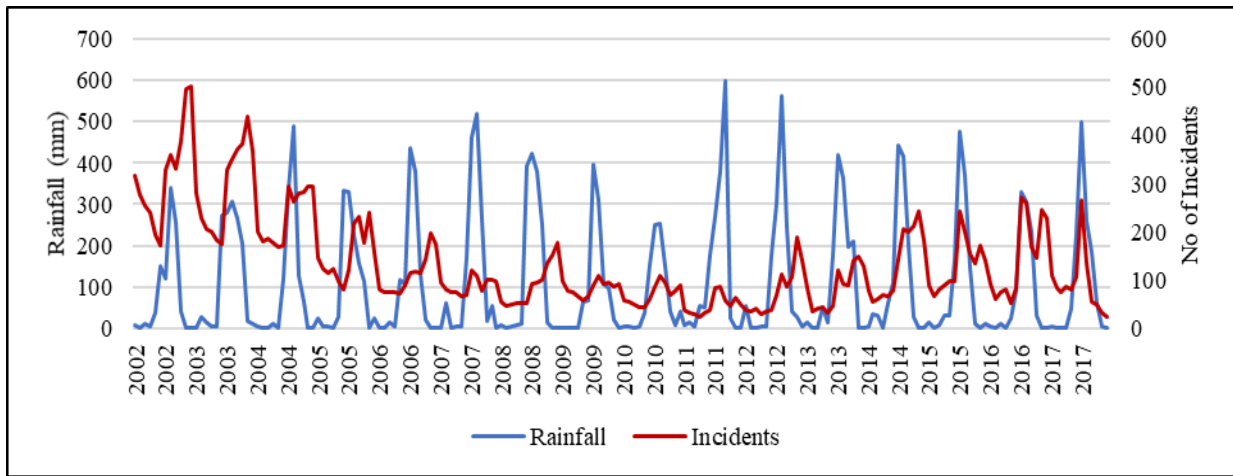
483

484 As shown in the table, periods of low temperature (28°-30°), low rainfall (0 -23mm), and low
485 relative humidity (68-82%) during the drier and cooler months of January and February are
486 characterized by lower cases of malaria. During the months of March-April-May, a period of low
487 precipitation, medium to higher temperature with lower relative humidity, there is an increase in
488 the number of incidences of malaria. This also agrees with the findings of Lee et al. [55], a study
489 conducted in the humid Arunachal Pradesh, India, that suggests decreasing precipitation and
490 increasing temperature resulted in increasing malaria incidence. With the arrival of the monsoon
491 and during the June-July-August-September, the period of high to very high rainfall, higher
492 temperature, and medium to high relative humidity, the malaria incidents were further on the
493 rise. Surprisingly, malaria incidents climbed even after the withdrawal of the monsoon, fall in the
494 temperature and humidity significantly during the November-December winter period. So, there
495 is possibly a lag-effect of the climatic phenomenon on the incidents.

496 Based on the results, it can be concluded that relative humidity and temperature showed a strong
497 association with malaria incidence, which is consistent with the study by Srimath-Tirumula [56],
498 in Vishakhapatnam, in India, which experiences similar climate compared to Sundargarh. In

499 contrast, rainfall showed a relatively weaker association, which is in line with the study by
500 Bomblies [3], which argues that during the rainy season, the breeding habitats of mosquitoes are
501 flushed away temporarily. Still, they start breeding again when the rains stop, and water becomes
502 stagnant, and the environmental condition is conducive for breeding.

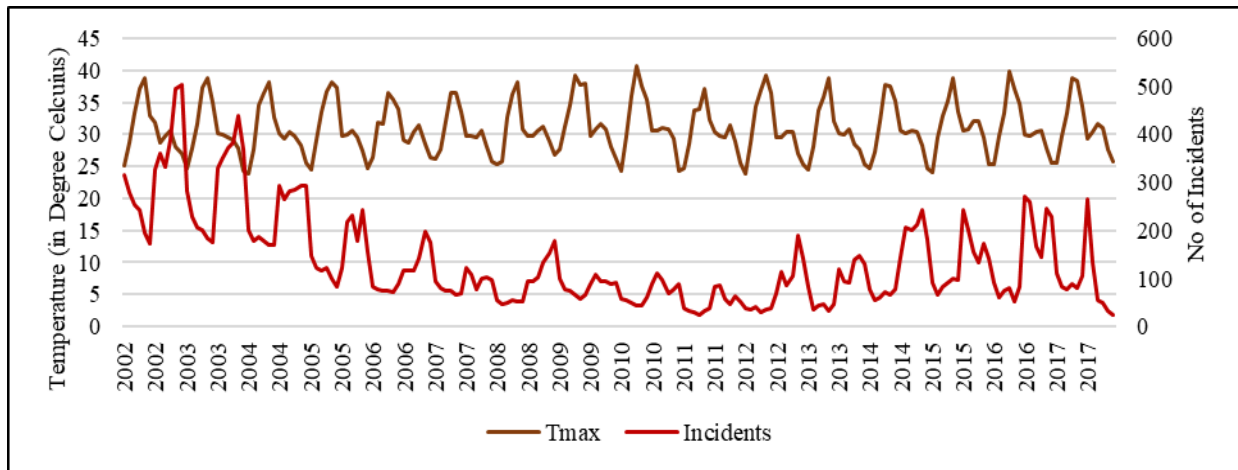
503 This study found that extremely high temperature is one of the crucial triggers of the higher
504 number of malaria incidents in Sundargarh district. Therefore, this agrees with the argument of
505 Smith et al. [57], when the temperatures increase, it reduces the time taken by the mosquito
506 parasite to complete its development. Furthermore, relative humidity also affects the
507 transmission of the malaria vector in agreement with [58], found out that mosquitoes survive
508 better under high humidity conditions. During high humid seasons, the number of malaria
509 incidents increases compared to the less humid conditions. The time series plots (Figure 12) for
510 the rainfall (a), temperature (b), and humidity (c) submit its direct association with the incidents
511 reported.



512

513

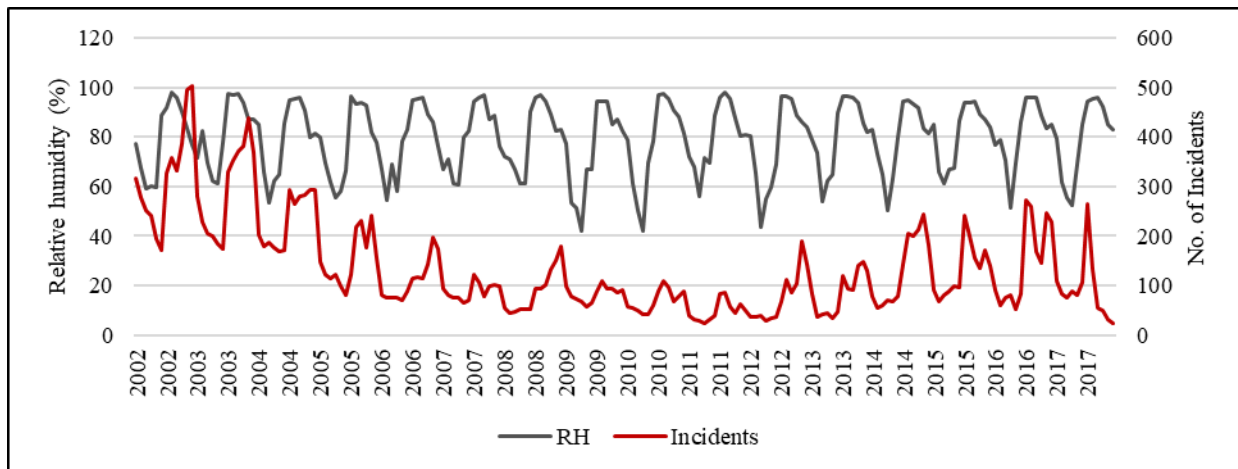
(a)



514

515

(b)



516

517

(c)

518 **Figure 12.** Time series plot for Monthly Rainfall (a), Maximum Temperature (b), and Relative

519 Humidity (c) in comparison with malaria incidents for the period of 2002-2017

520 It is concluded that relative humidity and temperature showed a significant relationship to

521 malaria incidences in the district, especially for some blocks, those are in flat terrain and near

522 dense vegetation like the forest. Additionally, rainfall also affects the transmission of malaria

523 vector incidence. However, the rate of vector transmission during the rainy season is relatively

524 lower, suggesting the influence of rain on malaria incidents may happen in a time lag mode as

525 well.

526 **Discussions**

527 A malaria early warning system and risk mapping tool is necessary to provide adequate support
528 to the public health workers to take preparedness measures and remain prepared for any possible
529 outbreaks in near real-time [59-61]. Several such attempts were proposed, such as a spatial
530 decision support system for Karnataka [62] or the operational system in Kenya [60], but not all
531 are very successful and useful. Statistical regression-based analysis like the multiple polynomial
532 regression, semiparametric Poisson distribution methods, distributed non-linear lag model,
533 hierarchical Bayesian framework and the time series regression models [14-19, 20-22, 56], or use
534 of ANN or machine learning as discussed by researchers, provides an opportunity to process the
535 dataset and establish an association between climate and the malaria incidents. Considering the
536 climate dimension only in malaria early warning is not adequate and requires a deep
537 understanding of the influence of all other facets, including climate in the establishment of an
538 effective and operational warning system.

539 **Model Selection and Model Evaluation**

540 Appropriate selection of model, algorithm, and model evaluation techniques are vital in machine
541 learning. The evaluation intends to estimate the performance of a model or algorithm on future
542 data. Running a learning algorithm over a training dataset with different hyperparameter settings
543 will result in different models [63]. Since we are typically interested in selecting the best-
544 performing model, estimation helps in choosing the best model to fit the purpose though, the
545 estimation of the absolute performance of a model is one of the most challenging tasks in
546 machine learning [63]. Working with small sample sizes in machine learning is acceptable but
547 choosing the correct sampling method is vital [64, 65]. Considering the sample size in this study
548 is smaller, and for parameter optimization, 10-fold cross-validation and Leave-One-Out cross-

549 validation are recommended as the best sampling mechanisms and generally would yield better
550 results [65].

551 Assessments made by researchers to evaluate different classifier performances [66]
552 recommended the use of leave-one-out cross-validation method as a preferred method of
553 prediction. Thus, this study put an effort to assess the 10-fold cross-validation method, which
554 incidentally also performed well. With more in-depth analysis, the reason would be linked to the
555 data sampling and training strategy implemented by the method in comparison to the others. A
556 list of explanations is provided below, which considers the mechanism with which the cross-
557 validation method works.

558 a) Utilized all the data samples for training and test and takes care of the multi-class issue
559 that arises in the percentage split method, where the sample sizes are static, and
560 generating multiple classes means a reduction in test sets.

561 b) We defined more metrics for the learning algorithm than other methods. If we have, n
562 samples there can $n-1$ models to predict one instance of the predictand.

563 c) Through model stacking and back-propagation models are processed in a pipeline
564 allowing model prediction by learning from the previous model in the forward direction
565 and feedback and model training in the backward direction. The model bias (error) is also
566 handled better in this process.

567 d) Finally, parameter fine-tuning, a process by which the parameters were tuned with an
568 independent validation set, that suggested the ideal number of trees in a classifier, hidden
569 layer size (activation function) in the Neural network.

570 The probable explanation for the model not performing well for some of the blocks and months
571 could be a factor that is external to the climate influence. Furthermore, the interventions in place

572 in parts of the district and other regions [67] and the socio-economic status of the significant
573 population in the eastern belt, being tribe and access to necessary facilities is limited as also
574 explained by Sundararajan et al. [13]. These factors influence the increase or decrease in the
575 cases, and not truly reflect the direct influence of climate. The other reflection from the analysis
576 is that the model's performance dips down significantly, especially during July, August. It picks
577 up during September October and then again shows lower accuracy during November and
578 December months (this is depicted in Figure 7). The poor performance of the model could be
579 because of the varying topography, which affects the intra-seasonal rainfall variability as well as
580 the spatial variation of the temperature and humidity could influence the results strongly. For
581 example, blocks such as Bonai, Koira, Gurundia, and Hemgiri are with higher elevation and
582 depict considerably lower accuracy for the prediction. In comparison, blocks with plain land and
583 forest cover had better performance.

584 **Conclusion**

585 The climatic condition of Odisha, especially the Sundargarh district, makes it vulnerable to
586 malaria y. Monsoon rainfall, maximum surface temperature ranging from 27° to 40° Celsius
587 during the summer, and relative humidity in the range of 60% to 85% provide a more favorable
588 climatic condition for the breeding of the malaria larva during the monsoon and post-monsoon
589 period. It was found that the increase in malaria incidents is significantly attributed to climatic
590 factors such as temperature, humidity, and monthly rainfall variability. Among the two classifier
591 models used, J48 has shown comparatively better skills over the MLP. J48 demonstrated less
592 error (RMSE = 0.6), better kappa = 0.63, and higher accuracy = 0.71), suggesting it to be a
593 suitable model.

594 J48 model has provided a greater insight into the predictability of malaria when compared in
595 seasonal scale. During the pre-monsoon (Mar-Apr-May) period it has accuracy = 0.68, kappa =
596 0.58 and RMSE = 0.67, for monsoon (Jun-July-Aug-Sep) period with accuracy = 0.73, kappa =
597 0.67 and RMSE = 0.70, post-monsoon (Oct-Nov) period with highest accuracy = 0.81, kappa =
598 0.68 and RMSE = 0.70 and during the winter (Dec-Jan-Feb) period with lowest accuracy = 0.64,
599 kappa = 0.59 and RMSE = 0.69, respectively. This suggests that the model performance is
600 particularly good during the monsoon and post-monsoon when the malaria incidents are at the
601 peak. Besides, non-climatic factors play a significant role in the malaria spreading, which was
602 reflected with a lower accuracy. However, climate being an extremely complex and variable
603 factor to predict, the results provided promising signal for the prediction of future malaria
604 incidents. Therefore, it is recommended that the public health department could adopt the J48
605 classifier machine learning technique in the malaria early warning system for the early detection
606 of malaria.

607 Even though the models have shown better performance in terms of predicting malaria incidence,
608 it is constrained by the non-availability of datasets for an elongated period. More finer scale
609 datasets (both climate and malaria cases) would have provided an opportunity for deeper analysis
610 to understand the phases and lags within a month as well. Furthermore, non-climatic factors such
611 as the demography, immunity within the population, the socio-economic structure of society,
612 availability of affordable public health facilities, and other environmental modifications
613 initiatives are strongly recommended to be factored in, while developing a malarial early
614 warning system.

615 **Declarations**

616 **Ethics approval and consent to participate**

617 Not applicable.

618 **Consent for publication**

619 Not applicable.

620 **Availability of data and materials**

621 The datasets generated and/or analyzed during the current study are not publicly available as they
622 are collected from government sources and due to the sensitivity but are available from the
623 corresponding author on reasonable request.

624 **Competing interests**

625 The authors declare that they have no competing interests.

626 **Funding**

627 Not applicable.

628 **Author Contributions**

629 PM and NKT conceptualized the study and drafted the manuscript. PM prepared the data set for
630 analysis. NKT contributed to the study design. NKT, IP, and SS provided critical input and
631 improvising the manuscript. All authors have read and agreed to the published version of the
632 manuscript.

633 **Acknowledgments**

634 The authors would like to acknowledge the contribution of the Directorate of Public Health,
635 Government of Odisha, for providing the malaria incident datasets. Also, the ECMWF for
636 keeping the climate datasets free and open for public access.

637 **Conflict of interest disclosure**

638 The authors declare that there is no conflict of interest regarding the publication of this paper.

639 **References**

- 640 1. Organization WH. World malaria report 2018. 2018. World Health Organization:
641 Geneva. 2020.
- 642 2. Kovats RS, Bouma MJ, Hajat S, Worrall E, Haines A. El Niño and health. *The Lancet*.
643 2003; 362(9394):1481-9.
- 644 3. Bomblies A. Modeling the role of rainfall patterns in seasonal malaria transmission.
645 *Climatic change*. 2012;112(3-4):673-85.
- 646 4. Githeko AK, Lindsay SW, Confalonieri UE, Patz JA. Climate change and vector-borne
647 diseases: a regional analysis. *Bulletin of the World Health Organization*. 2000;78:1136-
648 47.
- 649 5. Leal Filho W, Bönecke J, Spielmann H, Azeiteiro UM, Alves F, de Carvalho ML, et al.
650 Climate change and health: An analysis of causal relations on the spread of vector-borne
651 diseases in Brazil. *Journal of Cleaner Production*. 2018;177:589-96.
- 652 6. Van Lieshout M, Kovats R, Livermore M, Martens P. Climate change and malaria:
653 analysis of the SRES climate and socio-economic scenarios. *Global environmental*
654 *change*. 2004;14(1):87-99.
- 655 7. Holt RA, Subramanian GM, Halpern A, Sutton GG, Charlab R, Nusskern DR, et al. The
656 genome sequence of the malaria mosquito *Anopheles gambiae*. *science*.
657 2002;298(5591):129-49.
- 658 8. Kakmeni FMM, Guimapi RY, Ndjomatchoua FT, Pedro SA, Mutunga J, Tonnang HE.
659 Spatial panorama of malaria prevalence in Africa under climate change and interventions
660 scenarios. *International Journal of Health Geographics*. 2018;17(1):1-13.

- 661 9. Programme NVBDC. Malaria situation in India. Ministry of Health & Family Welfare,
662 Govt. of India: Directorate General of Health Services; [Available from:
663 <https://nvbdcp.gov.in/index4.php?lang=1&level=0&linkid=564&lid=3867.10>. Accessed
664 28 March 2020.
- 665 10. Organization WH. World malaria report 2019. World Health Organization: Geneva.
666 2019.
- 667 11. Mahakur PK, Nayak NC. Intrastate Income Inequalities in Odisha: Examining
668 Decomposition by Regions and Broad Sectors. Odisha Economy Discussion Series.
669 2019;1.
- 670 12. Kannan KP. Interrogating inclusive growth: Poverty and inequality in India. Routledge;
671 2017 Sep 19
- 672 13. Sundararajan R, Kalkonde Y, Gokhale C, Greenough PG, Bang A. Barriers to malaria
673 control among marginalized tribal communities: a qualitative study. PLoS One.
674 2013;8(12):e81966.
- 675 14. Chatterjee C, Sarkar RR. Multi-step polynomial regression method to model and forecast
676 malaria incidence. PLoS One. 2009;4(3):e4726.
- 677 15. Shimaponda-Mataa NM, Tembo-Mwase E, Gebreslasie M, Achia TN, Mukaratirwa S.
678 Modelling the influence of temperature and rainfall on malaria incidence in four endemic
679 provinces of Zambia using semiparametric Poisson regression. Acta tropica.
680 2017;166:81-91.
- 681 16. Guo C, Yang L, Ou C-Q, Li L, Zhuang Y, Yang J, et al. Malaria incidence from 2005–
682 2013 and its associations with meteorological factors in Guangdong, China. Malaria
683 Journal. 2015;14(1):116.

- 684 17. Arab A, Jackson MC, Kongoli C. Modelling the effects of weather and climate on
685 malaria distributions in West Africa. *Malaria Journal*. 2014;13(1):126.
- 686 18. Imai C, Armstrong B, Chalabi Z, Mangtani P, Hashizume M. Time series regression
687 model for infectious disease and weather. *Environmental Research*. 2015;142:319-27.
- 688 19. Rejeki DSS, Nurhayati N, Budi A, Murhandarwati EEH, KUSNANTO H. A Time Series
689 Analysis: Weather Factors, Human Migration and Malaria Cases in Endemic Area of
690 Purworejo, Indonesia, 2005–2014. *Iranian Journal of Public Health*. 2018;47(4):499.
- 691 20. Neter J, Kutner MH, Nachtsheim C, Wasserman W. *Applied linear statistical models*.
692 1996. WCB McGraw-Hill. 1996.
- 693 21. Basheer IA, Hajmeer M. Artificial neural networks: fundamentals, computing, design,
694 and application. *Journal of Microbiological Methods*. 2000;43(1):3-31.
- 695 22. Yao J, Tan CL, Poh H-L. Neural networks for technical analysis: a study on KLCI.
696 *International Journal of Theoretical and Applied Finance*. 1999;2(02):221-41.
- 697 23. OPERATIONS DOC, ODISHA. *District Census Handbook Sundargarh*. 2011; SERIES-
698 22. [Available from:
699 https://censusindia.gov.in/2011census/dchb/2105_PART_B_DCHB_SUNDARGARH.pdf
700 f]. Accessed 17 April 2020.
- 701 24. Servadio JL, Rosenthal SR, Carlson L, Bauer C. Climate patterns and mosquito-borne
702 disease outbreaks in South and Southeast Asia. *Journal of Infection and Public Health*.
703 2018;11(4):566-71.
- 704 25. Jolivet R, Grandin R, Lasserre C, Doin MP, Peltzer G. Systematic InSAR tropospheric
705 phase delay corrections from global meteorological reanalysis data. *Geophysical*
706 *Research Letters*. 2011;38(17).

- 707 26. Marcos R, González-Reviriego N, Torralba V, Soret A, Doblas-Reyes FJ.
708 Characterization of the near surface wind speed distribution at global scale: ERA-Interim
709 reanalysis and ECMWF seasonal forecasting system 4. *Climate Dynamics*. 2019;52(5-
710 6):3307-19.
- 711 27. Parker WS. Reanalyses and observations: What's the difference? *Bulletin of the*
712 *American Meteorological Society*. 2016;97(9):1565-72.
- 713 28. Bengtsson L, Hagemann S, Hodges KI. Can climate trends be calculated from reanalysis
714 data? *Journal of Geophysical Research: Atmospheres*. 2004;109(D11).
- 715 29. Holmes G, Donkin A, Witten IH. Weka: A machine learning workbench. In *Proceedings*
716 *of ANZIIS'94-Australian New Zealand Intelligent Information Systems Conference 1994*
717 *Nov 29 (pp. 357-361)*. IEEE.
- 718 30. Witten IH, Frank E. *Data mining: practical machine learning tools and techniques with*
719 *Java implementations*. *Acm Sigmod Record*. 2002 Mar 1;31(1):76-7.
- 720 31. Hornick M, inventor; Oracle International Corp, assignee. *Data mining agents for*
721 *efficient hardware utilization*. United States patent US 7,620,635. 2009 Nov 17.
- 722 32. Bui QT, Nguyen QH, Pham VM, Pham MH, Tran AT. Understanding spatial variations
723 of malaria in Vietnam using remotely sensed data integrated into GIS and machine
724 learning classifiers. *Geocarto International*. 2019 Oct 15;34(12):1300-14.
- 725 33. Sharma V, Kumar A, Lakshmi Panat D, Karajkhede G. Malaria outbreak prediction
726 model using machine learning. *International Journal of Advanced Research in Computer*
727 *Engineering & Technology (IJARCET)*. 2015 Dec;4(12).

- 728 34. Olayinka TC, Chiemeké SC. Predicting Paediatric Malaria Occurrence Using
729 Classification Algorithm in Data Mining. Journal of Advances in Mathematics and
730 Computer Science. 2019 Apr 1:1-0.
- 731 35. Mello-Román JD, Mello-Román JC, Gomez-Guerrero S, García-Torres M. Predictive
732 Models for the Medical Diagnosis of Dengue: A Case Study in Paraguay. Computational
733 and mathematical Methods in Medicine. 2019 Jul 29;2019.
- 734 36. Guo P, Liu T, Zhang Q, Wang L, Xiao J, Zhang Q, Luo G, Li Z, He J, Zhang Y, Ma W.
735 Developing a dengue forecast model using machine learning: A case study in China.
736 PLoS neglected tropical diseases. 2017 Oct 16;11(10):e0005973.
- 737 37. Atulbhai DK. COMPARISON AND COMBINATION OF MINING TECHNIQUES
738 FOR GENE ANALYSIS TO IDENTIFY DENGUE.
- 739 38. Shakil KA, Anis S, Alam M. Dengue disease prediction using weka data mining tool.
740 arXiv preprint arXiv:1502.05167. 2015 Feb 18.
- 741 39. Leo J, Luhanga E, Michael K. Machine Learning Model for Imbalanced Cholera Dataset
742 in Tanzania. The Scientific World Journal. 2019 Aug 1;2019.
- 743 40. Mahmud SH, Hossin MA, Ahmed MR, Noori SR, Sarkar MN. Machine Learning Based
744 Unified Framework for Diabetes Prediction. InProceedings of the 2018 International
745 Conference on Big Data Engineering and Technology 2018 Aug 25 (pp. 46-50).
- 746 41. Zia UA, Khan N. Predicting diabetes in medical datasets using machine learning
747 techniques. International Journal of Scientific & Engineering Research Volume. 2017
748 May;8(5).
- 749 42. Al Jarullah AA, editor Decision tree discovery for the diagnosis of type II diabetes. 2011
750 International Conference on Innovations in Information Technology; 2011: IEEE.

- 751 43. Sabarinathan V, Sugumaran V. Diagnosis of heart disease using decision tree.
752 International Journal of Research in Computer Applications & Information Technology.
753 2014;2(6):74-9.
- 754 44. Dangare CS, Apte SS. Improved study of heart disease prediction system using data
755 mining classification techniques. International Journal of Computer Applications.
756 2012;47(10):44-8.
- 757 45. Korting TS. C4.5 algorithm and multivariate decision trees. Image Processing Division,
758 National Institute for Space Research–INPE Sao Jose dos Campos–SP, Brazil. 2006.
- 759 46. Nicholson C. A Beginner's Guide to Multilayer Perceptrons (MLP) [Available from:
760 [https://pathmind.com/wiki/multilayer-](https://pathmind.com/wiki/multilayer-perceptron#three.%20https://pathmind.com/wiki/multilayer-perceptron#three)
761 [%20https://pathmind.com/wiki/multilayer-perceptron#three](https://pathmind.com/wiki/multilayer-perceptron#three)]. Accessed
762 20 April 2020.
- 763 47. Quinlan JR. C4.5: programs for machine learning; Elsevier; 2014.
- 764 48. Kaur G, Chhabra A. Improved J48 classification algorithm for the prediction of diabetes.
765 International Journal of Computer Applications. 2014;98(22).
- 766 49. McHugh ML. Interrater reliability: the kappa statistic. Biochemia medica: Biochemia
767 medica. 2012;22(3):276-82.
- 768 50. Landis JR, Koch GG. An application of hierarchical kappa-type statistics in the
769 assessment of majority agreement among multiple observers. Biometrics. 1977:363-74.
- 770 51. Kumar N, Khatri S, editors. Implementing WEKA for medical data classification and
771 early disease prediction. 2017 3rd International Conference on Computational
772 Intelligence & Communication Technology (CICT); 2017: IEEE.

- 773 52. Taylor KE. Summarizing multiple aspects of model performance in a single diagram.
774 Journal of Geophysical Research: Atmospheres. 2001;106(D7):7183-92.
- 775 53. Gupta S, Kumar D, Sharma A. Performance analysis of various data mining classification
776 techniques on healthcare data. International Journal of Computer Science & Information
777 Technology (IJCSIT). 2011 Aug;3(4):155-69.
- 778 54. Mangiafico S. Summary and analysis of extension program evaluation in R, version 1.15.
779 0. Rutgers Cooperative Extension, New Brunswick, NJ: [https://rcompanion](https://rcompanion.org/handbook/)
780 [org/handbook/](https://rcompanion.org/handbook/)[Google Scholar]. 2016.
- 781 55. Lee E, Burkhart J, Olson S, Billings AA, Patz JA, Harner EJ. Relationships of climate
782 and irrigation factors with malaria parasite incidences in two climatically dissimilar
783 regions in India. Journal of Arid Environments. 2016 Jan 1;124:214-24.
- 784 56. Srimath-Tirumula-Peddinti RCPK, Neelapu NRR, Sidagam N. Association of climatic
785 variability, vector population and malarial disease in district of Visakhapatnam, India: a
786 modeling and prediction analysis. PLoS One. 2015;10(6):e0128377.
- 787 57. Smith DL, Perkins TA, Tusting LS, Scott TW, Lindsay SW. Mosquito population
788 regulation and larval source management in heterogeneous environments. PloS One.
789 2013;8(8):e71247.
- 790 58. Goswami S, Saxena A, Singh KJ, Chandra S, Cleal CJ. An appraisal of the Permian
791 palaeobiodiversity and geology of the Ib-River Basin, eastern coastal area, India. Journal
792 of Asian Earth Sciences. 2018;157:283-301.
- 793 59. Connor SJ, Thomson MC, Flasse SP, Perryman AH. Environmental information systems
794 in malaria risk mapping and epidemic forecasting. Disasters. 1998;22(1):39-56.

- 795 60. Thomson M, Indeje M, Connor S, Dilley M, Ward N. Malaria early warning in Kenya
796 and seasonal climate forecasts. *The Lancet*. 2003;362(9383):580.
- 797 61. Thomson MC, Connor SJ. The development of malaria early warning systems for Africa.
798 *Trends in parasitology*. 2001;17(9):438-45.
- 799 62. Shekhar S, Yoo E, Ahmed S, Haining R, Kadannolly S. Analysing malaria incidence at
800 the small area level for developing a spatial decision support system: a case study in
801 Kalaburagi, Karnataka, India. *Spatial and spatio-temporal epidemiology*. 2017;20:9-25.
- 802 63. Raschka S. Model evaluation, model selection, and algorithm selection in machine
803 learning. arXiv preprint arXiv:181112808. 2018.
- 804 64. Lusa L. Joint use of over-and under-sampling techniques and cross-validation for the
805 development and assessment of prediction models. *BMC Bioinformatics*. 2015;16(1):1-
806 10.
- 807 65. Cawley GC, Talbot NL. On over-fitting in model selection and subsequent selection bias
808 in performance evaluation. *The Journal of Machine Learning Research*. 2010;11:2079-
809 107.
- 810 66. Bischl B, Mersmann O, Trautmann H, Weihs C. Resampling methods for meta-model
811 validation with recommendations for evolutionary computation. *Evolutionary
812 Computation*. 2012 Jun;20(2):249-75.
- 813 67. Sahu S, Gunasekaran K, Raju H, Vanamail P, Pradhan M, Jambulingam P. Response of
814 malaria vectors to conventional insecticides in the southern districts of Odisha State,
815 India. *The Indian Journal of Medical Research*. 2014;139(2):294.

Figures

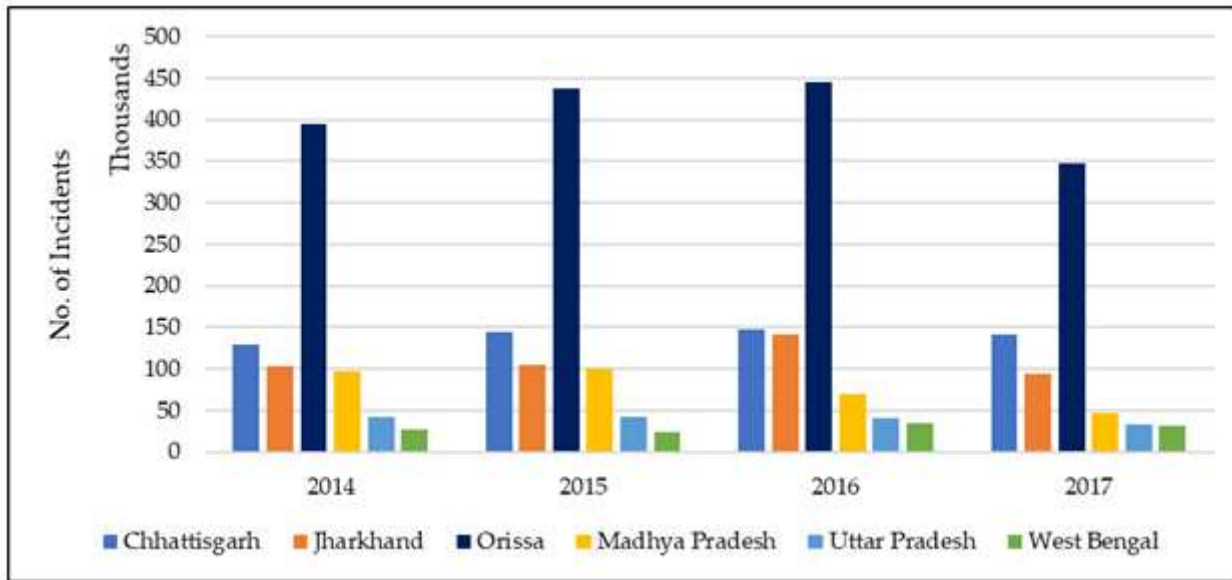


Figure 1

Malaria incidents across six different states of India for the period of 2014-2017.

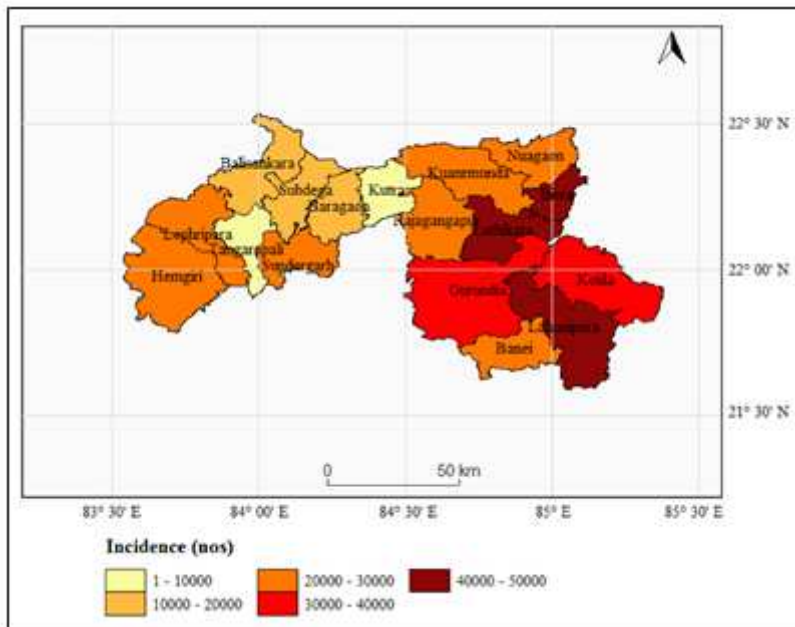


Figure 2

Cumulative malaria incidence map for the period 2002 to 2017 (Data Source: Directorate of Public Health, Odisha)

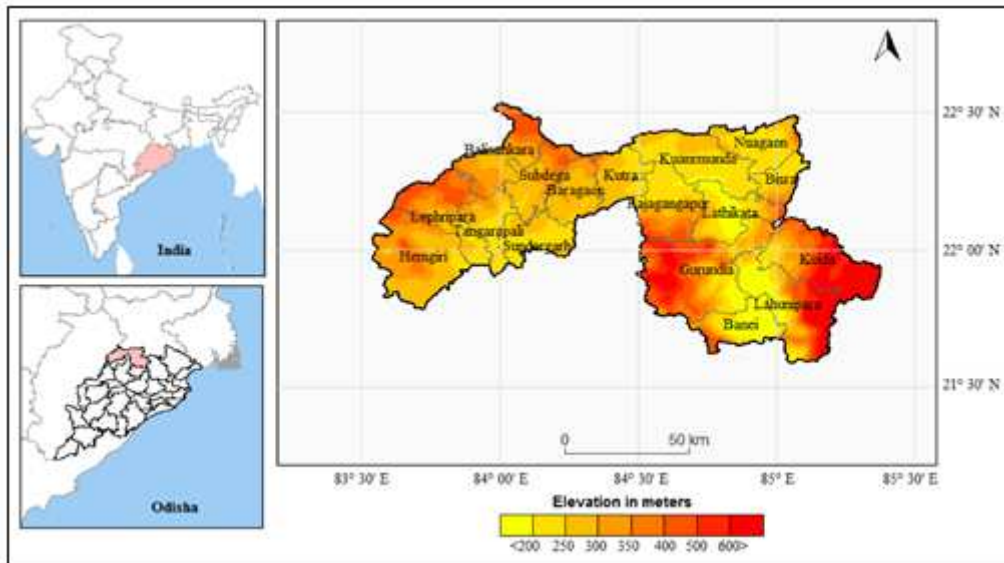
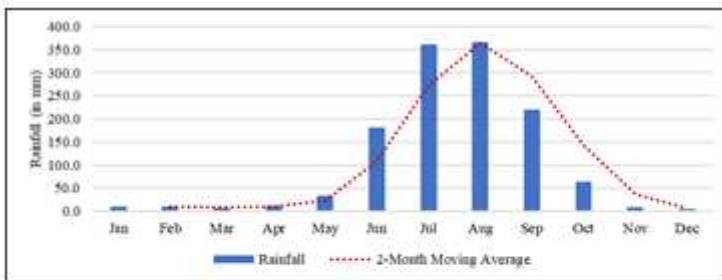
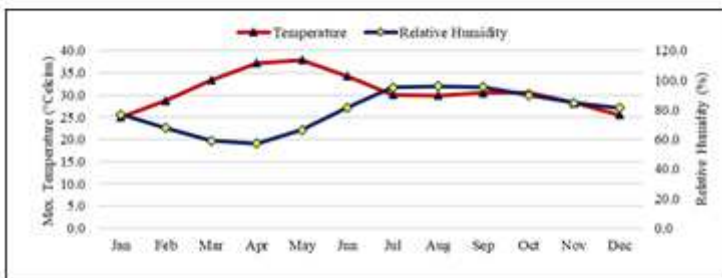


Figure 3

Study area with elevation for Sundargarh district in Odisha, India



(a)



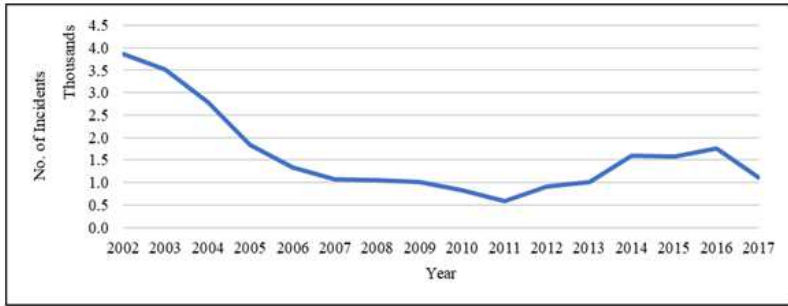
(b)



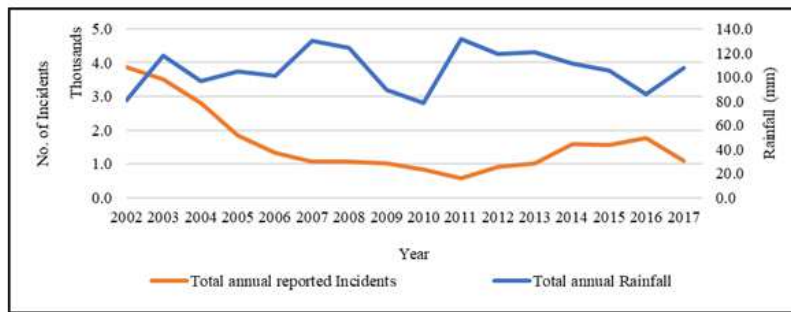
(c)

Figure 4

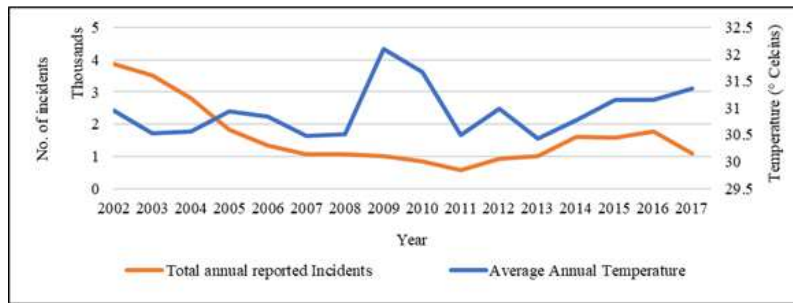
Average monthly rainfall (in mm) (a), maximum surface temperature (in °C), relative humidity (in %) (b), and incidents reported (c) for the period of 2002 to 2017.



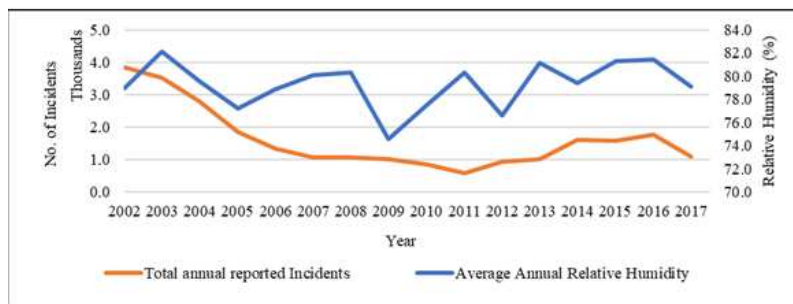
(a)



(b)



(c)



(d)

Figure 5

The trend of annual incidents (a); and comparison with average annual rainfall (b); temperature (c); and relative humidity (d) for the period of 2002 to 2017

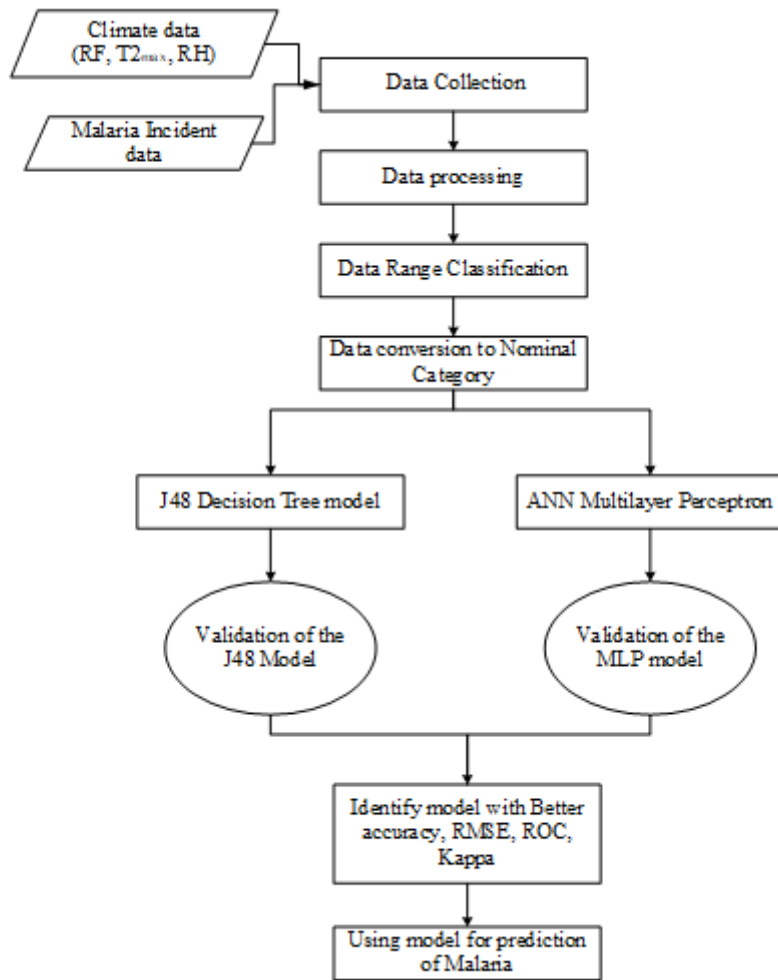


Figure 6

Detailed Methodology for Predictive Modeling of Malaria

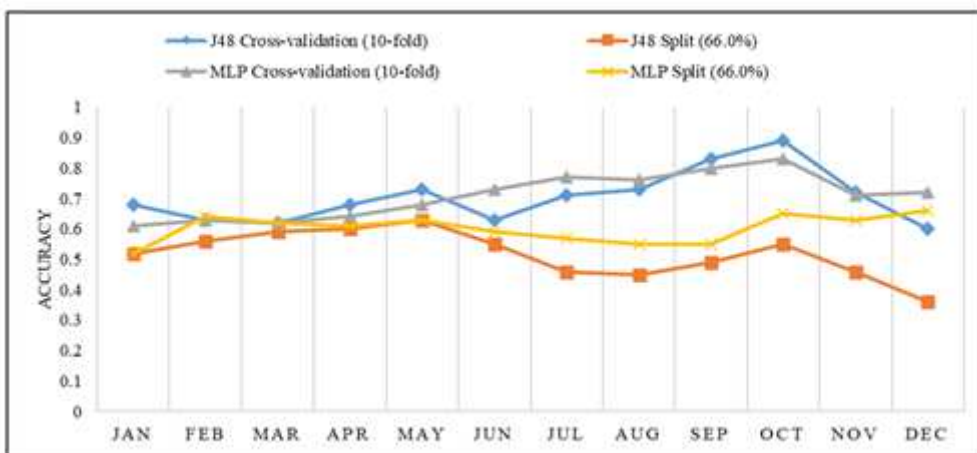


Figure 7

Month-wise Comparison of the Accuracy of the J48 (Cross-validation & Percent Split model) to the Multi-Layer Perceptron model for Accuracy for Sundargarh District

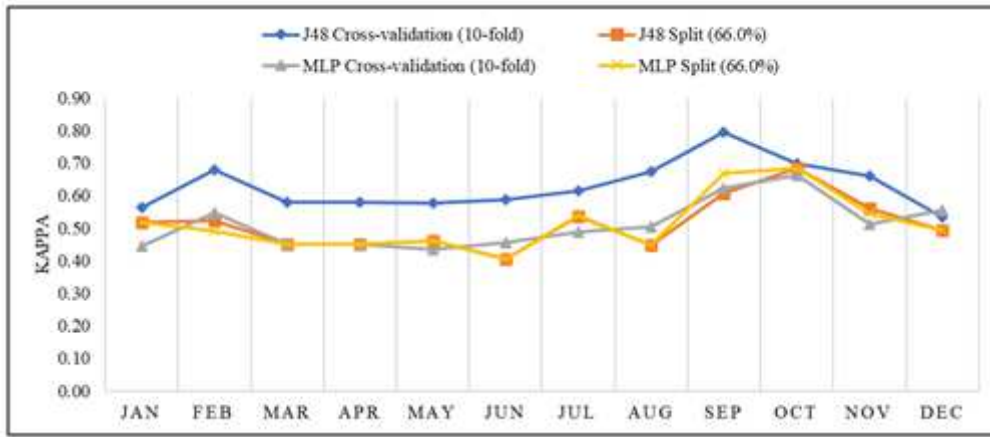


Figure 8

Month-wise Comparison of Kappa of the J48 (Cross-validation & Percent Split model) to the Multi-Layer Perceptron model for Accuracy for Sundargarh District

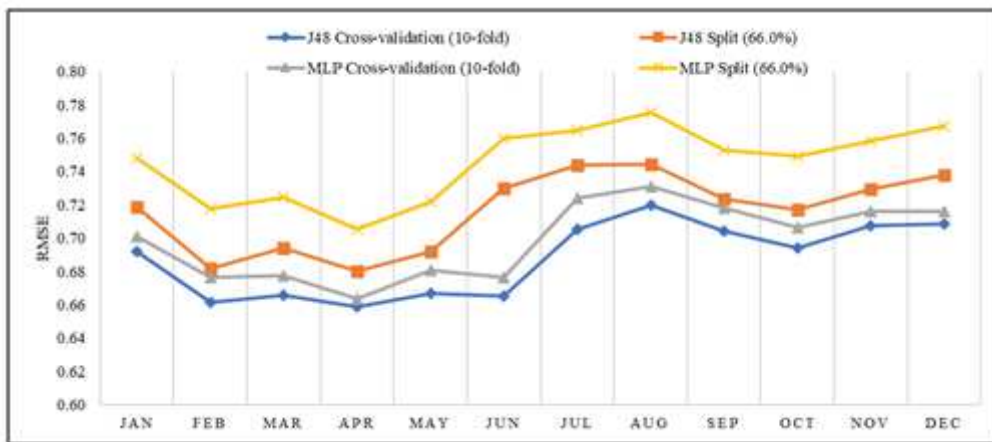
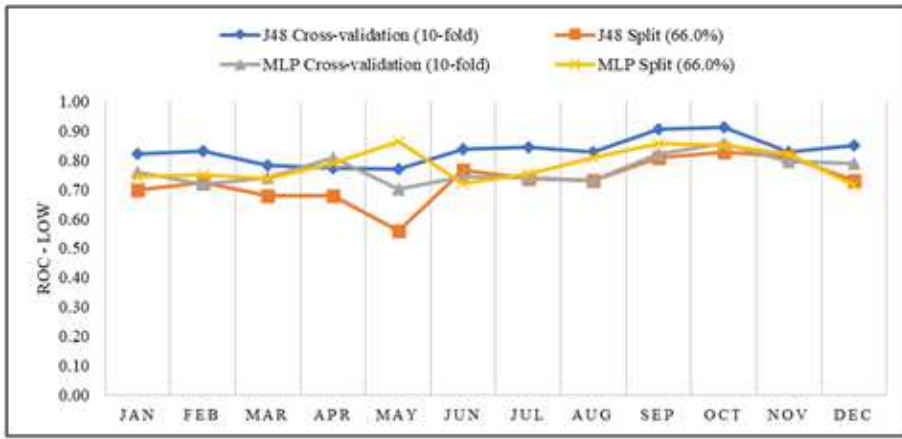
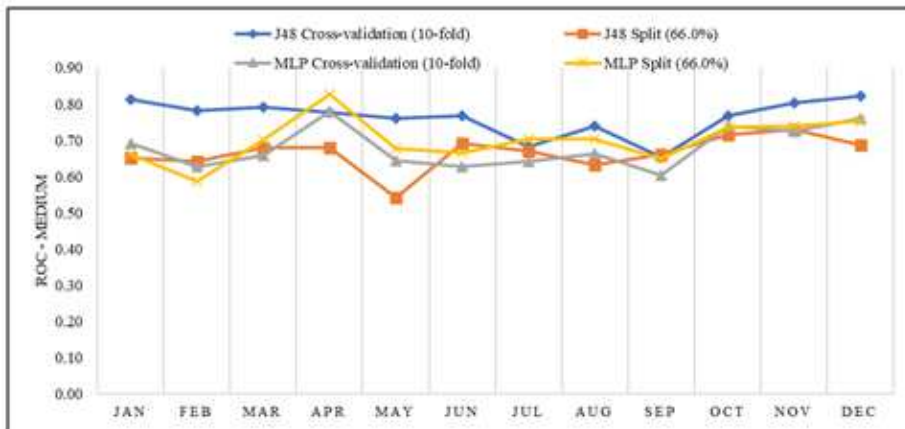


Figure 9

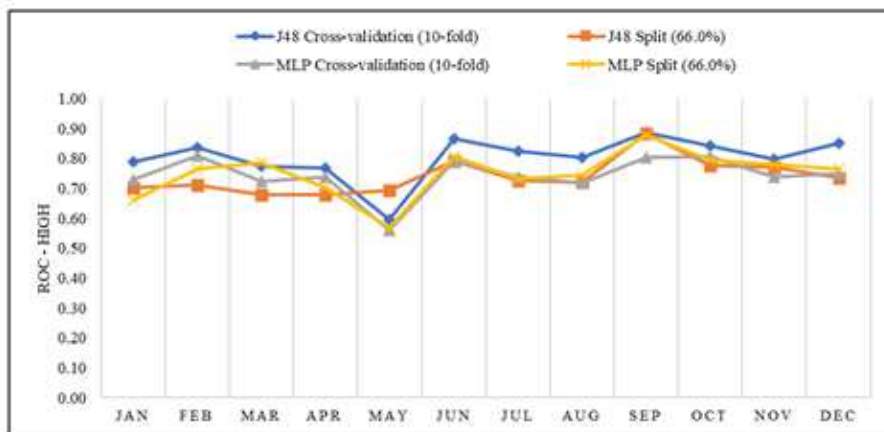
Month-wise Comparison of RMSE of the J48 (Cross-validation & Percent Split model) to the Multi-Layer Perceptron model for Accuracy for Sundargarh District



(a)



(b)



(c)

Figure 10

Month-wise Comparison of ROC for ((a) Low, (b) medium, and (c) High prediction category) of the J48 (Cross-validation & Percent Split model) to the Multi-Layer Perceptron model for Accuracy for Sundargarh District

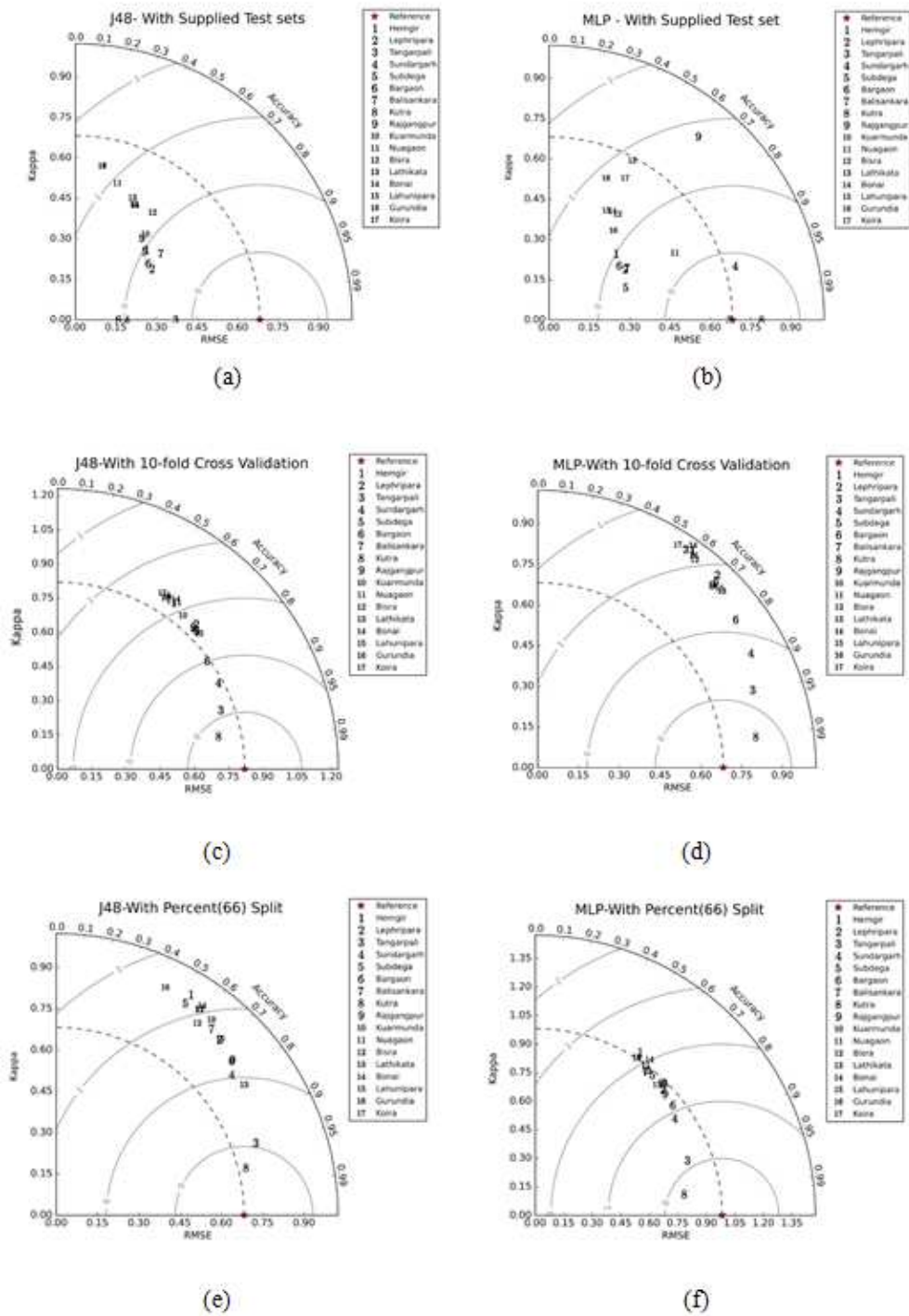
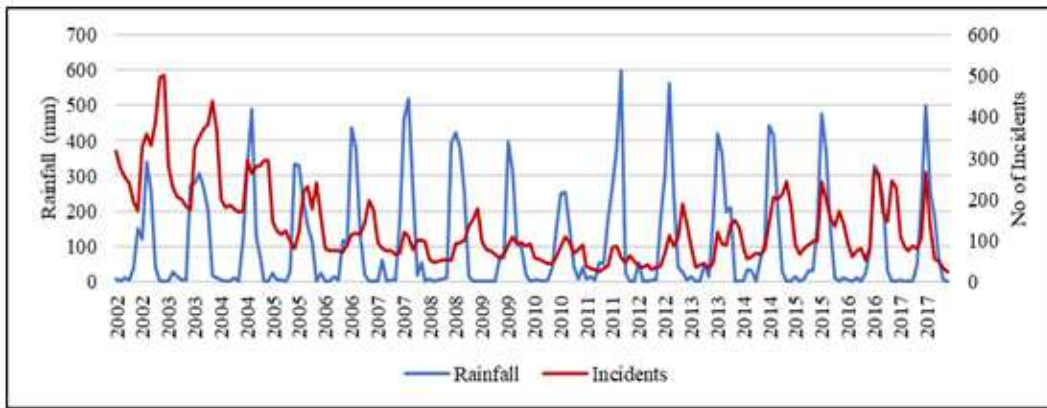
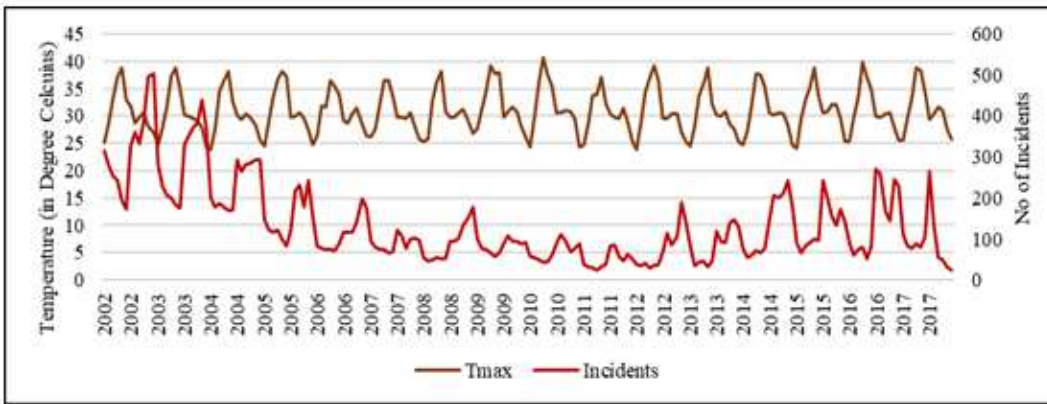


Figure 11

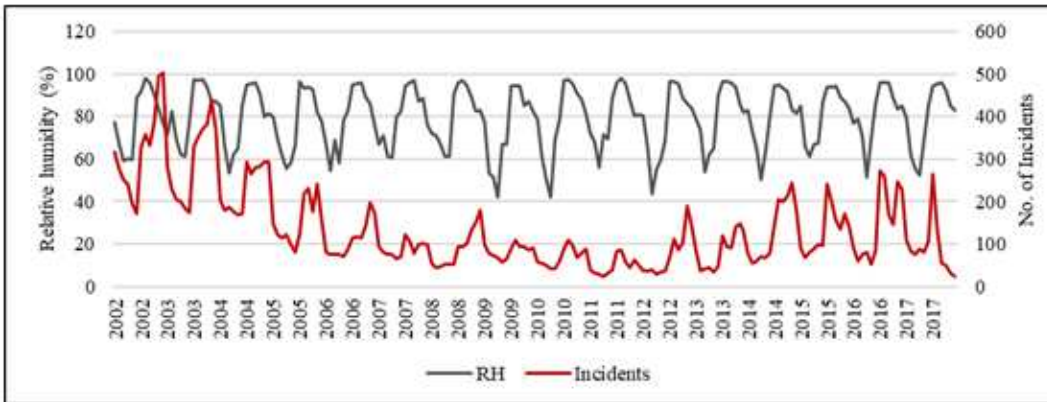
Block-wise accuracy (arc), RMSE (x-axis) and Kappa (Y-axis) for the J48 and the Multi-Layer Perceptron model with a Supplied test set, 10-fold Cross-validation & Percent Split classifiers



(a)



(b)



(c)

Figure 12

Time series plot for Monthly Rainfall (a), Maximum Temperature (b), and Relative Humidity (c) in comparison with malaria incidents for the period of 2002-2017