# Journal Pre-proof

Machine learning model for predicting malaria using clinical information

You Won Lee, Jae Woo Choi, Eun-Hee Shin

Please cite this article as: Y.W. Lee, J.W. Choi, E.-H. Shin, Machine learning model for predicting malaria using clinical information, *Computers in Biology and Medicine*, https://doi.org/10.1016/j.compbiomed.2020.104151.

This is a PDF file of an article that has undergone enhancements after acceptance, such as the addition of a cover page and metadata, and formatting for readability, but it is not yet the definitive version of record. This version will undergo additional copyediting, typesetting and review before it is published in its final form, but we are providing this version to give early visibility of the article. Please note that, during the production process, errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.

1 **Machine learning model for predicting malaria using clinical information**

2

3 You Won Lee[1] , Jae Woo Choi[2,3], Eun-Hee Shin[1,4,*]

4

5 [1]Department of Tropical Medicine and Parasitology, Seoul National University College of

6 Medicine and Institute of Endemic Diseases, Seoul 03080, Republic of Korea

7 [2]Department of Pharmacology, Yonsei University College of Medicine, Seoul 03722,

8 Republic of Korea

9 [3]Severance Biomedical Science Institute, Yonsei University College of Medicine, Seoul

10 03722, Republic of Korea

11 [4]Seoul National University Bundang Hospital, Seongnam 13620, Republic of Korea

12

13

14 **\*Corresponding author:** Eun-Hee Shin, PhD

15 Associate Professor, Department of Tropical Medicine and Parasitology, Seoul National

16 University College of Medicine

17 103 Daehak-ro, Jongno-gu, Seoul 03080, Republic of Korea

18 Tel.: +82-2-740-8344, E-mail: ehshin@snu.ac.kr

19

20

21 **Abstract**

22 Background: Rapid diagnosing is crucial for controlling malaria. Various studies have aimed

23 at developing machine learning models to diagnose malaria using blood smear images;

24 however, this approach has many limitations. This study developed a machine learning model

25 for malaria diagnosis using patient information.

26 Methods: To construct datasets, we extracted patient information from the PubMed abstracts

27 from 1956 to 2019. We used two datasets: a solely parasitic disease dataset and total dataset

28 by adding information about other diseases. We compared six machine learning models:

29 support vector machine, random forest (RF), multilayered perceptron, AdaBoost, gradient

30 boosting (GB), and CatBoost. In addition, a synthetic minority oversampling technique

31 (SMOTE) was employed to address the data imbalance problem.

32 Results: Concerning the solely parasitic disease dataset, RF was found to be the best model

33 regardless of using SMOTE. Concerning the total dataset, GB was found to be the best.

34 However, after applying SMOTE, RF performed the best. Considering the imbalanced data,

35 nationality was found to be the most important feature in malaria prediction. In case of the

36 balanced data with SMOTE, the most important feature was symptom.

37 Conclusions: The results demonstrated that machine learning techniques can be successfully

38 applied to predict malaria using patient information.

39

40 **Keywords:** Machine learning, Malaria, Diagnosis, Case reports, Patient information

41

**Introduction**

Malaria is a dangerous infection disease caused by various species of *Plasmodium* worldwide, which can be cured using drugs [1]. The World Health Organization (WHO)'s World Malaria report 2019 indicated 228 million cases of malaria, with 40,500 deaths, in more than 90 countries in 2018 [2, 3]. Early diagnosis of malaria is very important as it allows performing appropriate disease management and treatment [4-6]. Therefore, various diagnosis methods of malaria have been proposed so far, such as polymerase chain reaction (PCR), rapid diagnostic tests (RDTs), and microscopy [7-9]. Frickmann et al. evaluated a PCR assay corresponding to the differentiation of plasmodium [10]. Amaral established ribosomal- and non-ribosomal-targeting PCR assays for detecting low-density and mixed malaria [11]. Makuuchi evaluated RDTs by comparing their results with those of microscopy analysis [12]. However, these methods are generally expensive in terms of time and expert labor. Recently, machine learning-based diagnoses have been investigated to increase the diagnosis speed [5, 6, 13].

Various studies have been conducted for diagnosing malaria using machine learning [6, 14], most of which focused on the blood smear image approach [14]. Evidently, blood smear microscopic examination is the most reliable clue in parasitic disease diagnoses [15]; moreover, machine learning-based diagnosis reduces the required costs and professional labor while increasing the diagnosis accuracy [1]. However, supervised learning requires establishing appropriate labeling of images by experts to construct trained datasets, performing the so-called annotation [13, 16]. Moreover, diagnoses using microscopy methods considerably depend on the skills and experience of experts [1, 15]. Therefore, these methods require greater specificity and sensitivity of an expert [13, 15].

3

65    Meanwhile, another important indicator that needs to be considered in malaria diagnosis is

66    patient information, including symptomatology, nationality, age, gender, and travel history

67    [17, 18]. However, it is difficult to discriminate malaria infection from other parasitic

68    diseases [18], as usually, patients exhibit similar symptoms of malaria. Various effective

69    methods for machine learning diagnosis have been developed using patient information.

70    Spathis et al. considered age, gender, and symptomatology of a patient as variables for

71    diagnosing chronic obstructive pulmonary disease [19]. Terrada et al. classified and predicted

72    atherosclerosis using a machine learning approach, which trained the model on data including

73    age, gender, and symptoms [20]. Mello-Roman et al. predicted dengue using data on age,

74    gender, region, and symptomatology of a patient [21]. However, no study, so far, has

75    attempted to diagnose malaria using machine learning models trained on patient information.

76    Therefore, this paper proposes a machine learning model to predict malaria by using patient

77    information obtained from parasite case reports. We extract the data on nationality, disease,

78    gender, age, symptoms and body region of patients with symptoms. Then, we train six

79    machine learning models on these data.

80

81    **Methods**

82    **Dataset**

83    Using BioPython, we obtained the data corresponding to 56 parasitic disease reports provided

84    by the Center for Disease Control and Prevention (CDC) [22] and abstracts of case reports of

85    non-parasitic diseases (cancer, Alzheimer, rheumatoid disease, and diabetes) published from

86    1956 to 2019 by PubMed [23]. Based on CDC parasitic disease list, we classified 56 diseases

4

87  based on causative parasite genus or if the disease name was the same. For example, Hydatid

88  disease, Alveolar Echinococcosis, and Echinococosis were caused by the same genus, i.e.,

89  *Echinococcus*, and categorized as the same parasitic disease. Filaria, Filariasis,

90  Elephantiasis, and the infection of *Wuchereria bancrofti*, Brugia genus are regarded as the

91  same parasitic disease. Nonpathogenic intestinal protozoa (*Enteromonas*

92  *hominis, Retortamonas intestinalis*, and *Pentatrichomonas hominis*) can be present in feces

93  but they are not harmful and nonpathogenic. Therefore, we combined them into one disease.

94  Using this method, we were able to reorganize 56 parasitic diseases.

95  We selected non-parasitic diseases using the following standard. i) We chose diseases that

96  constitute the top ten causes of deaths worldwide [24] or diseases with more than 10000 cases

97  from 1956 to 2019. This approach was used because we wanted to have a diverse sample of

98  patients. ii) We chose a disease with a name without an organ name when collecting case

99  reports. For heart or brain diseases, if the name of the organ was included, it could overlap

100 with the case report of parasites. iii) We excluded infections because symptoms of infections

101 were similar to those of parasite infection diseases. Thus, we wanted to confirm the

102 applicability of our model in non-parasitic disease patients who had other symptoms. The

103 diseases that we used met these standards.

104 To extract relevant data, we ran queries using logical combinations such as operators "AND"

105 and "OR." We added our queries in Appendix A Supplementary Table 1. Parasitic diseases

106 have many related names. Therefore, we use "OR" and "AND". For example, another name

107 for sleeping sickness is African trypanosomiasis, and the causative parasite is Trypanosoma

108 genus. In addition, there are two types of trypanosomiasis: American trypanosomiasis

109 (Chagas disease) and African trypanosomiasis. We are only interested in African

5

110 trypanosomiasis and not American trypanosomiasis. Then, we used the following query:

111 *((Trypanosoma OR Sleeping Sickness OR trypanosomiasis) AND Africa) AND case report.*

112 If we are interested in American trypanosomiasis, then we used the following query:

113 *((Trypanosoma OR Chagas Disease OR trypanosomiasis) AND America) AND case report.*

114 We derived information regarding nation (meaning nationality or travel region of a patient),

115 disease, gender, age, symptom and body region of patients with symptoms using Python

116 scripts. The lists of body regions and symptoms were prepared by referring to the 10th edition

117 of International Classification of Disease.

118

119 **Dataset preprocessing**

120 Fig. 1 shows the data processing scheme. We removed missing variables or values from the

121 dataset, except for symptoms and body regions. If the symptoms or body regions had at least

122 one value, we did not remove these data. First, we constructed a dataset comprising only

123 parasitic disease patient information, and then, prepared the total dataset by adding

124 information about other diseases (Alzheimer, rheumatoid, cancer, and diabetes). All the

125 values were categorized using integers. Note that the prepared datasets incurred the data

126 imbalance problem. To address this problem, we applied the synthetic minority oversampling

127 technique (SMOTE) [25] provided by Scikit-learn [26].

128

129 **Model development**

130 We used various machine learning techniques to develop six models to diagnose malaria:

131 support vector machine (SVM) [27], random forest (RF) [28], multilayered perceptron (MLP)

132 [29], AdaBoost (Ada) [30], gradient boosting (GB) [31], and CatBoost (CB) [32].

133 SVM

134 SVM is a widely used supervised learning approach for classification or regression analysis.

135 It can be applied to transform training data into a high-dimensional feature space and

136 determine a linear optimal solution by separating a hyperplane that provides the smallest

137 distance between the hyperplane points and the largest margin between the classes [27, 33-

138 35].

139 RF

140 RF is an ensemble supervised learning method composed of multiple decision trees

141 corresponding to various subdatasets. Each tree calculates the results and obtains the average

142 of the prediction outcomes. This approach allows reducing variance in decision trees [28, 36,

143 37].

144 MLP

145 MLP is a supervised machine learning algorithm used for data classification tasks. It is

146 composed of three layers: an input layer, which includes input data; a hidden layer, which

147 computes complicated associations across the network; and an output layer, which generates

148 the final result. This process can be terminated when the error rate becomes sufficiently small.

149 We optimized the log-loss function using the stochastic gradient descent [29, 38].

150 Ada

151 Ada is an ensemble learning algorithm used to elevate a weak classifier to a strong one. First,

152 it trains a base classifier and assigns higher weights to the misclassified samples; thereafter, it

153 is applied to the next process. This iterative process continues until a stop condition is

7

154 reached or the error rate becomes sufficiently small [30, 39, 40].

155 GB

156 GB is an ensemble model based on decision trees. It minimizes the residual (negative

157 gradient) using gradient descents to classify data [31, 36, 41].

158 CB

159 CB is a modification of GB, and yields high performance in case of categorical features [32,

160 42].

161 These models were trained to execute effective malaria diagnosis.

162

163 **Model evaluation**

164 A 10-fold cross-validation (CV) was applied to avoid overfitting. In 10-fold CV, the data are

165 first randomly split into ten parts. Then, each subset is considered as the testing one and the

166 remaining subsets are used for training. The results of 10-fold CV are the averaged values of

167 accuracy obtained from the ten tests. We evaluated each model in terms of accuracy, precision,

168 recall, and F1-score. These parameters are defined as follows:

$$Accuracy = \frac{True\ positives + True\ neagtives}{True\ positives + True\ negatives + False\ positives + False\ negatives}$$

$$Precision = \frac{True\ positives}{True\ positives\ +\ False\ positives}$$

$$Recall = \frac{True\ positives}{True\ positives + False\ negatives}$$

$$F1 - score = \frac{2 * Precision * Recall}{Precision + Recall}$$

169

170    In addition, an AUC curve was drawn to compare the performance of each model.

171

172    **Feature importance**

173    We analyzed the feature importance using RF with optimized hyperparameters to evaluate

174    how each model works using yellowbrick [43]. The result is calculated based on the average

175    of the feature importance associated with each feature in terms of achieving high

176    performance for the model.

177

178    **Results**

179    **Data statistics**

180    In this study, information about 1,846 patients, obtained from case reports, was considered

181    (Table 1). Overall, 1,698 patients had parasitic diseases, 135 patients had malaria, and 148

182    patients had non-parasitic diseases. In the total dataset, the portion of malaria patients was

183    7.31% and that of the solely parasitic disease dataset was 8%. Table 1 provides information

184    about the patient demographics.

185

186    **Model performance**

187    Tables 2 and 3 and Fig. 2 describe the performance of the considered predictive models.

9

188  Concerning the solely parasitic disease dataset, the RF model achieved the best performance

189  with AUC of 73.2%. The worst model was Ada, with AUC of 59.6%. After applying SMOTE,

190  the AUC values of almost all models increased, except GB. In this case, RF achieved the best

191  performance (with AUC of 73.5%), while Ada demonstrated the worst performance (with

192  AUC of 68.8%). Within the total dataset, GB achieved the highest AUC (85.6%). The

193  performance of the models trained on the total dataset was higher compared to those trained

194  on the solely parasitic disease dataset. The values of accuracy, precision, recall, and F1-score

195  were also higher in the case of training on the total dataset. The model that showed the worst

196  results was SVM, with AUC of 77.6%. The AUC values of all models with SMOTE were

197  decreased; however, the values of accuracy, precision, recall, and F1-score were higher. RF

198  achieved the best performance among all classifiers (with AUC of 80.5%), while the worst

199  performing model was MLP (with AUC of 67.4%).

200

201  **Feature importance**

202  We calculated the feature importance using RF, which achieved the highest performance

203  overall (Fig. 3). Both in the total dataset and the solely parasitic disease dataset, the most

204  important feature was "nation," followed by "age" (Fig. 3A and C). However, in case of the

205  dataset with SMOTE, "symptom" was the most important feature, followed by "nation" (Fig.

206  3B and D).

207

208  **Discussion**

209  Recently, an increasing number of studies have been conducted on malaria diagnosis using

10

210 artificial intelligence (AI). Kim et al. and Wang et al. predicted malaria incidence by using a

211 seasonal climate dataset [44, 45]. Moreover, the methods based on AI for diagnosis using

212 blood smear images have been extensively investigated [1, 4, 6, 14]. Rajaraman et al. used

213 thin-blood smear images to construct deep neural ensemble models [4]. Molina et al.

214 introduced a machine learning model to discriminate infected blood cells from normal ones

215 [6].

216 In the present study, we used the parasitic disease patient information derived from the

217 abstracts of case reports provided by PubMed to train the models. Evidently, it is possible to

218 consider various databases for obtaining the epidemiology or symptom data on parasitic

219 diseases, such as Gideon [46] and CDC [22]. Moreover, the National Health and Nutrition

220 Examination Survey was used as a source for obtaining health and nutrition information

221 about patients. However, no database on parasitic disease patients provides information about

222 patients' nationality, age, symptoms, and gender. Even if such information was available, it

223 would not exhibit diversity in terms of regions or conditions of patients [47-49]. Therefore,

224 we constructed datasets based on the information obtained from the abstracts of all parasitic

225 disease case reports available in PubMed that were published from 1956 to 2019. Note that

226 the abstracts did not provide detailed information about patients; however, the available data

227 were sufficient to perform analysis to diagnose malaria using the methods considered in the

228 present study. Moreover, these data reflected the trend of overall parasitic disease patient

229 information with sufficient accuracy.

230 The performance estimates of almost all models trained on the solely parasitic disease dataset

231 were lower than those of the models trained on the total dataset and those trained on the data

232 with SMOTE. The observed results could be explained not only by a smaller dataset but also

11

233   by the characteristics of parasitic diseases. In clinical cases, the symptoms of a parasitic

234   disease are similar to those of malaria, and therefore, it was more difficult to discriminate

235   malaria using the solely parasitic disease dataset. For example, fever, which is a standard

236   symptom of malaria, can also indicate conditions such as toxoplasmosis [50] and pulmonary

237   eosinophilia [51]. Similarly, abdominal pain is a generic symptom for conditions such as

238   amoebic liver abscesses [52] and trichinellosis [53]. We hypothesized that SMOTE can

239   address this problem through oversampling; however, the AUC values were still lower than

240   those of the models trained on the total dataset.

241   According to the obtained results, RF achieved the best performance, except for the total

242   dataset without SMOTE. The remarkable performance of an RF model has also been reported

243   by other studies concerning various diseases [54-56]. An RF model has also been applied to

244   neuroimaging classification [54], the prediction of in-hospital cardiac arrest [56], and

245   biomarker prediction based on gene expression data [55]; in all these applications, the RF

246   model demonstrated great performance.

247   Meanwhile, the results of feature importance analysis indicated that nationality and age are

248   important factors to consider in diagnosing malaria using imbalanced data. Many previous

249   studies have reported that parasitic diseases, such as malaria, depend on the places in which

250   patients live or travel to [2, 57]. The obtained results suggested that nationality and the region

251   of traveling are crucial factors in the diagnosis of parasitic diseases.

252

253   **Limitations**

254   The limitations of our study are related to the small size of the datasets used and the limited

255 number of features without the process of feature selection. Moreover, the observed values of

256 precision, recall, and F1-score were lower than those reported previously. Specifically, the

257 dataset has high imbalance between parasitic and non-parasitic cases. We collected data from

258 more than 35519 non-parasitic patients, expecting to be able to obtain more than 1698

259 parasitic patients. However, the number of samples inevitably decreased when extracting only

260 patients whose information, such as the country, age, gender, symptoms, and body region of

261 patients with symptoms, was available. In particular, cancer, rheumatism, diabetes, and

262 Alzheimer's patients were able to provide less nationality information in case reports.

263 Therefore, we had no choice but to create a dataset with a small number of patients having

264 non-parasitic diseases. We considered that they could be improved by applying SMOTE.

265 However, even after the application of SMOTE, the value of precision, recall and F1-score

266 did not exceed 0.5. This can be addressed by increasing the number of patients and features in

267 the datasets.

268

269 **Conclusions**

270 This is the first study that aims to diagnose malaria using patient information. The novelty of

271 the utilized datasets lies in the fact that the data were obtained for parasitic disease patients

272 spread globally. We compared several machine learning models applied to malaria prediction

273 trained on parasitic disease patient data. The results showed that RF was the best model for

274 the diagnosis, indicating the possibility of diagnosing using only patient information with AI.

275

276

13

277 **Acknowledgments**

280

281 **Author contributions**

282 EHS provided the research idea. YWL and EHS conceived and designed the study. YWL and

283 JWC collected and analyzed data. YWL and JWC contributed materials and analysis tools.

284 YWL and EHS wrote the paper. EHS was responsible for the overall project administration

285 and acquiring of financial support.

286

287 **References**

288 [1] M. Poostchi, K. Silamut, R.J. Maude, S. Jaeger, G. Thoma, Image analysis and machine

289 learning for detecting malaria, Transl Res, 194 (2018) 36-55

290 https://doi.org/10.1016/j.trsl.2017.12.004.

291 [2] L. Zekar, T. Sharman, Malaria (Plasmodium Falciparum), StatPearls, Treasure Island (FL),

292 2020.

293 [3] W.H. Organization, World malaria report 2019, (2019).

294 [4] S. Rajaraman, S. Jaeger, S.K. Antani, Performance evaluation of deep neural ensembles

295 toward malaria parasite detection in thin-blood smear images, PeerJ, 7 (2019) e6977

296 https://doi.org/10.7717/peerj.6977.

297    [5] K. Torres, C.M. Bachman, C.B. Delahunt, J. Alarcon Baldeon, F. Alava, D. Gamboa Vilela,

298    S. Proux, C. Mehanian, S.K. McGuire, C.M. Thompson, T. Ostbye, L. Hu, M.S. Jaiswal, V.M.

299    Hunt, D. Bell, Automated microscopy for routine malaria diagnosis: a field comparison on

300    Giemsa-stained blood films in Peru, Malar J, 17 (2018) 339 https://doi.org/10.1186/s12936-

301    018-2493-0.

302    [6] A. Molina, S. Alferez, L. Boldu, A. Acevedo, J. Rodellar, A. Merino, Sequential

303    classification system for recognition of malaria infection using peripheral blood cell images, J

304    Clin Pathol, (2020) https://doi.org/10.1136/jclinpath-2019-206419.

305    [7] Z. Zheng, Z. Cheng, Advances in Molecular Diagnosis of Malaria, Adv Clin Chem, 80

306    (2017) 155-192 https://doi.org/10.1016/bs.acc.2016.11.006.

307    [8] P. Berzosa, A. de Lucio, M. Romay-Barja, Z. Herrador, V. Gonzalez, L. Garcia, A.

308    Fernandez-Martinez, M. Santana-Morales, P. Ncogo, B. Valladares, M. Riloha, A. Benito,

309    Comparison of three diagnostic methods (microscopy, RDT, and PCR) for the detection of

310    malaria parasites in representative samples from Equatorial Guinea, Malar J, 17 (2018) 333

311    https://doi.org/10.1186/s12936-018-2481-4.

312    [9] K.O. Mfuh, O.A. Achonduh-Atijegbe, O.N. Bekindaka, L.F. Esemu, C.D. Mbakop, K.

313    Gandhi, R.G.F. Leke, D.W. Taylor, V.R. Nerurkar, A comparison of thick-film microscopy,

314    rapid diagnostic test, and polymerase chain reaction for accurate diagnosis of Plasmodium

315    falciparum malaria, Malar J, 18 (2019) 73 https://doi.org/10.1186/s12936-019-2711-4.

316    [10] H. Frickmann, C. Wegner, S. Ruben, C. Behrens, H. Kollenda, R. Hinz, S. Rojak, N.G.

317    Schwarz, R.M. Hagen, E. Tannich, Evaluation of the multiplex real-time PCR assays

318    RealStar malaria S&T PCR kit 1.0 and FTD malaria differentiation for the differentiation of

15

319    Plasmodium species in clinical samples, Travel Med Infect Dis, 31 (2019) 101442

320    https://doi.org/10.1016/j.tmaid.2019.06.013.

321    [11] L.C. Amaral, D.R. Robortella, L.F.F. Guimaraes, J.E. Limongi, C.J.F. Fontes, D.B.

322    Pereira, C.F.A. de Brito, F.S. Kano, T.N. de Sousa, L.H. Carvalho, Ribosomal and non-

323    ribosomal PCR targets for the detection of low-density and mixed malaria infections, Malar J,

324    18 (2019) 154 https://doi.org/10.1186/s12936-019-2781-3.

325    [12] R. Makuuchi, S. Jere, N. Hasejima, T. Chigeda, J. Gausi, The correlation between

326    malaria RDT (Paracheck pf.(R)) faint test bands and microscopy in the diagnosis of malaria

327    in Malawi, BMC Infect Dis, 17 (2017) 317 https://doi.org/10.1186/s12879-017-2413-x.

328    [13] A. Rehman, N. Abbas, T. Saba, Z. Mehmood, T. Mahmood, K.T. Ahmed, Microscopic

329    malaria parasitemia diagnosis and grading on benchmark datasets, Microsc Res Tech, 81

330    (2018) 1042-1058 https://doi.org/10.1002/jemt.23071.

331    [14] S. Rajaraman, S.K. Antani, M. Poostchi, K. Silamut, M.A. Hossain, R.J. Maude, S.

332    Jaeger, G.R. Thoma, Pre-trained convolutional neural networks as feature extractors toward

333    improved malaria parasite detection in thin blood smear images, PeerJ, 6 (2018) e4568

334    https://doi.org/10.7717/peerj.4568.

335    [15] A. Mbanefo, N. Kumar, Evaluation of Malaria Diagnostic Methods as a Key for

336    Successful Control and Elimination Programs, Trop Med Infect Dis, 5 (2020)

337    https://doi.org/10.3390/tropicalmed5020102.

338    [16] K. Smith, F. Piccinini, T. Balassa, K. Koos, T. Danka, H. Azizpour, P. Horvath,

339    Phenotypic Image Analysis Software Tools for Exploring and Understanding Big Image Data

340    from    Cell-Based    Assays,    Cell    Syst,    6    (2018)    636-653

341 https://doi.org/10.1016/j.cels.2018.06.001.

342 [17] F. Jimenez-Morillas, M. Gil-Mosquera, E.J. Garcia-Lamberechts, I.-S. en representacion

343 de la seccion de enfermedades tropicales de, I.-S. Seccion de enfermedades tropicales de,

344 Fever in travellers returning from the tropics, Med Clin (Barc), 153 (2019) 205-212

345 https://doi.org/10.1016/j.medcli.2019.03.017.

346 [18] C. JY, Seo and Lee's Clinical Parasitology, Seoul National University Publishing

347 Council2011.

348 [19] D. Spathis, P. Vlamos, Diagnosing asthma and chronic obstructive pulmonary disease

349 with machine learning, Health Informatics J, 25 (2019) 811-827

350 https://doi.org/10.1177/1460458217723169.

351 [20] O. Terrada, B. Cherradi, A. Raihani, O. Bouattane, Classification and Prediction of

352 atherosclerosis diseases using machine learning algorithms, 2019 5th International

353 Conference on Optimization and Applications (ICOA), IEEE, 2019, pp. 1-5.

354 [21] J.D. Mello-Roman, J.C. Mello-Roman, S. Gomez-Guerrero, M. Garcia-Torres, Predictive

355 Models for the Medical Diagnosis of Dengue: A Case Study in Paraguay, Comput Math

356 Methods Med, 2019 (2019) 7307803 https://doi.org/10.1155/2019/7307803.

357 [22] C.f.D.C.a. Prevention, DPDx - Laboratory Identification of Parasites of Public Health

358 Concern 2020.

359 [23] P.J. Cock, T. Antao, J.T. Chang, B.A. Chapman, C.J. Cox, A. Dalke, I. Friedberg, T.

360 Hamelryck, F. Kauff, B. Wilczynski, M.J. de Hoon, Biopython: freely available Python tools

361 for computational molecular biology and bioinformatics, Bioinformatics, 25 (2009) 1422-

362 1423 https://doi.org/10.1093/bioinformatics/btp163.

363 [24] J.S. Rana, S.S. Khan, D.M. Lloyd-Jones, S. Sidney, Changes in Mortality in Top 10

364 Causes of Death from 2011 to 2018, J Gen Intern Med, (2020)

365 https://doi.org/10.1007/s11606-020-06070-z.

366 [25] N.V. Chawla, K.W. Bowyer, L.O. Hall, W.P.J.J.o.a.i.r. Kegelmeyer, SMOTE: synthetic

367 minority over-sampling technique, 16 (2002) 321-357.

368 [26] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P.

369 Prettenhofer, R. Weiss, V.J.t.J.o.m.L.r. Dubourg, Scikit-learn: Machine learning in Python, 12

370 (2011) 2825-2830.

371 [27] C. Cortes, V.J.M.l. Vapnik, Support-vector networks, 20 (1995) 273-297.

372 [28] L.J.M.l. Breiman, Random forests, 45 (2001) 5-32.

373 [29] X. Glorot, Y. Bengio, Understanding the difficulty of training deep feedforward neural

374 networks, Proceedings of the thirteenth international conference on artificial intelligence and

375 statistics, 2010, pp. 249-256.

376 [30] T. Hastie, S. Rosset, J. Zhu, H.J.S. Zou, i. Interface, Multi-class adaboost, 2 (2009) 349-

377 360.

378 [31] J.H.J.A.o.s. Friedman, Greedy function approximation: a gradient boosting machine,

379 (2001) 1189-1232.

380 [32] A.V. Dorogush, V. Ershov, A.J.a.p.a. Gulin, CatBoost: gradient boosting with categorical

381 features support, (2018).

382 [33] A. Gupta, B. Kahali, Machine learning-based cognitive impairment classification with

383 optimal combination of neuropsychological tests, Alzheimers Dement (N Y), 6 (2020)

18

384    e12049 https://doi.org/10.1002/trc2.12049.

385    [34] N. Liu, R. Zhao, L. Qiao, Y. Zhang, M. Li, H. Sun, Z. Xing, X. Wang, Growth Stages

386    Classification of Potato Crop Based on Analysis of Spectral Response and Variables

387    Optimization, Sensors (Basel), 20 (2020) https://doi.org/10.3390/s20143995.

388    [35] A. Gupta, R. Katarya, Social media based surveillance systems for healthcare using

389    machine learning: A systematic review, J Biomed Inform, 108 (2020) 103500

390    https://doi.org/10.1016/j.jbi.2020.103500.

391    [36] A. Dinh, S. Miertschin, A. Young, S.D. Mohanty, A data-driven approach to predicting

392    diabetes and cardiovascular disease with machine learning, BMC Med Inform Decis Mak, 19

393    (2019) 211 https://doi.org/10.1186/s12911-019-0918-5.

394    [37] C. Wang, X. Chen, L. Du, Q. Zhan, T. Yang, Z. Fang, Comparison of machine learning

395    algorithms for the identification of acute exacerbations in chronic obstructive pulmonary

396    disease,      Comput      Methods      Programs      Biomed,      188      (2020)      105267

397    https://doi.org/10.1016/j.cmpb.2019.105267.

398    [38] A.M. Ahmed, S.F. Aly, Egyptian License Plates Recognition System Using Morphologial

399    Operations and Multi Layered Perceptron, ICT in our lives-2019, 2019.

400    [39] B.X. Tran, G.H. Ha, L.H. Nguyen, G.T. Vu, M.T. Hoang, H.T. Le, C.A. Latkin, C.S.H.

401    Ho, R.C.M. Ho, Studies of Novel Coronavirus Disease 19 (COVID-19) Pandemic: A Global

402    Analysis      of      Literature,      Int      J      Environ      Res      Public      Health,      17      (2020)

403    https://doi.org/10.3390/ijerph17114095.

404    [40] L. Liu, C. Zhang, G. Zhang, Y. Gao, J. Luo, W. Zhang, Y. Li, Y. Mu, A study of aortic

405    dissection screening method based on multiple machine learning models, J Thorac Dis, 12

19

406    (2020) 605-614 https://doi.org/10.21037/jtd.2019.12.119.

407    [41] Y. Ye, Y. Xiong, Q. Zhou, J. Wu, X. Li, X. Xiao, Comparison of Machine Learning

408    Methods and Conventional Logistic Regressions for Predicting Gestational Diabetes Using

409    Routine Clinical Data: A Retrospective Cohort Study, J Diabetes Res, 2020 (2020) 4168340

410    https://doi.org/10.1155/2020/4168340.

411    [42] A. Gupta, A.S.R. Potty, D. Ganta, R.J. Mistovich, S. Penna, C. Cady, A.G. Potty,

412    Streamlining the KOOS Activities of Daily Living Subscale Using Machine Learning, Orthop

413    J Sports Med, 8 (2020) 2325967120910447 https://doi.org/10.1177/2325967120910447.

414    [43] B. Bengfort, R.J.J.o.O.S.S. Bilbro, Yellowbrick: Visualizing the scikit-learn model

415    selection process, 4 (2019) 1075.

416    [44] Y. Kim, J.V. Ratnam, T. Doi, Y. Morioka, S. Behera, A. Tsuzuki, N. Minakawa, N. Sweijd,

417    P. Kruger, R. Maharaj, C.C. Imai, C.F.S. Ng, Y. Chung, M. Hashizume, Malaria predictions

418    based on seasonal climate forecasts in South Africa: A time series distributed lag nonlinear

419    model, Sci Rep, 9 (2019) 17882 https://doi.org/10.1038/s41598-019-53838-3.

420    [45] M. Wang, H. Wang, J. Wang, H. Liu, R. Lu, T. Duan, X. Gong, S. Feng, Y. Liu, Z. Cui, C.

421    Li, J. Ma, A novel model for malaria prediction based on ensemble algorithms, PLoS One, 14

422    (2019) e0226910 https://doi.org/10.1371/journal.pone.0226910.

423    [46] G. Informatics, GIDEON, 2020.

424    [47] S. Mahmoudi, S. Mamishi, M. Banar, B. Pourakbari, H. Keshavarz, Epidemiology of

425    echinococcosis in Iran: a systematic review and meta-analysis, BMC Infect Dis, 19 (2019)

426    929 https://doi.org/10.1186/s12879-019-4458-5.

427    [48] M. Kotepui, K.U. Kotepui, Prevalence and laboratory analysis of malaria and dengue co-

428    infection: a systematic review and meta-analysis, BMC Public Health, 19 (2019) 1148

429    https://doi.org/10.1186/s12889-019-7488-4.

430    [49] D. Pierce, L. Merone, C. Lewis, T. Rahman, J. Croese, A. Loukas, M. McDonald, P.

431    Giacomin, R. McDermott, Safety and tolerability of experimental hookworm infection in

432    humans with metabolic disease: study protocol for a phase 1b randomised controlled clinical

433    trial, BMC Endocr Disord, 19 (2019) 136 https://doi.org/10.1186/s12902-019-0461-5.

434    [50] A.S. Kota, N. Shabbir, Congenital Toxoplasmosis, StatPearls, Treasure Island (FL), 2020.

435    [51] S.K. Jha, B. Karna, K. Mahajan, Tropical Pulmonary Eosinophilia, StatPearls, Treasure

436    Island (FL), 2020.

437    [52] T. Tharmaratnam, T. Kumanan, M.A. Iskandar, K. D'Urzo, P. Gopee-Ramanan, M.

438    Loganathan, T. Tabobondung, T.A. Tabobondung, S. Sivagurunathan, M. Patel, I. Tobbia,

439    Entamoeba histolytica and amoebic liver abscess in northern Sri Lanka: a public health

440    problem, Trop Med Health, 48 (2020) 2 https://doi.org/10.1186/s41182-020-0193-2.

441    [53] P. Rawla, S. Sharma, Trichinella Spiralis (Trichnellosis), StatPearls, Treasure Island (FL),

442    2020.

443    [54] S.I. Dimitriadis, D. Liparas, I. Alzheimer's Disease Neuroimaging, How random is the

444    random forest? Random forest algorithm on the service of structural imaging biomarkers for

445    Alzheimer's disease: from Alzheimer's disease neuroimaging initiative (ADNI) database,

446    Neural Regen Res, 13 (2018) 962-970 https://doi.org/10.4103/1673-5374.233433.

447    [55] L. Guo, Z. Wang, Y. Du, J. Mao, J. Zhang, Z. Yu, J. Guo, J. Zhao, H. Zhou, H. Wang, Y.

448    Gu, Y. Li, Random-forest algorithm based biomarkers in predicting prognosis in the patients

449    with    hepatocellular    carcinoma,    Cancer    Cell    Int,    20    (2020)    251

450     https://doi.org/10.1186/s12935-020-01274-z.

451     [56] R. Ueno, L. Xu, W. Uegami, H. Matsui, J. Okui, H. Hayashi, T. Miyajima, Y. Hayashi, D.

452     Pilcher, D. Jones, Value of laboratory results in addition to vital signs in a machine learning

453     algorithm to predict in-hospital cardiac arrest: A single-center retrospective cohort study,

454     PLoS One, 15 (2020) e0235835 https://doi.org/10.1371/journal.pone.0235835.

455     [57] F. Jimenez-Morillas, M. Gil-Mosquera, E.J. Garcia-Lamberechts, I.-S.t.d. department,

456     Fever in travellers returning from the tropics, Med Clin (Engl Ed), 153 (2019) 205-212

457     https://doi.org/10.1016/j.medcle.2019.03.013.

458

459     **Figure & table legends**

460     **Table 1.** Dataset review

461     **Table 2.** Solely parasite dataset

462     **Table 3**. Total dataset

463     **Fig 1**. Data processing

464     **Fig 2**. AUC curve: A) the solely parasitic disease dataset; B) solely parasitic disease dataset

465     with SMOTE; C) total dataset; D) total dataset with SMOTE.

466     **Fig 3**. Feature importance: A) the solely parasitic disease dataset; B) solely parasitic disease

467     dataset with SMOTE; C) total dataset; D) total dataset with SMOTE.

468     **Appendix A. Supplementary Table 1**. Query list of parasitic diseases

469

22

470

471

472

473

474

475    **Fig 1**. Data processing



Download abstract of case reports from
PubMed(1956~2019)
Number of parasitic disease case reports:25975
Number of non-parasitic disease case reports:35519

Extraction of patient information

Exclude missing value
Number of parasitic disease case reports:24277
Number of non-parasitic disease case reports:35371

**Total dataset**
Parasitic disease case reports
+ Non-parasitic disease case reports
(N=1846)

**Only parasitic disease dataset**
Parasitic disease case reports
(N=1698)

Categorized all value

SMOTE

Support vector machine, RandomForest,
Multi layer Perceptron, AdaBoost, GradientBoost, CatBoost

476

477

478

479

480

481

482

483

**Fig 2**. AUC curve: A) the solely parasitic disease dataset; B) solely parasitic disease dataset

with SMOTE; C) total dataset; D) total dataset with SMOTE.

486



487

488

489

490

491

492

493
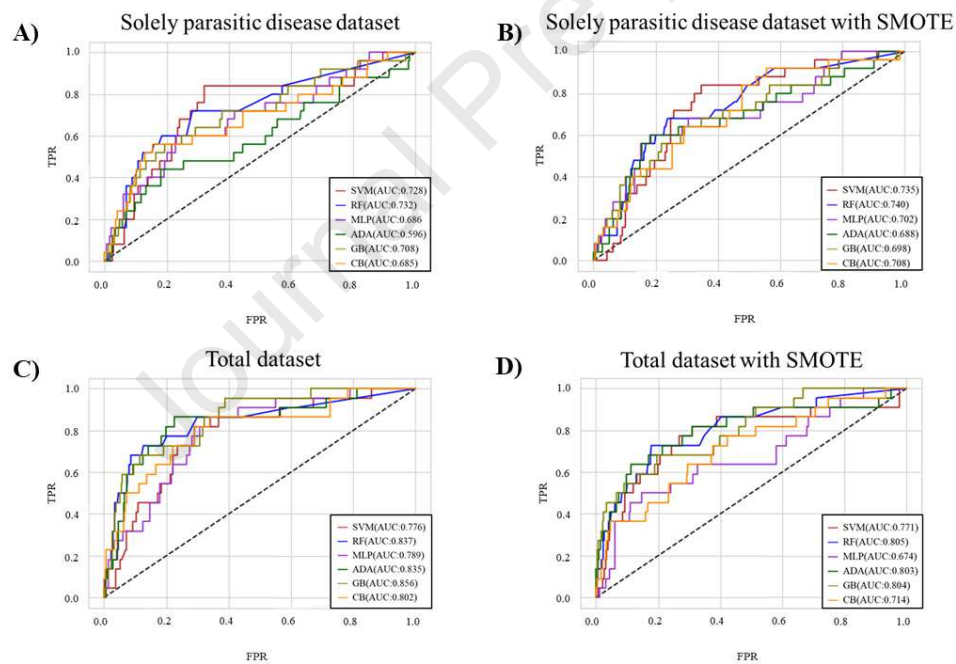
494

495
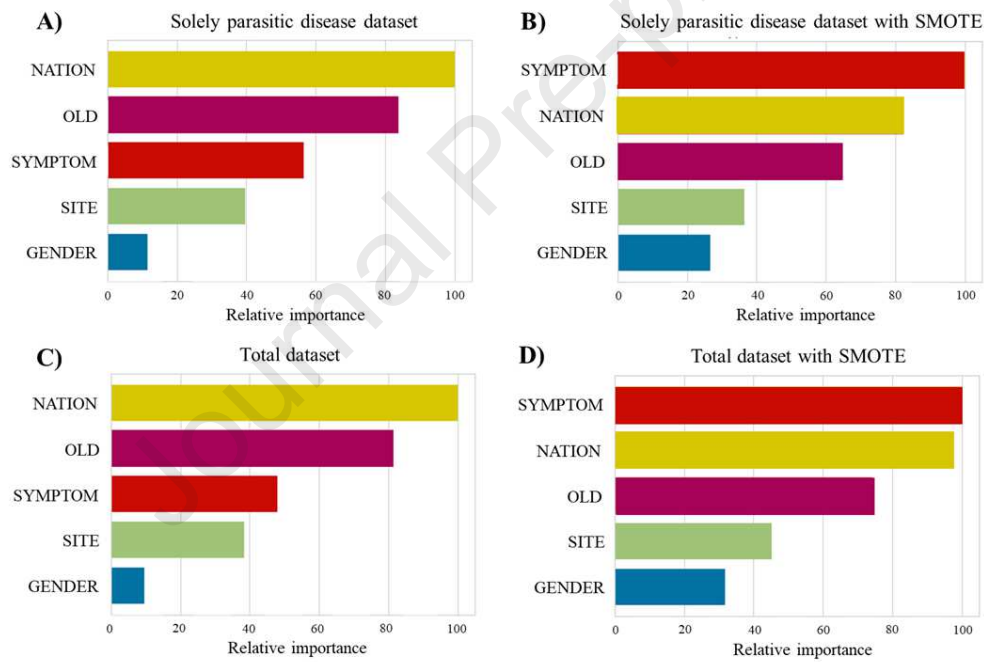
**Fig 3**. Feature importance: A) the solely parasitic disease dataset; B) solely parasitic disease

dataset with SMOTE; C) total dataset; D) total dataset with SMOTE.



498

499

500

501

502

25

503

504

505

506

507

508 **Table 1.** Dataset review

**Table 1. Dataset review**

| | Non - parasitic disease (N=148) | Only parasitic disease(N=1698) | | | Overall(N=1846) |
|---|---|---|---|---|---|
| | | Non - malaria (N=1563) | Malaria (N=135) | Total | |
| **Gender(n)** | | | | | |
| Male | 78 | 881 | 89 | 970 | 1048 |
| Female | 70 | 682 | 46 | 728 | 798 |
| **Age(n)** | | | | | |
| 1~20 | 6 | 334 | 17 | 351 | 357 |
| 21~40 | 24 | 562 | 56 | 618 | 642 |
| 41~60 | 54 | 421 | 50 | 471 | 525 |
| 61~80 | 56 | 233 | 12 | 245 | 301 |
| 81~ | 8 | 13 | 0 | 13 | 21 |
| **Nationality(n)** | | | | | |
| Africa | 16 | 251 | 69 | 320 | 336 |
| America | 16 | 309 | 14 | 323 | 339 |
| Asia | 87 | 591 | 37 | 628 | 715 |
| Europe | 27 | 364 | 13 | 377 | 404 |
| Ocearnia & Caribbean | 2 | 48 | 2 | 50 | 52 |
| **Symptomic body region(n, (%))** | 122(82.4) | 1360(87) | 78(57.8) | 1438(84.7) | 1560(83.7) |
| ABDOMEN | 21 | 476 | 24 | 500 | 521 |
| BACK | 7 | 32 | 7 | 39 | 46 |
| CHEST | 15 | 67 | 4 | 71 | 86 |

26

| | | | | | |
|---|---|---|---|---|---|
| EAR | 1 | 9 | 2 | 11 | 12 |
| EXTREMITIES | 4 | 22 | 1 | 23 | 27 |
| GASTROINTESTINAL | 8 | 68 | 2 | 70 | 78 |
| HAIR | 0 | 9 | 0 | 9 | 9 |
| HEAD | 0 | 28 | 0 | 28 | 28 |
| LYMPH NODE | 8 | 35 | 0 | 35 | 43 |
| MOUTH | 0 | 12 | 1 | 13 | 13 |
| NAIL | 1 | 4 | 0 | 4 | 5 |
| NECK | 3 | 41 | 1 | 42 | 45 |
| NEUROLOGICAL | 5 | 78 | 9 | 87 | 92 |
| OCULAR | 3 | 93 | 1 | 94 | 97 |
| PELVIS | 0 | 5 | 0 | 5 | 5 |
| PSYCHIATRIC | 0 | 7 | 6 | 13 | 13 |
| PULMONARY | 26 | 152 | 13 | 165 | 191 |
| RECTUM | 1 | 8 | 0 | 8 | 9 |
| SKIN | 15 | 164 | 5 | 169 | 184 |
| TOOTH | 1 | 1 | 0 | 1 | 2 |
| VAGINA | 1 | 17 | 0 | 17 | 18 |
| VISION | 2 | 32 | 2 | 34 | 36 |
| **Symptom(n, (%))** | 46(31.1) | 437(27.3) | 86(63.7) | 523(30.8) | 575(31.1) |
| ALOPECIA | 1 | 2 | 0 | 2 | 3 |
| APATHY | 0 | 4 | 0 | 4 | 4 |
| APHASIA | 1 | 6 | 0 | 6 | 7 |
| APNEA | 0 | 0 | 0 | 0 | 1 |
| APRAXIA | 0 | 2 | 0 | 2 | 2 |
| ARRHYTHMIA | 0 | 0 | 1 | 1 | 1 |
| ARTHRALGIA | 5 | 11 | 5 | 16 | 21 |
| ASTHENIA | 0 | 1 | 0 | 1 | 1 |
| ATAXIA | 1 | 12 | 3 | 15 | 16 |
| BACK PAIN | 7 | 18 | 0 | 18 | 25 |
| BLEEDING | 0 | 19 | 2 | 21 | 21 |
| BLINDNESS | 0 | 12 | 1 | 13 | 13 |
| BLURRED VISION | 0 | 6 | 0 | 6 | 6 |
| CHILLS | 1 | 8 | 7 | 15 | 16 |
| CHRONIC PAIN | 0 | 0 | 0 | 0 | 1 |
| CONFUSION | 2 | 7 | 2 | 9 | 11 |
| DEFORMITY | 0 | 4 | 0 | 4 | 4 |
| DEPRESSION | 0 | 3 | 1 | 4 | 4 |
| DISCHARGE | 2 | 15 | 1 | 16 | 18 |

27

| | | | | |
|---|---|---|---|---|
| DIZZINESS | 0 | 5 | 0 | 5 | 5 |
| DYSARTHRIA | 0 | 0 | 0 | 0 | 2 |
| FECAL INCONTINENCE | 0 | 0 | 0 | 0 | 1 |
| FEVER | 14 | 131 | 43 | 174 | 188 |
| HALLUCINATION | 0 | 2 | 0 | 2 | 2 |
| HEARING LOSS | 0 | 2 | 2 | 4 | 4 |
| HEARTBURN | 0 | 3 | 0 | 3 | 3 |
| HEMATEMESIS | 0 | 2 | 0 | 2 | 2 |
| INFERTILITY | 0 | 3 | 0 | 3 | 3 |
| IRRITABILITY | 0 | 1 | 1 | 2 | 2 |
| LACERATION | 0 | 2 | 0 | 2 | 2 |
| LHERMITTE'S SIGN | 0 | 1 | 0 | 1 | 1 |
| LOSS OF CONSCIOUSNESS | 0 | 13 | 1 | 14 | 14 |
| MALAISE | 0 | 10 | 5 | 15 | 15 |
| MYOCLONUS | 0 | 0 | 0 | 0 | 1 |
| NECK STIFFNESS | 0 | 3 | 1 | 4 | 4 |
| PARALYSIS | 0 | 3 | 0 | 3 | 3 |
| PARESIS | 0 | 5 | 0 | 5 | 5 |
| PELVIC PAIN | 0 | 7 | 0 | 7 | 7 |
| PETECHIA | 0 | 3 | 0 | 3 | 3 |
| PURPURA | 0 | 3 | 1 | 4 | 4 |
| RASH | 0 | 1 | 0 | 1 | 1 |
| SHIVERING | 0 | 3 | 4 | 7 | 7 |
| SHORT OF BREATH | 0 | 4 | 0 | 4 | 4 |
| SORE THROAT | 0 | 2 | 0 | 2 | 2 |
| SUICIDAL IDEATION | 0 | 0 | 3 | 3 | 3 |
| SWEATS | 1 | 6 | 1 | 7 | 8 |
| SWELLING | 6 | 71 | 0 | 71 | 77 |
| TINGLING | 0 | 2 | 0 | 2 | 2 |
| TREMOR | 3 | 0 | 1 | 1 | 4 |
| TRISMUS | 0 | 1 | 0 | 1 | 1 |
| URINARY RETENTION | 1 | 10 | 0 | 10 | 11 |
| VAGINAL DISCHARGE | 1 | 4 | 0 | 4 | 5 |
| VOMIT | 0 | 4 | 0 | 4 | 4 |

509

510

511

512

513

514

515  **Table 2.** Solely parasite dataset

Table 2. Solely parasitic disease dataset

| Model | Accuracy | Precision | Recall | F1-Score | CV-10 | AUC |
|---|---|---|---|---|---|---|
| SVM | 0.915 | 0.000 | 0.000 | 0.000 | 0.914 | 0.728 |
| RF | 0.903 | 0.250 | 0.160 | 0.195 | 0.906 | 0.732 |
| MLP | 0.909 | 0.286 | 0.160 | 0.205 | 0.916 | 0.686 |
| Ada | 0.894 | 0.211 | 0.160 | 0.182 | 0.905 | 0.596 |
| GB | 0.891 | 0.227 | 0.200 | 0.213 | 0.913 | 0.708 |
| CB | 0.909 | 0.250 | 0.120 | 0.162 | 0.919 | 0.685 |
| SMOTE+SVM | 0.874 | 0.091 | 0.080 | 0.085 | 0.914 | 0.735 |
| SMOTE+RF | 0.871 | 0.120 | 0.120 | 0.120 | 0.906 | 0.740 |
| SMOTE+MLP | 0.721 | 0.150 | 0.600 | 0.240 | 0.916 | 0.702 |
| SMOTE+Ada | 0.865 | 0.161 | 0.200 | 0.179 | 0.915 | 0.688 |
| SMOTE+GB | 0.885 | 0.208 | 0.200 | 0.204 | 0.912 | 0.698 |
| SMOTE+CB | 0.871 | 0.194 | 0.240 | 0.214 | 0.919 | 0.708 |

516

517  **Table 3**. Total dataset

Table 3. Total dataset

| Model | Accuracy | Precision | Recall | F1-Score | CV-10 | AUC |
|---|---|---|---|---|---|---|
| SVM | 0.938 | 0.000 | 0.000 | 0.000 | 0.921 | 0.776 |
| RF | 0.930 | 0.300 | 0.136 | 0.187 | 0.917 | 0.837 |
| MLP | 0.914 | 0.222 | 0.182 | 0.200 | 0.917 | 0.789 |
| Ada | 0.932 | 0.333 | 0.136 | 0.194 | 0.910 | 0.835 |
| GB | 0.930 | 0.300 | 0.136 | 0.187 | 0.908 | 0.856 |
| CB | 0.949 | 0.714 | 0.227 | 0.345 | 0.924 | 0.802 |
| SMOTE+SVM | 0.919 | 0.278 | 0.227 | 0.250 | 0.921 | 0.771 |

29

| | | | | | | |
|---|---|---|---|---|---|---|
| SMOTE+RF | 0.922 | 0.348 | 0.364 | 0.356 | 0.917 | 0.805 |
| SMOTE+MLP | 0.746 | 0.125 | 0.545 | 0.203 | 0.917 | 0.674 |
| SMOTE+Ada | 0.922 | 0.360 | 0.409 | 0.383 | 0.907 | 0.803 |
| SMOTE+GB | 0.927 | 0.400 | 0.455 | 0.426 | 0.907 | 0.804 |
| SMOTE+CB | 0.881 | 0.211 | 0.364 | 0.267 | 0.924 | 0.714 |

518

519

## Appendix A. Supplementary Table 1. Query list of parasitic disease

| Parasitic disease list | Query list |
|---|---|
| Chagas disease | ((Trypanosoma OR Chagas Disease OR trypanosomiasis) AND America) AND case report |
| Sleeping sickness | ((Trypanosoma OR Sleeping Sickness OR trypanosomiasis) AND Africa) AND case report |
| Acanthamoeba Infection | (Acantamoeba OR Granulomatous Amebic Encephalitis) AND case report |
| Angiostrongyliasis | (Angiostrongylus OR Angiostrongyliasis) AND case report |
| Anisakiasis | (Anisakis OR Anisakiasis OR Pseudoterranova) AND case report |
| Ascariasis | (Ascaris OR Ascariasis OR Intestinal Roundworms) AND case report |
| Babesiosis | (Babesia OR Babesiosis) AND case report |
| Balantidiasis | (Balantidium OR Balantidiasis) AND case report |
| Baylisascariasis | (Baylisascaris OR Baylisascariasis OR Raccoon Roundworm) AND case report |
| Bed Bugs | (Bed Bugs OR Tropical bedbug OR Triatomid bug OR Cimex OR Panstrongylus megistus OR Rhodnius OR Triatoma protracta) AND case report |
| Capillariasis | (Capillaria OR Capillariasis) AND case report |
| Cercarial Dermatitis | (Cercaria OR Cercarial Dermatitis OR Swimmer's Itch) AND case report |
| Clonorchiasis | (Clonorchis OR Clonorchiasis )AND case report |

| | |
|---|---|
| Cryptosporidiosis | (Cryptosporidiosis OR Cryptosporidium) AND case report |
| Cyclosporiasis | (Cyclospora OR Cyclosporiasis) AND case report |
| Cysticercosis | (Cysticercosis OR Neurocysticercosis OR Taenia OR Cysticercus) AND case report |
| Cystoisosporiasis | (Cystoisosporiasis OR Isospora OR Cystoisospora) AND case report |
| Dientamoeba fragilis Infection | (Dientamoeba fragilis) AND case report |
| Dipylidium caninum Infection | (Diphyllobothriasis OR Diphyllobothrium OR tapeworm) AND case report |
| Dirofilariasis | (Dirofilaria OR Dirofilariasis) AND case report |
| Echinococcosis, Hydatid Disease | (Echinococcus OR Echinococcosis OR Hydatid Disease OR Hydatidosis) AND case report |
| Nonpathogenic Intestinal Protozoa | (Endolimax nana OR Entamoeba coli OR Entamoeba dispar OR Entamoeba hartmanni OR Entamoeba polecki OR Nonpathogenic Intestinal Protozoa OR Harmless Intestinal Protozoa OR Iodamoeba buetschlii OR Entamoeba gingivalis) AND case report |
| Amoebiasis | (Entamoeba OR Amebiasis) AND case report |
| Enterobiasis | (Enterobiasis OR Pinworm OR Enterobius) AND case report |
| Fascioliasis | (Fasciola OR Fascioliasis) AND case report |
| Fasciolopsiasis | (Fasciolopsiasis OR Fasciolopsis) AND case report |
| Filariasis | (Filaria OR Filariasis OR Elephantiasis OR Wuchereria bancrofti OR Brugia) AND case report |
| Myiasis | (fly OR Myiasis OR Dermatobia hominis OR bot fly OR Cochliomyia hominovorax OR screwworm fly OR Chrysomya bezziana OR screwworm OR Cordylobia anthropophaga OR tumbu fly OR Cuterebra OR Oestrus OR Wohlfahrtia ) AND case report |
| Giardiasis | (Giardia OR Giardiasis) AND case report |
| Gnathostomiasis | (Gnathostoma OR Gnathostomiasis) AND case report |
| Dracunculiasis | (Guinea OR Dracunculiasis OR Dracunculus medinensis) AND case report |

| Heterophyiasis | (Heterophyes OR Heterophyiasis) AND case report |
|---|---|
| Ancylostomiasis/Hookworm | (Hook worm OR Ancylostomiasis OR Cutaneous larva migrans OR Ancylostoma OR Necator americanus) AND case report |
| Hymenolepiasis | (Hymenolepis OR Hymenolepiasis) AND case report |
| Leishmaniasis | (Leishmania OR Leishmaniasis OR Kala-azar) AND case report |
| Loiasis | (Loa loa OR Loiasis) AND case report |
| Lice Infestation | (Louse OR Body Lice OR Pediculosis OR Pthiriasis OR Pubic lice OR Pubic crab lice OR Head Lice OR Phthirus pubis OR Pediculus humanus corporis) AND case report |
| Malaria | (Malaria OR Plasmodium) AND case report |
| Microsporidiosis | (Microsporidiosis OR Microsporidia OR Anncaliia algerae OR Anncaliia connori OR Anncaliia vesicularum OR Brachiola OR Anncaliia OR Encephalitozoon cuniculi OR Encephalitozoon hellem OR Encephalitozoon intestinalis OR Septata intestinalis OR Tubulinosema acridophagus OR Enterocytozoon bieneusi OR Microsporidium ceylonensis OR Microsporidium africanum OR Nosema ocularum OR Pleistophora OR Trachipleistophora hominis OR Trachipleistophora anthropophthera OR Vittaforma corneae OR Tubulinosema acridophagus) AND case report |
| Mite Infestation | (mite OR Scabies OR Sarcoptes) AND case report |
| Naegleria Infection | (Naegleria OR brain eating amoeba OR primary amebic meningoencephalitis ) AND case report |
| Onchocerciasis | (Onchocerciasis OR Onchocerca OR River Blindness ) AND case report |
| Opisthorchiasis | (Opisthorchis OR Opisthorchiasis) AND case report |
| Paragonimiasis | (Paragonimus OR Paragonimiasis OR flatworm OR ) AND case report |
| Sappinia | (Sappinia OR amebic encephalitis) AND case report |
| Sarcocystosis | (Sarcocystosis OR Sarcocystis ) AND case report |
| Schistosomiasis | (Schistosomiasis OR Schistosoma OR Bilharzia) AND case report |

| | |
|---|---|
| Strongyloidiasis | (Strongyloidiasis OR Strongyloides) AND case report |
| Taeniasis | (Taenia OR Taeniasis OR Tapeworm OR cysticercosis ) AND case report |
| Toxocariasis | (Toxocara OR Ocular Larva Migrans OR Toxocariasis OR Roundworm OR Visceral Larva Migrans) AND case report |
| Toxoplasmosis | (Toxoplasma gondii OR Toxoplasmosis OR Toxoplasma) AND case report |
| Trichinosis | (Trichinella OR Trichinellosis OR Trichinosis) AND case report |
| Trichomoniasis | (Trichomonas OR Trichomoniasis) AND case report |
| Trichuriasis | (Trichuris OR Whipworm OR Trichuriasis ) AND case report |
| Balamuthia | Balamuthia AND case report |
| Chilomastix mesnili Infection | Chilomastix mesnili AND case report |

520

**Highlights**

- A machine learning model is proposed to predict malaria using patient information from parasite case reports.

- Feature importance analysis indicates that the nationality and region of travel are important factors to diagnose malaria.

- SMOTE application does not achieve any considerable improvement

**Conflict of interest**

The authors have no competing interests to declare.