



## OPEN ACCESS

## EDITED BY

Sebastian Vernal,  
University of São Paulo, Brazil

## REVIEWED BY

Luciana Silva Rodrigues,  
Rio de Janeiro State University, Brazil  
Maria Pena,  
Health Resources and Services Administration,  
United States

## \*CORRESPONDENCE

Marcelo Távora Mira  
✉ m.mira@pucpr.br  
Rafael Saraiva de Andrade Rodrigues  
✉ saraiva\_1988@hotmail.com

<sup>†</sup>These authors have contributed equally to this work

RECEIVED 01 June 2023

ACCEPTED 07 July 2023

PUBLISHED 26 July 2023

## CITATION

de Andrade Rodrigues RS, Heise EFJ, Hartmann LF, Rocha GE, Olandoski M, de Araújo Stefani MM, Latini ACP, Soares CT, Belone A, Rosa PS, de Andrade Pontes MA, de Sá Gonçalves H, Cruz R, Penna MLF, Carvalho DR, Fava VM, Bühner-Sékula S, Penna GO, Moro CMC, Nievola JC and Mira MT (2023) Prediction of the occurrence of leprosy reactions based on Bayesian networks. *Front. Med.* 10:1233220. doi: 10.3389/fmed.2023.1233220

## COPYRIGHT

© 2023 de Andrade Rodrigues, Heise, Hartmann, Rocha, Olandoski, de Araújo Stefani, Latini, Soares, Belone, Rosa, de Andrade Pontes, de Sá Gonçalves, Cruz, Penna, Carvalho, Fava, Bühner-Sékula, Penna, Moro, Nievola and Mira. This is an open-access article distributed under the terms of the [Creative Commons Attribution License \(CC BY\)](https://creativecommons.org/licenses/by/4.0/). The use, distribution or reproduction in other forums is permitted, provided the original author(s) and the copyright owner(s) are credited and that the original publication in this journal is cited, in accordance with accepted academic practice. No use, distribution or reproduction is permitted which does not comply with these terms.

# Prediction of the occurrence of leprosy reactions based on Bayesian networks

Rafael Saraiva de Andrade Rodrigues<sup>1\*†</sup>,  
Eduardo Ferreira José Heise<sup>1†</sup>, Luis Felipe Hartmann<sup>2</sup>,  
Guilherme Eduardo Rocha<sup>2</sup>, Marcia Olandoski<sup>1</sup>,  
Mariane Martins de Araújo Stefani<sup>3</sup>, Ana Carla Pereira Latini<sup>4</sup>,  
Cleverson Teixeira Soares<sup>4</sup>, Andrea Belone<sup>4</sup>,  
Patrícia Sammarco Rosa<sup>4</sup>, Maria Araci de Andrade Pontes<sup>5</sup>,  
Heitor de Sá Gonçalves<sup>5</sup>, Rossilene Cruz<sup>6</sup>,  
Maria Lúcia Fernandes Penna<sup>7</sup>, Deborah Ribeiro Carvalho<sup>2</sup>,  
Vinicius Medeiros Fava<sup>8</sup>, Samira Bühner-Sékula<sup>3</sup>,  
Gerson Oliveira Penna<sup>9</sup>, Claudia Maria Cabral Moro<sup>2</sup>,  
Julio Cesar Nievola<sup>10</sup> and Marcelo Távora Mira<sup>1,11\*</sup>

<sup>1</sup>School of Medicine and Life Sciences, Graduate Program in Health Sciences, Pontifícia Universidade Católica do Paraná – PUCPR, Curitiba, Paraná, Brazil, <sup>2</sup>Graduate Program in Health Technology, PUCPR, Curitiba, Paraná, Brazil, <sup>3</sup>Tropical Pathology and Public Health Institute, Federal University of Goiás, Goiania, Brazil, <sup>4</sup>Instituto Lauro de Souza Lima, Bauru, São Paulo, Brazil, <sup>5</sup>Dona Libânia Dermatology Centre, Ceará, Brazil, <sup>6</sup>Tropical Dermatology and Venerology Alfredo da Matta Foundation, Amazonas, Brazil, <sup>7</sup>Epidemiology and Biostatistics Department, Federal University Fluminense, Rio de Janeiro, Brazil, <sup>8</sup>Program in Infectious Diseases and Immunity in Global Health, Research Institute of the McGill University Health Centre, and The McGill International TB Centre, Departments of Human Genetics and Medicine, McGill University, Montreal, QC, Canada, <sup>9</sup>Tropical Medicine Centre, University of Brasília, and Fiocruz School of Government – Brasília, Brasília, Brazil, <sup>10</sup>Graduate Program in Informatics, PUCPR, Curitiba, Paraná, Brazil, <sup>11</sup>Pharmacy Program, School of Health and Biosciences, PUCPR, Curitiba, Paraná, Brazil

**Introduction:** Leprosy reactions (LR) are severe episodes of intense activation of the host inflammatory response of uncertain etiology, today the leading cause of permanent nerve damage in leprosy patients. Several genetic and non-genetic risk factors for LR have been described; however, there are limited attempts to combine this information to estimate the risk of a leprosy patient developing LR. Here we present an artificial intelligence (AI)-based system that can assess LR risk using clinical, demographic, and genetic data.

**Methods:** The study includes four datasets from different regions of Brazil, totaling 1,450 leprosy patients followed prospectively for at least 2 years to assess the occurrence of LR. Data mining using WEKA software was performed following a two-step protocol to select the variables included in the AI system, based on Bayesian Networks, and developed using the NETICA software.

**Results:** Analysis of the complete database resulted in a system able to estimate LR risk with 82.7% accuracy, 79.3% sensitivity, and 86.2% specificity. When using only databases for which host genetic information associated with LR was included, the performance increased to 87.7% accuracy, 85.7% sensitivity, and 89.4% specificity.

**Conclusion:** We produced an easy-to-use, online, free-access system that identifies leprosy patients at risk of developing LR. Risk assessment of LR for individual patients may detect candidates for close monitoring, with a potentially

positive impact on the prevention of permanent disabilities, the quality of life of the patients, and upon leprosy control programs.

#### KEYWORDS

leprosy, leprosy reactions, risk, Bayesian networks, artificial intelligence

## 1. Introduction

Leprosy is a chronic, disabling infectious disease caused by *Mycobacterium leprae* (*M. leprae*) (1) that affected 141,000 new individuals worldwide in 2021 – a number likely to be underestimated due to potential sub-notification caused by the COVID-19 pandemic – with most cases concentrated in India and Brazil (2). In the classical Ridley & Jopling (R&J) classification system, tuberculoid (TT) and lepromatous (LL) leprosy occupy opposite ends of a continuous disease spectrum that includes three borderline forms (BT, BB, and BL) (3). The TT + BT and BB + BL + LL cases roughly correspond to paucibacillary (PB) and multibacillary (MB) leprosy, according to the treatment-oriented World Health Organization (WHO) classification scheme, respectively (2, 4, 5). Today, it is widely accepted that exposure to *M. leprae* is necessary but not sufficient for the development of leprosy; different sets of host gene variants mediate susceptibility to leprosy in three different stages (6): (i) controlling infection *per se*, that is, the disease regardless of its clinical presentation, (ii) defining the clinical form of disease after the infection is established, and (iii) outlining the risk of developing leprosy reactions (LR) (7, 8).

Leprosy reactions are characterized by an intense and sudden (re) activation of the host inflammatory response that may be diagnosed concomitantly with leprosy, during or even after treatment (9–12). Upon diagnosis, LR requires immediate medical attention to prevent irreversible nerve damage, motor disability, and permanent anatomical deformities. In 2021, 6.04% of newly detected leprosy cases worldwide presented grade-2 disabilities in the diagnosis (2), often due to LR. Cohort studies estimate that, during leprosy, 16 to 56% of the patients will develop irreversible nerve damage, again, mainly due to reactional episodes (13–16). Over the past years, advances in genetic research improved our understanding of the molecular basis of leprosy pathogenesis, and several host genetic variations have been implicated in the control of LR episodes (17–19).

There are two major types of LR of distinct clinical presentation: type-1 (T1R) and type-2 reaction (T2R). T1R affects 10–30% of leprosy patients and occurs primarily within, but not limited to, the first 2 years after leprosy diagnosis (20, 21). Known risk factors for T1R are (i) borderline clinical groups BT-BL (22); (ii) age of leprosy onset, with older individuals being at higher risk (23, 24); (iii) positive bacillary index (25); (iv) an increased number of lesions at leprosy diagnosis (26, 27); (v) detection of *M. leprae* DNA in biopsies of lesions (24); and (vi) genetic/genomic studies have identified an association between T1R and genes *TLR1* (28), *TLR2* (29), *TLR3* (30), *TLR7* (30), *TLR10* (30), *NRAMP1/SCLC11A1* (31), *VRD* (32), *NOD2* (33), *TNFSF15/TNFSF8* (34, 35), lncRNA *ENSG00000235140* (36), *LRRK2* (19), and *PRKN* (19).

Leprosy T2R mainly affects patients classified within the BB-LL range (13, 37). Patients presenting bacterial index higher than 4+ in skin smears are at increased risk for T2R (38, 39). There is a wide

variation in the prevalence of T2R in different geographic and endemic settings. In Brazil, approximately 37% of BL and LL cases develop T2R, while in India, Nepal, and Thailand, the proportion is between 19–26% (40). A prospective study involving BL and LL patients from India followed for 11 years, showed that less than 10% of the individuals who developed T2R had a single episode, whereas 62% had chronic T2R (21). In Ethiopia, 63% of leprosy cases had more than one T2R episode, while 37% had a single event (41). Host genetics also seems to play a significant role in controlling the occurrence of T2R, and genes *C4B* (42), *TLR1* (43), *NRAMP1/SCLC11A1* (31), *NOD2* (33, 35), and *IL6* (12, 35) have been implicated as critical molecular players.

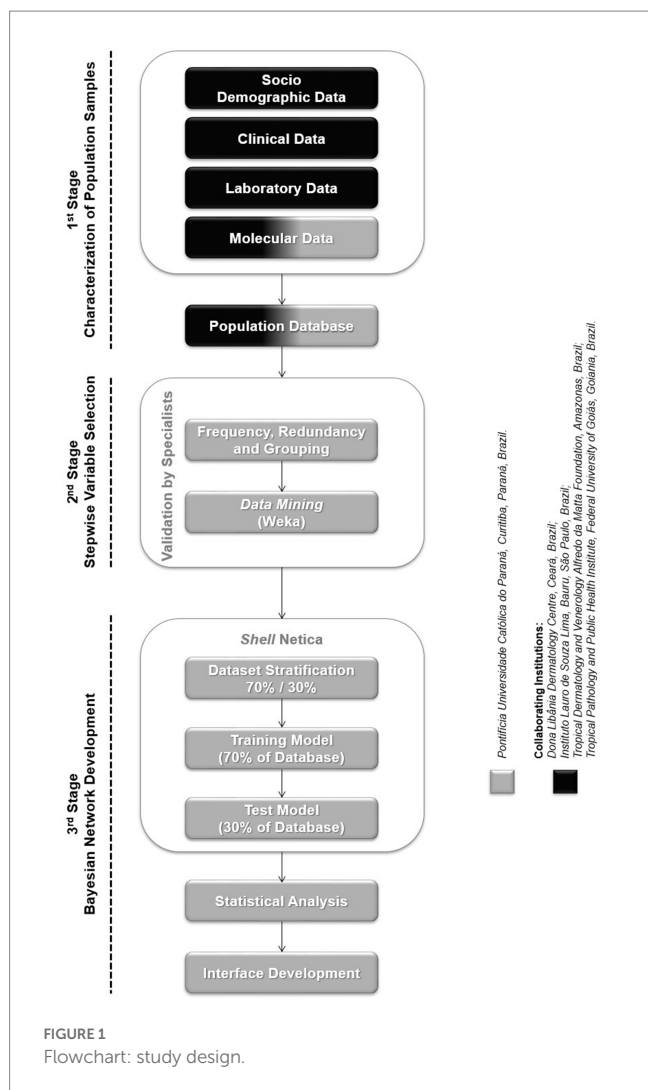
One of the challenges of translational medicine is to systematize the analysis of a large amount of patient data to predict a specific outcome. In addition, scientific results from basic research are often difficult to translate into daily medical practice. Artificial Intelligence (AI) methods seek to systematically address often large, complex data sets to provide a base for decision-making. Of particular interest in health care, Bayesian Networks (BN) are among the most successful techniques in processing and unraveling the relationship between a large number of variables, with risk estimation being the outcome (44).

A Bayesian Network (BN) is a graphical model of an outcome variable's posterior conditional probability distribution based on evidence. It contains nodes that represent the random variables and links between pairs of nodes, which represent the causal relationship of these nodes, together with a conditional probability distribution in each node. From the definition, one can deduce that any joint probability distribution may be represented by a Bayesian network, which shows its modeling power: any deterministic model is a particular case of a probabilistic model, and any probabilistic model may be represented as a Bayesian network (45, 46).

Several BN-based systems have been created using medical data, developed for different purposes, and applied to several health conditions such as cardiovascular diseases, liver diseases, cancer and Alzheimer's disease (47–57), including leprosy (44, 58–65). However, few initiatives aim to systematize a large amount of existing information of distinct nature to estimate the risk of the occurrence of a particular event. In the context of leprosy, creating a simple-to-use and flexible platform to predict the risk of LR based on patient data may help minimize the consequences of such aggressive events. Moreover, such a tool could improve leprosy control initiatives and public health systems. Here we present an AI system designed to predict the risk of a leprosy patient to develop LR using a complete or partial dataset of clinical, demographic, and host genetic data.

## 2. Materials and methods

A flowchart summarizing the three stages of the study and the procedures described next is available in Figure 1.



## 2.1. Population samples

This study used four pre-existing data sets from previous research initiatives of different/independent designs and contexts. The first database included in the study consisted of 409 leprosy patients diagnosed at the Reference Center for Diagnosis and Therapy located in Goiania, central-western Brazil, between February 2006 and March 2008, originally used for the genetic study that identified an association between T2R and variants of the *IL6* gene. A complete description of the Goiania population has been published elsewhere (12). Later, the Goiania population was used for an expanded investigation involving a larger number of candidate genes that detected an association between T1R and variants of the gene *TNFSF8* (34). Finally, an association between T1R and lncRNA *ENSG00000235140* (36) and *LRRK2* (unpublished data) was also detected in the Goiania sample. Two additional databases comprised 533 patients recruited at the Dermatological Center Dona Libânia, Fortaleza, northeast Brazil, and 137 patients diagnosed with leprosy at the Fundação Alfredo da Matta, Manaus, north Brazil. Enrolment of these two population samples was performed under a single protocol of a clinical study described previously (66) and conducted by the Tropical Medicine Center of the University of Brasília between

March 2007 and February 2012. Finally, a fourth database consisted of 371 patients diagnosed with leprosy at the Instituto Lauro de Souza Lima, Bauru, southeast Brazil, between March 2008 and January 2013, originally for a genetic study that detected an association between leprosy and variants of the *TLR1* (67) and *NOD2* (68) genes. For all databases, leprosy diagnosis/classification was defined after detailed dermatological and neurological examination by specialized leprologists, complemented by bacilloscopy and histopathology of skin lesions. All cases were classified following the R&J scheme (3). Patients were followed up for at least 2 years since diagnosis to monitor LR occurrence. Individuals who did not present LR at the initial diagnosis or during follow-up, were defined as non-reactional leprosy patients.

All patients were treated for leprosy according to WHO/MDT guidelines and for LR with the appropriate therapy. All subjects were evaluated for an extensive clinical, socioeconomic, and demographic information list.

## 2.2. Variable selection

The four databases included in this study were composed of clinical and laboratory parameters, most of them obtained for descriptive, epidemiological purposes unrelated to the occurrence of LR. Each one of the databases was subjected individually to a two-step, unbiased process aiming to identify those variables exerting the highest impact upon the risk of LR, thus, to be included in the system, as follows:

### 2.2.1. Frequency, redundancy, and grouping

The first selection step consisted of removing variables with low frequency (less than 15%) of occurrence and/or mutually correlated (redundant), consequently capturing the same information. In the case of redundant variables, the most frequent was selected to capture the information of the set.

### 2.2.2. Data mining

Data mining is one of the main stages of the knowledge extraction process from large databases, also known as KDD – Knowledge Discovery in the Databases (69). This AI method is defined as the process of discovering patterns in data to generate helpful information for the decision-making (70). WEKA (Waikato Environment for Knowledge Analysis) is an open-source program with a collection of algorithm implementations of various data mining techniques, such as pre-processing, classification, clustering, and visualization (71). This study used WEKA in the second variable selection step to identify those hierarchically important for LR occurrence in the population samples. The variables were selected using the C4.5 algorithm, which creates a decision tree and identifies the most relevant and non-redundant variants, thus reducing the number of attributes. The C4.5 selection is made according to the gain ratio, which is a normalization of the information gain, a parameter based on the entropy measure (originating from information theory) closely related to the maximum likelihood estimations (MLE) and usually used to make inferences about parameters of the underlying probability distribution from a given dataset (45, 72–75).

Four dermatologists/hansenologists with extensive experience in the area continuously validated the two-step variable selection

through a qualitative process based on their experience in the field of leprosy diagnosis. These specialists were also involved in conducting system performance assessments, evaluating usability, and organizing the workflow for integrating data from the four databases. By leveraging the knowledge and expertise of specialists, clinical decision systems can be effectively validated and optimized for real-world clinical use (76). Criteria for selecting the specialists were; (i) holding MD/Ph.D. degrees in dermatology/hansenology; (ii) having more than 10 years of experience in leprosy diagnosis; (iii) being representative of regions of Brazil with different levels of leprosy endemicity.

Finally, two datasets contributed with genotypic information: Goiania for genes *IL6*, *TNFSF8*, *LRRK2*, and *ENSG00000235140* and Bauru for *TLR1* and *NOD2*, all previously studied in these population samples.

### 2.3. System development

The system was created as a BN using Shell NETICA (Norsys Software Corporation) (77) with a customized dynamic interface considering the number of variables in the database. The system was designed to operate with complete or partial information, which is of critical importance considering the translational bias of the proposal and the fact that several leprosy centers may not have access to all the information included, particularly the molecular genetic data. The system loads a spreadsheet in which columns and lines refer to the variables and records, respectively. Each variable (column) is related to one node of the BN. The variables comprise demographic, clinical, laboratory, and genetic data (markers). For each one of the databases,

two groups were formed randomly to create the network: the test file, with 30% of patients, and the training file, with 70% of patients, both stored in an Excel file format.

The system's performance was assessed by its accuracy, sensitivity, specificity, and negative and positive predictive values. The patient's predicted outcome was defined by the class with higher risk, as estimated by the system. Predictive values were calculated using the prevalence of occurrence of reversal reactions observed for the studied population samples. The feature "importance" was also measured using the  $F_1$  score, which is the harmonic mean between positive predictive value (PPV) and sensitivity. The  $F_1$  score was calculated

accordingly to the equation  $F_1 \text{ score} = 2 * \frac{PPV * sensitivity}{PPV + sensitivity}$  using Python 3.7.9.

### 3. Results

Table 1 summarizes information on age, gender, and clinical form of leprosy according to the R&J classification system for T1R, T2R and non-reactional leprosy patients groups of all population samples. The mean age at diagnosis ranged from 40 to 59 years old, and males were consistently more frequent than females across all four population samples. Leprosy clinical form most frequently observed was BT (479, 33%) followed by BL (379, 26.1%), LL (346, 23.8%), BB (134, 9.2%), and TT (100, 6.9%). For the combined sample, 51% were non-reactional leprosy patients, 25.9 and 23.1% developed T1R or T2R, respectively. As expected, T1R was observed more often in BT + BB + BL cases, and T2R occurred more often in BL and LL individuals (Table 1).

TABLE 1 Distribution of sex, age at diagnosis, and clinical type of disease of leprosy-affected individuals with T1R, T2R, and non-reactional leprosy patients in each population sample.

	Patients, No. (%)														
	Goiania			Fortaleza			Manaus			Bauru			Combined		
Age, Years (Mean ± SD)	44.63 ± 16.67			45.15 ± 14.25			40.00 ± 15.39			59.00 ± 18.04			48.00 ± 17.29		
Sex															
Male	234 (57.1)			352 (66.0)			100 (72.9)			258 (69.5)			944 (65.1)		
Female	175 (42.9)			181 (34.0)			37 (27.1)			113 (30.5)			506 (34.9)		
Ridley&Jopling Classification	NRLP	T1R	T2R	NRLP	T1R	T2R	NRLP	T1R	T2R	NRLP	T1R	T2R	NRLP	T1R	T2R
TT	22	0	0	28	0	0	16	0	0	34	0	0	100	0	0
BT	124	79	0	164	24	0	36	4	0	18	30	0	342	137	0
BB	16	29	3	12	14	0	2	3	0	27	27	1	57	73	4
BL	26	46	8	47	71	66	12	28	10	12	20	33	97	165	117
LL	28	0	28	33	0	68	5	0	16	66	0	102	132	0	214
I	0	0	0	6	0	0	5	0	0	1	0	0	12	0	0
HI (Mean)	-	-	-	-	-	-	-	-	-	1.73	2.69	3.84	1.73	2.69	3.84
Proportion per Group	52.9	37.6	9.5	54.4	20.5	25.1	55.5	25.5	19.0	42.6	20.8	36.6	51.0	25.9	23.1
Total	409			533			137			371			1,450		

BB, borderline borderline; BL, borderline lepromatous; BT, borderline tuberculoid; HI, histological index; I, indeterminate leprosy; LL, lepromatous leprosy; NRLP, non-reactional leprosy patients; SD, standard deviation; TT, tuberculoid leprosy; T1R, type-1 reaction; T2R, type-2 reaction.



Our strategy for variable selection led to the inclusion of 34 demographic, clinical, laboratory, and genetic parameters (Supplementary Table S1) related to the occurrence of LR in the population samples (Table 2). Since the initial set of variables was not the same across the four datasets – thus, the variables selected by the two-step process and validated by the specialists were not necessarily the same – the prediction system was designed to include all variables selected in each population sample. Detailed information about the distribution of the included variables across the four different datasets is available in Supplementary Table S2.

The risk-prediction system was developed to allow the use of each of the four databases individually as a reference, as well as to use a single, combined dataset, thus favoring customization and facilitating the inclusion of new data sets. The system – named SEPAREH (from Portuguese: *Sistema Especialista Para Avaliação de Risco de Estado Reacional em Hanseníase*; in English: Specialist System for Evaluation of Risk of Occurrence of Reactional States in Leprosy) is designed to present a friendly graphical user interface (Figure 2), which allows the primary care professional to use it intuitively. Variation of the patient's risk of developing one of the two types of LR is shown in real time, as each available clinical and/or

genetic information is included in the interface. The platform can be accessed for free at <https://orfeu.ppgia.pucpr.br/separeh>.<sup>1</sup>

The overall sensitivity and specificity of the system, as estimated using the combined dataset of 1,450 patients, was 79.3% (95% CI 73.9–84.7) and 86.2% (95% CI 81.6–90.8), respectively. Accuracy reached 82.7% (95% CI 79.2–86.3), and positive and negative predicted values were 85.1% (95% CI 80.2–90.1) and 80.6% (95% CI 75.5–85.7), respectively.

To assess the importance of each of the variables individually, modeling was carried out after removing one at a time, and the impact on system performance was measured through changes in sensitivity, specificity, and F1. As summarized in Figure 3, the three attributes exerting the highest impact were R&J classification, combined genetic markers, and histological index. Interestingly, the highest estimates of accuracy, sensitivity, specificity, and both negative and positive predictive values were observed for the Bauru and the Goiania datasets, for which genotypic data was available, even higher than what was observed for the combined dataset of much larger sample size (the only exception being the positive predictive value for Bauru: 82.7% vs. 85.1% for the combined dataset) (Table 3).

## 4. Discussion

As an outcome of contact with its causative agent, leprosy is controlled by multiple environmental and socioeconomic factors and innate characteristics of both the host and pathogen. The specific contribution of each of these factors to the risk of developing leprosy and its endophenotypes is widely unknown. Today, LRs constitute a significant cause of disabilities associated with leprosy; thus, predicting patients at higher risk of developing LR at the time of leprosy diagnosis may help prevent permanent neural impairment. However, an accurate estimate of this risk demands analyzing a very complex set of variables, which is difficult – if not impossible – to perform by an unassisted primary healthcare professional. Here we present an easy-to-use, flexible, and automated system that identifies leprosy patients at increased risk of developing LR based on clinical, socio-economical, laboratory, and genetic data. Patients at high risk are candidates for close monitoring during and after treatment, aiming to prompt the management of these aggressive events, minimizing the likelihood of permanent disabilities. Our platform translates basic scientific data into a direct application that may immediately impact leprosy patients' quality of life and leprosy control programs' effectiveness.

The three features that exerted the highest impact on the system's performance were the R&J classification, the histological index, and the combined effect of the genetic markers (Figure 3). The R&J class is a well-accepted major risk factor for reversal reactions (7, 13, 21, 22, 37, 40, 41). As expected, simulations confirm that patients in the tuberculoid pole of the spectrum tend to have a higher chance of developing no reversal reaction (98% ~ when the classification is TT). As clinical form moves towards borderline, the probability of a T1R rises from <1 to 53% ~ when the category is BB and, finally, patients

TABLE 2 Demographical, clinical, laboratory, and genetic variables selected in the study.

Data	Variable information <sup>a</sup>
Socio-demographic	Sex
	Age group
	Ethnicity
Clinical	Multidrug therapy
	First signs and symptoms
	Ridley-Jopling classification
	Number of skin lesions
	Type of lesion
	Color of lesion
	Sensibility testing
Laboratory	Bacilloscopic index
	Histological index
	PGL-1
Genetic	<i>IL6</i> markers (4)
	<i>NOD2</i> marker (1)
	<i>TLR1</i> markers (2)
	<i>TNFSF8</i> markers (4)
	<i>ENSG00000235140</i> markers (4)
	<i>LRRK2</i> markers (3)
Family History	First degree <sup>b</sup>
	Second degree <sup>c</sup>
	Contact <sup>d</sup>

<sup>a</sup>Self-report in years since noticing the early signs and symptoms of leprosy.

<sup>b</sup>Father, mother, child, and sibs affected by leprosy.

<sup>c</sup>Cousins, nephews, uncles/aunts, grandparents, and grandchildren affected by leprosy.

<sup>d</sup>Close household contact affected by leprosy.

1 The access to the platform is limited to HTTPS protocol. In case of difficulty accessing the platform, please certify whether HTTPS is being used.

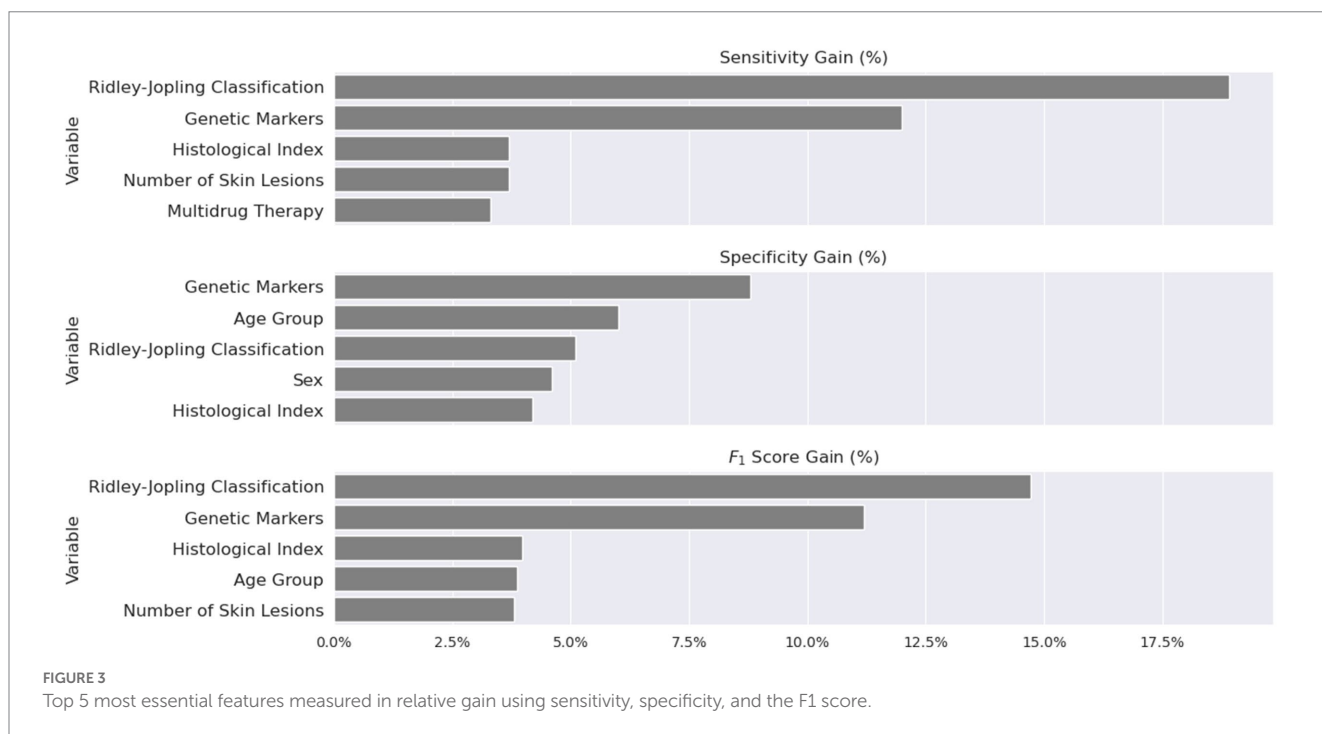
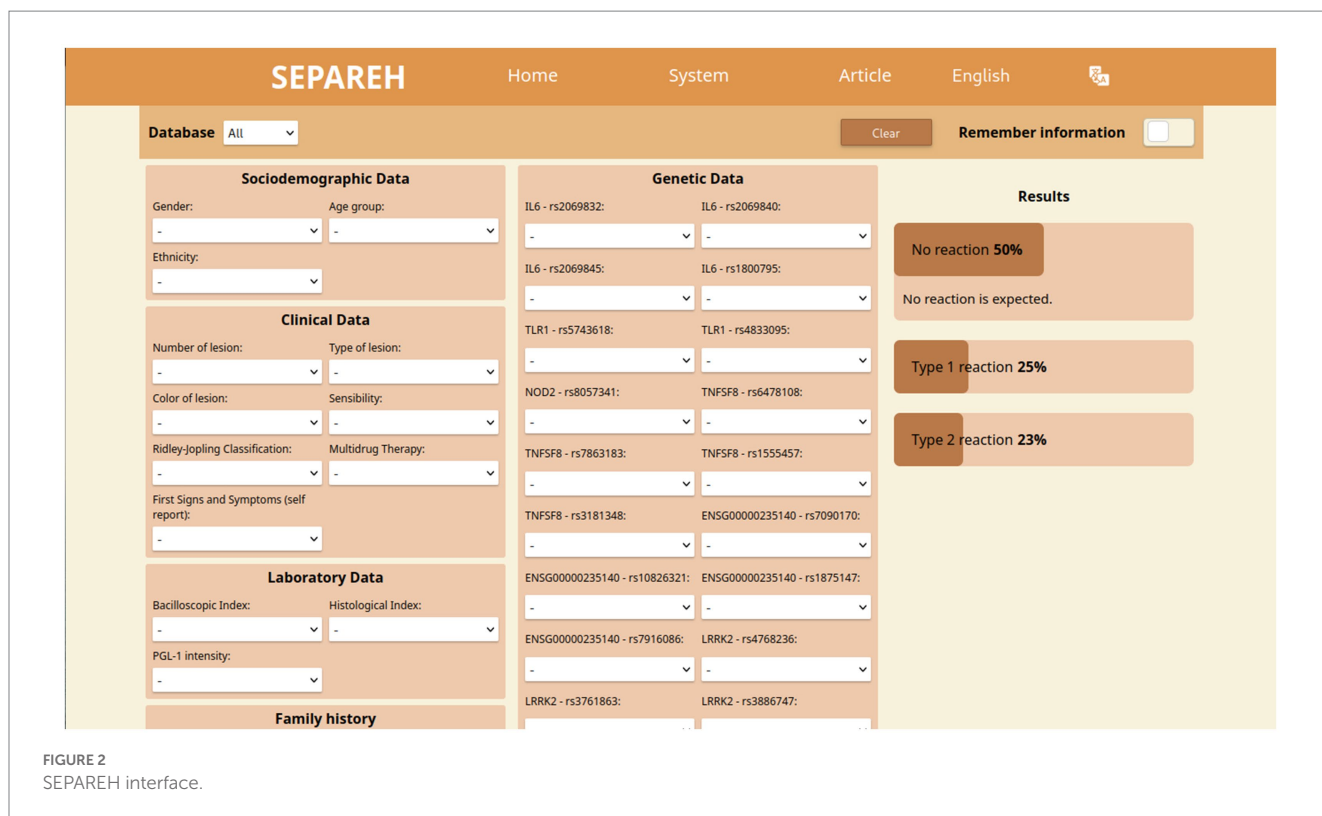


FIGURE 3 Top 5 most essential features measured in relative gain using sensitivity, specificity, and the F1 score.

at the lepromatous pole have a higher risk of developing T2R – more specifically, 61%~ when the type is LL. The second top-three parameter impacting the system is the histological index. An index equal to 2+ increases the risk of T1R to 56%~; values higher than 5+ shift the risk towards T2R – 45%~ when the histological index is 6+. This behavior is expected since an increase in the histological index is highly correlated with a higher bacterial load and, consequently, a

move toward the lepromatous pole of the disease. A histological index higher than 5+ is also a well-known risk factor for developing T2R (38, 39). Finally, genetic data seems critical to improving the system’s performance, which suggests that understanding the true, exact nature of LR depends on the description of the underlying genetic mechanisms.

We are aware of the study's limitations: we have had limited access to genetic information across the population samples; including genotypic data for additional, known LR susceptibility genes would likely positively impact the system's performance. In addition, the heterogeneity of the databases, originally obtained for independent studies of distinct designs, prevented a comprehensive analysis of the performance of the system, which we understand was yet quite remarkable, likely due to the ability of Bayesian methods to estimate risk using all available – even if partial – information. This is important considering that not all leprosy centers across the globe will have access to molecular data of all the patients; in these cases, the platform can still help estimate the risk of LR using only the clinical/laboratory and demographic data with fair sensitivity and specificity, as observed for the Fortaleza and Manaus datasets (Table 3). Of note: the heterogeneity of the dataset is known to enhance the quality of a trained model, since it tends to improve the generalization capturing a more comprehensive understanding of the problem and its nuances. Thus, the inclusion of diverse datasets is a known strategy to improve the performance of machine learning models. For example, in the field of Random Forests, the use of diverse datasets has been explored as a method to enhance the model's accuracy and robustness (78). This principle extends to various domains, including computer vision (79), and conversational

AI (80). For a comprehensive evaluation and refining of the system, datasets enrolled prospectively with these specific purposes will be necessary.

## 5. Conclusion

We produced SEPAREH as an easy-to-use, online, free-access system that identifies leprosy patients at higher risk of developing LR. We believe that SEPAREH can be useful to help primary healthcare services to establish a protocol for patient follow-up dedicated to improving early diagnosis and prevention of the devastating consequences of untreated LR. Ultimately, risk assessment of LR for individual patients may be of potential positive impact on the prevention of permanent disabilities, the quality of life of the patients, and upon leprosy control programs.

## Data availability statement

The raw data supporting the conclusions of this article will be made available by the authors, without undue reservation.

TABLE 3 Results obtained for each population sample.

Population sample	Two-by-two contingency				Results	95% CI
Combined		NRLP	LR	Total	Sensitivity = 79.3%	73.9–84.7%
	NRLP	187	45	232	Specificity = 86.2%	81.6–90.8%
	LR	30	172	202	PVP = 85.1%	80.2–90.1%
	Total	217	217	434	PVN = 80.6%	75.5–85.7%
					Accuracy = 82.7%	79.2–86.3%
Goiania		NRLP	LR	Total	Sensitivity = 85.7%	76.5–94.9%
	NRLP	59	8	67	Specificity = 89.4%	82.0–96.8%
	LR	7	48	55	PVP = 87.3%	78.5–96.1%
	Total	66	56	122	PVN = 88.0%	80.3–95.8%
					Accuracy = 87.7%	81.9–93.5%
Bauru		NRLP	LR	Total	Sensitivity = 82.7%	72.4–93.0%
	NRLP	51	9	60	Specificity = 85.0%	76.0–94.0%
	LR	9	43	52	PVP = 82.7%	72.4–93.0%
	Total	60	52	112	PVN = 85.0%	76.0–94.0%
					Accuracy = 83.9%	77.1–90.7%
Fortaleza		NRLP	LR	Total	Sensitivity = 78.1%	68.6–87.6%
	NRLP	62	16	78	Specificity = 71.3%	61.8–80.8%
	LR	25	57	82	PVP = 69.5%	59.5–79.5%
	Total	87	73	160	PVN = 79.4%	70.5–88.4%
					Accuracy = 74.3%	67.6–81.1%
Manaus		NRLP	LR	Total	Sensitivity = 77.8%	58.6–97.0%
	NRLP	18	4	22	Specificity = 78.3%	61.4–95.1%
	LR	5	14	19	PVP = 73.7%	53.9–93.5%
	Total	23	18	41	PVN = 81.8%	65.7–97.9%
					Accuracy = 78.0%	65.4–90.7%

LR, leprosy reactions; NRLP, non-reactional leprosy patients; PPV, positive predictive value; NPV, negative predictive value; CI, confidence interval.

## Ethics statement

This study was approved by the Brazilian Committee for Ethics in Research (CONEP) (protocol 1.722.447). All patients signed an informed consent to participate in the study; for patients <18 years old, the informed consent was signed by one of the parents or the legal guardian.

## Author contributions

RA, EH, DC, CM, and JN contributed to defining the AI-based protocol and data modeling. LH, GR, and EH developed the online platform. MA contributed to the recruitment and clinical characterization of the Tropical Pathology and Public Health Institute, Goiania, Goiás patients. AL, CS, AB, and PR contributed to the recruitment and clinical description of the Lauro de Souza Lima Institute, Bauru, São Paulo patients. MP and HS contributed to the recruitment and clinical characterization of the Dona Libânia Dermatology Centre patients, Fortaleza, Ceará. RC contributed to the recruitment and clinical description of the Alfredo da Matta Foundation patients, Manaus, Amazonas. MO contributed to the statistical analysis. VF contributed to the generation of the original genetic data. SB-S, GP, and MP contributed to coordinating the original study under which the population samples were recruited and characterized. RA, EH, CM, JN, and MM helped to draft the manuscript. MM is the principal investigator, the main responsible for the study design and execution, and provided senior supervision throughout the study. All authors contributed to the article and approved the submitted version.

## Funding

This work was supported by the Araucaria Foundation (Grant #41617.433.32610.10092013), the Brazilian Coordenação de Aperfeiçoamento de Pessoal de Nível Superior (CAPES), and the Leprosy Research Initiative (LRI)/Turing Foundation, grant ID# 704.16.31. MM is a Conselho Nacional de Desenvolvimento

## References

- Scollard DM, Adams LB, Gillis TP, Krahenbuhl JL, Truman RW, Williams DL. The continuing challenges of leprosy. *Clin Microbiol Rev.* (2006) 19:338–81. doi: 10.1128/CMR.19.2.338-381.2006
- World Health Organization. Guidelines for the diagnosis, treatment and prevention of leprosy. (2018). 1–6. Available at: <https://apps.who.int/iris/handle/10665/274127>
- Ridley DS, Jopling WH. Classification of leprosy according to immunity. A five-group system. *Int J Lepr Mycobact Dis.* (1966) 34:255–73.
- Croft RP, Nicholls PG, Steyerberg EW, Richardus JH, Cairns W, Smith S. A clinical prediction rule for nerve-function impairment in leprosy patients. *Lancet.* (2000) 355:1603–6. doi: 10.1016/S0140-6736(00)02216-9
- World Health Organization. *Guidelines for the diagnosis, treatment and prevention of leprosy.* (2018): 1–106. Available at: <https://apps.who.int/iris/handle/10665/274127>
- Abel L, Desein AJ. The impact of host genetics on susceptibility to human infectious diseases. *Curr Opin Immunol.* (1997) 9:509–16. doi: 10.1016/S0952-7915(97)80103-3
- Sauer ME, Salomão H, Ramos GB, D'Espindula HRS, Rodrigues RSA, Macedo WC, et al. Genetics of leprosy: expected and unexpected developments and perspectives. *Clin Dermatol.* (2015) 33:99–107. doi: 10.1016/j.clindermatol.2014.10.001
- Jacobson RR, Krahenbuhl JL. Leprosy. *Lancet.* (1999) 353:655–60. doi: 10.1016/S0140-6736(98)06322-3
- Alter A, Grant A, Abel L, Alcaïs A, Schurr E. Leprosy as a genetic disease. *Mamm Genome.* (2011) 22:19–31. doi: 10.1007/s00335-010-9287-1
- Sampaio LH, Stefani MMA, Oliveira RM, Sousa ALM, Ireton GC, Reed SG, et al. Immunologically reactive *M. leprae* antigens with relevance to diagnosis and vaccine development. *BMC Infect Dis.* (2011) 11:26. doi: 10.1186/1471-2334-11-26
- Walker SL, Lockwood DN. Leprosy type 1 (reversal) reactions and their management. *Lepr Rev.* (2008) 79:372–86. doi: 10.47276/lr.79.4.372
- Sousa AL, Fava VM, Sampaio LH, Martelli CMT, Costa MB, Mira MT, et al. Genetic and immunological evidence implicates interleukin 6 as a susceptibility gene for leprosy type 2 reaction. *J Infect Dis.* (2012) 205:1417–24. doi: 10.1093/infdis/jis208
- Britton WJ, Lockwood DN. Leprosy. *Lancet.* (2004) 363:1209–19. doi: 10.1016/S0140-6736(04)15952-7
- Reddy BN, Bansal RD. An epidemiological study of leprosy disability in a leprosy endemic rural population of Pondicherry (South India). *Indian J Lepr.* (1984) 56:191–9.
- Girdhar M, et al. Pattern of leprosy disabilities in Gorakhpur (Uttar Pradesh). *Indian J Lepr.* (1989) 61:503–13.
- Zhang G, Li W, Yan L, Yang Z, Chen X, Zheng T, et al. An epidemiological survey of deformities and disabilities among 14,257 cases of leprosy in 11 counties. *Lepr Rev.* (1993) 64:143–9.

Científico e Tecnológico (CNPq) productivity (PQ) researcher level 2, grant #304368/2018-0. MA is under a research fellowship grant from the Brazilian Research Council/CNPq (Grant #311986-2019-6).

## Acknowledgments

We are grateful to the patients and staff of the Reference Center for Diagnosis and Therapy, Goiania; Dermatological Center Dona Libânia, Fortaleza; Alfredo da Matta Foundation, Manaus; Instituto Lauro de Souza Lima, Bauru, for agreeing and their cooperation in this study. Thanks to Cássio Ghidella, MD dermatologist from Rondonia, for the valuable support.

## Conflict of interest

The authors declare that the research was conducted in the absence of any commercial or financial relationships that could be construed as a potential conflict of interest.

## Publisher's note

All claims expressed in this article are solely those of the authors and do not necessarily represent those of their affiliated organizations, or those of the publisher, the editors and the reviewers. Any product that may be evaluated in this article, or claim that may be made by its manufacturer, is not guaranteed or endorsed by the publisher.

## Supplementary material

The Supplementary material for this article can be found online at: <https://www.frontiersin.org/articles/10.3389/fmed.2023.1233220/full#supplementary-material>



17. Fava V, Orlova M, Cobat A, Alcais A, Mira M, Schurr E. Genetics of leprosy reactions: an overview. *Mem Inst Oswaldo Cruz.* (2012) 107:132–42. doi: 10.1590/S0074-02762012000900020
18. Cambri G, Mira MT. Genetic susceptibility to leprosy—from classic immune-related candidate genes to hypothesis-free, whole genome approaches. *Front Immunol.* (2018) 9:1674. doi: 10.3389/fimmu.2018.01674
19. Fava VM, Xu YZ, Lettre G, van Thuc N, Orlova M, Thai VH, et al. Pleiotropic effects for parkin and LRRK2 in leprosy type-1 reactions and Parkinson's disease. *Proc Natl Acad Sci U S A.* (2019) 116:15616–24. doi: 10.1073/pnas.1901805116
20. Scollard DM, Smith T, Bhoopat L, Theetrantong C, Rangdaeng S, Morens DM. Epidemiologic characteristics of leprosy reactions. *Int J Lepr Other Mycobact Dis.* (1994) 62:559–67.
21. Torres O, Suneetha S, Muzaffarullah S, Reddy R, Jain S, Pocaterra L, et al. Clinical course of erythema nodosum leprosum: an 11-year cohort study in Hyderabad, India. *Am J Trop Med Hyg.* (2006) 74:868–79. doi: 10.4269/ajtmh.2006.74.868
22. Van Brakel WH, Khawas IB, Lucas SB. Reactions in leprosy: an epidemiological study of 386 patients in West Nepal. *Lepr Rev.* (1994) 65:190–203. doi: 10.5935/0305-7518.19940019
23. Ranque B, Nguyen VT, Vu HT, Nguyen TH, Nguyen NB, Pham XK, et al. Age is an important risk factor for onset and sequelae of reversal reactions in Vietnamese patients with leprosy. *Clin Infect Dis.* (2007) 44:33–40. doi: 10.1086/509923
24. Costa MB, Martelli CMT, Pereira GAS, Stefani MMA, Narahashi K, Krahenbuhl JL, et al. Mycobacterium leprae DNA associated with type 1 reactions in single lesion paucibacillary leprosy treated with single dose rifampin, ofloxacin, and minocycline. *Am J Trop Med Hyg.* (2007) 77:829–33. doi: 10.4269/ajtmh.2007.77.829
25. Saunderson P, Gebre S, Byass P. Reversal reactions in the skin lesions of AMFES patients: incidence and risk factors. *Lepr Rev.* (2000) 71:309–17. doi: 10.5935/0305-7518.20000034
26. van Brakel WH, Khawas IB. Nerve function impairment in leprosy: an epidemiological and clinical study – part 2: results of steroid treatment. *Lepr Rev.* (1996) 67:104–18. doi: 10.5935/0305-7518.19960011
27. Kumar B, Dogra S, Kaur I. Epidemiological characteristics of leprosy reactions: 15 years experience from North India. *Int J Lepr Other Mycobact Dis.* (2004) 72:125–33. doi: 10.1489/1544-581X(2004)072<0125:ECOLRY>2.0.CO;2
28. Misch EA, Macdonald M, Ranjit C, Sapkota BR, Wells RD, Siddiqui MR, et al. Human TLR1 deficiency is associated with impaired mycobacterial signaling and protection from leprosy reversal reaction. *PLoS Negl Trop Dis.* (2008) 2:e231. doi: 10.1371/journal.pntd.0000231
29. Bochud PY, Hawn TR, Siddiqui MR, Saunderson P, Britton S, Abraham I, et al. Toll-like receptor 2 (TLR2) polymorphisms are associated with reversal reaction in leprosy. *J Infect Dis.* (2008) 197:253–61. doi: 10.1086/524688
30. Régo JL, de Lima Santana N, Machado PRL, Ribeiro-Alves M, de Toledo-Pinto TG, Castellucci LC, et al. Whole blood profiling of leprosy type 1 (reversal) reactions highlights prominence of innate immune response genes. *BMC Infect Dis.* (2018) 18:422. doi: 10.1186/s12879-018-3348-6
31. Teixeira MA, Silva NL, Ramos AL, Hatagima A, Magalhães V. NRAMP1 gene polymorphisms in individuals with leprosy reactions attended at two reference centers in Recife, northeastern Brazil. *Rev Soc Bras Med Trop.* (2010) 43:281–6. doi: 10.1590/S0037-86822010000300014
32. Sapkota BR, Macdonald M, Berrington WR, Misch EA, Ranjit C, Siddiqui MR, et al. Association of TNF, MBL, and VDR polymorphisms with leprosy phenotypes. *Hum Immunol.* (2010) 71:992–8. doi: 10.1016/j.humimm.2010.07.001
33. Berrington WR, Macdonald M, Khadge S, Sapkota BR, Janer M, Hagge DA, et al. Common polymorphisms in the NOD2 gene region are associated with leprosy and its reactive states. *J Infect Dis.* (2010) 201:1422–35. doi: 10.1086/651559
34. Fava VM, Cobat A, van Thuc N, Latini ACP, Stefani MMA, Belone AF, et al. Association of TNFSF8 regulatory variants with excessive inflammatory responses but not leprosy per se. *J Infect Dis.* (2015) 211:968–77. doi: 10.1093/infdis/jiu566
35. Fava VM, et al. Age-dependent association of TNFSF15/TNFSF8 variants and leprosy type 1 reaction. *Front Immunol.* (2017) 8:155. doi: 10.3389/fimmu.2017.00155
36. Fava VM, Manry J, Cobat A, Orlova M, van Thuc N, Moraes MO, et al. A genome wide association study identifies a lncRNA as risk factor for pathological inflammatory responses in leprosy. *PLoS Genet.* (2017) 13:e1006637. doi: 10.1371/journal.pgen.1006637
37. Eickelmann M, Steinhoff M, Metzke D, Tomimori-Yamashita J, Sunderkötter C. Erythema leprosum—after treatment of lepromatous leprosy. *J Dtsch Dermatol Ges.* (2010) 8:450–3. doi: 10.1111/j.1610-0387.2009.07294.x
38. Bex-Bleumink M, Berhe D. Occurrence of reactions, their diagnosis and management in leprosy patients treated with multidrug therapy; experience in the leprosy control program of the all Africa leprosy and rehabilitation training Center (ALERT) in Ethiopia. *Int J Lepr Other Mycobact Dis.* (1992) 60:173–84.
39. Manandhar R, LeMaster JW, Roche PW. Risk factors for erythema nodosum leprosum. *Int J Lepr Other Mycobact Dis.* (1999) 67:270–8.
40. Kahawita IP, Lockwood DN. Towards understanding the pathology of erythema nodosum leprosum. *Trans R Soc Trop Med Hyg.* (2008) 102:329–37. doi: 10.1016/j.trstmh.2008.01.004
41. Saunderson P, Gebre S, Byass P. ENL reactions in the multibacillary cases of the AMFES cohort in Central Ethiopia: incidence and risk factors. *Lepr Rev.* (2000) 71:318–24. doi: 10.5935/0305-7518.20000035
42. de Messias IJ, Santamaria J, Brenden M, Reis A, Mauff G. Association of C4B deficiency (C4B\*Q0) with erythema nodosum in leprosy. *Clin Exp Immunol.* (1993) 92:284–7. doi: 10.1111/j.1365-2249.1993.tb03393.x
43. Schuring RP, Hamann L, Faber WR, Pahan D, Richardus JH, Schumann RR, et al. Polymorphism N248S in the human toll-like receptor 1 gene is related to leprosy and leprosy reactions. *J Infect Dis.* (2009) 199:1816–9. doi: 10.1086/599121
44. Girardi DR, Moro CM, Bulegon H. SeyeS – support system for preventing the development of ocular disabilities in leprosy. *Conf Proc IEEE Eng Med Biol Soc.* (2010) 2010:6162–5.
45. Liu S, McGree J, Ge Z, Xie Y. Computational and statistical methods for analysing big data with applications. *JAAD Case Rep.* (2016) 18:1–194. doi: 10.1109/IEMBS.2010.5627769
46. Ben-Gal I. Bayesian networks. *Encyclopedia of Statistics in Quality & Reliability.* (2007) 1–6. doi: 10.1002/9780470061572.eqr089
47. Belle A, Kon MA, Najarian K. Biomedical informatics for computer-aided decision support systems: a survey. *ScientificWorldJ.* (2013) 2013:769639
48. Bellazzi R, Zupan B. Predictive data mining in clinical medicine: current issues and guidelines. *Int J Med Inform.* (2008) 77:81–97. doi: 10.1155/2013/769639
49. Lu Z, Mitchell RM, Smith RL, Karns JS, van Kessel JAS, Wolfgang DR, et al. Invasion and transmission of salmonella Kentucky in an adult dairy herd using approximate Bayesian computation. *BMC Vet Res.* (2013) 9:245. doi: 10.1186/1746-6148-9-245
50. Tylman W, Waszyrowski T, Napieralski A, Kamiński M, Trafidło T, Kulesza Z, et al. Real-time prediction of acute cardiovascular events using hardware-implemented Bayesian networks. *Comput Biol Med.* (2016) 69:245–53. doi: 10.1016/j.combiomed.2015.08.015
51. Twardy C, Nicholson A, Korb K, Mcneil J. Epidemiological data mining of cardiovascular Bayesian networks. *Electron J Health Inform.* (2006) 1:1–13.
52. Thornley S, Marshall RJ, Wells S, Jackson R. Using directed acyclic graphs for investigating causal paths for cardiovascular disease. *J Biom Biostat.* (2013) 4:1–6. doi: 10.4172/2155-6180
53. Fuster-Parra P, Tauler P, Bannasar-Veny M, Ligeza A, López-González AA, Aguiló A. Bayesian network modeling: a case study of an epidemiologic system analysis of cardiovascular risk. *Comput Methods Prog Biomed.* (2016) 126:128–42. doi: 10.1016/j.cmpb.2015.12.010
54. Jiao Y, Wang XH, Chen R, Tang TY, Zhu XQ, Teng GJ. Predictive models of minimal hepatic encephalopathy for cirrhotic patients based on large-scale brain intrinsic connectivity networks. *Sci Rep.* (2017) 7:11512. doi: 10.1038/s41598-017-11196-y
55. Zhang Y, Zhang T, Zhang C, Tang F, Zhong N, Li H, et al. Identification of reciprocal causality between non-alcoholic fatty liver disease and metabolic syndrome by a simplified Bayesian network in a Chinese population. *BMJ Open.* (2015) 5:e008204. doi: 10.1136/bmjopen-2015-008204
56. Refai A, Merouani HF, Aouras H, Aouras H. Maintenance of a Bayesian network: application using medical diagnosis. *Evol Syst.* (2016) 7:187–96. doi: 10.1007/s12530-016-9146-8
57. The Alzheimer's Disease Neuroimaging Initiative/Jin Y, Su Y, Zhou XH, Huang S. Heterogeneous multimodal biomarkers analysis for Alzheimer's disease via Bayesian network. *EURASIP J Bioinform Syst Biol.* (2016) 2016:12. doi: 10.1186/s13637-016-0046-9
58. Souza WV, Barcellos CC, Brito AM, Carvalho MS, Cruz OG, Albuquerque MFM, et al. Empirical bayesian model applied to the spatial analysis of leprosy occurrence. *Rev Saude Publica.* (2001) 35:474–80. doi: 10.1590/S0034-89102001000500011
59. Smith RL, Grohn YT. Use of approximate Bayesian computation to assess and fit models of Mycobacterium leprae to predict outcomes of the Brazilian control program. *PLoS One.* (2015) 10:e0129535. doi: 10.1371/journal.pone.0129535
60. Crump RE, Medley GF. Back-calculating the incidence of infection of leprosy in a Bayesian framework. *Parasit Vectors.* (2015) 8:534. doi: 10.1186/s13071-015-1142-5
61. Joshua V, Mehendale S, Gupte MD. Bayesian model, ecological factors & transmission of leprosy in an endemic area of South India. *Indian J Med Res.* (2016) 143:104–6. doi: 10.4103/0971-5916.178618
62. Zhang X, Yuan Z, Ji J, Li H, Xue F. Network or regression-based methods for disease discrimination: a comparison study. *BMC Med Res Methodol.* (2016) 16:100. doi: 10.1186/s12874-016-0207-2
63. Wang N, Wang Z, Wang C, Fu X, Yu G, Yue Z, et al. Prediction of leprosy in the Chinese population based on a weighted genetic risk score. *PLoS Negl Trop Dis.* (2018) 12:e0006789. doi: 10.1371/journal.pntd.0006789
64. Gama RS, Souza MLM, Sarno EN, Moraes MO, Gonçalves A, Stefani MMA, et al. A novel integrated molecular and serological analysis method to predict new cases of leprosy amongst household contacts. *PLoS Negl Trop Dis.* (2019) 13:e0007400. doi: 10.1371/journal.pntd.0007400
65. Tió-Coma M, Kielbasa SM, van den Eeden SJF, Mei H, Roy JC, Wallinga J, et al. Blood RNA signature RISK4LEP predicts leprosy years before clinical onset. *EBioMedicine.* (2021) 68:103379. doi: 10.1016/j.ebiom.2021.103379

66. Penna GO, Pontes MAA, Cruz R, Gonçalves HS, Penna MLF, Bühner-Sékula S. A clinical trial for uniform multidrug therapy for leprosy patients in Brazil: rationale and design. *Mem Inst Oswaldo Cruz.* (2012) 107:22–7. doi: 10.1590/S0074-02762012000900005
67. de Sales Marques C, Brito-de-Souza VN, Guerreiro LTA, Martins JH, Amaral EP, Cardoso CC, et al. Toll-like receptor 1 N248S single-nucleotide polymorphism is associated with leprosy risk and regulates immune activation during mycobacterial infection. *J Infect Dis.* (2013) 208:120–9. doi: 10.1093/infdis/jit133
68. Sales-Marques C, Salomão H, Fava VM, Alvarado-Arnez LE, Amaral EP, Cardoso CC, et al. NOD2 and CCDC122-LACC1 genes are associated with leprosy susceptibility in Brazilians. *Hum Genet.* (2014) 133:1525–32. doi: 10.1007/s00439-014-1502-9
69. Witten IH, Frank E. *Data mining practical machine learning tools and techniques, Ed. 2nd ed* Elsevier (2005).
70. Bellazzi R, Ferrazzi F, Sacchi L. Predictive data mining in clinical medicine: a focus on selected methods and applications. *WIREs. Data Min Knowl Disc.* (2011) 1:416–30. doi: 10.1002/widm.23
71. The University of Waikato, Weka 3: Data mining software in Java. Available at: <https://www.cs.waikato.ac.nz/ml/weka/>, (2018).
72. Salzberg SL C4.5: Programs for machine learning by J. Ross Quinlan. Morgan Kaufmann Publishers, Inc *Mach Learn.* (1993) 16:235–40. doi: 10.1007/BF00993309
73. Hall M, Frank E, Holmes G, Pfahringer B, Reutemann P, Witten IH. The WEKA data mining software: an update. *ACM SIGKDD Explorations Newsletter.* (2009) 11:10–8. doi: 10.1145/1656274.1656278
74. Ruggieri S. Efficient C4.5. *IEEE Trans Knowl Data Eng.* (2002) 14:438–44. doi: 10.1109/69.991727
75. Kotthoff L, et al. Auto-WEKA 2.0: automatic model selection and hyperparameter optimization in WEKA. *J Mach Learn Res.* (2016) 17:1–5.
76. Bates DW, Kuperman GJ, Wang S, Gandhi T, Kittler A, Volk L, et al. Ten commandments for effective clinical decision support: making the practice of evidence-based medicine a reality. *J Am Med Inform Assoc.* (2003) 10:523–30. doi: 10.1197/jamia.M1370
77. Norsys Software Corp. *Netica API Programmer's library. Reference manual. Version 4.18.* (2010). Available at: [http://www.norsys.com/netica-j/docs/NeticaJ\\_Man.pdf](http://www.norsys.com/netica-j/docs/NeticaJ_Man.pdf)
78. Hornung R. Diversity forests: using Split sampling to enable innovative complex Split procedures in random forests. *SN Comput Sci.* (2022) 3:1. doi: 10.1007/s42979-021-00920-1
79. Gontijo-Lopes R., Smullin S. J., Cubuk E. D., Dyer E. Affinity and diversity: quantifying mechanisms of data augmentation. arXiv [Preprint] (2020) arXiv:2002.08973 doi: 10.48550/arXiv.2002.08973
80. Aroyo L., et al., DICES dataset: diversity in conversational AI evaluation for safety. arXiv [Preprint] (2023) arXiv:2306.11247. doi: 10.48550/arXiv.2306.11247