



MODELOS NEURAIIS PARA APOIO AO DIAGNÓSTICO DE TUBERCULOSE
COM RESISTÊNCIA AOS MEDICAMENTOS

Luiz Henrique Ramos de Azevedo Évora

Dissertação de Mestrado apresentada ao Programa de Pós-graduação em Engenharia Elétrica, COPPE, da Universidade Federal do Rio de Janeiro, como parte dos requisitos necessários à obtenção do título de Mestre em Engenharia Elétrica.

Orientadores: José Manoel de Seixas

Afranio Lineu Kritski

Rio de Janeiro

Março de 2014

MODELOS NEURAIIS PARA APOIO AO DIAGNÓSTICO DE TUBERCULOSE
COM RESISTÊNCIA AOS MEDICAMENTOS

Luiz Henrique Ramos de Azevedo Évora

DISSERTAÇÃO SUBMETIDA AO CORPO DOCENTE DO INSTITUTO
ALBERTO LUIZ COIMBRA DE PÓS-GRADUAÇÃO E PESQUISA DE
ENGENHARIA (COPPE) DA UNIVERSIDADE FEDERAL DO RIO DE
JANEIRO COMO PARTE DOS REQUISITOS NECESSÁRIOS PARA A
OBTENÇÃO DO GRAU DE MESTRE EM CIÊNCIAS EM ENGENHARIA
ELÉTRICA.

Examinada por:

Prof. José Manoel de Seixas, D.Sc.

Prof. Afranio Lineu Kritski, D.Sc.

Prof. Antonio Mauricio Miranda de Sá, D.Sc.

Profa. Fernanda Carvalho de Queiroz Mello, D.Sc.

Prof. Guilherme de Alencar Barreto, D.Sc.

RIO DE JANEIRO, RJ – BRASIL

MARÇO DE 2014

Évora, Luiz Henrique Ramos de Azevedo

Modelos Neurais para Apoio ao Diagnóstico de Tuberculose com Resistência aos Medicamentos/Luiz Henrique Ramos de Azevedo Évora. – Rio de Janeiro: UFRJ/COPPE, 2014.

XV, 98 p.: il.; 29, 7cm.

Orientadores: José Manoel de Seixas

Afranio Lineu Kritski

Dissertação (mestrado) – UFRJ/COPPE/Programa de Engenharia Elétrica, 2014.

Referências Bibliográficas: p. 82 – 91.

1. Tuberculose. I. Seixas, José Manoel de *et al.*
- II. Universidade Federal do Rio de Janeiro, COPPE, Programa de Engenharia Elétrica. III. Título.

Resumo da Dissertação apresentada à COPPE/UFRJ como parte dos requisitos necessários para a obtenção do grau de Mestre em Ciências (M.Sc.)

MODELOS NEURAIIS PARA APOIO AO DIAGNÓSTICO DE TUBERCULOSE
COM RESISTÊNCIA AOS MEDICAMENTOS

Luiz Henrique Ramos de Azevedo Évora

Março/2014

Orientadores: José Manoel de Seixas

Afranio Lineu Kritski

Programa: Engenharia Elétrica

A tuberculose (TB) continua a ser um dos principais problemas de saúde no mundo. Existem casos em que a doença apresenta resistência aos medicamentos utilizados em seu tratamento, sendo esta situação um agente dificultador em relação à sua cura, exigindo maiores cuidados e atenção sobre o paciente. Os métodos existentes apresentam deficiências em sua taxa de detecção, tempo de resposta, ou possuem altos custos e necessitam de infraestrutura para sua implementação. Neste trabalho, baseando-se em respostas de um questionário de anamnese, é proposto o desenvolvimento de modelos neurais para a classificação de pacientes entre diagnósticos positivos para TB (incluindo TB sensível, TB droga-resistente e TB multirresistente), assim como diagnósticos negativos. Para a classificação entre TB resistente e sensível, é visto um valor de $0,7308 \pm 0,0813$. Entre os diagnósticos de resistência, a sensibilidade foi de $0,7067 \pm 0,0967$. Estudos sobre grupos de baixo, médio e alto risco, utilizando mapas auto-organizáveis, são realizados, avaliando-se com isso, as contribuições de cada variável utilizada para o resultado final.

Abstract of Dissertation presented to COPPE/UFRJ as a partial fulfillment of the requirements for the degree of Master of Science (M.Sc.)

NEURAL MODELS TO SUPPORT THE DIAGNOSIS OF DRUG-RESISTANT TUBERCULOSIS

Luiz Henrique Ramos de Azevedo Évora

March/2014

Advisors: José Manoel de Seixas

Afranio Lineu Kritski

Department: Electrical Engineering

Tuberculosis (TB) remains as one of the major health problems in the world. There are cases in which the illness is resistant to the drugs used in its treatment, making it difficult to treat, and demanding more attention on the patient. All methods available have deficiencies in its detection rates, response time, or have a higher cost and needs a more complex infrastructure to be implemented. In this work, by using responses retrieved from an anamnesis questionnaire, it is proposed the development of neural models to classify patients into positive and negative diagnosis for TB (including drug-sensitive, drug-resistant and multi-drug-resistant TB). For the classification between resistant and sensitive TB, a sensibility of $0,7308 \pm 0,0813$ was seen. Between resistant TB diagnoses, a sensitivity of $0,8205 \pm 0,0676$ was reached. Studies on risk groups, using self-organizing maps (SOM) are made, evaluating the contribution of each variable to the final results.

Sumário

Lista de Figuras	viii
Lista de Tabelas	xiv
1 Introdução	1
1.1 Motivação	2
1.2 Objetivos	6
1.3 Organização do documento	6
2 Tuberculose Pulmonar	8
2.1 Introdução	8
2.2 Fatores de Risco	13
2.3 Sinais e Sintomas	18
2.4 Diagnóstico	19
2.4.1 Modelos Estatísticos	24
3 Base de Dados	30
4 Método	34
4.1 Análise das variáveis utilizadas	34
4.2 Redes neurais para apoio ao diagnóstico	37
4.2.1 Redes neurais multicamadas	37
4.2.1.1 Pré-processamento	37
4.2.1.2 Treinamento	40
4.2.1.3 Avaliação do desempenho	42

4.2.2	Mapas Auto-organizáveis	46
4.2.2.1	Pré-processamento	46
4.2.2.2	Treinamento	46
4.2.2.3	Identificação de agrupamentos e avaliação das variáveis	47
5	Resultados	49
5.1	Estudo entre diagnósticos de TB positivo e TB negativo	49
5.2	Estudo entre diagnósticos de TB resistente e TB sensível	57
5.3	Estudo entre diagnósticos de TB-DR e TB-MDR	64
5.4	Estudo entre diagnósticos de TB resistente, TB sensível e TB negativo	71
5.4.1	Rede neural dedicada	71
5.4.2	Rede neural composta por redes especialistas	74
6	Conclusão	78
	Referências Bibliográficas	82
A	Redes neurais multicamadas	92
B	Mapas Auto-organizáveis (SOM)	95

Lista de Figuras

2.1	Estimação da taxa de incidência de tuberculose em 2011. Fonte OMS . . .	10
2.2	Coeficiente de incidência de tuberculose no Brasil em 2012. Fonte Ministério da Saúde	10
2.3	Taxa de casos de tuberculose multirresistentes quando o paciente já havia sido tratado da doença anteriormente. Fonte OMS	12
2.4	Porcentagem de pacientes com diagnóstico positivo para tuberculose e com HIV por país. Fonte OMS	15
4.1	Distribuição das idades dos pacientes no (a) conjunto completo de dados, (b) com diagnóstico de TB positivo, (c) TB droga-resistente e (d) multirresistente.	35
4.2	Ilustração da técnica de validação cruzada pelo método k-fold.	39
4.3	Exemplo de particionamento de um conjunto de dados utilizando a técnica de k-means. Os pontos pretos representam os eventos, e os pontos azuis os centróides após algumas iterações.	40
4.4	Representação da evolução do erro de classificação de um treinamento por iteração e escolha do ponto de parada com a melhor generalização.	42
4.5	Exemplo de curva ROC para medir a relação entre sensibilidade e falsos positivos.	43
4.6	Distribuição de pacientes com diagnósticos positivos e negativos e exemplo de ponto de corte para definição da sensibilidade e especificidade.	44

5.1	Valores de SP médios para cada topologia, pelos conjuntos de dados de teste e treino, considerando as redes obtidas pela validação cruzada. As barras de erro são representadas pelos valores RMS dos índices SP. A topologia com 4 neurônios é a que apresenta melhor média para o conjunto de teste.	50
5.2	Representações da curva ROC para a rede de melhor SP e sensibilidade, evolução do treino para a rede de melhor SP e distribuição dos valores de saída da rede neural com maior sensibilidade para classificar pacientes entre os diagnósticos de TB positivo e TB negativo. . .	51
5.3	Estudo de relevância das variáveis na rede neural multicamadas obtida para os diagnósticos TB positivo e TB negativo. O índice SP da rede de operação é subtraído do índice SP alcançado pela rede quando se fixa cada variável com o seu valor médio.	52
5.4	U-Matrix do mapa SOM e o resultado da sua clusterização pela técnica de k-means em 3 agrupamentos para pacientes com diagnósticos TB positivo e TB negativo.	53
5.5	Definição dos agrupamentos em grupos de baixo, médio e alto risco, de acordo com a distribuição de pacientes.	54
5.6	Mapa de componentes referente ao mapa SOM para pacientes com diagnósticos TB positivo e TB negativo.	54
5.7	Clusterização do mapa SOM gerado para os pacientes com diagnósticos de TB positivo e TB negativo. (a) Índice de Davis-Bouldin por número de clusters (b) Mapa SOM clusterizado com 10 agrupamentos	56
5.8	Rotulação do mapa SOM entre os diagnósticos de TB positivo e TB negativo. (a) Frequência de diagnósticos em cada cluster (b) Rotulação dos clusters por diagnóstico	56

5.9	Valores de SP médios para cada topologia, pelos conjuntos de dados de teste e treino, considerando as redes obtidas pela validação cruzada. As barras de erro são representadas pelo valor RMS dos índices SP. A topologia com 8 neurônios é a que apresenta melhor média para o conjunto de teste.	58
5.10	Representações da curva ROC para a rede de melhor SP e sensibilidade, evolução do treino para a rede de melhor SP e distribuição dos valores de classificação para a rede neural com maior sensibilidade para classificar pacientes entre os diagnósticos de TB resistente e TB sensível.	59
5.11	Estudo de relevância das variáveis na rede neural multicamadas obtida para os diagnósticos TB resistente e TB sensível. O índice SP da rede de operação é subtraído do índice SP alcançado pela rede quando se define cada variável com o seu valor médio.	59
5.12	U-Matrix do mapa SOM e o resultado da sua clusterização pela técnica de k-means em 3 agrupamentos para dados de pacientes com TB resistente e TB sensível.	61
5.13	Definição dos agrupamentos em grupos de baixo, médio e alto risco, de acordo com a distribuição de pacientes.	61
5.14	Mapa de componentes referente ao mapa SOM para pacientes com diagnósticos TB resistente e TB sensível.	62
5.15	Clusterização do mapa SOM gerado para os pacientes com diagnósticos de TB resistente e TB sensível. (a) Índice de Davis-Bouldin por número de clusters (b) Mapa SOM clusterizado com 15 agrupamentos	63
5.16	Rotulação do mapa SOM entre os diagnósticos de TB resistente e TB sensível. (a) Frequência de diagnósticos em cada cluster (b) Rotulação dos clusters por diagnóstico	64

5.17	Valores de SP médios para cada topologia, pelos conjuntos de dados de teste e treino, considerando as redes obtidas pela validação cruzada. As barras de erro são representadas pelo valor RMS dos índices SP. A topologia com 11 neurônios é a que apresenta melhor média para o conjunto de teste.	65
5.18	Representações da curva ROC para a rede de melhor SP e sensibilidade, evolução do treino para a rede de melhor SP e distribuição dos valores de classificação para a rede neural com maior sensibilidade para classificar pacientes entre os diagnósticos de TB-DR e TB-MDR.	66
5.19	Estudo de relevância das variáveis na rede neural multicamadas obtida para os diagnósticos TB resistente e TB sensível. O índice SP da rede de operação é subtraído do índice SP alcançado pela rede quando se define cada variável com o seu valor médio.	67
5.20	U-Matrix do mapa SOM e o resultado da sua clusterização pela técnica de k-means em 3 agrupamentos.	68
5.21	Definição dos agrupamentos em grupos de baixo, médio e alto risco, de acordo com a distribuição de pacientes.	68
5.22	Mapa de componentes referente ao mapa SOM para pacientes com diagnósticos TB-DR e TB-MDR.	69
5.23	Clusterização do mapa SOM gerado para os pacientes com diagnósticos de TB-DR e TB-MDR. (a) Índice de Davis-Bouldin por número de clusters (b) Mapa SOM clusterizado com 6 agrupamentos	70
5.24	Rotulação do mapa SOM entre os diagnósticos de TB-DR e TB-MDR. (a) Frequência de diagnósticos em cada cluster (b) Rotulação dos clusters por diagnóstico	70

5.25	Valores de SP médios para cada topologia, pelos conjuntos de dados de teste e treino, considerando as redes obtidas pela validação cruzada. As barras de erro são representadas pelo valor RMS dos índices SP. A topologia com 13 neurônios é a que apresenta melhor média para o conjunto de teste.	72
5.26	Representações da evolução do treino e distribuição dos valores de classificação para a rede neural de operação que classifica pacientes entre os diagnósticos de TB resistente, TB sensível e TB negativo. . .	73
5.27	Estudo de relevância das variáveis na rede neural multicamadas obtida para os diagnósticos TB resistente, TB sensível e TB negativo. O índice SP da rede de operação é subtraído do índice SP alcançado pela rede quando se define cada variável com o seu valor médio. . . .	74
5.28	Modelo de redes neurais em cascata, composta por redes especialistas, para simular cenário de triagem médica.	75
5.29	Desempenho do valor preditivo positivo (VPP) e valor preditivo negativo (VPN), por prevalência, para o modelo de triagem.	77
5.30	Desempenho do valor preditivo positivo (VPP) e valor preditivo negativo (VPN), por prevalência, para o modelo de diagnóstico.	77
A.1	Exemplo básico de rede neural multicamada, com camada de entrada, escondida e de saída.	93
B.1	Diagrama de um mapa auto-organizável bi-dimensional, onde o vetor de entrada x é aplicado ao mapa com pesos sinápticos w . O neurônio vencedor é representado de cor clara, e pode-se observar a área de vizinhança, afetada pela sua ativação, possuindo nesse caso uma profundidade de dois neurônios.	96

B.2 Cada unidade da U-Matrix representa a distância entre um neurônio e seus vizinhos, calculo representado à esquerda. Um exemplo de diagrama de um mapa auto-organizável é mostrado à direita, o qual mostra uma região de afastamento horizontal representada em cor mais escura.	97
--	----

Lista de Tabelas

3.1	Variáveis utilizadas e sua respectiva codificação.	32
3.2	Quantidade de pacientes por diagnóstico e local de triagem médica.	32
3.3	Frequência de respostas dadas para as variáveis utilizadas por todo o conjunto de pacientes.	33
4.1	Razão de chance para diagnósticos de TB multirresistente em comparação com TB polirresistente, TB resistente (a qual inclui todas as formas de resistência) contra TB sensível, e TB positivo (incluindo os casos sensíveis e resistentes) contra diagnósticos negativos para TB.	36
4.2	Valores de parâmetros utilizados para o treinamento das redes neurais.	42
5.1	Desempenho alcançado pelo classificador para pacientes com diagnósticos TB positivo e TB negativo, considerando a topologia com 4 neurônios na camada escondida.	50
5.2	Desempenho alcançado pelo classificador, por localização da triagem, para pacientes com diagnósticos TB positivo e TB negativo.	52
5.3	Desempenho alcançado pelo classificador para pacientes com diagnósticos TB resistente e TB sensível, considerando a topologia com 8 neurônios na camada escondida.	57
5.4	Desempenho alcançado pelo classificador, por localização da triagem, para pacientes com diagnósticos TB resistente e TB sensível.	60

5.5	Desempenho alcançado pelo classificador para pacientes com diagnósticos TB-DR e TB-MDR, considerando a topologia com 11 neurônios na camada escondida.	65
5.6	Desempenho alcançado pelo classificador, por localização da triagem, para pacientes com diagnósticos TB-DR e TB-MDR.	67
5.7	Desempenho alcançado pelo classificador para pacientes com diagnósticos TB resistente, TB sensível e TB negativo, considerando a topologia com 13 neurônios na camada escondida.	72
5.8	Desempenho alcançado pelo classificador composto por redes especialistas em cascata para modelos de triagem e diagnóstico.	76

Capítulo 1

Introdução

Na medicina, as situações básicas de diagnóstico, tratamento, prognóstico e prevenção estão sempre presentes. Em cada uma dessas situações, para se obter um bom resultado é preciso identificar o problema e realizar um estudo cuidadoso sobre as informações disponíveis.

O diagnóstico é parte da consulta ou do atendimento pela equipe de saúde que tem como objetivo a identificação de uma eventual doença. Um conjunto de dados é analisado pelo profissional de saúde para, então, ser sintetizado em uma ou mais doenças. Os dados podem ser provenientes dos sinais e sintomas apresentados, do histórico clínico, de exame físico e de exames complementares. Com a identificação da doença, é então realizado o planejamento para o tratamento, além do prognóstico, tendo como base o quadro apresentado.

O diagnóstico médico é uma tarefa fundamental na prática da medicina e para o direcionamento do doente para um tratamento adequado. A habilidade de realizar um diagnóstico correto exige não só conhecimento médico e experiência do profissional de saúde, mas também um elaborado raciocínio clínico. Pelo grande volume de informações que um médico necessita para exercer sua profissão, sistemas de apoio à decisão se apresentam como uma poderosa ferramenta no entendimento e processamento de informações e das incertezas médicas associadas.

Um Sistema de Apoio à Decisão (SAD) é uma classe de Sistemas de Informação (SI) que se refere a um modelo genérico de tomada de decisão. Esse modelo possi-

bilita a análise de um grande volume de variáveis para que seja possível a indicação de uma resposta a uma determinada questão.

Uma das áreas de maior atuação dos SAD é a medicina. Conforme afirma a Organização Mundial da Saúde, a quantidade de informações referentes à área da saúde dobra a cada três anos, surgindo novos métodos de diagnósticos e terapêutica, princípios químicos, inovações da área da biologia molecular e da genética, entre outros avanços.

Além disso, um grande volume de informações disponíveis para o profissional de saúde sobre um paciente pode dificultar o julgamento clínico ou reconhecimento de padrões entre uma variedade de pacientes. Portanto, um SAD aplicado na medicina tem papel fundamental na análise, diagnóstico e processamento de informações sobre um paciente, o que, aliado aos conhecimentos dos médicos, auxilia na identificação de possíveis soluções.

É desejável que os sistemas de apoio aos procedimentos de triagem e diagnóstico médico apresentem uma alta sensibilidade, ou seja, um alto nível de acerto entre os pacientes que apresentem a doença, assim como uma alta especificidade, que representa a correta identificação de pacientes sem a doença. O sistema deve ser de fácil utilização e implementação, o que pode proporcionar a agilização do procedimento médico, além de não acarretar em altos custos e mantendo sua eficácia.

1.1 Motivação

A potencialidade da utilização de sistemas probabilísticos de apoio ao diagnóstico, na área médica, é vista principalmente, pelos fatores econômicos e sociais, além da limitação de sensibilidade observada em diversos testes diagnósticos, como é o caso da tuberculose (TB). Cerca de 8,6 milhões de pessoas foram diagnosticadas apenas no ano de 2012 no mundo, onde 1,3 milhões de mortes foram consequências desse cenário [1]. Apesar de existirem tratamentos e medicamentos eficazes, a TB é considerada uma das principais doenças que aflige a humanidade, constituindo em um sério problema de saúde pública.

Mais especificamente no Brasil, foram notificados 70.047 casos novos em 2012, colocando o país no 17º lugar em relação aos 22 países que apresentam maior incidência em tuberculose [2]. Algumas unidades federativas ainda apresentam uma taxa de TB bem mais elevadas que a média nacional, como é o caso do Rio de Janeiro, onde em 2012 só nessa localidade foram notificados 13.053 casos [3]. Um dos fatores que mais favoreceu o aumento da incidência da tuberculose foi a coinfeção pelo vírus da imunodeficiência humana (HIV). Essa associação HIV/TB constitui um sério problema de saúde pública, que leva ao aumento do número de mortes pela doença. Outros fatores como a falta de sistemas públicos de saúde eficientes, a desigualdade social, o crescimento da população marginalizada rural e urbana, são relacionados com o aumento da incidência de TB.

A TB pulmonar é uma doença infecto-contagiosa transmitida pelo ar. É causada por um microrganismo denominado *Mycobacterium tuberculosis*, o qual se propaga por meio do ar, por meio de gotículas contendo os bacilos expelidos por um doente ao tossir ou espirrar. A inalação destas gotículas por pessoas saudáveis, acarreta na infecção tuberculosa e o risco de desenvolvimento da doença. Após penetrar no organismo pela via respiratória, o agente causador da tuberculose ainda pode disseminar-se e instalar-se em qualquer órgão, chamando-se tuberculose extra-pulmonar, sendo mais frequentes na pleura e nos linfonodos.

Uma grande ameaça à saúde pública são os casos que apresentam resistência às drogas utilizadas normalmente no tratamento de TB pulmonar. Esse fator pode ser acarretado por regimes de tratamentos inadequados, os quais podem ser ocasionados pelo abandono total ao tratamento; interrupção do tratamento; uso irregular das drogas, como doses incorretas por erro de prescrição, como no caso de não ser levada em consideração a interação medicamentosa; ou mesmo por recusa do doente. Ainda, um paciente que desenvolva a TB em sua forma resistente, poderá transmiti-la para outros indivíduos, o que agrava ainda mais a situação.

A resistência ao medicamento pode se apresentar em sua forma monorresistente, onde o paciente apresenta resistência a somente um medicamento anti-TB; polirre-

sistente, onde apresenta resistência a mais de um medicamento anti-TB, mas não a combinação de Isoniazida e Rifampicina; multiresistente (MDR), na qual a resistência é sobre a combinação Isoniazida e Rifampicina; e extensivamente resistente (XDR), onde o bacilo é resistente até a medicamentos utilizados no tratamento do tipo MDR.

Atualmente, em torno de 3,7% de pacientes com tuberculose no mundo possuem a doença em sua forma resistente. Caso sejam considerados apenas pacientes que já haviam sido tratados, o número sobe para 20%. No ano de 2012, estima-se que 450 mil pessoas tenham sido diagnosticadas com TB-MDR, sendo mais da metade dos casos ocorridos na China, Índia e Rússia [4]. Apesar dos casos de TB-XDR serem considerados mais raros, estimativas sugerem que cerca de 9% dos casos de TB-MDR possuem características de TB-XDR. O número ainda pode ser maior, considerando a dificuldade de países mais pobres no diagnóstico de tais casos.

O diagnóstico da TB pulmonar é realizado com base nos sinais e sintomas relatados pelo paciente, associados ao uso de testes diagnósticos. A baciloscopia e a cultura para a micobactéria são indicados como testes fundamentais para o diagnóstico da tuberculose pulmonar. A baciloscopia direta do escarro é o exame prioritário para os casos suspeitos de TB pulmonar, já que permite descobrir a mais importante fonte de infecção, o paciente bacilífero. Por ser um método simples e seguro, é praticado em todos os serviços de saúde que possuem um laboratório. Como uma desvantagem, a baciloscopia possui uma baixa sensibilidade, entre 40% e 80% [5]. Outra característica da baciloscopia é que ela não consegue identificar as cepas de tuberculose resistentes ao tratamento e também não as discrimina de outras micobactérias atípicas.

Já a cultura é o teste mais sensível para o diagnóstico de TB pulmonar, sendo considerado o padrão-ouro. Entretanto, a cultura possui o inconveniente de demorar de 4 a 8 semanas até a obtenção do resultado, sendo assim, pouco utilizada no processo de tomada de decisão clínica [5]. Além disso, a cultura exige infra-estrutura maior em biossegurança, dificultando sua implementação em locais de poucos recur-

sos. Desse modo, a cultura é indicada a suspeitos de tuberculose que apresentem diagnósticos negativos na baciloscopia, além de ser indicada para os casos extrapulmonares, HIV positivos e em crianças. É possível ainda a diferenciação de cepas sensíveis e resistentes aos medicamentos disponíveis para o tratamento, sendo com isso recomendada também para casos suspeitos de resistência bacteriana aos medicamentos.

Métodos mais caros e que necessitam de maior estrutura também podem ser utilizados. A radiografia de tórax é utilizada como exame auxiliar, apresentando taxas de sensibilidade de 70% a 80% na presença das anormalidades típicas em adultos imunocompetentes, e 95% se considerarmos outras anormalidades [6]. Sua grande desvantagem é a logística necessária para se realizar esse procedimento, assim como seu custo. O teste molecular Xpert® MTB/RIF também é uma opção, funcionando não só para a detecção do bacilo da tuberculose, mas para a triagem de cepas resistentes a rifampicina [7]. O teste realizado pode fornecer resultados em um laboratório em menos de 2 horas, sem necessitar de tratamento da amostra ou de pessoas especializadas em biologia molecular. Apesar de alcançar bons níveis em sensibilidade na detecção da doença, sua desvantagem principal é o alto custo e a necessidade de estrutura específica no local do exame, sendo sua utilização ainda restrita a poucos laboratórios de referência.

Portanto, na ausência do resultado da cultura, e ou pela deficiência na sensibilidade de outros métodos ou sua indisponibilidade, muitos casos acabam sendo diagnosticados com base nos sintomas clínicos, exames radiológicos e outros testes laboratoriais, o que reduz a eficiência no diagnóstico. Desse modo, o diagnóstico da TB tende a produzir atrasos na identificação de um paciente contaminado com a doença, retardando o tratamento e facilitando a transmissão do bacilo em uma comunidade.

Portanto, o desenvolvimento de um modelo computacional que, ao ser alimentado pelos dados fornecidos ainda na triagem pelo paciente, possibilite identificar um diagnóstico de tuberculose pulmonar, assim como suas formas resistentes, tende a

ser de grande auxílio na prática médica. Além disso, auxiliar no entendimento de como cada sintoma ou resposta dada por um paciente pode se relacionar com a doença, assim como características relacionadas a regionalidade, são também de grande valia no processo de tomada de decisão.

1.2 Objetivos

Este trabalho objetiva o desenvolvimento de um modelo matemático, por meio de técnicas de inteligência computacional, que permita indicar com rapidez o diagnóstico de um paciente com suspeita de tuberculose resistente aos tratamentos medicamentosos. O modelo deverá ser de fácil utilização por equipes de enfermagem e aplicável a regiões com poucos recursos. Para isso, técnicas de redes neurais supervisionadas, como redes multicamadas, e não-supervisionadas, como mapas auto-organizáveis (SOM), são utilizadas. Um conjunto seletivo de sinais e sintomas dos pacientes é analisado, escolhidos pela sua relevância em relação a doença. Os modelos desenvolvidos identificam as variáveis mais determinantes na eficiência de detecção de TB resistente e multirresistente, estabelecendo as probabilidades de existência de TB e multirresistência, além de associar, de forma complementar, grupos de risco (baixo, médio e alto) e uma forma visual de representação de informações.

1.3 Organização do documento

No próximo capítulo, é apresentada uma breve introdução a doença em estudo, definido suas principais características, sintomas e seu processo de diagnóstico. É realizada uma revisão da literatura sobre diferentes trabalhos realizados em sistemas de apoio à decisão na área médica e técnicas utilizadas no diagnóstico de TB por meio de inteligência computacional.

No capítulo 3 é apresentada a base de dados a ser utilizada nos estudos. Já no capítulo 4 é descrito o método do trabalho, sendo descritas as técnicas utilizadas para atingir o objetivo deste trabalho.

No capítulo 5, são apresentados os resultados obtidos. As conclusões e discussões sobre os resultados são apresentadas no capítulo 6, assim como perspectivas futuras para a continuidade dos trabalhos de pesquisa.

Capítulo 2

Tuberculose Pulmonar

Nesse capítulo, será apresentada uma introdução sobre a Tuberculose Pulmonar, além de descrever o tipo da doença que apresenta resistência aos medicamentos, sendo esses os temas objeto do estudo desse trabalho. Os fatores de risco são descritos, além da análise dos sintomas apresentados por pacientes com a doença e as formas e avaliações de diagnóstico existentes. Um estudo bibliográfico é realizado, apresentando os métodos mais utilizados e um resumo de trabalhos já produzidos nessa área.

2.1 Introdução

A TB é uma doença que aflige os seres humanos desde os primórdios da história [8]. Acredita-se que o primeiro registro de estudos sobre esse mal foi escrito por Demócrito, em 460 A.C., um trabalho que hoje considera-se perdido, tendo apenas seu título "*On those who are attacked by cough after illness*" sobrevivido ao tempo. Já alguns séculos depois, encontram-se registros de Hipócrates em que os procedimentos para lidar com a doença eram expostos, princípios que utilizamos até os dias de hoje, como manter o doente em um ambiente arejado e limpo [9].

Atualmente, a TB continua sendo um importante problema de saúde no mundo e que exige o desenvolvimento de estratégias para controlá-la, que envolvem aspectos humanitários, econômicos e de saúde pública. Apesar de já existirem recursos

tecnológicos capazes de promover seu controle, ainda não há perspectiva de se obter, em futuro próximo, sua eliminação como problema de saúde pública [10].

A tuberculose pulmonar é uma doença infecciosa e contagiosa, causada por um microrganismo denominado *Mycobacterium tuberculosis*, chamado também de bacilo de Koch, o qual se propaga por meio do ar, por meio de gotículas contendo os bacilos expelidos por um doente ao tossir, espirrar ou falar em voz alta. Quando estas gotículas são inaladas por pessoas saudáveis, provocam a infecção tuberculosa e o risco de desenvolver a doença [11].

Estima-se que, em 2012, 8,6 milhões de pessoas foram diagnosticadas com TB, com 1,3 milhões de mortes decorrentes desse cenário [12]. Com esses números, a TB se apresenta em segundo lugar, enquanto doença com mais alta taxa de mortalidade, considerando-se casos onde existe um único agente infeccioso, perdendo apenas para o HIV/AIDS [13]. Um fator comum entre essas duas doenças são seus efeitos profundos no sistema imunológico, que são capazes de enfraquecer as respostas imunológicas do hospedeiro [14].

Em grande parte do mundo desenvolvido tem se observado um declínio continuado da taxa de incidência de tuberculose. Entretanto, países menos desenvolvidos ainda apresentam, em suas populosas regiões, a manutenção de elevadas taxas de TB. Em áreas muito afetadas pela AIDS, como a África subsaariana, as taxas de casos de TB chegam a ser de 50 a 100 vezes maiores que na América do Norte ou Europa. Estimulado pela imunodeficiência causada pelo HIV, o aumento do número dos casos é visto como uma tendência para as próximas décadas [15]. A Figura 2.1 mostra as taxas mundiais de incidência de tuberculose no mundo, que dá um panorama desse cenário.

Com a notificação de 70.047 casos novos no Brasil em 2012, o país ainda não conseguiu alcançar a meta, estipulada pela OMS, de curar 85% dos casos novos bacilíferos, mesmo havendo queda no índice de proporção de cura para os casos novos de TB bacilífera. Em 2010, a proporção de cura foi de 73,4% e em 2011 regrediu para 71,6% [2]. Algumas unidades federativas apresentam coeficientes de

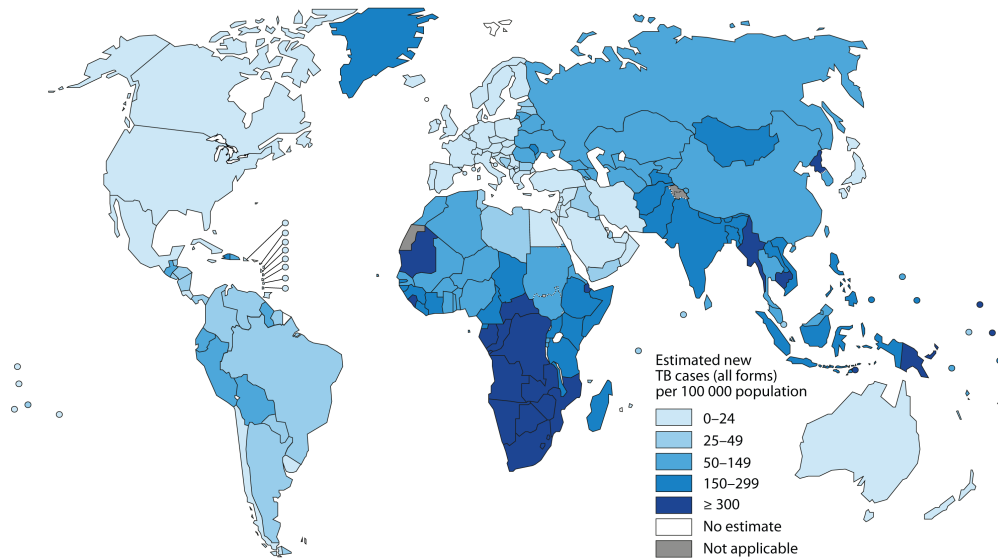


Figura 2.1: Estimação da taxa de incidência de tuberculose em 2011. Fonte OMS [13].

contaminação por número de habitantes bem elevados em comparação com a média nacional, como visto na Figura 2.2, onde os estados do Rio de Janeiro e do Amazonas apresentam as situações mais graves. Somente no ano de 2012, foram notificados, no Rio de Janeiro, 13.053 casos com confirmação de tuberculose [3].

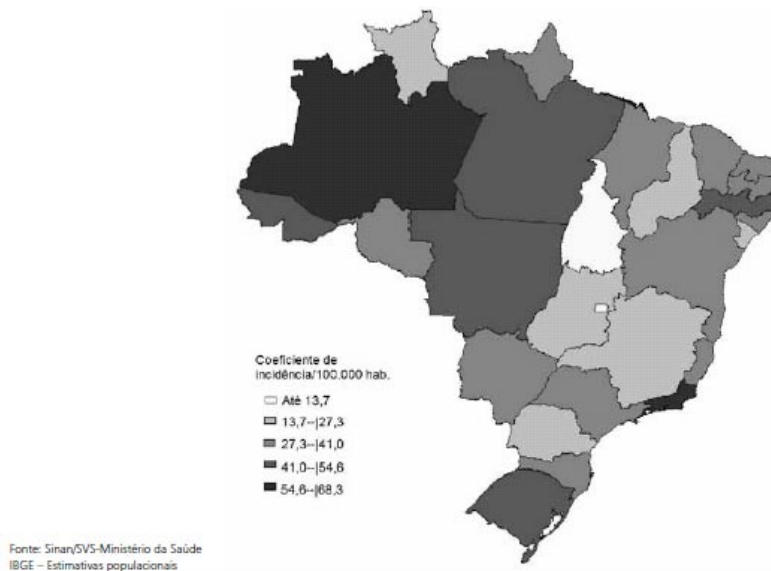


Figura 2.2: Coeficiente de incidência de tuberculose no Brasil em 2012. Fonte Ministério da Saúde[2].

Existem casos em que a doença se apresenta de forma resistente às drogas utilizadas normalmente no seu tratamento. Esse fator é visto como uma grande ameaça à saúde pública mundial, sendo consequência de um número de ações, como regi-

mes de tratamentos inadequados que podem ser ocasionados pelo abandono total ao tratamento; interrupção do tratamento; uso irregular das drogas, como doses incorretas por erro de prescrição ou no caso de não ser levada em consideração a interação medicamentosa; ou mesmo por recusa do doente [16]. Essencialmente, os casos de resistência aos medicamentos ocorrem em áreas onde os programas de controle de tuberculose são de baixa efetividade. Ainda, um paciente que desenvolva a TB em sua forma resistente, poderá transmití-la para outros indivíduos, o que agrava ainda mais a situação [17].

Os tipos mais básicos de resistência aos medicamentos são denominadas tuberculose monorresistente, onde é vista resistência a apenas um medicamento anti-TB, e a tuberculose polirresistente (ou droga-resistente), que apresenta resistência a mais de um medicamento anti-TB, mas que não seja a combinação de Isoniazida e Rifampicina. Ambos os tipos não impedem que o tratamento seja realizado pela primeira linha de medicação [18].

Apresentando um nível de resistência mais complexo em termos de tratamento, foi conceitualmente denominada *multidrugresistant tuberculosis* (MDR-TB) nos Estados Unidos, ao final dos anos 80, o tipo de tuberculose que apresentava bacilos resistentes às drogas Rifampicina e à Isoniazida simultaneamente [19]. O órgão de saúde pública norte-americano, o CDC (*Centers for Disease Control*), publicou em 1992 a trajetória da tuberculose multirresistente no país. Durante o período de 1982-1986, apresentou-se um número de 0,5% dos casos na condição de resistência aos medicamentos, havendo um aumento dessa proporção em 1991 para 3,1%. Esse aumento, que levou os EUA a registrarem 28.000 casos acima do esperado, entre os anos de 1980 e 1992, fez o problema ser reconhecido como de saúde pública. Com os norte-americanos elevando sua preocupação em relação a essa nova apresentação da doença, o mundo, em especial os países europeus, retomaram seus cuidados com esse cenário [20].

Quando identificada a ocorrência da TB-MDR, o tratamento da doença passa a ser realizado por medicamentos de segunda linha, os quais são menos atraentes

inicialmente pelos seus efeitos colaterais, necessidade de um período maior de tratamento e custo mais elevado. Entretanto, existem ainda estirpes de tuberculose que são resistentes não só ao tratamento de primeira linha, mas também a alguns tratamentos de segunda linha. Definida como *Extensively Drug-Resistant Tuberculosis* (TB-XDR), esse tipo de infecção foi apontada em março de 2006 pela Organização Mundial da Saúde (OMS) e o CDC, como uma emergente e grave ameaça à Saúde Pública e ao controle da TB [21]. Seu nível de resistência torna seu tratamento extremamente complicado, se não impossível, em ambientes com recursos limitados.

Atualmente, cerca de 3,7% de pacientes com tuberculose no mundo possuem a doença em sua forma resistente. Esse número é mais elevado se considerados os pacientes que já haviam sido tratados, cerca de 20%. Estima-se que 450 mil pessoas foram notificadas com casos de MDR-TB no mundo em 2012. Dentre esses casos, mais da metade do total ocorreram na China, Índia e Rússia. Em 2009, 48% dos pacientes que apresentaram a forma resistente da doença obtiveram sucesso no tratamento [4]. Na Figura 2.3 é possível observar a distribuição dos casos de tuberculose com resistência, quando o paciente já havia sido tratado anteriormente da doença, fator esse de alta relevância nesse cenário.

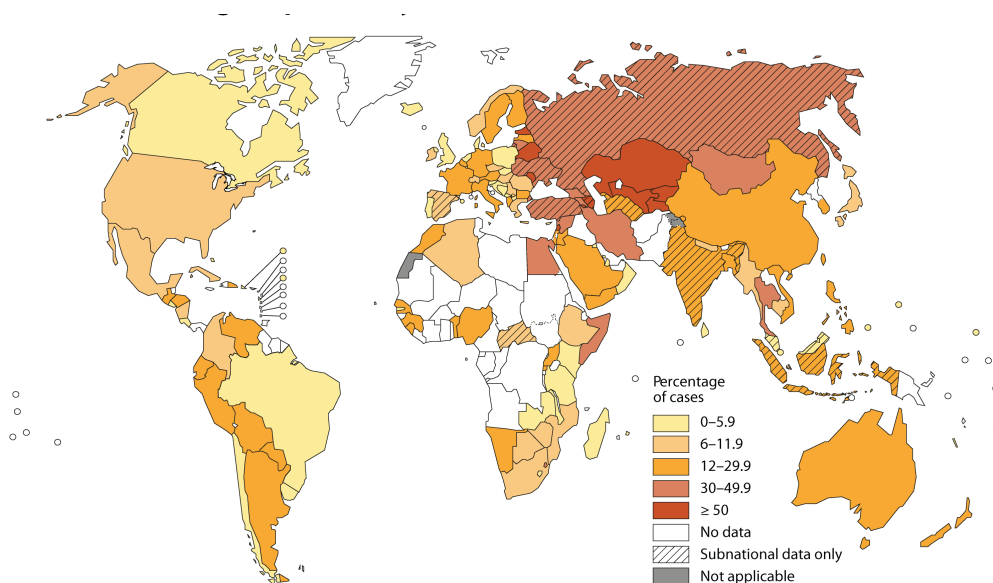


Figura 2.3: Taxa de casos de tuberculose multirresistentes quando o paciente já havia sido tratado da doença anteriormente. Fonte OMS [13].

Já casos de XDR-TB são mais raros, apesar de que 77 países reportaram pelo me-

nos um caso até o final de 2011. As informações divulgadas pelos países sugerem que cerca de 9% dos casos de MDR-TB possuem as características de XDR-TB, porém, devido à deficiência dos países mais pobres em diagnosticar casos de MDR/XDR-TB, estima-se que apenas uma pequena proporção desses casos são detectados e tratados apropriadamente [22].

No ano 2000, foi estabelecido o *Stop TB Partnership* [23], compreendendo um conjunto de 1.000 parceiros institucionais, incluindo organizações internacionais, não-governamentais e governamentais, que possui como objetivo a eliminação da tuberculose como problema de saúde pública, buscando estratégias para aumentar a efetividade do tratamento diretamente observado (TDO). Nesse âmbito, foi traçado um plano global com o objetivo de que, entre os anos de 2011 e 2015, cerca de 1 milhão de pacientes com MDR-TB deverão ser identificados e colocados em tratamento. Outra meta é alcançar uma taxa de 75% de sucesso no tratamento. Como base para esses números, em 2011, cerca de 18% dos casos entraram em tratamento e a taxa de 75% de sucesso no tratamento foi atingida por 30 países (considerando casos com tratamento iniciado em 2009)[23].

A estratégia de TDO foi adotada por 187 dos 193 membros da OMS com níveis altos de alcance na população, estimando-se que, em 2005, 89% das pessoas no mundo viviam em áreas que implementam essa política. Passados 10 anos do lançamento da estratégia de TDO, foi observada uma melhora na proporção de sucesso do tratamento anti-TB, entretanto, o aumento da cobertura desta estratégia não levou a uma melhoria na detecção dos casos de TB, fator mais evidenciado para os infectados pelo tipo resistente da doença [24].

2.2 Fatores de Risco

Dentre os fatores de risco relacionados a TB, pessoas portadoras da Síndrome da Imunodeficiência Adquirida (AIDS), que fazem tratamento de quimioterapia e receptores de transplantes, sob medicação para evitar rejeição, possuem uma maior probabilidade de contrair a doença. Um risco maior também afeta crianças, idosos,

indivíduos com má nutrição, profissionais de saúde, mineiros portadores de silicose, indivíduos com dependência de álcool, pessoas que utilizam medicamentos com corticoides e portadores de outras doenças, como a diabetes [25].

Algumas pessoas desenvolvem a tuberculose logo após serem infectadas, em um período de semanas, antes que o sistema imunológico possa combater a bactéria causadora. Em outros casos, a doença pode se desenvolver anos depois, aproveitando-se de uma situação em que a imunidade da pessoa esteja enfraquecida. Cerca de 5 a 10% das pessoas infectadas, que não receberam tratamento para a infecção de tuberculose latente, irão desenvolver a doença em alguma época de suas vidas [26].

Desde o surgimento da AIDS, em 1981, vem sendo observado, tanto em países desenvolvidos como nos em desenvolvimento, um número crescente de notificações de casos de TB em pessoas infectadas pelo HIV. Os índices de mortalidade em pacientes com AIDS e TB são maiores que de pacientes apenas com AIDS, sem TB. Sendo assim, a infecção pelo HIV, se não é o mais grave fator de risco para o desenvolvimento da tuberculose, é o mais característico, em termos de morbidade [27]. Em alguns casos, pessoas com tuberculose latente, que são infectadas pelo vírus do HIV, aumentam 100 vezes o risco da reativação da TB [28]. No Brasil, 53,3% dos novos casos de TB diagnosticados em 2012 realizaram o teste de anti-HIV, ano em que a coinfeção TB-HIV registrada no país foi de 9,7% [2]. A Figura 2.4 mostra um mapa com a porcentagem dos casos notificados com diagnóstico positivo para TB que apresentam HIV, o que ajuda a evidenciar essa situação, onde um grande número de países apresenta um número maior que 50%, inclusive o Brasil, com uma taxa de 55%.

A tuberculose, no Brasil, predomina no sexo masculino, na relação de 2 para 1, em relação ao feminino. Apesar de ocorrer com maior força na faixa etária do adulto jovem, a doença também se apresenta nas faixas etárias dos idosos [29], possuindo características diferentes. Nesses casos, a história de tuberculose no passado se torna predominante, e a associação com diabetes, doenças cardiovasculares e doenças pulmonares é maior. A dispneia e o emagrecimento são mais prevalentes, os eventos

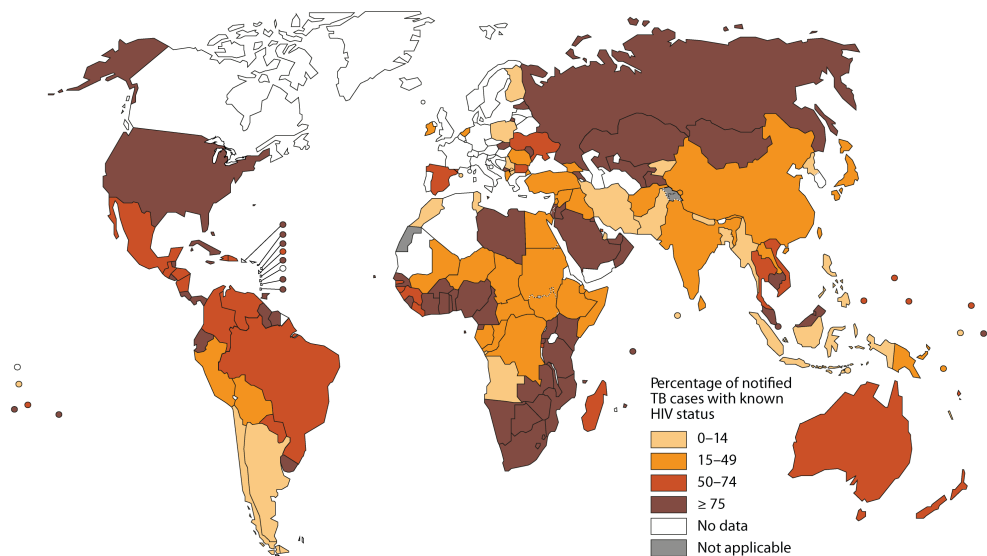


Figura 2.4: Porcentagem de pacientes com diagnóstico positivo para tuberculose e com HIV por país. Fonte OMS[13].

adversos associados aos medicamentos antituberculose são mais intensos, o índice de cura é menor, e a mortalidade é significativa [30][31].

Uma estimativa precisa da carga global de tuberculose em crianças é difícil, principalmente por causa dos desafios sobre a apuração e diagnóstico, além dos programas de saúde fracos em muitos países com uma alta carga de tuberculose [32]. Especialmente quando possuem menos de 5 anos de idade, existe uma dificuldade em expelir catarro, exigindo a realização de outros métodos para retirar o material necessário para os exames, como aspiração nasofaríngea ou gástrica [33]. Por estas razões, as crianças são muitas vezes excluídas dos estudos de prevalência da tuberculose [34], o que dificulta ainda mais a coleta de informações confiáveis. Estima-se que as crianças menores de 15 anos contribuem com 15% a 20% da carga global de tuberculose [35]. Poucos países têm relatado dados sobre tuberculose na infância, e as taxas de notificação de casos relatados variam entre 3% e 25% [36]. O risco de infecção de tuberculose em crianças depende da duração, e da proximidade de exposição a uma fonte infecciosa, que é muitas vezes um adulto [37], e sua progressão após a infecção primária é altamente variável, dependendo da idade e do estado imunológico [38]. A vacina BCG é recomendada para todas as crianças logo após o nascimento em ambientes de tuberculose-predominante, mas apresenta uma eficácia

preventiva média de 50% contra a tuberculose pulmonar [32].

Outras doenças, que podem se apresentar como co-morbidades, como o câncer, a insuficiência renal crônica, e o diabete mellitus, podem determinar uma maior demora na suspeita e na determinação do diagnóstico, o que ajuda na contribuição para uma maior transmissão de TB em unidades de saúde, não só entre pacientes, mas entre profissionais de saúde [39]. Ainda, a OMS considera o profissional de saúde com risco ocupacional relevante, levando em consideração os que atuam em atividades assistenciais de controle da TB em países em desenvolvimento [40].

Em relação a tuberculose multirresistente, estudos apontam as condições de moradia, como falta de esgoto no domicílio, além de alcoolismo associado ao tabagismo, número de tratamentos anteriores, tratamentos irregulares e cavidades pulmonares como importantes fatores para o seu desenvolvimento [41][42]. Em geral, os fatores de risco relacionados à forma resistente da doença são similares aos da não-resistente, devido às suas proximidades na formas de transmissão e desenvolvimento do agente causador.

Após penetrar no organismo pela via respiratória, o agente causador da tuberculose pode disseminar-se e instalar-se em qualquer órgão. Com isso, existe o risco de casos (cerca de 10 a 15% do total infectado) em que a tuberculose é observada em outro local do organismo que não o pulmão, chamando-se tuberculose extrapulmonar. As principais formas de tuberculose extrapulmonar são: tuberculose pleural, empiema pleural tuberculoso, tuberculose ganglionar periférica, tuberculose meningoencefálica, tuberculose pericárdica e tuberculose óssea [43]. Geralmente, quando se está detectando a TB extrapulmonar, o pulmão pode também já estar afetado [44]

Um outro importante fator de risco para a disseminação da doença, e desenvolvimento de suas formas resistentes, é o abandono ao tratamento. Ao não seguir corretamente o tratamento, o portador de TB continua doente, e permanece como fonte de contágio, aumentando-se ainda a probabilidade do surgimento de uma resistência medicamentosa. A não adesão ou abandono do tratamento podem estar

associados a diversos fatores, podendo ser relacionados ao paciente (uso irregular da medicação e/ou não ingestão da mesma, baixo nível socioeconômico, hábitos de vida e internações por outras doenças), ao medicamento (pelos efeitos colaterais e o relativo longo tempo de duração do tratamento) e ao desempenho do serviço de saúde (falhas na orientação do paciente, prescrições medicamentosas inadequadas, falta de fornecimento da medicação e falhas no agendamento de consultas) [45]. No Brasil, a proporção de abandono é considerada alta, atingindo em algumas capitais, cerca de 25%, em média, dos pacientes tratados [46]. Com uma política de tratamento observado, é possível ver uma melhora nos índices de abandono. Como um exemplo, na região Centro-Oeste e no estado de Tocantins, a proporção de abandono de um grupo auto-administrado e de um grupo de pacientes observados apresentaram números de 11% e 2,2% respectivamente, significando uma grande melhora [47]. Como comparação, a proposta da OMS é que sejam mantidos valores para essas taxas inferior a 5% do total de casos [48].

O ambiente hospitalar, caso não adote procedimentos adequados ao fluxo de pacientes com tuberculose, também se torna um fator de risco para contágio da doença. Salas de espera cheias e/ou partilhadas por muitas especialidades, além de diversas etapas que possam ser criadas para o acesso do paciente ao médico, são um terreno fértil à micobactéria. Para evitar tais situações, medidas devem ser tomadas, como atendimento com hora marcada para evitar aglomerações, e salas de espera em locais bem ventilados, se possível ao ar livre [49]. Pacientes com suspeita de tuberculose, ou com diagnóstico positivo ainda em período infectante, devem ter prioridade no atendimento, para tornar o tempo na unidade de saúde o menor possível. Na âmbito hospitalar, pacientes considerados de risco devem ficar em isolamento respiratório, em quartos individuais, ou em condições excepcionais, ter mais de um paciente por quarto, desde que os pacientes estejam com diagnóstico confirmado e não possuam o tipo resistente da doença [50].

2.3 Sinais e Sintomas

A tuberculose é uma doença que, nem sempre, quando uma pessoa é infectada, há a ocorrência de sintomas. Quando surge um caso bacilífero em uma pessoa que habita com familiares, e outros membros da família são infectados pelo contato com o doente, nem todos poderão vir a apresentar sintomas ou mesmo adoecer. Isso acontece pela ocorrência de uma imunidade natural, mediada por fatores genéticos, ainda sem explicação pela ciência. É comum um paciente realizar uma radiografia do tórax por algum motivo qualquer, e o resultado do exame mostrar lesões cicatriciais típicas da tuberculose, mesmo com a pessoa negando ter estado doente anteriormente. Ainda é possível uma radiografia torácica com infiltrado típico de tuberculose, sem o paciente apresentar qualquer sintomatologia e com exame direto de escarro negativo[51].

Em grande parte dos casos, se o sistema imunológico da pessoa estiver em condições normais, a bactéria não causará a doença, ficando em estado latente, talvez até por anos. Caso a defesa do organismo fique debilitada, a infecção pode evoluir de um estado latente para um estado de atividade, causando a doença, chamada de TB pós-primária. Na possibilidade da doença se desenvolver logo após o primeiro contato com o bacilo, é então chamada de TB primária.

Uma infecção pulmonar e seus sintomas apresentam muita relação com o germe a qual estão relacionados. Uma pneumonia por estreptococos ou estailococos inicia-se repentinamente com tosse, expectoração purulenta, febre elevada, dispnéia, entre outros. A velocidade dos sintomas explica-se pela rapidez com a qual os germes se multiplicam, apresentando taxas exponenciais a cada 2 horas. Ao realizar a comparação com a bactéria causadora da tuberculose, observa-se uma taxa muito mais lenta, se multiplicando a cada 12 a 20 horas, com a aparição dos sintomas muito a posteriori do momento da infecção. O risco de uma pessoa contaminada só perceber que esta doente após 30-60 dias é grande, acontecendo apenas quando o estado sintomático já está em evidência, o escarro bacilífero, e com tempo suficiente para já ter contaminado outras pessoas [52].

Os sintomas e sinais mais comuns em pacientes com tuberculose consistem de fatores como a tosse persistente produtiva (muco e eventualmente sangue) ou não, febre, sudorese noturna e emagrecimento. É possível encontrar linfadenomegalias no exame físico, que pode ser relacionado tanto a presença de TB extrapulmonar, quanto a existência de coinfeção pelo HIV [53]. Outros marcadores clínicos e indicadores da doença são fraqueza, anorexia, dor torácica moderada, imunodepressão por qualquer causa, contato recente com caso bacilífero e história de tuberculose.

Ao contrário da pneumonia, a febre não é elevada, ocorrendo no período da tarde ou no início da noite, além de apresentar na madrugada uma intensa sudorese. O emagrecimento do doente é rápido e intenso, com casos onde a perda de peso chega a até 6 kg em dois meses. A fraqueza é muito comum e significativa, mas a dor torácica nem sempre é presente e pode ocorrer de forma moderada e contínua. Escarros com sangue, embora com uma frequência mais baixa, são os sintomas que mais rapidamente levam a pessoa ao atendimento médico. A dispnéia é na maioria das vezes discreta ou inexistente. A presença de candidose oral é comum e pode auxiliar no diagnóstico, assim como rouquidão presente por mais de um mês.

2.4 Diagnóstico

A realização de um diagnóstico positivo de um paciente com tuberculose pulmonar, principalmente se possui o tipo bacilífero (o qual elimina bacilos no ar ambiente ao tossir), e tratá-lo corretamente com sucesso, é eliminar uma fonte de infecção. Reduzir as fontes de infecção diminui a taxa de transmissão da doença, e consequentemente, o problema da tuberculose na comunidade. Portanto, quanto mais precoce é feito o diagnóstico menor a chance de disseminação da doença [54].

Considerada um importante fator para o diagnóstico, a tosse pode ocorrer em uma variedade de outras doenças, além de TB pulmonar, como em infecções agudas respiratórias, asma e doença pulmonar obstrutiva crônica (DPOC). Por isso, modelos de orientação para o momento do início da avaliação diagnóstica da TB pulmonar, em áreas com prevalência moderada da doença, utilizam o critério "tempo de tosse",

e não só sua presença, para auxiliar na definição de um paciente como suspeito de TB. Com isso, um suspeito de portador de TB pulmonar é definido como um indivíduo que apresente tosse por duas a três semanas (ou mais) de duração [55].

Pesquisas recentes sugerem que perguntar ativamente sobre tosse para todos os adultos que comparecem à unidades de saúde com sintomas respiratórios, independentemente do tempo de duas semanas como ponto de corte, pode aumentar significativamente o número de casos diagnosticados de TB, especialmente em regiões com elevada prevalência da doença, como o Brasil. Estudos demonstraram que 11% a 18% dos indivíduos com mais de 5 anos de idade, que procuraram unidades de saúde, tinham queixas respiratórias, sendo cerca de 5% classificados suspeito de tuberculose devido à tosse sem motivo aparente, por mais de duas a três semanas [56]. Com isso, é visto como um importante passo para o diagnóstico a investigação de TB em adultos e crianças que apresentem tosse, uma vez que nem todos os pacientes com problemas respiratórios recebem uma avaliação adequada para TB, como indicado em vários estudos [57][58].

Após a identificação dos sintomas característicos, a confirmação do diagnóstico é dada através da identificação da bactéria em material da lesão. Antigamente isso só era possível através de exames bacteriológicos, em especial a cultura. Hoje em dia, com o desenvolvimento de técnicas imunológicas e de métodos de imagem, outros recursos, como a radiografia de tórax, também podem ser utilizados para legitimar o diagnóstico.

Através desses aspectos, é visto que o diagnóstico de tuberculose é baseado em análise de dados de anamnese de exames físicos de pacientes, em conjunto com exames em laboratório, chegando-se com isso a uma conclusão final que determina se a pessoa tem ou não a doença [59]. Entretanto, os exames e observações envolvem incertezas, devido às diferentes características dos pacientes, erro na observação de sintomas ou outros motivos. Com isso, as conclusões obtidas pelos diagnósticos podem ser avaliadas pelas seguintes possibilidades:

- Teste sensível: O resultado do diagnóstico é positivo e o paciente tem a doença.

- Teste falso-positivo: O resultado do diagnóstico é positivo e o paciente não tem a doença.
- Teste falso-negativo: O resultado do diagnóstico é negativo e o paciente tem a doença.
- Teste específico: O resultado do diagnóstico é negativo e o paciente não tem a doença.

Na fase inicial de triagem, um teste com muita sensibilidade é considerado útil para detectar a presença da doença em pacientes. Já testes com grande especificidade excluem a possibilidade de presença da doença em pessoas que não estão infectadas, com utilidade na fase de diagnóstico, após a realização da triagem. O resultado dessa avaliação é importante para auxiliar na prioridade a ser dada entre os pacientes, pois casos com maiores indicações de serem positivos podem levar até mesmo a um possível isolamento do paciente, para evitar a transmissão da doença [60].

Em relação aos métodos de diagnóstico, a baciloscopia é o método mais utilizado mundialmente, e também considerado o principal, para o diagnóstico da tuberculose, com suas características de ser simples, rápido e barato [61]. Ao permitir identificar o indivíduo bacilífero, possui indicação para todos os pacientes sintomáticos respiratórios que apresentam alterações pulmonares na radiografia de tórax. A baciloscopia identifica a bactéria através da análise da amostra clínica (escarro), o qual é preparado para permitir a leitura e a interpretação da lâmina pelo microscopista. A sensibilidade da baciloscopia varia entre 40 e 80%, variação que se deve a fatores como o tipo da lesão, tipo e número de amostras, a atenção e persistência do microscopista e a presença de co-infecção com HIV. A negatividade desse método pode ser influenciada por um número insuficiente de bacilos na amostra, por um estágio ainda inicial da doença ou pela imunidade do indivíduo, que pode manter os níveis da infecção sob controle [5]. Outra característica da baciloscopia é que ela não diferencia as cepas de tuberculose resistentes ao tratamento e também não as discrimina de outras micobactérias atípicas [23].

Já a cultura é o método considerado padrão-ouro para a confirmação da tuberculose [5]. Esse exame permite a multiplicação e o isolamento dos bacilos a partir da semente da amostra clínica, em meios de cultura específicos para micobactérias. Sua sensibilidade para o diagnóstico de tuberculose são altas, detectando de 70% a 89% dos casos [62]. Através de identificação por meio de testes bioquímicos ou moleculares, e a partir da cultura, é possível realizar a diferenciação de cepas sensíveis e resistentes aos medicamentos disponíveis para o tratamento [61]. Apesar de considerada padrão-ouro, a cultura possui o inconveniente de possuir uma demora de 4 a 8 semanas até a obtenção do resultado, sendo com isso pouco utilizada no processo de tomada de decisão clínica [5]. Além disso, a cultura exige infraestrutura maior em biossegurança, que torna difícil sua implementação em locais de poucos recursos [63]. Desse modo, a cultura é indicada a suspeitos de tuberculose que apresentem diagnósticos negativos na baciloscopia, além de ser indicada para os casos extra-pulmonares e também para pacientes com HIV positivos e em crianças. Sua indicação ainda é feita para casos de suspeita de resistência bacteriana aos medicamentos [64].

A radiografia de tórax é utilizada como exame auxiliar, apresentando taxas de sensibilidade de 70% a 80% na presença das anormalidades típicas em adultos imunocompetentes, e 95% se considerarmos outras anormalidades, entretanto, em torno de 10% dos pacientes com AIDS a radiografia pode ser normal [6]. Além de não auxiliar na caracterização de TB resistente, uma grande desvantagem é a logística necessária para se realizar esse procedimento, assim como seu custo.

Outro método existente é o Xpert® MTB/RIF (Cepheid Inc., CA, USA), um teste de amplificação de ácido nucleico, o qual funciona não só para a detecção do bacilo da tuberculose, mas para a triagem de cepas resistentes a medicamentos, mais especificamente, a Rifampicina. O teste realizado pode fornecer resultados em um laboratório em menos de 2 horas, sem necessitar de tratamento da amostra ou de pessoas especializadas em biologia molecular [7]. A OMS recomenda o uso do Xpert® MTB/RIF como primeiro exame em pessoas com suspeita de TB multirre-

sistente (MDR) e na população com HIV [65]. Como desvantagens, para o teste ser realizado é preciso uma fonte de energia elétrica estável e ininterrupta, assim como um computador acoplado para a análise dos dados. O instrumento necessita de, pelo menos, uma calibração por ano, realizada por um técnico treinado. De acordo com estudos realizados, a sensibilidade alcançada por esse teste pode variar, para grupo de amostras com resultado positivo para baciloscopia e para cultura, entre 98% e 100%. Ainda, o teste foi capaz de detectar até 78% de casos negativos à baciloscopia (falso-negativo). Em relação a especificidade, os resultados obtidos variaram de 90,9% a 100% em relação a cultura [66]. Para tentar aumentar a disponibilidade do teste, a OMS obteve uma redução do preço do Xpert® MTB/RIF em 75% para países de baixa e média renda. Entretanto, seu custo ainda permanece alto para a realidade de alguns locais e como barreira para sua utilização em larga escala [67].

Em geral, para diagnosticar casos de tuberculose multirresistente, na maioria dos estudos publicados o critério utilizado para sua definição foram os testes de sensibilidade em cultura. O método prevalente é o das proporções, descrito originalmente por Canetti et al [68] em 1963, que consiste em colocar iguais quantidades com diferentes diluições com medida aproximada do inócuo no controle e meio com o fármaco a ser testado. A resistência ou sensibilidade é estimada comparando o número de colônias que crescem no meio com o fármaco com o número que cresce no meio controle. Os tipos de resistência a TB podem ser resumidos como [69]:

- Natural: decorrente de mutação espontânea, independentemente de exposição prévia a fármacos.
- Inicial: observada no momento em que o paciente se apresenta para tratamento, com resistência a um ou mais fármacos. Inclui os pacientes com resistência primária ou adquirida, sobre os quais não se conhece informações de tratamentos anteriores.
- Primária: observada em pacientes sabidamente não tratados antes, infectados por uma fonte doente com forma resistente.

- Adquirida ou secundária: resultante de uso prévio de medicação de forma inadequada.

O diagnóstico de TB multirresistente ainda é considerado demorado, pelo tempo exigido pelos métodos diagnósticos e pelas limitações e inconsistências nos testes de sensibilidade. Os únicos fármacos testados regularmente para identificação da resistência, e que têm comprovadamente 100% de confiabilidade, são a rifampicina e a isoniazida. Estudos relataram a comparação entre os diagnósticos de 16 laboratórios de referência em micobacteriologia no mundo, e uma grande disparidade entre os resultados alcançados em relação a resistência de fármacos foi vista [70]. Portanto, para auxiliar na tomada de decisão nesses casos quanto ao regime a ser utilizado em pacientes já tratados de TB, deve-se considerar além do teste de sensibilidade, o histórico terapêutico desses pacientes [71].

Como pode ser observado, os métodos básicos descritos para realizar o diagnóstico de TB não apresentam uma grande sensibilidade na detecção da doença. Para conseguir uma taxa mais elevada, são necessários maiores gastos, disponibilidade de equipamentos, além de um intervalo de tempo maior até a obtenção do resultado final. Essas características dificultam suas aplicações em regiões remotas, de difícil acesso, ou com recursos econômicos limitados. Além disso, o tempo necessário para o diagnóstico pode ser considerado demasiado em relação ao objetivo de evitar a disseminação da infecção, representando uma maior gravidade em regiões de grande densidade populacional. Ainda não existe perspectiva em um futuro próximo do desenvolvimento de um teste simples, rápido, barato e de fácil execução, que possa substituir a microscopia em todas as regiões do mundo [67].

2.4.1 Modelos Estatísticos

Nesse cenário, não só para TB, mas como para outras diversas doenças, tem-se observado que o processo de diagnóstico pode ser auxiliado e/ou melhorado através da combinação de testes clínicos com modelos estatísticos. Os modelos são vistos como sistemas que, a partir de bases de dados consistentes, conseguem representar o

problema relacionado ao diagnóstico e, portanto, auxiliar na tomada de decisão dos profissionais de saúde. Atualmente, os modelos estatísticos comumente utilizados no apoio ao diagnóstico de doenças, inclusive TB, são:

- Regressão Logística;
- Redes Bayesianas;
- Árvores de Decisão;
- Redes Neurais Artificiais.

A regressão logística tem como objetivo produzir, a partir de um conjunto de observações, um modelo que permita a predição de valores de variáveis categóricas, as quais são frequentemente binárias, a partir de uma série de variáveis explicativas contínuas e/ou binárias. Essa técnica é bastante utilizada na área da medicina para identificação de doenças de interesse. A limitação dessa técnica se encontra quando não existe uma grande base de dados a ser utilizada na geração do modelo [72].

Em relação às redes bayesianas, elas consistem em representar, em forma de grafos, um conjunto de variáveis aleatórias e suas dependências condicionais, mostrando, desse modo, a relação probabilística entre doenças e sintomas. Através dessa técnica, é possível identificar as relações de causa e efeito entre as variáveis de estudo [73], o que pode ser utilizado em ferramentas para apoiar o diagnóstico médico. Apesar disso, a complexidade em se especificar uma rede bayesiana a torna difícil para aplicações em que se encontrem muitas variáveis, já que a técnica requer os valores de todas as probabilidades condicionais e, antecipadamente, os valores das variáveis explicativas.

Já as árvores de decisão representam uma técnica que particiona, de forma recursiva, um conjunto de dados, até que cada subconjunto do particionamento contenha casos iguais, formando um modelo graficamente estruturado, compacto e que descreve, de forma clara, a classificação do conjunto de variáveis [74]. Em problemas mais complexos, árvores de decisão maiores tendem a ser necessárias, onde podem

surgir sub-árvores idênticas em caminhos diferentes. Como consequência, quanto maior o número de decisões a serem tomadas, mais nós são necessários percorrer e menos acurados são os resultados.

A técnica de rede neural artificial é utilizada para o reconhecimento de padrões, onde os modelos gerados se inspiram na estrutura e no funcionamento do cérebro humano, baseando-se em interconexões de unidades, chamadas de neurônios. As redes neurais podem ser treinadas de forma supervisionada, onde o conhecimento das respostas do modelo para cada entrada do sistema é necessário, ou não-supervisionado, na qual as respostas são completamente desconhecidas do sistema [75]. Essa técnica é uma realidade na área médica, sendo eficiente em diversas áreas da medicina, inclusive no auxílio ao diagnóstico, prognóstico, e reconhecimento de padrões biológicos, exames clínicos e imagens médicas. Uma desvantagem dessa abordagem seria o ajuste de parâmetros necessário para a realização do modelo, que depende de métodos heurísticos onde a solução é única para cada problema.

Os modelos descritos acima podem auxiliar no apoio ao diagnóstico de tuberculose, na triagem de pacientes, prognóstico e no direcionamento para o tratamento, facilitando a ação dos profissionais da área médica. A rapidez na detecção, assim como eficiência, custo-efetivo, e facilidade na interpretação dos resultados, são fatores fundamentais para a aplicabilidade dos modelos, se tornando de grande utilidade em lugares com alta carga de TB. Apesar de não terem sido encontrados trabalhos que utilizam modelos estatísticos para apoiar o diagnóstico de tuberculose com resistência à medicamentos, diversas metodologias são sugeridas na literatura para o diagnóstico de TB, sendo algumas apresentadas a seguir.

Bock et al. (1996) [76] utilizaram um modelo logístico multivariado para analisar os dados de 295 pacientes colocados em isolamento respiratório durante 3 meses, para tentar descobrir variáveis importantes para o diagnóstico de tuberculose. Como resultado, 5 variáveis foram identificadas com relevância preditiva, sendo elas achados radiológicos no lóbulo superior dos pulmões ou cavitações, histórico de contato com doentes, prova tuberculínica positiva e a não utilização da terapia preventiva

com isoniazida. Se tivessem utilizados esses fatores para decidir sobre o isolamento dos pacientes, o número de pessoas sem tuberculose colocadas nessa situação teria caído de 253 para 95, porém 8 dos 42 pacientes com tuberculose teriam sido avaliados erroneamente.

Gordin et al. (1996) [77] avaliaram 1.373 casos de tuberculose ocorridos entre janeiro de 1992 a junho de 1994 em diferentes regiões dos Estados Unidos para determinar o impacto da AIDS em relação a tuberculose multirresistente. Através de uma análise por regressão logística multivariada, o HIV mostrou ser um fator de risco para a TB-MDR, independente da localização geográfica, históricos de tratamentos anteriores, idade e etnia.

Samb et al. (1997) [78] utilizaram regressão logística multivariada para identificar sintomas importantes para o diagnóstico de TB, tendo como resultado a tosse por mais de 21 dias, dores no peito por mais de 15 dias, expectoração e dispnéia. Com duas quaisquer das quatro variáveis alcança-se sensibilidade de 85% e especificidade de 67%, já com três das quatro variáveis esses valores se alteram para 49% e 86%, para sensibilidade e especificidade, respectivamente.

Veropoulos et al. (1998) [79] utilizaram redes neurais para desenvolver um sistema para analisar imagens de radiografias e identificar o bacilo responsável pela TB através do reconhecimento de padrões. Uma acurácia de 97,9% foi atingida, com sensibilidade de 94,1% e especificidade de 99,1%. No ano seguinte, El-Solh et al. (1999)[80] desenvolveram uma rede neural que utilizava informações clínicas e radiográficas para previsão de TB pulmonar ativa. Foram utilizadas as variáveis idade, diabetes mellitus, AIDS, dor torácica, emagrecimento, tosse, sudorese, febre, dispnéia e achados em imagens de radiografias.

Mello (2001) [81], através de um modelo de regressão logística multivariada e árvore de classificação, conseguiu para a detecção de TB sensibilidade de 65,9% e especificidade de 60,1% com o primeiro método, e sensibilidade de 64,2% e especificidade de 60,6% com as árvores. Para isso foram utilizadas informações sobre achados radiológicos, presença de escarro, emagrecimento e idade de pacientes atendidos na

Rede de Saúde Municipal do Rio de Janeiro.

Santos (2003) [82] e Santos et al.(2003) [83] utilizaram modelos com redes neurais artificiais e árvores de classificação para analisar a presença de tuberculose paucibacilar entre 126 pacientes do Hospital Universitário Clementino Fraga Filho, onde 59 eram casos positivos. As variáveis utilizadas compreenderam características demográficas como gênero, idade e renda familiar, clínicas incluindo tosse, febre, sudorese, emagrecimento, anorexia, e fatores de risco como diabetes, alcoolismo, AIDS, entre outros. Com as redes neurais foi obtida uma sensibilidade de 73% e uma especificidade de 67%, e com as árvores de classificação sensibilidade de 40% e a especificidade de 70%. A partir da filtragem dos dados utilizados considerando sua relevância dentro dos modelos, um conjunto de 12 variáveis apresentou resultados de sensibilidade em 83% e uma especificidade de 71%, sendo elas: AIDS, internação hospitalar, TB extrapulmonar, fumante, anorexia, dispneia, emagrecimento, febre, sudorese, hemoptise, tosse e idade.

Benfu et al. (2009) [84] utilizaram um modelo neural para diagnosticar TB utilizando um conjunto total de 29 variáveis, dentre informações pessoais, histórico médico, achados radiográficos e exames laboratoriais. O modelo obteve sensibilidade de 88,9% e especificidade de 100%, considerando um total de 589 pacientes sendo 291 casos com diagnóstico positivo. Aguiar et al (2012)[85] utilizaram árvores de classificação com variáveis de entrada como radiografias de tórax, perda de peso e dispneia para predição de TB em pacientes hospitalizados, além de auxiliar em um uso mais racional de salas de isolamento, sem aumentar o risco de transmissão. A base de dados tinha um total de 215 pacientes, obtendo como resultados uma especificidade de 76% e 60% de sensibilidade.

Cascão (2011) [86] desenvolveu modelos de redes neurais para a triagem e diagnóstico de TB pulmonar. Foram utilizadas redes neurais não-supervisionadas, no caso mapas auto-organizáveis, para analisar as variáveis relevantes e obter uma triagem dos pacientes em um posto de saúde. Métodos para o diagnóstico, além da análise de grupos de risco, também foram realizados. O melhor modelo alcançou

um resultado de 81,4% de sensibilidade e 61,3% de especificidade.

Seixas et al. (2012) [87] propuseram um sistema de suporte à decisão baseado em redes neurais para auxiliar no direcionamento de pacientes para áreas de isolamento. O sistema identifica a probabilidade do paciente ter tuberculose pulmonar e utilizou registros de 290 pessoas que foram admitidas em salas de isolamento no Hospital. Seixas et al. (2013) [88] utilizaram redes neurais para auxiliar em um diagnóstico não-invasivo de TB pleural. O histórico clínico e a presença do vírus HIV de 137 pacientes foram utilizados como informações, alcançando um resultado para a prevalência de uma acurácia de mais de 90%.

Capítulo 3

Base de Dados

A base de dados utilizada nas análises desse trabalho foi constituída de informações de pacientes de diferentes unidades de saúde e hospitais do Brasil, coletadas dentre os anos de 2011 a 2013, sendo elas: no Ceará, o Hospital de Messejana Doutor Carlos Alberto Studart Gomes; em São Paulo, o Instituto Clemente Ferreira; no Rio Grande do Sul, a Unidade Básica de Saúde do bairro de Partenon em Porto Alegre; e, no Rio de Janeiro, o Centro de Referência Professor Hélio Fraga (Fiocruz), o Hospital Municipal Raphael de Paula Souza, o Instituto Estadual de Doenças do Tórax Ary Parreiras e unidade de saúde do bairro de Paciência. Outros pacientes que apresentavam triagem realizada em outros estados, como Minas Gerais, Paraíba, Paraná e Rio Grande do Norte, também foram incluídos. Porém, pela baixa frequência nesses casos, apenas uma ocorrência por estado, não foram considerados para o estudo de regionalidade, a ser posteriormente descrito.

Através de formulários, foram coletadas informações clínicas dos pacientes por um questionário padrão sobre os sinais e sintomas da doença, fatores de risco para tuberculose, história de tuberculose e de outras doenças prévias e dados sócio-demográficos, além de ter sido incluído o acompanhamento sobre o resultado de exames realizados. Devido ao objetivo desse estudo, somente os formulários de pacientes que já apresentavam, na triagem, suspeita de tuberculose com resistência, foram selecionados.

Para ser possível o desenvolvimento e avaliação dos modelos matemáticos pre-

sententes nesse trabalho, foi preciso estabelecer um diagnóstico conclusivo de cada paciente. Portanto, através de análises sobre as informações contidas nos dados de cada paciente, cada diagnóstico foi estabelecido pelas seguintes condições:

- TB negativo: Paciente apresentou resultado de exame de cultura negativo.
- TB sensível: Paciente apresentou resultado de exame de cultura positivo.
- TB-DR: Paciente apresentou resultado de cultura positivo, e teste de sensibilidade apontou resistência para drogas anti-TB, exceto a combinação de Rifampicina e Isoniazida.
- TB-MDR: Paciente apresentou resultado de cultura positivo, e teste de sensibilidade apontou resistência para drogas anti-TB, incluindo a combinação de Rifampicina e Isoniazida.

Além disso, dentre todas as informações existentes para cada paciente, foi selecionado um conjunto de variáveis de acordo com o seu nível de dependência com a TB pulmonar. As variáveis selecionadas foram baseadas em estudos sobre modelos de classificação entre pacientes com diagnóstico positivo e negativo para tuberculose. Com isso, é possível analisar a possibilidade de discriminação dos tipos resistentes da doença ao se utilizar o mesmo conjunto de informações desses modelos de diagnósticos. A lista de variáveis utilizadas neste trabalho, e sua codificação correspondente para o desenvolvimento dos modelos, são descritas na Tabela 3.1.

Já a Tabela 3.2 exibe a quantidade de pacientes por estado onde foi realizada a triagem, assim como sua distribuição por diagnósticos. Com um total de 560 pacientes suspeitos de TB resistente, é possível observar que a maior quantidade está no estado do Rio de Janeiro, seguido pelo Ceará. Outra observação pode ser feita em relação a grande parcela de pacientes que apresentam algum tipo de tuberculose, em comparação aos que possuem diagnóstico negativo. Isso se deve principalmente pela escolha da utilização de dados que correspondem à pacientes que já chegaram a unidade de saúde ou hospital apresentando suspeita de tuberculose com resistência.

Variável	Codificação
Idade	anos
Emagrecimento	
Tosse	
Hemoptise	
Sudorese noturna	ausência = -1, presença = 1, ignorado = 0
Febre	
Perda de apetite	
Dispneia	
Tabagismo	
Ex-fumante	não = -1, sim = 1, ignorado = 0

Tabela 3.1: Variáveis utilizadas e sua respectiva codificação..

Pacientes	CE	RJ	RS	SP	Outros	Total
TB Negativo	38	162	32	12	4	248
TB Sensível	21	34	20	3	0	78
TB-DR	38	43	21	18	0	120
TB-MDR	33	54	5	22	0	114
Total	130	293	78	55	4	560

Tabela 3.2: Quantidade de pacientes por diagnóstico e local de triagem médica.

É possível encontrar respostas dos pacientes a certos sintomas descritos como ignorados. Esses casos podem ocorrer pelo fato do paciente ter escolhido não responder a pergunta por se sentir constrangido, não saber ou não ter certeza da resposta. A frequência das respostas para cada variável, com exceção da idade, é vista na Tabela 3.3, onde se pode observar que existem respostas ignoradas para 8 das 9 variáveis, porém em baixa representatividade.

Variável	Sim (%)	Não (%)	Ignorado (%)
Emagrecimento	56,4	42,8	0,53
Tosse	81,4	17,5	1,07
Catarro com sangue	21,96	76,96	1,07
Suor noturno	52,32	43,92	3,75
Febre	53,03	45,17	1,78
Perda de apetite	55,71	43,21	1,07
Falta de ar	54,28	43,57	2,14
Fumante	44,82	54,1	1,07
Ex-fumante	27,14	72,85	0

Tabela 3.3: Frequência de respostas dadas para as variáveis utilizadas por todo o conjunto de pacientes.

Capítulo 4

Método

Essa seção apresenta as metodologias utilizadas para cumprir os objetivos propostos. Inicialmente, é feita uma análise das variáveis a serem utilizadas nos estudos, para então ser descrito o procedimento para o apoio ao diagnóstico de pacientes através de modelos matemáticos. Em seguida, o método para a análise dos dados de pacientes por similaridade estatística é exposto, o qual visa a identificação de padrões de interesse no conjunto de dados, assim como o estudo das variáveis por relevância.

4.1 Análise das variáveis utilizadas

Foram selecionadas 10 variáveis explicativas, sendo uma numérica e as outras qualitativas. Para a variável idade, a qual é numérica, histogramas foram gerados para analisar sua distribuição por conjunto de diagnóstico, verificando sua média e desvio padrão, como mostra a Figura 4.1. É possível observar que a média das distribuições não apresenta grandes variações entre os conjuntos, mas os pacientes com diagnóstico de TB apresentam médias menores que a do conjunto total de pacientes, tendência que é acompanhada pelos pacientes com diagnósticos de TB-DR e TB-MDR.

Para as variáveis qualitativas, foi feito um estudo para calcular a razão de chance (em inglês, *odds ratio*) [89] das mesmas. Esse cálculo possui uma simples interpre-

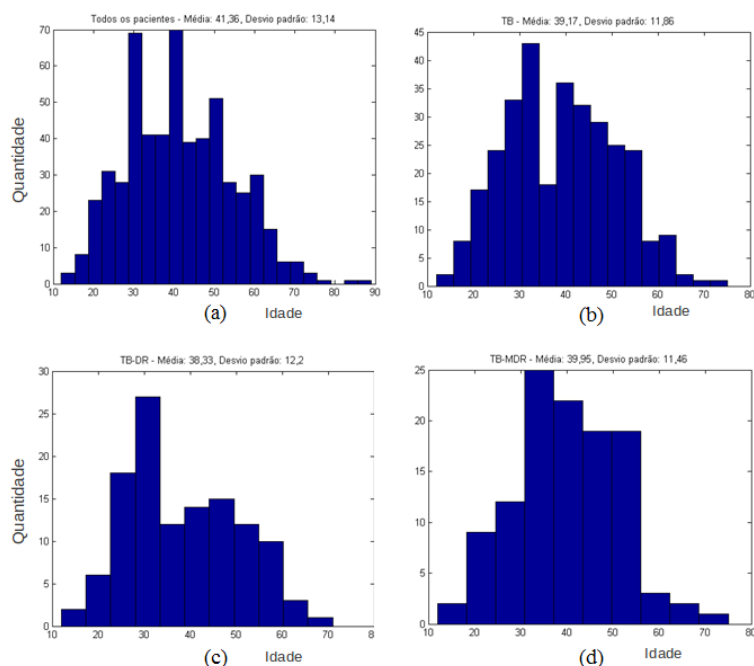


Figura 4.1: Distribuição das idades dos pacientes no (a) conjunto completo de dados, (b) com diagnóstico de TB positivo, (c) TB droga-resistente e (d) multirresistente.

tação, a qual indica o nível do relacionamento de um sintoma, ou fator, com uma doença. Uma razão de chances com valor igual a 1 significa que o sintoma em estudo possui probabilidade igual de ocorrer nos dois grupos analisados, enquanto um valor maior que 1 significa que o sintoma tem maior probabilidade de ocorrer no caso considerado como foco da análise. Caso o valor encontrado seja menor que 1, o sintoma tem maior probabilidade de ocorrer no caso fora do foco da análise. Foram realizados estudos entre os diagnósticos de TB multirresistente em comparação com TB droga-resistente, TB resistente (a qual inclui todas as formas de resistência) contra TB sensível, e TB positivo (incluindo os casos sensíveis e resistentes) contra diagnósticos negativos para TB. Os resultados podem ser vistos na Tabela 4.1.

Ao analisar os resultados das chances entre TB-MDR e TB-DR, é possível observar que os sintomas tosse, hemoptise, sudorese noturna, perda de apetite, além do fato do paciente ser ex-fumante, são bem próximas do valor 1, o que significa uma igual probabilidade de ser presente nos dois grupos. Porém, emagrecimento exibe uma maior probabilidade de acontecer em casos de TB-MDR, enquanto dispneia e o fator tabagismo apresentam uma tendência ao grupo de TB-DR.

OR	TB-MDR x TB-DR
Emagrecimento	1.826255
Tosse	1.090909
Hemoptise	1.063291
Sudorese noturna	1.187291
Febre	1.067041
Perda de appetite	0.913462
Dispneia	0.580091
Fumante	0.697218
Ex-fumante	0.958580
OR	TB resistente x TB sensível
Emagrecimento	0.573477
Tosse	1.194444
Hemoptise	1.515698
Sudorese noturna	0.517241
Febre	0.824798
Perda de appetite	0.754870
Dispneia	0.511930
Fumante	0.634415
Ex-fumante	0.842318
OR	TB positivo x TB negativo
Emagrecimento	4.079304
Tosse	6.226824
Hemoptise	2.254111
Sudorese noturna	2.942415
Febre	3.521368
Perda de appetite	3.069711
Dispneia	1.319244
Fumante	0.536164
Ex-fumante	2.001465

Tabela 4.1: Razão de chance para diagnósticos de TB multirresistente em comparação com TB polirresistente, TB resistente (a qual inclui todas as formas de resistência) contra TB sensível, e TB positivo (incluindo os casos sensíveis e resistentes) contra diagnósticos negativos para TB.

Enquanto isso, pelos resultados da razão de chances entre os grupos de TB resistente e TB sensível, apenas os sintomas tosse e hemoptise apresentam valores maiores que 1, ou seja, mais prováveis no grupo de tuberculose resistente. Isso pode indicar que os pacientes com casos com resistência tendem a ter esses sintomas como seus maiores indícios para seu diagnóstico. Já pelos resultados das análises entre os grupos de TB positivo e TB negativo, é visto a ocorrência dos sintomas característicos de TB muito mais prováveis de ocorrer em pacientes com algum tipo da doença. Esse resultado mostra aderência ao que é esperado pelas práticas médicas.

4.2 Redes neurais para apoio ao diagnóstico

Um dos principais objetivos deste trabalho é apontar qual o diagnóstico mais provável para um paciente em fase de triagem médica, apoiando com isso, a decisão médica. A técnica de redes neurais multicamadas foi escolhida para realizar essa tarefa, sendo descrito os detalhes e as escolhas feitas nas fases de pré-processamento dos dados, treinamento e na avaliação do desempenho dos resultados.

Além disso, deve-se investigar padrões entre as estatísticas dos pacientes, como grupos de risco e as respectivas contribuições de cada variável utilizada no resultado. Para isso, técnicas de inteligência computacional para agrupamento não-supervisionado dos dados podem ser aplicadas, como mapas auto-organizáveis, os quais foram utilizados nesse estudo, e os detalhes de seu processo de treinamento e avaliação dos resultados são abaixo discutidos.

4.2.1 Redes neurais multicamadas

4.2.1.1 Pré-processamento

O objetivo final do problema de classificação com redes neurais multicamadas é encontrar os pesos das sinapses em que os dados de entrada gerem saídas que sejam classificadas da forma mais correta possível. Para isso, existe a fase de treinamento,

onde os pesos são ajustados até chegarem a seu valor final. Entretanto, alguns cuidados devem ser tomados antes de realizar esse procedimento.

Se a rede for treinada com o conjunto total de dados, a configuração dos melhores pesos pode se tornar tendenciosa a acertar os dados aos quais a rede foi apresentada. Com isso, a classificação de novos dados que sejam apresentados à rede, os quais não foram utilizados para treino, pode ter sua eficiência prejudicada. Para evitar essa situação, deve-se separar os conjuntos de dados em conjunto de treino, utilizado para o ajuste dos pesos das sinapses da rede, e teste, para avaliar a qualidade de classificação.

O *over-fitting*, o qual é um problema de generalização, ocorre basicamente quando após um certo período de treinamento, a rede se especializa em um conjunto de dados e perde a capacidade de generalização. Para evitar esse problema, pode ser usado também um conjunto de dados para validação, com a finalidade de verificar a eficiência da rede quanto a sua capacidade de generalização no decorrer do processo de treinamento. Esse fator pode ser utilizado como critério de parada da rede. Porém, na área médica, é comum termos uma quantidade limitada de dados, o que dificulta a divisão dos dados entre treino, teste e validação. Nesses casos, como alternativa, utiliza-se o conjunto de teste como conjunto de validação, que foi o procedimento utilizado nesse trabalho.

Além disso, os dados apresentam flutuações estatísticas que devem ser levadas em consideração. Caso contrário, ao realizar a separação simples dos dados entre treino e teste, é possível que alguma determinada característica de um grupo de eventos fique representada apenas em um desses conjuntos, prejudicando o processo de aprendizado do modelo. Para amenizar esse problema, é utilizada a técnica da validação cruzada. Essa técnica, quando aplicada às redes neurais, as quais possuem vários parâmetros de configuração a serem escolhidos, pode ser utilizada para externalizar as possíveis flutuações estatísticas. Dentre suas opções de implementação, nesse trabalho foi utilizada a técnica de *k-fold*.

A técnica de validação cruzada *k-fold* consiste na divisão do conjunto total de

dados entre subconjuntos de quantidade similar entre si. Com isso, é possível formar diversas combinações entre os subconjuntos para obter diferentes arranjos para treino e teste. Um ilustração da técnica pode ser vista na Figura 4.2. Para o presente estudo, foram utilizados um total de 12 subconjuntos, sendo 9 para treino e 3 para teste. O procedimento é realizado até que todas as combinações de subconjuntos entre treino e teste sejam utilizadas, no caso sem repetição. Nesse caso, um limite de 50 combinações possíveis foram utilizadas para as análises.

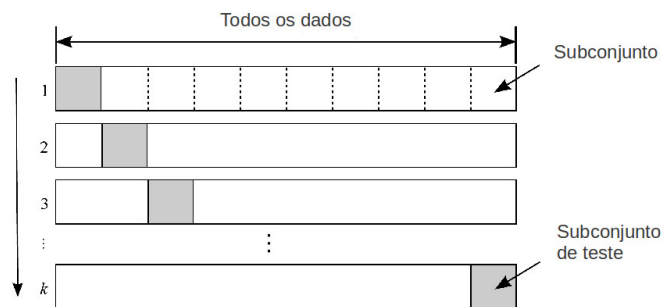


Figura 4.2: Ilustração da técnica de validação cruzada pelo método k-fold.

Os dados foram distribuídos entre os subconjuntos após a aplicação da técnica de *k-means*. Essa técnica é um método particional, amplamente utilizada pela comunidade científica, que utiliza o conceito de centroides. Em um conjunto de dados, definem-se N centroides (ou pontos de coordenada) espalhados aleatoriamente, e calcula-se a distância euclidiana quadrática dos eventos a cada um desses pontos. Os eventos mais próximos de cada centróide, em comparação com os outros, são considerados pertencentes ao mesmo agrupamento. Com isso, os pontos dos centróides são recalculados para serem o baricentro dos agrupamentos. O processo é repetido até não haver variação significativa da localização dos centróides. Ao utilizar essa técnica, visa-se separar previamente o conjunto de dados por sua similaridade estatística, para distribuí-lo de forma equilibrada e tentar garantir a variedade estatística de cada subconjunto da validação cruzada, e com isso, dos conjuntos de treino e teste. A Figura 4.3 mostra um exemplo de particionamento realizado por *k-means*, onde os pontos pretos representam os eventos, e o ponto azul o centróide final após algumas iterações.

Ainda na fase de pré-processamento dos dados, a idade dos pacientes foi norma-

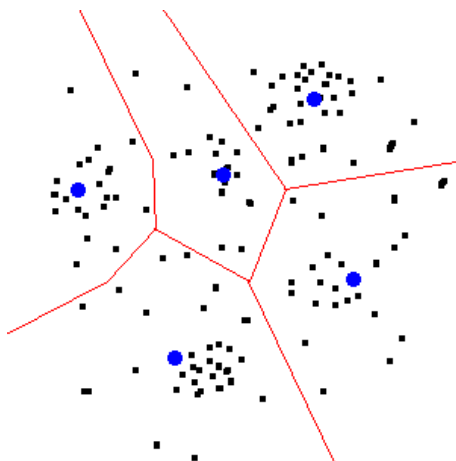


Figura 4.3: Exemplo de particionamento de um conjunto de dados utilizando a técnica de k-means. Os pontos pretos representam os eventos, e os pontos azuis os centróides após algumas iterações.

lizada para ter desvio padrão unitário e média zero, se tornando desse modo mais equilibrada com relação aos valores das outras variáveis, as quais são categóricas, com valores 1, 0 e -1. Esse procedimento auxilia a rede em seu treinamento e tende a gerar resultados melhores. Conjuntos de treino, que possuam classes de dados com quantidades de eventos desbalanceadas, também podem afetar o treinamento. Para diminuir esse problema, eventos escolhidos de forma aleatória foram replicados dentro do conjunto de menor representatividade, tendo como objetivo igualar as quantidades totais de cada classe.

4.2.1.2 Treinamento

Para escolher a melhor configuração de uma rede neural multicamadas, diversas redes devem ser treinadas e avaliadas, na busca por um resultado que alcance a melhor generalização possível. A escolha do número de neurônios da camada escondida pode ser realizada através de experimentos, por não possuir um método exato para sua escolha na literatura. Portanto, as topologias, nesse trabalho, foram variadas para possuir dentre 1 a 20 neurônios na camada escondida, existindo preferência por redes com poucos neurônios, as quais tendem a possuir uma melhor generalização, reduzindo o problema de *overfitting* quando há restrição estatística no conjunto de dados.

O algoritmo utilizado para treinamento foi o *Resilient Backpropagation*, o qual está sujeito a mínimos locais. Para evitar esse cenário, foram realizadas diversas inicializações para cada classificador proposto, e selecionando-se apenas o melhor resultado. Para esse trabalho, foram feitas 100 inicializações para cada configuração de topologia. O algoritmo ainda requer o ajuste de parâmetros, como número de iterações do algoritmo, critério de parada, inicialização dos pesos e taxa de aprendizado. Esses parâmetros influenciam diretamente na qualidade da generalização da rede treinada.

A inicialização dos pesos sinápticos é considerada um fator de importância no processo de treinamento, sendo aqui utilizados valores aleatórios entre -0,2 e 0,2. Já a taxa de aprendizado escolhida depende da função a ser aproximada. A escolha de uma taxa muito pequena torna o treinamento lento, já valores muito grandes podem acarretar em sua divergência. Para os modelos gerados nesse trabalho, o valor inicial de 0,01 foi utilizado. Quando o sinal da derivada da função objetivo, em relação a um peso sináptico, alternava por diversas iterações consecutivas, a taxa era diminuída em 30%. Se o sinal permanece o mesmo, a taxa de aprendizado era aumentada em 5%.

Para alcançar a melhor generalização possível, o critério de parada do treinamento da rede deve considerar sua capacidade de generalização. Para isso, a cada iteração de treinamento, a rede é colocada em operação com o conjunto de dados de validação, sendo o valor do erro de classificação monitorado. Ao final do número determinado de iterações, os pesos utilizados serão aqueles que obtiveram o menor valor de MSE de validação, seguindo o critério *Save the Best*. A Figura 4.4 mostra um exemplo de evolução de treinamento e o critério de parada.

A camada de saída das redes foi projetada para possuir apenas um nó, no caso de um discriminador para duas classes, ou três nós, no caso de três classes. Os valores de classificação dos nós de saída são variados entre 1 (indica que o paciente possui a doença) e -1 (o paciente não apresenta a doença). No caso de um discriminador com três nós na camada da saída, a classe é indicada pelo nó que apresentou o maior

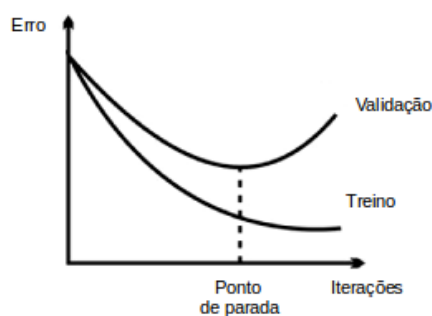


Figura 4.4: Representação da evolução do erro de classificação de um treinamento por iteração e escolha do ponto de parada com a melhor generalização.

valor de probabilidade entre os nós de saída. Ainda, por sempre possuir saídas nessa faixa de valores, a função de transferência utilizada foi a tangente hiperbólica.

A Tabela 4.2 mostra a lista dos parâmetros utilizados para o treinamento das redes neurais.

Parâmetro	Valor
Número máximo de épocas	10.000
Gradiente mínimo	1e-10
Taxa de aprendizagem	0,01
Número de inicializações	100

Tabela 4.2: Valores de parâmetros utilizados para o treinamento das redes neurais.

4.2.1.3 Avaliação do desempenho

Para realizar a avaliação do desempenho da classificação do modelo, deve-se aplicar o conjunto de dados como informação de entrada na rede, obtendo-se com isso índices, como a taxa de acerto do classificador. Entretanto, para um classificador neural, é recomendável que se aplique a fração do conjunto de dados que não foi utilizada em sua fase de treinamento, ou seja, o conjunto de teste. Com isso, evita-se que a rede alcance taxas mais altas de acerto sobre a classe esperada apenas pelo fato de sua generalização não ter sido tão eficaz.

Além da acurácia, ou seja, a taxa de classificação correta atingida pela aplicação do conjunto de teste, outras estatísticas também são utilizadas para medir desempenho, como a sensibilidade e especificidade. O valor preditivo positivo (VPP), o

qual é a probabilidade de um indivíduo ter uma determinada doença dado que ele é positivo a um teste diagnóstico, e o valor preditivo negativo (VPN), que é definido como a probabilidade de um indivíduo não ser doente dado que tem teste diagnóstico negativo, também são avaliados. Todos esses índices são indicados pelos seus valores médios e com suas respectivas incertezas, representando os valores que devem ser esperados como desempenho de classificação.

Um modo de se analisar a relação entre sensibilidade e especificidade de um modelo é através da curva ROC (*Receiver Operating Characteristic*) [90], a qual permite ver os valores para os quais existe maior otimização da sensibilidade em função da especificidade, que corresponde ao ponto que se encontra mais próximo do canto superior esquerdo do diagrama visto na Figura 4.5.

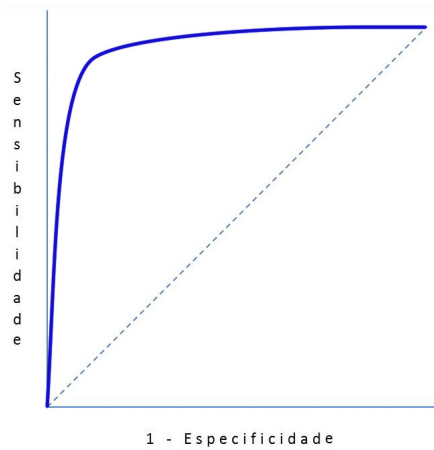


Figura 4.5: Exemplo de curva ROC para medir a relação entre sensibilidade e falsos positivos.

O valor do ponto de corte que define a sensibilidade e especificidade a serem utilizadas pode ser selecionado arbitrariamente entre os valores possíveis para a variável de decisão, acima do qual o paciente é classificado como positivo (teste positivo, paciente com a doença) e abaixo do qual é classificado como negativo (teste de diagnóstico negativo, ausência de doença), como visto na Figura 4.6. Para cada ponto de corte são calculados valores de sensibilidade e especificidade, que podem então serem dispostos no gráfico como um ponto da curva ROC. Um classificador perfeito corresponderia a uma linha horizontal no topo do gráfico, porém dificilmente é alcançado. Na prática, curvas consideradas boas estarão entre a linha diagonal e a linha

perfeita, onde quanto maior a distância da linha diagonal, melhor o sistema. A linha diagonal indica uma classificação aleatória, ou seja, um sistema que aleatoriamente seleciona saídas como positivas ou negativas.

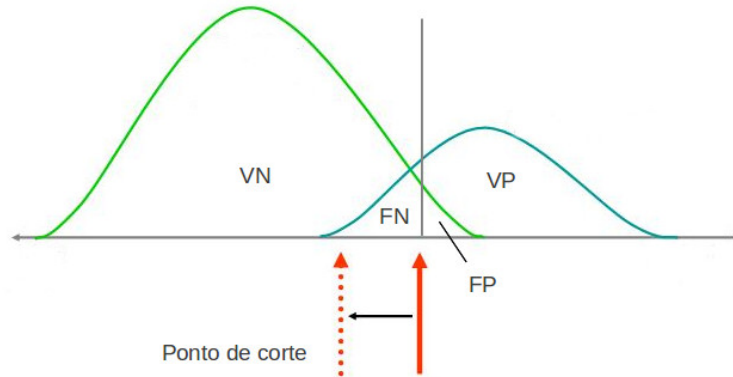


Figura 4.6: Distribuição de pacientes com diagnósticos positivos e negativos e exemplo de ponto de corte para definição da sensibilidade e especificidade.

Os resultados de cada índice e o ponto de corte escolhido devem ser avaliados de acordo com o objetivo final da rede. Nesse trabalho, um dos critérios utilizados para auxiliar na escolha da melhor rede é o índice SP, cujo cálculo é definido por [91]

$$SP = \sqrt{\frac{S + E}{2}} \times \sqrt{S \times E} \quad (4.1)$$

onde S representa a sensibilidade e E a especificidade. Esse índice visa representar o ponto de corte ótimo para um classificador [92]. Para cada ponto de corte um índice SP percentual é calculado, e onde esse valor for máximo indica o ponto de eficiência alta.

Como descrito anteriormente, são experimentadas topologias que variam entre 1 e 20 neurônios na camada escondida, e para cada uma dessas configurações, 50 redes são modeladas por combinações diferentes de conjuntos de dados obtidas pela validação cruzada. Multiplicando-se a isso um total de 100 inicializações realizadas para cada rede, das quais apenas as que obtiveram o melhor índice SP são consideradas, ao todo 100 mil redes são geradas.

Através desses resultados, a topologia que apresenta a melhor média do índice

SP, ao se propagar o conjunto de teste, é a escolhida como o número ótimo de neurônios na camada escondida. Os índices de desempenho para as 50 redes com essa topologia, obtidas pela validação cruzada, são avaliados para representar os valores de sensibilidade, especificidade, acurácia e índice SP médios, assim como suas incertezas, as quais são referentes ao seu valor RMS.

Entretanto, a utilização de um conjunto de redes, em uma situação operacional, pode ser uma tarefa complexa. Por isso, uma só rede, dentre o conjunto que apresenta a melhor topologia, é selecionada para um cenário de operação. Para isso, é aplicado todo o conjunto de dados (treino e teste) nas redes de interesse, sendo a rede de operação escolhida como a que apresentou o maior índice SP, sensibilidade ou especificidade, cada uma com um propósito diferenciado. Para efeitos de triagem, a rede com melhor sensibilidade é a indicada, pois possui uma maior taxa de classificação correta dos pacientes doentes, auxiliando em seu encaminhamento posterior. Já a rede com maior especificidade é indicada para a fase de diagnóstico, pois classifica melhor os pacientes não doentes, evitando, por exemplo, gastos desnecessários com medicamentos. A rede com melhor SP apresenta o ponto ótimo de equilíbrio entre sensibilidade e especificidade, sendo utilizada nesse trabalho para avaliação de desempenho do modelo.

Após selecionar a rede de operação com o melhor índice SP, é analisado como cada variável contribuiu para o resultado obtido. Para isso, cada variável é analisada separadamente, sendo ela definida como seu valor médio em todas as entradas de pacientes, o que retira seu impacto no resultado, para então o conjunto completo de dados ser propagado pela rede que havia obtido o melhor SP. Os valores da diferença entre o SP da rede de operação e o SP alcançado pelo conjunto gerado para a análise das variáveis é avaliado, sendo visto nesse caso que quanto menor o valor obtido, menor é a sua relevância para o resultado

4.2.2 Mapas Auto-organizáveis

4.2.2.1 Pré-processamento

Assim como para o treinamento das redes neurais multicamadas, a idade dos pacientes é normalizada com variância unitária e média zero, para tornar seus valores próximos aos que representam as variáveis categóricas. O resultado do treinamento tende a ser melhor através desse procedimento.

4.2.2.2 Treinamento

No SOM, existem, além da função de vizinhança, quatro parâmetros que devem ser decididos antes de iniciar um treinamento: o número de neurônios, a dimensão, o formato e o tipo de grade do mapa. A granularidade do mapa é definida através do número de neurônios, o qual deve ser o maior possível, onde o tamanho da vizinhança irá controlar a suavização e generalização do mapa. Porém, um mapa muito grande pode prejudicar o treinamento, pois muito neurônios podem nunca ser ativados, o que aumenta o peso computacional do treinamento, além de apresentar uma granularidade incompatível ao conjunto de dados.

Antes do treinamento, os pesos sinápticos dos neurônios do mapa são inicializados de forma aleatória, com os pesos selecionados através de uma distribuição uniforme com valores máximos e mínimos dados pelos respectivos valores da base de dados após o pré-processamento. Todos os mapas treinados foram bidimensionais, sendo a grade utilizada hexagonal, que faz a distância ser igual entre os neurônios vizinhos. Além disso, a função de vizinhança escolhida é a gaussiana, a qual possui características de interesse, como: simetria topológica com relação ao ponto máximo, ou seja, do neurônio vencedor em relação a sua vizinhança; a amplitude da função diminui monotonicamente com o aumento da distância ao ponto máximo, condição necessário para a convergência; e possui independência da localização do neurônio vencedor.

As dimensões do mapa são determinadas a partir do seguinte critério [93]:

$$\begin{cases} n_x n_y = 5\sqrt{n_d} \\ \frac{n_y}{n_x} = \sqrt{0.75} \frac{\lambda_1}{\lambda_2} \end{cases}$$

onde n_x é o número de colunas e n_y o de linhas na grade, n_d o número de amostras disponíveis para treinamento e λ_1 e λ_2 os dois maiores autovalores da análise PCA [94] para esse conjunto de dados.

4.2.2.3 Identificação de agrupamentos e avaliação das variáveis

A qualidade dos agrupamentos pode ser medida através de métricas, sendo o índice de Davies-Bouldin [95] frequentemente utilizado para esse objetivo, o qual independe do número de agrupamentos e do método de partição dos dados. Esse índice visa avaliar não só o grau de espalhamento dentro de um agrupamento quanto o grau de separação entre eles. O grau de separação entre agrupamentos é avaliado pelo somatório das distâncias entre os centroides. Já o espalhamento de um agrupamento é dado como uma média da distância entre o centroide e os dados correspondentes. O valor do índice é inversamente proporcional à separação entre agrupamentos e diretamente ao espalhamento dentro de um agrupamento. Portanto, quanto menor o valor do índice, melhor o agrupamento dos dados.

Para o agrupamento dos dados, o algoritmo *k-means* foi executado sobre o conjunto formado pelos vetores de peso dos neurônios do SOM, os quais possuem a dimensionalidade dos dados utilizados para o treinamento. O *k-means* foi executado para agrupamentos que variaram de 1 a 30 clusters, sendo repetido até 20 vezes com inicializações diferentes para cada agrupamento, onde a melhor é escolhida pela menor soma de erro quadrático alcançada. Para cada agrupamento, o índice de Davis-Bouldin é calculado. A variação do valor do índice entre os diferentes números de clusters é analisada, sendo o menor valor alcançado uma indicação para o número ótimo de clusters.

No caso da análise de grupos de risco, o mapa é agrupado em três clusters, independentemente do índice de Davis-Bouldin. Cada um dos clusters é rotulado como

de baixo, médio ou alto risco, de acordo com o diagnóstico dos pacientes localizados no agrupamento. Como exemplo, quanto mais pacientes com TB resistente um agrupamento possuir, maior o risco desse grupo em relação a um conjunto que apresente mais TB sensível ou diagnósticos negativos.

Para a mineração de dados e análise das variáveis marcadoras com relação ao caso observado, um rótulo deve ser atribuído a cada agrupamento formado, o qual indicará que grupo de interesse está sendo agrupado naquela determinada região do mapa. Para isso, utiliza-se o critério de frequência, que avalia qual grupo de pacientes, os quais os agrupamentos corretos já são previamente conhecidos, foi mais associado ao agrupamento indicado pelo SOM quando os mesmos são projetados sobre o mapa. Portanto, os dados mapeados em um agrupamento são rotulados como a classe que é mais frequente nessa região do mapa. Quanto maior a diferença entre a frequência das classes na região, maior a tendência do agrupamento estar rotulado corretamente. Ao final, as regiões do mapa se encontram subdivididas entre as classes de interesse, e índices de classificação como sensibilidade, especificidade, acurácia, podem ser calculados.

Ainda, uma das principais características do SOM é a relação topológica entre o mapa treinado e os dados utilizados. Com isso, é possível utilizar esse fator para entender como cada variável influencia no mapa gerado pelos Planos de Componentes. Esses planos são constituídos pelos valores médios de cada variável utilizada na classificação, projetados na treliça de neurônios do mapa [96], o que fornece uma visualização da distribuição espacial no mapa de uma determinada variável.

Ao analisar o plano dos componentes, em relação à clusterização obtida no mapa, é possível avaliar a relação topológica entre as variáveis e os agrupamentos formados. Isso auxilia na análise da influência de cada variável nos grupos, e com isso, ter uma indicação de qual fator pode ser mais influente dentro do agrupamento de interesse.

Capítulo 5

Resultados

Neste capítulo, são apresentados os resultados referentes a aplicação da metodologia descrita no capítulo anterior. Diferentes análises são realizadas sobre os diagnósticos dos pacientes, onde classificadores utilizando a técnica de redes neurais multicamadas são projetados, assim como o estudo da relevância de variáveis e grupos de risco são feitos através da técnica de mapas auto-organizáveis.

5.1 Estudo entre diagnósticos de TB positivo e TB negativo

Para pacientes com o diagnóstico de TB positivo e TB negativo, um discriminador neural para realizar a classificação entre esses dois diagnósticos foi proposto. A Figura 5.1 mostra o índice SP médio alcançado pelas redes, por topologia, ao se aplicar na rede os conjuntos de teste e treino em separados, assim como sua incerteza (representada pelo valor RMS das 50 seleções da validação cruzada) e o valor máximo obtido. A topologia ótima é escolhida pela melhor média do índice SP para o conjunto de teste, nesse caso, representada pela estrutura com 4 neurônios na camada escondida. Já a Tabela 5.1 descreve os valores obtidos referentes ao desempenho do modelo neural, incluindo sua incerteza.

Para um cenário de operação, redes com maiores índices de sensibilidade e especi-

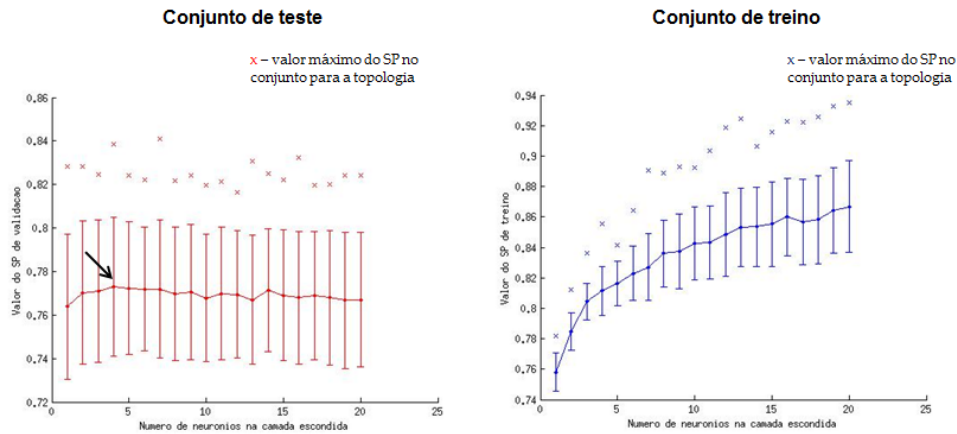


Figura 5.1: Valores de SP médios para cada topologia, pelos conjuntos de dados de teste e treino, considerando as redes obtidas pela validação cruzada. As barras de erro são representadas pelos valores RMS dos índices SP. A topologia com 4 neurônios é a que apresenta melhor média para o conjunto de teste.

Desempenho do modelo neural	
Índice SP	$0,7730 \pm 0,0318$
Sensibilidade	$0,8205 \pm 0,0676$
Especificidade	$0,7286 \pm 0,0548$
Acurácia	$0,7780 \pm 0,0327$

Tabela 5.1: Desempenho alcançado pelo classificador para pacientes com diagnósticos TB positivo e TB negativo, considerando a topologia com 4 neurônios na camada escondida.

ficidade são selecionadas. A Figura 5.2 mostra a curva ROC obtida para as redes de melhor SP e sensibilidade. Também é possível observar a evolução do treinamento para a rede de melhor SP, o qual teve um ponto de parada em 401 épocas com um MSE igual a 0,909. A distribuição dos valores de saída de classificação para a rede com maior sensibilidade apresenta um comportamento que atinge o alvo -1 para o diagnóstico de TB negativo, porém, para TB positivo, o pico é localizado próximo ao valor de 0,6. Isso indica uma dificuldade na caracterização dessa classe.

Após selecionar a rede de operação com o melhor índice SP, foi analisado como cada variável contribuiu para o resultado obtido. A Figura 5.3 mostra os valores da diferença entre o SP da rede de operação e o SP alcançado pelo conjunto gerado para a análise das variáveis. Nesse caso, quanto menor o valor obtido, menor sua

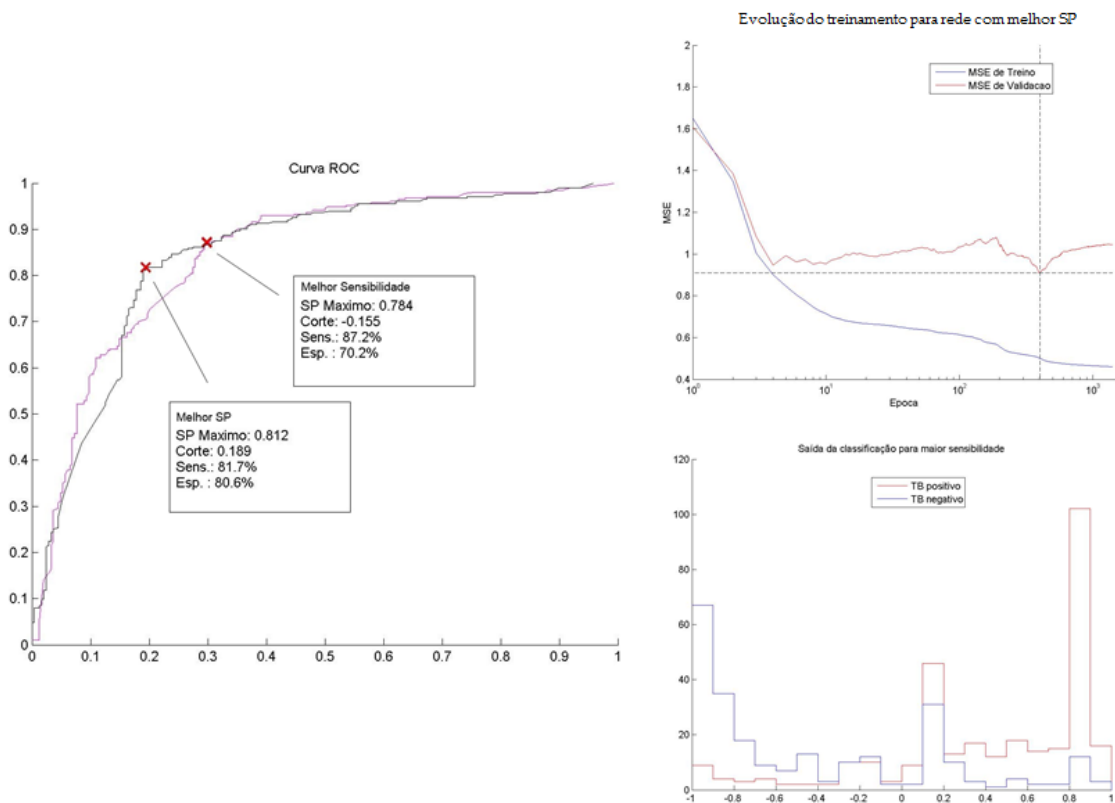


Figura 5.2: Representações da curva ROC para a rede de melhor SP e sensibilidade, evolução do treino para a rede de melhor SP e distribuição dos valores de saída da rede neural com maior sensibilidade para classificar pacientes entre os diagnósticos de TB positivo e TB negativo.

relevância para o resultado. Pode-se observar que os sintomas catarro com sangue e falta de ar apresentaram as menores contribuições para o classificador, o que pode indicar que esses fatores não possuem um grande impacto na discriminação entre os diagnósticos de TB positivo e TB negativo.

A Tabela 5.2 exibe os valores de SP, sensibilidade, especificidade e acurácia para os pacientes com triagem realizadas nos estados do Rio de Janeiro, Ceará, São Paulo e Rio Grande do Sul, utilizando, para isso, a rede que obteve melhor índice SP. É possível observar que a sensibilidade para os pacientes no Rio Grande do Sul apresentou valores bem abaixo do conjunto completo, apesar de um leve aumento em sua especificidade. Já em São Paulo a situação se inverte, apresentando uma especificidade baixa, enquanto teve um aumento de sensibilidade. Os outros valores, apesar de pequena variação, não apresentam discrepâncias tão grandes em relação aos valores alcançados pelo conjunto de todos os pacientes.

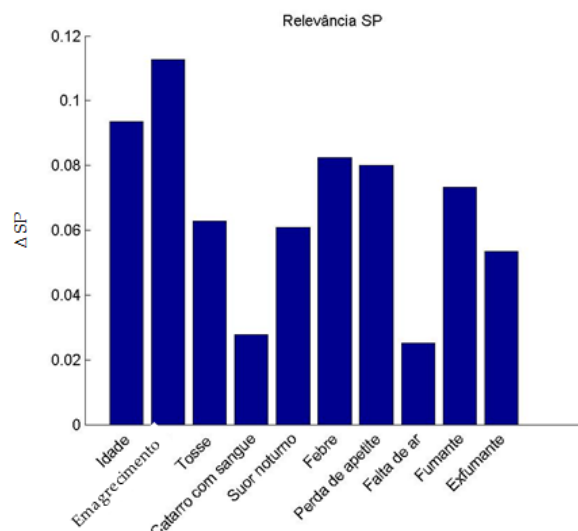


Figura 5.3: Estudo de relevância das variáveis na rede neural multicamadas obtida para os diagnósticos TB positivo e TB negativo. O índice SP da rede de operação é subtraído do índice SP alcançado pela rede quando se fixa cada variável com o seu valor médio.

Local de triagem	Índice SP	Sensibilidade (%)	Especificidade (%)	Acurácia (%)
Rio de Janeiro	0.7839 ± 0.0270	0.8252 ± 0.0445	0.7448 ± 0.0541	0.7773 ± 0.0279
Ceará	0.7629 ± 0.0333	0.8070 ± 0.0809	0.7232 ± 0.0747	0.7825 ± 0.0463
São Paulo	0.6569 ± 0.0594	0.7086 ± 0.1083	0.6164 ± 0.1325	0.6614 ± 0.0564
Rio Grande do Sul	0.7724 ± 0.0854	0.7100 ± 0.1620	0.8560 ± 0.1897	0.6093 ± 0.0876
Todos	0,7730 ± 0,0318	0,8205 ± 0,0676	0,7286 ± 0,0548	0,7780 ± 0,0327

Tabela 5.2: Desempenho alcançado pelo classificador, por localização da triagem, para pacientes com diagnósticos TB positivo e TB negativo.

Para o estudo de grupos de risco, é aplicada a técnica de mapas auto-organizáveis, a qual permite visualizar, em duas dimensões, os agrupamentos formados por pacientes. A Figura 5.4 mostra a U-Matrix obtida pelo treinamento do mapa SOM, o qual apresentou a dimensão de 13 por 9 neurônios. Com isso, o algoritmo k-means foi utilizado para identificar três agrupamentos, sendo repetido 20 vezes com inici-

alizações diferentes, onde a melhor configuração foi escolhida pela menor soma do erro quadrático.

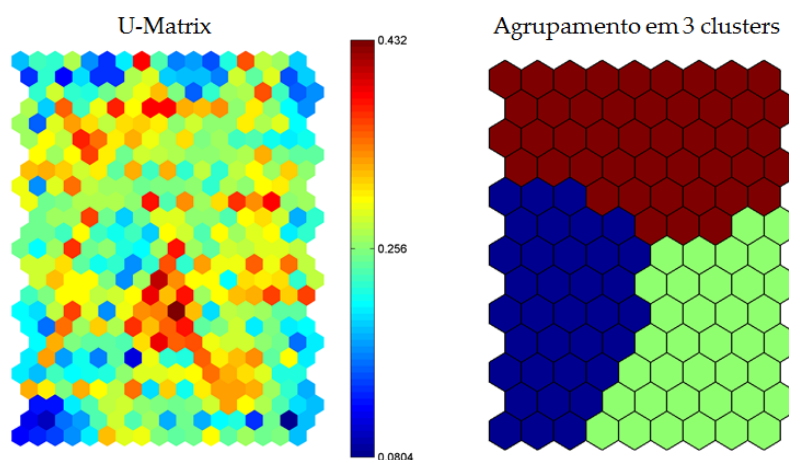


Figura 5.4: U-Matrix do mapa SOM e o resultado da sua clusterização pela técnica de k-means em 3 agrupamentos para pacientes com diagnósticos TB positivo e TB negativo.

Com os três agrupamentos definidos, foram rotulados grupos de baixo, médio e alto risco, os quais visam mapear os riscos de um paciente possuir qualquer um dos tipos de TB considerados, avaliando, para isso, os pacientes (e seus diagnósticos) que foram mapeados dentro de cada região. O grupo de alto risco foi rotulado dessa forma pela quantidade de pacientes com TB positivo presente nesse setor do mapa, possuindo um total de 48,4% contra 26% e 25,6% das outras regiões. Já o grupo de baixo risco foi avaliado pela quantidade de pacientes com diagnóstico negativo para TB, sendo 42% do total contra 33% do grupo considerado de médio risco. É possível observar que o grupo de alto risco apresenta uma grande quantidade de pacientes com TB-DR e TB-MDR. Ainda, caso os grupos de médio e alto risco sejam considerados como um só, uma alta detecção de casos com resistência é alcançada. A Figura 5.5 exhibe os grupos de risco identificados no mapa, assim como a distribuição de pacientes em cada agrupamento.

Com os grupos de risco definidos, é possível avaliar a contribuição de cada variável através da análise do mapa de componentes, o qual pode ser visto na Figura 5.6. Por meio de inspeção visual dos mapas, algumas características podem ser observadas.

Quando vista a distribuição da idade no mapa, vemos uma concentração de pes-

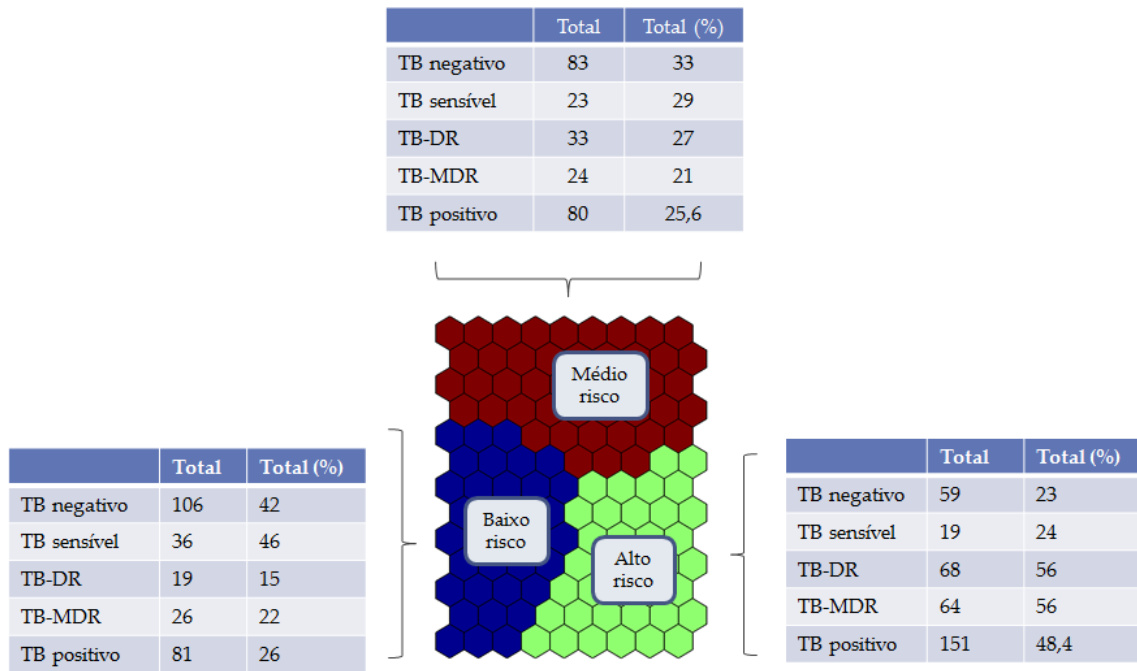


Figura 5.5: Definição dos agrupamentos em grupos de baixo, médio e alto risco, de acordo com a distribuição de pacientes.

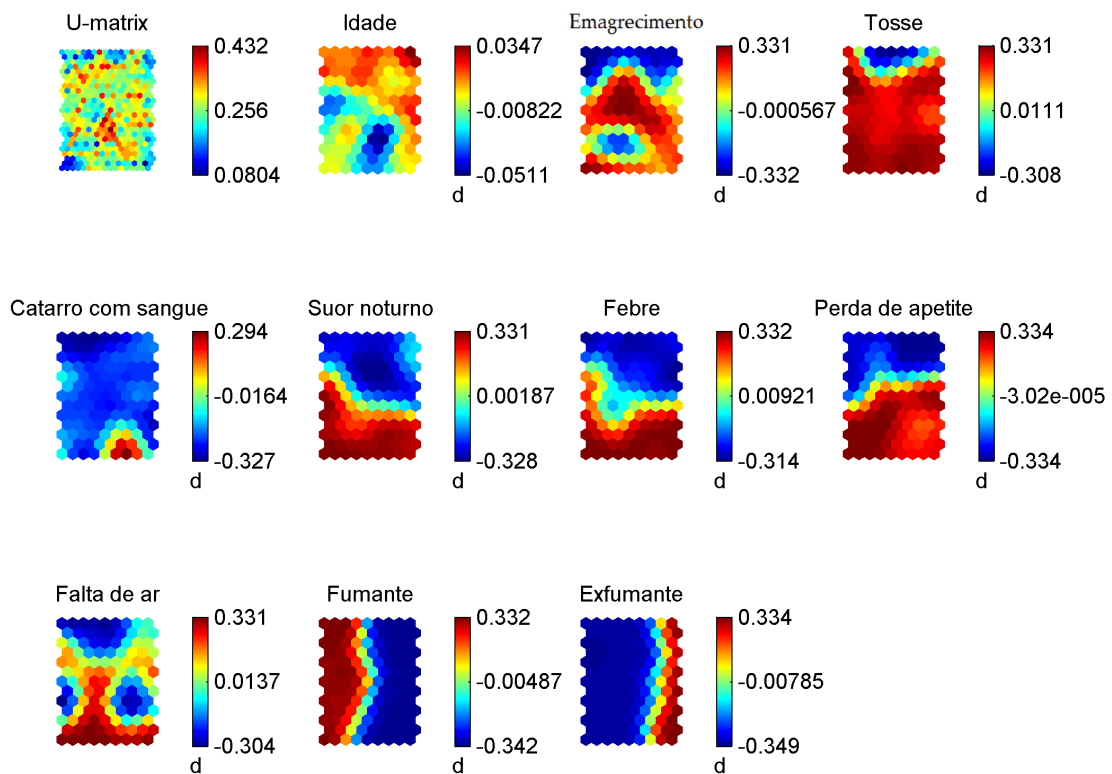


Figura 5.6: Mapa de componentes referente ao mapa SOM para pacientes com diagnósticos TB positivo e TB negativo.

soas de idade mais elevada na parte superior do mapa, o que corresponde a área de médio risco, enquanto as idades mais baixas apresentam um forte agrupamento

na região de alto risco. O plano que representa o emagrecimento mostra uma distribuição da presença ou não do sintoma em diferentes regiões do mapa, com uma maior concentração de pessoas sem o sintoma na região de médio risco. Apesar disso, poucas conclusões podem ser tomadas em relação a essa informação.

Já a variável tosse aparece concentrada em uma região muito específica do mapa, dentro da região de médio risco, indicando que pacientes que não apresentam esse sintoma são casos bem distintos em relação ao resto da população estudada.

Ao compararmos os sintomas suor noturno e febre, é possível observar que a distribuição de ambos são bem similares, com uma grande concentração na parte superior do mapa, sendo esta característica um forte indício de que esses sintomas podem estar estatisticamente correlacionadas.

Casos que apresentam perda de apetite aparecem localizados na parte inferior do mapa, onde se localizam as regiões de baixo e alto risco, o que não traz informações muito relevantes para essa análise. Pessoas com catarro com sangue são mais concentradas na região de alto risco, o que pode levar a consideração de que esse fator é possui uma associação relevante com riscos maiores. A presença ou não de Falta de ar, de acordo com a distribuição do mapa, não é conclusiva, pois é possível observar concentrações em partes distintas do mapa e em diferentes grupos de risco. Esse fator pode ser um indício de concordância com a análise de relevância realizada pelo índice SP através da rede neural multicamada, que também indicou que esse sintoma poderia ser pouco relevante nesse caso. Já a distribuição de fumantes e ex-fumantes apresentam distribuições similares, o que pode indicar um correlacionamento entre esses aspectos, e que também condiz com a realidade.

Após a análise com grupos de risco, o mapa foi clusterizado com um número livre de agrupamentos, o qual foi indicado pelo índice de Davis-Bouldin. A Figura 5.7 mostra o valor do índice DB para mapas com 1 até 30 clusters, onde é possível observar que não há variação significativa no índice entre o intervalo de 10 a 23 grupos. Para evitar uma granularidade muito fina, o que pode prejudicar a identificação de agrupamentos, 10 clusters foram utilizados na análise.

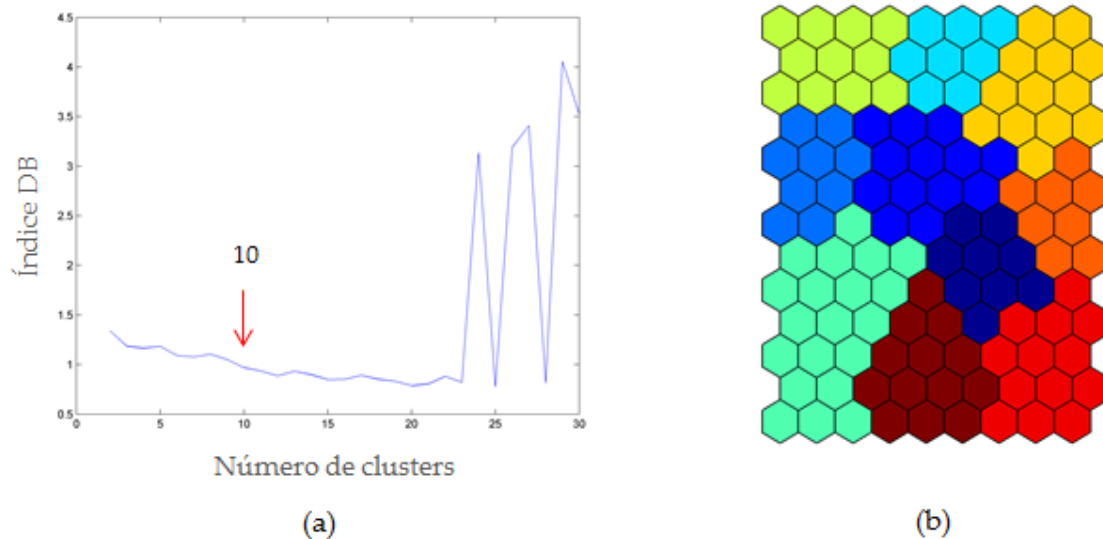


Figura 5.7: Clusterização do mapa SOM gerado para os pacientes com diagnósticos de TB positivo e TB negativo. (a) Índice de Davis-Bouldin por número de clusters (b) Mapa SOM clusterizado com 10 agrupamentos

Em seguida à clusterização, os agrupamentos são rotulados de acordo com o diagnóstico mais frequente na região. A Figura 5.8 exibe um mapa de frequência, representado por um código de cores que varia de 100% de quantidade de TB negativo em um cluster, até 100% de pacientes com TB positivo, além do mapa já rotulado entre os dois diagnósticos alvos.

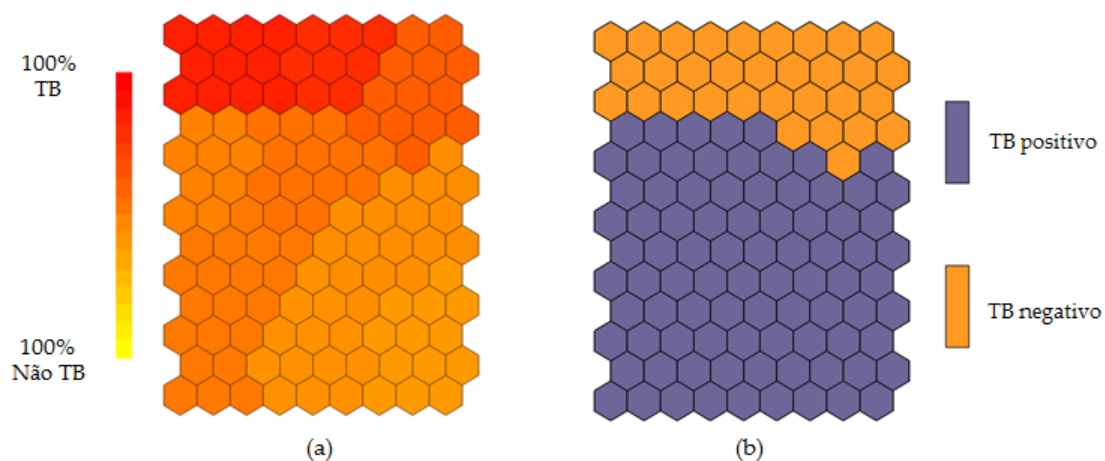


Figura 5.8: Rotulação do mapa SOM entre os diagnósticos de TB positivo e TB negativo. (a) Frequência de diagnósticos em cada cluster (b) Rotulação dos clusters por diagnóstico

Por essa rotulação, um índice SP de 0,656 foi obtido, com sensibilidade de 84%, especificidade de 48,8% e acurácia com valor de 68,9%. Ao comparar visualmente o

mapa rotulado com o plano de componentes, é possível observar que os sintomas de suor noturno, febre, perda de apetite e tosse apresentam também uma distribuição horizontal, compatível com a que foi encontrada pela rotulação, o que pode indicar uma contribuição para a definição das classes. Além disso, para esses sintomas, pacientes que não os possuem acabam sendo localizados na parte superior do mapa, a qual boa parte foi rotulada como TB negativa, indo ao encontro da prática médica. Já a falta de ar não apresenta distribuição compatível com nenhum dos rótulos, o que indica pouca relevância para a classificação, e entra em concordância com as análises anteriores.

5.2 Estudo entre diagnósticos de TB resistente e TB sensível

Um discriminador neural foi também proposto para realizar a classificação entre pacientes com diagnósticos de TB resistente e TB sensível. A Figura 5.9 exibe os valores médios do índice SP, por topologia, para os conjuntos de treino e teste, indicando sua incerteza (valor RMS) e valor máximo. A topologia de 8 neurônios é selecionada como a melhor, utilizando como base para a escolha, a média do SP para os dados de teste. Os índices de desempenho esperado para este modelo neural são descritos na Tabela 5.3, avaliando-se para isso, as 50 redes com a topologia escolhida.

Desempenho do modelo neural	
Índice SP	$0,7627 \pm 0,0402$
Sensibilidade	$0,7308 \pm 0,0813$
Especificidade	$0,7986 \pm 0,0791$
Acurácia	$0,7478 \pm 0,0533$

Tabela 5.3: Desempenho alcançado pelo classificador para pacientes com diagnósticos TB resistente e TB sensível, considerando a topologia com 8 neurônios na camada escondida.

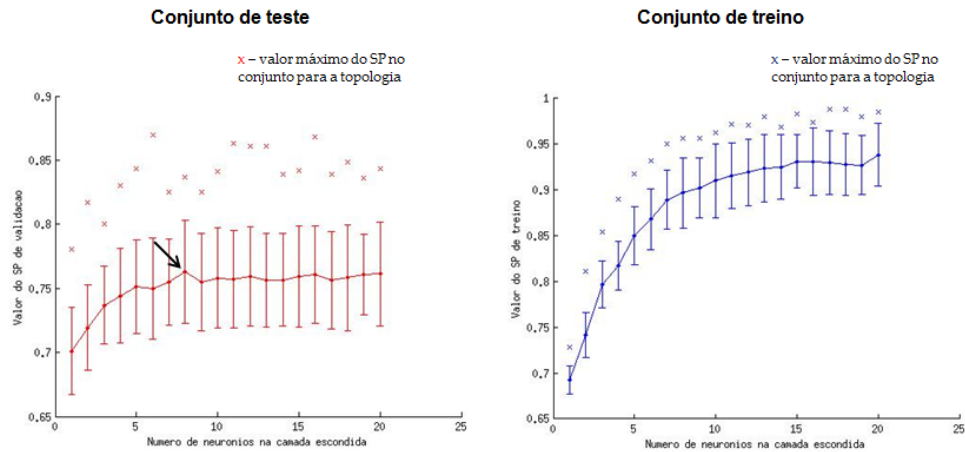


Figura 5.9: Valores de SP médios para cada topologia, pelos conjuntos de dados de teste e treino, considerando as redes obtidas pela validação cruzada. As barras de erro são representadas pelo valor RMS dos índices SP. A topologia com 8 neurônios é a que apresenta melhor média para o conjunto de teste.

Considerando um cenário de operação, a Figura 5.10 exhibe a curva ROC obtida, assim como é ilustrada a evolução do treinamento para a rede que apresentou o melhor SP, que teve um ponto de parada em 565 épocas com valor de MSE igual a 0,8842. Os valores de saída para a rede com melhor sensibilidade se encontram bem concentrados nos extremos dos valores (-1 e 1), o que sugere uma boa classificação entre os diagnósticos, pois os alvos são atingidos.

A partir da rede que obteve o maior índice SP, foram analisadas as contribuições de cada variável no classificador. A Figura 5.11 mostra a relevância de cada variável. Ao observar os valores, a variável que indica ter uma menor relevância para a classificação entre TB resistente e TB sensível é a tosse. As outras variáveis, apesar de algumas variações, apresentam maiores diferenças, com destaque para uma maior contribuição da falta de ar, a qual havia apresentado a menor contribuição na análise anterior entre TB positivo e TB negativo.

A Tabela 5.4 indica os valores de SP, sensibilidade, especificidade e acurácia para os estados avaliados, utilizando para isso, a rede que apresentou melhor índice SP. Não é possível observar grandes discrepâncias em relação dos índices em relação ao conjunto completo de dados, com exceção do Rio Grande do Sul, que alcança uma especificidade de 100%, além de um pequeno aumento de sensibilidade, o que eleva

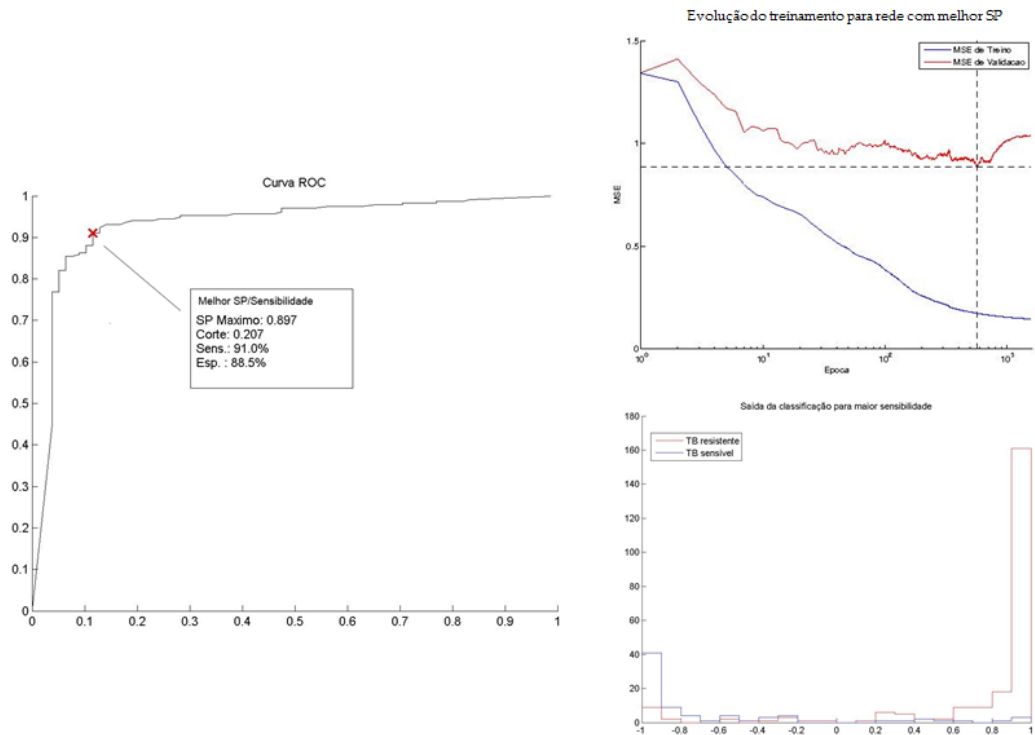


Figura 5.10: Representações da curva ROC para a rede de melhor SP e sensibilidade, evolução do treino para a rede de melhor SP e distribuição dos valores de classificação para a rede neural com maior sensibilidade para classificar pacientes entre os diagnósticos de TB resistente e TB sensível.

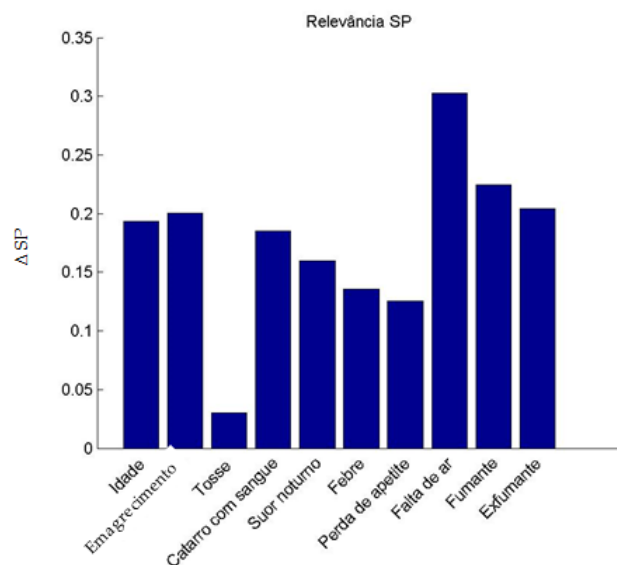


Figura 5.11: Estudo de relevância das variáveis na rede neural multicamadas obtida para os diagnósticos TB resistente e TB sensível. O índice SP da rede de operação é subtraído do índice SP alcançado pela rede quando se define cada variável com o seu valor médio.

também o seu valor de SP.

Local de triagem	Índice SP	Sensibilidade (%)	Especificidade (%)	Acurácia
Rio de Janeiro	0.7543 ± 0.1106	0.7377 ± 0.1228	0.7771 ± 0.1505	0.7479 ± 0.1066
Ceará	0.8314 ± 0.1127	0.8265 ± 0.0829	0.8438 ± 0.1861	0.8304 ± 0.0807
São Paulo	0.7259 ± 0.1050	0.6884 ± 0.1501	0.7730 ± 0.1415	0.7152 ± 0.1109
Rio Grande do Sul	0.8429 ± 0.1150	0.7887 ± 0.1868	0.9133 ± 0.1623	0.6327 ± 0.1270
Todos	0,7627 ± 0,0402	0,7308 ± 0,0813	0,7986 ± 0,0791	0,7478 ± 0,0533

Tabela 5.4: Desempenho alcançado pelo classificador, por localização da triagem, para pacientes com diagnósticos TB resistente e TB sensível.

Já no estudo de grupos de risco, a técnica de mapas auto-organizáveis foi utilizada, sendo responsável por representar de forma bidimensional os agrupamentos formados por pacientes. A Figura 5.12 mostra a U-Matrix gerada através do treinamento do mapa SOM, apresentando dimensões de 11 por 8 neurônios. A partir desse mapa, três agrupamentos são indentificados através da aplicação do algoritmo k-means, técnica que é repetida 20 vezes, sendo a melhor configuração selecionada pelo menor erro quadrático.

Os grupos de risco são definidos entre baixo, médio e alto risco em relação a TB resistente, utilizando o mapa e agrupamentos gerados anteriormente. Essas definições foram tomadas de acordo com o mapeamento dos pacientes com seus respectivos diagnósticos, sendo mostradas na Figura 5.13. O grupo de alto risco foi definido pelo número de pacientes localizados nessa região que apresentam TB em sua forma resistente, representando um valor de 42,7% do total. O grupo de baixo risco também foi avaliado pelo número de pacientes com TB resistente, o qual tem um total de 23,5%, contra 33,7% do grupo de risco médio. Ao avaliar os conjuntos de alto e médio risco, é possível ver uma grande quantidade de casos

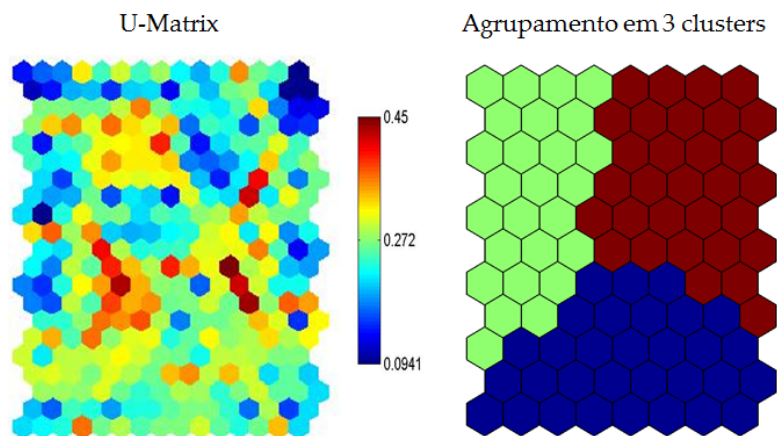


Figura 5.12: U-Matrix do mapa SOM e o resultado da sua clusterização pela técnica de k-means em 3 agrupamentos para dados de pacientes com TB resistente e TB sensível.

com resistência, entretanto, com baixa especificidade pelo número de casos sensíveis também presentes.

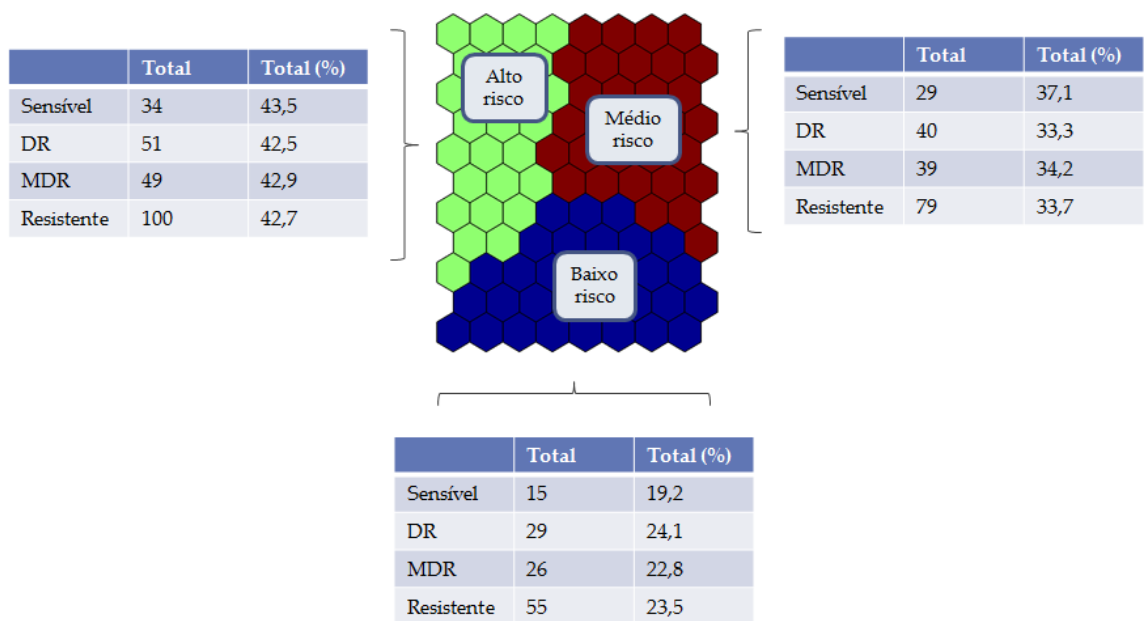


Figura 5.13: Definição dos agrupamentos em grupos de baixo, médio e alto risco, de acordo com a distribuição de pacientes.

Ao serem definidos grupos de risco, é possível avaliar as variáveis e suas contribuições, a partir de suas distribuições, pela análise visual do mapa de componentes, o qual pode ser visto na Figura 5.14.

Ao observar a distribuição da variável idade, não é possível identificar uma relação relevante aos grupos de risco, pois seu mapa apresenta apenas uma concentração

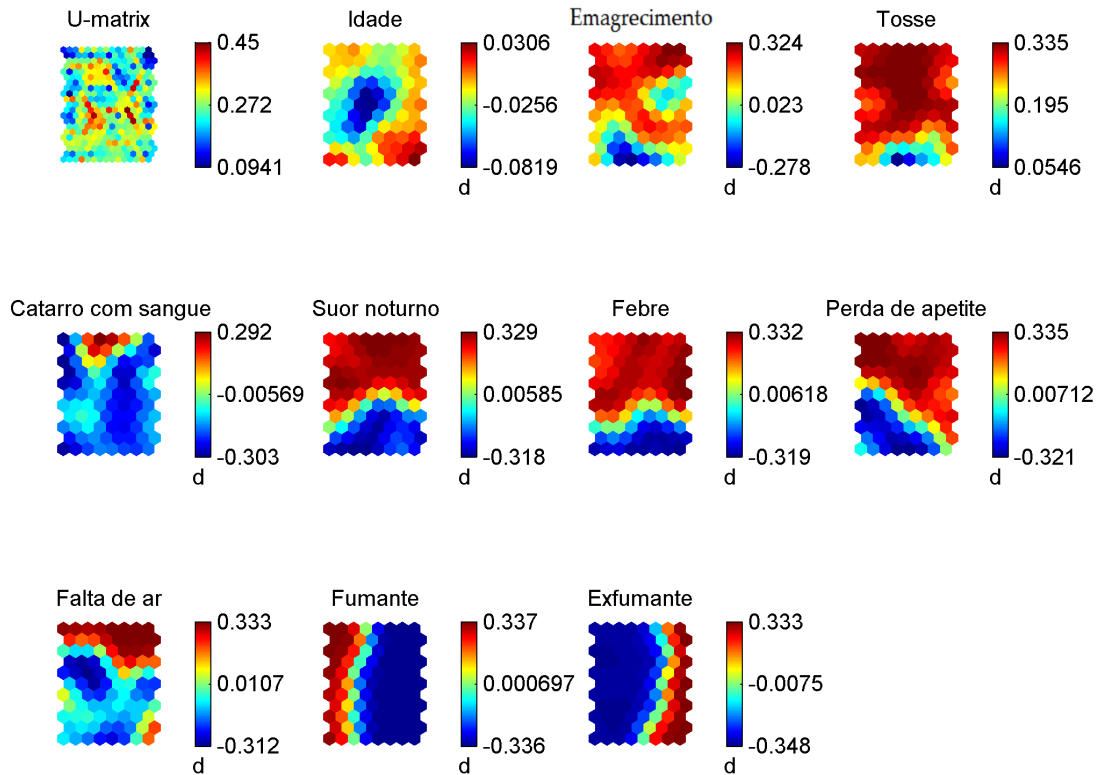


Figura 5.14: Mapa de componentes referente ao mapa SOM para pacientes com diagnósticos TB resistente e TB sensível.

mais visível na região central dos pacientes com idade mais baixa. Entretanto, um agrupamento de pacientes com idades mais elevadas é visto localizado na região inferior direita, o que poderia ser compatível com o grupo de baixo risco.

Casos que não apresentam o sintoma de emagrecimento se localizam em regiões correspondentes aos grupos de baixo e médio risco, podendo indicar um relacionamento maior da ausência desse fator com a TB sensível. Os pacientes sem tosse aparecem concentrados em uma região muito específica do mapa, indicando não só a distinção desse grupo em relação aos outros pacientes, já que a incidência desse sintoma é alta em ambos os diagnósticos, mas uma possível relação com o grupo de baixo risco. O espalhamento dos pacientes com tosse no mapa é um indício que contribui com a análise anterior de relevância por SP, na qual esse sintoma se destacou como o que menos contribui para uma correta classificação entre diagnósticos.

Já os sintomas de febre e suor noturno, como no caso anterior, quando comparados, apresentam distribuições similares, com ambas possuindo concentrações na

parte inferior do mapa, o que pode fortemente significar a existência de uma correlação estatística entre elas. A presença de Catarro com sangue se localiza concentrado na parte superior do mapa, região onde estão localizados os grupos de médio e alto risco. A Perda de apetite apresenta uma distribuição que localiza os pacientes com sintoma no canto superior direito, enquanto os que não possuem no canto inferior esquerdo. Apesar disso, não é visto informações relevantes em relação aos grupos de risco.

A falta de ar se apresenta, em sua maioria, na parte superior do mapa, porém sua ausência se apresenta espalhada em boa parte das regiões, então não é possível tirar conclusões dessa análise. Os mapas de Fumante e Ex-fumante apresentam distribuições correlacionadas, o que é o esperado, mas sem relacionamento com os grupos de risco.

Em sequência a análise de grupos de risco, uma clusterização do mapa com um número de agrupamentos guiado pelo índice de Davis-Bouldin foi realizado. A partir dos valores desse índice para diversas variações de número de agrupamentos, como visto na Figura 5.15, um número de 15 clusters foi escolhido, sendo esse o que apresentou menor valor de Davis-Bouldin.

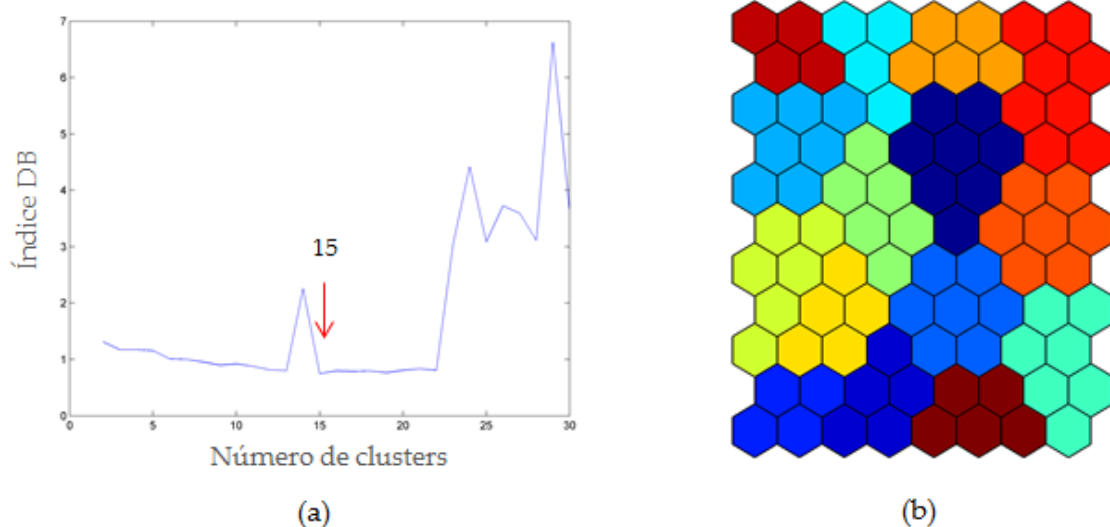


Figura 5.15: Clusterização do mapa SOM gerado para os pacientes com diagnósticos de TB resistente e TB sensível. (a) Índice de Davis-Bouldin por número de clusters (b) Mapa SOM clusterizado com 15 agrupamentos

Com a clusterização em 15 agrupamentos, rótulos são atribuídos de acordo com a incidência dos diagnósticos em cada cluster. Seguindo esse procedimento, não foi possível identificar um cluster com maioria de pacientes com TB sensível, como mostra a Figura 5.16. Esse fato pode ser explicado pelo fato do conjunto de pacientes utilizados no estudo serem selecionados por possuírem suspeita de TB resistente, o que pode contribuir para a baixa estatística em relação aos casos de TB sensível, e com isso, dificultar a formação de agrupamentos para esses casos.

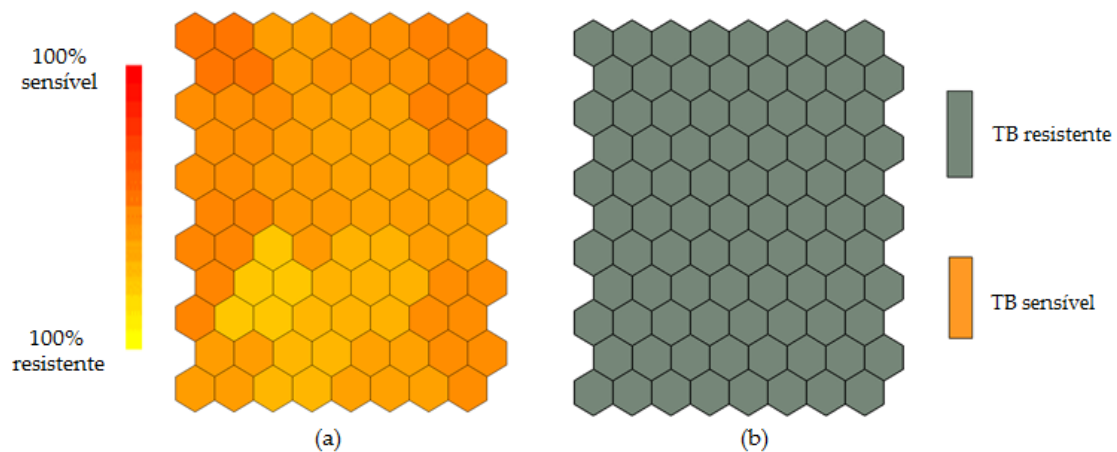


Figura 5.16: Rotulação do mapa SOM entre os diagnósticos de TB resistente e TB sensível. (a) Frequência de diagnósticos em cada cluster (b) Rotulação dos clusters por diagnóstico

5.3 Estudo entre diagnósticos de TB-DR e TB-MDR

Após a realização dos modelos neurais anteriores, foi proposto um discriminador entre pacientes com o diagnóstico de TB-DR e TB-MDR. A Figura 5.17 exibe os valores do índice SP por topologia, obtidos pelos conjuntos de teste e treino, assim como sua incerteza (valor RMS) e valor máximo. A melhor topologia é obtida pela maior média de SP alcançada pelo conjunto de teste, sendo nesse caso a configuração com 11 neurônios na camada escondida. O modelo neural teve seu desempenho avaliado pelo índice SP, sensibilidade, especificidade e acurácia, e os resultados são descritos na Tabela 5.5.

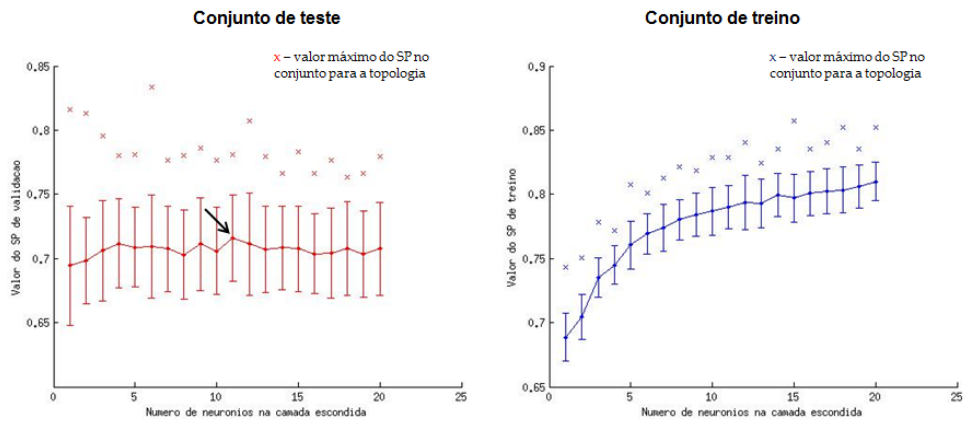


Figura 5.17: Valores de SP médios para cada topologia, pelos conjuntos de dados de teste e treino, considerando as redes obtidas pela validação cruzada. As barras de erro são representadas pelo valor RMS dos índices SP. A topologia com 11 neurônios é a que apresenta melhor média para o conjunto de teste.

Desempenho do modelo neural	
Índice SP	$0,7159 \pm 0,0336$
Sensibilidade	$0,7067 \pm 0,0967$
Especificidade	$0,7308 \pm 0,0948$
Acurácia	$0,7183 \pm 0,0336$

Tabela 5.5: Desempenho alcançado pelo classificador para pacientes com diagnósticos TB-DR e TB-MDR, considerando a topologia com 11 neurônios na camada escondida.

A rede selecionada para um ambiente de operação são selecionadas pelos melhores índices SP, sensibilidade e especificidade. A Figura 5.18 exhibe as curvas ROC para as redes de melhor SP e sensibilidade, assim como a evolução do treinamento da rede com melhor SP, que teve sua parada com 48 épocas e com MSE de 1,0214. A distribuição dos valores de saída da rede com melhor sensibilidade também pode ser vista, a qual apresenta a dificuldade na identificação das duas classes, apresentando picos nos valores de -0,3 e 0,5.

Um estudo sobre a contribuição de cada variável, através da rede neural com melhor SP, também foi realizado. A Figura 5.19 mostra os valores da diferença entre o SP alcançado pela propagação dos dados originais e os dados com as variáveis definidas para a análise. Como anteriormente, apenas a variável tosse apresenta

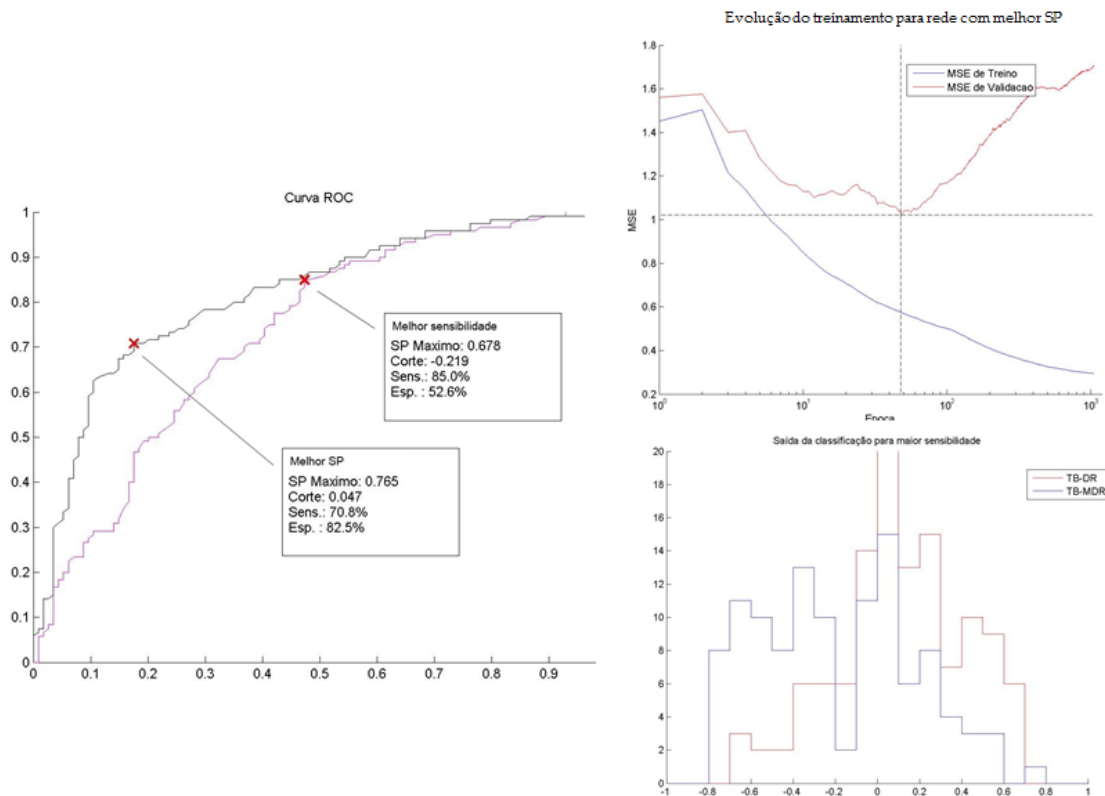


Figura 5.18: Representações da curva ROC para a rede de melhor SP e sensibilidade, evolução do treinamento para a rede de melhor SP e distribuição dos valores de classificação para a rede neural com maior sensibilidade para classificar pacientes entre os diagnósticos de TB-DR e TB-MDR.

indícios de que pode não ser relevante para a classificação entre os diagnóstico de TB-DR e TB-MDR, já que as outras, apesar de pequenas variações, se mantém bastante similares entre si.

Também com a rede que teve o melhor SP, um estudo sobre as diferenças entre os índices de diferentes regiões de triagem pode ser visto na Tabela 5.6, que apresenta valores de SP, sensibilidade, especificidade e acurácia para cada estado. O Ceará apresenta uma queda dos valores de sensibilidade e especificidade, o que leva também a uma queda no índice SP. Já o Rio Grande do Sul apresenta valores maiores para SP, sensibilidade, e alcançando 100% de especificidade. São Paulo possui valores mais elevados de sensibilidade, porém com queda de especificidade, situação inversa a do Rio de Janeiro.

A técnica de mapas auto-organizáveis foi utilizada para fazer o estudo de grupos de risco para TB-MDR, a qual possibilita a representação em duas dimensões dos

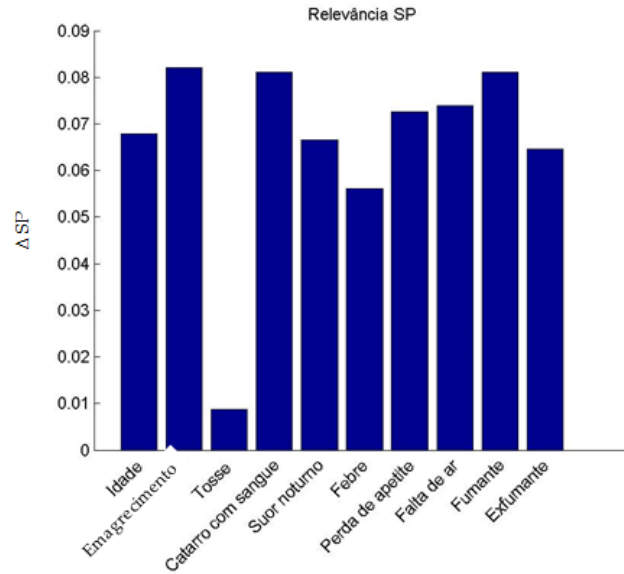


Figura 5.19: Estudo de relevância das variáveis na rede neural multicamadas obtida para os diagnósticos TB resistente e TB sensível. O índice SP da rede de operação é subtraído do índice SP alcançado pela rede quando se define cada variável com o seu valor médio.

Local de triagem	Índice SP	Sensibilidade (%)	Especificidade (%)	Acurácia (%)
Rio de Janeiro	0.6399 ± 0.0488	0.6581 ± 0.1271	0.6330 ± 0.1243	0.6441 ± 0.0474
Ceará	0.6257 ± 0.0473	0.6311 ± 0.1156	0.6309 ± 0.1282	0.6310 ± 0.0446
São Paulo	0.6569 ± 0.0594	0.7086 ± 0.1083	0.6164 ± 0.1325	0.6614 ± 0.0564
Rio Grande do Sul	0.7724 ± 0.0854	0.7100 ± 0.1620	0.8560 ± 0.1897	0.6093 ± 0.0876
Todos	0,7159 ± 0,0336	0,7067 ± 0,0967	0,7308 ± 0,0948	0,7183 ± 0,0336

Tabela 5.6: Desempenho alcançado pelo classificador, por localização da triagem, para pacientes com diagnósticos TB-DR e TB-MDR.

agrupamentos formados pelos pacientes. A U-Matrix gerada pode ser vista na Figura 5.20, que representa um mapa com tamanho de 10 por 8 neurônios, sendo obtida através de treinamento do SOM. Os três agrupamentos, que serão as bases para os grupos de risco, são definidos pela aplicação do k-means, o qual é inicializado 20

vezes, e o melhor agrupamento escolhido pelo menor erro quadrático.

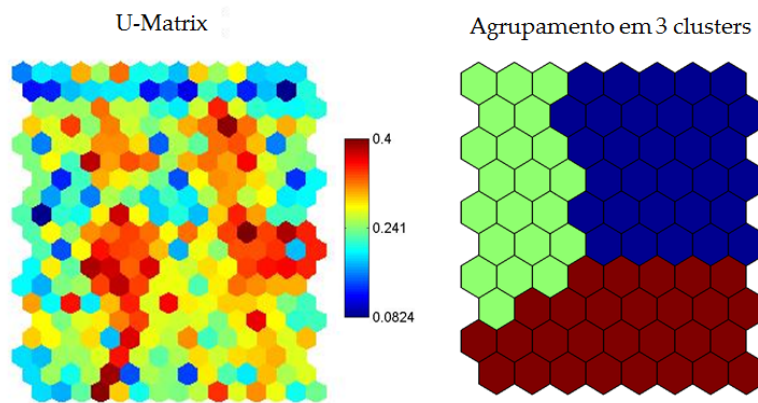


Figura 5.20: U-Matrix do mapa SOM e o resultado da sua clusterização pela técnica de k-means em 3 agrupamentos.

Os grupos de baixo, médio e alto risco devem ser então rotulados a partir do agrupamento anterior, analisando a distribuição dos pacientes no mapa. Apesar disso, como visto na Figura 5.21, a distribuição entre pacientes diagnosticados com TB-DR e TB-MDR não é significativa nos clusters para fazer qualquer afirmação nesse âmbito, não sendo possível, portanto, identificar grupos de risco.

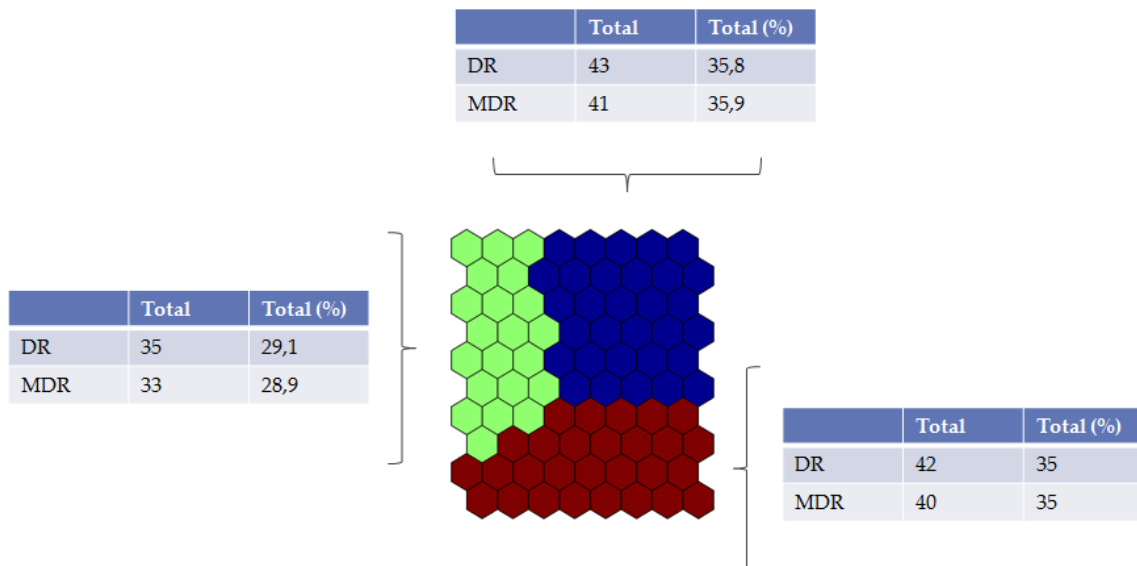


Figura 5.21: Definição dos agrupamentos em grupos de baixo, médio e alto risco, de acordo com a distribuição de pacientes.

Na falta de grupos de risco, ao analisar o mapa de componentes, visto na Figura 5.22, é possível observar que os sintomas suor noturno e febre apresentam distribuições semelhantes, apresentando características que indicam uma correlação

estatística entre eles. Do maneira semelhante, as variáveis fumante e ex-fumante também apresentam similaridades, e com isso, um possível correlacionamento. Pacientes com a ausência de tosse e emagrecimento se encontram localizados em regiões muito específicas do mapa, podendo representar pacientes que são distintos da maioria do conjunto. Já catarro com sangue, perda de apetite, falta de ar e idade não apresentam características relevantes o suficiente para se fazer alguma afirmação nesse tipo de análise.

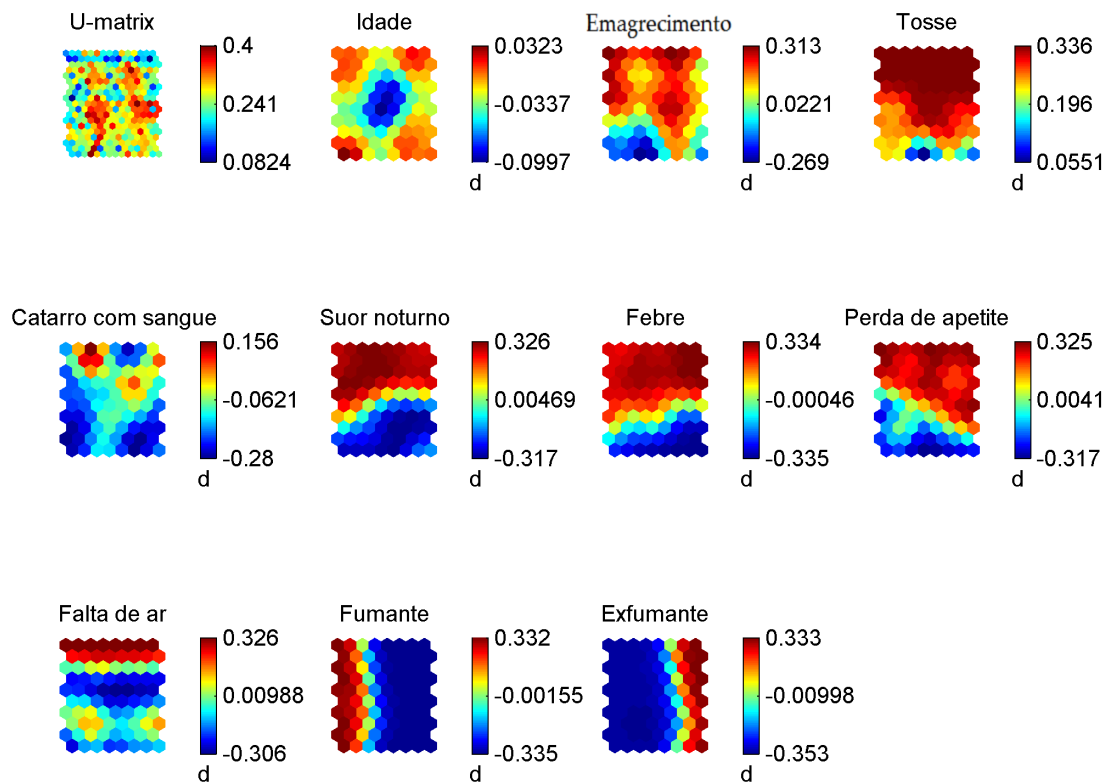


Figura 5.22: Mapa de componentes referente ao mapa SOM para pacientes com diagnósticos TB-DR e TB-MDR.

O mapa foi então clusterizado sem restrições quanto a quantidade de agrupamentos, onde o melhor valor é indicado pelo índice de Davis-Bouldin. A variação do valor desse índice para configurações que apresentam de 1 a 30 clusters é mostrada na Figura 5.23, onde é observado pouca variação desde a configuração com 6 clusters até um total de 16. Com isso, é escolhido um número mínimo de 6 agrupamentos, que tende a evitar uma granularidade alta, o que pode prejudicar a identificação dos agrupamentos.

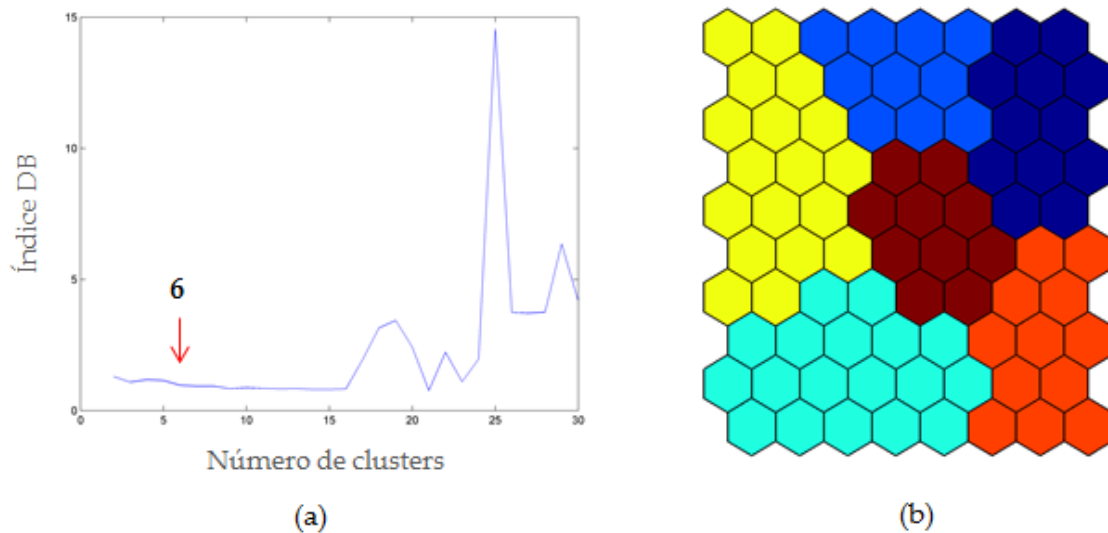


Figura 5.23: Clusterização do mapa SOM gerado para os pacientes com diagnósticos de TB-DR e TB-MDR. (a) Índice de Davis-Bouldin por número de clusters (b) Mapa SOM clusterizado com 6 agrupamentos

A rotulação de cada agrupamento foi realizada pela comparação da frequência de pacientes com diagnósticos diferentes em cada região do mapa, sendo o diagnóstico que aparece mais, o responsável pelo rótulo. A Figura 5.24 mostra o mapa de frequência através da variação de cores, que representam os extremos entre as presenças dos diferentes diagnósticos. É possível observar que apesar de regiões diferentes terem sido rotuladas pelos diagnósticos, não existe um cluster tão bem definido em direção a uma classe.

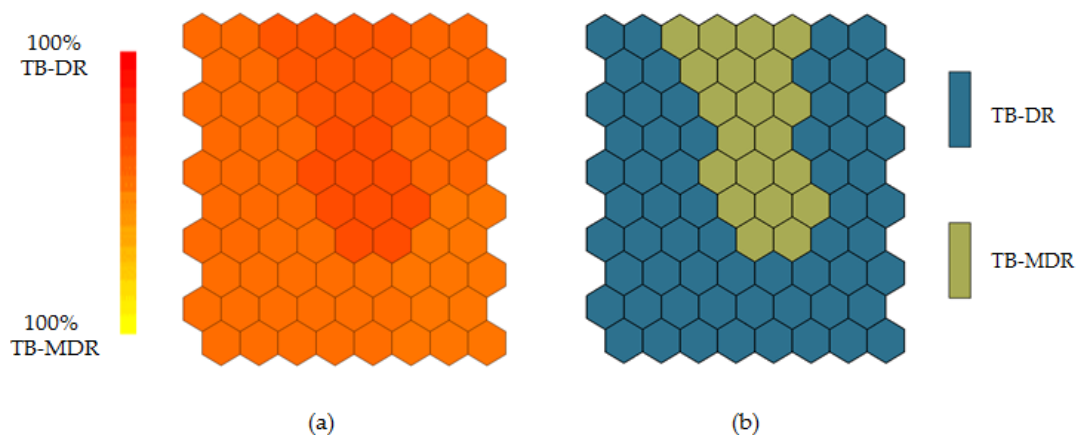


Figura 5.24: Rotulação do mapa SOM entre os diagnósticos de TB-DR e TB-MDR. (a) Frequência de diagnósticos em cada cluster (b) Rotulação dos clusters por diagnóstico

Com a análise por SOM, são observados índices de sensibilidade de 83,3%, especificidade de 25% e acurácia de 55,1%. Pela análise através do mapa das componentes, é possível ver indícios de que a ausência de Tosse e Perda de apetite se localizam em sua maioria no rótulo de TB-MDR. Apesar disso, pela distribuição das variáveis no mapa, assim como a rotulação realizada, não é possível fazer conclusões fortes sobre a relevância e contribuição das variáveis para esse caso.

5.4 Estudo entre diagnósticos de TB resistente, TB sensível e TB negativo

Após realizar os estudos, descritos na seção anterior, sobre classificadores especialistas utilizando redes neurais multicamadas, assim como sobre a relevância das variáveis e grupos de risco para cada classificação, foi proposto o desenvolvimento de um modelo neural único para ser responsável por classificar diretamente um paciente entre os todos os diagnósticos analisados. Para isso, a utilização de uma rede neural com 3 nós na camada de saída é analisada, assim como um modelo que utiliza as redes especialistas modeladas anteriormente em cascata, resultados os quais são descritos a seguir.

5.4.1 Rede neural dedicada

Assim como nos estudos anteriores, para realizar o processo de treinamento do discriminador, os dados foram divididos entre conjuntos de treino e teste, e a técnica de validação cruzada foi aplicada. A Figura 5.25 exibe os valores de SP para cada topologia ao se propagar pelas redes os conjuntos de treino e teste separadamente, sendo a incerteza gerada pelo valor RMS sobre as 50 redes de cada topologia. A configuração com 13 neurônios foi escolhida como ótima com base na melhor média de SP alcançada pela aplicação do conjunto de teste.

Para uma rede neural com três nós de saída, os índices de desempenho avaliados são o índice SP, sensibilidade, especificidade e acurácia, sendo seus valores descritos

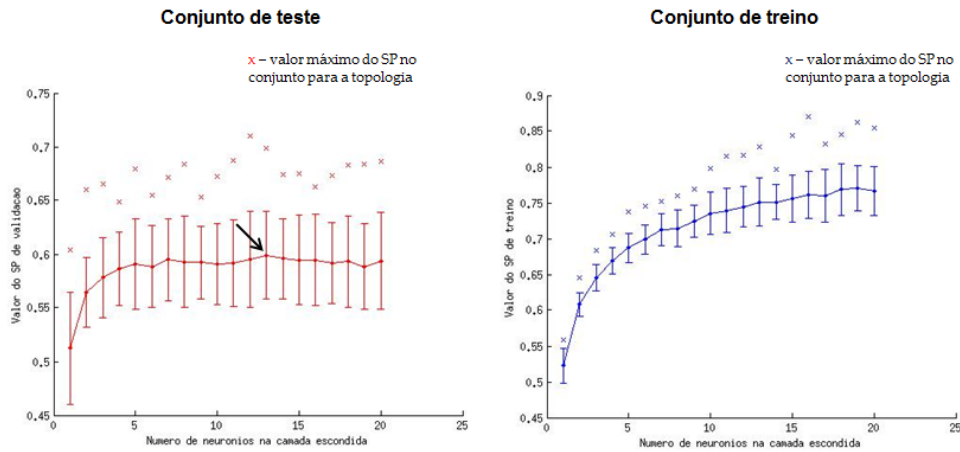


Figura 5.25: Valores de SP médios para cada topologia, pelos conjuntos de dados de teste e treino, considerando as redes obtidas pela validação cruzada. As barras de erro são representadas pelo valor RMS dos índices SP. A topologia com 13 neurônios é a que apresenta melhor média para o conjunto de teste.

na Tabela 5.7. A sensibilidade e especificidade são calculadas considerando a classe alvo o diagnóstico de tuberculose resistente. Suas incertezas são calculadas pelo valor RMS das 50 redes geradas com a topologia de 13 neurônios.

Desempenho do modelo neural	
Índice SP	$0,5989 \pm 0,0408$
Sensibilidade	0.5760 ± 0.0903
Especificidade	0.8334 ± 0.0235
Acurácia	$0,7183 \pm 0,0336$

Tabela 5.7: Desempenho alcançado pelo classificador para pacientes com diagnósticos TB resistente, TB sensível e TB negativo, considerando a topologia com 13 neurônios na camada escondida.

Para selecionar a rede de operação, o conjunto completo de dados (treino e teste) foi propagado pelas redes com a topologia considerada ótima, sendo escolhida a que teve maior valor de SP, nesse caso, a que apresentou um valor de 0,7275. Na Figura 5.26 é possível observar a evolução do treinamento da rede, que teve um ponto de parada em 60 épocas com um MSE de 0,752. A distribuição dos valores de saída do classificador também pode ser vista, onde, quando o alvo é TB resistente, apesar das distribuições das outras classes estarem em sua maior parte com valores negativos, e

apresentando picos próximos às extremidades, os valores para os dados de pacientes com resistência se encontram concentrados próximos ao valor 0. Situação similar é vista quando o alvo é o diagnóstico de TB sensível, o que não é uma situação esperada como ideal, onde as classes alvos e não-alvos possuiriam suas distribuições localizadas em extremidades opostas, sugerindo alguma dificuldade na classificação direta entre as três classes. Ainda, quando o alvo é TB negativo, a distribuição dos pacientes alvos são bem espalhadas dentre os limites estabelecidos, mas apesar disso, os outros diagnósticos apresentam picos bem próximos a extremidade negativa, o que se assemelha ao comportamento desejado.

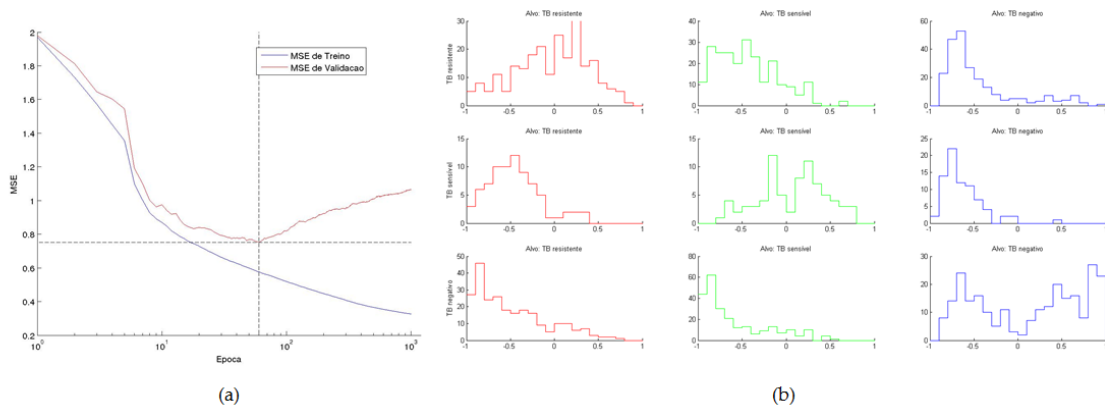


Figura 5.26: Representações da evolução do treino e distribuição dos valores de classificação para a rede neural de operação que classifica pacientes entre os diagnósticos de TB resistente, TB sensível e TB negativo.

Uma análise de relevância das variáveis foi realizada, sendo os resultados observados na Figura 5.30, onde as variáveis emagrecimento, catarro com sangue, falta de ar e fumante, apresentam valores negativos. Esse comportamento pode sugerir que essas variáveis são confundidoras, ou seja, dificultam a distinção das classes pela rede. Apesar disso, como os valores são muito baixos, não se pode fazer uma afirmação apenas com base nessa análise, sendo necessário para isso o estudo das redes após a remoção desses sintomas. As outras variáveis apresentam certa contribuição para o resultado, com maior destaque para idade e suor noturno.

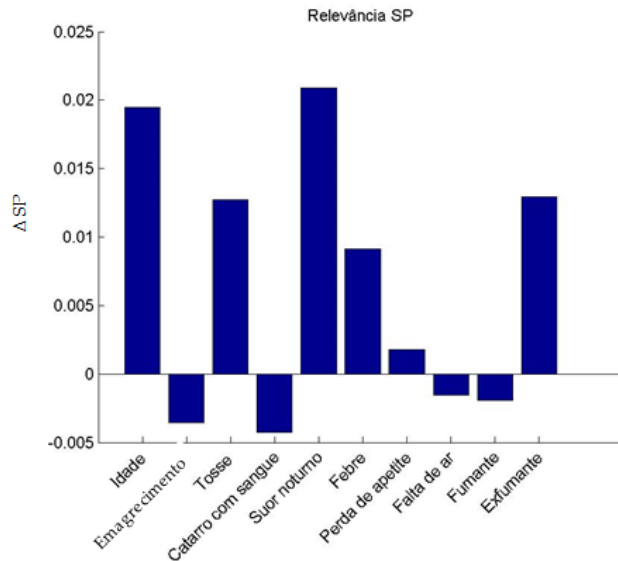


Figura 5.27: Estudo de relevância das variáveis na rede neural multicamadas obtida para os diagnósticos TB resistente, TB sensível e TB negativo. O índice SP da rede de operação é subtraído do índice SP alcançado pela rede quando se define cada variável com o seu valor médio.

5.4.2 Rede neural composta por redes especialistas

Outra proposta para a identificação das três classes através de um único modelo, foi a utilização das redes especialistas modeladas anteriormente em cascata, ou seja, o resultado da classificação de uma rede serve como indicação dos pacientes a serem classificados pela rede seguinte. As redes utilizadas foram as que obtiveram maiores taxas de sensibilidade, visando com isso, a operação em triagem, e especificidade, para a fase de diagnóstico. Como primeiro passo, o conjunto total de pacientes é propagado pela rede discriminadora para diagnósticos de TB positivo e TB negativo, sendo aqueles que foram indicados por essa rede como positivos para TB, aplicados na rede seguinte que classifica entre TB resistente e TB sensível. Esse processo é repetido, onde agora os pacientes indicados como TB resistente servem como entrada para a rede classificadora de TB-DR e TB-MDR.

A Figura 5.28 exibe uma simulação da situação de triagem, onde os pacientes são classificados entre os diagnósticos, em diferentes fases que se tornam cada vez mais específicas em direção a tuberculose que apresenta multirresistência aos medicamentos. É possível observar que na primeira fase, pouco menos de um terço dos

pacientes com diagnóstico negativo para TB são classificados erroneamente como TB positivo. Esses pacientes são propagados para a segunda fase de triagem, onde mais da metade desse conjunto é classificado como TB resistente. Além disso, é possível observar que apenas uma pequena taxa de pacientes com TB sensível à medicamentos é classificada como pacientes possuindo TB resistente. Já, na última fase de classificação, os pacientes classificados de forma equivocada, e que chegaram a esse estágio, são classificados em sua maioria como TB droga-resistente, que seria o tipo menos complexo da forma com resistência. Esse resultado diminuiria, por exemplo, o tratamento de doentes com medicamentos mais caros, os quais são fornecidos à pacientes com suspeita de TB-MDR. Ao final, índices de desempenho são obtidos para cada passo, sendo seus valores mostrados na Tabela 5.8.

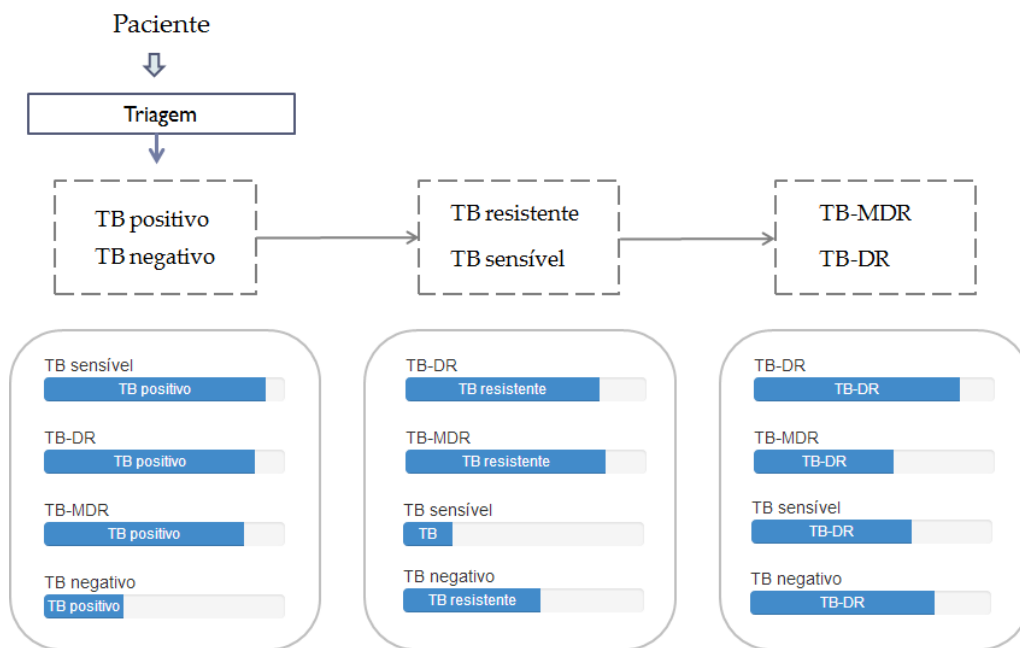


Figura 5.28: Modelo de redes neurais em cascata, composta por redes especialistas, para simular cenário de triagem médica.

Para o modelo de triagem, em comparação com os resultados alcançados anteriormente, é possível ver uma queda nos valores de sensibilidade e especificidade na etapa de classificação entre TB resistente e TB sensível, que agora possuem valores de 82,0% e 79,1%, e como rede especialista alcançava valores de 91% e 88,4% respectivamente. Com isso, uma diminuição do SP também é vista, caindo de 0,897 para o valor de 0,805. Um aumento de sensibilidade é visto na etapa de classificação

Triagem	Índice SP	Sensib. (%)	Especif. (%)	Acur. (%)
TB pos x TB neg	0.765	87,1	66,5	78,0
TB res x TB sens	0.805	82,0	79,1	81,2
TB-DR x TB-MDR	0.646	87,0	45,5	67,0
Diagnóstico	Índice SP	Sensib. (%)	Especif. (%)	Acur. (%)
TB pos x TB neg	0.809	81,7	80,2	81,0
TB res x TB sens	0.825	74,7	90,7	78,8
TB-DR x TB-MDR	0.580	39,1	80,8	60,5

Tabela 5.8: Desempenho alcançado pelo classificador composto por redes especialistas em cascata.

entre TB-DR e TB-MDR, onde a sensibilidade subiu de 85,0% para 87,0%. Já a especificidade diminuiu, possuindo 45,5% contra o valor de 52,6% alcançado anteriormente. O valor de SP obteve um decréscimo, apresentando nesse caso, um valor de 0,646 em comparação com 0,678 do estudo anterior.

Já para o modelo de diagnóstico, a sensibilidade apresenta quedas tanto para o estágio de classificação entre TB resistente e TB sensível, quanto para TB-DR e TB-MDR. Diminuição de valores também são vistas para a especificidade, entretanto, mantendo os índices acima de 90% no segundo estágio, e 80% no último, o qual discrimina entre os tipos de resistência.

Outro fator avaliado é a análise do desempenho das etapas em função do Valor Preditivo Positivo (VPP), o que traduz a quantidade de pacientes que de fato foram diagnosticados com a doença alvo e foram previstos corretamente pelo modelo como suspeitos, e do Valor Preditivo Negativo (VPN), que é análogo ao VPP, porém para casos sem a doença. Os valores de VPP e VPN foram estimados para as 3 etapas de classificação, considerando prevalências (probabilidades pré-teste) que variam de 2% a 50%.

No modelo de triagem, para o estágio intermediário, para classificação de TB resistente e sensível, para a população com uma prevalência de 20%, um VPN de 94,6% é encontrado, Já no último estágio de classificação, para a população com a mesma prevalência, alcança-se um VPN de 92,2%.

Analisando o modelo de diagnóstico, para a segunda etapa de classificação, e com prevalência de 20%, um total de 66,9% dos pacientes são classificados com a doença ao passarem pela triagem, já na última etapa, o valor de VPP é de 33,7%, ou seja, a cada 3 pacientes, um é diagnosticado positivamente para a doença de interesse.

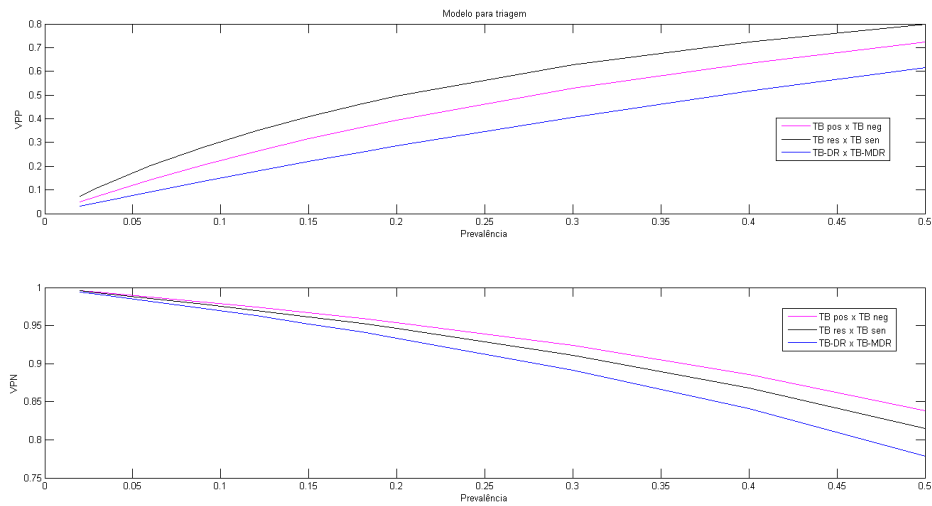


Figura 5.29: Desempenho do valor preditivo positivo (VPP) e valor preditivo negativo (VPN), por prevalência, para o modelo de triagem.

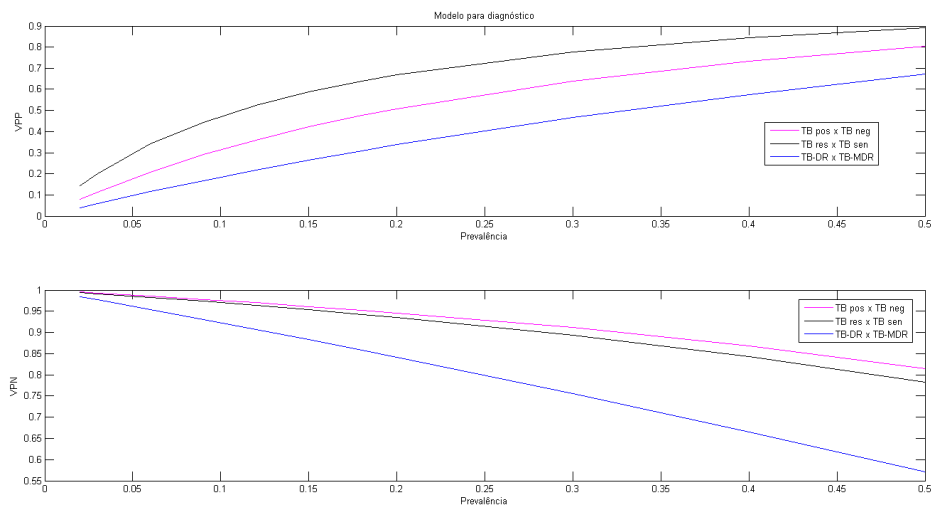


Figura 5.30: Desempenho do valor preditivo positivo (VPP) e valor preditivo negativo (VPN), por prevalência, para o modelo de diagnóstico.

Capítulo 6

Conclusão

Como descrito anteriormente, os métodos tradicionais para diagnóstico de tuberculose pulmonar apresentam limitações. A baciloscopia apresenta sensibilidade de 40% a 60%, enquanto a cultura, apesar de apresentar sensibilidade de 70% a 80%, seu resultado necessita de 4 a 8 semanas de espera, período no qual a doença pode ser agravada, além de poder ocorrer a transmissão para outros indivíduos.

Já os tipos da doença que apresentam resistência aos medicamentos indicados para o tratamento, não conseguem ser detectados pela baciloscopia, a qual não discrimina as cepas características dessas variações da doença. Através da cultura é possível identificar a TB resistente, porém, com as mesmas características de demora no resultado, além de precisar de uma infra-estrutura maior em biossegurança, o que aumenta a dificuldade de sua implementação em lugares com poucos recursos.

Outros métodos, como a radiografia de tórax, podem ser também utilizados, mas a um custo mais elevado, assim como apresenta uma logística necessária mais complexa. Testes mais recentes, como o Xpert® MTB/RIF, conseguem bons índices de sensibilidade, além de identificar resistência à medicamentos, porém, seus custos são elevados, assim como possui necessidades específicas de infra-estrutura, sendo características que dificultam na sua implantação em larga escala.

Vista essa situação, o objetivo desse trabalho foi fornecer uma ferramenta para o apoio a tomada de decisão relativa a triagem de tuberculose pulmonar e suas variações resistentes à medicamentos, utilizando, para isso, informações de fácil obtenção

através de anamnese. A modelagem de redes neurais multicamadas foi proposta para realizar a classificação dos pacientes entre os diagnósticos de TB positivo, incluindo TB sensível, TB-DR e TB-MDR, e TB negativo. Para um melhor entendimento dos sintomas e variáveis utilizadas nas análises, técnicas de redes neurais não supervisionadas, no caso, os mapas auto-organizáveis, foram utilizadas para a identificação de agrupamentos existentes entre os pacientes, assim como a avaliação de grupos de risco.

Para o modelo que classifica os pacientes entre TB positivo e TB negativo, a sensibilidade alcançada foi de $0,8205 \pm 0,0676$ com especificidade de $0,7286 \pm 0,0548$. Através da análise de relevância por influência no índice SP da rede, os sintomas catarro com sangue e falta de ar foram o que apresentaram menor contribuição aparente. A falta de ar também foi apontada com uma menor relevância pelo estudo do SOM, entretanto, catarro com sangue apresentou comportamento não só compatível ao grupo de alto risco, mas com o grupo rotulado no mapa como TB positivo. Já a falta de ar ainda não apresentou correspondência à rotulação realizada no mapa entre os dois diagnósticos, fazendo com que esse conjunto de análises infra que o sintoma falta de ar não é relevante para o diagnóstico proposto. Outras variáveis ainda apresentam comportamentos que indicam sua correlação, como suor noturno e febre, enquanto pacientes com idade elevada puderam ser vistos com maiores risco de apresentar tuberculose. Os sintomas de suor noturno, febre, perda de apetite e tosse apresentaram características relacionadas com a rotulação dos grupos em TB positivo, indicando a possibilidade de serem relevantes a identificação da classe.

A classificação do discriminador neural para TB resistente e TB sensível apresentou valores de sensibilidade de $0,7308 \pm 0,0813$, mais baixa em relação a classificação de TB positivo e TB negativo, além de uma especificidade de $0,7986 \pm 0,0791$. Pela análise de relevância das variáveis, através de sua contribuição para o índice SP da rede, foi possível ver que o sintoma tosse é o que aparenta menos contribuição, enquanto falta de ar se destaca inversamente. A tosse ainda se apresenta espalhada no mapa SOM, sendo mais um indício de que não é relevante para essa classificação,

assim como a falta de ar, o que não ajuda a confirmar a suspeita anterior de sua contribuição mais elevada. Os sintomas febre e suor noturno continuam apresentando comportamentos correlacionados, e pacientes com idade elevada parecem se associar mais, nesse caso, ao grupo de baixo risco.

O classificador entre TB-DR e TB-MDR obteve índices de sensibilidade de $0,7067 \pm 0,0967$ e especificidade de $0,7308 \pm 0,0948$, apresentando incertezas maiores que as encontradas nos classificadores anteriores. Dentre os sintomas avaliados, a tosse apresentou a menor contribuição na identificação das classes, segundo a análise realizada pelo SP da rede. Pela análise com SOM, os sintomas suor noturno e febre continuam a apresentar comportamentos aparentemente correlacionados. A tosse se encontra espalhada por boa parte do mapa, ajudando a validar o indício de sua baixa contribuição na classificação, mesmo comportamento apresentado pelo sintoma de emagrecimento. Sendo um sintoma característico de tuberculose, é natural que a tosse não seja um fator de decisão, tanto na classificação de TB resistente e TB sensível, como entre TB-DR e TB-MDR.

Após o estudo com as redes especialistas, modelos únicos para a classificação entre os três diagnósticos foram propostos, um sendo uma rede neural com uma camada de saída possuindo três neurônios, e outro, utilizando as redes especialistas para classificar pacientes em cascata. O índice SP, para a rede com 3 nós de saída, alcançou um valor de $0,5989 \pm 0,0408$, valor bem abaixo dos alcançados anteriormente, indicando uma maior dificuldade em identificar os três diagnósticos sem a realização de uma triagem prévia dos pacientes. As redes em cascata apresentaram valores próximos aos que foram encontrados na análise individual de cada rede, sendo uma boa opção para a análise por um modelo unificado entre todos os diagnósticos.

Em medicina, um método de diagnóstico que apresente uma sensibilidade maior que 80% é considerado um método sensível. Com isso, os resultados encontrados estão dentro desse limite, considerando os valores de incerteza obtidos. Como uma grande vantagem do modelo neural gerado está a possibilidade de se identificar os tipo de tuberculose droga-resistente e multirresistente, diagnósticos de difícil ob-

tenção por métodos usuais, seja pelo custo, equipamentos necessários ou tempo de espera para o resultado.

Existem estudos que apresentam resultados para a classificação entre os diagnósticos positivo e negativo para TB, utilizando modelos estatísticos, que alcançam 100% de sensibilidade e 80% de especificidade. Para esse trabalho, era considerado razoável a obtenção de valores menores a esses, pois são utilizados apenas pacientes que apresentam suspeita de TB resistente, assim como a própria inclusão de pacientes com o diagnóstico de tuberculose com resistência em conjunto com sensíveis. Adiciona-se a isso a presença de pacientes de diferentes localizações do Brasil, que podem apresentar características específicas a região de triagem, e que não foram consideradas para as classificações realizadas.

Como trabalhos futuros, é vista a possibilidade de se aprofundar na análise sobre as regiões de triagem dos pacientes. Modelos neurais para cada região podem ser treinados para avaliar se os resultados conseguem ser melhorados caso os modelos sejam especializados, considerando as características regionais de cada conjunto de pacientes. Outras arquiteturas de modelo único de classificação ainda podem ser testadas, como o treinamento de uma rede que utilize como entrada a saída da camada escondida de todas as redes especialistas, formando ao final uma rede unificada com cinco camadas. O uso de comitês para decisão entre os diagnósticos das redes pode ser desenvolvido, assim como a opção de rejeição de pacientes para uso durante o treinamento da rede neural multicamadas, a qual tentaria diminuir o nível de confusão entre os diagnósticos. Mais variáveis ainda podem ser incluídas no modelo para observar os resultados, características ao tipo de tuberculose resistente, como os achados radiológicos, tempo em que o paciente apresenta os sintomas, contato com doentes com TB-MDR e TB-DR, ou outros sintomas e/ou coinfeções.

Referências Bibliográficas

- [1] WHO. *Global Tuberculosis Report 2013*. World Health Organization, WHO, 2013.
- [2] SECRETARIA DE VIGILÂNCIA EM SAÚDE, MINISTÉRIO DA SAÚDE. *Boletim Epidemiológico Tuberculose - 2013*. Volume 44 N° 2 -2013.
- [3] “BRASIL, Ministério da Saúde. Tuberculose: Casos confirmados notificados no Sistema de Informação de Agravos de Notificação (Sinan)”. Acessado em 18 Ago. 2013.
<http://dtr2004.saude.gov.br/sinanweb/tabnet/tabnet?sinannet/>.
- [4] WORLD HEALTH ORGANIZATION. *Multidrug-resistant tuberculosis - 2013 Update*. WHO, 2013.
- [5] A.A., F., K.C., Q., K.P., T., et al. *Os fatores associados à tuberculose pulmonar e a baciloscopia: uma contribuição ao diagnóstico nos serviços de saúde pública*. Rev Bras Epidemiol 2005; 8(2): 142-9., 2005.
- [6] CONDE, M.B., MELLO, F., LIMA, M.A., GUERRA, R.L., MIRANDA, S.S., GALVÃO, T.S., PINHEIRO, V.G., CARVALHO, N.B. *Tuberculose Pulmonar: Diagnóstico - Técnicas Convencionais*. Sociedade Brasileira de Pneumologia e Tisiologia, 2011.
- [7] “Foundation for Innovative New Diagnostics (FIND), 2011. Frequently Asked Questions on Xpert MTB/RIF assay.” Acessado em 30 Dez. 2013.
http://www.finddiagnostics.org/export/sites/default/Xpert_FAQs.pdf.
- [8] DONOGHUE, H. “Human tuberculosis – an ancient disease, as elucidated by ancient microbial biomolecules”. In: *Microbes and Infection*, p. 1156–1162, dez. 2009.
- [9] HALL, F. *Tuberculosis and allied diseases*. 1 ed. Michigan, Kalamazoo, Mich., The Yonkerman company, 1975.

- [10] PROGRAMA NACIONAL DE CONTROLE DA TUBERCULOSE. *Manual de Recomendações para o Controle da Tuberculose no Brasil*. Ministério da Saúde, Secretaria de Vigilância em Saúde, 2010.
- [11] DEPARTAMENTO DE ATENÇÃO BÁSICA. *Manual Técnico para o Controle da Tuberculose*. Ministério da Saúde, Secretaria de Políticas de Saúde, 2002.
- [12] UNITED NATIONS. *The Millennium Development Goals Report*. World Health Organization, 2013.
- [13] WHO. *Global Tuberculosis Report 2012*. World Health Organization, WHO, 2012.
- [14] PAWLOWSKI, A., JANSSON, M., SKÖLD, M., et al. “Tuberculosis and HIV Co-Infection”. In: *PLoS Pathog*, n. 8(2): e1002464. doi:10.1371/journal.ppat.1002464, fev. 2012.
- [15] SCHAAF, H. S., ZUMLA, A. *Tuberculosis. A Comprehensive Clinical Reference*. 1 ed. , Saunders, 2009.
- [16] NATAL, S. *Emergência da resistência às drogas*. Bol. Pneumol. Sanit. vol.10 no.2 Rio de Janeiro Dec. 2002, 2002.
- [17] “Multidrug-resistant tuberculosis, World Health Organization”. Acessado em 18 Ago. 2013. <http://www.who.int/tb/challenges/mdr/en/index.html>.
- [18] PROGRAMA NACIONAL DE CONTROLO DA TUBERCULOSE - MOÇAMBIQUE. *Manual de Diagnóstico e Tratamento de Tuberculose Resistente e Multi-Droga Resistente*. Organização Mundial de Saúde, 2009.
- [19] DALCOLMO, M. P., DE NORONHA ANDRADE, M. K., PICON, P. D. *Tuberculose multirresistente no Brasil: histórico e medidas de controle*. Rev Saúde Pública 2007;41(Supl. 1):34-42, 2007.
- [20] DALCOLMO, M. P., FORTES, A., MELO, F. F. D., et al. *Estudo de efetividade de esquemas alternativos para o tratamento da tuberculose multirresistente no Brasil*. J. Pneumologia, Abr 1999, vol.25, no.2, p.70-77. ISSN 0102-3586, 1999.
- [21] TEIXEIRA, G. M. *XDR-TB - uma grave e emergente ameaça à saúde pública*. Rev. Bras. Pneumol. Sanit. vol.15 no.1 Rio de Janeiro Dec., 2007.

- [22] WORLD HEALTH ORGANIZATION. *WHO - Frequently asked questions - XDR-TB*. WHO, 2013.
- [23] STOP TB PARTNERSHIP. *The global plan to stop TB*. Stop TB Partnership, 2011.
- [24] OBERMEYER, Z., ABBOTT-KLAFTER, J., MURRAY, C. J. L. *Has the DOTS Strategy Improved Case Finding or Treatment Success? An Empirical Assessment*. 25. PLoS ONE. 2008;3(3):e1721., 2008.
- [25] KRITSKI, A. L., CONDE, M. B., MUSY, G. R. *Tuberculose: do ambulatório a enfermaria*. Atheneu, 2006.
- [26] “Risk Factors, Center of Disease Control and Prevention”.
Acessado em 18 Ago. 2013.
<http://www.cdc.gov/tb/topic/basics/risk.htm>.
- [27] MUNIZ, J. N., RUFFINO-NETTO, A., VILLA, T. C. S., et al. *Aspectos epidemiológicos da co-infecção tuberculose e vírus da imunodeficiência humana em Ribeirão Preto (SP), de 1998 a 2003*. J. bras. pneumol. vol.32 no.6 São Paulo Nov./Dez, 2006.
- [28] CHAISSON, R. E., MARTINSON, N. A. *Tuberculosis in Africa — Combating an HIV-Driven Crisis*. N Engl J Med 2008; 358:1089-1092, 2008.
- [29] HIJJAR, M. A., PROCÓPIO, M. J., DE FREITAS, L. M. R., et al. *Epidemiologia da tuberculose: importância no mundo, no Brasil e no Rio de Janeiro*. Pulmão RJ, pp. 310–314, 2005.
- [30] CANTALICE FILHO, J.P, SANT’ANNA, C.C., BÓIA, M.N. *Aspectos clínicos da tuberculose pulmonar em idosos atendidos em hospital universitário do Rio de Janeiro, RJ, Brasil*. J Bras Pneumol 2007;33(6):699-706, 2007.
- [31] MACIEL, E.L., GOLUB, J.E., PERES, R.L., HADAD, D.J., FÁVERO, J.L., MOLINO, L.P. *Delay in diagnosis of pulmonary tuberculosis at a primary health clinic in Vitoria, Brazil*. Int J Tuberc Lung Dis. 2010;14(11):1403-10, 2010.
- [32] NEWTON S.M. AND BRENT A.J. AND ANDERSON S. AND WHITTAKER E. AND KAMPMANN B. *Paediatric tuberculosis*. Lancet Infect Dis 2008;8:498-510, 2008.
- [33] COULTER, J.B. *Diagnosis of pulmonary tuberculosis in young children*. Ann Trop Paediatr 2008;28:3-12., 2008.

- [34] WORLD HEALTH ORGANIZATION. *Tuberculosis prevalence surveys: a handbook*. Geneva, Switzerland: World Health Organization, 2011.
- [35] B.J., M., A., G., J.R., S., et al. *Tuberculosis in women and children*. Lancet 2010;375:2057-9., 2010.
- [36] WORLD HEALTH ORGANIZATION. *Guidance for national tuberculosis programmes on the management of tuberculosis in children*. Geneva, Switzerland: World Health Organization, 2006.
- [37] C., L., J., S., K., F. *Risk factors for tuberculosis infection in children in contact with infectious tuberculosis cases in The Gambia, West Africa*. Pediatrics 2003;111:e608-14., 2003.
- [38] M., S., A.G., S., N., A., et al. *Latent tuberculosis in children: diagnosis and management*. Indian J Pediatr 2011;78:464-8., 2011.
- [39] DE QUEIROZ MELLO, F. C. *Modelos Preditivos para o Diagnóstico da Tuberculose Pulmonar Paucibacilar*. Tese de D.Sc., Faculdade de Medicina/UFRJ, Rio de Janeiro, RJ, Brasil, 2001.
- [40] GRANICH R. ET AL. *Guidelines for the prevention of tuberculosis in health facilities in resource-limited settings*. Geneva, World Health Organization, 1999:1-51 (document WHO/CDC/TB/99.269), 1999.
- [41] BARROSO, E. C., MOTA, R. S., SANTOS, R. O., et al. *Fatores de risco para tuberculose multirresistente adquirida*. J Pneumol 2003;29(2):89-97, 2003.
- [42] DE SOUZA, M. B., DE FIGUEIREDO ANTUNES, C. M., GARCIA, G. F. *Perfil de sensibilidade e fatores de risco associados à resistência do Mycobacterium tuberculosis, em centro de referência de doenças infecto-contagiosas de Minas Gerais*. J. bras. pneumol. vol.32 no.5 São Paulo Set./Out, 2006.
- [43] BRASIL. MINISTÉRIO DA SAÚDE. SECRETARIA DE VIGILÂNCIA EM SAÚDE. PROGRAMA NACIONAL DE CONTROLE DA TUBERCULOSE. *Manual de Recomendações para o Controle da Tuberculose no Brasil*. Brasília, 2010.
- [44] LOPES, A. J., CAPONE, D., MOGAMI, R., et al. *Tuberculose extrapulmonar: aspectos clínicos e de imagem*. Pulmão RJ 2006;15(4):253-261, 2006.
- [45] DA SILVA, C. C. A. V., ANDRADE, M. S., CARDOSO, M. D. *Fatores associados ao abandono do tratamento de tuberculose em indivíduos acompanhados em unidades de saúde de referência na cidade do Recife, Estado*

de Pernambuco, Brasil, entre 2005 e 2010. *Epidemiol. Serv. Saúde* vol.22 no.1 Brasília Mar., 2013.

- [46] DE MELO MENDES, A., FENSTERSEIFER, L. M. *Tuberculose: porque os pacientes abandonam o tratamento?* *Bol. Pneumol. Sanit.* vol.12 no.1 Rio de Janeiro Apr., 2004.
- [47] MINISTÉRIO DA SAÚDE. *Plano Estratégico para o Controle da Tuberculose, Brasil 2007-2015.* 2006.
- [48] VIEIRAI, A. A., RIBEIRO, S. A. *Abandono do tratamento de tuberculose utilizando-se as estratégias tratamento auto-administrado ou tratamento supervisionado no Programa Municipal de Carapicuíba, São Paulo, Brasil.* *J. bras. pneumol.* vol.34 no.3 São Paulo Mar., 2008.
- [49] WORLD HEALTH ORGANIZATION. *Guidelines for the Prevention of Tuberculosis in Health Care Facilities in Resource-Limited Settings.* WHO, 1999.
- [50] CENTERS FOR DISEASE CONTROL AND PREVENTION - CDC. *Guidelines for preventing the transmission of Mycobacterium tuberculosis in health care facilities.* *MMWR* 1994;43(RR-13):1-131, 1994.
- [51] DE SIQUEIRA, H. R. *Enfoque Clínico da Tuberculose Pulmonar.* *Pulmão RJ* 2012;21(1):15-18, 2012.
- [52] BORGDORF, M.W., FLOYD, K, BROEKMANS, J.F. *Interventions to reduce tuberculosis mortality and transmission in low and middle income countries.* *Bull World Health Health Organ.* 2002;80(3):217-27, 2002.
- [53] MINISTÉRIO DA SAÚDE. *Tratamento diretamente observado da tuberculose na atenção básica.* 2010.
- [54] CAMPOS, H. S. *Diagnóstico da Tuberculose.* *Pulmão RJ* 2006;15(2):92-99, 2006.
- [55] ENARSON, D. A., RIEDER, H. L., ARNADOTTIR, T., et al. *Management of tuberculosis: a guide for low income countries.* Paris: International Union Against Tuberculosis and Lung Disease, 2000.
- [56] BASTOS, L., FONSECA, L., MELLO, F., et al. *Prevalence of pulmonary tuberculosis among respiratory symptomatic subjects in an out-patient primary health unit.* International Union Against Tuberculosis and Lung Disease, 2007.

- [57] PRASAD, R., NAUTIYAL, R.G., MUKHERJI, P.K., JAIN, A., SINGH, K., AHUJA, R.C. *Diagnostic evaluation of pulmonary tuberculosis: what do doctors of modern medicine do in India?* The International Journal of Tuberculosis and Lung Disease, Volume 7, Number 1, January 2003 , pp. 52-57(6), 2003.
- [58] SINGLA, N., SHARMA, P.P., SINGLE, R., JAIN, R.C. *Survey of knowledge, attitudes and practices for tuberculosis among general practitioners in Delhi, India.* Int J Tuberc Lung Dis. 1998 May;2(5):384-9., 1998.
- [59] FRIEDMAN, H. H. *Manual de Diagnóstico Clínico.* 3 ed. Rio de Janeiro, 1985.
- [60] AGUIAR, F. S. *Redes neurais artificiais e árvores de classificação e regressão para previsão do diagnóstico de tuberculose pulmonar em pacientes hospitalizados.* Tese de D.Sc., Faculdade de Medicina/UFRJ, Rio de Janeiro, RJ, Brasil, 2013.
- [61] A, V. R., PAGE-SHIPPI, L, S., et al. *Xpert(®) MTB/RIF for point-of-care diagnosis of TB in high-HIV burden, resource-limited countries: hype or hope?* Expert Rev Mol Diagn. 2010; 10(7):937-46., 2010.
- [62] SCHIRM, J., OOSTENDORP, L. A., MULDER, J. G. *Comparasion of amplificador, in house PCR and conventional culture for detection of mycobacterium in clinical samples.* Journal Clinical Microbiology, pp. 3321–3324, 1995.
- [63] “World Health Organization. Health topics - Tuberculosis. 2010”.
Acessado em 30 Dez. 2013.
<http://www.who.int/topics/tuberculosis/es/>.
- [64] A., K. *Emergência de tuberculose resistente: renovado desafio.* J Bras Pneumol. 2010; 36(2):157-158, 2010.
- [65] WORLD HEALTH ORGANIZATION. *Automated Real-time Nucleic Acid Amplification Technology for Rapid and Simultaneous Detection of Tuberculosis and Rifampicin Resistance: Xpert MTB/RIF System.* WHO., 2011.
- [66] “Boletim Brasileiro de Avaliação de Tecnologias em Saúde, Anvisa.”
Acessado em 30 Dez. 2013.
<http://portal.anvisa.gov.br/wps/wcm/connect/>.
- [67] WALLIS, R.S., PAI M., MENZIES D. *Biomarkers and diagnostics for tuberculosis: progress, needs, and translation into practice.* Lancet 2010; published online May 19. DOI:10.1016/S0140-6736(10)60359-5., 2010.

- [68] CANETTI, G., RIST, N., GROSSET, J. *Measurement of sensitivity of the tuberculous bacillus to anti-bacillary drugs by the method of proportions. Methodology, resistance criteria, results and interpretation.* Rev Tuberc Pneumol (Paris). 1963; 27:217-72, 1963.
- [69] DALCOLMO, M. P., DE NORONHA ANDRADE, M. K., PICON, P. D. *Tuberculose multirresistente no Brasil: histórico e medidas de controle.* Rev Saúde Pública 2007;41(Supl. 1):34-42, 2007.
- [70] MITCHISON, D.A. *Drug Resistance in tuberculosis.* Eur Respir J. 2005;25(2):376-9, 2005.
- [71] PICON, P.D., RIZZON, C.F.C., FREITAS, T.M., AZEVEDO, S.N.B., GUTIERREZ, R.S. *Resultados do tratamento.* Tuberculose: Epidemiologia, Diagnóstico e tratamento em clínica e saúde pública. Rio de Janeiro: Editora Médica e Científica Ltda; 1993. p.504-23, 1993.
- [72] AGRESTI, A. *An Introduction to Categorical Data Analysis.* Wiley, 2007.
- [73] KEVIN B. KORB, A. E. N. *Bayesian Artificial Intelligence.* Chapman Hall /CRC, 2003.
- [74] MITCHELL, T. *Machine Learning.* McGraw Hill, 1997.
- [75] HAYKIN, S. *Neural Networks and Learning Machines.* Prentice-Hall, Inc., 2008.
- [76] BOCK, N.N., MCGOWAN, J.E. JR, AHN, J., TAPIA, J., BLUMBERG, H.M. *Clinical predictors of tuberculosis as a guide for a respiratory isolation policy.* Am J Respir Crit Care Med. 1996 Nov;154(5):1468-72, 1996.
- [77] GORDIN, F. M., NELSON, E. T., MATTS, J. P., et al. *The impact of human immunodeficiency virus infection on drug-resistant tuberculosis.* American Journal of Respiratory and Critical Care Medicine, Vol. 154, No. 5 (1996), pp. 1478-83., 1996.
- [78] SAMB, B., HENZEL, D., DALEY, C.L., MUGUSI, F., NIYONGABO, T., MLIKA-CABANNE, N., KAMANFU, G. *Methods for diagnosing tuberculosis among in-patients in eastern Africa whose sputum smears are negative.* Int J Tuberc Lung Dis. 1997 Feb;1(1):25-30., 1997.
- [79] VEROPOULOS, K., CAMPBELL, C. *The Automated Identification of Tubercle Bacilli using Image Processing and Neural Computing Techniques.* Proceedings of the 8th International Conference on Artificial Neural Networks, vol 2, pp 797-802, 1998.

- [80] EL-SOLH, A. A., HSIAO, C.-B., GOODNOUGH, S., et al. *Predicting Active Pulmonary Tuberculosis Using an Artificial Neural Network*. Chest. 1999;116(4):968-973. doi:10.1378/chest.116.4.968, 1999.
- [81] MELLO, F. C. Q. *Modelos preditivos para tuberculose pulmonar paucibacilar*. Tese de D.Sc., Faculdade de Medicina / UFRJ, Rio de Janeiro, RJ, Brasil, 2001.
- [82] SANTOS, A. M. *Redes Neurais e Árvores de Classificação Aplicadas ao Diagnóstico da Tuberculose Pulmonar Paucibacilar*. Tese de D.Sc., COPPE/UFRJ, Rio de Janeiro, RJ, Brasil, 2003.
- [83] SANTOS, A. M., PEREIRA, B., SEIXAS, J. M. *Neural networks: An Application for Predicting Smear Negative Pulmonary Tuberculosis*". Advances in Statistical Methods for the Health Sciences, pp. 279–289, 2007.
- [84] BENFU, Y., HONGMEI, S., YE, S., et al. *Study on the artificial neural network in the diagnosis of smear negative pulmonary tuberculosis*. WRI World Congress on Computer Science and Information, v. 5, pp. 584 – 588, 2009, 2009.
- [85] AGUIAR F. S. AND ALMEIDA L. L. AND RUFFINO-NETO A. AND KRITSKI A. L. AND MELLO F. C. AND WERNECK G. L. *Classification and Regression Tree (CART) Model to Predict Pulmonary Tuberculosis in Hospitalized Patients*. BMC Pulmonary Medicine, 2012.
- [86] CASCÃO, L. V. C. *Modelos de inteligência computacional para apoio da triagem de pacientes e diagnóstico clínico de tuberculose pulmonar*. Tese de M.Sc., COPPE/UFRJ, Rio de Janeiro, RJ, Brasil, 2011.
- [87] SEIXAS, J. M., SOUZA FILHO, J. B. D. O. E., MELLO, F. C. D. Q., et al. *An Intelligent System for Managing the Isolation of Patients Suspected of Pulmonary Tuberculosis*. International Conference on Intelligent Data Engineering and Automated Learning, 2012, Natal, Rio Grande do Norte. 2012. v. 7435. p. 818-825, 2012.
- [88] SEIXAS, J. M., FARIA, J, SOUZA FILHO, J. B., VIEIRA, A. F., KRITSKI, A., TRAJMAN, A. *Artificial neural network models to support the diagnosis of pleural tuberculosis in adult patients*. The International Journal of Tuberculosis and Lung Disease, Volume 17, Number 5, 1 May 2013 , pp. 682-686(5), 2013.
- [89] EDWARDS, A. *The measure of association in a 2x2 table*. Journal of the Royal Statistical Society, v. 126, pp. 109-114, 2009.

- [90] TREES, H. *Detection, Estimation and Modulation Theory, Part III*. Nova York, John Wiley Sons, Inc, 2001.
- [91] SIMAS FILHO, E. *Análise Não-Linear de Componentes Independentes para uma Filtragem Online Baseada em Calorimetria de Alta Energia e com Fina Segmentação*. Tese de Doutorado, COPPE/UFRJ, Rio de Janeiro, 2010.
- [92] CIODARO, T., DEVA, D., DE SEIXAS, J., ET AL. *Online particle detection with Neural Networks based on topological calorimetry information*. In: Journal of Physics: Conference Series, v. 368, p. 012030. IOP Publishing, 2012.
- [93] “SOM implementation in SOM Toolbox”.
Julho 2013. Disponível em:
<http://www.cis.hut.fi/somtoolbox/documentation/somalg.shtml>.
- [94] JOLLIFFE, I.T. *Principal Component Analysis*. Springer, 2nd ed. 2002, XXIX, 489 p., 2002.
- [95] DAVIES, D., BOULDIN, D. *A cluster separation measure*. Pattern Analysis and Machine Intelligence, IEEE Transactions on, n. 2, pp. 224–227, 1979.
- [96] KOHONEN, T. *Self-Organizing Maps*. Springer, 2000.
- [97] WASSERMAN, P. *Neural computing: theory and practice*. E.U.A., Van Nostrand Reinhold Co., 1989.
- [98] BISHOP, C., OTHERS. *Neural networks for pattern recognition*. 1995.
- [99] ICHIKAWA, Y., SAWA, T. *Neural network application for direct feedback controllers*. Neural Networks, IEEE Transactions on, v. 3, n. 2, pp. 224–231, 1992.
- [100] DENG, W., CHEN, R., GAO, J., ET AL. *A novel parallel hybrid intelligence optimization algorithm for a function approximation problem*. Computers Mathematics with Applications, 2011.
- [101] LAPEDES, A., FARBER, R. *Nonlinear signal processing using neural networks: Prediction and system modelling*. 1987.
- [102] RUMELHART, D. E.; HINTON, G. E., WILLIAMS, R. J. *Learning representations by back-propagating errors*. Nature 323 (6088): 533–536, 1986.

- [103] KOHONEN, T. *Self-Organizing Maps*. Secaucus, NJ, USA, Springer-Verlag New York, Inc., 1997. ISBN: 3-540-62017-6, 1997.
- [104] SIMAS FILHO, E. *Análise Não-Linear de Componentes Independentes para uma Filtragem Online Baseada em Calorimetria de Alta Energia e com Fina Segmentação*. Tese de Doutorado, COPPE/UFRJ, Rio de Janeiro, 2010.
- [105] ULTSCH, A. *Self-Organizing Neural Networks for Visualization and Classification*. Information and Classification. Springer.
- [106] COSTA, J. A. F. *Uma Nova Abordagem para Visualização e Detecção de Agrupamentos em Mapas de Kohonen Baseado em Gradientes das Componentes*. Learning and Nonlinear Models, Journal of the Brazilian Neural Network Society, v. 9, pp. 20–31, 2011.
- [107] VESANTO, J., ALHONIEMI, E. *Clustering of the Self-Organizing Map*. IEEE Transactions on Neural Networks, v. 11, pp. 586–600, 2000.

Apêndice A

Redes neurais multicamadas

As redes neurais multicamadas são uma classe de redes neurais artificiais (RNA) [97], as quais são modelos matemáticos criados com inspiração nas características do cérebro humano, capazes de absorver conhecimento e generalizá-lo. A generalização consiste na capacidade da rede em aprender através de um conjunto reduzido de exemplo e, a posteriori, fornecer respostas coerentes a dados não apresentados anteriormente.

Para isso, a ideia básica é construir um modelo composto por um número de unidades de processamento muito simples, que são chamadas de neurônios, com conexões entre eles. O processamento básico de informação na rede ocorre nos neurônios, enquanto a transmissão da informação é realizada pelas conexões, denominadas sinapses ou pesos sinápticos.

Uma diferença relevante entre as redes neurais e os métodos clássicos, é que não é preciso a formulação de um modelo matemático a partir dos sinais. Com as redes neurais, os conjuntos de dados são diretamente utilizados pelo classificador, e a modelagem matemática é representada pelos pesos das sinapses obtidos pelo treinamento.

A topologia de uma RNA se refere à organização dos neurônios e os tipos de conexões permitidas entre eles, e redes neurais multicamadas são uma importante classe nesse aspecto. Nela, uma rede irá consistir de um conjunto de nós de entrada que constituem a camada de entrada, uma ou mais camadas escondidas intermediá-

rias, e uma camada de saída. Com isso, o sinal de entrada deve se propagar através da rede por cada camada, em sequência, até chegar a saída. A Figura A.1 mostra a estrutura de uma rede multicamada padrão.

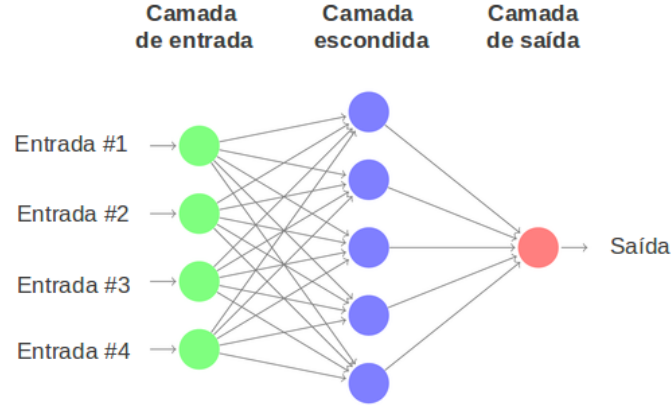


Figura A.1: Exemplo básico de rede neural multicamada, com camada de entrada, escondida e de saída.

Grande parte do processamento é realizado pelas camadas escondidas, através de conexões ponderadas, as quais são utilizadas como extratoras de características. Considerando-se que uma rede neural apresenta uma camada escondida de J neurônios, um neurônio j recebe as entradas multiplicadas pelos pesos sinápticos, os quais são ainda somados a uma constante, chamada de *bias*. Seja $x = (x_1, x_2, \dots, x_p)^T$ o vetor de entrada da rede, onde $x_i = (x_{i1}, x_{i2}, \dots, x_{in})^T$ para $i = 1, \dots, p$, e n corresponde ao tamanho do conjunto de dados, o valor dessa operação é dado por

$$a_j = \omega_{0j} + \sum_{i=1}^p \omega_{ij}x_i \quad (\text{A.1})$$

onde ω_{ij} é o peso sináptico da conexão entre o j -ésimo neurônio da camada escondida e sua entrada x_i , e ω_{0j} seu *bias*. O valor resultante é ainda passado pela função de ativação, a qual possui a finalidade de ativar ou inibir o próximo neurônio. Com isso, a saída v_j do neurônio j da camada escondida é dada por

$$v_j = \phi_j(a_j) = \phi_j\left(\omega_{0j} + \sum_{i=1}^p \omega_{ij}x_i\right) \quad (\text{A.2})$$

onde $\phi_j(\cdot)$ é a função de ativação. O resultado da função de ativação é passado para

os neurônios da próxima camada, repetindo o processo, até alcançar a camada de saída.

Esse modelo tem sido aplicado com sucesso para resolver problemas difíceis e diversos, como reconhecimento de padrões [98], controle e identificação de sistemas [99], processamento de sinais [100] e aproximação de funções [101], sendo treinado de maneira supervisionada através do algoritmo chamado de *error back-propagation* [102]. Esse algoritmo se baseia na regra de aprendizado por correção através do erro. Basicamente, o aprendizado consiste em aplicar os dados a camada de entrada da rede, e observar sua propagação através de todas as camadas. Ao final, um conjunto de valores de saída é produzido como resposta do sistema, mantendo durante todo esse procedimento os pesos das sinapses constantes. Em seguida, os valores de saída são subtraídos da resposta a qual era desejada, produzindo-se assim, um valor de erro. O erro encontrado é então propagado de volta pela rede, da camada de saída até a entrada, contra o sentido natural das sinapses, para realizar ajustes em seus pesos, de forma a tornar a resposta original mais próxima a desejada. Esse procedimento é repetido até se chegar a um valor de erro satisfatório.

As redes multicamadas apresentam características específicas. Por exemplo, cada neurônio apresenta, em seu modelo, uma função de ativação não-linear, a qual é suave, ou seja, diferenciável em todos os pontos. Caso não houvesse a não-linearidade, o modelo poderia ser resumido em apenas uma camada de neurônios, ou seja, capaz de resolver apenas problemas linearmente separáveis. Além disso, a rede contém uma ou mais camadas escondidas de neurônios, as quais são responsáveis por extrair características de ordem elevada.

Apêndice B

Mapas Auto-organizáveis (SOM)

Os Mapas Auto-organizáveis (SOM, da sigla em inglês para *Self-Organizing Maps*) são um tipo de rede neural onde seu treinamento é feito de forma não-supervisionada, ou seja, sem a necessidade de saber, a priori, qual a classificação correta dos dados. Através dessa técnica, consegue-se produzir, em uma representação discreta de baixa dimensionalidade, um mapeamento não-linear do espaço de entrada [103]. Com isso, é possível ter uma representação gráfica dos dados, auxiliando em sua interpretação. Por essas características, e por utilizar uma função de vizinhança que preserva as propriedades topológicas do espaço de entrada, essa técnica é muito aplicada para mineração de dados e análise de grupos de risco.

O objetivo do aprendizado do SOM visa fazer com que diferentes setores do mapa respondam de forma semelhante a certos padrões de entrada. A formação do mapa ocorre através de três processos [104]: competição, cooperação e adaptação. Na competição, para cada vetor de entrada, haverá apenas um neurônio vencedor. Já na cooperação, os neurônios adjacentes irão ser excitados, conforme uma função de vizinhança. Por fim, a adaptação reforça a resposta do neurônio vencedor, e de seus vizinhos, ao padrão de entrada, após o ajuste dos pesos sinápticos.

Ao considerar os vetores de entrada como $x = [x_1, x_2, \dots, x_n]^T$, e os pesos sinápticos dos neurônios definidos como $w_i = [w_{1i}, w_{2i}, \dots, w_{ni}]^T$, já que os neurônios são todos conectados com as entradas, a atualização do vetor de pesos do neurônio vencedor é feita sequencialmente pela equação:

$$w_i(t + 1) = w_i(t) + \eta(t)h_{ij}(t)(x(t) - w_i(t)) \quad (\text{B.1})$$

onde $\eta(t)$ é uma taxa de aprendizagem monotonicamente decrescente e $h_{ij}(t)$ é a função de vizinhança, a qual deve ter seu valor máximo se afastando do neurônio vencedor com área de atuação decrescente com o tempo. A Figura B.1 mostra um exemplo de mapa, o qual é formado por neurônios conectados entre si em uma visualização bi-dimensional.

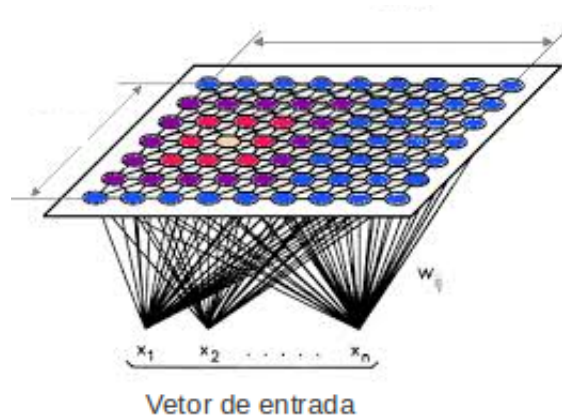


Figura B.1: Diagrama de um mapa auto-organizável bi-dimensional, onde o vetor de entrada x é aplicado ao mapa com pesos sinápticos w . O neurônio vencedor é representado de cor clara, e pode-se observar a área de vizinhança, afetada pela sua ativação, possuindo nesse caso uma profundidade de dois neurônios.

O SOM pode ser interpretado de duas formas. Na primeira, considera-se que os pesos sinápticos dos neurônios, enquanto ponteiros para o espaço de entrada, formam uma aproximação da distribuição dos eventos utilizados no treinamento. Uma maior quantidade de neurônios irá apontar para regiões com alta concentração dos dados e menos para as outras regiões. Alternativamente, pode-se interpretar com o sentido de que, durante o treinamento, os neurônios de certa vizinhança se movem para uma certa direção, pois dados similares tendem a ativar neurônios adjacentes. Com isso, um mapa semântico é formado, que mapeia eventos semelhantes próximos um dos outros, e caso contrário, distantes entre si. Esse mapeamento pode ser visto através da U-Matrix do SOM [105].

A representação pela U-Matrix visa mostrar graficamente a distância entre um neurônio e seus vizinhos. Com isso, aplica-se a mesma métrica utilizada durante o

treinamento para calcular distâncias entre pesos de neurônios vizinhos. O resultado é uma matriz, a qual pode ser interpretada como uma imagem, na qual as coordenadas de cada pixel são derivadas das coordenadas dos neurônios no mapa, sendo a intensidade de cada pixel correspondente a distância calculada.

Através desse fator, é possível identificar regiões onde os neurônios são mais próximos, o que pode significar uma formação de agrupamento, ou distantes, o que representaria uma dissimilaridade ou uma fronteira entre agrupamentos [106]. Portanto, a U-Matrix se torna uma ferramenta que pode ser de grande ajuda na procura pela formação de agrupamentos, tanto matematicamente como visualmente. A Figura B.2 mostra, à esquerda, um exemplo do cálculo da distância da U-Matrix, e à direita, um exemplo de U-Matrix, onde é possível identificar um agrupamento na parte inferior do mapa e outro na parte superior, os quais são separados por uma região de fronteira horizontal, representada pela região mais escura do mapa.

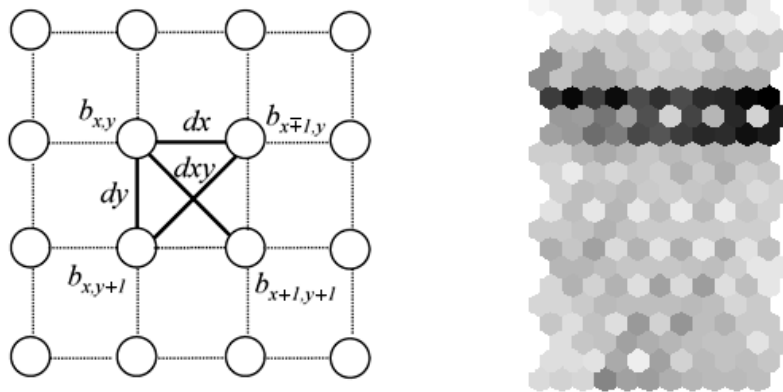


Figura B.2: Cada unidade da U-Matrix representa a distância entre um neurônio e seus vizinhos, cálculo representado à esquerda. Um exemplo de diagrama de um mapa auto-organizável é mostrado à direita, o qual mostra uma região de afastamento horizontal representada em cor mais escura.

Apesar de útil, a inspeção visual só pode ser utilizada para uma análise qualitativa. Para descrições quantitativas dos dados, os grupos de interesse devem ser identificados no mapa. Para isso, o mapeamento dos dados realizado pelo SOM, após o treinamento, pode ser aplicado a outros métodos de clusterização, onde seus neurônios irão ser definidos em grupos distintos. Dessa forma, os eventos passam a ser os neurônios do mapa, reduzindo o custo computacional no caso de base de

dados volumosas, e reduz o impacto de eventos atípicos que poderiam impactar na clusterização [107].