



**Universidade do Estado do Rio de Janeiro**

Centro de Tecnologia e Ciências

Faculdade de Engenharia

Noemi da Paixão Pinto

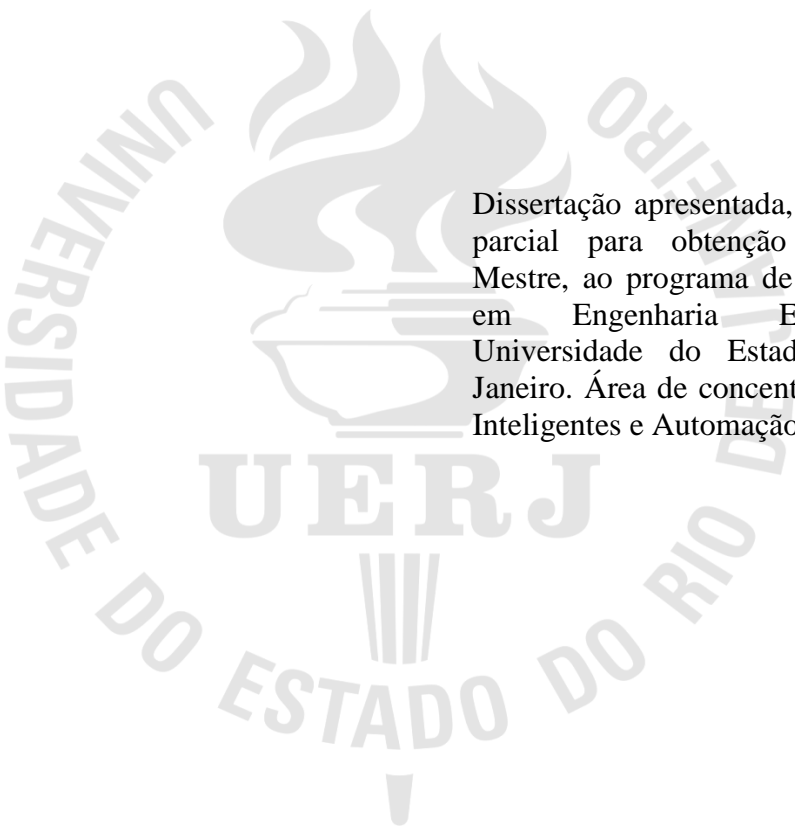
**Deteccção de Alterações Respiratórias na Fibrose Cística Através da Técnica  
de Oscilações Forçadas e Algoritmos de Aprendizado de Máquinas**

Rio de Janeiro

2018

Noemi da Paixão Pinto

**Deteccão de Alterações Respiratórias na Fibrose Cística Através da Técnica de Oscilações Forçadas e Algoritmos de Aprendizado de Máquinas**



Dissertação apresentada, como requisito parcial para obtenção do título de Mestre, ao programa de Pós-Graduação em Engenharia Eletrônica da Universidade do Estado do Rio de Janeiro. Área de concentração: Sistemas Inteligentes e Automação.

Orientadores: Prof. Dr. Jorge Luís Machado do Amaral

Prof. Dr. Pedro Lopes de Melo

Rio de Janeiro

2018

CATALOGAÇÃO NA FONTE  
UERJ / REDE SIRIUS / BIBLIOTECA CTC/B

P659 Pinto, Noemi da Paixão.  
Detecção de alterações respiratórias na fibrose cística através da técnica de oscilações forçadas e algoritmos de aprendizado de máquinas / Noemi da Paixão Pinto. – 2018.  
114f.  
  
Orientadores: Jorge Luís Machado do Amaral, Pedro Lopes de Melo.  
Dissertação (Mestrado) – Universidade do Estado do Rio de Janeiro, Faculdade de Engenharia.  
  
1. Engenharia eletrônica - Teses. 2. Aprendizado do computador - Teses. 3. Teoria bayesiana de decisão estatística - Teses. 4. Algoritmos genéticos - Teses. I. Amaral, Jorge Luís Machado do. II. Melo, Pedro Lopes de. III. Universidade do Estado do Rio de Janeiro, Faculdade de Engenharia. IV. Título.

CDU 004.891

Bibliotecária: Júlia Vieira – CRB7/6022

Autorizo, apenas para fins acadêmicos e científicos, a reprodução total ou parcial desta tese, desde que citada a fonte.

---

Assinatura

---

Data

Noemi da Paixão Pinto

**Deteção de Alterações Respiratórias na Fibrose Cística Através da Técnica de Oscilações Forçadas e Algoritmos de Aprendizado de Máquinas**

Dissertação apresentada, como requisito parcial para obtenção do título de Mestre, ao programa de Pós-Graduação em Engenharia Eletrônica da Universidade do Estado do Rio de Janeiro. Área de concentração: Sistemas Inteligentes e Automação.

Aprovado em 15 de maio de 2018.

Banca Examinadora:

---

Prof. Dr. Jorge Luís Machado do Amaral (Orientador)  
Faculdade de Engenharia – UERJ

---

Prof. Dr. Pedro Lopes de Melo (Orientador)  
Instituto de Biologia – UERJ

---

Prof. PhD. Ana Cristina Bicharra Garcia  
Centro de Ciências Exatas e Tecnologia – UNIRIO

---

Prof. Dr. Nayat Sánchez Pi  
Instituto de Matemática e Estatística – UERJ

Rio de Janeiro

2018

## **DEDICATÓRIA**

Dedico este trabalho ao Laboratório de Redes Industriais e Sistemas de Automação (LARISA) e ao Laboratório de Instrumentação Biomédica da UERJ (LIB-UERJ) pelos esforços em se manterem ativos, mesmo em meio à crise vivida no Estado do Rio de Janeiro, e por abrir espaço para o desenvolvimento de métodos para o diagnóstico e estudo de doenças que afetam o sistema respiratório, como a fibrose cística.

## **AGRADECIMENTOS**

Agradeço a Deus por se mostrar presente em minha vida, sendo meu refúgio e fortaleza, e pelo privilégio em continuar os estudos através do mestrado. Aos meus amigos e pais, Augusto e Rose, pelo incentivo demonstrado em todas as etapas de minha vida, sempre acompanhados de muito bom humor e palavras de ânimo. À minha amiga e irmã, Laís, que me influenciou e incentivou a seguir nessa área. Por todas as palavras de ânimo que, também acompanhadas de muito bom humor, sempre trouxeram alívio nos mais diversos momentos. Ao meu amigo e companheiro de todas as horas, Sávio, que me incentivou a entrar no mestrado, me ajudou a concluir essa etapa e sempre me deu palavras de ânimo para continuar. À minha amiga, Patrícia, que conheci através do mestrado, e aos amigos Adriano e Alexandre que sempre me ajudaram por meio de explicações, palavras de ânimo e incentivo durante essa caminhada. Aos companheiros de laboratório Hugo, Everton, George e Anderson que me receberam muito bem e me fizeram sentir parte do LARISA. Pelas conversas que fizeram toda a diferença em meio às tarefas. Ao professor Pedro Lopes pela orientação, dedicação e oportunidade de continuar desenvolvendo um trabalho com aplicação na área de biomédica. Ao professor Jorge Amaral pela oportunidade, orientação, dedicação, incentivo e muita paciência demonstrados ao longo deste projeto. À FAPERJ e ao CNPq pelo apoio financeiro desse projeto.

E formou o Senhor Deus o homem do pó da terra, e soprou em suas narinas o fôlego da vida;  
e o homem foi feito alma vivente.

*Genesis 2:7 ACF*

## RESUMO

PINTO, Noemi P. *Detecção de alterações respiratórias na fibrose cística através da técnica de oscilações forçadas e algoritmos de aprendizado de máquinas*. 114f. 2018. Dissertação (Mestrado em Engenharia Eletrônica) – Faculdade de Engenharia, Universidade do Estado do Rio de Janeiro, Rio de Janeiro, 2018.

Quando começou a ser estudada, a fibrose cística levava recém-nascidos a óbito em seu primeiro ano de vida. Entretanto, devido a avanços no tratamento, esses pacientes têm chegado até a fase adulta. Exames como teste de suor e espirometria, vêm sendo usados na tentativa de detectar a doença em sua fase inicial, porém esses métodos não têm sido eficientes. Sendo assim, um novo método vem sendo estudado para avaliar as propriedades mecânicas do sistema respiratório: a técnica de oscilações forçadas (FOT). A fim de comprovar a eficácia dessa nova técnica, este trabalho propõe o uso de algoritmos de aprendizado de máquinas para auxiliar a investigação e diagnóstico de alterações respiratórias na fibrose cística. Os dados fornecidos pela FOT foram aplicados nos algoritmos: *K Nearest Neighbor* (K-NN), *Radial Support Vector Machine* (RSVM), *Adaboost* (ADAB) e *Random Forest* (RF). Com o objetivo de manter uma boa acurácia e aumentar a interpretabilidade dos resultados obtidos, esses dados também foram submetidos a um algoritmo de Redes Bayesianas sintetizadas com algoritmo genético (RBGAOT). Dos experimentos realizados, a reatância respiratória fornecida pela FOT, foi o atributo que apresentou melhor desempenho individual (AUC=0,85). No experimento com oito atributos o algoritmo RBGAOT apresentou melhor desempenho (AUC=0,88). Com a aplicação dos métodos produto cruzado e seleção de variáveis, o K-NN e ADAB foram os algoritmos que tiveram melhores resultados (AUC=0,89). Os experimentos realizados mostraram que o uso de algoritmos de aprendizado de máquina aumentou a acurácia no diagnóstico de alterações respiratórias da fibrose cística. Já a inferência sobre as redes construídas pelo RBGAOT gerou um aumento na interpretabilidade das relações existentes entre as variáveis fornecidas pela FOT.

Palavras-chave: Fibrose cística; Técnica de oscilações forçadas; FOT; Aprendizado de máquina; Redes Bayesianas; Algoritmo genético; AUC.



## ABSTRACT

PINTO, Noemi P. *Detection of respiratory changes in cystic fibrosis by forced oscillation technique and machine learning algorithms*. 114f. 2018. Dissertação (Mestrado em Engenharia Eletrônica) – Faculdade de Engenharia, Universidade do Estado do Rio de Janeiro, Rio de Janeiro, 2018.

When the cystic fibrosis studies began, it used to lead newborns to death after their first year of life. However, due to advances in treatment of cystic fibrosis, these patients have reached adulthood. Medical exams such as sweat test and spirometry, have been used as an attempt to diagnose the disease on its first stage, but these methods have not been efficient. Therefore, a new method is being studied to evaluate the mechanical properties of the respiratory system: the Forced Oscillation Technique (FOT). To prove the efficiency of this new technique, the present work proposes the use of machine learning algorithms to help the investigation and diagnosis of respiratory changes in cystic fibrosis. The data provided by FOT were used on the following algorithms: K Nearest Neighbor (K-NN), Radial Support Vector Machine (RSVM), Adaboost (ADAB) and Random Forest (RF). With the purpose of keeping a good accuracy and increase the interpretability of the results, this data was submitted to Bayesian Network synthesized by genetic algorithm (RBGAOT). From the experiments performed, the respiratory reactance provided by the FOT was the feature selection that presented the best individual performance (AUC=0.85). On the experiment with eight features, the RBGAOT had the best performance (AUC=0.88). When the methods of cross product and feature selection were applied, the K-NN and ADAB were the algorithms with the best results (AUC=0.89). The experiments realized showed that the use of machine learning algorithms increased the accuracy on the diagnosis of respiratory changes in cystic fibrosis. The inference about the networks constructed by RBGAOT generated an increase in the interpretability of the existing relation between the variables provided by the FOT.

Keywords: Cystic fibrosis; Forced oscillation technique; FOT; Machine learning; Bayesian Networks; Genetic algorithm; AUC.

## LISTA DE ILUSTRAÇÃO

Figura 1 – Fluxograma de recomendações para o diagnóstico da fibrose cística.....	21
Figura 2 – Diagrama em blocos básico do sistema .....	22
Figura 3 – Indivíduo realizando ensaios pela técnica de oscilações forçadas .....	26
Figura 4 – Exemplo da configuração 1-NN.....	28
Figura 5 – Exemplo das configurações: (a) 3-NN e (b) 5-NN .....	29
Figura 6 – Exemplo de fronteira de decisão e os vetores de suporte do algoritmo SVM .....	33
Figura 7 – Exemplo de hiperplanos na classificação SVM.....	35
Figura 8 – Exemplo de hiperplanos na classificação SVM com margens suaves.....	36
Figura 9 – Conjunto de dados: (a) em um espaço unidimensional e não linearmente separável e (b) em um novo espaço bidimensional e linearmente separável .....	37
Figura 10 – Elementos de representação das Redes Bayesianas .....	38
Figura 11 – Topologia de uma Rede Bayesiana simples.....	39
Figura 12 – Tipos de inferências bayesianas: (a) Causal; (b) Diagnóstico; (c) Intercausal; ....	41
Figura 13 – Rede Bayesiana construída para o problema de câncer de pulmão .....	41
Figura 14 – Tabelas de distribuição de probabilidade conjunta para o problema de câncer no pulmão .....	42
Figura 15 – Estrutura de Rede Bayesiana selecionada para o problema <i>Contratar</i> .....	46
Figura 16 – Tabelas de distribuição de probabilidade conjunta do exemplo <i>Contratar</i> .....	47
Figura 17 – Fluxograma básico de um algoritmo genético .....	49
Figura 18 – Fluxograma resumido do modelo proposto.....	51
Figura 19 – Divisão para validação cruzada com 10 pastas .....	54
Figura 20 – Matriz confusão das possíveis classificações de uma instância.....	55
Figura 21 – Representação de uma Rede Bayesiana em matriz esparsa .....	58
Figura 22 – Comparação dos parâmetros da FOT de indivíduos do grupo controle e do grupo teste.....	68
Figura 23 – Curvas ROC dos parâmetros da FOT .....	70
Figura 24 – Curvas ROC do experimento com todos os parâmetros da FOT .....	72
Figura 25 – Análise da sensibilidade com especificidade em 75% e 90% no experimento com oito atributos.....	73
Figura 26 – Curvas ROC do experimento com oito atributos cruzados.....	75

Figura 27 – Análise da sensibilidade com valores de especificidade em 75% e 90% no experimento com oito atributos cruzados.....	76
Figura 28 – Curvas ROC do experimento com seleção de atributos da FOT .....	77
Figura 29 – Análise da sensibilidade com valores de especificidade em 75% e 90% no experimento com seleção de atributos da FOT .....	78
Figura 30 – Curvas ROC do experimento com cinco parâmetros da FOT cruzados .....	80
Figura 31 – Análise da sensibilidade com valores de especificidade em 75% e 90% no experimento com cinco parâmetros da FOT cruzados .....	81
Figura 32 – Curvas ROC do experimento com atributos do produto cruzado selecionados....	82
Figura 33 – Análise da sensibilidade com valores de especificidade em 75% e 90% no experimento com atributos do produto cruzado selecionado .....	83
Figura 34 – Resumo dos maiores valores de AUC obtidos durante os experimentos.....	84
Figura 35 - Resumo dos maiores valores de sensibilidade com especificidade fixada em 75%, obtidos durante os experimentos .....	85
Figura 36 – Resumo dos maiores valores de sensibilidade com especificidade fixada em 90%, obtidos durante os experimentos .....	85
Figura 38 – Estrutura da rede com oito atributos de entrada.....	87
Figura 39 – Estrutura da rede com cinco atributos de entrada .....	93
Figura 40 – Estrutura da rede 1 gerada com cinco atributos de entrada.....	109
Figura 41 – Estrutura da rede 2 gerada com cinco atributos de entrada.....	112

## LISTA DE TABELAS

Tabela 1 – Incidência de fibrose cística em diferentes regiões .....	19
Tabela 2 – Parâmetros fornecidos pela FOT .....	25
Tabela 3 – Algoritmo básico da configuração 1-NN.....	28
Tabela 4 – Algoritmo básico do <i>Random Forest</i> .....	31
Tabela 5 – Algoritmo básico do Adaboost .....	32
Tabela 6 – Variáveis e seus possíveis estados no problema do câncer de pulmão.....	42
Tabela 7 – Variáveis e possíveis estados do problema <i>Contratar</i> .....	46
Tabela 8 – Resultados do treinamento do algoritmo K-NN .....	56
Tabela 9 – Resultados do treinamento do algoritmo ADAB.....	57
Tabela 10 – Resultados do treinamento do algoritmo RF .....	57
Tabela 11 – Exemplo de tabela de DPC .....	59
Tabela 12 – Parâmetros fornecidos pela FOT .....	66
Tabela 13 – Pontos de corte para discretização dos parâmetros da FOT, média e desvio padrão .....	69
Tabela 14 – Comportamento geral das características do grupo controle e do grupo teste .....	69
Tabela 15 – Desempenho individual dos parâmetros da FOT na classificação de pacientes... ..	70
Tabela 16 – Resultado dos oito parâmetros da FOT submetidos aos classificadores .....	71
Tabela 17 – Comparação dos valores de AUC dos classificadores no experimento com todos os atributos da FOT .....	72
Tabela 18 – Resultado do experimento com oito atributos cruzados submetidos aos classificadores.....	74
Tabela 19 – Comparação dos valores da AUC dos classificadores no experimento com oito atributos cruzados .....	75
Tabela 20 – Resultado do experimento com a seleção de cinco atributos submetidos aos classificadores.....	77
Tabela 21 – Comparação dos valores de AUC dos classificadores no experimento com seleção de atributos da FOT .....	78
Tabela 22 – Resultado do experimento com produto cruzado dos atributos selecionados da FOT submetidos aos classificadores.....	79
Tabela 23 – Comparação dos valores de AUC dos classificadores no experimento com cinco parâmetros da FOT cruzados .....	80

Tabela 24 – Resultado do experimento com atributos do produto cruzado selecionados e submetidos aos classificadores .....	82
Tabela 25 – Comparação dos valores de AUC dos classificadores no experimento com atributos do produto cruzado selecionados.....	83
Tabela 26 – Probabilidades à priori da variável <i>classe</i> com oito atributos de entrada.....	87
Tabela 27 – DPC para a variável $Z_{4Hz}$ da rede com oito atributos de entrada.....	88
Tabela 28 – DPC para a variável $R_m$ da rede com oito atributos de entrada.....	88
Tabela 29 – DPC para a variável $E_{din}$ da rede com oito atributos de entrada.....	89
Tabela 30 – DPC para a variável $X_m$ da rede com oito atributos de entrada.....	90
Tabela 31 – DPC para a variável $R_o$ da rede com oito atributos de entrada.....	90
Tabela 32 – DPC para a variável $C_{din}$ da rede com oito atributos de entrada.....	91
Tabela 33 – DPC para a variável $S$ da rede com oito atributos de entrada.....	92
Tabela 34 – Probabilidades à priori da variável <i>classe</i> da rede com cinco atributos de entrada .....	94
Tabela 35 – DPC para a variável $R_o$ com cinco atributos de entrada .....	94
Tabela 36 – DPC para a variável $C_{din}$ da rede com cinco atributos de entrada .....	95
Tabela 37 – DPC para a variável $R_m$ da rede com cinco atributos de entrada .....	96
Tabela 38 – DPC para a variável $X_m$ da rede com cinco atributos de entrada .....	96
Tabela 39 – DPC para a variável $Z_{4Hz}$ da rede com cinco atributos de entrada .....	97
Tabela 40 – Probabilidades à priori da variável <i>classe</i> da rede 1 com cinco atributos de entrada .....	109
Tabela 41 – Probabilidades à priori da variável $R_o$ da rede 1 com cinco atributos de entrada .....	110
Tabela 42 – DPC da variável $C_{din}$ da rede 1 com cinco atributos de entrada .....	110
Tabela 43 – DPC da variável $X_m$ da rede 1 com cinco atributos de entrada .....	111
Tabela 44 – DPC da variável $Z_{4Hz}$ da rede 1 gerada com cinco atributos de entrada.....	111
Tabela 45 – DPC da variável $R_o$ da rede 1 gerada com cinco atributos de entrada .....	111
Tabela 46 – Probabilidades à priori da variável <i>classe</i> da rede 2 com cinco atributos de entrada .....	112
Tabela 47 – DPC para a variável $R_o$ da rede 2 gerada com cinco atributos de entrada.....	113
Tabela 48 – DPC para a variável $C_{din}$ da rede 2 com cinco atributos de entrada .....	113
Tabela 49 – DPC para a variável $X_m$ da rede 2 com cinco atributos de entrada .....	113
Tabela 50 – DPC para a variável $Z_{4Hz}$ da rede 2 com cinco atributos de entrada .....	114
Tabela 51 – DPC para a variável $R_m$ da rede 2 com cinco atributos de entrada.....	114

## LISTA DE ABREVIATURAS E SIGLAS

ADAB	Adaboost
AE	Algoritmos evolucionários
AG	Algoritmo genético
ANOVA	Análise de Variância
AUC	<i>Area Under the ROC curve</i>
BDeu	<i>Bayesian Dirichlet Equivalent Uniform</i>
CFTR	<i>Cystic Fibrosis Transmembrane Conductance Regulator</i>
DAG	<i>Directed Acyclic Graphs</i>
DPC	Distribuições de probabilidade conjunta
DPN	Diferença de potencial nasal
DPOC	Doença respiratória obstrutiva crônica
FOT	<i>Forced Oscillation Technique</i>
GAOT	<i>Genetic Algorithms for optimization</i>
K-NN	<i>K Nearest Neighbor</i>
LIB	Laboratório de Instrumentação Biomédica da UERJ
MPF	Melhor parâmetro da FOT
PGM	<i>Probabilistic Graphical Model</i>
PNT	Pneumotacômetro
RBGAOT	Redes Bayesianas sintetizadas por algoritmos genéticos
ROC	<i>Receiver Operating Characteristic</i>
RF	<i>Random Forest</i>
RSVM	<i>Radial Support Vector Machine</i>
LSVM	<i>Linear Support Vector Machine</i>
TP	Transdutor

## SUMÁRIO

<b>INTRODUÇÃO</b> .....	15
<b>1. FIBROSE CÍSTICA</b> .....	19
<b>2. TÉCNICA DE OSCILAÇÕES FORÇADAS</b> .....	22
<b>3. ALGORITMOS DE APRENDIZADO DE MÁQUINAS</b> .....	27
3.1. <i>K-Nearest Neighbor</i> .....	27
3.2. <i>Random Forests</i> .....	30
3.3. Adaboost .....	31
3.4. <i>Support Vector Machines</i> .....	33
3.5. Redes Bayesianas .....	38
3.5.1. Distribuição de Probabilidade Conjunta .....	40
3.5.2. Tipos de Inferência Bayesiana .....	40
3.5.3. Aprendizagem e Construção de uma Rede .....	45
3.5.4. Exemplo de Aplicação .....	45
3.5.5. Vantagens e Desvantagens das Redes Bayesianas .....	48
3.6. Algoritmos Genéticos .....	48
<b>4. MODELO PROPOSTO</b> .....	50
4.1. Dados de Entrada .....	51
4.2. Seleção de Atributos .....	52
4.3. Treinamento do Modelo .....	53
4.4. Medida de desempenho .....	54
4.5. Classificadores .....	56
4.6. Redes Bayesianas sintetizadas com Algoritmos Genéticos .....	58
4.6.1. Discretização dos Dados .....	60
4.6.2. RBGAOT .....	62
4.6.3. Representação do cromossomo .....	62
4.6.4. População Inicial .....	63
4.6.5. Função de Avaliação .....	63
4.6.6. Função de Seleção .....	64
4.6.7. Operadores Genéticos .....	64
<b>5. ESTUDO DE CASO</b> .....	66
5.1. Descrição do Conjunto de Dados .....	66

5.2. Experimento Individual dos Atributos .....	69
5.3. Experimento com Oito Atributos .....	71
5.4. Experimento com Oito Atributos Cruzados .....	74
5.5. Experimento com Cinco Atributos Seleccionados .....	76
5.6. Experimento com Cinco Atributos Cruzados.....	79
5.7. Experimento com Seleção de Cinco Atributos do Produto Cruzados.....	81
5.8. Inferência sobre Redes Bayesianas .....	86
5.8.1. Rede com Oito Atributos .....	86
5.8.2. Rede com Seleção de Cinco Atributos.....	93
<b>CONCLUSÃO</b> .....	99
<b>REFERÊNCIAS</b> .....	102
<b>APÊNDICE A – COMBINAÇÕES DO PRODUTO CRUZADO</b> .....	108
<b>APÊNDICE B – INFERÊNCIA SOBRE ESTRUTURAS DE REDES BAYESIANAS</b> .....	109
1. Inferência sobre a Rede 1 com cinco atributos de entrada .....	109
2. Inferência sobre a Rede 2 com cinco atributos de entrada .....	112



## INTRODUÇÃO

A fibrose cística é uma doença genética que inicialmente era diagnosticada em recém-nascidos. Essas crianças eram levadas a óbito ainda no primeiro ano de vida, apresentando problemas, principalmente, no sistema respiratório. Porém, nos últimos anos houve avanço no tratamento e diagnóstico da doença, fazendo com que esses pacientes chegassem à idade adulta (LIMA et al., 2010). A espirometria é um dos principais métodos usados atualmente para o diagnóstico da fibrose cística, porém por ser um exame mais simples, não caracteriza em detalhes o sistema respiratório e não permite um melhor entendimento dos processos da doença. Sendo assim, a busca por novas técnicas tem sido uma grande motivação na pesquisa para aprimorar a identificação dessa doença.

Dentre os novos métodos pesquisados, a técnica de oscilações forçadas (FOT - *Forced Oscillation Technique*) tem sido estudada para avaliar as propriedades mecânicas do sistema respiratório (LIMA et al., 2015). A utilização dos parâmetros obtidos pela FOT, associada aos métodos de aprendizado de máquinas, trouxe importantes avanços no diagnóstico de doenças respiratórias (AMARAL et al., 2013; AMARAL et al., 2015; AMARAL et al., 2017).

Diversas publicações têm mostrado que é possível aplicar os algoritmos de aprendizado de máquinas no diagnóstico e estudo de doenças respiratórias. O artigo (AMARAL et al., 2013) descreve o uso de classificadores para aumentar a acurácia na identificação de mudanças no sistema respiratório de pacientes com tabagismo. Usando como medida de desempenho a área sob a curva ROC (AUC – *Area Under the Receiver Operating Characteristic Curve*), os algoritmos usados foram: classificadores logísticos lineares, *K*-NN, redes neurais e SVM. Dentre os testes apresentados, os melhores desempenhos foram obtidos pelo *K*-NN e SVM com valores de AUC iguais a 0,91. Esses resultados caracterizam alta taxa de acerto na classificação e comprovam que o uso de algoritmos de aprendizado de máquinas aumentou a acurácia na identificação de alterações no sistema respiratório geradas pelo tabagismo.

Já o artigo (AMARAL et al., 2015), propõe técnicas para classificar automaticamente os níveis de obstrução das vias aéreas de portadores de doença pulmonar obstrutiva crônica (DPOC). Os algoritmos *K*-NN, RF, LSVM e RSVM, foram usados durante os experimentos e avaliados de acordo com a AUC. Os classificadores *K*-NN e RF apresentaram melhor desempenho, com valores de AUC maiores que 0,9 na maioria dos procedimentos realizados. Esses resultados comprovam que o uso de algoritmos de aprendizado de máquinas pode

ajudar na categorização da obstrução das vias aéreas da DPOC e auxiliar os médicos na análise da progressão da doença.

O artigo (AMARAL et al., 2017), propõe o desenvolvimento de classificadores automáticos para simplificar o uso clínico e aumentar a precisão da FOT no diagnóstico de obstrução das vias aéreas em pacientes portadores de asma. Os algoritmos K-NN, RF, AdaBoost (ADAB) e classificador no espaço de dissimilaridade (FDSC – *Feature-based Dissimilarity Space Classifier*) foram usados e avaliados durante os experimentos, de acordo com a AUC. Os classificadores ADAB e K-NN apresentaram melhor desempenho com valores de AUC variando entre 0,88 e 0,91 durante os testes realizados. De acordo com os resultados observados, os classificadores usados podem ajudar no diagnóstico da obstrução das vias aéreas em pacientes asmáticos, auxiliando os médicos na identificação da obstrução das vias aéreas.

Embora o uso de aprendizado de máquina em associação com os parâmetros da FOT apresente elevado potencial no diagnóstico de alterações respiratórias na fibrose cística, essa associação ainda não foi investigada. Desta maneira, esse trabalho se insere na linha de pesquisa dos trabalhos citados, propondo o uso de algoritmos de aprendizado de máquinas para aprimorar ainda mais o diagnóstico e aplicação da FOT em doenças respiratórias. Dessa forma, novas informações podem auxiliar a equipe médica na investigação e diagnóstico de alterações respiratórias em portadores de fibrose cística, através dos parâmetros fornecidos pela FOT.

Atualmente, esses parâmetros são usados separadamente para realizar o diagnóstico de doenças respiratórias, sendo o atributo que apresentar maior número de acertos, o critério selecionado para identificação da doença. Nos trabalhos dessa linha de pesquisa citados anteriormente, o conjunto de dados usado para treinar os algoritmos de aprendizado de máquinas era composto por diversas amostras obtidas pela FOT, apresentando melhor desempenho do que o método atual. Apesar desses artigos comprovarem uma boa acurácia no diagnóstico de doenças respiratórias, ainda são necessários estudos que aumentem a interpretação dos resultados fornecidos por esses métodos.

Geralmente, a acurácia é a medida de desempenho usada em algoritmos de aprendizado de máquinas. Porém, há casos onde a interpretação do processo de classificação é mais, ou tão importante, quanto à previsão feita pelo modelo. Quando isso ocorre, o algoritmo escolhido precisa realizar o estudo do conjunto de dados disponível, gerando informações novas sobre o problema e facilitando o entendimento do usuário final (BRATKO, 1997). Essa característica de expressar o comportamento de um sistema de forma compreensível é

chamada de interpretabilidade e está relacionada a fatores ligados a estrutura do modelo, porém não possui uma medida padrão para ser avaliada (GACTO et al., 2011).

Mesmo que seja uma técnica de simples execução, a análise do sistema respiratório pela FOT é de difícil compreensão. Por isso é necessário treinamento e experiência da equipe médica para interpretar as curvas de resistência, reatância e os diversos parâmetros provenientes de outras medidas obtidas por essa técnica (AMARAL et al., 2013). Logo, optar por um algoritmo de aprendizado de máquinas que forneça interpretabilidade do resultado, pode agregar mais informações para o estudo e diagnóstico da fibrose cística. Sendo assim, este trabalho também propõe o uso do algoritmo de Redes Bayesianas que fornece interpretação de seus resultados por meio de grafos, realizando a tomada de decisões a partir do raciocínio baseado em probabilidades. Com estruturas gráficas, essas redes possibilitam a representação e o raciocínio sobre um domínio incerto, tornando possível lidar com a falta de informação.

A estrutura de uma Rede Bayesiana é formada por nós, que representam as variáveis do problema e são interligadas por arcos, cuja única limitação é a exigência que os grafos formados sejam acíclicos dirigidos (DAG – *Directed Acyclic Graphs*). As ligações entre as variáveis podem ser quantificadas com o cálculo das tabelas de probabilidades condicionais (SANTANA et al., 2007). Esse algoritmo tem como vantagem o fato de conseguir lidar com grande quantidade de atributos de entrada e mesmo assim apresentar essas tabelas de forma mais compacta, já que cada variável é influenciada apenas pelas variáveis diretamente ligadas a ela.

O aprendizado das Redes Bayesianas pode ser dividido em duas etapas: o aprendizado da topologia, considerado um problema complexo, e o aprendizado das probabilidades condicionais. Ambos podem ser feitos por um especialista, mas também há possibilidade de realizá-los de forma automática por meio de outros algoritmos. Nesse caso, as estruturas são geralmente construídas com base no conjunto de dados e o modelo obtido é usado para prever novos resultados (GONÇALVES, 2017).

Dentre os algoritmos usados para a busca de uma estrutura de Rede Bayesiana, podem ser destacados o *K2* e o *B* (PIFER, 2006). O algoritmo *K2* inicia sua busca com uma estrutura simples onde todas as variáveis são consideradas independentes. A cada iteração a entropia da rede é calculada e os arcos são adicionados à medida que a entropia é minimizada. As probabilidades condicionais são obtidas diretamente do conjunto de dados (HERSKOVITS et al., 1991). Já o algoritmo *B* é inicializado como o algoritmo *K2*, porém os nós e arcos são acrescentados de acordo com a diferença entre os valores de qualidade. Esse processo é

repetido até que a qualidade não aumente mais ou até que a rede esteja completa (CASTILLO et al., 1996).

A busca realizada por esses algoritmos é considerada NP-difícil devido à grande quantidade de estruturas DAG que podem descrever as relações entre suas variáveis (LARRAÑAGA et al., 1996). Motivados por essa limitação, estratégias vem sendo estudadas para a seleção de estruturas e pesquisas têm mostrado que o uso de algoritmos evolucionários (AE), fornece resultados eficientes para essa busca (TONDA et al., 2012; MYERS et al., 1999; MURUZÁBAL et al., 2007; LARRAÑAGA et al., 1996). Em geral, os AE realizam uma busca probabilística baseada nos princípios da evolução natural das espécies. Essa técnica pode ser aplicada mesmo em casos onde há muitos atributos de entrada e um conjunto de dados limitado (TONDA et al., 2012). Os AE se dividem em três principais tipos: algoritmos genéticos, programação evolutiva e estratégias de evolução (GABRIEL et al., 2008).

Este trabalho também propõe o uso de algoritmos genéticos (AG) para estimar as estruturas de Redes Bayesianas que melhor representam as relações existentes entre as características fornecidas pela FOT. Os AG são inspirados na teoria de seleção natural das espécies, proposta por Darwin, e usada para busca e otimização de problemas complexos. Dessa forma, mais informações podem ser geradas para auxiliar a equipe médica no estudo e diagnóstico de anormalidades respiratórias na fibrose cística.

Os três primeiros capítulos a seguir são destinados à revisão teórica dos principais assuntos que se baseia este trabalho. O primeiro capítulo descreve a fibrose cística, abordando suas causas, sintomas e atuais métodos de diagnóstico. O segundo capítulo apresenta os parâmetros fornecidos pela FOT e as vantagens em fazer uso dessa técnica. O terceiro capítulo descreve de forma inicial os algoritmos de aprendizado de máquinas escolhidos para aplicar os dados obtidos na FOT. O quarto capítulo apresenta o modelo proposto para este trabalho, que conta com os algoritmos descritos no capítulo 3, além da aplicação de algoritmos genéticos para gerar estruturas de Redes Bayesianas. Já o capítulo 5, mostra os resultados dos experimentos realizados com o uso das técnicas de seleção de variáveis e produto cruzado, bem como a inferência realizada nas estruturas geradas.

## 1. FIBROSE CÍSTICA

A mucoviscidose, ou fibrose cística, é uma doença hereditária autossômica recessiva que atinge pessoas de ambos os sexos, sendo mais incidente na população caucasiana. É causada por mutações no gene localizado no braço longo do cromossomo sete e responsável por codificar uma proteína chamada CFTR (*Cystic Fibrosis Transmembrane Conductance Regulator*) (MOTA et al., 2015; CASTELLANI et al., 2008). A CFTR é um canal responsável por regular e participar do transporte de eletrólitos por meio das membranas celulares dos sistemas respiratório, digestivo e do aparelho reprodutor, sendo o sistema respiratório o mais afetado (DALCIN et al., 2008).

Nos recém-nascidos, os primeiros sintomas podem ser observados através de alterações nas pequenas vias aéreas e tosses crônicas, logo nos primeiros meses. Já pacientes menores de 18 anos, podem apresentar quadros de pneumonia recorrentes (RIBEIRO, 2002). De uma forma geral, por se tratar de uma doença progressiva, a fibrose cística causa um aumento na obstrução do fluxo de ar, bem como o aumento da frequência respiratória e da dificuldade expiratória. Esses sintomas causam incômodo durante o sono e uma diminuição na tolerância às atividades físicas, à fisioterapia e até mesmo atividades normais realizadas no cotidiano, diminuindo assim, a expectativa de vida dos pacientes.

Quando a fibrose cística começou a ser estudada, pacientes eram levados ao óbito no primeiro ano de vida. Porém, devido ao avanço no diagnóstico e tratamento da doença, esses indivíduos têm chegado até a fase adulta. Atualmente, cerca de 70.000 pacientes estão registrados em todo o mundo. As incidências da doença variam de acordo com a região, conforme a Tabela 1, considerando pacientes da população branca de recém-nascidos (vivos) (MOTA et al., 2015).

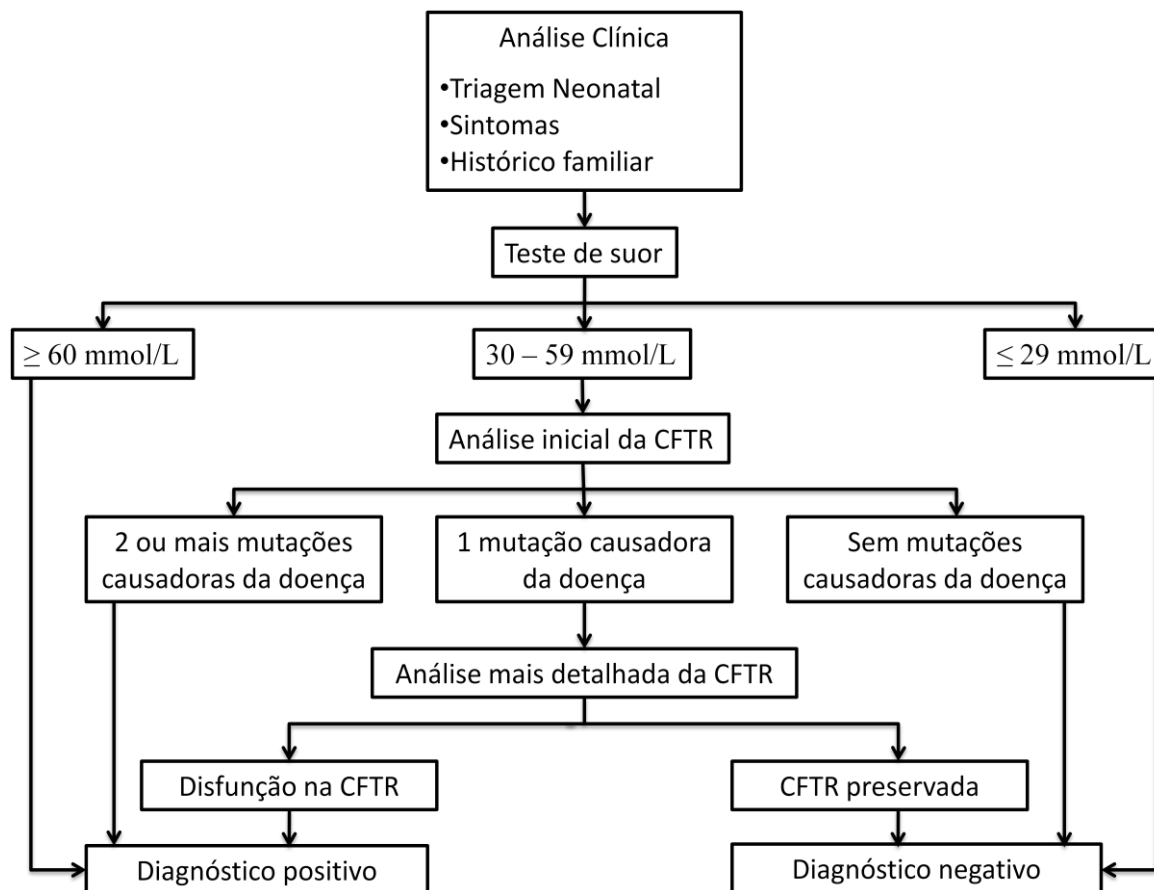
**Tabela 1 – Incidência de fibrose cística em diferentes regiões**  
(Extraído de: LIMA et al., 2010)

<b>Região ou País</b>	<b>Incidência</b>
Europa	1:2000 a 1:3000
Latino-Americanos	1:4000 a 1:10000
Africa do Sul	1:7056
Japão	1:350.000

Com o passar dos anos, os estudos sobre a doença foram avançando e tornaram possíveis análises mais profundas sobre o gene defeituoso, bem como a avaliação do estado do paciente, tornando possível também a expansão e avanços no tratamento. Desde então, a expectativa de vida dos portadores de fibrose cística tem se tornado cada vez maior.

Em 1938, o primeiro trabalho com a descrição da doença relatou que a expectativa de vida dos recém-nascidos era menor que um ano de idade (ANDERSEN, 1938). Já registros norte-americanos realizados em 2007 mostraram que 43% dos portadores de Fibrose Cística tinham mais de 18 anos e a idade média dos pacientes chegava a 36,5 anos de idade (DALCIN et al., 2008). Já em 2014, a expectativa média de vida dos pacientes era de 39,3 anos (LIMA et al., 2010; *Cystic Fibrosis Foundation Patient Registry*, 2014).

Atualmente, é recomendado que o diagnóstico da fibrose cística seja feito com base em três parâmetros. O primeiro parâmetro é a análise clínica que inclui: triagem neonatal, sintomas característicos da doença e histórico familiar. O segundo parâmetro é a concentração de cloreto de sódio obtido através do teste de suor, onde indivíduos com resultado menor que 30 mmol/L, são considerados indivíduos com poucas chances de portar a doença e indivíduos com resultados maiores que 60 mmol/L possuem grandes chances de portar a doença. Já indivíduos com concentração de cloreto na faixa de 30 a 59 mmol/L, devem passar pela análise do terceiro parâmetro: a avaliação da proteína CFTR. No caso da descoberta de duas ou mais mutações nessa proteína, há fortes indícios de se tratar de um portador da doença. Se for identificada apenas uma mutação, é feita uma análise mais profunda da disfunção da CFRT através da medição da diferença de potencial nasal (DPN). O fluxograma da Figura 1 mostra um resumo das recomendações para o diagnóstico da fibrose cística (FARRELL et al., 2017).

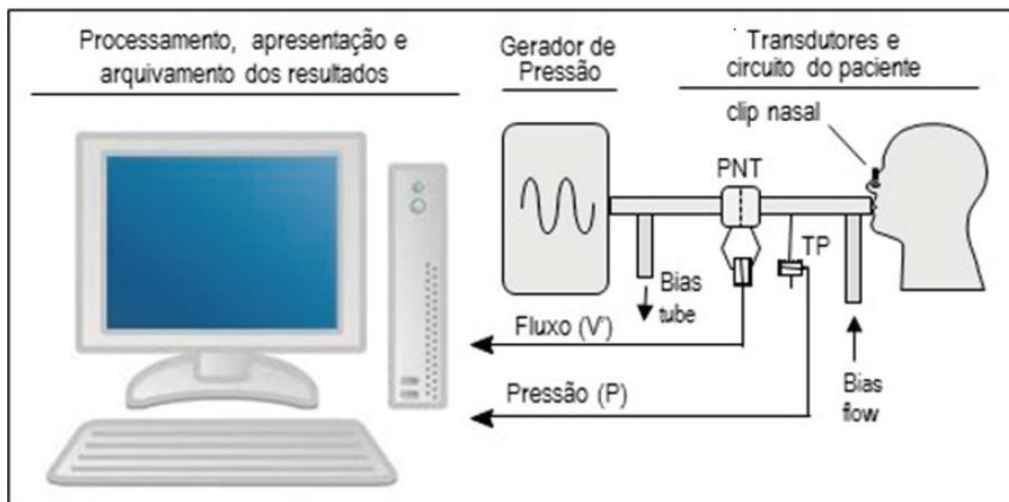


**Figura 1 – Fluxograma de recomendações para o diagnóstico da fibrose cística**  
(Adaptado de FARRELL et al., 2017)

Outra ferramenta usada para complementar a investigação e diagnóstico de alterações respiratórias na fibrose cística é a espirometria, avaliando as alterações respiratórias dos pacientes por meio de fluxos e volumes respiratórios. Dentre os resultados obtidos, especificamente o volume expiratório forçado no primeiro segundo, é o mais usado na tentativa de prever e identificar sintomas característicos em pacientes de diferentes idades. Porém, a espirometria não tem sido suficiente para diagnosticar alterações respiratórias em portadores dessa doença e uma nova técnica, que avalia as propriedades resistivas e reativas do sistema respiratório, vem sendo estudada: a técnica de oscilações forçadas (LIMA et al., 2010).

## 2. TÉCNICA DE OSCILAÇÕES FORÇADAS

Com o intuito de desenvolver uma análise mecânica do sistema respiratório de forma mais simples, mas que ainda apresentasse novas informações, no ano de 1956 foi proposto o uso da técnica de oscilações forçadas (FOT – *Forced Oscillation Technique*) (DUBOIS et al., 1956). Durante um ensaio da FOT, o paciente deve permanecer sentado, fazer o uso de um *clip* nasal e respirar de forma espontânea, enquanto um fluxo constante (*bias flow*) renova o ar inspirado pelo mesmo. Pequenas oscilações de pressão são geradas por um aparelho externo, normalmente um alto falante, e aplicadas às vias aéreas do paciente, que permanece respirando espontaneamente. Essa pressão  $P$  é medida por um transdutor (TP) e estimula um fluxo oscilatório ( $V'$ ), medido por um pneumotacômetro (PNT). O valor da resistência total do sistema respiratório, chamada impedância de entrada ( $Z_{rs}$ ), é calculado pelos sinais obtidos por TP e PNT (Figura 2) (MELO et al., 2000).



**Figura 2 – Diagrama em blocos básico do sistema**  
(MELO, 2015)



Para diminuir o tempo de execução dos ensaios da FOT, um sistema responsável por analisar impedâncias realiza o processamento dos sinais obtidos pelas oscilações nas faixas de frequências desejadas. Pela transformada de Fourier, é possível decompor os sinais  $P$  e  $V'$  e também realizar uma avaliação das alterações do módulo de  $Z_{rs}$  em diversas frequências, conforme equação (1) (LIMA et al., 2010).

$$Z_{rs}(f) = \frac{FFT(P)}{FFT(V')} \quad (1)$$

Sendo,

$FFT(P)$ : Transformada de Fourier da pressão  $P$

$FFT(V')$ : Transformada de Fourier do fluxo  $V'$

$f$ : frequência desejada

As funções senoidais provenientes da decomposição dos sinais  $P$  e  $V'$  podem ser representadas com suas respectivas componentes, conforme as equações (2) e (3):

$$P = P_m \text{sen}(\omega t) \quad (2)$$

$$V' = V'_m \text{sen}(\omega t + \varphi) \quad (3)$$

Sendo:

$P_m$ : amplitude do sinal  $P$

$V'_m$ : amplitude do sinal  $V'$

$\omega$ : frequência angular igual a  $2\pi f$

$\varphi$ : diferença de fase entre os sinais  $P$  e  $V'$

A variável  $Z_{rs}$  representa toda a carga mecânica, que inclui os efeitos das propriedades resistivas, elásticas e inertivas do sistema respiratório. Normalmente durante um ensaio da FOT, as impedâncias  $Z_{rs}$  são representadas por componentes reais e imaginários, descritos respectivamente pela resistência respiratória ( $R_{rs}$ ) e pela reatância respiratória ( $X_{rs}$ ) (MELO, 2015), conforme equação (4):

$$Z_{rs} = \sqrt{R_{rs}^2 + X_{rs}^2} \quad (4)$$

Os componentes de  $R_{rs}$  e  $X_{rs}$  podem ser derivados da seguinte forma (MACLEOD et al., 2001):

$$R_{rs} = |Zrs| \cos \varphi \quad (5)$$

$$X_{rs} = |Zrs| \sin \varphi \quad (6)$$

A energia cinética usada durante a aceleração do fluxo aéreo é descrita por meio da inertância respiratória ( $I_{rs}$ ). Essa variável é normalmente desprezada em análises realizadas em baixas frequências, sendo o sistema respiratório modelado apenas por um componente resistivo e um complacente. Já em frequências mais elevadas, como ocorre na FOT, a  $I_{rs}$  torna-se relevante, permitindo a obtenção de informações mais detalhadas sobre as características mecânicas do aparelho respiratório com base na reatância (MELO, 2015).

A faixa de frequência de 4 a 32Hz, é a mais utilizada durante os ensaios da FOT. Nesse intervalo, a resistência respiratória caracteriza a dissipação total da energia do sistema, que abrange a soma dos efeitos vindos de resistências relacionadas a quatro fatores: ao tecido pulmonar, à parede torácica, às vias aéreas e à redistribuição do fluxo do gás nos pulmões. A reatância respiratória ( $X_{rs}$ ) caracteriza o armazenamento de energia potencial do sistema que está associada à complacência respiratória ( $C_{rs}$ ), sendo o armazenamento de energia cinética associado à inertância  $I_{rs}$ . As propriedades elásticas estão associadas à complacência ( $C_{din}$ ) e à elastância dinâmica ( $E_{din}$ ), sendo o parâmetro  $E_{din}$ , o inverso de  $C_{din}$  (equação (7)) (MELO et al., 2000). A relação entre  $X_{rs}$ ,  $I_{rs}$  e  $C_{rs}$ , é descrita na equação (8):

$$E_{din} = \frac{1}{C_{din}} \quad (7)$$

$$X_{rs} = \omega I_{rs} - j \frac{1}{\omega C_{rs}} \quad (8)$$

Sendo,  $\omega=2\pi f$  e  $j=\sqrt{-1}$ .

Devido à complexidade do sistema respiratório, é comum que as componentes de  $Z_{rs}$  não estejam na mesma fase. Porém na frequência de ressonância ( $F_r$ ), os efeitos da complacência e inertância são iguais e, conseqüentemente,  $X_{rs}$  é igual a zero (MIRANDA et al., 2013).

É possível analisar a resistência e a reatância de forma mais detalhada usando diversas faixas de frequências. Mesmo sendo um método mais lento ainda é vantajoso, pois os valores obtidos representam as médias dos resultados do sistema respiratório durante vários períodos de ventilação. Como não há consenso na literatura sobre as características a serem avaliadas nesses valores médios, há grupos que realizam a análise pelo método de regressão linear para uma faixa de frequências de 4 a 16Hz. Dessa forma, é possível determinar a resistência no intercepto em 0Hz ( $R_o$ ) e o coeficiente angular da curva de resistência ( $S$ ) (LIMA et al., 2015; AMARAL et al., 2017).

O parâmetro  $R_o$  estima como as resistências newtonianas associadas às vias aéreas e aos tecidos, bem como sua resistência tardia proveniente da distribuição do gás, reagem em frequências baixas. Já o parâmetro  $S$  está associado à alteração na distribuição do fluxo de gás dentro do sistema respiratório de acordo com a frequência utilizada (MIRANDA et al., 2013). A Tabela 2 mostra um resumo de todas as características fornecidas pela FOT.

**Tabela 2 – Parâmetros fornecidos pela FOT**

<b>Parâmetro</b>	<b>Descrição do Parâmetro</b>
$R_o$	Resistência no Intercepto
$R_m$	Resistência Média
$X_m$	Reatância Média
$C_{din}$	Complacência Dinâmica
$S$	Inclinação da Curva de Resistência
$Z_{4Hz}$	Impedância em 4Hz
$F_r$	Frequência de Ressonância
$E_{din}$	Elastância Dinâmica

Em suma, a FOT possui duas principais vantagens. Primeiramente, esse método permite uma análise mais detalhada do sistema respiratório, fornecendo parâmetros que não podem ser obtidos pela espirometria e demais técnicas usadas. Sendo assim, a FOT apresenta forte potencial para diagnósticos e contribui para melhor compreensão dos processos da doença. Outra grande vantagem desse método é a fácil execução do exame para o profissional, dependendo apenas da cooperação do paciente que deve respirar de forma espontânea, conforme mostrado na Figura 3.

Por ser uma técnica nova que investiga as alterações mecânicas no sistema respiratório, são necessários os mais diversos estudos e testes para que se comprove sua eficácia. Assim sendo, diversas pesquisas têm sido feitas para mostrar que é possível aplicar a FOT para auxílio da equipe médica (RIBEIRO et al., 2018; MARINHO et al., 2017; LACERDA et al., 2017).



**Figura 3 – Indivíduo realizando ensaios pela técnica de oscilações forçadas (MELO et al., 2015)**

### 3. ALGORITMOS DE APRENDIZADO DE MÁQUINAS

Com o intuito de auxiliar a equipe médica no diagnóstico de alterações respiratórias na fibrose cística por meio da FOT, foi usada a técnica de aprendizado de máquinas. Cinco algoritmos foram escolhidos para realização dos testes: *K-Nearest Neighbor*, *Random Forest*, *Adaboost*, *Support Vector Machine* e Redes Bayesianas. Cada um desses algoritmos foi descrito neste capítulo, com suas principais vantagens e desvantagens.

#### 3.1. *K-Nearest Neighbor*

O algoritmo dos  $K$  vizinhos mais próximos ( $K$ -NN – *K Nearest Neighbor*) é considerado um dos mais simples algoritmos de aprendizado de máquinas. Esse algoritmo tem o aprendizado por instância, onde o conjunto de treinamento é armazenado durante o estágio de aprendizado. Quando uma nova instância precisa ser classificada, o algoritmo encontra as  $K$  instâncias de treinamento mais próximas, usando uma função de similaridade, que normalmente é a distância euclidiana. Uma instância  $x$  pode ser representada como um vetor de atributos  $c_r$ , conforme a equação (9) (FACELI et al., 2011):

$$c_r(x) = (c_1(x), c_2(x), c_3(x), \dots, c_n(x)) \quad (9)$$

Sendo  $c_n(x)$  o valor de cada atributo do vetor  $c_r(x)$ .

O método  $K$ -NN assume que todas as instâncias estão dentro de um espaço  $n$ -dimensional  $\mathfrak{R}^n$ . Dessa forma, o cálculo da distância euclidiana entre duas instâncias  $x_i$  e  $x_j$  é feito conforme a equação (10):

$$d(x_i, x_j) = \sqrt{\sum_{r=1}^n (c_r(x_i) - c_r(x_j))^2} \quad (10)$$

Durante o treinamento da configuração mais simples, considerando apenas 1 vizinho mais próximo (1-NN), o algoritmo aprende um conjunto de dados ( $D$ ) com seus respectivos rótulos. Para atribuir uma classe a uma amostra de teste ( $a$ ) não rotulada, é feito o cálculo da

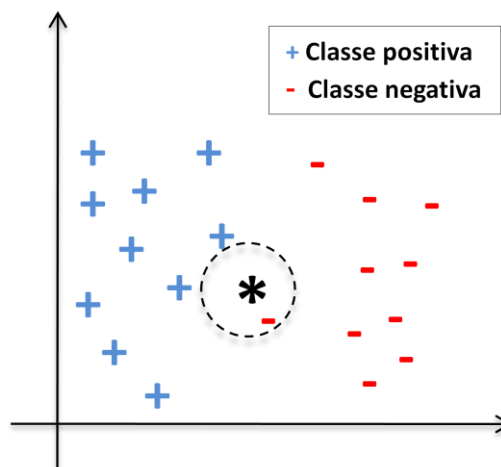
distância entre o vetor de características de  $a$  e o vetor de características de todas as instâncias do conjunto  $D$ , armazenadas pelo algoritmo. O vetor de característica com menor distância classifica a amostra de teste. A Tabela 3 mostra um algoritmo básico da configuração 1-NN (SMOLA et al., 2008). A Figura 4 mostra um exemplo da configuração 1-NN aplicada a um problema de classificação, onde uma amostra pode receber o rótulo positivo ou negativo. O asterisco representa a amostra a ser classificada e, nesse exemplo, devido à proximidade com um ponto da classe negativa, a amostra em questão é classificada como negativa.

**Tabela 3 – Algoritmo básico da configuração 1-NN**  
(Adaptado de: FACELI et al., 2011)

---

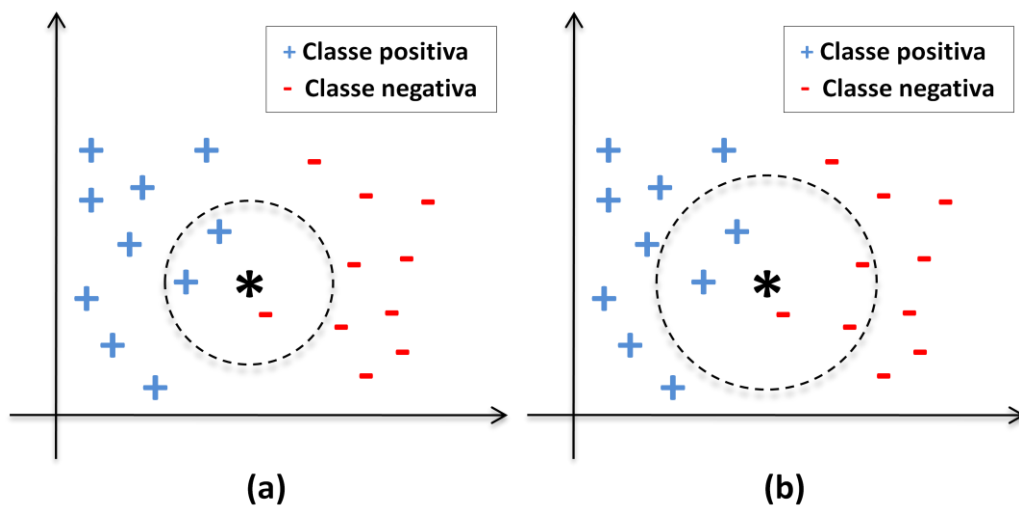
Conjunto de dados para treinamento: $D = \{(x_i, y_i)\}, i = 1 \dots n$
Uma amostra que se deseja classificar: $a = \{x_a, y_a = ?\}$
Distância entre as instâncias: $d(x_i, x_j)$
Resultado da classificação da amostra a: $y_a$
$d_{\min} \leftarrow +\infty$
Para $i=1:n$ faça
Se $d(x_i, x_j) < d_{\min}$
$d_{\min} \leftarrow d(x_i, x_j)$
$\text{ind} \leftarrow i$
Fim
Fim
$y_a = y_{\text{ind}}$

---



**Figura 4 – Exemplo da configuração 1-NN**  
(Adaptado de: FACELI et al., 2011)

Na Figura 5, um exemplo das configurações 3-NN e 5-NN mostram que dependendo do número de vizinhos considerados, a mesma amostra poderia ser classificada com diferentes rótulos. Na configuração 3-NN (Figura 5 (a)), a amostra é classificada como positiva, pois está próxima a dois pontos dessa classe, mas apenas a 1 ponto da classe negativa. Já na configuração 5-NN (Figura 5 (b)), a classificação seria negativa devido à proximidade da amostra com três pontos dessa classe. A única restrição para a escolha do  $K$  é que seja um valor ímpar, para realizar a classificação de uma amostra sem que haja empate entre as classes.



**Figura 5 – Exemplo das configurações: (a) 3-NN e (b) 5-NN**  
(Adaptado de: FACELI et al., 2011)

Dentre as vantagens do algoritmo  $K$ -NN, pode-se destacar seu treinamento simples que se resume ao armazenamento dos dados de consulta em sua memória. Outra vantagem é sua aplicabilidade em problemas mais complexos, já que realiza aproximações para cada nova instância a ser classificada. Como desvantagem, pode ser citado o custo computacional em grandes conjuntos de treinamento, pois é necessário calcular a distância entre a amostra a ser classificada e cada um desses pontos. Outra desvantagem é a sensibilidade quanto à quantidade de atributos usados no problema, já que o número de atributos define a dimensão do espaço  $\mathcal{R}^n$ .

### 3.2. *Random Forests*

As florestas aleatórias (*Random Forests*) são comitês de árvores de decisão e se baseiam em dois conceitos principais (BREIMAN, 2001). O primeiro é a seleção aleatória dos atributos de entrada para a formação de diversos subconjuntos, que serão submetidos às árvores de decisões. Esse processo contribui para a redução da correlação entre as diversas árvores construídas (COSTA, 2012). O segundo conceito é o *bagging* (*Bootstrap Aggregation*), usado para criação desses subconjuntos através da amostragem por *bootstrap*, onde o mesmo número de amostras do conjunto original é selecionado com repetição para cada novo subconjunto (LIAW et al., 2002). Dessa forma, podem existir tanto amostras repetidas, quanto amostras não inclusas durante o treinamento. O resultado das diversas árvores criadas é combinado, reduzindo então, a variância do resultado final fornecido pelo algoritmo.

O funcionamento do algoritmo *Random Forest* (Tabela 4) tem início com a seleção aleatória de um subconjunto  $Z$ , formado por amostras dos dados de treinamento com o total de  $p$  atributos. Em seguida, uma árvore  $T_b$  é construída em três etapas: seleção aleatória de  $m$  dos  $p$  atributos ( $m \ll p$ ), escolha do melhor ponto de corte dentre os atributos selecionados e divisão de um nó em dois nós filhos, com base nesse ponto de corte. Esse procedimento é repetido para cada novo nó até alcançar o tamanho mínimo de nós ( $n_{min}$ ). Com as árvores de decisão construídas, é possível configurar o algoritmo para regressão ou classificação, nesse caso, a previsão final dada pelo algoritmo será a previsão dada pela maioria das árvores individuais (HASTIE et al., 2008).

Como vantagem, o algoritmo *Random Forest* consegue lidar com um grande número de variáveis de entrada mantendo sua rapidez na construção das redes. Como a escolha dos atributos de entrada é aleatória, as árvores construídas são descorrelacionadas e, conseqüentemente, outra vantagem é uma diminuição na variância da combinação das árvores. As desvantagens do método estão na sensibilidade a muitos ruídos na base de dados e na dificuldade de interpretação do modelo (COSTA, 2012).



**Tabela 4 – Algoritmo básico do *Random Forest***  
(HASTIE et al., 2008)

---

Para  $b = 1, \dots, B$  faça

Seleciona um subconjunto  $Z$  com dados de treinamento

Constrói uma árvore de decisão seguindo as três etapas:

1. Seleciona  $m$  atributos
2. Define o melhor atributo dentre  $m$  para ponto de corte
3. Divide o nó em dois nós filhos

Fim

Saída: Conjunto de árvores  $\{T_b\}$

- Para classificação: Sendo  $\hat{C}_b(\mathbf{x})$  a classe de um ponto  $\mathbf{x}$  a ser classificado, tem-se:
 
$$\hat{C}_b(\mathbf{x}) = \text{maioria dos votos } \{\hat{C}_b(\mathbf{x})\}_1^B$$
- Para regressão:  $\hat{f}_{RF}^B(\mathbf{x}) = \frac{1}{B} \sum_{b=1}^B T_b(\mathbf{x})$

---

### 3.3. Adaboost

*Adaptive Boosting* (Adaboost) é uma técnica de aprendizado de máquina cujo objetivo é criar um classificador forte (alta acurácia), através da combinação de vários classificadores simples (baixa acurácia). Esses classificadores são treinados em sequência e a cada novo modelo os ajustes são feitos aumentando a probabilidade dos pontos classificados de forma errada, aparecerem no próximo conjunto de treinamento (MARGINEANTU et al., 1997).

Essa probabilidade é calculada da seguinte forma: em um conjunto de treinamento  $D = \{x_i, y_i\}$ ,  $x_i$  representa o vetor com os dados de entrada no sistema e  $y_i$  representa seus respectivos rótulos  $\in \{-1, +1\}$ . A cada iteração ( $t$ ), é calculada uma distribuição ( $D_t$ ):

$$D_{t+1}(i) = \frac{D_t(i) \exp(-\alpha_t y_i h_t(x_i))}{\sum_{i=1}^n D_t(i)} \quad (11)$$

Com a distribuição calculada, um algoritmo simples é aplicado para encontrar uma hipótese ( $h_t$ ). Pelo erro  $\varepsilon_t$ , o peso do classificador simples ( $\alpha_t$ ) pode ser calculado através da equação (12), onde as hipóteses com menor valor de erro recebem maior peso  $\alpha_t$ .

$$\alpha_t = \frac{1}{2} \ln \left( \frac{1 - \varepsilon_t}{\varepsilon_t} \right) \quad (12)$$

Em seguida, obtém-se o resultado final  $H(x)$  pela aplicação da função sinal na combinação ponderada de todas as hipóteses  $h_t$  (equação (13)). O algoritmo básico que descreve o funcionamento do Adaboost está na Tabela 5 (SCHAPIRE, 2013).

$$F(x) = \sum_{t=1}^T \alpha_t h_t(x) \rightarrow H(x) = \text{sign} \left( \sum_{t=1}^T \alpha_t h_t(x) \right) \quad (13)$$

**Tabela 5 – Algoritmo básico do Adaboost**  
(Adaptado de: SCHAPIRE, 2013)

---

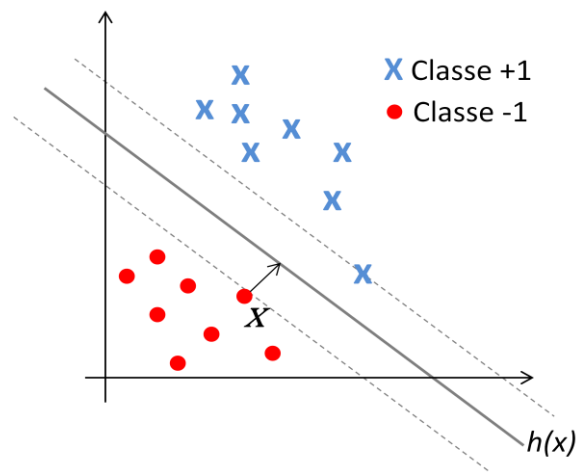
Conjunto de dados para treinamento: $D = \{(x_i, y_i)\}$
Inicializa: $D_1(i) = \frac{1}{m}$ para $i=1, \dots, m$
Para $t = 1, \dots, T$ faça
Treina o classificador através de $D_t$
Calcula as hipóteses $h_t$
Seleciona $h_t$ com menor erro $\varepsilon_t$
Escolha do $\alpha_t$
Atualiza o valor de $D_t$ para $i=1, \dots, T$
Fim
Saída $H(x)$

---

O Adaboost possui como vantagem sua simples implementação, já que o algoritmo utiliza classificadores simples e estes, sucessivamente, vão se especializando em acertar a classificação que os classificadores anteriores fizeram de forma errônea. A única restrição a ser feita é que os classificadores simples devem ter um desempenho superior aos classificadores aleatórios. Outra vantagem é sua boa generalização, sendo adequada para qualquer problema de classificação. Dentre as desvantagens do Adaboost estão o risco de *overfitting* durante o treinamento e sua sensibilidade ao lidar com dados ruidosos.

### 3.4. Support Vector Machines

O algoritmo *Support Vector Machines* (SVM) é baseado na teoria de aprendizagem estatística no qual, são aplicados princípios matemáticos para auxiliar a seleção de um classificador específico ( $h$ ), por meio de seu desempenho e complexidade, a partir de um conjunto de treinamento ( $D$ ) (FACELI et al., 2011). O SVM tem como ideia principal a criação de um hiperplano como uma fronteira de classificação, onde a margem, definida como a distância entre um ponto  $x$  e essa fronteira, é maximizada (Figura 6). Os limiares dessa fronteira são conhecidos como vetores de suporte (KUNCHEVA, 2014).



**Figura 6 – Exemplo de fronteira de decisão e os vetores de suporte do algoritmo SVM**  
(Adaptado de: KUNCHEVA, 2014)

Quanto maior a margem selecionada, melhor será a capacidade de generalização do SVM (KUNCHEVA, 2014). Considerando um conjunto de treinamento linearmente separável  $D = \{x_i, y_i\}$ , onde  $x_i$  são as entradas e  $y_i$  os rótulos das classes representados por dois valores possíveis: -1 ou 1, a fronteira de classificação é dada por um hiperplano representado pela equação (14):

$$h(x) = w \cdot x + b \quad (14)$$

Onde:

$w$ : vetor normal ao hiperplano

$b$ : número escalar

$w \cdot x$ : produto escalar

A equação (14) divide o espaço dos dados de entrada em duas regiões:

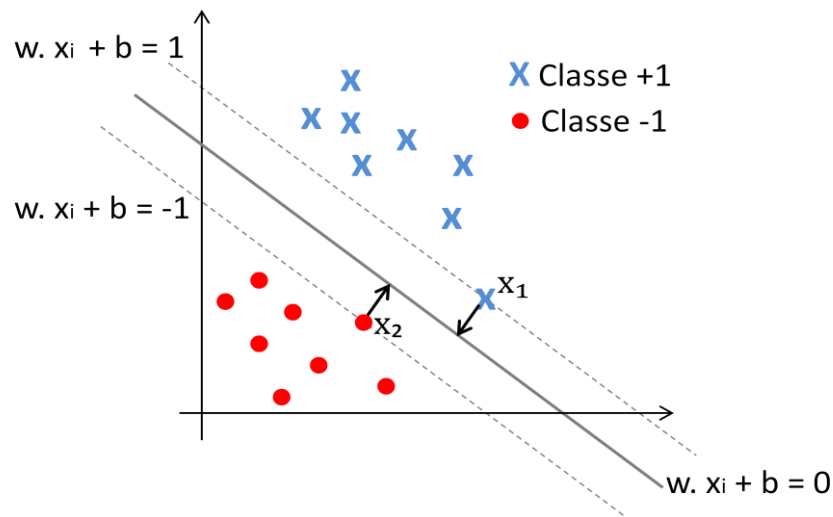
$$w \cdot x + b > 0 \quad e \quad w \cdot x + b < 0$$

Para classificar um ponto  $x$ , é necessário aplicar uma função sinal ( $sgn$ ) em  $h(x)$ :

$$g(x) = sgn(h(x)) = \begin{cases} +1, & w \cdot x + b > 0 \\ -1, & w \cdot x + b < 0 \end{cases} \quad (15)$$

Ao multiplicar o vetor  $w$  e a constante  $b$  na equação (14) por uma mesma variável, é possível obter diversos hiperplanos correspondentes. Sendo assim,  $w$  e  $b$  geram hiperplanos onde:  $h(x) = w \cdot x + b \geq 0$ , quando  $y_i = +1$  e  $h(x) = w \cdot x + b < 0$ , quando  $y_i = -1$ . Dessa forma, podem-se escrever as equações no sistema (16) e observar sua ilustração da Figura 7 (SMOLA et al., 2008):

$$\begin{cases} w \cdot x_i + b \geq +1 & \text{se } y_i = +1 \\ w \cdot x_i + b \leq -1 & \text{se } y_i = -1 \end{cases} \quad (16)$$



**Figura 7 – Exemplo de hiperplanos na classificação SVM**  
(Adaptado de: FACELI et al., 2011)

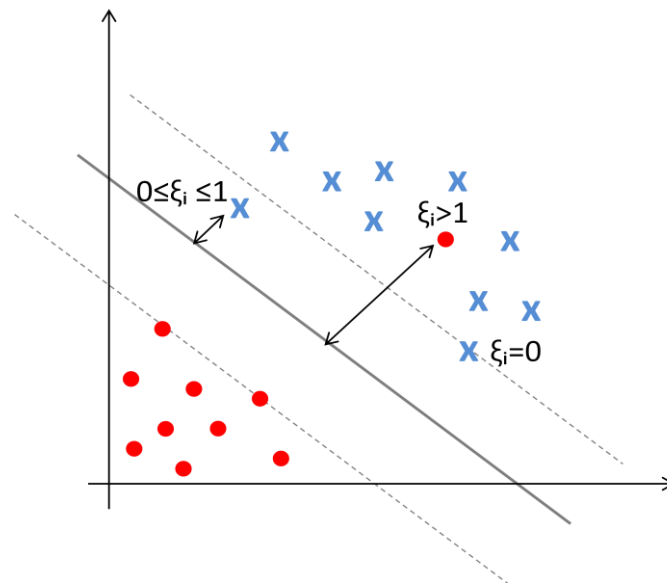
Para realizar o cálculo da margem na Figura 7, é necessário subtrair a equação do hiperplano de  $x_2$  da equação do hiperplano de  $x_1$ :

$$\begin{cases} w \cdot x_1 + b = +1 \\ w \cdot x_2 + b = -1 \end{cases} \rightarrow w(x_1 - x_2) = 2 \rightarrow \|x_1 - x_2\| = \frac{2}{\|w\|} \quad (17)$$

Para maximizar a margem devem-se minimizar os pesos. A fim de facilitar esse cálculo foram feitas alterações matemáticas para que um problema de maximização fosse reduzido a um problema de minimização de uma função quadrática, que representa a função de custo:

$$\min \|w\| \rightarrow \min \frac{1}{2} \|w\|^2 \quad (18)$$

O algoritmo SVM linear também é chamado de SVM com margens rígidas, pois são estabelecidas restrições para certificar que pontos do conjunto de treinamento não estejam entre as margens que separam as classes do problema. Porém, para aplicações em problemas reais, dificilmente são encontrados dados linearmente separáveis. Nesses casos, é possível realizar a adição de variáveis de folga ( $\xi_i$ ) que suavizam as restrições lineares, permitindo que alguns pontos de treinamento estejam entre os hiperplanos (Figura 8) e também permitem erros na classificação, como ruídos.

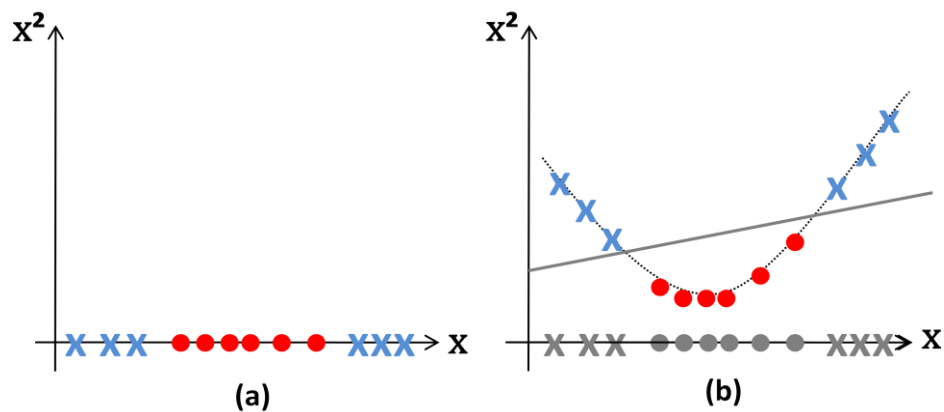


**Figura 8 – Exemplo de hiperplanos na classificação SVM com margens suaves**  
(Adaptado de: KUNCHEVA, 2014)

Dessa forma, há diferentes situações que podem ser observadas na classificação de uma amostra. Quando  $\xi_i > 1$ , a amostra está fora da região de separação e do lado incorreto de sua classificação. Quando  $0 < \xi_i \leq 1$ , a amostra está classificada corretamente, mas entre as margens de separação. Quando  $\xi_i = 0$ , a amostra está sobre as margens de separação (FACELLI et al., 2011). As equações do sistema (16) podem ser reescritas conforme o sistema (19). Por flexibilizar restrições, essa configuração é chamada de SVM com margens suaves.

$$\begin{cases} w \cdot x_i + b \geq +1 - \xi_i \text{ se } y_i = +1 \\ w \cdot x_i + b \leq -1 + \xi_i \text{ se } y_i = -1 \end{cases} \quad (19)$$

Existem problemas cujos dados não podem ser divididos por um hiperplano. Nesses casos, realiza-se o mapeamento do conjunto de treinamento para um espaço de dimensão maior, chamado espaço de características. Nesse novo espaço, espera-se que esses dados sejam linearmente separáveis (Figura 9). Após o mapeamento, é usada a configuração do SVM linear com margens suaves para lidar com classificações erradas ou ruídos (HASTIE et al., 2008).



**Figura 9 – Conjunto de dados: (a) em um espaço unidimensional e não linearmente separável e (b) em um novo espaço bidimensional e linearmente separável**  
(Adaptado de: KUNCHEVA, 2014)

Esse mapeamento em dimensões mais altas é feito através de funções de *Kernels*. O uso dessas funções na representação de espaços de características é realizado quando não se tem o conhecimento do mapeamento a ser feito. Há três principais tipos de funções *Kernel*: Polinomial, Gaussiana e *Radial Basis Function* (RBF) (FACELI et al., 2011; KUNCHEVA, 2014).

Uma das vantagens do algoritmo SVM é a sua eficiência em encontrar a melhor solução possível, já que sua função objetivo é convexa e possui apenas um mínimo global. Outra vantagem é poder aplicar o SVM em problemas que possuem elevado número de atributos, bem como em problemas de regressão. Como desvantagens, o SVM não permite uma interpretação da decisão tomada pelo algoritmo e também se mostra sensível quanto à escolha de seus parâmetros. O conjunto de treinamento do SVM precisa conter apenas dados numéricos, sendo necessário converter atributos discretos.

### 3.5. Redes Bayesianas

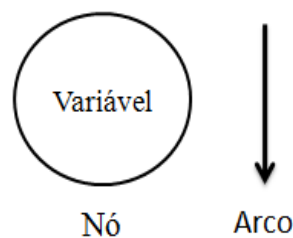
O nome destas redes é derivado da regra de Bayes, estabelecida por Thomas Bayes, que mostra como um efeito  $E$  transforma a probabilidade à priori  $P(C_j)$  em uma probabilidade à posteriori  $P(C_j/E)$  (equação (20)), alterando assim, a estimativa inicial  $C_j$  com base na nova informação fornecida por  $E$ .

$$P(C_j|E) = \frac{P(E|C_j)P(C_j)}{P(E)} \quad (20)$$

O termo  $P(E/C_j)$  é a probabilidade condicional do efeito  $E$  ser observado, dado a causa  $C_j$ .  $P(E)$  é um fator de normalização dado por um somatório, como mostrado na equação (21), e pode ser desprezado (SILVA et al., 2016). Os cálculos dessas probabilidades representam as relações causais entre as variáveis do problema, permitindo um aumento na interpretabilidade (SANTANA et al., 2007).

$$P(E) = \sum_{j=1}^n P(E|C_j)P(C_j) \quad (21)$$

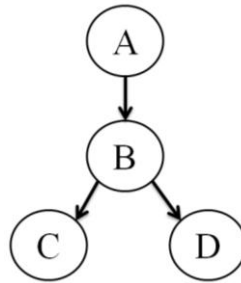
A representação em uma Rede Bayesiana mostra de forma simplificada as relações de causalidade entre as variáveis de um sistema. Essa representação é feita por nós, que correspondem às variáveis do problema, e por arcos que correspondem às conexões entre essas variáveis (Figura 10), mostrando a dependência direta entre elas (MARQUES et al., 2003).



**Figura 10 – Elementos de representação das Redes Bayesianas**



Na nomenclatura usada para os elementos das Redes Bayesianas, alguns termos são comuns para indicar a hierarquia dos nós da rede. Os termos pai e filho mostram a dependência direta entre dois ou mais nós por meio de um arco. O nó de onde parte o arco é chamado nó pai. O nó em que o arco chega, é chamado nó filho. Na rede da Figura 11, *A* é dito pai do nó *B*. Por sua vez, o nó *B* é dito filho do nó *A*. Da mesma forma, o nó *B* é chamado pai de *C* e *D*, e ambos são filhos do mesmo nó *B*. Os nós que não estão diretamente ligados pelos arcos podem ser chamados de nós antecedentes ou nós descendentes. Usando também como exemplo a rede da Figura 11, *A* é dito antecedente de *C* e *D*, bem como *C* e *D* são descendentes de *A* (ARA-SOUZA, 2010).



**Figura 11 – Topologia de uma Rede Bayesiana simples**

Outros termos usados são: nó raiz, nó folha e nó intermediário. Os nós raízes representam a origem do problema e não possuem pais. Os nós folhas mostram o resultado final do problema e não possuem filhos. Os nós que não são raízes e nem folha, são chamados de nós intermediários.

Há uma propriedade relativa a nós pais, filhos e descendentes, chamada propriedade de Markov, que diz: “não existem dependências diretas no sistema que está sendo modelado, que não sejam explicitamente mostradas nos arcos”, ou seja, uma variável é condicionalmente independente de todos os outros nós da rede que não sejam seus antecedentes (NEAPOLITAN, 2003). Atender essa propriedade é importante no uso das Redes Bayesianas, pois elas simplificam o cálculo das relações existentes entre as variáveis do problema, considerando apenas os nós que exercem influência sobre seus descendentes (KORB et al., 2011).

### 3.5.1. Distribuição de Probabilidade Conjunta

As relações entre os nós podem ser quantificadas através do cálculo da probabilidade condicional. É necessário olhar para cada um dos nós, ou variáveis do sistema, e analisar todas as possíveis combinações de valores dos nós pais em relação aos seus nós filhos. Vale ressaltar que a soma das probabilidades deve somar 1 para cada um dos possíveis estados de uma variável. O cálculo dessas probabilidades dá origem à tabela de distribuição de probabilidade conjunta (MARQUES et al., 2003).

Se uma Rede Bayesiana satisfaz a propriedade de Markov, em que cada nó depende apenas dos seus nós pais, então o cálculo da distribuição de probabilidade pode ser escrito como na equação (22):

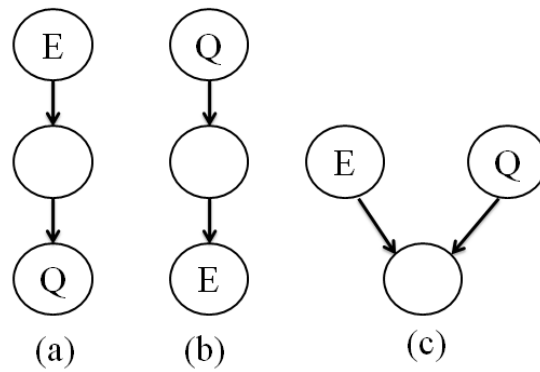
$$P(x_1, x_2, \dots, x_n) = \prod_i P(x_i | \text{Pais}(X_i)) \quad (22)$$

Uma tabela de distribuição de probabilidade conjunta representa a descrição completa do domínio de um sistema. Sendo assim, até os nós raízes possuem uma tabela que represente suas probabilidades à priori.

### 3.5.2. Tipos de Inferência Bayesiana

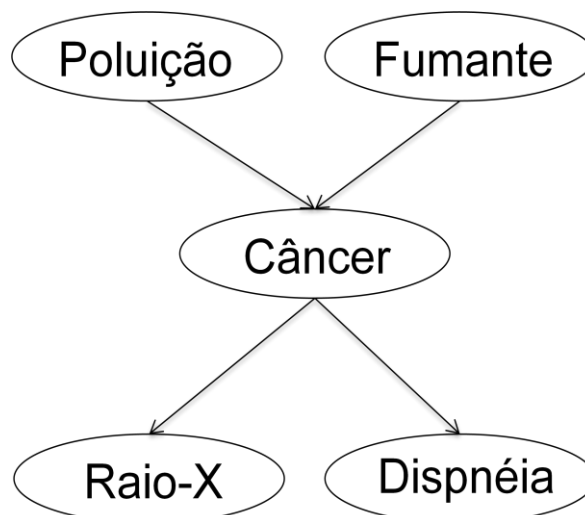
Para o raciocínio em problemas através de métodos probabilísticos, é necessário o uso do cálculo da probabilidade à posteriori sobre uma variável consulta (*query*), dada uma evidência forte (*Hard Evidence* ou *Evidence*). Dependendo do objetivo desejado, há pelo menos três formas de inferir sobre uma Rede Bayesiana e calcular a probabilidade  $P(\text{Query}|\text{Evidence})$ .

Quando o objetivo é descobrir as causas do problema, o raciocínio é feito da causa em direção ao efeito (Figura 12 (a)), seguindo a mesma direção dos arcos da rede. Quando o objetivo é o diagnóstico, o raciocínio é feito a partir dos efeitos em direção a causa (Figura 12(b)). Esse raciocínio segue o sentido oposto ao dos arcos da estrutura da rede. Se o objetivo for descobrir as causas de um efeito em comum, deve ser usado o raciocínio intercausal, onde são analisadas as variáveis *query* ( $Q$ ) e *hard evidence* ( $E$ ) (Figura 12 (c)) (MARQUES et al., 2003).



**Figura 12 – Tipos de inferências bayesianas: (a) Causal; (b) Diagnóstico; (c) Intercausal;**  
(Adaptado de: MARQUES et al., 2003)

O exemplo a seguir mostra como realizar a inferência diagnóstica em um problema (KORB et al., 2010). Um laboratório deseja estudar as relações entre as principais causas e efeitos de câncer no pulmão. Considerando que as duas causas principais que afetam as chances de um paciente desenvolver a doença, são: exposição a altos níveis de poluição e ser fumante. Uma vez que o paciente seja diagnosticado com câncer, o exame de Raio-X, geralmente, tem resultado positivo e o paciente apresenta certa dificuldade na respiração, chamada de dispneia. Ordenando as variáveis com base na opinião de um especialista, uma possível Rede Bayesiana foi construída, conforme Figura 13.

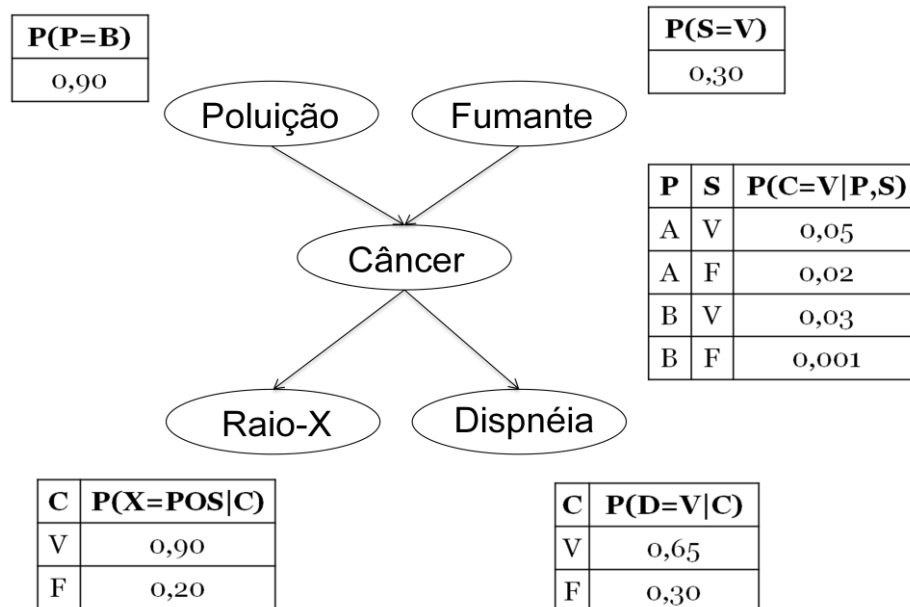


**Figura 13 – Rede Bayesiana construída para o problema de câncer de pulmão**  
(Adaptado de: KORB et al., 2010)

Cada um dos cinco nós dessa rede pode assumir dois estados (Tabela 6). Dessa forma, com a estrutura construída, e usando como base um banco de dados com registros de vários pacientes, é possível calcular as tabelas de distribuição de probabilidade conjunta para as variáveis (Figura 14).

**Tabela 6 – Variáveis e seus possíveis estados no problema do câncer de pulmão**

Variável	Possíveis estados
<i>Poluição (P)</i>	Baixo nível ( <i>B</i> ) ou Alto nível ( <i>A</i> )
<i>Fumante (F)</i>	Verdadeiro ( <i>V</i> ) ou Falso ( <i>F</i> )
<i>Câncer (C)</i>	Verdadeiro ( <i>V</i> ) ou Falso ( <i>F</i> )
<i>Raio-X (X)</i>	Verdadeiro ( <i>V</i> ) ou Falso ( <i>F</i> )
<i>Dispneia (D)</i>	Verdadeiro ( <i>V</i> ) ou Falso ( <i>F</i> )



**Figura 14 – Tabelas de distribuição de probabilidade conjunta para o problema de câncer no pulmão**

(Adaptado de: KORB et al., 2010)

A fim de obter mais informações em uma Rede Bayesiana, é possível realizar a inferência diagnóstica, analisando a estrutura no sentido contrário aos arcos. Por exemplo, para calcular a probabilidade de um indivíduo ser fumante ( $S=V$ ), dado que possui dispneia ( $D=V$ ), é feito pela aplicação da Regra de Bayes, conforme equação a seguir:

$$P(S = V|D = V) = \frac{P(D = V|S = V)P(S = V)}{P(D = V)} \quad (23)$$

Cada termo da equação (23) deve ser calculado separadamente, com os valores obtidos pelas tabelas de DPC da Figura 14. Como a variável *Câncer* está entre os nós *Dispneia* e *Fumante*, o cálculo do termo  $P(D=V/S=V)$  precisa levar em consideração a probabilidade de um indivíduo ter câncer ( $C=V$ ), ou não ( $C=F$ ):

$$\begin{aligned} P(D = V|S = V) &= P(D = V|C = V).P(C = V|S = V) + \\ &+ P(D = V|C = F).P(C = F|S = V) \end{aligned} \quad (24)$$

Porém, o nó *Câncer* também recebe influência do nó *Poluição*. Sendo assim, o cálculo do termo  $P(C=V/S=V)$  também considera a probabilidade do paciente ter sido exposto a altos níveis de poluição ( $P=A$ ), ou não ( $P=B$ ):

$$\begin{aligned} P(C = V|S = V) &= P(C = V|P = A, S = V).P(P = A) \\ &+ P(C = V|P = B, S = V).P(P = B) \end{aligned} \quad (25)$$

$$P(C = V|S = V) = 0,05.0,1 + 0,03.0,9$$

$$P(C = V|S = V) = 0,032$$

O termo  $P(C=F/S=V)$  pode ser encontrado usando o resultado da equação (25):

$$P(C = F|S = V) = 1 - P(C = V|S = V) \quad (26)$$

$$P(C = F|S = V) = 1 - 0,032$$

$$P(C = F|S = V) = 0,968$$

Com esses valores é possível encontrar o resultado da equação (24):

$$P(D = V|S = V) = 0,65 \cdot 0,032 + 0,30 \cdot 0,968$$

$$P(D = V|S = V) = 0,311$$

O próximo passo é calcular a probabilidade de ter dispneia. Esse cálculo também engloba os possíveis estados da variável *Câncer*, conforme a equação a seguir:

$$P(D = V) = P(D = V|C = V) \cdot P(C = V) + P(D = V|C = F) \cdot P(C = F) \quad (27)$$

Para o cálculo da probabilidade de ter câncer, é necessário considerar também as variáveis *Poluição* e *Fumante*:

$$\begin{aligned} P(C = V) &= P(C = V|P = A, S = V) \cdot P(P = A) \cdot P(S = V) + \\ &+ P(C = V|P = A, S = F) \cdot P(P = A) \cdot P(S = F) + \\ &+ P(C = V|P = B, S = V) \cdot P(P = B) \cdot P(S = V) + \\ &+ P(C = V|P = B, S = F) \cdot P(P = B) \cdot P(S = F) \end{aligned} \quad (28)$$

$$P(C = V) = 0,0116$$

Com esses termos calculados, é possível obter o resultado da equação (27):

$$P(D = V) = P(D = V|C = V) \cdot P(C = V) + P(D = V|C = F) \cdot P(C = F)$$

$$P(D = V) = 0,65 \cdot 0,0116 + 0,30 \cdot (1 - 0,0116)$$

$$P(D = V) = 0,304$$

Com todos os termos da equação (23), a probabilidade de um indivíduo ser fumante, dado que ele possui dispneia, é:

$$P(S = V|D = V) = \frac{0,311 \cdot 0,30}{0,304} = 0,307$$

### 3.5.3. Aprendizagem e Construção de uma Rede

A aprendizagem Bayesiana visa fornecer uma estrutura que melhor represente o problema e que facilite a obtenção de informações. Esse processo pode ser dividido em duas partes. Na primeira parte ocorre a aprendizagem da topologia da rede, ordenando as variáveis do problema e suas relações de causalidade. Na segunda parte acontece a aprendizagem dos parâmetros numéricos, quando é feito o cálculo das probabilidades condicionais.

Um especialista pode analisar e definir as duas etapas da aprendizagem Bayesiana, tomando por base apenas seu conhecimento prévio. Porém, tanto a estrutura da rede, quanto as tabelas de distribuição de probabilidade, podem ser obtidas através de um conjunto de dados. Nesse caso, um algoritmo é usado para gerar a Rede Bayesiana de forma automática (GONÇALVES, 2017).

### 3.5.4. Exemplo de Aplicação

Uma Rede Bayesiana pode ser definida como a representação compacta da distribuição de probabilidade conjunta do domínio de um problema. Para um especialista, essas estruturas mostram de forma simples e gráfica as relações de causalidade das variáveis que compõe um sistema (MARQUES et al., 2003). O exemplo de aplicação de Redes Bayesianas a seguir, mostra uma das diversas possibilidades em representar o problema abordado.

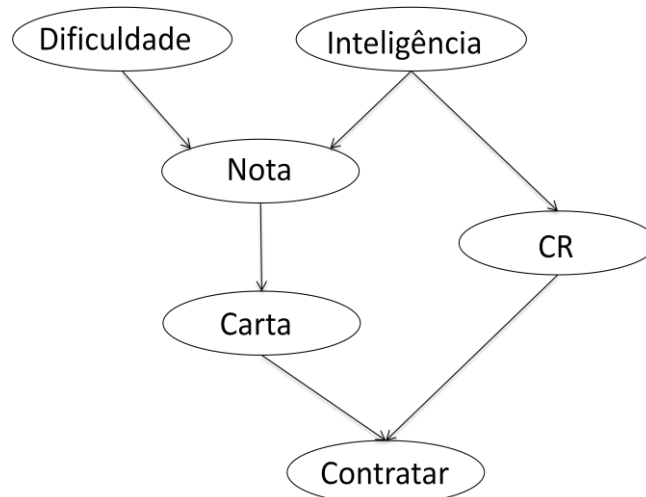
Considerando o processo seletivo de uma companhia, onde é necessário contratar funcionários com alto nível de inteligência, optou-se por selecionar alunos recém-formados. Como não há uma forma direta de testar a inteligência dos candidatos, a companhia decidiu efetuar uma análise baseando-se em três critérios: a nota em uma disciplina específica, coeficiente de rendimento e uma carta de recomendação.

Pelos critérios escolhidos, algumas observações podem ser feitas. A nota de um aluno em uma disciplina específica, e de interesse para a empresa, depende da inteligência do estudante e da dificuldade em cursar a matéria. O coeficiente de rendimento no decorrer da faculdade também depende da inteligência do aluno. No caso da carta de recomendação, é provável que o professor não se lembre do desempenho de todos os seus alunos, logo, a carta pode ser redigida com base nas notas dos estudantes. No total, seis variáveis e seus possíveis estados podem ser listados, conforme Tabela 7. Com a seleção das variáveis relevantes para o

domínio do problema, é necessário ordená-las de acordo com suas relações de causalidade. A Figura 15 mostra uma possível estrutura de Rede Bayesiana para esse exemplo.

**Tabela 7 – Variáveis e possíveis estados do problema *Contratar***

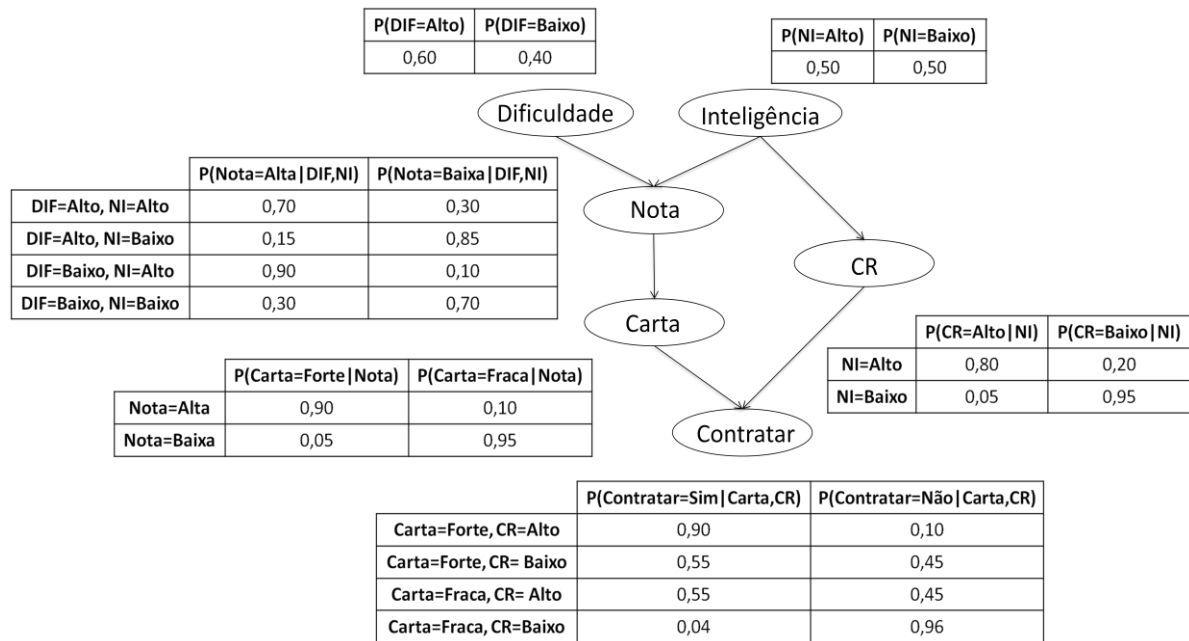
Variáveis	Possíveis Estados
Nível de inteligência do aluno (NI)	Alto ou Baixo
Nota em determinada disciplina (Nota)	Alta ou Baixa
Nível de dificuldade na disciplina (DIF)	Alto ou Baixo
Coefficiente de Rendimento (CR)	Alto ou Baixo
Carta de recomendação (Carta)	Forte ou Fraca
Contratar	Sim ou Não



**Figura 15 – Estrutura de Rede Bayesiana selecionada para o problema *Contratar***

Com a estrutura definida, é necessário montar as tabelas de distribuição de probabilidade conjunta. Nesse exemplo, os valores escolhidos simulam casos onde há um conjunto de dados para ser usado no cálculo dessas probabilidades. A Figura 16 mostra as seis tabelas para a estrutura da Rede Bayesiana gerada.





**Figura 16 – Tabelas de distribuição de probabilidade conjunta do exemplo *Contratar***

Os nós raízes *Dificuldade* e *Inteligência* possuem tabelas com probabilidades à priori, pois não há nós que exerçam influência sobre eles. A variável *CR* é influenciada apenas pela variável *Inteligência*, já que um bom desempenho acadêmico (*CR=Alto*) está ligado a um bom nível de inteligência do aluno (*NI=Alto*). A variável *Carta* recebe influência apenas da variável *Nota*, mostrando que há alta probabilidade de um aluno ter uma carta com fortes recomendações, dado que foi observada uma nota alta em determinada disciplina.

As variáveis *Nota* e *Contratar* são influenciadas por dois nós pais, apresentando assim, tabelas com mais probabilidades a serem definidas. A variável *Nota* mostra que há probabilidade de 0,90 de um aluno ter bom desempenho em uma disciplina (*Nota=Alta*), dado que a matéria não é difícil (*DIF=Baixo*) e ele possui alto nível de inteligência (*NI=Alto*).

Já a variável *Contratar* recebe influência dos nós *Carta* e *CR*. Através dos diferentes estados que os nós pais podem assumir, é possível observar a alta probabilidade existente em duas situações. Um aluno com carta de recomendação forte e alto *CR*, possui probabilidade de 0,90 de ser contratado. Já um aluno com carta de recomendação fraca e baixo *CR*, possui probabilidade de 0,96 em não ser contratado.

Dessa forma, as Redes Bayesianas permitem inferir sobre o domínio de um problema, representando graficamente e quantificando as relações entre suas variáveis. Vale ressaltar que a estrutura feita para esse exemplo é apenas uma das formas de ordenar os nós.

### **3.5.5. Vantagens e Desvantagens das Redes Bayesianas**

Sendo um método que utiliza o raciocínio probabilístico, as Redes Bayesianas permitem tomar decisões mesmo com grande quantidade de dados e informações insuficientes. Também permite expressar as relações causais entre as variáveis de forma visual e de fácil entendimento. Outra vantagem do método é a apresentação de tabelas de probabilidade condicional compactas, já que cada variável recebe influência apenas dos seus nós pais.

Como desvantagem, o algoritmo de Redes Bayesianas exige um bom conhecimento do problema para construir uma base de dados probabilística, mesmo com informações incompletas, o que pode exigir gastos. Outra desvantagem está na eliminação de dados para compactar o problema, o que pode eliminar parâmetros importantes.

### **3.6. Algoritmos Genéticos**

Inspirados na teoria de seleção natural das espécies, proposta por Darwin, os algoritmos genéticos são técnicas usadas para busca e otimização de problemas complexos. A estratégia usada por esses algoritmos é a geração de uma população inicial composta por indivíduos que representam as possíveis soluções do problema. Esses indivíduos são codificados em estruturas chamadas de cromossomos que passam pelas gerações e evoluem de acordo com o princípio de seleção e sobrevivência dos mais aptos.

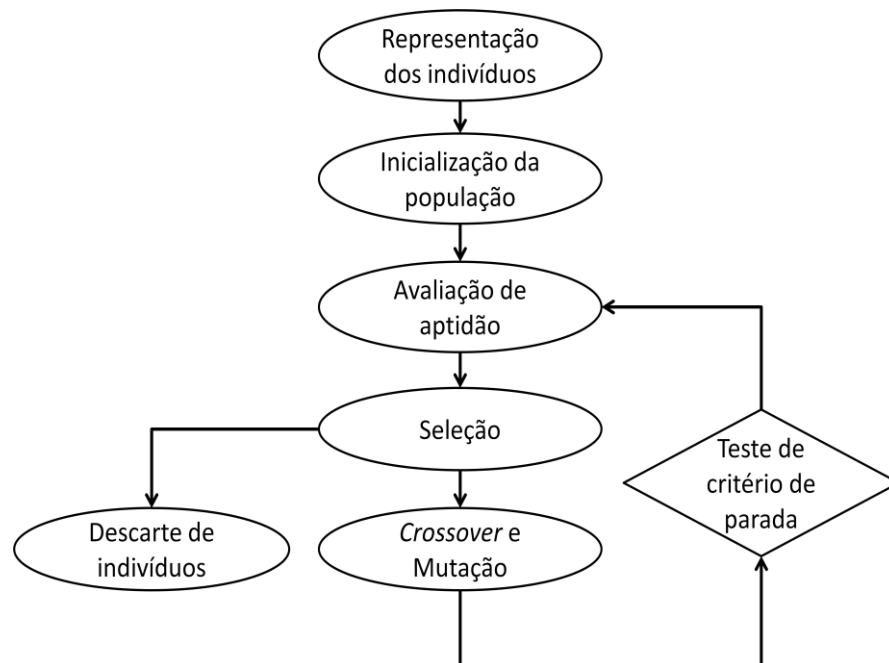
Na natureza, observa-se a competição de indivíduos por recursos básicos à sobrevivência. Os indivíduos que não tem sucesso em obter esses recursos possuem uma probabilidade menor em ter seus genes transferidos para as próximas gerações, e consequentemente, tem menos chance de deixar descendentes. Já os indivíduos que tem sucesso, possuem uma probabilidade maior em se manter nas próximas gerações, e dessa forma, produzir novos indivíduos com características mais adequadas ao seu meio ambiente. De forma análoga, a população de indivíduos representa o espaço de busca que contém possíveis soluções. As gerações são representadas pelos ciclos e o meio ambiente é o problema a ser resolvido (ROSA et al., 2009).

Todos os indivíduos da população são avaliados por uma função e recebem uma medida de aptidão, que reflete o quão boa uma solução é para o problema. Para que ocorra a geração de descendentes, um conjunto de indivíduos é selecionado com base na sua aptidão,

que serão, posteriormente, submetidos aos operadores de cruzamento (*crossover*) ou mutação. Esse processo é repetido até que o critério de parada seja atingido.

Dessa forma, os algoritmos genéticos otimizam problemas fornecendo a melhor solução possível de acordo com a aplicação desejada, mas não garante que haja convergência para uma solução ótima. A Figura 17 mostra um fluxograma com as principais etapas desse método.

Como vantagens, os Algoritmos Genéticos são robustos, podem ser usados em conjunto com outras técnicas e ser aplicados em diversos tipos de problemas. Como desvantagens, esse método apresenta dificuldade em encontrar a ótima solução exata e também possui a necessidade em ter um grande número de avaliações de função de aptidão, ocasionando em um desempenho mais lento (PINHO et al., 2013; LACERDA et al., 1999).



**Figura 17 – Fluxograma básico de um algoritmo genético**  
(Adaptado de: ROSA et al., 2009)

## 4. MODELO PROPOSTO

O modelo desenvolvido neste projeto tem como objetivo aprimorar a detecção de anormalidades respiratórias, decorrentes da fibrose cística, por meio das características fornecidas pela FOT. Outro objetivo é a geração de estruturas de Redes Bayesianas que melhor descrevam as relações existentes entre essas características.

Inicialmente, os dados da FOT podem passar pelo processo de seleção de atributos. Em seguida, essas características selecionadas também podem passar pelo método do produto cruzado, quando é gerado um novo conjunto de dados composto por colunas extras com o cálculo do produto dessas características. O exemplo descrito na equação (29) mostra a saída resultante da aplicação do produto cruzado em um conjunto de dados ( $X$ ), inicialmente composto por dois atributos,  $X_1$  e  $X_2$ . Dessa forma, é possível apresentar ao modelo dados em uma dimensão mais alta, como tentativa de aumentar sua acurácia<sup>1</sup>.

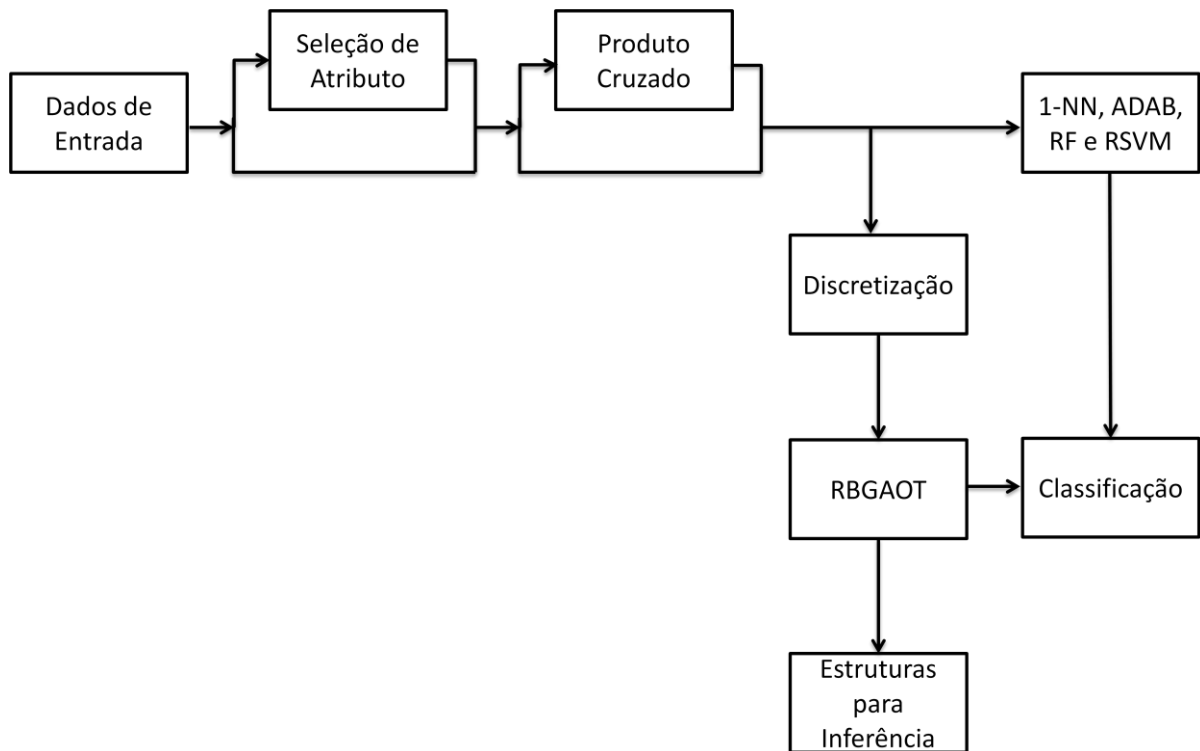
$$X = X_1^2 + X_1X_2 + X_2^2 \quad (29)$$

Com os atributos de entrada definidos, o conjunto de dados é submetido a quatro modelos de algoritmos de aprendizado de máquina: *1-Nearest Neighbor (1-NN)*, *Adaboost (ADAB)*, *Random Forest (RF)* e *Radial Support Vector Machine (RSVM)*. Esses dados também são submetidos às Redes Bayesianas sintetizadas por algoritmo genético (RBGAOT), usado para construção e seleção de estruturas, capazes de fornecer maior interpretabilidade das características mecânicas do sistema respiratório. Antes de passar pelo RBGAOT, os dados precisam ser discretizados a fim de viabilizar o cálculo das distribuições de probabilidade usado nesse método.

No presente estudo, todos os algoritmos de aprendizado de máquinas fornecem uma classificação como resultado final e tem seu desempenho calculado por meio da AUC. No caso do RBGAOT, também são obtidas estruturas que apresentaram melhor resultado durante a classificação e podem ser analisadas. Todo o modelo proposto, desenvolvido no *software* Matlab R2016a, foi descrito nos itens a seguir, de acordo com seu fluxograma resumido (Figura 18).

---

<sup>1</sup> Disponível em: <<https://jakevdp.github.io/PythonDataScienceHandbook/05.04-feature-engineering.html>>, Acessado em: 19/06/2018.



**Figura 18 – Fluxograma resumido do modelo proposto**

#### 4.1. Dados de Entrada

O conjunto de dados usado neste projeto foi obtido através de exames realizados por um sistema de técnica de oscilações forçadas (FOT), desenvolvido no Laboratório de Instrumentação Biomédica da UERJ (LIB-UERJ). O procedimento para realizar os exames pela FOT em cada indivíduo, consistiu em três medições com intervalo de um minuto e duração de 16 segundos. A fim de evitar o vazamento do ar e induzir a respiração normal pelo bocal do equipamento, foi necessário que os indivíduos fizessem uso de *clip* nasal e permanecessem sentados durante o exame (MIRANDA et al., 2013).

O aparelho usado forneceu a impedância do sistema respiratório em uma faixa de frequências de 4 a 32Hz, que foi medida com incrementos de 2Hz. Através de um bocal, um alto-falante gerou oscilações de pressão com amplitude de 1 cmH<sub>2</sub>O no sistema respiratório do paciente, durante sua respiração espontânea. Um pneumotacômetro e um transdutor de pressão foram usados para medir esses sinais de fluxo e de pressão, respectivamente, próximos à boca do paciente (LIMA et al., 2015).

## 4.2. Seleção de Atributos

Durante o projeto de construção de um classificador, a seleção de atributos de entrada é aplicada com o intuito de escolher as características que melhor descrevem o problema e, dessa forma, melhorar o desempenho do algoritmo. O uso dessa estratégia também permite a redução da complexidade do modelo, já que ocorre uma redução nos atributos e, conseqüentemente, uma diminuição do número de parâmetros que precisa ser estimado. Outras vantagens da técnica são o aumento da velocidade de execução do algoritmo, a possibilidade em visualizar os dados e a obtenção de uma melhor compreensão do processo que gera os resultados obtidos (GUYON et al., 2003).

Há duas formas principais de realizar a seleção de atributos. A primeira forma é através de um especialista que destaca os parâmetros que melhor descrevem o problema, com base em sua própria experiência. O outro método é feito de forma automática e pode utilizar técnicas de filtragem, *Wrapper* ou o método embutido. A filtragem realiza a classificação ordenada das características antes que os dados sejam submetidos ao algoritmo. Essa ordenação é feita com base em um critério escolhido, como coeficientes de correlação ou testes estatísticos. A técnica de *Wrapper* também é usada antes da classificação, fazendo uso de algoritmos de aprendizado de máquina para avaliar subconjuntos criados a partir do conjunto de dados (HORTA et al., 2010). O subconjunto que apresentar melhor desempenho tem seus atributos selecionados para o classificador. Já o método embutido realiza a seleção de variáveis durante o treinamento. Ele é usado especificamente em alguns algoritmos, como as árvores de decisões, que selecionam atributos no seu próprio processo de construção de um modelo (AMARAL et al., 2013). O RBGAOT também pode ser considerado um caso de algoritmo com seleção embutida, visto que mesmo submetendo um conjunto de atributos às Redes Bayesianas, variáveis podem ser descartadas durante a construção das estruturas geradas.

Neste projeto, o método *Wrapper* foi usado para a seleção de atributos por ser um método heurístico e guiado em sua busca pelo conjunto de atributos que maximize a média da AUC. Essa busca foi realizada de forma direta (*forward*), onde os atributos são acrescentados um por vez com base em um critério, até completar o subconjunto. O critério escolhido foi a taxa de acerto no algoritmo  $K$ -NN, com  $K$  igual a 1 (1-NN). O treinamento do classificador 1-NN foi feito através da validação cruzada *leave-one-out*, onde uma amostra  $n$  é testada com base nas  $n-1$  amostras restantes (RODRIGUES et al., 2017). A seleção feita pelo 1-NN foi

implementada pela função *featself* disponível na *toolbox Pattern Recognition* (prtools) do *software* Matlab (DUIN, 2007).

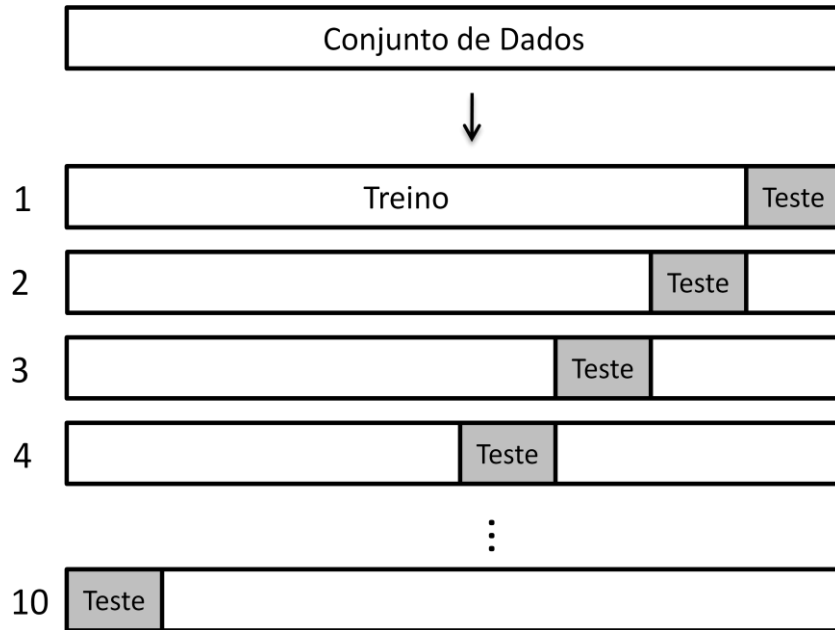
Como o uso da AUC é recomendado em diagnósticos médicos (METZ, 1978; HANLEY et al., 1982), essa medida também foi testada como critério de seleção de atributos. Entretanto, as variáveis selecionadas por esse método não apresentaram desempenho superior ao das variáveis selecionadas pelo algoritmo 1-NN durante os experimentos. A seleção de atributos também foi feita pela consulta a um especialista, que elegeu o mesmo conjunto de variáveis obtido pelo algoritmo 1-NN.

### 4.3. Treinamento do Modelo

A técnica de validação cruzada foi usada durante o treinamento do modelo. Devido a pouca quantidade de amostras disponíveis e a grande quantidade de atributos fornecidos pela FOT, optou-se pelo método de validação cruzada por  $k$ -pastas, onde uma parte dos dados é separada para treino e o restante é destinado para teste do modelo. Durante esse processo, os dados são divididos em  $k$  pastas, sendo geralmente uma pasta usada para testar e  $k-1$  pastas usadas para o aprendizado do modelo.

Esse processo é repetido  $k$  vezes e a cada iteração são usadas diferentes pastas para o conjunto de treino e teste, gerando então,  $k$  medidas de erro. A média desses  $k$  erros é a medida de generalização do modelo. Dessa forma, é possível evitar uma estimativa muito otimista fornecida por alguma das  $k$  partições, como poderia ocorrer na validação cruzada feita pela técnica *hold-out* (HASTIE, 2008). No exemplo da Figura 19 é possível observar a divisão para  $k$  igual a 10, valor escolhido para o uso da validação cruzada neste projeto.

A métrica usada para a seleção dos classificadores e a configuração escolhida para cada modelo estão descritas dos itens 4.4 a 4.6.



**Figura 19 – Divisão para validação cruzada com 10 pastas**

#### 4.4. Medida de desempenho

A medida de desempenho usada para a seleção dos melhores modelos de classificação foi feita com base na área sob a curva ROC. A AUC é uma ferramenta normalmente usada para diagnóstico médico (METZ, 1978; HANLEY et al., 1982), que também fornece informações sobre a eficácia de algoritmos de aprendizado de máquina (HUANG et al, 2005).

Considerando um conjunto de dados com duas classes, positiva ( $p$ ) e negativa ( $n$ ), os rótulos obtidos na classificação desses dados podem ser respectivamente,  $P$  e  $N$ . A matriz confusão da Figura 20 fornece quatro possibilidades para uma instância ao ser classificada. Se a instância for positiva e classificada como positiva, é um caso de verdadeiro positivo ( $VP$ ). Se for classificada como negativa pelo algoritmo, é um caso de falso negativo ( $FN$ ). Se a instância for negativa e classificada como negativa, é um caso verdadeiro negativo ( $VN$ ). Se for classificada como positiva, é um caso de falso positivo ( $FP$ ) (FAWCETT, 2006).



		<u>Classe Verdadeira</u>	
		<i>p</i>	<i>n</i>
<u>Resultado da Classificação</u>	<i>P</i>	Verdadeiro Positivo (VP)	Falso Positivo (FP)
	<i>N</i>	Falso Negativo (FN)	Verdadeiro Negativo (VN)

**Figura 20 – Matriz confusão das possíveis classificações de uma instância**  
(Adaptado de: FAWCETT, 2006)

Pelos valores da matriz na Figura 20, podem ser calculadas diversas métricas, sendo a diagonal principal as classificações corretas do modelo. Para este projeto, além da AUC, também foram usadas a sensibilidade e a especificidade. A equação (30) mostra o cálculo da sensibilidade, que corresponde à probabilidade de ter uma classificação positiva quando a instância é positiva. Já a equação (31), mostra a probabilidade de ter uma classificação negativa quando a instância é negativa.

$$\text{Sensibilidade} = \frac{VP}{VP + FN} \quad (30)$$

$$\text{Especificidade} = \frac{VN}{VN + FP} \quad (31)$$

Para construir a curva ROC, são usados os valores de sensibilidade no eixo *y* e o resultado da subtração de 1-especificidade no eixo *x*, caracterizando a relação entre os benefícios (verdadeiros positivos) e os custos (falsos positivos) de um modelo (FAWCETT, 2006). Caso não seja possível selecionar visualmente um classificador, a área sob a curva ROC é calculada e, pelo seu valor, é possível definir o classificador com melhor desempenho.

## 4.5. Classificadores

Dentre os classificadores descritos no capítulo 3, quatro foram escolhidos e implementados pela *toolbox* *prtools*<sup>2</sup>, com base em trabalhos realizados nessa linha de pesquisa (AMARAL et al., 2013; AMARAL et al., 2015; AMARAL et al., 2017): *K-Nearest Neighbor* (K-NN), *Adaboost* (ADAB), *Random Forest* (RF) e *Radial Support Vector Machine* (RSVM), cujos parâmetros foram definidos de acordo com o desempenho durante a validação cruzada.

No algoritmo K-NN, o valor de *K* foi definido de acordo com o erro encontrado durante o treinamento usando a área sob a curva ROC (AUC) como medida de desempenho ( $E_{AUC}$ ). Foram avaliados diferentes valores de *K*, sendo *K* igual a 1 (1-NN) a configuração que apresentou melhor desempenho, conforme Tabela 8.

**Tabela 8 – Resultados do treinamento do algoritmo K-NN**

<b>K</b>	<b><math>E_{AUC}</math></b>
1	0,1252
3	0,1768
5	0,1467
7	0,1554
9	0,1746
11	0,1807
13	0,1866
15	0,1582
17	0,1501
19	0,1563

A quantidade de árvores de decisões usadas como classificadores simples no algoritmo *ADAB*, foi selecionada de acordo com o erro de AUC encontrado durante o treinamento. Conforme Tabela 9, o número de árvores de decisões usado que apresentou melhor desempenho foi 200.

<sup>2</sup> *Prtools: Toolbox for Pattern Recognition*. Disponível em: <<http://prtools.org/>>, Acessado em: 16/03/2018.

**Tabela 9 – Resultados do treinamento do algoritmo ADAB**

Número de Árvores de Decisão	$E_{AUC}$
50	0,1238
100	0,1296
150	0,1192
200	0,1019
250	0,1120

O algoritmo RF foi implementado de acordo com o erro encontrado durante o treinamento, utilizando a AUC como medida do desempenho do modelo testado. Nesse caso, foram avaliadas a quantidade de subgrupos de atributos formados e a quantidade de árvores geradas. De acordo com a Tabela 10, a configuração que apresentou menor erro possui 50 árvores geradas e tamanho do subconjunto de atributos igual a 1.

**Tabela 10 – Resultados do treinamento do algoritmo RF**

Árvores geradas	10	20	50	100	150
Subgrupos					
1	0,1664	0,1152	0,0970	0,1268	0,1361
2	0,1245	0,1383	0,1030	0,1375	0,1139
3	0,1190	0,1363	0,1155	0,1235	0,1079
5	0,1472	0,1304	0,1147	0,1142	0,1219
7	0,1401	0,1246	0,1257	0,1148	0,1053

No caso do algoritmo RSVM, foi necessário definir dois parâmetros: o desvio padrão da base radial ( $r$ ) e o parâmetro de regularização ( $C$ ). A busca por esses parâmetros foi realizada por uma validação cruzada interna durante o treinamento<sup>3</sup>.

Com o intuito de extrair informações a respeito das relações entre os atributos para obter uma explicação da classificação, utilizou-se um classificador baseado em Redes Bayesianas. De acordo com trabalhos já realizados (TONDA et al., 2012; LARRAÑAGA et al., 1996), é possível fazer uso de algoritmos evolutivos para o aprendizado dessas redes.

<sup>3</sup> Disponível em: <<http://www.37steps.com/prhtml/prtools/rbsvc.htm>>, Acessado em: 16/03/2018

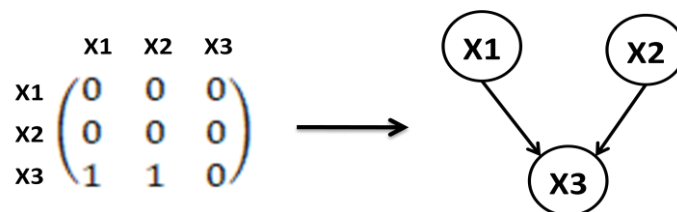
No artigo “*Bayesian Network Structure Learning from Limited Datasets through Graph Evolution*” (TONDA et al., 2012), é proposto o uso de um algoritmo evolutivo para a aprendizagem das estruturas de Redes Bayesianas, tomando por base um conjunto de dados com número de amostras limitado. Esse trabalho apresenta a possibilidade de trabalhar usando diretamente estruturas gráficas. Sua função *fitness* é baseada na métrica de informação *Akaike*, considerando a precisão e a complexidade da estrutura fornecida durante o treinamento do modelo.

Já o artigo (LARRAÑAGA et al., 1996), mostra o uso de algoritmo genético para realizar a busca da melhor estrutura de Rede Bayesiana, com base em um conjunto de dados. Esse trabalho usa matrizes para representar suas soluções e propõe o uso de um operador responsável por corrigir redes que não forem DAG. A função *fitness* calcula a métrica *K2*, onde a estrutura com maior valor de probabilidade conjunta, dado um conjunto de treinamento, é selecionada.

Com base nesses trabalhos, a técnica de algoritmos genéticos foi utilizada para realizar a busca pela estrutura com melhor desempenho. A descrição desse método escolhido está nos itens a seguir.

#### 4.6. Redes Bayesianas sintetizadas com Algoritmos Genéticos

As Redes Bayesianas foram implementadas pela *toolbox Probabilistic Graphical Model 9.2.3* (PGM<sup>4</sup>), onde a primeira etapa consiste na leitura de um grafo direcionado e acíclico (MENSXMACHINA, 2011). No *software* Matlab essas estruturas são representadas por matrizes esparsas binárias, onde elementos iguais a 0 são eliminados e os elementos iguais a 1 representam as ligações entre as variáveis do problema, conforme Figura 21:



**Figura 21 – Representação de uma Rede Bayesiana em matriz esparsa**

<sup>4</sup> *Toolbox* PGM: *Probabilistic Graphical Model 9.2.3*. Disponível em: <<http://mensxmachina.org/en/software/pgm-toolbox/>>, Acessado em: 16/03/2018

Durante essa primeira etapa, podem ser apresentadas ao algoritmo estruturas inválidas de Redes Bayesianas, como redes que não sejam DAG ou redes que não possuam a variável classe. Essas matrizes são identificadas através do ajuste de linhas e colunas, realizado pela *toolbox* PGM, retornando redes com apenas uma linha. Dessa forma, elas são identificadas e recebem valores de AUC igual à zero, para que sejam descartadas durante as próximas iterações.

Em seguida, o aprendizado das distribuições de probabilidade conjunta (DPC) é feito pelo algoritmo BDeu (*Bayesian Dirichlet Equivalent Uniform*), desenvolvido por Heckerman (ONISKO et al., 2001) e baseado na métrica *Bayesian Dirichlet* desenvolvida por Cooper e Herskovits (COOPER et al., 1991). Essa pontuação corresponde ao logaritmo da probabilidade à posteriori de uma rede  $B_s$ , dado um conjunto de dados  $A$ , logo, a pontuação é obtida pelo cálculo de  $\log(P(B_s|A))$ . A métrica BDeu deve corresponder à capacidade de uma rede em capturar a probabilidade conjunta dos dados e prever novas amostras, apresentando então, sua relação direta com a capacidade de inferência da estrutura analisada (BROWN et al., 2004).

Durante essa segunda etapa, as tabelas de DPC são calculadas e contém: o nome da variável analisada, os diversos valores que seus respectivos nós pais podem assumir e as probabilidades condicionadas aos nós pais. No exemplo da Tabela 11, a variável analisada é  $X3$ , os valores que as variáveis podem assumir são positivo ( $p$ ) ou negativo ( $n$ ) e os nós pais são  $X1$  e  $X2$ . Com a estrutura definida e as tabelas de DPC calculadas, a função *bayesnet* realiza a construção da Rede Bayesiana.

**Tabela 11 – Exemplo de tabela de DPC**

	$P(X3=p X1,X2)$	$P(X3=n X1,X2)$
$X1=p, X2=p$	0,75	0,25
$X1=n, X2=p$	0,05	0,95
$X1=p, X2=n$	0,85	0,15
$X1=n, X2=n$	0,10	0,90

A terceira etapa disponível na *toolbox* PGM é um mecanismo de inferência que pode ser usado para a classificação de um conjunto de teste, através de uma Rede Bayesiana já construída. O algoritmo *Junction Tree* (BARBER, 2003) é aplicado com o intuito de fornecer uma ótima sequência de decomposição, ou marginalização, dessa rede. A estrutura resultante desse algoritmo é capaz de calcular a distribuição de probabilidade marginal de cada amostra

de teste. Essa distribuição corresponde às probabilidades de uma nova amostra pertencer a cada uma das classes do problema.

Considerando como exemplo um conjunto de dados com duas classes: positiva e negativa, durante a terceira etapa são calculadas as probabilidades marginais de uma amostra de teste pertencer à classe positiva ou a classe negativa. Esses valores são importantes para a classificação com Redes Bayesianas, sendo que a maior probabilidade marginal encontrada define o rótulo dessa amostra.

As três etapas descritas, são a base para o uso das Redes Bayesianas pela *toolbox* PGM. Porém, estratégias externas podem ser aplicadas para trabalhar com o conjunto de dados, gerar diversas estruturas e selecionar aquela que melhor descreve o problema. Neste trabalho, duas estratégias foram escolhidas: a discretização dos dados de entrada e a aplicação de um algoritmo genético no treinamento das redes.

#### **4.6.1. Discretização dos Dados**

As tabelas de distribuição de probabilidade conjunta de uma Rede Bayesiana quantificam as relações existentes entre suas variáveis, abordando os diversos estados que seus nós pais podem assumir. Quando o conjunto de treinamento é composto por valores contínuos, é necessário escolher uma abordagem para esses dados. Neste trabalho, foi escolhido o método de discretização para possibilitar o uso das Redes Bayesianas, devido sua simples implementação e interpretação na leitura das tabelas de distribuição de probabilidade conjunta, por parte da equipe médica.

A discretização pode ser definida como o processo de transformação de uma variável contínua em uma variável discreta. Dessa forma, as amostras passam a ser apresentadas em intervalos cujos limiares são pontos de corte que podem ser definidos por diversos tipos de cálculos. O uso de dados discretizados facilita até mesmo a análise de características do problema por parte dos usuários e especialistas.

Todos os métodos usados para discretização podem ser classificados como supervisionados e não supervisionados. No método supervisionado, as classes dos atributos são levadas em consideração, já no caso não-supervisionado, a discretização é feita considerando apenas os valores dos atributos (CARVALHO, 2010).

O método escolhido para discretizar os dados fornecidos pela FOT é supervisionado e usa o conceito de entropia. Esse método tem como objetivo determinar um ponto de corte que seja capaz de gerar os mais puros subconjuntos possíveis, com base no maior ganho de informação (FAYYAD et al., 1993):

$$Ganho(A) = E(\text{conjunto Atual}) - \sum E(\text{subconjuntos}) \quad (32)$$

Sendo:

$E()$  a entropia

*Conjunto atual* os dados a serem discretizados

*Subconjuntos* os intervalos criados para o atributo  $A$

De acordo com a equação (32), para aumentar o ganho de informação é necessário minimizar o somatório das entropias relativas aos subconjuntos (MERSCHMANN, 2007):

$$E(\text{subconjunto}) = - \sum_{i=1}^c p_i \log_2(p_i) \quad (33)$$

Sendo:

$C$  a quantidade de classes do problema

$p_i$  a probabilidade da classe  $i$  ocorrer

Os pontos de corte são calculados até que o critério de parada seja alcançado e haja uma diminuição dos intervalos em cada atributo. O princípio usado como critério de parada é o MDL (*Minimum Description Length*) que compara se a criação de um novo intervalo aumenta o ganho de informação. Quando não houver mais aumento, o ponto de corte que fornecer o maior ganho de informação é selecionado (FAYYAD et al., 1993).

### 4.6.2. RBGAOT

A estratégia escolhida para realizar a aprendizagem da estrutura das Redes Bayesianas foi o uso de algoritmos genéticos, por meio da *toolbox Genetic Algorithms for optimization* (GAOT). Para implementar o uso dessa técnica na criação e seleção da melhor estrutura que descreva as relações entre as variáveis do problema, foi feita a junção das *toolboxes* de Redes Bayesianas e Algoritmo Genético, chamada RBGAOT.

O RBGAOT gera de forma aleatória diversas redes representadas em matrizes de adjacência que apresentam possíveis soluções ao problema. Cada rede é construída com base nessas matrizes e tem suas distribuições de probabilidade calculadas pela *toolbox* de Redes Bayesianas (PGM). Uma vez que a rede já esteja com todas as suas características definidas, é possível analisar seu desempenho classificatório.

Há cinco elementos principais que precisam ser definidos para o uso do RBGAOT: representação do cromossomo, criação de uma população inicial, função de aptidão, função de seleção e operadores genéticos.

### 4.6.3. Representação do cromossomo

Um cromossomo equivale a cada indivíduo da população em um algoritmo genético e é composto por uma sequência de genes. No RBGAOT, um cromossomo corresponde à estrutura de uma Rede Bayesiana com  $n$  variáveis e genes formados por dígitos binários. Uma rede pode ser representada por uma matriz de adjacência  $C$  de tamanho  $n \times n$ , cujos elementos são descritos de acordo com as ligações existentes, ou não, entre  $j$  e  $i$ , conforme a seguir:

$$c_{ij} = \begin{cases} 1, & \text{se } j \text{ é pai de } i \\ 0, & \text{se não existe ligação entre } j \text{ e } i \end{cases} \quad (34)$$

Dessa forma, as ligações entre as variáveis são expressas em uma matriz que, por sua vez, pode ser decomposta coluna a coluna para gerar um vetor (LARRAÑAGA et al., 1996), conforme o exemplo a seguir:

$$C = \begin{pmatrix} c_{11} & \cdots & c_{1n} \\ \vdots & \ddots & \vdots \\ c_{n1} & \cdots & c_{nn} \end{pmatrix}$$

$$\text{Vetor} = c_{11} c_{21} c_{31} \dots c_{n1} \quad c_{12} c_{22} c_{32} \dots c_{n2} \quad \dots \quad c_{1n} c_{2n} c_{3n} \dots c_{nn}$$



#### 4.6.4. População Inicial

Um indivíduo na população inicial corresponde a uma estrutura de Rede Bayesiana que pode ser selecionada para formar a próxima geração ou passar pelos operadores de mutação e *crossover*. A criação de uma população inicial  $P_{ij}$ , é feita de forma aleatória com uma distribuição uniforme ( $U$ ), conforme a equação (35) (CIVICIOGLU, 2013):

$$P_{ij} = U(A_j, B_j), \text{ para } i = 1, 2, \dots, N \text{ e } j = 1, 2, \dots, D \quad (35)$$

Sendo:

$A_j$  e  $B_j$  o limite mínimo e máximo de cada indivíduo do vetor  $P$ , respectivamente.

$N$  o tamanho da população

$D$  a dimensão dos indivíduos da população

No RBGAOT, a população inicial foi criada com 15 indivíduos formados por valores entre 0 e 1, e 20 gerações. Em seguida, toda a população teve seus genes aproximados para valores binários, conforme o sistema da equação (36):

$$P_{ij} = \begin{cases} 1, & U(A_j, B_j) > 0,5 \\ 0, & U(A_j, B_j) < 0,5 \end{cases} \quad (36)$$

#### 4.6.5. Função de Avaliação

A função de avaliação ou *fitness* é usada para determinar a aptidão de cada indivíduo gerado durante a busca pela melhor solução, no RBGAOT. Cada vetor, que representa um indivíduo gerado, é recebido por essa função e convertido para uma matriz esparsa, conforme o processo descrito no item 4.6.3. Uma vez que se tenha a estrutura em formato de matriz, o RBGAOT faz uso da *toolbox* de Redes Bayesianas para o treinamento e teste da estrutura gerada.

As duas saídas fornecidas pela função de avaliação desse algoritmo são: o valor da AUC da estrutura testada e um vetor com a probabilidade das amostras do grupo de teste obtidas durante a classificação. Essas probabilidades serão usadas na construção da curva ROC.

#### 4.6.6. Função de Seleção

O RBGAOT realiza a seleção de indivíduos de forma probabilística pelo método de roleta, onde os mais aptos têm maior probabilidade de serem escolhidos para formar a próxima geração. Também foi usado o *ranking* por normalização geométrica, para ordenação dos indivíduos ( $i$ ) de acordo com a probabilidade  $P(i)$ , definida conforme equação (37). Essa técnica evita que indivíduos com aptidão muito acima da média sejam sempre escolhidos, levando o algoritmo a uma convergência prematura (HOUCK et al., 1995).

$$P(i) = q_t(1 - q)^{r-1}, \quad q_t = \frac{q}{1-(1-q)^T} \quad (37)$$

Sendo:

$i$  o indivíduo ou possível solução

$q$  a probabilidade de selecionar o melhor indivíduo

$r$  o *rank* do indivíduo (onde 1 é o melhor)

$T$  o tamanho da população

#### 4.6.7. Operadores Genéticos

Os operadores genéticos são mecanismos básicos de busca usados pelo algoritmo genético e tem como função criar novos indivíduos com base na população já existente. Um dos principais operadores é o *crossover*, que usa dois indivíduos pais para gerar dois novos indivíduos filhos através do cruzamento de seus cromossomos.

Apesar da *toolbox* GAOT disponibilizar diversos tipos de *crossover*, o cruzamento simples foi o que apresentou melhor desempenho, onde dois indivíduos pais ( $X$  e  $Y$ ) de uma população de tamanho  $m$  formam dois novos indivíduos ( $X'$  e  $Y'$ ). Para isso, um número aleatório  $r$  é criado por uma distribuição uniforme de 1 até  $m$ , para ser aplicado como ponto de corte, conforme as equações (38) e (39) (HOUCK et al., 1995).

$$x'_i = \begin{cases} x_i, & \text{se } i < r \\ y_i, & \text{caso contrário} \end{cases} \quad (38)$$

$$y'_i = \begin{cases} y_i, & \text{se } i < r \\ x_i, & \text{caso contrário} \end{cases} \quad (39)$$

Quanto maior for o valor da taxa de *crossover* escolhida, maior será a quantidade de novas estruturas promissoras, uma vez que ele combina as características de pais com alta aptidão. Entretanto, isso pode levar a uma convergência prematura da evolução. Caso a taxa de *crossover* seja muito baixa, o algoritmo poderá demorar a convergir para uma solução aceitável. Sendo assim, a taxa de 0,6 foi escolhida para o *crossover* no RBGAOT.

O operador de mutação também é muito usado nos algoritmos genéticos. Seu objetivo é alterar o cromossomo de um indivíduo  $X$  da população e gerar apenas uma nova solução  $X'$ . Dentre os operadores de mutação disponíveis na *toolbox* usada, a mutação binária apresentou melhor desempenho, realizando alterações com base em uma probabilidade ( $p_m$ ). Os genes dos novos indivíduos são definidos conforme a equação (40) (HOUCK et al., 1995):

$$x'_i = \begin{cases} 1 - x_i, & \text{se } U(0,1) < p_m \\ x_i, & \text{caso contrário} \end{cases} \quad (40)$$

A taxa de mutação pode evitar que o algoritmo fique estagnado em uma solução, permitindo que a busca seja realizada em mais pontos no espaço de soluções. Valores mais altos tendem a tornar essa busca aleatória. Dentre os valores testados na faixa de 0,005 e 0,1, a taxa de mutação de 0,01 gerou os melhores resultados e foi o valor escolhido para o RBGAOT.

## 5. ESTUDO DE CASO

Neste capítulo, foram realizados experimentos para testar os cinco algoritmos de aprendizado de máquina descritos no capítulo 3. Após uma descrição mais detalhada, os parâmetros fornecidos pela FOT foram submetidos a experimentos individuais e em conjunto. Outros testes foram feitos com aplicação de métodos como produto cruzado e seleção de variáveis. Em seguida, foram selecionadas estruturas geradas pelas Redes Bayesianas sintetizadas com Algoritmo Genético, para análise das ligações entre as variáveis usadas e suas tabelas de distribuição de probabilidade conjunta.

### 5.1. Descrição do Conjunto de Dados

O conjunto de dados usado neste projeto foi obtido por um sistema de oscilações forçadas (FOT), desenvolvido no Laboratório de Instrumentação Biomédica da UERJ. Os exames foram realizados em 23 indivíduos do grupo controle e 27 portadores de Fibrose Cística, que formam um grupo de teste. Em cada exame foram feitas três medidas, o que totalizou um conjunto de dados de 150 instâncias para os experimentos. As informações fornecidas pela FOT estão na tabela a seguir:

**Tabela 12 – Parâmetros fornecidos pela FOT**

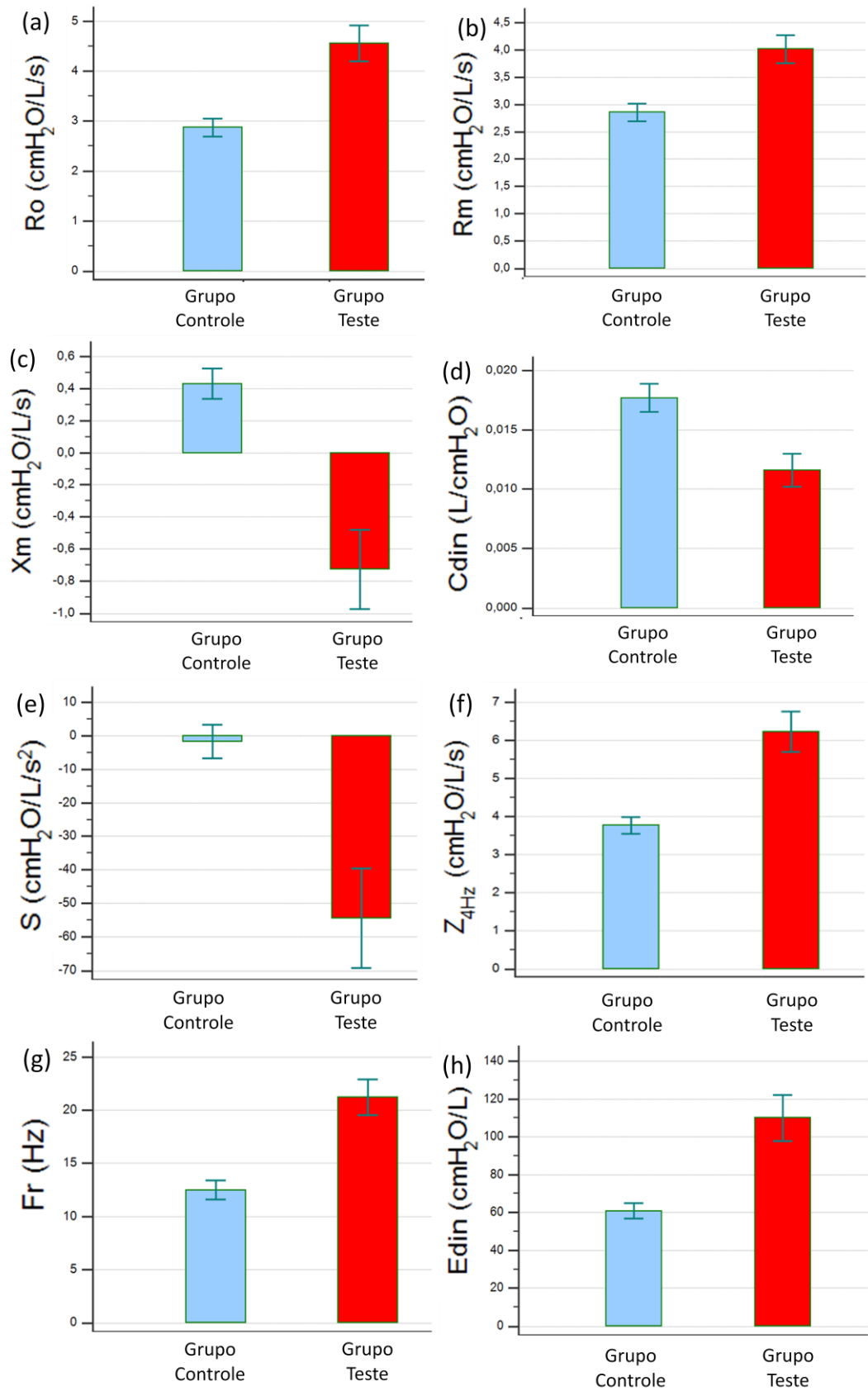
Parâmetro	Descrição do Parâmetro
$R_o$	Resistência no Intercepto
$R_m$	Resistência Média
$X_m$	Reatância Média
$C_{din}$	Complacência Dinâmica
$S$	Inclinação da Curva de Resistência
$Z_{4Hz}$	Impedância em 4Hz
$F_r$	Frequência de Ressonância
$E_{din}$	Elastância Dinâmica

As características de indivíduos pertencentes ao grupo controle e teste, foram comparadas na Figura 22. Os gráficos com barras mostram as médias das variáveis calculadas em um intervalo de confiança de 95%, cujos valores de desvio padrão são indicados acima ou abaixo das barras. Por exemplo, o valor médio da resistência  $R_o$  no grupo controle é

$2,87 \pm 0,76$ . Já no grupo de teste, a média de  $R_o$  sobe para  $4,56 \pm 1,61$ . Através da Análise de Variância (ANOVA), todos os parâmetros da FOT mostraram diferença significativa nos seus respectivos valores de média ( $p < 0,001$ ).

De acordo com a Figura 22, houve um aumento na média das variáveis  $R_o$ ,  $R_m$ ,  $Z_{4Hz}$ ,  $F_r$  e  $E_{din}$  dos indivíduos do grupo teste, se comparado aos do grupo controle. Ou seja, indivíduos portadores de fibrose cística geralmente possuem valores mais altos de resistências ( $R_o$  e  $R_m$ ), impedância ( $Z_{4Hz}$ ), frequência de ressonância ( $F_r$ ) e elastância ( $E_{din}$ ), se comparado a não portadores da doença. Já as variáveis  $X_m$ ,  $C_{din}$  e  $S$  do grupo teste, apresentaram uma diminuição em sua média. Logo, conclui-se que portadores de fibrose cística possuem valores mais negativos de reatância ( $X_m$ ) e inclinação da curva de resistência ( $S$ ), e valores menores de complacência ( $C_{din}$ ).

Para submeter os dados da FOT nas Redes Bayesianas, todas as amostras do conjunto de dados foram discretizadas e para cada característica da Tabela 12 foi estabelecido um ponto de corte (Tabela 13). Os valores abaixo desse ponto foram rotulados como 1, representando valores mais baixos que a variável pode assumir. Já os valores acima do ponto de corte foram rotulados como 2, representando os valores mais altos da variável. No caso da variável *classe*, indivíduos do grupo controle receberam o rótulo 0, e indivíduos do grupo teste receberam o rótulo 1. Com base nessas informações, o comportamento geral das características da FOT pode ser resumido conforme a Tabela 14.



**Figura 22 – Comparação dos parâmetros da FOT de indivíduos do grupo controle e do grupo teste**

**Tabela 13 – Pontos de corte para discretização dos parâmetros da FOT, média e desvio padrão**

Parâmetro	Ponto de corte	Média ( $\pm$ Desvio Padrão)
R <sub>o</sub>	3,31	3,78 $\pm$ 1,54
R <sub>m</sub>	3,21	3,48 $\pm$ 1,13
X <sub>m</sub>	0,18	-0,19 $\pm$ 1,04
C <sub>din</sub>	0,014	0,014 $\pm$ 0,007
S	-10,25	-30,15 $\pm$ 57,61
Z <sub>4Hz</sub>	4,44	5,10 $\pm$ 2,25
F <sub>r</sub>	14,19	17,20 $\pm$ 7,50
E <sub>din</sub>	72,54	87,40 $\pm$ 48,80

**Tabela 14 – Comportamento geral das características do grupo controle e do grupo teste**

	R <sub>o</sub>	R <sub>m</sub>	Z <sub>4Hz</sub>	F <sub>r</sub>	E <sub>din</sub>	X <sub>m</sub>	C <sub>din</sub>	S	Classe
<b>Grupo Controle</b>	1	1	1	1	1	2	2	2	0
<b>Grupo Teste</b>	2	2	2	2	2	1	1	1	1

## 5.2. Experimento Individual dos Atributos

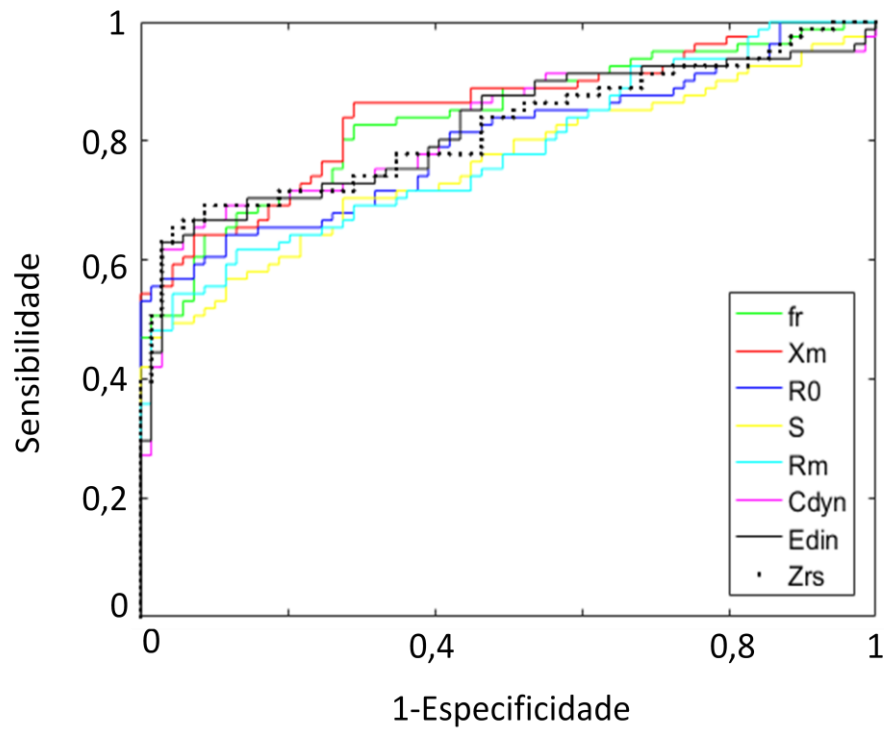
Cada atributo fornecido pela FOT foi submetido à análise individual para que seu desempenho em classificar pacientes fosse testado. Todos os parâmetros tiveram suas medidas de AUC, erro padrão ( $E_{AUC}$ ) e intervalo de confiança de 95% (IC 95%) calculados (DELONG et al., 1988), conforme Tabela 15.

A reatância X<sub>m</sub> e a frequência F<sub>r</sub> foram os parâmetros que apresentaram melhor desempenho individual com valores de AUC iguais a 0,85 e 0,84, respectivamente. Os demais parâmetros apresentaram valores de AUC entre 0,76 e 0,82. Sendo assim, o desempenho de todos os atributos analisados separadamente se enquadra na faixa de acurácia moderada (0,70 a 0,90), sendo observada diferença significativa apenas entre os valores de AUC de S e F<sub>r</sub> ( $p < 0,01$ ) e entre S e X<sub>m</sub> ( $p < 0,005$ ).

As curvas ROC com o desempenho de cada atributo foram traçadas, mostrando que a área sob a curva é maior na faixa final do eixo x, onde uma maior quantidade de falsos positivos é aceita (Figura 23).

**Tabela 15 – Desempenho individual dos parâmetros da FOT na classificação de pacientes**

	AUC	$E_{AUC}$	IC 95%
$F_r$ (Hz)	0,84	0,03	0,77-0,89
$X_m$ (cmH <sub>2</sub> O/L/s)	0,85	0,03	0,78-0,90
$R_o$ (cmH <sub>2</sub> O/L/s)	0,80	0,04	0,73-0,86
$S$ (cmH <sub>2</sub> O/L/s <sup>2</sup> )	0,76	0,04	0,68-0,83
$R_m$ (cmH <sub>2</sub> O/L/s)	0,78	0,04	0,70-0,84
$C_{din}$ (L/cmH <sub>2</sub> O)	0,82	0,04	0,75-0,88
$E_{din}$ (cmH <sub>2</sub> O/L)	0,82	0,04	0,75-0,88
$Z_{4Hz}$ (cmH <sub>2</sub> O/L/s)	0,82	0,04	0,75-0,88



**Figura 23 – Curvas ROC dos parâmetros da FOT**



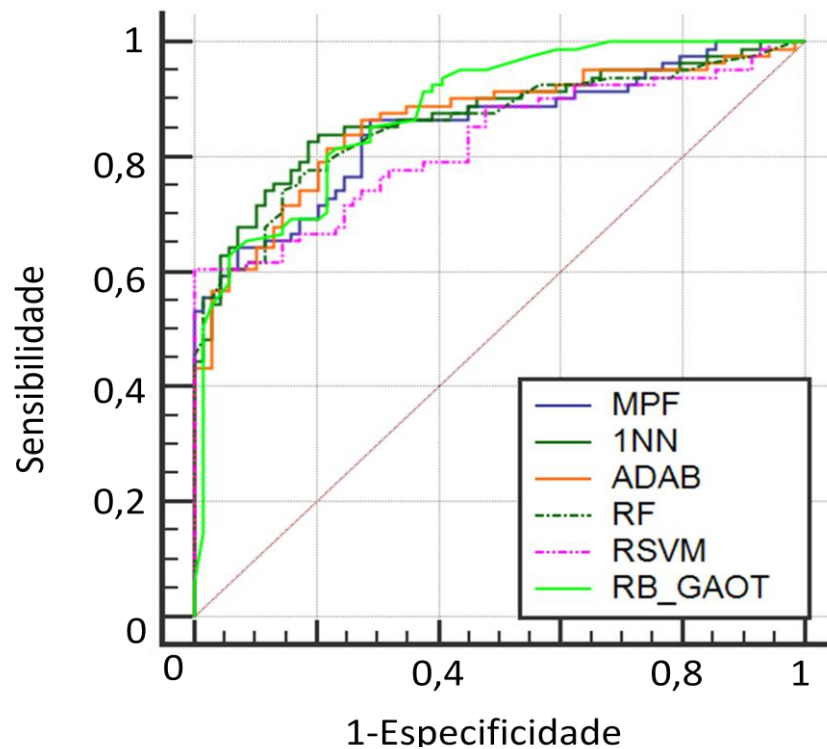
### 5.3. Experimento com Oito Atributos

Para esse experimento, os oito atributos da FOT foram usados nos seguintes algoritmos de aprendizado de máquina: *1-Nearest Neighbor* (1-NN), *Adaboost* (ADAB), *Random Forest* (RF), *Radial Support Vector Machine* (RSVM) e Redes Bayesianas sintetizadas com Algoritmos Genéticos (RBGAOT). Esses cinco classificadores também foram comparados com o melhor parâmetro da FOT (MPF), a variável  $X_m$ . Pela Tabela 16, pode-se observar que o algoritmo RBGAOT apresentou melhor desempenho com AUC igual a 0,88. O segundo melhor resultado foi obtido pelo 1-NN com valor de AUC igual a 0,87.

Além da AUC, foram calculadas as probabilidades de ter um resultado positivo quando o indivíduo for portador da doença, denominada sensibilidade (Sens). Também foram calculadas as probabilidades de ter um resultado negativo quando o indivíduo não portar a doença, denominada especificidade (Esp). O intervalo de confiança está abaixo dos respectivos valores de Sens, Esp e AUC (DELONG et al, 1988). Na Figura 24, pode-se observar que a área sob a curva ROC dos classificadores é maior na faixa final dos eixos, onde é aceito maior quantidade de casos falsos positivos.

**Tabela 16 – Resultado dos oito parâmetros da FOT submetidos aos classificadores**

	Sens (%)	Esp (%)	AUC	$E_{AUC}$
<b>MPF</b>	86,42 (77,0-93,0)	71,01 (58,8-81,3)	0,85 (0,78-0,90)	0,03
<b>1-NN</b>	82,72 (72,7-90,2)	81,16 (69,9-89,6)	0,87 (0,81-0,92)	0,03
<b>ADAB</b>	81,48 (71,3-89,2)	78,26 (66,7-87,3)	0,86 (0,79-0,91)	0,03
<b>RF</b>	74,07 (63,1-83,2)	85,51 (75,0-92,8)	0,86 (0,79-0,91)	0,03
<b>RSVM</b>	60,49 (49,0-71,2)	100 (94,8-100,0)	0,82 (0,75-0,88)	0,04
<b>RBGAOT</b>	80,25 (69,9-88,3)	78,26 (66,7-87,3)	0,88 (0,82-0,93)	0,03



**Figura 24 – Curvas ROC do experimento com todos os parâmetros da FOT**

A Tabela 17 mostra, em pares, a comparação das áreas sobre as curvas ROC de todos os métodos e o erro padrão em um intervalo de confiança de 95% (DELONG et al, 1988), sendo a interseção entre linha e coluna a diferença entre dois classificadores. Nesse experimento, não foi observada diferença significativa em nenhum dos quinze pares analisados ( $p > 0,05$ ). Isto pode ter ocorrido devido ao tamanho do conjunto de dados usado que conta apenas com 150 instâncias, limitando assim, a quantidade de amostras usada para os testes.

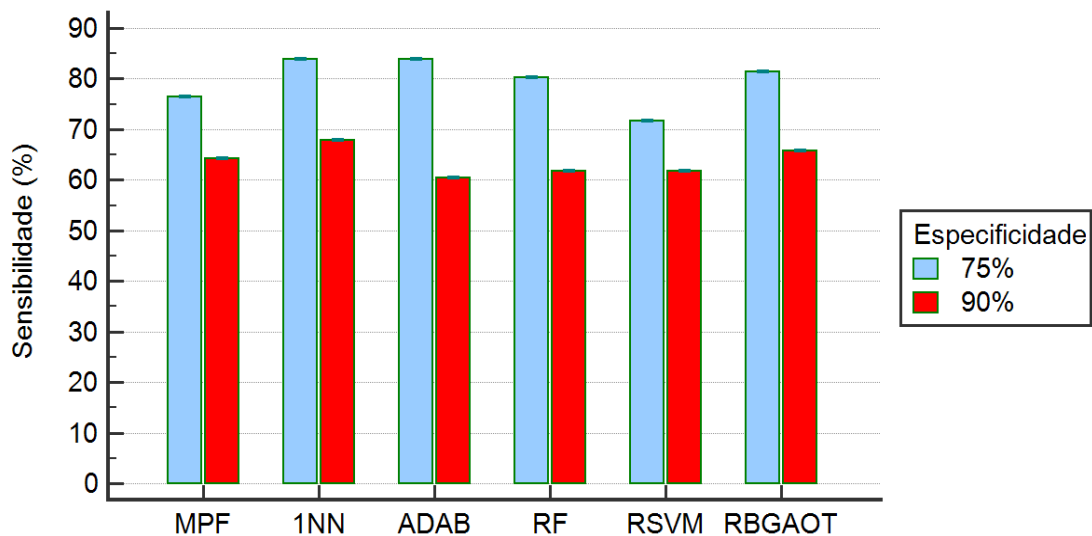
**Tabela 17 – Comparação dos valores de AUC dos classificadores no experimento com todos os atributos da FOT**

	<b>1-NN</b>	<b>ADAB</b>	<b>RF</b>	<b>RSVM</b>	<b>RBGAOT</b>
<b>MPF</b>	0,02±0,035	0,01±0,031	0,01±0,028	0,03±0,029	0,03±0,036
<b>1-NN</b>	-	0,01±0,022	0,01±0,024	0,05±0,035	0,01±0,037
<b>ADAB</b>	-	-	0,003±0,014	0,04±0,033	0,02±0,038
<b>RF</b>	-	-	-	0,03±0,031	0,02±0,034
<b>RSVM</b>	-	-	-	-	0,06±0,039

É possível comparar a sensibilidade ao observar a especificidade com valores fixos. Essa análise limita os falsos positivos, fornecendo medidas em situações onde o algoritmo dificilmente erra a classificação dos portadores da doença. Nesse experimento, foram escolhidas especificidades com valores fixados em 75%, representando um valor moderado e 90%, representando alta especificidade (Figura 25). Esses valores limitam respectivamente, em 25 e 10% os casos de falsos positivos.

Com a especificidade fixada em 75%, observou-se melhora no desempenho dos classificadores 1-NN, ADAB, RF e RBGAOT, fazendo com que eles alcançassem valores de sensibilidade acima de 80%. De forma geral, tanto o MPF, quanto os classificadores estão dentro da faixa de sensibilidade moderada (70 a 90%).

Já com a especificidade fixada em 90%, observou-se uma diminuição nos valores da sensibilidade, fazendo com que todos os classificadores ficassem abaixo da faixa moderada. Os melhores resultados a uma especificidade de 75% ocorrem devido à maior tolerância a falsos positivos, se comparado a 90%.



**Figura 25 – Análise da sensibilidade com especificidade em 75% e 90% no experimento com oito atributos**

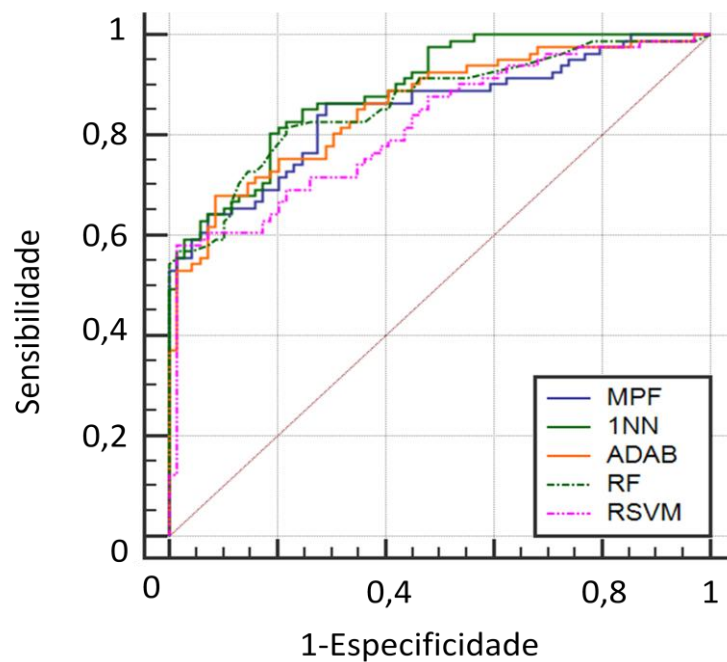
## 5.4. Experimento com Oito Atributos Cruzados

Nesse experimento, o produto cruzado dos oito parâmetros da FOT foi aplicado como entrada dos classificadores. Foram geradas 36 combinações que, somadas à variável *classe*, totalizaram um conjunto com 37 atributos. Para representar possíveis soluções no RBGAOT, são necessárias matrizes de tamanho 37x37. Durante a marginalização da rede, o método *Junction Tree* (BARBER, 2003), fornecido pela *toolbox* PGM, realiza diversos processos que geram um alto custo computacional, fazendo com que o algoritmo não conseguisse convergir. No caso dos demais algoritmos, não houve falhas e o experimento pode ser realizado. As combinações geradas nesse experimento estão disponíveis no Apêndice A.

No algoritmo 1-NN houve um aumento no valor da AUC de 0,87 para 0,89 com relação ao experimento anterior. Os demais classificadores não mostraram mudança no valor da AUC, porém foram observadas alterações nos valores da sensibilidade e especificidade, conforme Tabela 18. Com base nesses valores, a curva ROC pode ser calculada para os quatro classificadores que finalizaram o experimento, além do melhor parâmetro da FOT. Por meio da Figura 26, observa-se que a curva ROC dos classificadores é maior no final dos eixos, pois esse trecho possibilita maior quantidade de falsos positivos.

**Tabela 18 – Resultado do experimento com oito atributos cruzados submetidos aos classificadores**

	Sens (%)	Esp (%)	AUC	E <sub>AUC</sub>
<b>MPF</b>	86,42 (77,0-93,0)	71,01 (58,8-81,3)	0,85 (0,78-0,90)	0,03
<b>1-NN</b>	80,25 (69,9-88,3)	81,16 (69,9-89,6)	0,89 (0,83-0,94)	0,03
<b>ADAB</b>	67,90 (56,6-77,8)	91,30 (82,0-96,7)	0,86 (0,79-0,91)	0,03
<b>RF</b>	81,48 (71,3-89,2)	78,26 (66,7-87,3)	0,86 (0,80-0,91)	0,03
<b>RSVM</b>	58,02 (46,5-68,9)	98,55 (92,2-100,0)	0,82 (0,75-0,86)	0,03
<b>RBGAOT</b>	-	-	-	-



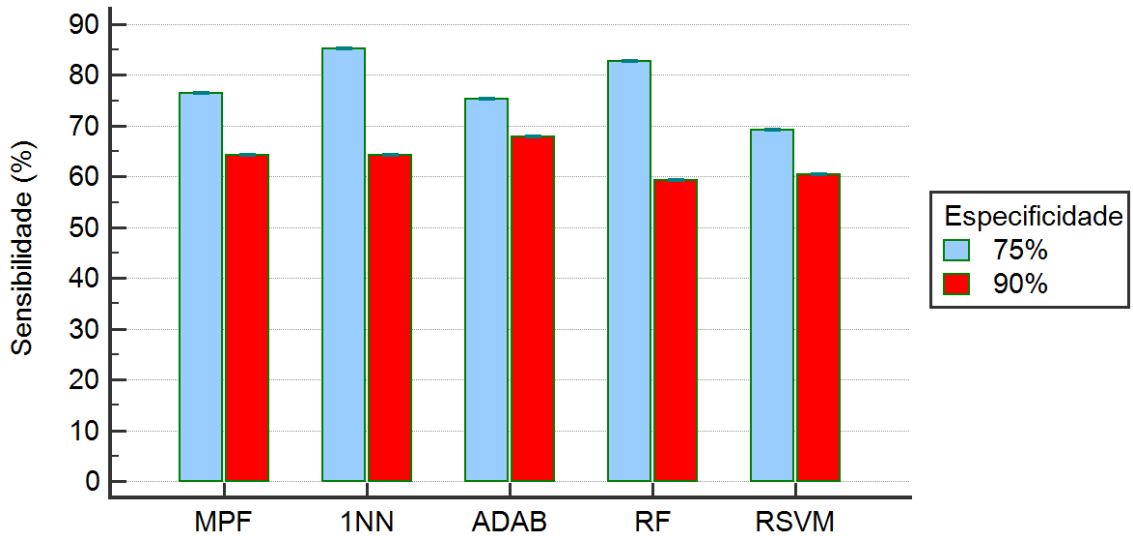
**Figura 26 – Curvas ROC do experimento com oito atributos cruzados**

As curvas ROC obtidas pelos quatro algoritmos foram comparadas, conforme a Tabela 19. A diferença e o erro padrão de cada par foram calculados em um intervalo de confiança de 95% (DELONG et al, 1988). Dentre os dez pares, não foram encontradas diferenças significativas ( $p > 0,05$ ).

**Tabela 19 – Comparação dos valores da AUC dos classificadores no experimento com oito atributos cruzados**

	<b>1-NN</b>	<b>ADAB</b>	<b>RF</b>	<b>RSVM</b>
<b>MPF</b>	0,04±0,03	0,01±0,03	0,02±0,026	0,03±0,031
<b>1-NN</b>	-	0,03±0,024	0,03±0,025	0,07±0,031
<b>ADAB</b>	-	-	0,007±0,014	0,04±0,028
<b>RF</b>	-	-	-	0,05±0,028

A Figura 27 mostra a comparação da sensibilidade com valores de especificidade em 75% e 90%. No primeiro caso, os classificadores 1-NN e RF apresentaram valores de sensibilidade acima de 80%, já o RSVM ficou abaixo da faixa moderada (70 a 90%). Com a especificidade a 90%, todos os classificadores apresentaram valores abaixo de 70%, assim como no experimento anterior.



**Figura 27 – Análise da sensibilidade com valores de especificidade em 75% e 90% no experimento com oito atributos cruzados**

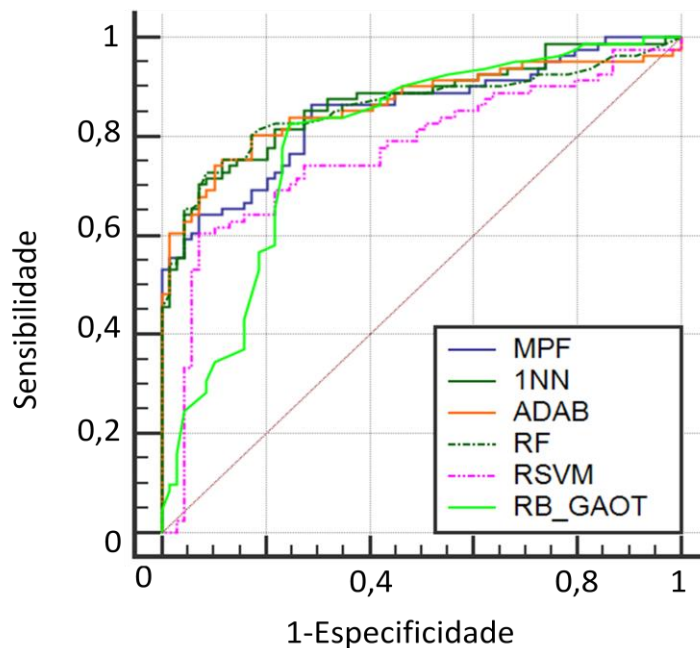
### 5.5. Experimento com Cinco Atributos Selecionados

Na tentativa de melhorar o desempenho dos classificadores, cinco atributos fornecidos pela FOT foram selecionados para o experimento. Essa seleção foi feita pela aplicação da função *featsel* que utiliza o algoritmo 1-NN para classificar os atributos do conjunto de dados. A validação cruzada pelo método *leave-one-out* também foi usada para a seleção dos melhores atributos, sendo, a cada iteração, uma amostra escolhida para teste e as demais para treino. O critério para seleção é o resultado dessa classificação. Da mesma forma, foi testado o uso da AUC como critério de avaliação, porém as variáveis selecionadas não apresentaram resultado superior às selecionadas pelo 1-NN, que foram:  $R_o$ ,  $R_m$ ,  $X_m$ ,  $C_{din}$  e  $Z_{4Hz}$ . Essas cinco variáveis também estão de acordo com a seleção feita por um especialista.

De acordo com a Tabela 20, o algoritmo 1-NN apresentou melhor desempenho com AUC no valor de 0,87. Já o RSVM e o algoritmo RBGAOT apresentaram desempenho inferior com AUC iguais a 0,77 e 0,79, respectivamente. Os algoritmos ADAB e RF permaneceram com AUC iguais a 0,86, porém com alterações nos valores da sensibilidade e especificidade. A Figura 28 mostra a curva ROC com o desempenho dos cinco classificadores, além do melhor parâmetro da FOT. Observa-se que na faixa final dos eixos, onde maior quantidade de falsos positivos é aceita, a curva ROC desses classificadores é maior.

**Tabela 20 – Resultado do experimento com a seleção de cinco atributos submetidos aos classificadores**

	<b>Sens (%)</b>	<b>Esp (%)</b>	<b>AUC</b>	<b>E<sub>AUC</sub></b>
<b>MPF</b>	86,42 (77,0-93,0)	71,01 (58,8-81,3)	0,85 (0,78-0,90)	0,03
<b>1-NN</b>	70,37 (59,2-80,0)	92,75 (83,9-97,6)	0,87 (0,81-0,92)	0,03
<b>ADAB</b>	74,07 (63,1-83,2)	89,86 (80,2-95,8)	0,86 (0,80-0,91)	0,03
<b>RF</b>	72,84 (61,8-82,1)	91,30 (82,0-96,7)	0,86 (0,79-0,91)	0,03
<b>RSVM</b>	60,49 (49,0-71,2)	92,75 (83,9-97,6)	0,77 (0,69-0,83)	0,04
<b>RBGAOT</b>	82,72 (72,7-90,2)	75,36 (63,5 - 84,9)	0,79 (0,72-0,86)	0,04



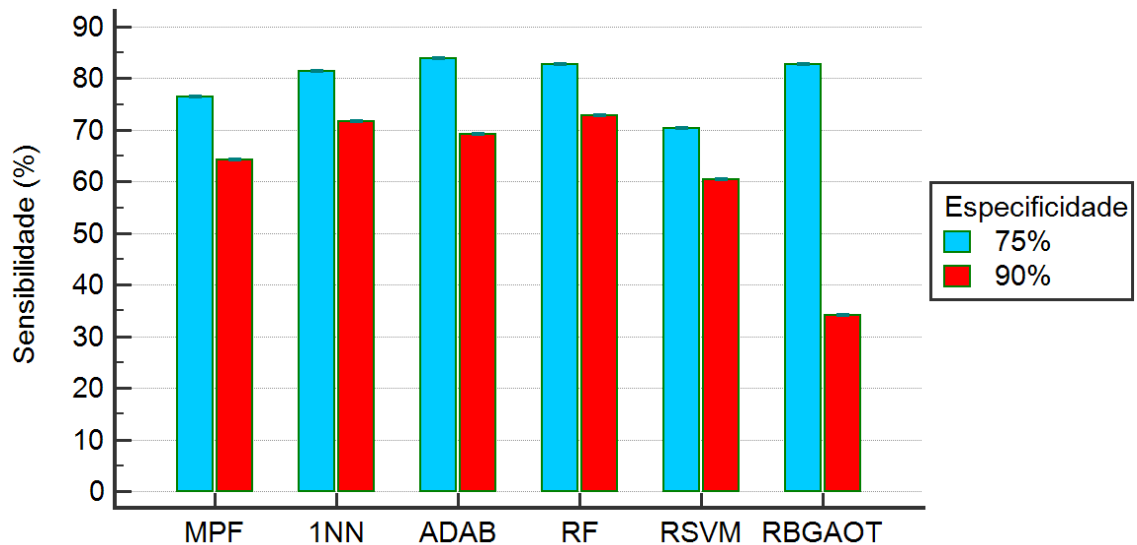
**Figura 28 – Curvas ROC do experimento com seleção de atributos da FOT**

A Tabela 21 mostra a diferença entre os valores das curvas ROC dos cinco algoritmos testados e o MPF. Os cálculos foram feitos em pares, com os erros padrões em um intervalo de confiança de 95% (DELONG et al, 1988). Dentre as quinze áreas comparadas, não foram encontradas diferenças significativas ( $p > 0,05$ ).

**Tabela 21 – Comparação dos valores de AUC dos classificadores no experimento com seleção de atributos da FOT**

	<b>1-NN</b>	<b>ADAB</b>	<b>RF</b>	<b>RSVM</b>	<b>RBGAOT</b>
<b>MPF</b>	0,02±0,033	0,01±0,033	0,01±0,032	0,08±0,036	0,06±0,042
<b>1-NN</b>	-	0,01±0,020	0,01±0,021	0,10±0,037	0,08±0,047
<b>ADAB</b>	-	-	0,003±0,012	0,10±0,037	0,07±0,046
<b>RF</b>	-	-	-	0,09±0,035	0,07±0,045
<b>RSVM</b>	-	-	-	-	0,03±0,051

A análise da sensibilidade realizada com valores fixados em 75% e 90% de especificidade pode ser feita na Figura 29. Em 75%, todos os classificadores apresentaram sensibilidade acima de 80%, exceto o RSVM. Já com especificidade de 90%, os algoritmos 1-NN e RF conseguiram apresentar valores na faixa de sensibilidade de 71,60% e 72,84% respectivamente, estando os demais modelos abaixo da faixa moderada (70 a 90%).



**Figura 29 – Análise da sensibilidade com valores de especificidade em 75% e 90% no experimento com seleção de atributos da FOT**



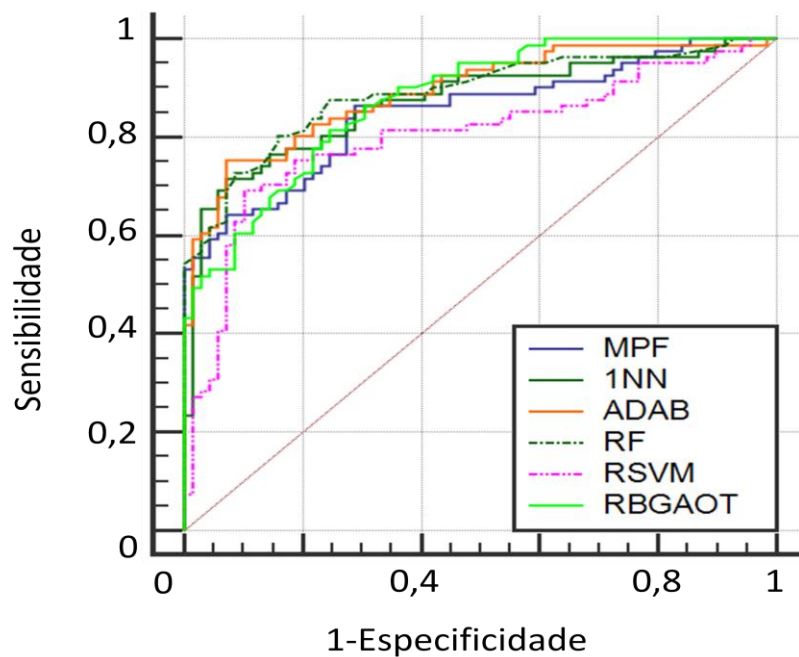
## 5.6. Experimento com Cinco Atributos Cruzados

Nesse experimento, o produto cruzado dos atributos selecionados no item 5.5, foram usados como entrada nos classificadores. Após a aplicação desse método, foram obtidas 15 combinações que, somadas a variável *classe*, totalizaram 16 variáveis de entrada para os cinco algoritmos. Nesse caso, para representar uma possível solução no RBGAOT é necessário o uso de uma matriz de tamanho 16x16. Com esse valor, foi possível passar pelo algoritmo *Junction Tree* e realizar a marginalização das redes, sem que o RBGAOT apresentasse falhas. As combinações geradas nesse experimento também estão disponíveis no Apêndice A.

Os algoritmos RF e ADAB apresentaram melhor desempenho com valores de AUC iguais a 0,89. Mesmo com a aplicação do produto cruzado, o classificador 1-NN apresentou mudança nos valores da sensibilidade e especificidade, porém não apresentou mudança no valor da AUC. No caso do RSVM e do RBGAOT, os valores da AUC aumentaram respectivamente para 0,80 e 0,88, em relação ao experimento anterior. A Figura 30 mostra as curvas ROC calculadas para os classificadores e o MPF, com base nos valores da Tabela 22. Através dessas curvas é possível observar que a curva ROC de todos os classificadores é maior na faixa final dos eixos, onde a quantidade de falsos positivos é maior.

**Tabela 22 – Resultado do experimento com produto cruzado dos atributos selecionados da FOT submetidos aos classificadores**

	Sens (%)	Esp (%)	AUC	E <sub>AUC</sub>
<b>MPF</b>	86,42 (77,0-93,0)	71,01 (58,8-81,3)	0,85 (0,78-0,90)	0,03
<b>1-NN</b>	71,60 (60,5-81,1)	92,75 (83,9-97,6)	0,87 (0,81-0,92)	0,03
<b>ADAB</b>	69,14 (57,9-78,9)	92,75 (83,9-97,6)	0,89 (0,83-0,94)	0,03
<b>RF</b>	82,72 (72,7-90,2)	82,61 (71,6-90,7)	0,89 (0,83-0,94)	0,03
<b>RSVM</b>	83,95 (74,1-91,2)	72,46 (60,4-82,5)	0,80 (0,73-0,86)	0,03
<b>RBGAOT</b>	56,79 (45,3-67,8)	98,55 (92,2-100,0)	0,88 (0,81-0,92)	0,03



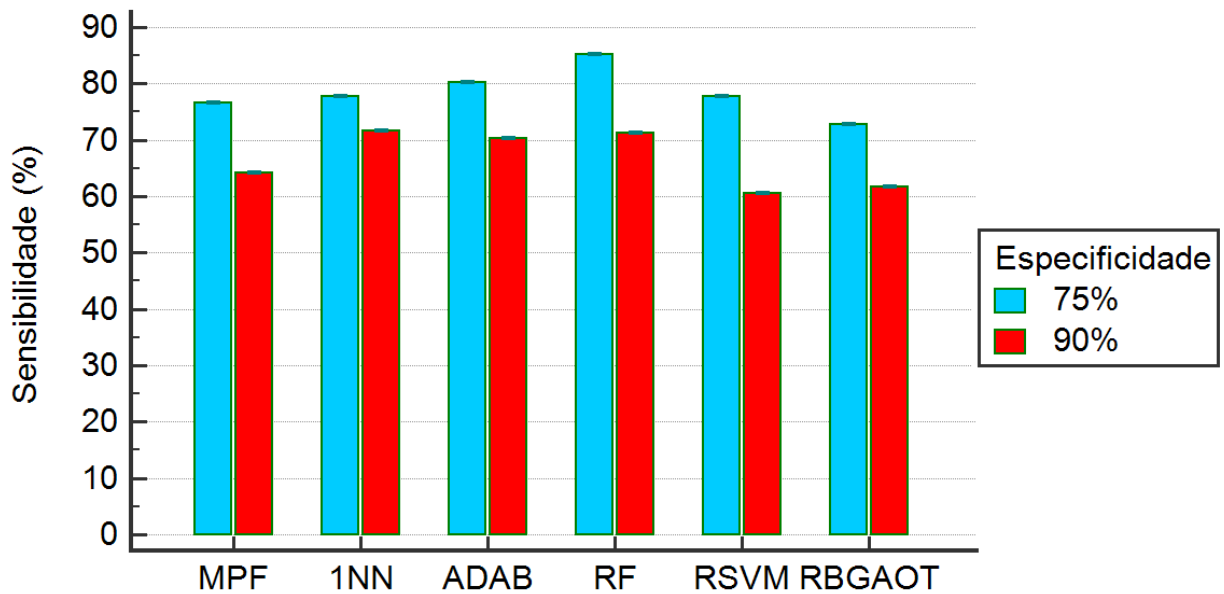
**Figura 30 – Curvas ROC do experimento com cinco parâmetros da FOT cruzados**

A diferença entre os valores da AUC dos classificadores e seu respectivo erro padrão calculado em um intervalo de confiança de 95% (DELONG et al, 1988), podem ser observados na Tabela 23. Dentre os valores encontrados, não foram observadas diferenças significativas ( $p > 0,05$ ).

**Tabela 23 – Comparação dos valores de AUC dos classificadores no experimento com cinco parâmetros da FOT cruzados**

	<b>1-NN</b>	<b>ADAB</b>	<b>RF</b>	<b>RSVM</b>	<b>RBGAOT</b>
<b>MPF</b>	0,02±0,041	0,04±0,038	0,04±0,039	0,05±0,049	0,03±0,041
<b>1-NN</b>	-	0,01±0,022	0,01±0,022	0,04±0,027	0,03±0,037
<b>ADAB</b>	-	-	0,02±0,017	0,03±0,031	0,01±0,040
<b>RF</b>	-	-	-	0,05±0,027	0,04±0,038
<b>RSVM</b>	-	-	-	-	0,02±0,043

A análise dos valores de sensibilidade com a especificidade fixada em 75% e 90% é mostrada na Figura 31. No primeiro caso, observou-se que a sensibilidade dos classificadores apresentou valores na faixa de sensibilidade moderada (70 a 90%), sendo que o RF chegou a 85%. No caso da especificidade em 90%, os algoritmos 1-NN, ADAB e RF conseguiram alcançar valores de sensibilidade acima de 70%, mesmo nesse caso com restrição de falsos positivos a 10%.



**Figura 31 – Análise da sensibilidade com valores de especificidade em 75% e 90% no experimento com cinco parâmetros da FOT cruzados**

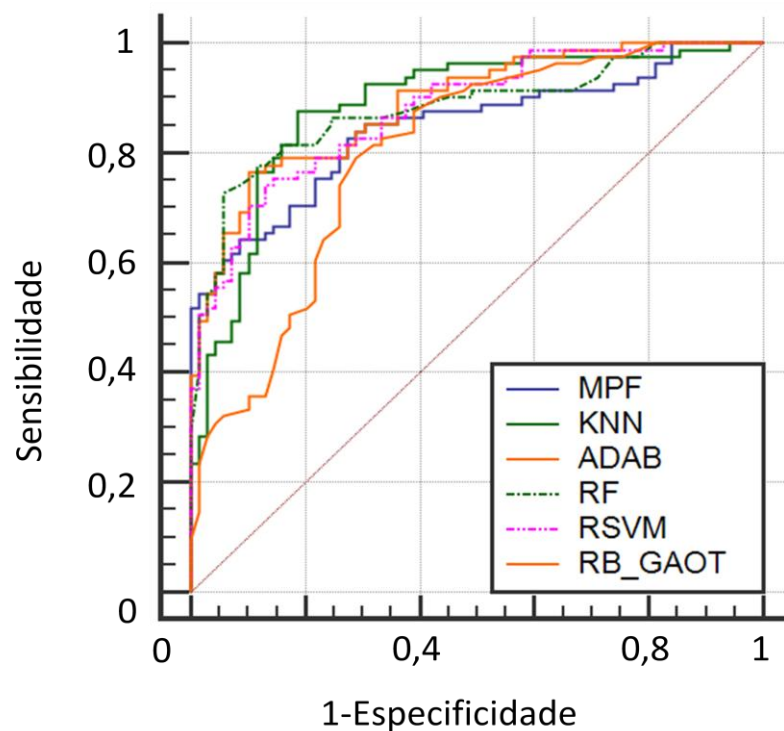
## 5.7. Experimento com Seleção de Cinco Atributos do Produto Cruzados

Nesse experimento, a seleção de cinco atributos também foi feita pela função *featself*, como no experimento do item 5.5, com base no desempenho do algoritmo 1-NN e no método de validação cruzada *leave-one-out*. A diferença é que, nesse caso, a seleção foi feita após a aplicação do produto cruzado.

De acordo com a Tabela 24, todos os algoritmos apresentaram bom desempenho, sendo o 1-NN e ADAB com valores de AUC iguais a 0,89 e o RF e RSVM com valores de AUC iguais a 0,88. Já o RBGAOT apresentou menor desempenho com AUC igual a 0,80. O maior valor de sensibilidade foi obtido no 1-NN, já os maiores valores de especificidade foram obtidos nos classificadores ADAB e RF. A Figura 32 mostra a curva ROC com o desempenho dos classificadores, além do melhor parâmetro da FOT.

**Tabela 24 – Resultado do experimento com atributos do produto cruzado selecionados e submetidos aos classificadores**

	Sens (%)	Esp (%)	AUC	E <sub>AUC</sub>
<b>MPF</b>	86,42 (77,0-93,0)	71,01 (58,8-81,3)	0,85 (0,78-0,90)	0,03
<b>1-NN</b>	87,65 (78,5-93,9)	81,16 (69,9-89,6)	0,89 (0,82-0,93)	0,03
<b>ADAB</b>	76,54 (65,8-85,2)	89,86 (80,2-95,8)	0,89 (0,83-0,94)	0,03
<b>RF</b>	72,84 (61,8-82,1)	94,20 (85,8-98,4)	0,88 (0,82-0,93)	0,03
<b>RSVM</b>	74,07 (63,1-83,2)	86,96 (76,7-93,9)	0,88 (0,83-0,93)	0,03
<b>RBGAOT</b>	79,01 (68,5-87,3)	71,01 (58,8-81,3)	0,80 (0,73-0,86)	0,04



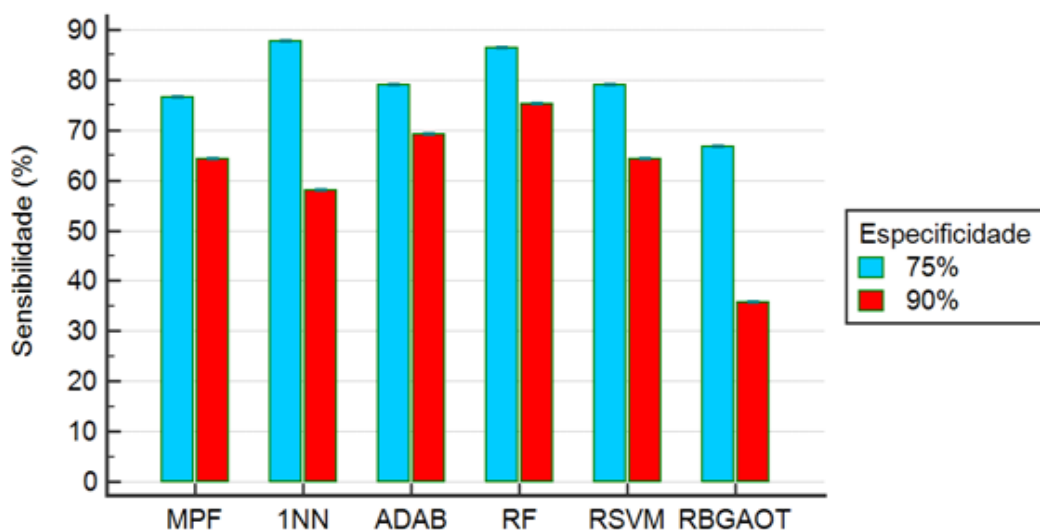
**Figura 32 – Curvas ROC do experimento com atributos do produto cruzado selecionados**

A Tabela 25 mostra a diferença entre os valores das curvas ROC dos algoritmos testados e o melhor parâmetro da FOT. Os cálculos foram feitos em pares, com os erros padrões em um intervalo de confiança de 95% (DELONG et al, 1988). Dentre as quinze áreas comparadas, houve diferença significativa entre o 1-NN e o RBGAOT ( $p < 0,05$ ), e entre ADAB e o RBGAOT ( $p < 0,05$ ). Nas demais combinações não foram encontradas diferenças ( $p > 0,05$ ).

**Tabela 25 – Comparação dos valores de AUC dos classificadores no experimento com atributos do produto cruzado selecionados**

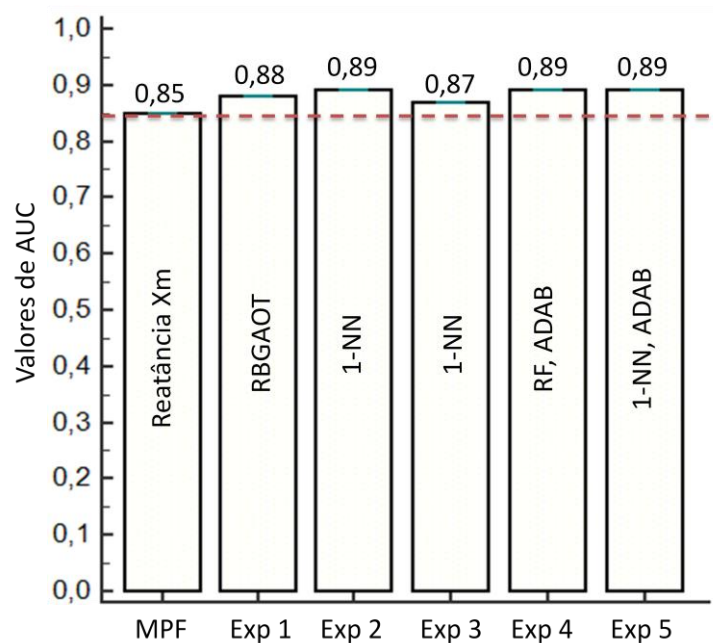
	1-NN	ADAB	RF	RSVM	RBGAOT
MPF	0,04±0,035	0,04±0,028	0,04±0,028	0,03±0,031	0,05±0,045
1-NN	-	0,003±0,023	0,004±0,023	0,008±0,024	0,09±0,043
ADAB	-	-	0,008±0,013	0,01±0,021	0,09±0,043
RF	-	-	-	0,004±0,023	0,08±0,045
RSVM	-	-	-	-	0,08±0,045

A análise da sensibilidade feita com valores fixados em 75% e 90% de especificidade pode ser vista na Figura 33. Em 75%, todos os classificadores apresentaram sensibilidade na faixa moderada, exceto o RBGAOT, tendo o 1-NN e o RF alcançado valores de sensibilidade acima de 80%. Já com especificidade em 90%, apenas o algoritmo RF conseguiu apresentar valor de sensibilidade acima de 70%, estando os demais modelos abaixo da faixa moderada (70 a 90%).



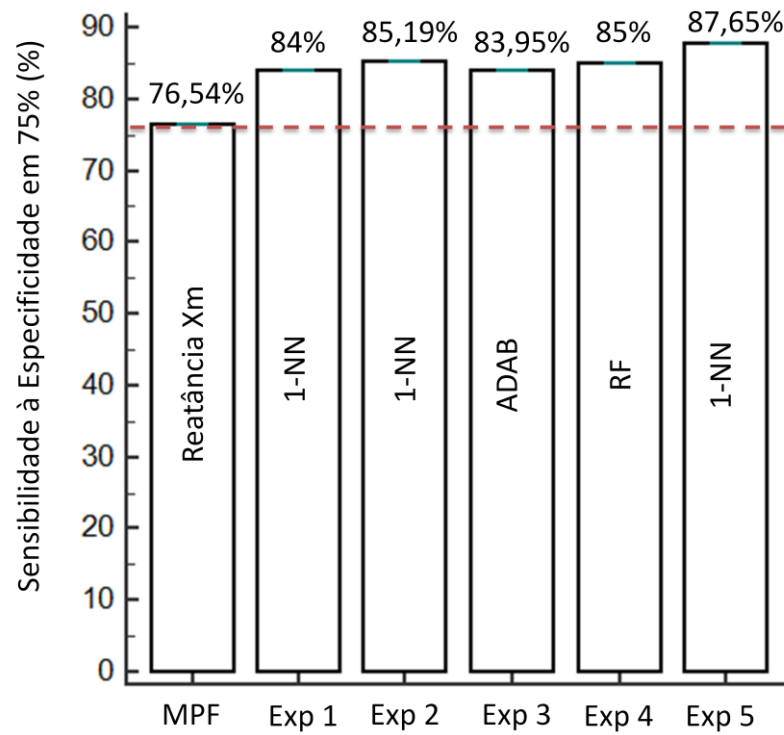
**Figura 33 – Análise da sensibilidade com valores de especificidade em 75% e 90% no experimento com atributos do produto cruzado selecionado**

A Figura 34 mostra um resumo com os classificadores que apresentaram melhor desempenho em cada um dos experimentos realizados, comparando-os com o melhor parâmetro da FOT. Durante os testes, os valores de AUC dos melhores algoritmos variaram entre 0,87 e 0,89 (acurácia moderada). Esses resultados mostram que o uso de algoritmos de aprendizado de máquinas teve melhor desempenho do que a classificação com a melhor característica da FOT, a reatância  $X_m$ .

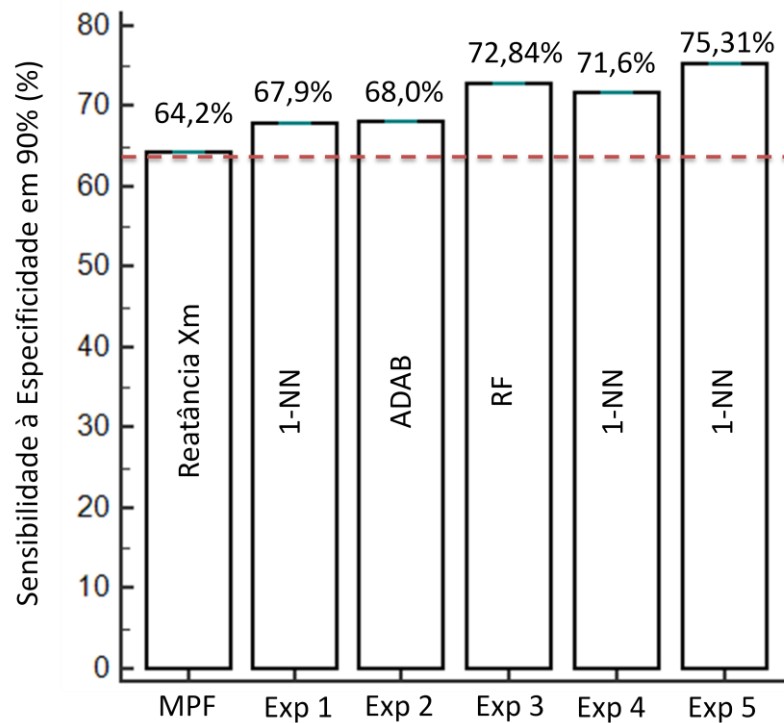


**Figura 34 – Resumo dos maiores valores de AUC obtidos durante os experimentos**

Os gráficos da Figura 35 e Figura 36 apresentam um resumo com o desempenho dos melhores algoritmos de cada experimento, usando como base a especificidade fixada em 75% e 90%, respectivamente. No primeiro caso, os valores de sensibilidade encontrados foram superiores, se comparados ao melhor parâmetro da FOT, variando de 83,95% a 85,19%. Já no segundo caso, onde é simulada uma situação mais restrita para o algoritmo, os valores de sensibilidade variaram entre 67,9% e 72,84%. Mesmo com a especificidade em 90%, simulando uma situação onde o algoritmo aceita apenas 10% de casos de falsos positivos, o uso de algoritmos de aprendizado de máquinas fez com que no terceiro e quarto experimentos com variáveis selecionadas, houvesse valores de sensibilidade dentro da faixa moderada.



**Figura 35 - Resumo dos maiores valores de sensibilidade com especificidade fixada em 75%, obtidos durante os experimentos**



**Figura 36 – Resumo dos maiores valores de sensibilidade com especificidade fixada em 90%, obtidos durante os experimentos**

## 5.8. Inferência sobre Redes Bayesianas

As Redes Bayesianas fornecem grafos que mostram as ligações de dependência entre as variáveis do problema. Sendo assim, as redes construídas com base em matrizes, conforme item 4.6, que representam as melhores soluções geradas, são selecionadas pelo RBGAOT. Além de apresentar boa acurácia, essas estruturas devem permitir uma análise gráfica das relações existentes entre as características fornecidas pela FOT.

Algumas dessas redes foram selecionadas para análise com base na quantidade de ligações existentes entre suas variáveis. Logo, quanto menor o número de arcos existentes entre os nós da rede, mais simples será sua representação e, conseqüentemente, mais simples serão suas tabelas de distribuição de probabilidade conjunta (DPC). Nos itens 5.8.1 e 5.8.2 foram realizadas inferências sobre as redes selecionadas de acordo com o número de atributos de entrada usados em suas respectivas construções.

### 5.8.1. Rede com Oito Atributos

Com o intuito de analisar algumas das redes geradas pelo RBGAOT com todos os atributos da FOT, foi feita uma seleção com base em dois parâmetros. Primeiramente observou-se o menor número de ligações entre as variáveis, resultando em tabelas de DPC mais simples. Em seguida, foram selecionadas redes cujas tabelas de DPC apresentaram menor ocorrência de probabilidades iguais a 0,5, já que esse valor não permite inferir de forma mais precisa sobre uma situação. Como as tabelas de DPC são construídas de forma automática, elas produzem combinações que não condizem com a biomecânica. A fim de destacar apenas as combinações que descrevem situações possíveis de ocorrer, essas probabilidades foram destacadas em cada uma das tabelas.

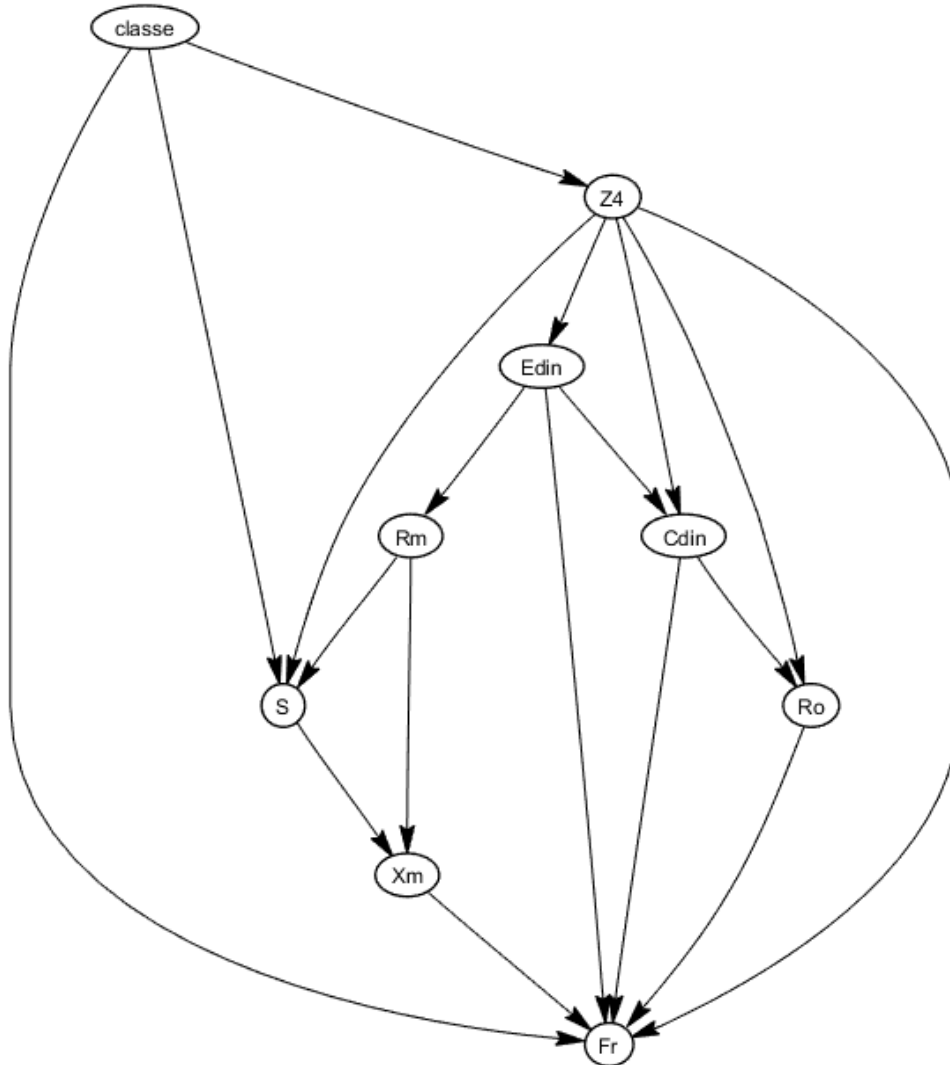
A rede representada na Figura 37 tem estrutura composta pelos oito atributos fornecidos pela FOT:  $R_o$ ,  $R_m$ ,  $X_m$ ,  $C_{din}$ ,  $S$ ,  $Z_{4Hz}$ ,  $F_r$  e  $E_{din}$ , além da variável *classe*. Essa estrutura apresenta nove tabelas de DPC, cujas informações foram analisadas e comparadas com os gráficos da Figura 22.

A Tabela 26 mostra as probabilidades à priori da variável *classe*, o único nó raiz dessa rede. Por essas probabilidades, é possível observar que o conjunto de dados usado para a construção da estrutura é formado por 45% de indivíduos do grupo controle e 55% do grupo teste.



**Tabela 26 – Probabilidades à priori da variável *classe* com oito atributos de entrada**

$P(\text{classe} = 0)$	$P(\text{classe} = 1)$
0,45	0,55



**Figura 37 – Estrutura da rede com oito atributos de entrada**

A variável  $Z_{4Hz}$  recebe influência apenas da variável *classe*. Pela Tabela 27, pode-se perceber que há alta probabilidade de um indivíduo ter baixa impedância ( $Z_{4Hz}=1$ ) dado que não é portador da doença (*classe*=0). Já um indivíduo portador de fibrose cística (*classe*=1), tem maior probabilidade de apresentar alta impedância ( $Z_{4Hz} = 2$ ). Comparando com o gráfico da Figura 22(f), em média, a impedância  $Z_{4Hz}$  é de fato menor para indivíduos que não possuem fibrose cística (grupo controle) e maior para indivíduos que possuem a doença (grupo teste).

**Tabela 27 – DPC para a variável  $Z_{4Hz}$  da rede com oito atributos de entrada**

	$P(Z_{4Hz} = 1 classe)$	$P(Z_{4Hz} = 2 classe)$
<i>classe</i> = 0	0,93	0,07
<i>classe</i> = 1	0,34	0,66

A variável  $R_m$  é influenciada apenas pela elastância ( $E_{din}$ ) e, de acordo com a Tabela 28, possui comportamento diretamente proporcional a ela. Ao observar um indivíduo com baixa elastância, há probabilidade de 0,91 desse indivíduo também apresentar baixa resistência  $R_m$ .

**Tabela 28 – DPC para a variável  $R_m$  da rede com oito atributos de entrada**

	$P(R_m = 1 E_{din})$	$P(R_m = 2 E_{din})$
$E_{din} = 1$	0,91	0,09
$E_{din} = 2$	0,34	0,66

Pelos gráficos da Figura 22(b) e da Figura 22(h), tanto a resistência  $R_m$  como a elastância  $E_{din}$ , possuem valores mais baixos para indivíduos do grupo controle e valores mais altos para indivíduos do grupo teste. Logo, em média, ambos os parâmetros apresentam comportamento similar, conforme indicado pela Tabela 28.

A variável  $E_{\text{din}}$  também possui comportamento diretamente proporcional à impedância  $Z_{4\text{Hz}}$ . Pela Tabela 29 pode-se observar que no caso de  $Z_{4\text{Hz}}=1$ , a elastância mostra alta probabilidade de ser baixa ( $E_{\text{din}}=1$ ). De forma similar, quando um indivíduo possui  $Z_{4\text{Hz}}=2$ , há alta probabilidade de ter  $E_{\text{din}}=2$ .

**Tabela 29 – DPC para a variável  $E_{\text{din}}$  da rede com oito atributos de entrada**

	$P(E_{\text{din}}=1   Z_{4\text{Hz}})$	$P(E_{\text{din}}=2   Z_{4\text{Hz}})$
$Z_{4\text{Hz}} = 1$	0,95	0,05
$Z_{4\text{Hz}} = 2$	0,14	0,86

Os gráficos da Figura 22(f) e da Figura 22(h) mostram que a elastância  $E_{\text{din}}$  e a impedância  $Z_{4\text{Hz}}$ , em média, também apresentam um comportamento equivalente, concordando assim, com as informações da Tabela 29.

A variável  $X_m$  recebe influências de outras duas variáveis:  $S$  e  $R_m$ . A reatância  $X_m$  mostra comportamento diretamente proporcional à inclinação da curva de resistência  $S$ . Analisando as condições a seguir retiradas da Tabela 30, é possível comprovar que há alta probabilidade de um indivíduo ter  $X_m$  com valor baixo, tendo observado  $S$  com valor baixo:

$$P(X_m = 1 | S = 1, R_m = 1) = 0,76$$

$$P(X_m = 1 | S = 1, R_m = 2) = 0,87$$

De forma similar, há maior probabilidade em observar alta reatância  $X_m$ , tendo observado uma maior inclinação  $S$  da curva de resistência, independentemente do valor de  $R_m$ :

$$P(X_m = 2 | S = 2, R_m = 1) = 0,96$$

$$P(X_m = 2 | S = 2, R_m = 2) = 0,55$$

Apresentando ao algoritmo a combinação improvável onde  $S=2$  (característica do grupo controle) e  $R_m=2$  (característica do grupo teste), observa-se a dificuldade do RBGAOT calcular um valor que favoreça uma discriminação mais clara, resultando em uma probabilidade próxima a 0,5, conforme a linha 4 da Tabela 30:

**Tabela 30 – DPC para a variável  $X_m$  da rede com oito atributos de entrada**

	$P(X_m = 1   S, R_m)$	$P(X_m = 2   S, R_m)$
$S = 1, R_m = 1$	0,76	0,24
$S = 2, R_m = 1$	0,05	0,95
$S = 1, R_m = 2$	0,87	0,13
$S = 2, R_m = 2$	0,45	0,55

As variáveis  $X_m$  e  $S$ , que possuem valores baixos, geralmente caracterizam indivíduos do grupo teste, conforme Figura 22(c) e Figura 22(e). Esse comportamento também é observado na Tabela 30.

As variáveis  $C_{din}$  e  $Z_{4Hz}$  exercem influência sobre a variável  $R_o$ , de acordo com a rede da Figura 37. Há maior probabilidade em observar baixa resistência  $R_o$ , dado que foram observados baixos valores de  $Z_{4Hz}$ , independente do valor de  $C_{din}$ . O mesmo comportamento é observado para altas resistências, conforme Tabela 31:

**Tabela 31 – DPC para a variável  $R_o$  da rede com oito atributos de entrada**

	$P(R_o = 1   C_{din}, Z_{4Hz})$	$P(R_o = 2   C_{din}, Z_{4Hz})$
$C_{din} = 1, Z_{4Hz} = 1$	0,64	0,36
$C_{din} = 2, Z_{4Hz} = 1$	0,97	0,03
$C_{din} = 1, Z_{4Hz} = 2$	0,19	0,81
$C_{din} = 2, Z_{4Hz} = 2$	0,34	0,66

Analisando as informações dos gráficos da Figura 22(a) e Figura 22(f), a impedância  $Z_{4Hz}$  e a resistência  $R_o$  de fato apresentam valores baixos para o grupo controle e valores altos para o grupo teste. Portanto, com base na biomecânica, as linhas 1 e 4 da Tabela 31 descrevem situações difíceis de ocorrer. Isso justifica valores de probabilidades mais próximos, mostrando maior indecisão.

A variável  $C_{din}$  é influenciada por  $E_{din}$  e  $Z_{4Hz}$ . Ao analisar a Tabela 32, obtém-se probabilidades de 0,98 e 0,81 do indivíduo possuir valor de complacência maior ( $C_{din}=2$ ), dado que foi observada baixa elastância ( $E_{din}=1$ ). O mesmo comportamento ocorre quando  $C_{din}=1$ . Pode-se concluir que de acordo com essa rede,  $C_{din}$  é inversamente proporcional a  $E_{din}$ .

No caso das linhas 2 e 3 na Tabela 32, há duas situações improváveis de ocorrer, já que a impedância é diretamente proporcional à elastância. Porém, apenas no caso onde  $E_{din} = 2$  e  $Z_{4Hz} = 1$ , foi encontrada uma probabilidade menor, pois no caso onde  $E_{din} = 1$  e  $Z_{4Hz} = 2$ , o valor da probabilidade continuou alto.

**Tabela 32 – DPC para a variável  $C_{din}$  da rede com oito atributos de entrada**

	$P(C_{din}=1 E_{din}, Z_{4Hz})$	$P(C_{din}=2 E_{din}, Z_{4Hz})$
$E_{din} = 1, Z_{4Hz} = 1$	0,02	0,98
$E_{din} = 2, Z_{4Hz} = 1$	0,65	0,35
$E_{din} = 1, Z_{4Hz} = 2$	0,19	0,81
$E_{din} = 2, Z_{4Hz} = 2$	0,97	0,03

Esse comportamento inversamente proporcional da complacência e da elastância pode ser observado também nos gráficos da Figura 22(d) e Figura 22(h), comprovando assim, as informações obtidas na Tabela 32.

A variável  $S$  é influenciada pelas variáveis  $R_m$ ,  $Z_{4Hz}$  e *classe*. De acordo com as combinações a seguir, há maior probabilidade de um indivíduo ter inclinação da curva de resistência com alto valor ( $S=2$ ) se ele não for portador da doença (*classe*=0):

$$P(S = 2|R_m = 1, Z_{4Hz} = 1, classe = 0) = 0,96$$

$$P(S = 2|R_m = 2, Z_{4Hz} = 1, classe = 0) = 0,72$$

$$P(S = 2|R_m = 1, Z_{4Hz} = 2, classe = 0) = 0,72$$

Mesmo em situações improváveis, representadas nas linhas 2 e 3, as informações da Tabela 33 concordam com o comportamento da inclinação da curva de resistência, descrito na Figura 22 (e), onde o valor alto de  $S$  é uma característica do grupo controle. Já um valor baixo de  $S$ , caracteriza indivíduos portadores de fibrose cística.

Dentre as combinações improváveis contidas na Tabela 33, há duas em que o algoritmo calculou probabilidades iguais a 0,5. A primeira é a  $P(S | R_m=2, Z_{4Hz}=2, classe=0)$ , que descreve um indivíduo com valores altos de resistência e impedância, e mesmo assim não é portador da doença ( $classe=0$ ). A segunda é  $P(S | R_m=2, Z_{4Hz} =1, classe=1)$ , que mostra um indivíduo com baixa impedância e portador da doença ( $classe=1$ ). Geralmente, isto se deve ao fato do algoritmo gerar todas as combinações possíveis, até mesmo situações como essas representadas nas linhas 4 e 6, onde o RBGAOT não foi capaz de inferir e calcular as probabilidades com valores distantes.

A linha cinco da Tabela 33 descreve uma situação onde o indivíduo possui baixos valores de  $R_m$  e  $Z_{4Hz}$ , mas é doente ( $classe=1$ ). Mesmo nessa combinação improvável, o algoritmo calculou uma alta probabilidade do indivíduo ter valor de  $S$  elevado, característica de quem não é portador de fibrose cística. Essa alta probabilidade pode ter sido influenciada pelo fato do indivíduo apresentar uma combinação de duas características de um não portador da doença ( $R_m = 1$  e  $Z_{4Hz} = 1$ ), contra uma característica de quem é portador da fibrose cística ( $classe=1$ ).

**Tabela 33 – DPC para a variável  $S$  da rede com oito atributos de entrada**

	$P(S =1 R_m, Z_{4Hz}, classe)$	$P(S =2 R_m, Z_{4Hz}, classe)$
$R_m = 1, Z_{4Hz} = 1, classe = 0$	0,04	0,96
$R_m = 2, Z_{4Hz} = 1, classe = 0$	0,28	0,72
$R_m = 1, Z_{4Hz} = 2, classe = 0$	0,28	0,72
$R_m = 2, Z_{4Hz} = 2, classe = 0$	0,50	0,50
$R_m = 1, Z_{4Hz} = 1, classe = 1$	0,03	0,97
$R_m = 2, Z_{4Hz} = 1, classe = 1$	0,50	0,50
$R_m = 1, Z_{4Hz} = 2, classe = 1$	0,54	0,46
$R_m = 2, Z_{4Hz} = 2, classe = 1$	0,77	0,23

A tabela de DPC da frequência  $F_r$  que, por ser nó folha dessa rede e ser influenciada por outras seis variáveis, exige uma tabela com  $2^6$  linhas para representar todas as possíveis combinações ligadas a essa variável. Devido à complexidade desse caso, o algoritmo forneceu apenas probabilidades iguais a 0,5.

Pela análise dessa rede com oito atributos, observou-se que a estrutura obtida pelo algoritmo RBGAOT apresentou informações probabilísticas das relações entre os parâmetros da FOT. Em geral, os valores mais altos de probabilidades encontrados nas tabelas de DPC, foram obtidos em combinações coerentes com a biomecânica. Essas informações foram conferidas com os valores médios da Figura 22, e apresentaram o mesmo comportamento descrito por esses gráficos, confirmando a consistência do RBGAOT.

### 5.8.2. Rede com Seleção de Cinco Atributos

Dentre as estruturas geradas com cinco atributos de entrada, uma rede foi selecionada para inferência. Essa escolha também foi feita com base no menor número de ligações entre as variáveis e pela menor ocorrência de probabilidades iguais a 0,5. Os cinco atributos usados na construção das redes foram:  $R_o$ ,  $R_m$ ,  $X_m$ ,  $C_{din}$  e  $Z_{4Hz}$ , além da variável *classe*.

A rede da Figura 38 gerou um total de seis tabelas de DPC, que foram analisadas e comparadas com as informações da Figura 22. A Tabela 34 apresenta as probabilidades à priori da variável *classe*, onde a probabilidade de um indivíduo não ser portador da doença é igual a 0,49. Já a probabilidade do indivíduo ser portador da doença é igual a 0,51.

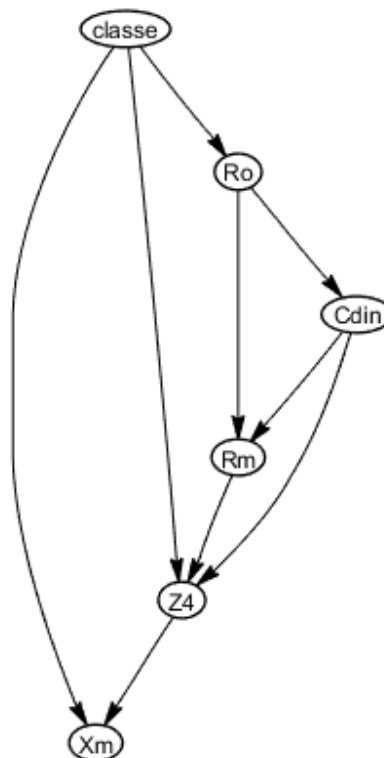


Figura 38 – Estrutura da rede com cinco atributos de entrada

**Tabela 34 – Probabilidades à priori da variável *classe* da rede com cinco atributos de entrada**

$P(\text{classe} = 0)$	$P(\text{classe} = 1)$
0,49	0,51

De acordo com a rede representada na Figura 38, a variável  $R_o$  é influenciada apenas pela variável *classe*. Sendo assim, a Tabela 35 mostra as probabilidades de um indivíduo ter baixa ou alta resistência ( $R_o = 1$  ou  $2$ ), dado que foi observada a sua classe.

**Tabela 35 – DPC para a variável  $R_o$  com cinco atributos de entrada**

	$P(R_o = 1   \text{classe})$	$P(R_o = 2   \text{classe})$
<i>classe</i> = 0	0,94	0,06
<i>classe</i> = 1	0,47	0,53

Usando como base as condições, retiradas da Tabela 35, é possível concluir que há alta chance de um indivíduo apresentar baixo valor de  $R_o$ , dado que pertence a *classe* 0. Também há chances do indivíduo ter alto valor de  $R_o$ , dado que pertence a *classe* 1. Esse comportamento também é observado na Figura 22(a), onde a baixa resistência é característica do grupo controle, e a alta resistência é característica do grupo teste.

$$P(R_o = 1 | \text{classe} = 0) = 0,94$$

$$P(R_o = 2 | \text{classe} = 1) = 0,53$$

Pela Tabela 35, também foi observado que as probabilidades calculadas para *classe*=1, não possuem diferença grande entre seus valores, como ocorre na *classe*=0. Por apresentarem valores próximos a 0,5, o algoritmo mostra maior dificuldade em discriminar portadores da doença:

$$P(R_o = 1 | \text{classe} = 1) = 0,47$$

$$P(R_o = 2 | \text{classe} = 1) = 0,53$$



A variável  $C_{din}$  depende apenas da variável  $R_o$ . Pela Tabela 36 é possível observar que há uma probabilidade de 0,87 do indivíduo ter o valor de complacência alto ( $C_{din}=2$ ), dado que foi observado um baixo valor de resistência ( $R_o=1$ ). Já os indivíduos que apresentam maior valor de  $R_o$ , possuem probabilidade de 0,84 de ter complacência menor ( $C_{din}=1$ ). Esse comportamento inversamente proporcional da complacência e da resistência, também é visto na Figura 22(a) e na Figura 22(d).

**Tabela 36 – DPC para a variável  $C_{din}$  da rede com cinco atributos de entrada**

	$P(C_{din}=1   R_o)$	$P(C_{din}=2   R_o)$
$R_o = 1$	0,13	0,87
$R_o = 2$	0,84	0,16

No caso da variável  $R_m$ , há duas outras variáveis que a influenciam:  $C_{din}$  e  $R_o$ . Pela Tabela 37, é possível analisar as relações a seguir e perceber que independente do valor de  $C_{din}$ , há probabilidade da resistência  $R_m$  apresentar um valor baixo, dado que também é observado um valor baixo de  $R_o$ :

$$P(R_m = 1 | C_{din} = 1, R_o = 1) = 0,79$$

$$P(R_m = 1 | C_{din} = 2, R_o = 1) = 0,98$$

O mesmo ocorre com a resistência  $R_m$ , dado que  $R_o$  apresenta um valor alto:

$$P(R_m = 2 | C_{din} = 1, R_o = 2) = 0,82$$

$$P(R_m = 2 | C_{din} = 2, R_o = 2) = 0,77$$

A relação diretamente proporcional da resistência média ( $R_m$ ) e da resistência no intercepto ( $R_o$ ) descrita na Tabela 37, concordam com os gráficos da Figura 22(a) e Figura 22(b). Essa relação se mantém mesmo nas inconsistências biomecânicas  $C_{din}=1$  e  $R_m=1$ , ou  $C_{din}=2$  e  $R_o=2$ , como ocorre respectivamente nas linhas 1 e 4.

**Tabela 37 – DPC para a variável  $R_m$  da rede com cinco atributos de entrada**

	$P(R_m=1   C_{din}, R_o)$	$P(R_m=2   C_{din}, R_o)$
$C_{din} = 1, R_o = 1$	0,79	0,21
$C_{din} = 2, R_o = 1$	0,98	0,02
$C_{din} = 1, R_o = 2$	0,18	0,82
$C_{din} = 2, R_o = 2$	0,23	0,77

A variável  $X_m$  é influenciada pelas variáveis  $Z_{4Hz}$  e *classe*. Conforme as probabilidades da Tabela 38, há alta chance de um indivíduo apresentar alta reatância  $X_m$ , dado que  $Z_{4Hz}$  é baixo, independente do valor da variável *classe*. Esse comportamento é observado mesmo na situação improvável, onde há baixa impedância  $Z_{4Hz}$  em um portador de fibrose cística:

$$P(X_m = 2 | Z_{4Hz} = 1, classe = 0) = 0,98$$

$$P(X_m = 2 | Z_{4Hz} = 1, classe = 1) = 0,94$$

Do mesmo modo, há probabilidade do indivíduo apresentar baixa reatância  $X_m$ , dado que foi observado alto valor de impedância  $Z_{4Hz}$ . Mesmo com a situação improvável de um indivíduo não possuir a doença (*classe*=0) e ter alta impedância ( $Z_{4Hz}$ =2), esse comportamento se manteve, apresentando apenas, valor de probabilidade mais baixo:

$$P(X_m = 1 | Z_{4Hz} = 2, classe = 0) = 0,64$$

$$P(X_m = 1 | Z_{4Hz} = 2, classe = 1) = 0,80$$

Os gráficos da Figura 22 (c) e da Figura 22 (f), comprovam o comportamento inversamente proporcional da reatância  $X_m$  e da impedância  $Z_{4Hz}$ , descrito na Tabela 38.

**Tabela 38 – DPC para a variável  $X_m$  da rede com cinco atributos de entrada**

	$P(X_m=1   Z_{4Hz}, classe)$	$P(X_m=2   Z_{4Hz}, classe)$
$Z_{4Hz} = 1, classe = 0$	0,02	0,98
$Z_{4Hz} = 2, classe = 0$	0,64	0,36
$Z_{4Hz} = 1, classe = 1$	0,06	0,94
$Z_{4Hz} = 2, classe = 1$	0,80	0,20

Já a variável  $Z_{4Hz}$  é influenciada por três outras variáveis:  $R_m$ ,  $C_{din}$  e *classe*. As condições a seguir, retiradas da Tabela 39, mostram que há maior probabilidade de um indivíduo apresentar alta impedância  $Z_{4Hz}$ , dado que ele é doente. Pela Figura 22(f), observa-se que de fato indivíduos portadores da doença (*classe*=1) possuem valores altos de impedância.

$$P(Z_{4Hz} = 2 | R_m = 1, C_{din} = 1, classe = 1) = 0,77$$

$$P(Z_{4Hz} = 2 | R_m = 2, C_{din} = 1, classe = 1) = 0,97$$

$$P(Z_{4Hz} = 2 | R_m = 2, C_{din} = 2, classe = 1) = 0,81$$

Do ponto de vista biomecânico, as únicas situações que descrevem combinações prováveis são as linhas 3 e 6. Para esses casos, foram encontrados altos valores de probabilidade:

$$P(Z_{4Hz} = 1 | R_m = 1, C_{din} = 2, classe = 0) = 0,99$$

$$P(Z_{4Hz} = 2 | R_m = 2, C_{din} = 1, classe = 1) = 0,97$$

**Tabela 39 – DPC para a variável  $Z_{4Hz}$  da rede com cinco atributos de entrada**

	$P(Z_{4Hz}=1   R_m, C_{din}, classe)$	$P(Z_{4Hz}=2   R_m, C_{din}, classe)$
$R_m = 1, C_{din} = 1, classe = 0$	0,28	0,72
$R_m = 2, C_{din} = 1, classe = 0$	0,50	0,50
$R_m = 1, C_{din} = 2, classe = 0$	0,99	0,01
$R_m = 2, C_{din} = 2, classe = 0$	0,72	0,28
$R_m = 1, C_{din} = 1, classe = 1$	0,23	0,77
$R_m = 2, C_{din} = 1, classe = 1$	0,03	0,97
$R_m = 1, C_{din} = 2, classe = 1$	0,92	0,08
$R_m = 2, C_{din} = 2, classe = 1$	0,19	0,81

Para o caso de um indivíduo com as características:  $R_m = 2$ ,  $C_{din} = 1$  e *classe* = 0, o algoritmo não conseguiu definir os cálculos das probabilidades, resultando em  $P(Z_{4Hz}|R_m=2,C_{din}=1,classe=0)$  igual a 0,5. Normalmente, isto ocorre em situações improváveis como essa, em que um indivíduo possui alta resistência respiratória ( $R_m=2$ ) e complacência próxima à zero ( $C_{din}=1$ ), porém não possui a doença (*classe*=0).

Assim como no item 5.8.1, a estrutura gerada pelo algoritmo RBGAOT apresentou probabilidades que descrevem as relações existentes entre os parâmetros da FOT. Nas combinações que representam situações coerentes com a biomecânica, também foram obtidos os maiores valores de probabilidades. Essas relações estão de acordo com os gráficos da Figura 22 e podem fornecer informações mesmo quando são submetidas a combinações improváveis, porém com menores valores de probabilidades. A inferência sobre outras redes pode ser lida no Apêndice B.

## CONCLUSÃO

Este projeto seguiu a linha de pesquisa de trabalhos já realizados (AMARAL et al., 2013; AMARAL et al., 2015; AMARAL et al., 2017), com o uso dos dados fornecidos pela FOT em algoritmos de aprendizado de máquinas, mostrando que essa associação também foi eficiente na detecção de alterações respiratórias na fibrose cística. Durante os experimentos, os classificadores: *1-Nearest Neighbor*, *Adaboost*, *Radial Support Vector Machine*, *Random Forest* e Redes Bayesianas, apresentaram valores de AUC maiores do que os valores obtidos pelo melhor parâmetro da FOT, indicando assim, maior acurácia no diagnóstico.

Dentre os testes realizados, a reatância respiratória ( $X_m$ ) foi o atributo que apresentou melhor desempenho individual (AUC=0,85). No primeiro experimento, foram usados oito atributos de entrada fornecidos pela FOT, sendo o RBGAOT o algoritmo que apresentou melhor resultado (AUC=0,88).

No experimento seguinte, o produto cruzado dos oito atributos da FOT foi usado com o intuito de melhorar o desempenho dos algoritmos, fornecendo um conjunto de dados em uma dimensão mais alta. Foram geradas 36 combinações, que somadas a variável classe, totalizaram 37 atributos de entrada nos algoritmos testados. Como resultado, o 1-NN teve melhor desempenho com AUC igual a 0,89. Já o RBGAOT não convergiu nesse teste, devido aos cálculos realizados durante a marginalização da rede realizados pelo algoritmo *Junction Tree*. Essa limitação também pode ser observada em outros trabalhos com Redes Bayesianas, como no artigo (SILANDER et al., 2012), onde o número máximo de variáveis suportadas pelo modelo é 30.

No terceiro experimento foi feita a seleção de atributos de entrada com base na acurácia que, além de coincidir com a seleção realizada por um especialista, foi a técnica de seleção que indicou variáveis com melhor desempenho. Ao todo, cinco variáveis foram selecionadas:  $R_o$ ,  $R_m$ ,  $X_m$ ,  $C_{din}$  e  $Z_{4Hz}$ , sendo o 1-NN o algoritmo que apresentou melhor desempenho, com AUC igual a 0,87. Já o algoritmo RBGAOT obteve seu desempenho mais baixo, com AUC igual a 0,79.

Durante o quarto experimento, foi aplicado o produto cruzado nos cinco atributos selecionados no teste anterior, gerando um total de 15 combinações que, somada a variável classe, totalizaram 16 atributos na entrada dos classificadores. O ADAB e RF foram os algoritmos que tiveram melhores resultados com valores de AUC iguais a 0,89.

No quinto experimento, a seleção de atributos foi feita dentre as 36 combinações resultantes da aplicação do método do produto cruzado nos oito atributos fornecidos pela FOT. Os classificadores que apresentaram melhor desempenho foram 1-NN e ADAB, com valores de AUC iguais a 0,89, seguidos do RF e RSVM, com valores de AUC iguais a 0,88.

Em todos os experimentos, pelo menos um algoritmo de aprendizado de máquina apresentou sensibilidade acima de 80%, ao observar uma especificidade fixada em 75%. Em uma situação mais restrita para o algoritmo, com especificidade fixada em 90%, pelo menos dois algoritmos alcançaram a faixa da sensibilidade moderada (70 a 90%) nos experimentos com seleção de atributos. Esse é um bom resultado já que essa análise representa uma situação onde o algoritmo é limitado a 10% de falsos positivos. Em ambos os casos, os resultados obtidos superaram a sensibilidade obtida pelo melhor parâmetro da FOT.

Além da acurácia no diagnóstico, a interpretabilidade também foi analisada pelas Redes Bayesianas, construídas e selecionadas por Algoritmos Genéticos. Mesmo sendo treinada com um conjunto de dados limitado, essa técnica se mostrou eficiente, apresentando probabilidades condicionais capazes de descrever o comportamento das características do sistema respiratório de um indivíduo portador de fibrose cística. Vale ressaltar que as redes geradas com cinco atributos da FOT (terceiro experimento), apresentaram maior facilidade em sua inferência, porém menor valor de AUC (AUC=0,79). Já as redes geradas com oito atributos (primeiro experimento), apresentaram melhor desempenho (AUC=0,88), porém possuem uma análise mais difícil devido à maior quantidade de variáveis e ligações entre elas.

O presente trabalho mostrou que o uso de Redes Bayesianas fornece interpretabilidade ao resultado obtido, mostrando as relações existentes entre as variáveis que descrevem a biomecânica do sistema respiratório. Por meio das estruturas geradas, é possível quantificar e compreender melhor como essas variáveis se relacionam, mantendo ainda boa acurácia na detecção de alterações respiratórias em portadores de fibrose cística. Sendo assim, novas informações são geradas e, somadas aos métodos atuais, podem ser usadas para auxílio da equipe médica no estudo da doença.

Uma das limitações observadas no algoritmo RBGAOT foi sua sensibilidade durante o experimento com 37 atributos de entrada, fazendo com que ele não convergisse. Uma melhoria proposta para esse problema é o uso de outro método que realize a marginalização da rede e tenha menor custo computacional. Outra limitação está no descarte das estruturas que não sejam DAG ou que não possuam a variável classe. Uma possível solução para estes casos é a elaboração de uma rotina que realize o reparo dessas redes geradas, transformando-

as em estruturas válidas para o problema. Outras metaheurísticas também podem ser testadas para criação e seleção de estruturas de Redes Bayesianas, além do algoritmo genético.

Este trabalho foi desenvolvido em ambiente Matlab por se tratar de um protótipo. Entretanto, como melhoria futura, será feita a implementação do algoritmo RBGAOT no *software* Python, devido à facilidade em usar um ambiente aberto e gratuito. Há outras duas propostas que visam melhorar a forma de analisar o resultado obtido pelas Redes Bayesianas. Uma delas é a disponibilização de uma plataforma na internet, onde outros pesquisadores possam obter modelos ao inserir seus dados. Outra proposta é o desenvolvimento de um algoritmo que forneça os valores de probabilidade obtidos pela inferência diagnóstica nas redes de forma automática.

## REFERÊNCIAS

- AMARAL JLM, LOPES AJ, FARIA ACD, MELO PL. *Machine learning algorithms and forced oscillation measurements to categorise the airway obstruction severity in chronic obstructive pulmonary disease*, *Computer Methods and Programs in Biomedicine*, Elsevier 118, p. 186-197, 2015;
- AMARAL JLM, LOPES AJ, JANSEN JM, FARIA ACD, MELO PL. *An improved method of early diagnosis of smoking-induced respiratory changes using machine learning algorithms*, *Computer Methods and Programs in Biomedicine*, Elsevier 112, p. 441-454, 2013;
- AMARAL JLM, LOPES AJ, VEIGA J, FARIA ACD, MELO PL. *High-accuracy detection of airway obstruction in asthma using machine learning algorithms and forced oscillation measurements*, *Computer Methods and Programs in Biomedicine*, Elsevier 144, p. 113–125, 2017;
- ANDERSEN DH, *Cystic fibrosis of the pancreas and its relation to celiac disease clinical and pathologic study*, *American Journal of Diseases of Children* 56(2):344-399, 1938;
- ARA-SOUZA AL. *Redes Bayesianas: uma introdução aplicada a Credit Scoring*, Centro de Ciências Exatas e Tecnológicas – Universidade Federal de São Carlos, 2010;
- BARBER D. *Probabilistic Modelling and Reasoning The Junction Tree Algorithm*, Universidade de Edinburgh, 2003;
- BRATKO I. *Machine Learning: Between Accuracy and Interpretability, Learning Networks and Statistics*, *International Centre for Mechanical Sciences*, Vol. 382, p. 164-177, Springer, Vienna, 1997;
- BREIMAN L. *Random Forests*, *Kluwer Academic Publishers, Machine Learning*, v.45, p. 5–32, 2001;
- BROWN LE; TSAMARDINOS I; ALIFERIS CF. *A Novel Algorithm for Scalable and Accurate Bayesian Network Learning*; Departamento de Informática Biomédica da Universidade de Vanderbilt, 2004;
- CARVALHO MA. *Discretização de Atributos Contínuos em Sistemas de Informação Utilizando Algoritmos Genéticos para a Aplicação da Teoria dos Conjuntos Aproximados*, Universidade Federal de Itajubá, 2010;
- CASTELLANI C., CUPPENS H., MACEK M.J., CASSIMAN J.J., KEREM E., DURIE P., TULLIS E., ASSAEL B.M., BOMBIERI C., BROWN A., CASALS T., CLAUSTRES M., CUTTING G.R., DEQUEKER E., DODGE J., DOULL I., FARRELL P., FEREC C., GIRODON E., JOHANNESOM M., KEREM B., KNOWLES M., MUNCK A., PIGNATTI P.F., RADOJKOVIC D., RIZZOTTI P., SCHWARZ M., STUHRMANN M., TZETIS M., ZIELENSKI J., ELBORN J.S; *Consensus on the use and interpretation of cystic fibrosis mutation analysis in clinical practice*; *Journal of Cystic Fibrosis* 7, p.179–196; 2008;
- CASTILLO E; GUTIÉRREZ JM; ALI SH. *Sistemas expertos y modelos de redes probabilísticas*, *Academia Española de Ingeniería*, Madrid, 1996;



CHAVES BB. Estudo do Algoritmo Adaboost de Aprendizagem de Máquina Aplicado a Sensores e Sistemas Embarcados, Escola Politécnica da Universidade de São Paulo, 2012;

CIVICIOGLU P. *Backtracking Search Optimization Algorithm for numerical optimization problems*, *Applied Mathematics and Computation*, Vol. 219, p. 8121–8144, Elsevier Science, 2013;

COLLINS M; SCHAPIRE RE; SINGER Y. *Logistic Regression, AdaBoost and Bregman Distances*, *Machine Learning*, v. 48, p. 253-285, Kluwer Academic Publishers, 2002;

COOPER G; HERSKOVITS E. *A Bayesian method for the induction of probabilistic networks from data*, *Technical Report SMI-91-1, Section on Medical Informatics*, Universidade de Stanford, 1991;

COSTA HSRM. Estudo comparativo de abordagens ao problema de débito de transações bancárias em contas com saldo insuficiente, Departamento de Matemática Aplicada Faculdade de Ciências da Universidade do Porto, 2012;

*Cystic Fibrosis Foundation Patient Registry, Annual Data Report*, Bethesda, Maryland, 2014;

DALCIN PLT, SILVA FAA; Fibrose cística no adulto: aspectos diagnósticos e terapêuticos; *Jornal Brasileiro de Pneumologia*, p. 107-117; 2008;

DUBOIS AB, BRODY AW, LEWIS DH, BURGESS BF. *Oscillation mechanics of lungs and chest in man*, *Journal of applied physiology*, 8:587-594, 1956;

DUIN RPW; JUSZCZAK P; PACLIK P; PEKALSKA E; RIDDER DMJ; TAX DMJ; VERZAKOV S. *PRTools 4.1, A Matlab Toolbox for Pattern Recognition*, Universidade de Tecnologia Delft, Holland, 2007;

FACELI K; LORENA AC; GAMA J; CARVALHO ACPLF. *Inteligência Artificial: uma Abordagem de Aprendizado de Máquina*, LTC, 2011;

FARRELL PM, WHITE TB, REN CL, HEMPSTEAD SE, ACCURSO F, DERICHS N, HOWENSTINE M, MCCOLLEY SA, ROCK M, ROSENFELD M, SERMET-GAUDELUS I, SOUTHERN KW, MARSHALL BC, SOSNAY PR, *Diagnosis of Cystic Fibrosis: Consensus Guidelines from the Cystic Fibrosis Foundation*, *The Journal of Pediatrics* Volume 181S, 2017;

FAWCETT T. *An Introduction to ROC Analysis*, *Pattern Recognition Letters*, V. 27, N. 8, p. 861–874, 2006;

FAYYAD UM; IRANI KB. *Multi-Interval Discretization of Continuous-Valued Attributes for Classification Learning*, *Machine Learning*, p. 1022-1027, 1993;

GABRIEL PHR; DELBEM ACB. *Fundamentos de Algoritmos Evolutivos*, Notas Didáticas do ICMC-USP, N. 75, p. 35, 2008;

GACTO MJ; ALCALÁ R; HERRERA F. *Interpretability of linguistic fuzzy rule-based systems: An overview of interpretability measures*, *Information Sciences*, Elsevier, V. 181, p. 4340–4360, DOI:10.1016, 2011;

GONÇALVES AR. *Redes Bayesianas*; Universidade de Campinas; Disponível em: <<http://www-users.cs.umn.edu/~andre/arquivos/pdfs/bayesianas.pdf>> Acessado em: 03/01/2018;

GUYON I; ELISSEEFF A. *An Introduction to Variable and Feature Selection*, *Journal of Machine Learning Research*, Vol 3, 1157-1182, 2003;

HANLEY JA; MCNEIL BJ, *The Meaning and Use of the Area under a Receiver Operating Characteristic*; *Diagnostic Radiology*; Vol. 143, N. 1; 1982;

HASTIE T; TIBSHIRANI R; FRIEDMAN J. *The Elements of Statistical Learning: Data Mining, Inference, and Prediction*, Springer, 2008;

HERSKOVITS E; COOPER G. *Kutató: An Entropy-Driven System for Construction of Probabilistic Expert Systems from Databases*, Knowledge Systems Laboratory, Medical Computer Science - Stanford University, 1991;

HORTA RAM; CARVALHO FA; ALVES FJS; JORGE MJ. *Comparação de Técnicas de Seleção de Atributos para Previsão de Insolvência de Empresas Brasileiras no Período 2005-2007*. 34º Encontro da ANPAD, 2010;

HOUCK CR; JOINES JA; KAY MG. *A Genetic Algorithm for Function Optimization: A Matlab Implementation*, Universidade do Estado da Carolina do Norte, NCSU-IE, TR 95–09, 1995;

HUANG J; LING CX. *Using AUC and Accuracy in Evaluating Learning Algorithms*, *IEEE Transaction Knowledge and Data Engineering*, Vol. 17, N. 3, p. 299–310, 2005;

KORB KB; NICHOLSON AE. *Introducing Bayesian Network, Bayesian Artificial Intelligence*, 2ª Edição, Cap. 2, p. 29-54, CRC Press, 2010;

KUNCHEVA LI. *Combining Pattern Classifiers: Methods and Algorithms*, Wiley Interscience, New Jersey, 2014;

LACERDA EGM; CARVALHO ACPLF. *Introdução aos Algoritmos Genéticos, Sistemas Inteligentes: Aplicações a Recursos Hídricos e Ciências Ambientais*, Capítulo 3, p. 87-148, Editora da Universidade Federal do Rio Grande do Sul, 1999;

LACERDA LS; LOPES AJ; CARVALHO ARS; GUIMARÃES ARM; FIRMIDA MC; CASTRO MCS; MOGAMI R; MELO PL. *The Role of Multidetector Computed Tomography and the Forced Oscillation Technique in Assessing Lung Damage in Adults With Cystic Fibrosis*, *Respiratory Care*, Vol 63 Issue 3, PubMed: 29208759, 2017;

- LARRAÑAGA P; POZA M; YURRAMENDI Y; MURGA RH; KUIJPERS CMH. *Structure Learning of Bayesian Networks by Genetic Algorithms: Performance Analysis of Control Parameters*; IEEE Transactions on pattern analysis and machine intelligence, Vol. 18, N. 9, p. 912-926, 1996;
- LIAW A; WIENER M. *Classification and Regression by Random Forest*; R. News, Vol. 2/3, p.18–22, 2002;
- LIMA AN, FARIA ACD, LOPES AJ et al. Técnica de oscilações forçadas na avaliação funcional de pacientes com fibrose cística com idade superior a 18 anos, Pulmão RJ 2010;
- LIMA AN, FARIA CDF, LOPES AJ, JANSEN JM, MELO PL. *Forced oscillations and respiratory system modeling in adults with cystic fibrosis*, BioMedical Engineering OnLine, DOI 10.1186/s12938-015-0007-7, 2015;
- LORENA AC; CARVALHO ACPLF. Uma Introdução às *Support Vector Machines*, Revista de Informática Teórica e Aplicada, Volume 14 - Número 2, 2007;
- MACLEOD D, BIRCH M. *Respiratory input impedance measurement: forced oscillation methods*, Medical & Biological Engineering & Computing, Vol 39 p. 505-516, 2001;
- MANDAL S; SINHA RK; MITTAL K. *Comparative Analysis of Backtrack Search Optimization Algorithm with other Evolutionary Algorithms for Global Continuous Optimization*, International Journal of Computer Science and Information Technologies, V. 6, N. 3, p. 3237-3241, 2015;
- MARGINEANTU DD; DIETTERICH TG. *Pruning Adaptive Boosting*, Machine Learning: Proceedings of the Fourteenth International Conference, p. 211-218, 1997;
- MARINHO CL; MAIOLI MCP; AMARAL JLM; LOPES AJ; MELO PL. *Respiratory resistance and reactance in adults with sickle cell anemia: Correlation with functional exercise capacity and diagnostic use*, PLoS One 12 (12): e0187833, 2017;
- MARQUES RL; DUTRA I. Redes Bayesianas: o que são, para que servem, algoritmos e exemplos de aplicações; Coppe Sistemas – UFRJ, 2003;
- MELO PL, GIANELLA-NETO WM A. Avaliação da mecânica ventilatória por oscilações forçadas: Fundamentos e aplicações clínicas. Jornal brasileiro de pneumologia: publicação oficial da Sociedade Brasileira de Pneumologia e Tisiologia; 26:194-206; 2000;
- MELO PL. Técnica de oscilações forçadas na prática pneumológica: Princípios e exemplos de potenciais aplicações, Pulmão RJ, Vol 24 p. 42-48, 2015;
- MENSXMACHINA, *Toolbox Probabilistic Graphical Model 9.2.3*, Universidade de Creta, Departamento de Ciência da Computação, Campus Voutes, 2011, Disponível em: <<http://mensxmachina.org/en/software/pgm-toolbox/>>. Acessado em: 08/02/2018;
- MERSCHMANN LHC. Classificação Probabilística Baseada em Análise de Padrões, Universidade Federal Fluminense, 2007;

METZ CE. *Basic Principles of ROC Analysis, Seminars in Nuclear Medicine*, Vol. 8, N. 4, 1978;

MIRANDA IA, FARIA ACD, LOPES AJ, JANSEN JM, MELO PL. *On the Respiratory Mechanics Measured by Forced Oscillation Technique in Patients with Systemic Sclerosis*; Plos One, Vol. 8(4): e61657, doi:10.1371/journal.pone.0061657, 2013;

MITCHELL TM. *Machine Learning*, McGraw-Hill; 1997;  
MOTA LR, SOUZA EL, ROCHA PHSA, FONSECA MJ, SANTOS JF, LAGE VMGB, LIMA RLLF. Estudos genéticos sobre a Fibrose Cística no Brasil: uma revisão sistemática, *Revista de Ciências Médicas e Biológicas*, 2015;

MURUZÁBAL J; COTTA C. *A Study on the Evolution of Bayesian Network Graph Structures, Advances in Probabilistic Graphical Models*, p. 193-213, 2007;

MYERS J, LASKEY K, DEJONG K. *Learning Bayesian Networks from Incomplete Data using Evolutionary Algorithms*, 1st Annual Conference on Genetic and Evolutionary Computation, V. 1, p. 458-465, 1999;

NEAPOLITAN, RE. *Learning Bayesian Networks, Prentice Hall Series in Artificial Intelligence, Northeastern Illinois University*, 2003;

ONISKO A; DRUZDZEL MJ; WASYLUK H. *Learning Bayesian Network Parameters from Small Data Sets: Application of Noisy-OR Gates, International Journal of Approximate Reasoning*, Elsevier, Vol. 27, p. 165-182, 2001;

PINHO AF; MONTEVECHI JAB; MARINS FAS; MIRANDA RC. Algoritmos Genéticos: Fundamentos e Aplicações, *Meta-Heurísticas em Pesquisa Operacional*, Capítulo 2, p. 21-32, DOI: 10.7436/2013.mhpo.02, 2013;

PGM – *Probabilistic Graphical Model toolbox*, Mens x Machina, Departamento da Ciência da Computação, Universidade de Creta, Grécia, Disponível em:  
<<http://mensxmachina.org/en/software/pgm-toolbox>> Acessado em: 10/02/2018;

PIFER AC. Estudo comparativo de métricas de pontuação para aprendizagem estrutural de Redes Bayesianas, Centro de Tecnologia – Programa de Pós-Graduação em Engenharia Elétrica, Universidade Federal do Rio Grande do Norte, 2006;

REIS FJC, DAMACENO N. Fibrose Cística, *Jornal de Pediatria, Sociedade Brasileira de Pediatria* Vol 74 Supl 1, 1998;

RIBEIRO FCV; LOPES AJ; MELO PL. *Reference values for respiratory impedance measured by the Forced Oscillation Technique in adult men and women*, *The Clinical Respiratory Journal*, DOI: 10.1111/crj.12783, PMID: 29470844, 2018;

RIBEIRO JD, RIBEIRO MAGO, RIBEIRO AF, *Controvérsias na fibrose cística – do pediatra ao especialista*, *Jornal de Pediatria* Vol 78 Supl 2, 2002;

RODRIGUES YE; MANICA E; ZIMMER ER; PASCOAL TA; MATHOTAARACHCHI SS; ROSA-NETO P. *Wrappers Feature Selection in Alzheimer's Biomarkers Using kNN and*

*SMOTE Oversampling*. Sociedade Brasileira de Matemática Aplicada e Computacional, Tendências em Matemática Aplicada e Computacional, Vol. 18, N. 1, p. 15-34, doi: 10.5540, tema 2017.018.01.0015, 2017;

ROSA TO; LUZ HS. Conceitos Básicos de Algoritmos Genéticos: Teoria e Prática, Anais do XI Encontro de Estudantes de Informática do Tocantins, p. 27-37, 2009;

SANTANA AL; REGO LP; FRANCÊS CRL; CARVALHO SV; VIJAYKUMAR NL. Aplicação de Modelos Markovianos para a Análise Temporal e Melhoria da Interpretabilidade de Redes Bayesianas, 39º SBPO – A pesquisa Operacional e o Desenvolvimento, p. 456 -465, 2007;

SCHAPIRE RE. *Explaining AdaBoost, Empirical Inference*, p. 37–52, Springer, 2013;

SILANDER T; MYLLYMÄKI P. A Simple Approach for Finding the Globally Optimal Bayesian Network Structure, *Proceedings of the Twenty-second Annual, Conference on Uncertainty in Artificial Intelligence*, 2012;

SILVA WT; LADEIRA M. Mineração de Dados em Redes Bayesianas, Universidade de Brasília, Capítulo 6, 2016;

SMOLA A; VISHWANATHAN SVN. *Introduction to Machine Learning*, Cambridge University Press, 2008;

TONDA A; LUTTON E; REUOLLION R; SQUILLERO G; WUILLEMIN PH. *Bayesian Network Structure Learning from Limited Datasets through Graph Evolution*, 15º European Conference on Genetic Programming, EuroGP 2012, Malaga – Spain, 7244, p. 254-265, 2012;

ZHU J; ZOU H; SAHARON R; HASTIE T. *Multi-class AdaBoost*, *Statistics and Its Interface*, v. 2, p. 349–360, 2009;

## APÊNDICE A – Combinações do Produto Cruzado

A Tabela 40 mostra as 36 combinações obtidas ao aplicar o método do produto cruzado nos oito parâmetros fornecidos pela FOT. Essas combinações foram usadas como atributos de entrada nos experimentos realizados nos itens 5.4 e 5.7. Já a Tabela 41 mostra as 15 combinações geradas a partir dos cinco atributos da FOT selecionados durante o experimento do item 5.6.

**Tabela 40 – 36 Combinações geradas pelo produto cruzado no experimento dos itens 5.4 e 5.7 com seleção de cinco atributos**

	$F_r$	$X_m$	$R_o$	$S$	$R_m$	$C_{din}$	$E_{din}$	$Z_{4Hz}$
$F_r$	$F_r.F_r$	$F_r.X_m$	$F_r.R_o$	$F_r.S$	$F_r.R_m$	$F_r.C_{din}$	$F_r.E_{din}$	$F_r.Z_{4Hz}$
$X_m$	-	$X_m.X_m$	$X_m.R_o$	$X_m.S$	$X_m.R_m$	$X_m.C_{din}$	$X_m.E_{din}$	$X_m.Z_{4Hz}$
$R_o$	-	-	$R_o.R_o$	$R_o.S$	$R_o.R_m$	$R_o.C_{din}$	$R_o.E_{din}$	$R_o.Z_{4Hz}$
$S$	-	-	-	$S.S$	$S.R_m$	$S.C_{din}$	$S.E_{din}$	$S.Z_{4Hz}$
$R_m$	-	-	-	-	$R_m.R_m$	$R_m.C_{din}$	$R_m.E_{din}$	$R_m.Z_{4Hz}$
$C_{din}$	-	-	-	-	-	$C_{din}.C_{din}$	$C_{din}.E_{din}$	$C_{din}.Z_{4Hz}$
$E_{din}$	-	-	-	-	-	-	$E_{din}.E_{din}$	$E_{din}.Z_{4Hz}$
$Z_{4Hz}$	-	-	-	-	-	-	-	$Z_{4Hz}.Z_{4Hz}$

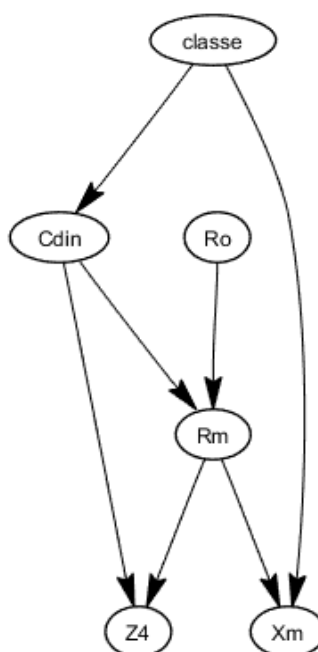
**Tabela 41 – 15 Combinações geradas pelo produto cruzado no experimento do item 5.6 com seleção de cinco atributos**

	$X_m$	$Z_{4Hz}$	$R_m$	$C_{din}$	$R_o$
$X_m$	$X_m.X_m$	$X_m.Z_{4Hz}$	$X_m.R_m$	$X_m.C_{din}$	$X_m.R_o$
$Z_{4Hz}$	-	$Z_{4Hz}.Z_{4Hz}$	$Z_{4Hz}.R_m$	$Z_{4Hz}.C_{din}$	$Z_{4Hz}.R_o$
$R_m$	-	-	$R_m.R_m$	$R_m.C_{din}$	$R_m.R_o$
$C_{din}$	-	-	-	$C_{din}.C_{din}$	$C_{din}.R_o$
$R_o$	-	-	-	-	$R_o.R_o$

## APÊNDICE B – Inferência sobre estruturas de Redes Bayesianas

### 1. Inferência sobre a Rede 1 com cinco atributos de entrada

A rede da Figura 39 possui uma característica diferente das outras redes apresentadas. Além da variável *classe* (Tabela 42), a variável  $R_o$  também foi usada como nó raiz, e portanto, também possui uma tabela de probabilidade à priori (Tabela 43). Isso ocorre devido à aleatoriedade dos indivíduos gerados na população inicial ou criados através dos operadores de mutação e *crossover* do algoritmo genético. Esses indivíduos são apresentados como possíveis soluções ao problema e selecionados de acordo com o valor da AUC que possuem, ou seja, o algoritmo genético não leva em consideração a presença de um ou mais nós raízes, desde que a estrutura gerada tenha um bom desempenho. Mesmo com essa particularidade, é possível inferir sobre as tabelas de DPC obtidas por essa rede.



**Figura 39 – Estrutura da rede 1 gerada com cinco atributos de entrada**

**Tabela 42 – Probabilidades à priori da variável *classe* da rede 1 com cinco atributos de entrada**

$P(\text{classe} = 0)$	$P(\text{classe} = 1)$
0,50	0,50

**Tabela 43 – Probabilidades à priori da variável  $R_o$  da rede 1 com cinco atributos de entrada**

$P(R_o = 1)$	$P(R_o = 2)$
0,69	0,31

A variável  $C_{din}$  é influenciada apenas pela variável *classe*. Pela Tabela 44, pode-se observar que há probabilidade de 0,92 do indivíduo ter complacência maior, dado que não possui a doença (*classe*=0). Para o caso de valores mais baixos ( $C_{din}$ =1), há uma possibilidade de 0,61 do paciente pertencer à classe 1.

**Tabela 44 – DPC da variável  $C_{din}$  da rede 1 com cinco atributos de entrada**

	$P(C_{din} = 1 classe)$	$P(C_{din} = 2 classe)$
<i>classe</i> = 0	0,08	0,92
<i>classe</i> = 1	0,61	0,39

De acordo com a rede 1, a variável  $X_m$  recebe influencias das variáveis  $R_m$  e *classe*. Pela Tabela 45, são retiradas as condições a seguir mostrando que na maioria das combinações há menor probabilidade de um indivíduo apresentar  $X_m=1$ , dado que foi observado um baixo valor da resistência ( $R_m=1$ ), independente do valor da classe:

$$P(X_m = 1|R_m = 1, classe = 0) = 0,04$$

$$P(X_m = 1|R_m = 1, classe = 1) = 0,25$$

Outra observação pode ser feita levando em consideração a variável *classe*. Também há menor chance de um indivíduo apresentar baixa reatância, dado que não é portador da doença (*classe*=0):

$$P(X_m = 1|R_m = 1, classe = 0) = 0,04$$

$$P(X_m = 1|R_m = 2, classe = 0) = 0,36$$



**Tabela 45 – DPC da variável  $X_m$  da rede 1 com cinco atributos de entrada**

	$P(X_m = 1   R_m, classe)$	$P(X_m = 2   R_m, classe)$
$R_m = 1, classe = 0$	0,04	0,96
$R_m = 2, classe = 0$	0,36	0,64
$R_m = 1, classe = 1$	0,25	0,75
$R_m = 2, classe = 1$	0,78	0,22

A variável  $Z_{4Hz}$  é influenciada por  $R_m$  e  $C_{din}$ . Pela Tabela 46 pode-se observar que há maior probabilidade em ter alta impedância  $Z_{4Hz}$ , dado que foi observada complacência igual a 1. Outra observação a ser feita é sobre a combinação  $R_m=2$  e  $C_{din}=2$ . Trata-se de uma situação difícil de ocorrer, pois normalmente um indivíduo com alta resistência  $R_m$  é portador da doença e possui complacência  $C_{din}=1$ . Mesmo assim o algoritmo calcula uma probabilidade de 0,65 para um paciente ter alta impedância  $Z_{4Hz}$ , dado que foram observados  $R_m$  e  $C_{din}$  iguais a 2.

**Tabela 46 – DPC da variável  $Z_{4Hz}$  da rede 1 gerada com cinco atributos de entrada**

	$P(Z_{4Hz} = 1   R_m, C_{din})$	$P(Z_{4Hz} = 2   R_m, C_{din})$
$R_m = 1, C_{din} = 1$	0,22	0,78
$R_m = 2, C_{din} = 1$	0,05	0,95
$R_m = 1, C_{din} = 2$	0,97	0,03
$R_m = 2, C_{din} = 2$	0,35	0,65

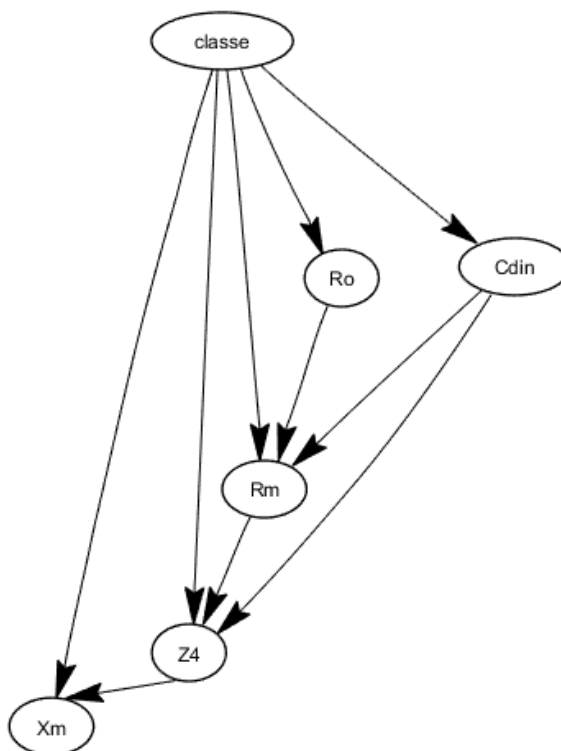
A Tabela 47 mostra as relações de  $R_m$ ,  $C_{din}$  e  $R_o$ , de acordo com a rede 1. As probabilidades contidas nessa tabela reafirmam a relação direta entre  $R_m$  e  $R_o$ , mostrando a alta probabilidade de  $R_m$  ser baixo, dado que  $R_o$  também é baixo. Da mesma forma, observa-se alta probabilidade de um paciente ter  $R_m=2$ , dado que foi observado  $R_o=2$ .

**Tabela 47 – DPC da variável  $R_o$  da rede 1 gerada com cinco atributos de entrada**

	$P(R_m = 1   C_{din}, R_o)$	$P(R_m = 2   C_{din}, R_o)$
$C_{din} = 1, R_o = 1$	0.88	0.12
$C_{din} = 2, R_o = 1$	0.98	0.02
$C_{din} = 1, R_o = 2$	0.18	0.82
$C_{din} = 2, R_o = 2$	0.19	0.81

## 2. Inferência sobre a Rede 2 com cinco atributos de entrada

A rede da Figura 40 possui seis tabelas de DPC. As probabilidades à priori da variável *classe* estão na Tabela 48 e mostram a probabilidade de um indivíduo pertencer à classe 0 ou a classe 1.



**Figura 40 – Estrutura da rede 2 gerada com cinco atributos de entrada**

**Tabela 48 – Probabilidades à priori da variável *classe* da rede 2 com cinco atributos de entrada**

$P(\text{classe} = 0)$	$P(\text{classe} = 1)$
0,50	0,50

Conforme a rede da Figura 40, a variável  $R_o$  é influenciada apenas pela variável *classe*. Pela Tabela 49 é possível concluir que há probabilidade de um indivíduo ter baixa resistência  $R_o$ , dado que pertence a *classe 0*. Da mesma forma, há chances do indivíduo ter alto valor de  $R_o$ , dado que pertence a *classe 1*.

**Tabela 49 – DPC para a variável  $R_o$  da rede 2 gerada com cinco atributos de entrada**

	$P(R_o = 1 classe)$	$P(R_o = 2 classe)$
$classe = 0$	0,94	0,06
$classe = 1$	0,47	0,53

Assim como observado na rede 1, a variável  $C_{din}$  é dependente apenas da variável  $classe$ . De acordo com as condições a seguir, retiradas da Tabela 50, há alta probabilidade de um indivíduo ter complacência alta ( $C_{din}=2$ ), dado que não possui a doença ( $classe=0$ ). Também há maior probabilidade do paciente ter baixa complacência ( $C_{din}=1$ ), dado que possui a doença ( $classe=1$ ).

**Tabela 50 – DPC para a variável  $C_{din}$  da rede 2 com cinco atributos de entrada**

	$P(C_{din} = 1 classe)$	$P(C_{din} = 2 classe)$
$classe = 0$	0,06	0,94
$classe = 1$	0,61	0,39

As variáveis  $Z_{4Hz}$  e  $classe$  exercem influência sobre a variável  $X_m$ . Conforme as probabilidades da Tabela 51, independente da classe, há alta probabilidade de um indivíduo ter alta reatância ( $X_m=2$ ), dado que  $Z_{4Hz}$  é baixo:

**Tabela 51 – DPC para a variável  $X_m$  da rede 2 com cinco atributos de entrada**

	$P(X_m = 1  Z_{4Hz}, classe)$	$P(X_m = 2  Z_{4Hz}, classe)$
$Z_{4Hz} = 1, classe = 0$	0,02	0,98
$Z_{4Hz} = 2, classe = 0$	0,64	0,36
$Z_{4Hz} = 1, classe = 1$	0,06	0,94
$Z_{4Hz} = 2, classe = 1$	0,80	0,20

A variável  $Z_{4Hz}$  é influenciada pelas variáveis:  $R_m$ ,  $C_{din}$  e *classe*. De acordo com a Tabela 52, há maior probabilidade de um indivíduo apresentar alta impedância  $Z_{4Hz}$ , dado que é portador da doença:

**Tabela 52 – DPC para a variável  $Z_{4Hz}$  da rede 2 com cinco atributos de entrada**

	$P(Z_{4Hz} = 1   R_m, C_{din}, classe)$	$P(Z_{4Hz} = 2   R_m, C_{din}, classe)$
$R_m = 1, C_{din} = 1, classe = 0$	0,28	0,72
$R_m = 2, C_{din} = 1, classe = 0$	0,50	0,50
$R_m = 1, C_{din} = 2, classe = 0$	0,99	0,01
$R_m = 2, C_{din} = 2, classe = 0$	0,72	0,28
$R_m = 1, C_{din} = 1, classe = 1$	0,23	0,77
$R_m = 2, C_{din} = 1, classe = 1$	0,03	0,97
$R_m = 1, C_{din} = 2, classe = 1$	0,92	0,08
$R_m = 2, C_{din} = 2, classe = 1$	0,19	0,81

As variáveis  $C_{din}$ ,  $R_o$  e *classe* influenciam  $R_m$ . As probabilidades da Tabela 53 mostram que  $R_m$  é diretamente proporcional a  $R_o$ , independente dos valores assumidos pelas variáveis *classe* e  $C_{din}$ .

O algoritmo apresentou probabilidade igual a 0,5 para o caso de um indivíduo com as seguintes características:  $C_{din}=1$ ,  $R_o=2$  e *classe*=0. Normalmente, isto ocorre devido à combinação difícil de um indivíduo ter complacência baixa ( $C_{din}=1$ ), alta resistência ( $R_o=2$ ) e não possuir a doença (*classe*=0). Em geral, essas são características de portadores de fibrose cística (*classe*=1).

**Tabela 53 – DPC para a variável  $R_m$  da rede 2 com cinco atributos de entrada**

	$P(R_m = 1   C_{din}, R_o, classe)$	$P(R_m = 2   C_{din}, R_o, classe)$
$C_{din} = 1, R_o = 1, classe = 0$	0,72	0,28
$C_{din} = 2, R_o = 1, classe = 0$	0,99	0,01
$C_{din} = 1, R_o = 2, classe = 0$	0,50	0,50
$C_{din} = 2, R_o = 2, classe = 0$	0,28	0,72
$C_{din} = 1, R_o = 1, classe = 1$	0,80	0,20
$C_{din} = 2, R_o = 1, classe = 1$	0,97	0,03
$C_{din} = 1, R_o = 2, classe = 1$	0,17	0,83
$C_{din} = 2, R_o = 2, classe = 1$	0,19	0,81