

UNIVERSIDADE FEDERAL DE ITAJUBÁ
PROGRAMA DE PÓS GRADUAÇÃO EM
CIÊNCIA E TECNOLOGIA DA COMPUTAÇÃO

Estudo e Avaliação de Classificadores para um Sistema
Anti-Spam

Marcelo Vinícius Cysneiros Aragão

Itajubá, junho de 2018

UNIVERSIDADE FEDERAL DE ITAJUBÁ
PROGRAMA DE PÓS GRADUAÇÃO EM
CIÊNCIA E TECNOLOGIA DA COMPUTAÇÃO

Marcelo Vinícius Cysneiros Aragão

Estudo e Avaliação de Classificadores para um Sistema
Anti-Spam

Dissertação submetida ao Programa de Pós-Graduação em Ciência e
Tecnologia da Computação como parte dos requisitos para obtenção
do Título de Mestre em Ciência e Tecnologia da Computação

Área de Concentração: Sistemas de Computação

Orientador: Prof. Dr. Otávio Augusto Salgado Carpinteiro

Coorientador: Prof. Dr. Edmilson Marmo Moreira

Junho de 2018

Itajubá - MG

Agradecimentos

Agradeço aos amigos e familiares pela paciência, tolerância e apoio emocional durante esta jornada turbulenta e gratificante.

Também registro minha gratidão aos colegas e mentores que me ouviram, criticaram, opinaram e, acima de tudo, me ajudaram a manter a cabeça erguida e a mente aberta, de modo que pudesse fazer o melhor trabalho possível, diante das circunstâncias às quais fui sujeito nos últimos anos.

Finalmente, um saudoso agradecimento com sabor de dedicatória ao meu avô Giuseppe (*in memoriam*), por ter me mostrado que ser “mestre” não consiste apenas em receber um título por cumprir determinados requisitos acadêmicos, mas também em manifestar diariamente a incessável vontade de aprender e consequente responsabilidade de ensinar.

Grande astro! Que seria da tua felicidade se te faltassem aqueles a quem iluminas?

Friedrich Wilhelm Nietzsche

Resumo

O advento das redes de dados móveis vem popularizando o acesso à Internet, aumentando o número de usuários e o volume de troca de *e-mails*. Desde 2005, porém, pelo menos metade deste tráfego é composto por *spam*, ou seja, mensagens não-solicitadas enviadas em massa a fim de promover um produto ou serviço, podendo incluir anúncios, esquemas de pirâmide, doações, correntes, entre outros.

Este tráfego massivo de mensagens indesejadas traz prejuízos aos usuários, tais como o uso excessivo e desnecessário da largura de banda, impacto ambiental, perda de produtividade, exposição de conteúdo impróprio a públicos sensíveis, prejuízos financeiros e legais e redução da efetividade de propagandas legítimas.

Uma técnica largamente empregada por servidores de *e-mail*, a fim de lidar com este problema são as *listas negras*, um banco de dados em tempo real que determina se um endereço IP está enviando *e-mails* que podem ser considerados *spam*. Este tipo de sistema apresenta desvantagens, tais como código proprietário, parcialidade questionável, limitações dos protocolos IPv4/IPv6 e lentidão na classificação.

Outra abordagem ao problema de classificação de *e-mails* consiste em usar os dados contidos nas próprias mensagens e um modelo de aprendizado de máquina, a fim de construir um sistema de decisão automático, inteligente, imparcial e adaptativo, capaz de identificar padrões e se adaptar às características comuns a ambos os tipos de *e-mail*, apresentando alta precisão de classificação.

Portanto, o objetivo deste trabalho é o estudo e a aplicação de métodos de aprendizado de máquina à classificação de *e-mails*, a fim de obter modelos utilizáveis em sistemas anti-*spam* do mundo real. Para isto, são também abordadas técnicas de seleção de características e redução de dimensionalidade, necessária e complementar, respectivamente, ao pré-processamento dos dados utilizados pelos algoritmos.

Após realizar experimentos em diversos *corpora*, os resultados confirmaram que modelos de aprendizado de máquina, aliados a métodos de seleção de características baseados em informação mútua e às técnicas de redução de dimensionalidade por correlação, podem ser aplicados com sucesso à categorização de *e-mails*, detectando *spams* com alto grau de confiabilidade para o usuário.

Palavras-chave: aprendizado de máquina. categorização de texto. anti-spam.

Abstract

The advent of mobile data networks has been popularizing Internet access, increasing the number of *e-mail* users and the volume of exchange of these messages. Since 2005, however, at least half of this traffic is made up of *spam*. i.e., unsolicited messages sent in bulk to promote a product or service, which may include advertisements, pyramid schemes, donations, chains, among others.

This massive traffic of unwanted messages causes harm to users, such as excessive and unnecessary bandwidth use, environmental impact, loss of productivity, exposure of sensitive audiences to inappropriate content, financial and legal damage, and reduction of the effectiveness of legitimate advertisements.

One technique widely employed by *e-mail* servers to deal with this problem is *blacklisting*, a real-time database that determines whether an IP address is sending *e-mail* that may be considered *spam*. This type of system has disadvantages such as proprietary code, questionable partiality, limitations of IPv4/IPv6 protocols and slowness in classification.

Another approach to the problem of *e-mail* classification is to use the data contained in the messages themselves and a machine learning method to build an automatic, intelligent, unbiased and adaptive decision system capable of identifying patterns and adapting features common to both types of *e-mail*, presenting high classification accuracy.

Therefore, the objective of this work is the study and application of machine learning methods to the classification of *e-mails*, in order to obtain models usable in real-world anti-*spam* systems. For this, techniques of characteristics selection and dimensionality reduction, necessary and complementary, respectively, to the pre-processing of the data used by the algorithms will also be approached.

After performing experiments on several *corpora*, the results confirmed that machine learning models, coupled with mutual information-based characteristics selection methods and correlation-based dimensionality reduction techniques, can be applied successfully to the categorization of *e-mails*, detecting *spam* with a high degree of reliability for the user.

Keywords: machine learning. text categorization. anti-spam.

Sumário

Lista de Figuras

Lista de Tabelas

Glossário	p. 19
1 Introdução	p. 23
1.1 Contexto da pesquisa e justificativas	p. 23
1.2 Propostas existentes	p. 27
1.3 Deficiências e necessidade de uma nova abordagem	p. 29
1.4 Objetivos e contribuições da dissertação	p. 30
1.5 Estrutura da dissertação	p. 31
2 Fundamentação Teórica	p. 33
2.1 Aprendizado de máquina	p. 34
2.2 Categorização de texto	p. 36
2.3 Seleção de características	p. 37
2.3.1 Estatística Qui-quadrado (CHI2)	p. 38
2.3.2 Frequência no Documento (FD)	p. 38
2.3.3 Informação Mútua (MI)	p. 39

2.4	Redução de dimensionalidade	p. 40
2.5	Balanceamento de classes	p. 42
2.6	Algoritmos de classificação	p. 43
2.6.1	Métodos Probabilísticos	p. 44
2.6.1.1	Naïve Bayes (NB)	p. 44
2.6.1.2	Averaged One-Dependence Estimators (AODE)	p. 46
2.6.2	Redes Neurais Artificiais	p. 48
2.6.2.1	Single Layer Perceptron (SLP)	p. 50
2.6.2.2	Radial Basis Function Network (RBF)	p. 52
2.6.3	Árvores de Decisão	p. 56
2.6.3.1	Reduced-Error Pruning Tree (REPT)	p. 58
2.6.3.2	Boosted Reduced-Error Pruning Tree (AB-M1)	p. 59
2.6.4	Máquina de Vetores de Suporte	p. 61
2.6.4.1	Linear Support Vector Machine (L-SVM)	p. 63
2.6.4.2	Non-Linear Support Vector Machine (NL-SVM)	p. 65
3	Revisão Bibliográfica	p. 67
4	Metodologia	p. 77
4.1	Modelo proposto	p. 77
4.2	Ferramentas utilizadas	p. 79
4.3	Conjuntos de dados utilizados	p. 81
4.3.1	Ling Spam	p. 81

4.3.2	Spam Assassin	p. 82
4.3.3	TREC	p. 83
4.3.4	Unifei 2017	p. 84
4.3.5	Unifei 2018	p. 88
4.4	Métricas de desempenho	p. 92
4.4.1	Precisão	p. 93
4.4.2	Revocação	p. 93
4.4.3	Escore F_1	p. 94
4.4.4	Tempo de Treinamento	p. 94
4.4.5	Tempo de Classificação	p. 94
5	Experimentos e Resultados	p. 95
5.1	Seleção de características e padrões zerados	p. 96
5.2	Redução de dimensionalidade	p. 102
5.3	Treinamento e classificação	p. 108
6	Conclusão	p. 117
6.1	Discussão dos resultados e considerações finais	p. 117
6.2	Sugestões para novos trabalhos	p. 120
7	Anexos	p. 121
8	Apêndices	p. 123
	Referências	p. 135

Lista de Figuras

1	Usuários de <i>E-mail</i>	p. 23
2	Tráfego Diário de <i>E-mails</i>	p. 23
3	Percentual de <i>Ham</i> e <i>Spam</i> ao Ano.	p. 24
4	Funcionamento da SBL Advisory.	p. 27
5	Aprendizado não-supervisionado.	p. 35
6	Aprendizado supervisionado.	p. 35
7	Maldição da dimensionalidade.	p. 40
8	Comparação entre dependências dos classificadores NB e AODE. . .	p. 47
9	Neurônio artificial de McCulloch–Pitts (MCP).	p. 48
10	Exemplos de funções de ativação.	p. 49
11	Rede Perceptron de camada única.	p. 50
12	Funções binárias e separabilidade linear.	p. 50
13	Rede de função de base radial.	p. 52
14	Iterações inicial, final e campos receptivos de uma rede RBF. . . .	p. 54
15	Árvore de decisão apresentando resposta à oferta de emprego. . .	p. 57
16	Exemplo de poda de árvore de decisão.	p. 58
17	Ilustração do processo de funcionamento do AdaBoost.	p. 59
18	Hiperplano, vetores de suporte e margem de separação.	p. 61

19	Aumento de dimensionalidade por meio do <i>kernel trick</i>	p. 62
20	Agrupamentos de dados usando máquinas de vetores de suporte. .	p. 62
21	Representação gráfica das margens de separação rígida e suave. . .	p. 63
22	Ilustração da influência do parâmetro C na margem de uma SVM.	p. 64
23	Diagrama de blocos da metodologia proposta.	p. 77
24	Tela principal do WEKA.	p. 79
25	Tela do gerenciador de pacotes do WEKA.	p. 80
26	Visualização 2D da base Unifei 2017 (método CHI2).	p. 85
27	Visualização 2D da base Unifei 2017 (método FD).	p. 86
28	Visualização 2D da base Unifei 2017 (método MI).	p. 87
29	Visualização 2D da base Unifei 2018 (método CHI2).	p. 89
30	Visualização 2D da base Unifei 2018 (método FD).	p. 90
31	Visualização 2D da base Unifei 2018 (método MI).	p. 91
32	Padrões zerados na base Ling Spam.	p. 97
33	Padrões zerados na base Spam Assassin.	p. 97
34	Padrões zerados na base TREC.	p. 97
35	Padrões zerados na base Unifei 2017.	p. 98
36	Padrões zerados na base Unifei 2018.	p. 98
37	Padrões zerados em todas as bases.	p. 98
38	Escore F_1 vs métodos de seleção de características (Ling Spam). .	p. 99
39	Escore F_1 vs métodos de seleção de características (Spam Assassin).	p. 99
40	Escore F_1 vs métodos de seleção de características (TREC). . . .	p. 99

41	Escore F_1 vs métodos de seleção de características (Unifei 2017). . .	p. 100
42	Escore F_1 vs métodos de seleção de características (Unifei 2018). . .	p. 100
43	Escore F_1 vs métodos de seleção de características em todas as bases.	p. 100
44	Redução de dimensionalidade na base Ling Spam.	p. 102
45	Redução de dimensionalidade na base Spam Assassín.	p. 103
46	Redução de dimensionalidade na base TREC.	p. 104
47	Redução de dimensionalidade na base Unifei 2017.	p. 105
48	Redução de dimensionalidade na base Unifei 2018.	p. 106
49	<i>Ranking</i> dos classificadores na base Ling Spam.	p. 109
50	<i>Ranking</i> dos classificadores na base Spam Assassín.	p. 110
51	<i>Ranking</i> dos classificadores na base TREC.	p. 111
52	<i>Ranking</i> dos classificadores na base Unifei 2017.	p. 112
53	<i>Ranking</i> dos classificadores na base Unifei 2018.	p. 113
54	Distribuição dos métodos de seleção de características no <i>ranking</i> . . .	p. 114
55	Distribuição das quantidades de características no <i>ranking</i>	p. 114

Lista de Tabelas

1	Exemplos de tipos de texto e sugestões de classificação.	p. 36
2	Conjunto de dados antes do balanceamento.	p. 42
3	Amostras replicadas a fim de balancear o conjunto de dados. . . .	p. 42
4	Exemplos de funções de <i>kernel</i>	p. 65
5	Subconjuntos do <i>corpus</i> Spam Assassin.	p. 82
6	Quantidades de padrões zerados nos conjuntos de dados	p. 96
7	Redução de dimensionalidade na base Ling Spam.	p. 102
8	Redução de dimensionalidade na base Spam Assassin.	p. 103
9	Redução de dimensionalidade na base TREC.	p. 104
10	Redução de dimensionalidade na base Unifei 2017.	p. 105
11	Redução de dimensionalidade na base Unifei 2018.	p. 106
12	<i>Ranking</i> dos classificadores na base Ling Spam.	p. 109
13	<i>Ranking</i> dos classificadores na base Spam Assassin.	p. 110
14	<i>Ranking</i> dos classificadores na base TREC.	p. 111
15	<i>Ranking</i> dos classificadores na base Unifei 2017.	p. 112
16	<i>Ranking</i> dos classificadores na base Unifei 2018.	p. 113
17	Parâmetros de configuração usados no WEKA.	p. 124
18	Resultados do classificador NB.	p. 125

19	Resultados do classificador AODE.	p. 126
20	Resultados do classificador SLP.	p. 127
21	Resultados do classificador RBF.	p. 128
22	Resultados do classificador REPT.	p. 129
23	Resultados do classificador AB-M1.	p. 130
24	Resultados do classificador SVM (linear).	p. 131
25	Resultados do classificador SVM (polinomial).	p. 132
26	Detalhes do <i>hardware</i> utilizado nos experimentos.	p. 133
27	Detalhes do <i>software</i> utilizado nos experimentos.	p. 133

Glossário

1D-LBP	<i>One-Dimension Linear Binary Pattern</i>
AB-M1	<i>Boosted Reduced-Error Pruning Tree</i>
ADTrees	<i>Alternating Decision Trees</i>
AODE	<i>Averaged One-Dependence Estimators</i>
API	<i>Application Programming Interface</i>
ASCII	<i>American Standard Code for Information Interchange</i>
AT&T	<i>American Telephone and Telegraph</i>
BN	<i>BayesNet</i>
BNS	<i>Bi-Normal Separation</i>
CFS	<i>Correlation-based Feature Selection</i>
CUDA	<i>Compute Unified Device Architecture</i>
DF	<i>Document Frequency</i>
DNS	<i>Domain Name System</i>
DNSBL	<i>Domain Name System-based Blackhole List</i>
EAF	<i>Empirical Attainment Function</i>
EANT2	<i>Evolutionary Acquisition of Neural Topologies, Version 2</i>
ENORA	<i>Evolutionary NON-dominated Radial slots based Algorithm</i>
EQM	Erro Quadrático Médio
EUA	Estados Unidos da América
FD	Frequência no Documento
FLDA	<i>Fisher Linear Discriminant Analysis</i>
FRES	<i>Flesch Reading Ease Score</i>
FT	<i>Functional Tree</i>
GNARL	<i>GeNeralized Acquisition of Recurrent Links</i>

GPESC	Grupo de Pesquisas em Engenharia de Sistemas e de Computação
GPU	<i>Graphics Processing Unit</i>
GRBF	<i>Generalized Radial Basis Functions</i>
HTML	<i>HyperText Markup Language</i>
IB	<i>Information Bottleneck</i>
ID3	<i>Iterative Dichotomiser 3</i>
IG	<i>Information Gain</i>
IP	<i>Internet Protocol</i>
IPv4	<i>Internet Protocol version 4</i>
IPv6	<i>Internet Protocol version 6</i>
ISP	<i>Internet Service Provider</i>
KNN	<i>k-Nearest Neighbors</i>
L-SVM	<i>Linear Support Vector Machine</i>
LMS	<i>Least Mean Squares</i>
LNO	<i>Largest Number of Objects</i>
LOF	<i>Local Outlier Factor</i>
LTS	<i>Largest Total Strength</i>
LVF	<i>Las Vegas Filter</i>
MAP	<i>Maximum A Posteriori</i>
MCC	<i>Matthews Correlation Coefficient</i>
MCP	<i>McCulloch–Pitts</i>
MFD	<i>Most Frequent Decision</i>
MI	<i>Mutual Information</i>
MIME	<i>Multipurpose Internet Mail Extensions</i>
MIT	<i>Massachusetts Institute of Technology</i>
MLP	<i>Multilayer Perceptron</i>
NB	<i>Naïve Bayes</i>
NEAT	<i>NeuroEvolution of Augmenting Topologies</i>
NL-SVM	<i>Non-Linear Support Vector Machine</i>
NSA-PSO	<i>Negative Selection Algorithm - Particle Swarm Optimization</i>

NSGA-II	<i>Nondominated Sorting Genetic Algorithm II</i>
Open-MaLBAS	<i>Open Machine Learning-Based Anti-Spam</i>
P-SVM	<i>Polynomial Support Vector Machine</i>
PCA	<i>Principal Component Analysis</i>
RBF	<i>Radial Basis Function</i>
RDF	<i>Resource Description Framework</i>
REPT	<i>Reduced-Error Pruning Tree</i>
RF	<i>Random Forest</i>
ROC	<i>Receiver Operating Characteristic</i>
RT	<i>Random Tree</i>
SBL	<i>Spamhaus Block List</i>
SLP	<i>Single Layer Perceptron</i>
SMOG	<i>Simple Measure of Gobbledygook</i>
SOM	<i>Self-Organizing Map</i>
SPEA2	<i>Strength Pareto Evolutionary Algorithm 2</i>
SSF	<i>Simplified Silhouette Filter</i>
STA	Servidor Técnico-Administrativo
SVC	<i>Support Vector Clustering</i>
SVM	<i>Support Vector Machine</i>
SVM-RFE	<i>Support Vector Machine Recursive Feature Elimination</i>
TF	<i>Term Frequency</i>
TF-IDF	<i>Term Frequency - Inverse Document Frequency</i>
TREC	<i>Text REtrieval Conference</i>
t-SNE	<i>t-Distributed Stochastic Neighbor Embedding</i>
UCE	<i>Unsolicited Commercial E-mail</i>
UNIFEI	Universidade Federal de Itajubá
UTF-8	<i>8-bit Unicode Transformation Format</i>
WEKA	<i>Waikato Environment for Knowledge Analysis</i>

1 Introdução

1.1 Contexto da pesquisa e justificativas

Nos dias de hoje, com o advento das redes de dados móveis, um número cada vez maior de pessoas possui acesso à Internet de qualquer lugar. Seja no trabalho, em casa ou até mesmo no trânsito, é possível receber e enviar *e-mails* a qualquer momento. Isto implica em um aumento na quantidade de usuários de *e-mail* e, conseqüentemente, na quantidade de trocas de mensagens, conforme pode ser confirmado nos gráficos da figuras 1 e 2 (RADICATI, 2018), respectivamente.

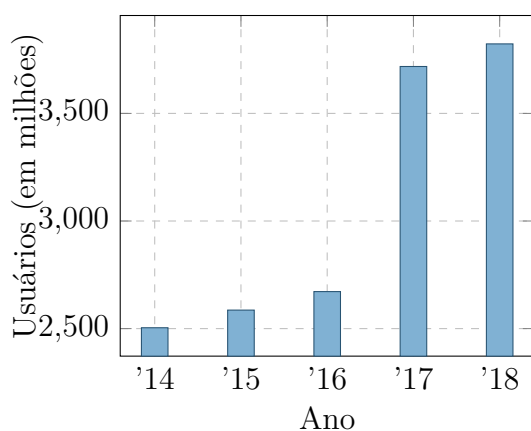


Figura 1: Usuários de *E-mail*.

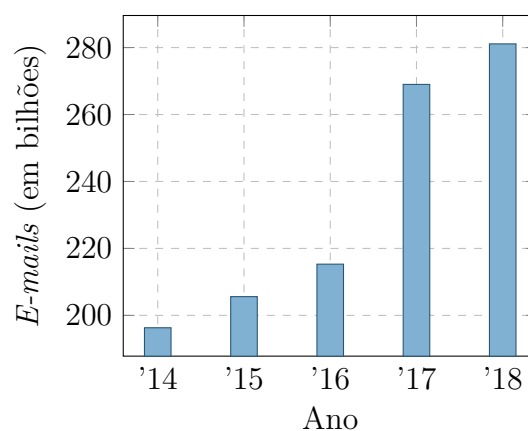


Figura 2: Tráfego Diário de *E-mails*.

Infelizmente não são somente a quantidade de usuários e a facilidade de acesso que vêm aumentando. De acordo com Symantec Corporation (2005), Symantec Corporation (2012) e Symantec Corporation (2018), pelo menos metade do tráfego de *e-mails* desde 2005 é composta por *spam* (vide figura 3), ou seja, mensagens

não solicitadas cujo destinatário não é um conhecido, cliente ou um membro voluntariamente inscrito em uma lista de distribuição. Neste mesmo sentido, é comum fazer referência a *e-mails* legítimos (não *spam*) por meio do termo *ham*.

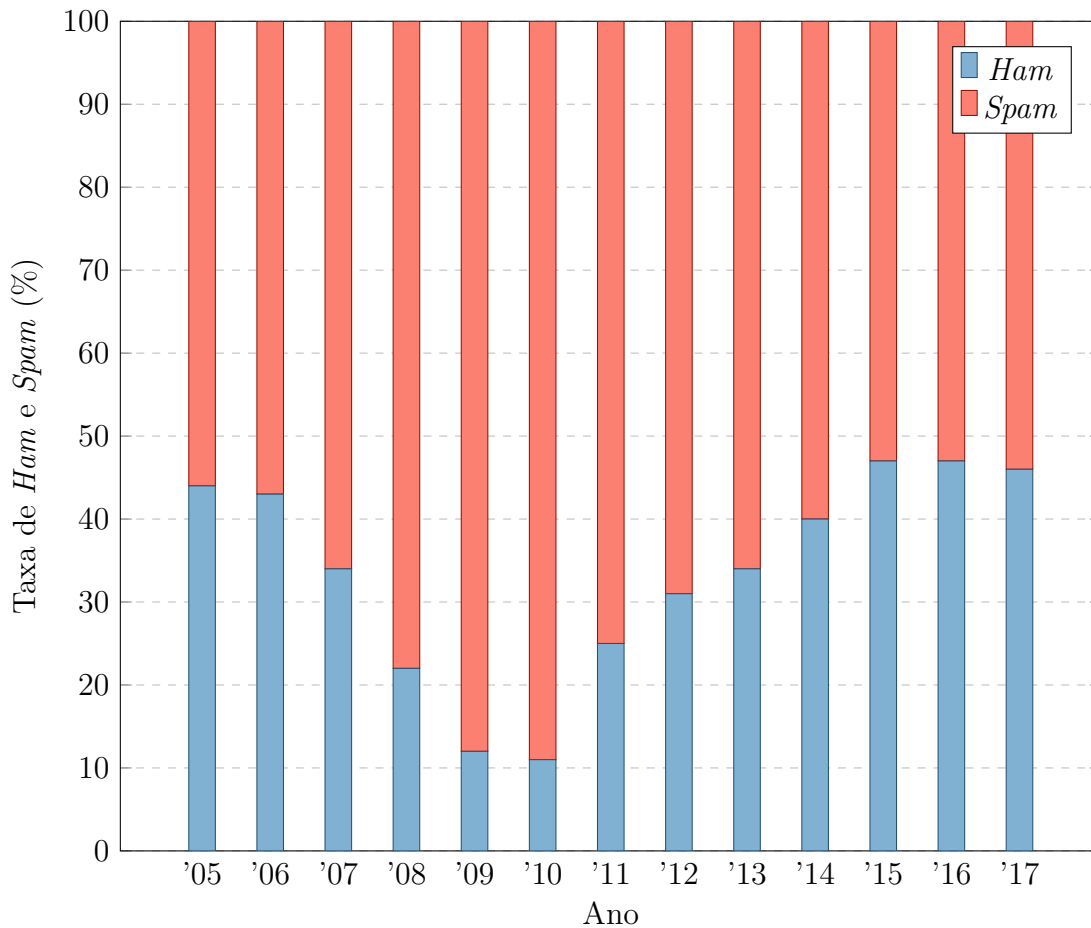


Figura 3: Percentual de *Ham* e *Spam* ao Ano.

Geralmente, *spams* são enviados a um grande número de pessoas para promover um produto ou serviço e, por esta razão, muitas vezes são chamados de *e-mail* comercial não solicitados (UCE, do inglês *Unsolicited Commercial E-mail*). Isso inclui anúncios, esquemas de pirâmide, doações, correntes de mensagens, *e-mail* político, assessoria no mercado de ações e avisos únicos.

Este tráfego massivo de mensagens indesejadas provoca diversas consequências negativas, podendo ser citadas:

- Uso excessivo da largura de banda, provocando lentidão de ISPs (do inglês *Internet Service Provider*) e sobrecarga de servidores de *e-mail*.
- Impacto ambiental: considerando o custo energético para o processamento de *spam*, a energia gasta poderia ser utilizada para outros fins (McAfee Inc.; ICF International, 2009).
- Perda de produtividade: o recebimento de *spam* diminui a produtividade de todos que utilizam *e-mails* como ferramenta de trabalho, pois um tempo extra é dedicado à tarefa de recebimento e leitura, podendo também uma mensagem importante ser apagada por engano, ignorada ou lida com atraso dado o volume de *spams* recebidos. Esta perda também é diretamente associada a prejuízos financeiros.
- Conteúdo impróprio ou ofensivo: como os *spams* são disseminados de forma aleatória, é bem provável que o usuário algumas vezes julgue impróprio e ofensivo o conteúdo da mensagem recebida, podendo até expor crianças a material inadequado.
- Prejuízos financeiros e legais causados por fraude: *spams* estão sendo utilizados para induzir usuários a acessar páginas clonadas de instituições financeiras ou a instalar programas maliciosos, projetados para furtar dados pessoais e financeiros, causando grandes prejuízos aos usuários.
- Redução da efetividade de propagandas legítimas.

Diante dos problemas apresentados, nota-se a importância de filtrar este tipo de mensagem. Uma solução resultaria em um alívio nos sistemas de telecomunicações e em melhor experiência do usuário na utilização de seu *e-mail*.

Um usuário experiente de computador consegue facilmente identificar um *spam*, seja por meio de seu conteúdo (palavras-chave, *links* e/ou imagens), seja por

mera vigilância epistêmica¹(SPERBER et al., 2010) adquirida com o tempo e, possivelmente, com experiências ruins com este tipo de mensagem (exemplo: ser alvo de algum tipo de fraude). Entretanto, para um computador, este processo não é tão trivial. Isto porque o cérebro humano, devido à sua natureza, trabalha de forma altamente paralelizada e possui grande capacidade de aprendizado. Já o computador apresenta comportamento linear, e a identificação de *spams* envolve a criação de regras complexas e algoritmos custosos para processamento, análise e categorização de texto.

Este processo de categorização é normalmente feito utilizando modelos de aprendizado de máquina, sejam estes supervisionados (com amostras rotuladas) ou não-supervisionados (com amostras de treinamento sem classificação). Estes modelos são capazes de identificar padrões e de se ajustar a características comuns a ambos os tipos de *e-mail*, resultando em classificações cada vez mais precisas, com o decorrer das épocas de treinamento. Entretanto, cada modelo apresenta defeitos e qualidades particulares que influenciam na aplicabilidade em sistemas de classificação e processos decisórios em tempo real.

Este trabalho, portanto, visa apresentar e comparar diversos modelos de aprendizado de máquina aplicados à classificação de mensagens eletrônicas, cujas amostras foram obtidas a partir de conjuntos de dados disponibilizados sob domínio público. Tais comparações serão feitas sob a luz de métricas — precisão, revocação, escore F_1 e tempos de treinamento e classificação — bem disseminadas no contexto acadêmico e relevantes ao tema (POWERS, 2011).

Por meio deste projeto, busca-se disseminar o entendimento de diversos modelos de aprendizado de máquina — conceito que chama cada vez mais atenção devido à sua intrínseca adaptabilidade — e evidenciar sua aplicação na solução de problemas práticos e de grande relevância nos dias de hoje.

¹Conjunto de mecanismos cognitivos para filtragem de informação comunicada recebida: calibração de confiança, raciocínio, avaliação de argumentos e questionamento da plausibilidade de reivindicações.

1.2 Propostas existentes

Uma técnica frequentemente empregada por servidores de *e-mail* para lidar com *spam* é por meio do uso de listas negras (JUNG; SIT, 2004) (ZHANG et al., 2009).

O papel de uma DNSBL (do inglês *Domain Name System-based Blackhole List*, comumente conhecida como *blacklist* ou *block list*) é fornecer uma opinião a qualquer usuário que questionar se um determinado endereço IP (do inglês *Internet Protocol*) atende à política de aceitação de *e-mails* de entrada definida no servidor em questão. Tais listas podem ser consultadas no momento em que um novo *e-mail* é recebido pelo servidor, com o objetivo de manter mensagens maliciosas longe das caixas de entrada dos usuários.

Um exemplo de *blacklist* é a Spamhaus Block List (SBL) Advisory (The Spamhaus Project, 2018a), ilustrada na figura 4.

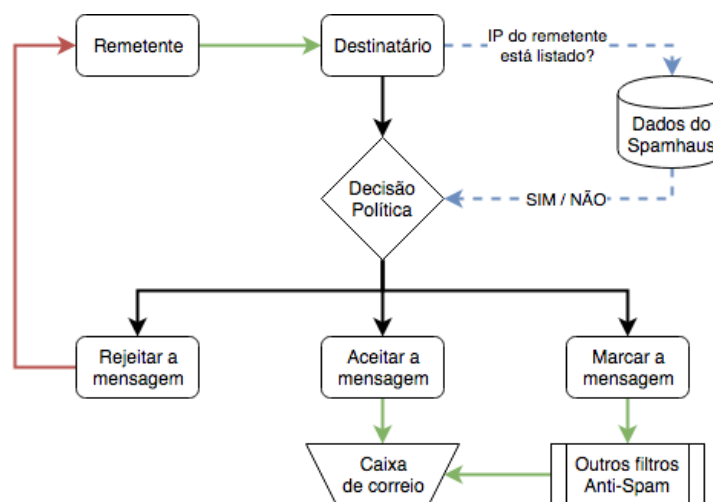


Figura 4: Funcionamento da SBL Advisory (The Spamhaus Project, 2018b).

Como uma etapa inicial adicional, podem ser consultadas listas brancas (*whitelists*), compostas por domínios ou endereços IP previamente aprovados e com permissão de entrega e que, portanto, são isentos de serem consultados em uma lista negra ou submetidos a posteriores verificações de presença de *spam*.

Uma outra abordagem ao problema de classificação de *e-mails* consiste em usar os dados contidos na própria mensagem (isto é, remetente, destinatário, assunto, cabeçalho, corpo, anexos etc.) para construir um sistema de decisão automático, inteligente e imparcial.

Para isto, é necessária a coleta de conjuntos de dados (ou seja, amostras reais de *e-mails*), sobre os quais um ou mais algoritmos de extração de características e aprendizado de máquina analisarão padrões inerentes a cada tipo de mensagem (legítima ou não) e tentarão derivar um conjunto de regras ou uma fronteira (e conseqüentemente uma função) de discernibilidade. O modelo obtido com o treinamento deste método poderá, então, ser usado para classificar novas mensagens com base no conhecimento obtido por meio da análise das mensagens anteriores.

Diversos autores propuseram o uso de modelos de aprendizado supervisionado, sobre *e-mails* previamente classificados da base de dados (ou seja, cada amostra já fora rotulada como *ham* ou *spam*). Utilizando diversos algoritmos, são construídos modelos capazes de inferir sobre a categoria de uma nova mensagem. Dentre as propostas existentes, pode-se citar o método probabilístico Naïve Bayes (ANDROUTSOPOULOS et al., 2000a) (MEYER; WHATELEY, 2004), redes neurais artificiais (CARPINTEIRO et al., 2006), máquina de vetores de suporte (DRUCKER et al., 1999) (ARAGÃO et al., 2016) (KUMARESAN; SARAVANAKUMAR; BALAMURUGAN, 2017) e *k-nearest neighbors* (BARIGOU; BELDJILALI; ATMANI, 2014).

Outros trabalhos adotaram modelos de aprendizado não-supervisionado, que também são capazes de apresentar taxas de precisão de classificação razoáveis, como na comparação entre máquina de vetores de suporte e mapas auto-organizáveis, apresentada em Berger e Merkl (2005).

1.3 Deficiências e necessidade de uma nova abordagem

Sistemas anti-*spam* baseados em listas negras e brancas apresentam diversas desvantagens, dentre as quais:

- Soluções pagas e de código proprietário: apesar das soluções proprietárias apresentarem altas taxas de captura de *spam* (SNYDER, 2015), muitas custam caro (da ordem de dezenas de milhares de dólares ao ano), comprometendo o orçamento de diversas instituições e/ou impedindo filtragem decente de *e-mails*. Além disso, licenças de código fechado impedem o estudo, personalização e extensão de funcionalidades por desenvolvedores independentes.
- Parcialidade questionável: diversos provedores oferecem serviços de consulta e remoção de endereços de *blacklists* mediante pagamento. Isto compromete toda a parcialidade e confiabilidade destes sistemas, já que virtualmente qualquer *spammer* (ou seja, pessoa ou organização que envia *spam*) pode deixar de sê-lo a qualquer momento.
- Fadadas à exaustão devido aos limites do IPv4/IPv6:
 - Se cada endereço IPv4 usar 4 B e seu nome de domínio em uma tabela DNS usar 8 B, seriam necessários $2^{32} * (4 + 8)B = 48 \text{ GiB}$ para armazenar estes dados. É um volume factível de armazenamento; entretanto, esta quantidade de endereços está próxima de se esgotar (ALI, 2012).
 - Se cada endereço IPv6 usar 16 B e seu nome de domínio em uma tabela DNS usar 8 B, seriam necessários $2^{128} * (16 + 8)B = 7.428 \times 10^{27} \text{ TiB}$ para armazenar todos estes dados, o que é impraticável nos dias de hoje.
- Lentas (internet/latência): cada consulta a um servidor de *blacklist* está sujeita a uma latência, ou seja, o tempo que toda requisição leva para atingi-lo, ser processada e um resultado ser retornado ao cliente. Considerando um *round trip time* de 1 ms, um servidor de *e-mails* que processa 1 milhão de mensagens por dia gastaria 16 minutos e 40 segundos apenas para consultar o banco de endereços da *blacklist*. Caso seja feito um *cache* local dos dados, os problemas de armazenamento levantados no item anterior entram em cena.

1.4 Objetivos e contribuições da dissertação

O objetivo geral deste trabalho é aplicar modelos de aprendizado de máquina propostos por diversos autores à classificação de *e-mails*, a fim de obter modelos utilizáveis em sistemas anti-*spam* do mundo real.

Como objetivos e contribuições específicos, tem-se os seguintes pontos:

- pesquisar, estudar e comparar métodos de seleção de características e de redução de dimensionalidade, a fim de eleger aqueles que mais bem atendam às necessidades da proposta;
- elaborar um *framework* de *software* científico genérico, capaz de treinar diversos modelos de aprendizado de máquina e salvá-los para posterior reuso (e, em alguns casos, refinamento);
- expandir a funcionalidade do sistema anti-*spam* de código aberto Open-MaLBAS (do inglês *Open Machine-Learning-Based Anti-Spam*) (FERREIRA, 2018), fornecendo uma gama de métodos de classificação facilmente configuráveis e intercambiáveis;
- efetuar simulações em diversos conjuntos de dados e manter um histórico de resultados para fins de comparação;
- atingir desempenho satisfatório em termos de precisão, revocação, escore F_1 e tempos de treinamento e classificação, por meio da modificação de parâmetros de configuração de cada método;
- consolidar os resultados experimentais, comparando-os àqueles de outros trabalhos da comunidade acadêmica;
- apontar potenciais melhorias nos modelos propostos.

1.5 Estrutura da dissertação

A fundamentação teórica necessária para o desenvolvimento deste estudo está apresentada no capítulo 2. Nele são abordados conceitos de aprendizado de máquina, contemplando suas diferentes estratégias. Também apresenta conceitos de extração de características, filtragem de atributos, balanceamento de classes e elementos da teoria necessária sobre os algoritmos de classificação.

O capítulo 3 revisa, de forma cronológica, diversos trabalhos existentes sobre classificação de *e-mails*, abordando os modelos e métodos utilizados bem como os resultados alcançados em cada um.

Estão descritos no capítulo 4 os detalhes sobre os modelos avaliados neste estudo, desde o tratamento dos dados até a sua apresentação aos métodos de aprendizado de máquina, detalhando todos os seus componentes, inclusive ferramentas externas. Também descreve os conjuntos de dados utilizados. Por fim, esse capítulo aponta os resultados obtidos com a aplicação de cada modelo, inclusive sob a forma de *rankings* visando facilitar a interpretação dos resultados experimentais.

No capítulo 5, são apresentadas visualizações, análises e conclusões sobre os resultados obtidos, agrupados em três seções de forma a facilitar o entendimento de cada aspecto do trabalho.

Concluindo a dissertação, o capítulo 6 sumariza o problema e a proposta para resolvê-lo, bem como faz uma análise qualitativa dos resultados. Por fim, são apresentadas propostas para novos trabalhos.

Adicionalmente, é possível encontrar nos anexos uma cópia da licença sob a qual o código e os conjuntos de dados foram disponibilizados, e nos apêndices há tabelas com todos os resultados experimentais obtidos neste trabalho e detalhes do *hardware* e dos *softwares* utilizados.

2 Fundamentação Teórica

As seções a seguir visam explicar os principais conceitos necessários ao entendimento da proposta deste trabalho. Primeiramente, são introduzidos o conceito e os tipos de aprendizado de máquina, base para toda a aquisição de conhecimento e construção de sistemas computacionais inteligentes. Também visita-se o conceito de classificação de texto, essencial à categorização de quaisquer exemplos de dados multimídia compostos por caracteres (em vez de valores numéricos reais), como por exemplo os *e-mails*, objeto de estudo desta proposta.

Em seguida, as três técnicas de seleção de características adotadas no trabalho são apresentadas, bem como o método de redução de dimensionalidade e de balanceamento de classes aplicado aos conjuntos de dados. A finalidade de tais técnicas é refinar as amostras de treinamento, removendo informações redundantes e garantindo uma base de conhecimento apropriada à construção dos modelos de classificadores.

Na sequência, cada técnica usada para classificação das mensagens é categorizada quanto à sua natureza (métodos probabilísticos, redes neurais artificiais, árvores de decisão e máquina de vetores de suporte) e explicada. Finalmente as métricas para avaliação dos resultados são apresentadas e definidas matematicamente.

2.1 Aprendizado de máquina

O aprendizado é um processo que acompanha o homem desde o início de sua existência. Quando representados de forma coerente e inteligível, estímulos sensoriais desencadeiam todo um processo de organização, associação e geração de conhecimento. O artefato deste processo é a capacidade natural de inferência, conferindo ao homem as capacidades de raciocínio, discernimento e tomada de decisão, mesmo diante de temas sobre os quais nunca havia lidado.

Desde o final da década de 1950 (SAMUEL, 1959), vem-se tentando imitar o complexo processo de aprendizado em ambientes computacionais, a fim de conceber sistemas inteligentes e que possam auxiliar (ou até mesmo substituir) o homem em situações que demandam tomadas sistemáticas de decisão. Esta área do conhecimento recebeu o nome de aprendizado de máquina e é um assunto cada vez mais em voga na comunidade acadêmica e no meio científico.

As tarefas de aprendizagem de máquinas podem ser classificadas em três categorias, dependendo da presença e do tipo de *feedback* de aprendizado disponível para um sistema de aprendizagem (RUSSELL; NORVIG, 2013):

- Aprendizagem não-supervisionada: o agente aprende padrões na entrada, embora não seja fornecido nenhum *feedback* explícito. A tarefa mais comum de aprendizagem não supervisionada é o agrupamento: a detecção de grupos de exemplos de entrada potencialmente úteis.
- Aprendizagem por reforço: o agente aprende a partir de uma série de reforços — recompensas ou punições. Cabe ao agente decidir qual das ações anteriores ao reforço foram as maiores responsáveis por isso.
- Aprendizagem supervisionada: o agente observa exemplos de pares de entrada e saída, e aprende uma função que faz o mapeamento da entrada para a saída.

A figura 5 ilustra um processo de aprendizado não-supervisionado, onde faz-se o agrupamento dos vetores de entrada com base em alguma medida de similaridade entre suas características. Para um novo dado, a saída do modelo preditivo pode ser uma probabilidade, uma indicação do grupo mais provável ou uma representação melhorada deste dado, com base no conhecimento adquirido durante o treinamento.

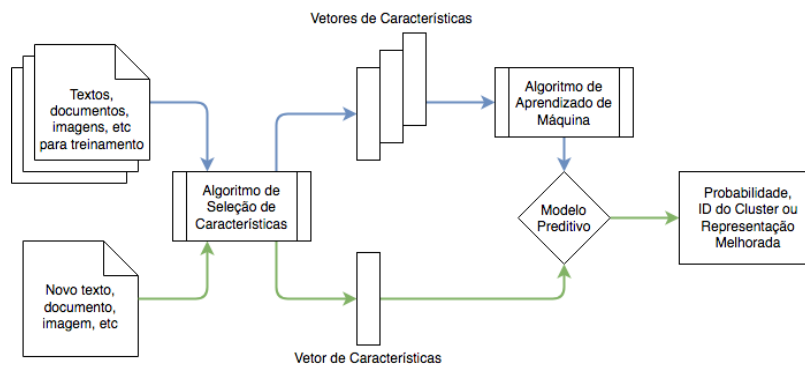


Figura 5: Aprendizado não-supervisionado.

A figura 6 ilustra um processo de aprendizado supervisionado, onde cada par vetor-rótulo de entrada é analisado de forma a obter uma função capaz de mapear um novo dado de entrada a um dos rótulos conhecidos pelo algoritmo. Para um novo dado, portanto, a saída do modelo preditivo é um rótulo que indica a qual classe este dado provavelmente pertence, com base no conhecimento adquirido durante o treinamento.

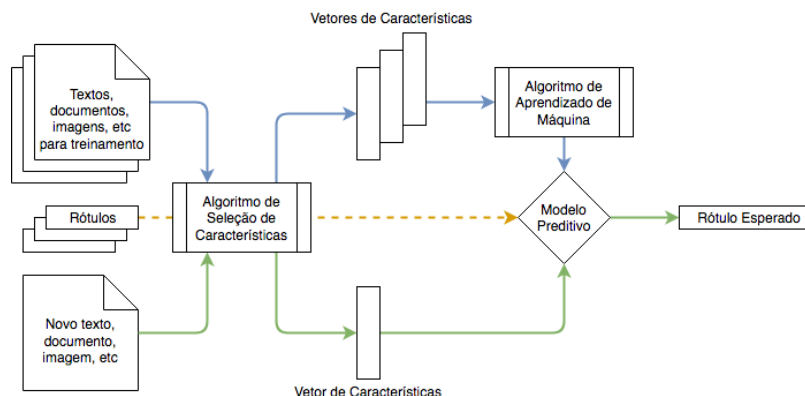


Figura 6: Aprendizado supervisionado.

2.2 Categorização de texto

A categorização de texto é tarefa de atribuir categorias predefinidas a documentos de texto livre. Ela pode ser útil nos mais diversos contextos e aplicações do mundo real. A tabela 1 mostra como este conceito pode ser aplicado a certos tipos de texto (YANG; JOACHIMS, 2008):

Tabela 1: Exemplos de tipos de texto e sugestões de classificação.

Tipo de texto	Exemplos de categorização
Notícias	Assunto, tópico ou código geográfico
Trabalhos acadêmicos	Domínios técnicos e subdomínios
Prontuários médicos	Tipo de doença, procedimentos cirúrgicos e medicação
<i>E-mail</i>	<i>Ham</i> , <i>spam</i> , pessoais ou de negócios
Redes sociais	Postagens casuais, difamatórias e <i>fake news</i>

Em vez de manualmente classificar documentos ou criar regras de classificação automática, é possível utilizar métodos de aprendizado de máquina para aprender regras de classificação automática com base em documentos de treinamento que foram previamente categorizados por humanos.

Para isto, é necessário primeiramente representar estes documentos (que, no contexto deste trabalho, consistem em mensagens de *e-mail*) em uma forma passível de processamento computacional. Esta representação, que consiste em um conjunto de características do documento, pode ser simplificada por meio de métodos de seleção de características. Alguns destes métodos são apresentadas na seção 2.3. Um estudo mais aprofundado pode ser encontrado em Forman e George (2003).

Em seguida, entram em cena os modelos de aprendizado de máquina que, por meio de treinamento (usando vetores de características extraídas de amostras de texto), são capazes de classificar textos, *e-mails* etc. Vários destes modelos são apresentados na seção 2.6. Um levantamento mais abrangente pode ser encontrado em Aggarwal e Zhai (2012).

2.3 Seleção de características

Os problemas de categorização de texto normalmente apresentam espaços de característica de alta dimensionalidade, já que são compostos por todos os termos únicos (palavras, por exemplo) que ocorrem nos documentos do conjunto. Este grande volume de dados de entrada, que muitas vezes inclui características redundantes ou irrelevantes, frequentemente atinge a casa de centenas ou milhares de características, podendo dificultar (ou até mesmo impedir) o treinamento e a geração de modelos com certos algoritmos de aprendizado de máquina.

Diante deste problema, mostra-se necessário efetuar uma redução do conjunto de características, de modo que os algoritmos não gastem quantidades impraticáveis de tempo e recursos (armazenamento e processamento) computacionais. Além disso, deseja-se escolher um subconjunto de características que retenha a maior quantidade possível de informação do conjunto original, de forma a não prejudicar a capacidade de generalização e precisão de categorização do modelo.

Em Yang e Pedersen (1997), são revisitados e comparados diversos métodos estatísticos de seleção de características, especificamente aplicados à categorização de textos, tais como “seleção de termos baseada na frequência no documento, ganho de informação, informação mútua, teste χ^2 e força do termo”. É importante ressaltar que os termos não são necessariamente palavras, mas qualquer elemento presente no *e-mail* após o pré-processamento (ex: *tokens* para caracterizar *tags* HTML (do inglês *HyperText Markup Language*), números, pontuação, *links* etc.).

Na literatura, também é possível encontrar métodos de extração de características que visam criar conjuntos de novas características definidas em função das originais. Apesar de tais métodos fugirem ao escopo deste trabalho, podem ser citadas as abordagens baseadas em análise semântica latente (DEERWESTER et al., 1990) (HOFMANN, 1999), mapeamento isométrico de características (TENENBAUM; SILVA; LANGFORD, 2000) e incorporação estocástica de vizinhos (HINTON; ROWEIS, 2002) (MAATEN; HINTON, 2008).

A seguir, serão explicados os três métodos de seleção de características — estatística χ^2 , frequência no documento e informação mútua — utilizados neste trabalho.

2.3.1 Estatística Qui-quadrado (CHI2)

A estatística χ^2 (ou simplesmente CHI2) mede a dependência entre um termo t e uma classe c em particular. É definida pela equação (2.1).

$$CHI2(t, c) = \frac{n \cdot F(t)^2 \cdot (p_c(t) - P_c)^2}{F(t) \cdot (1 - F(t)) \cdot P_c \cdot (1 - P_c)} \quad (2.1)$$

Na qual:

- n é o número total de documentos na coleção;
- $p_c(t)$ é a probabilidade condicional c para documentos que contenham t ;
- P_c é a fração global de documentos que contém a classe c ;
- $F(t)$ é a fração global de documentos que contem o termo t .

Uma característica deste método é que os valores $CHI2(t, c)$ são normalizados e, portanto, são comparáveis entre termos de uma mesma categoria.

2.3.2 Frequência no Documento (FD)

É um método estatístico que mede a frequência que um termo t ocorre em uma classe c , em particular. É definido pela equação (2.2).

$$FD(t, c) = \frac{C(t, c)}{\sum_{t_i \in T(c)} C(t_i, c)} \quad (2.2)$$

Na qual:

- o numerador é a quantidade de ocorrências de t nos documentos da classe c ;
- o denominador é o somatório das quantidades de ocorrências de todos os termos nos documentos da classe c .

2.3.3 Informação Mútua (MI)

Derivada da teoria da informação (COVER; THOMAS, 1991), esta medida consiste na quantidade de informação que um termo agrega em relação a uma dada classe. A informação mútua pontual $M(t, c)$ entre o termo t e a classe c é baseada no nível de coocorrência entre a classe c e o termo t , dada pela equação (2.3).

$$MI(t, c) = \log \left(\frac{F(t) \cdot p_c(t)}{F(t) \cdot P_c} \right) = \log \left(\frac{p_c(t)}{P_c} \right) \quad (2.3)$$

Na qual:

- $F(t) \cdot P_c$ é a coocorrência esperada da classe x e termo t com base na independência mútua;
- $F(t) \cdot p_c(t)$ é a coocorrência verdadeira (na prática, esse valor pode ser muito maior ou menor que o esperado, dependendo do nível de correlação entre a classe x e o termo t).

Claramente, o termo t é positivamente correlacionado à classe c quando $MI(t, c) > 0$, e negativamente correlacionada quando $MI(t, c) < 0$.

Note que as equações (2.1), (2.2) e (2.3) são computadas para uma única classe. Para criar um valor que abranja todas as classes, podem ser usados os valores médios ou máximos, como nas equações (2.4), (2.5), (2.6), (2.7), (2.8) e (2.9):

$$CHI2_{m\u00e9dio}(t) = \sum_{c \in C} P_i \cdot \chi^2(t, c) \quad (2.4) \quad CHI2_{m\u00e1ximo}(t) = \max_{c \in C} \chi^2(t, c) \quad (2.5)$$

$$FD_{m\u00e9dio}(t) = \frac{1}{|C|} \sum_{c \in C} FD(t, c) \quad (2.6) \quad FD_{m\u00e1ximo}(t) = \max_{c \in C} FD(t, c) \quad (2.7)$$

$$MI_{m\u00e9dio}(t) = \sum_{c \in C} P_i \cdot MI(t, c) \quad (2.8) \quad MI_{m\u00e1ximo}(t) = \max_{c \in C} MI(t, c) \quad (2.9)$$

2.4 Redução de dimensionalidade

Um problema central no aprendizado de máquina é identificar um conjunto representativo de atributos (características) para construir um modelo de classificação para uma tarefa específica. Em modelos preditivos, atributos em excesso podem confundir o classificador, tornando a redução de dimensionalidade uma etapa importante no processo de aprendizado.

A “maldição da dimensionalidade” refere-se a vários fenômenos que surgem ao analisar e organizar dados em espaços de alta dimensão (muitas vezes com centenas ou milhares de dimensões) que não ocorrem em espaços de baixa dimensão, como o espaço físico tridimensional da experiência cotidiana. (DONOHO, 2000)

No aprendizado de máquina, muitos problemas exigem grandes quantidades de dados de treinamento para garantir que existam várias amostras com valores distintos. Com um número fixo de amostras de treinamento, o poder preditivo de um classificador ou regressor inicialmente aumenta à medida que o número de dimensões (características) usados é aumentado, mas depois decai. Isto é conhecido como “fenômeno de Hughes” (HUGHES, 1968), e é ilustrado na figura 7.

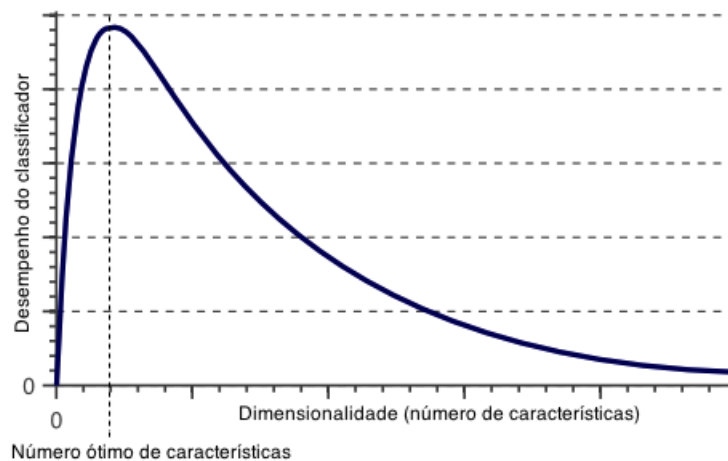


Figura 7: Maldição da dimensionalidade (traduzido) (SPRUYT, 2014).

Diversas técnicas de redução de dimensionalidade foram propostas, tais como:

- Análise de componentes principais (PCA, do inglês *Principal Component Analysis*): utiliza a ortogonalização de vetores para converter um conjunto de variáveis possivelmente correlacionadas num conjunto de variáveis linearmente não correlacionadas (chamadas de componentes principais). (F.R.S., 1901)
- Consistência de subconjunto (LVF, do inglês *Las Vegas Filter*): baseia-se em um algoritmo de busca aleatório que demonstra alta robustez (em termos de atributos redundantes ou correlações de ordem alta) quando utilizado com conjuntos de dimensionalidade alta. (LIU; SETIONO, 1996)
- Correlação de subconjunto (CFS, do inglês *Correlation-based Feature Selection*): baseia-se na hipótese central de que bons conjuntos de características possuem características altamente correlacionadas com a classe, ainda que não correlacionadas entre si. (HALL, 2000)
- Máquina de vetores de suporte (SVM-RFE, do inglês *Support Vector Machine Recursive Feature Elimination*): faz uso de SVMs com eliminação recursiva de características para efetuar redução de dimensionalidade, descartando atributos redundantes e produzindo subconjuntos compactos e, ao mesmo tempo, com grande capacidade de predição. (GUYON et al., 2002)
- Filtro de silhueta simplificada (SSF, do inglês *Simplified Silhouette Filter*): consiste em agrupar atributos em subconjuntos de cardinalidade estimada automaticamente, podendo considerar correlações entre características e classes. (COVÕES; HRUSCHKA, 2011)

Neste trabalho, foi adotada a técnica CFS devido a sua grande capacidade de redução de dimensionalidade em uma variedade de conjuntos de dados (com características tanto discretas quanto contínuas), ao mesmo tempo que mantém ou aumenta o desempenho dos algoritmos de aprendizado. (HALL, 2000)

O algoritmo CFS primeiro calcula uma matriz de correlações entre características-classes e características-características a partir do conjunto de treinamento e, em seguida, realiza uma busca pelo subconjunto com maior nota. Para esta última tarefa, foi utilizado o algoritmo evolucionário multiobjetivo ENORA (do inglês *Evolutionary Non-dominated Radial slots based Algorithm*) (JIMENEZ et al., 2002).

2.5 Balanceamento de classes

A finalidade desta etapa é assegurar que haja a mesma quantidade de amostras de *hams* e *spams*. Isto porque um conjunto de dados desbalanceado (com quantidades diferentes de instâncias de cada classe) pode fazer com que o método de aprendizado de máquina enfatize mais uma classe do que outra, obtendo resultados enviesados. Para resolver esta situação, pode-se selecionar aleatoriamente mensagens no conjunto de dados de menor cardinalidade e replicá-las até que os conjuntos possuam o mesmo tamanho. Considere o conjunto de dados da tabela 2, onde k é o índice da amostra, x_i^k é o valor da i -ésima característica da k -ésima amostra e y^k é a classe da k -ésima amostra.

Tabela 2: Conjunto de dados antes do balanceamento.

k	x_0^k	x_1^k	x_2^k	x_3^k	y^k
1	0,24	0,59	0,81	0,04	<i>ham</i>
2	0,66	0,24	0,86	-0,11	<i>ham</i>
3	0,48	0,15	0,76	0,68	<i>ham</i>
4	0,53	0,39	0,34	0,14	<i>ham</i>
5	0,08	0,56	0,01	0,55	<i>ham</i>
6	1,14	0,25	0,81	1,13	<i>spam</i>
7	-0,44	-0,54	-0,10	0,78	<i>spam</i>
8	0,36	0,59	0,30	0,91	<i>spam</i>

Neste caso, a classe *ham* possui cinco amostras, enquanto a *spam* possui apenas três. Para resolver esta situação, escolhem-se aleatoriamente $5 - 3 = 2$ amostras de *spam*, que são replicadas no conjunto de dados visando balancear as cardinalidades (vide tabela 3). Assim, as classes passam a quantidades iguais de amostras, equalizando as influências de cada uma perante o algoritmo de aprendizado.

Tabela 3: Amostras replicadas a fim de balancear o conjunto de dados.

k	x_0^k	x_1^k	x_2^k	x_3^k	y^k
9	-0,44	-0,54	-0,10	0,78	<i>spam</i>
10	1,14	0,25	0,81	1,13	<i>spam</i>

2.6 Algoritmos de classificação

A classificação consiste em, utilizando um conjunto de dados com amostras rotuladas e um algoritmo de aprendizado de máquina, criar um modelo capaz de atribuir classes a dados inéditos. Existem diversos modelos de classificação, e uma grande variedade é estudada de forma abrangente em Bishop (2006). No escopo deste trabalho, fazem-se relevantes os métodos discriminados na listagem a seguir, divididos em quatro categorias:

- Métodos probabilísticos (subseção 2.6.1): consiste em classificadores capazes de prever, dada uma observação de entrada, uma distribuição de probabilidade em um conjunto de classes, em vez de apenas produzir a classe mais provável a que a observação deveria pertencer.
- Redes neurais artificiais (subseção 2.6.2): são modelos computacionais inspirados pelo sistema nervoso central de um animal, capazes de realizar o aprendizado de máquina bem como o reconhecimento de padrões.
- Árvores de decisão (subseção 2.6.3): é uma ferramenta de suporte à decisão que usa um grafo ou modelo de decisões em árvore e suas possíveis consequências; são comumente usados na pesquisa de operações para ajudar a identificar uma estratégia com maior probabilidade de atingir um objetivo, sendo também uma ferramenta popular na aprendizagem de máquinas.
- Máquina de vetores de suporte (subseção 2.6.4): é um classificador discriminativo formalmente definido por um hiperplano separador; em outras palavras, com dados de treinamento rotulados (aprendizagem supervisionada), o algoritmo produz um ótimo hiperplano que categoriza novos exemplos.

Algoritmos que não foram abordados neste trabalho mas que merecem destaque incluem regressão logística, modelos ocultos de Markov, aprendizado por quantização vetorial (KOHONEN, 1995) e *gradient boosting* (FRIEDMAN, 2001).

2.6.1 Métodos Probabilísticos

2.6.1.1 Naïve Bayes (NB)

Naïve Bayes é um algoritmo para classificação conhecido tanto por sua simplicidade teórica e facilidade de implementação, quanto também pela eficácia (especialmente na categorização de textos), já que é capaz de construir modelos precisos e fazer previsões rapidamente. Pode ser utilizado em diversos cenários, tais como filtragem de *spam*, análise sentimental, categorização de notícias, diagnósticos médicos, previsão do tempo e reconhecimento facial e de caracteres.

O algoritmo aprende a probabilidade de um objeto com determinadas características pertencer a uma determinada classe ou categoria; por esta razão, é tido como um “classificador probabilístico”. Além disso, é dito “bayesiano” por ter como base o teorema de Bayes e “ingênuo” por supor que a ocorrência de uma determinada característica independe da ocorrência e/ou influência de outras características.

Proposto pelo matemático Thomas Bayes em 1763 e aprimorado pelo matemático, astrônomo e físico Pierre-Simon Laplace em 1812, o teorema (dado pela equação (2.10)) descreve a probabilidade de ocorrência de um evento com base no conhecimento prévio de condições que possam estar relacionadas ao evento.

$$P(A|B) = \frac{P(B|A) P(A)}{P(B)} \quad (2.10)$$

em que:

- A e B são eventos e $P(B) \neq 0$;
- $P(A|B)$ é a probabilidade da ocorrência do evento A dado que B ocorreu;
- $P(A)$ e $P(B)$ são as probabilidades de ocorrência de A e B , respectivamente;
- $P(B|A)$ é a probabilidade da ocorrência do evento B dado que A ocorreu.

Dado um vetor de dados $\vec{x} = \{x_1, \dots, x_n\}$ com n características, o classificador prediz a classe C_k para \vec{x} de acordo com a equação (2.11) (MCGONAGLE, 2018):

$$p(C_k|\vec{x}) = p(C_k|x_1, \dots, x_n) \text{ para } k = 1, \dots, K \quad (2.11)$$

Usando o teorema de Bayes, a equação (2.11) pode ser fatorada como:

$$p(C_k|\vec{x}) = \frac{p(\vec{x}|C_k)p(C_k)}{p(\vec{x})} = \frac{p(x_1, \dots, x_n|C_k)p(C_k)}{p(x_1, \dots, x_n)} \quad (2.12)$$

Usando a regra da cadeia, o fator $p(x_1, \dots, x_n|C_k)$ no numerador da equação (2.12) pode ser decomposto ainda mais em:

$$p(x_1, \dots, x_n|C_k) = p(x_1|x_2 \dots x_n, C_k)p(x_2|x_3 \dots x_n, C_k) \dots p(x_{n-1}|x_n, C_k)p(x_n|C_k) \quad (2.13)$$

Aplicando a suposição ingênua de independência condicional (ou seja, assumindo que a característica x_i é independente da característica x_j para $i \neq j$ dada a classe C_k) à equação (2.13), tem-se:

$$p(x_1, \dots, x_n|C_k) = p(x_i|C_k) \Rightarrow p(x_1, \dots, x_n|C_k) = \prod_{i=1}^n p(x_i|C_k) \quad (2.14)$$

Desta forma,

$$\begin{aligned} p(C_k|x_1, \dots, x_n) &\propto p(C_k, x_1, \dots, x_n) \\ &\propto p(C_k)p(x_1, \dots, x_n|C_k) \\ &\propto p(C_k)p(x_1|C_k)p(x_2|C_k) \dots p(x_n|C_k) \\ &\propto p(C_k) \prod_{i=1}^n p(x_i|C_k) \end{aligned} \quad (2.15)$$

O classificador NB combina o modelo descrito pelas equações anteriores com uma regra de decisão, como a estimativa de probabilidade máxima *a posteriori*, que consiste na escolha da hipótese mais provável. O classificador correspondente é a função que atribui a classe $\hat{C} = C_k$ para algum k como descrito na equação (2.16):

$$\hat{C} = \operatorname{argmax}_{k \in \{1, \dots, K\}} p(C_k) \prod_{i=1}^n p(x_i | C_k). \quad (2.16)$$

Além do NB, há outras abordagens que fazem uso da suposição de independência. No trabalho de McCallum e Nigam (1998), por exemplo, são comparados os modelos multivariado de Bernoulli e multinomial aplicados à classificação de texto.

2.6.1.2 Averaged One-Dependence Estimators (AODE)

Conforme explicado na seção 2.6.1.1, o classificador NB baseia-se na suposição ingênua de independência entre atributos. Entretanto, em conjuntos obtidos de contextos em que os atributos são fortemente relacionados, tal premissa pode acarretar em altas taxas de erro quando da classificação de amostras. Neste sentido, Sahami (1996) introduz a noção de estimadores de x dependências, por meio dos quais a probabilidade do valor de cada atributo é condicionada pela classe e, no máximo, outros x atributos, induzindo um classificador “menos ingênuo” que o NB.

Obviamente, a introdução de dependências entre atributos aumenta exponencialmente a quantidade de poder computacional necessária para estimar as probabilidades condicionais. Entretanto, estes dados podem ser pré-computados uma única vez (ou seja, antes do treinamento propriamente dito), minimizando assim o *overhead* causado por essa suposição menos ingênua.

Neste trabalho, será adotado $x = 1$, ou seja, o foco será em “estimadores de uma dependência” (AODE). Neste cenário, um problema rapidamente vem à tona: um classificador deste tipo requer que cada atributo dependa de outro atributo, e tal atributo específico deve ser necessariamente selecionado. A solução para este impasse consiste em considerar todos os classificadores de uma dependência e agregar as previsões de todos eles, desde que estejam qualificados para tal.

Sucintamente, um classificador é qualificado para contribuir para a previsão final somente se possuir um único atributo que seja pai de todos os outros e se houver no mínimo m amostras no conjunto de treinamento com o atributo pai preenchido, ou seja, com um valor não-nulo. Este último critério visa incluir modelos para os quais as estimativas básicas de probabilidades são imprecisas. Além disso, no mesmo trabalho, adota-se $m = 30$, “por ser uma quantidade mínima de amostras amplamente utilizada para propósitos de inferência estatística” (WEBB; BOUGHTON; WANG, 2005).

As figuras 8a e 8b ilustram o funcionamento do NB e do AODE, respectivamente. Na primeira, nota-se que as probabilidades de classe $P(c)$ dependem diretamente dos atributos x_1, x_2, \dots, x_n , e nenhuma probabilidade entre eles é considerada. Já na segunda, $P(c)$ está condicionada aos atributos que, por sua vez, levam em consideração as probabilidades conjuntas em relação a um único outro atributo, caracterizando assim um estimador de uma dependência.

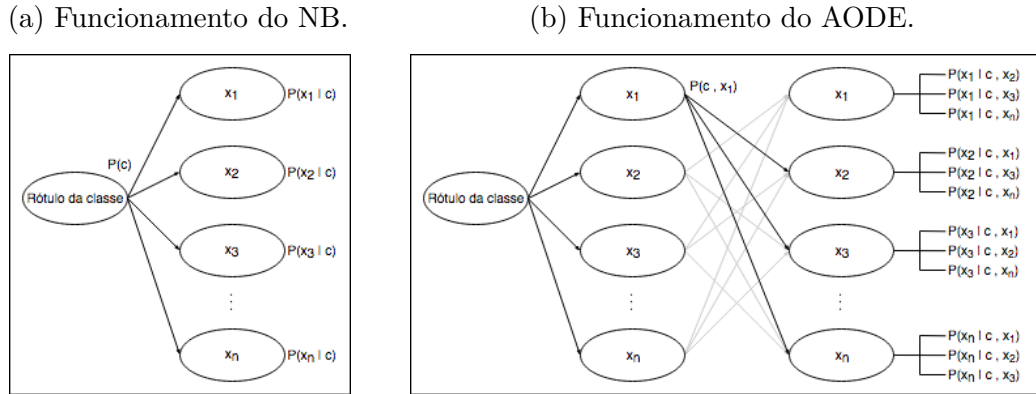


Figura 8: Comparação entre dependências dos classificadores NB e AODE.

O funcionamento do AODE consiste, portanto, em estimar as probabilidades de uma amostra \vec{x} com valores de características x_1, x_2, \dots, x_n pertencer a cada classe y_i e, em seguida, eleger a classe final \hat{y} como aquela que maximiza a equação (2.17). Nela, o termo $F(x_i)$ representa a quantidade de exemplos que possuem o valor de atributo x_i , e é usado para assegurar o limite m de amostras necessárias para um dado classificador poder contribuir para a estimativa de probabilidade (WEBB; BOUGHTON; WANG, 2005).

$$\hat{y} = \operatorname{argmax}_y \left(\sum_{i:1 \leq i \leq n \wedge F(x_i) \geq m} P(y, x_i) \prod_{j=1}^n P(x_j | y, x_i) \right) \quad (2.17)$$

É importante destacar que, se não existir um classificador que satisfaça a condição de qualificação, a consideração de uma dependência não será levada em conta e, portanto, o método AODE comportar-se-á como o NB.

2.6.2 Redes Neurais Artificiais

As redes neurais artificiais são modelos matemáticos compostos por unidades discretas de processamento, cujo funcionamento baseia-se no de um neurônio biológico. Quando dispostos em uma ou mais camadas neurais associadas por meio de conexões direcionadas, tais elementos passam a caracterizar uma rede neural, estruturas frequentemente comparadas ao cérebro humano.

A figura 9 apresenta um modelo simples do neurônio artificial (MCCULLOCH; PITTS, 1943). Se a combinação linear de suas entradas exceder um limiar pré determinado, o neurônio produzirá um potencial excitatório em sua saída; caso contrário, o neurônio produzirá um potencial inibitório.

Cada unidade de uma rede neural possui uma entrada fictícia $a_0 = 1$ como um peso $w_{0,j}$ associado a ela. Uma ligação da unidade i para a unidade j serve para propagar a ativação a_i de i para j , e o peso desta ligação (chamado de peso sináptico, já que as “sinapses” são as regiões onde ocorre a transmissão de impulsos nervosos em um neurônio biológico), representado por $w_{i,j}$, determina a força e o sinal da conexão. Cada unidade j primeiro calcula uma soma ponderada de suas entradas e, em seguida, aplica uma função de ativação g a essa soma para obter a saída $a_j = g(in_j) = g(\sum_{i=0}^n w_{i,j}a_i)$. (RUSSELL; NORVIG, 2013)

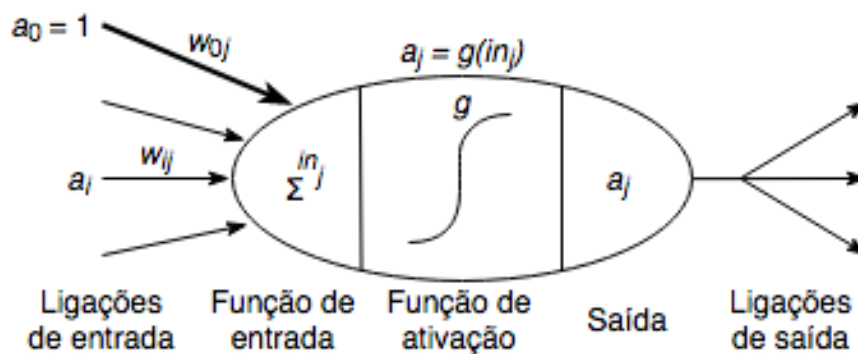


Figura 9: Neurônio artificial de McCulloch–Pitts (MCP) (traduzido).

A função g é tipicamente uma das listadas nos gráficos da figura 10.

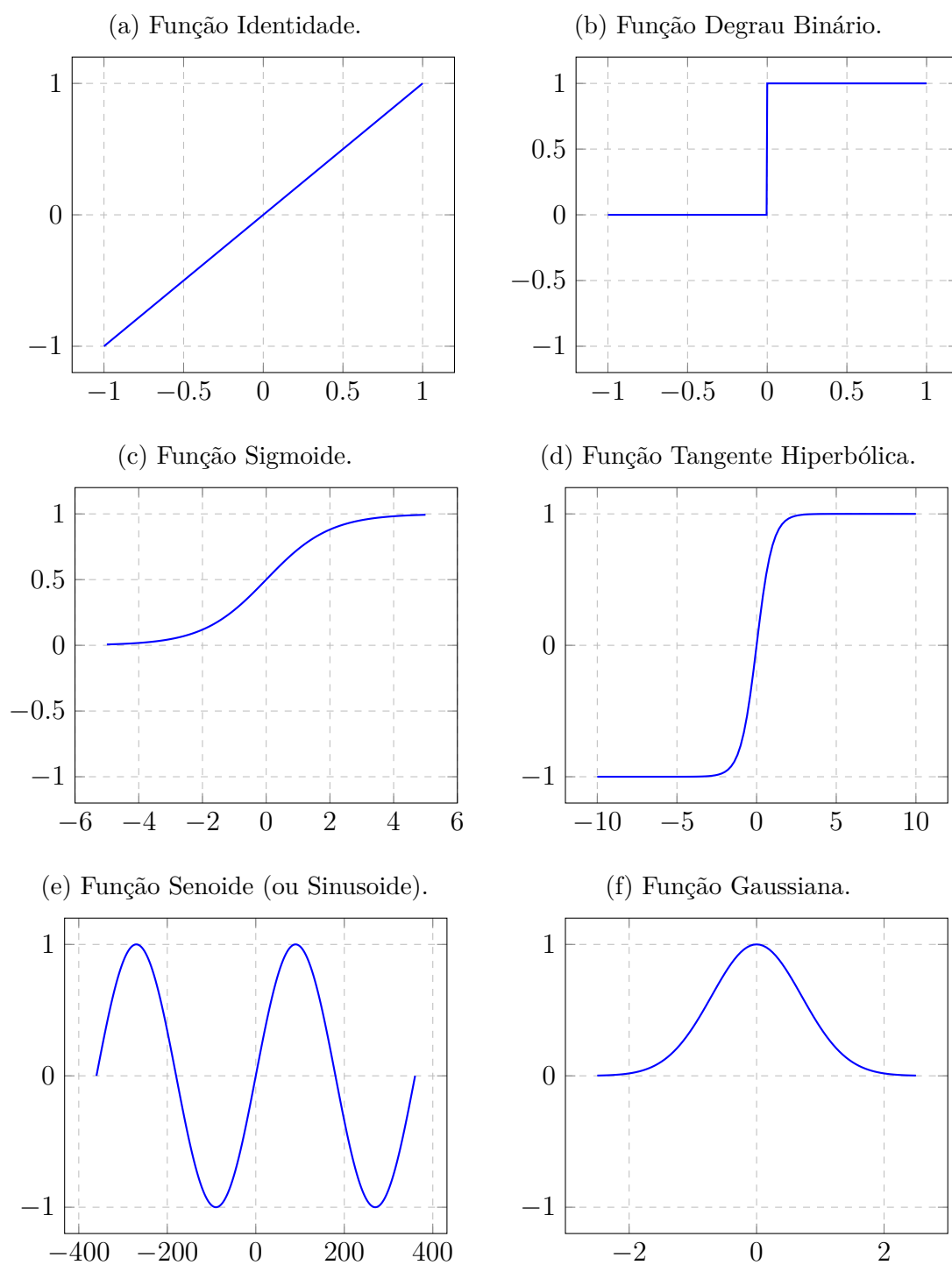


Figura 10: Exemplos de funções de ativação.

2.6.2.1 Single Layer Perceptron (SLP)

Uma rede neural com todas as entradas conectadas diretamente às saídas é dita simples, de camada única ou rede Perceptron. A figura 11 mostra um exemplo deste tipo de rede com duas entradas (x_1 e x_2) e uma saída (y).

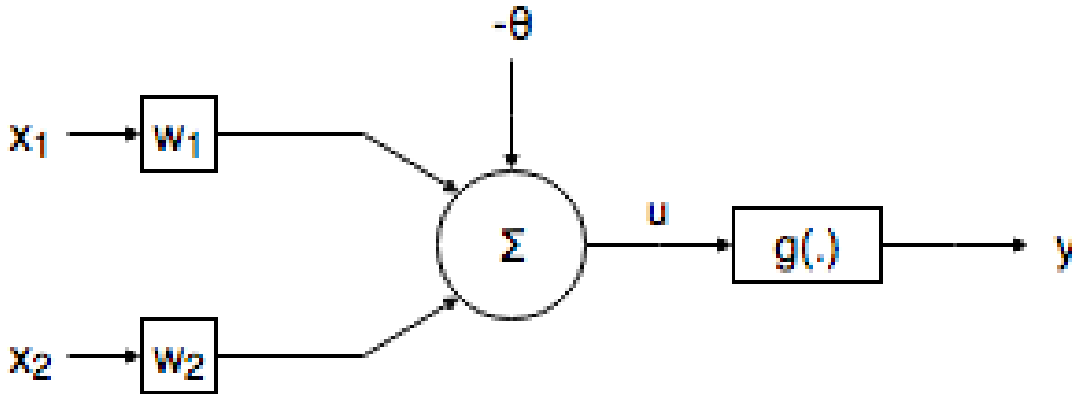


Figura 11: Rede Perceptron de camada única.

A capacidade de cálculo do Perceptron está limitada pela separabilidade linear (vide figuras 12a e 12b) das classes. Sua incapacidade em lidar com classes distribuídas de tal maneira que seja impossível criar um hiperplano para separá-las (vide figura 12c) foi um revés significativo para as redes neurais na década de 1960. Nestes casos, faz-se necessário utilizar abordagens mais sofisticadas, como uma rede de função de base radial, que é discutida adiante na seção 2.6.2.2.

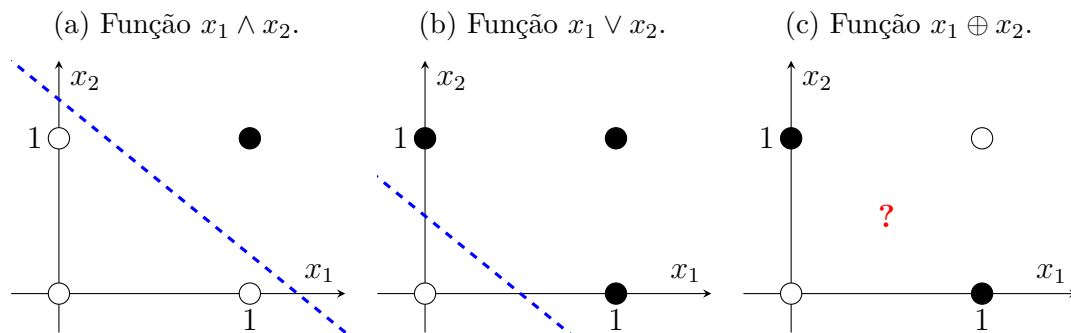


Figura 12: Funções binárias e separabilidade linear.

O treinamento de uma rede Perceptron consiste em ajustar os valores dos pesos e limiar, e baseia-se na Teoria Hebbiana (HEBB, 1949). Resumidamente, se a saída produzida pelo Perceptron difere da desejada, seus pesos sinápticos e limiares serão então ajustados proporcionalmente aos valores dos sinais de entrada. Este processo é repetido para todas as amostras de treinamento até que a saída produzida pela rede seja similar à saída desejada de cada amostra (ou até que outro critério de parada, por exemplo quantidade de épocas de treinamento, seja alcançado).

Em termos matemáticos, as regras de ajuste dos pesos sinápticos \vec{w} e do limiar θ do neurônio podem ser expressas pela equação (2.18):

$$\vec{w}_i^{actual} = \vec{w}_i^{anterior} + \eta \cdot (d^{(k)} - y) \cdot \vec{x}^{(k)} \quad (2.18)$$

Na qual:

- $\vec{w} = [\theta \ w_1 \ w_2 \ \dots \ w_n]^\top$ é o vetor de pesos sinápticos;
- $\vec{x}^{(k)} = [-1 \ x_1^{(k)} \ x_2^{(k)} \ \dots \ x_n^{(k)}]^\top$ é a k -ésima amostra de treinamento;
- $d^{(k)}$ é o valor desejado para a k -ésima amostra de treinamento;
- y é o valor de saída produzido pelo Perceptron;
- η é uma constante que define a taxa de aprendizagem da rede.

A taxa de aprendizagem η “exprime o quão rápido o processo de treinamento da rede será conduzido rumo à convergência (estabilização) e, normalmente, são adotados valores entre 0 e 1” (SILVA; SPATTI; FLAUZINO, 2010).

Outra técnica frequentemente adotada para treinar este tipo de rede baseia-se no algoritmo LMS (do inglês *Least Mean Squares*) (WIDROW; HOFF, 1960). Esta abordagem visa ajustar os pesos e limiar do neurônio de forma a minimizar o erro da rede (ou seja, as diferenças entre as saídas desejadas e as respostas do combinador linear) na classificação das amostras de treinamento. O erro da rede é normalmente medido usando a função EQM (Erro Quadrático Médio), dada pela equação (2.19).

$$E(\vec{w}) = \frac{1}{2} \sum_{k=1}^p (d^{(k)} - (\vec{w}^\top \cdot \vec{x}^{(k)} - \theta))^2 \quad (2.19)$$

2.6.2.2 Radial Basis Function Network (RBF)

As redes Perceptron de múltiplas camadas (WHITE; ROSENBLATT, 1963) são caracterizadas pela presença de pelo menos uma camada intermediária (escondida) de neurônios, situada entre as camadas de entrada e de saída. Este tipo de rede pode ser aplicado em diversos tipos de problemas, tais como aproximação universal de funções, reconhecimento de padrões, identificação e controle de processos, previsão de séries temporais e otimização de sistemas (SILVA; SPATTI; FLAUZINO, 2010).

As redes denominadas de funções de base radial (BROOMHEAD; LOWE, 1988), convencionalmente conhecidas como RBF (do inglês *Radial Basis Function*), assim como as Perceptron de camadas múltiplas, pertence à arquitetura *feedforward*, já que o fluxo de informações na sua estrutura se inicia na camada de entrada, percorre a camada intermediária única (neurônios com funções de ativação $g_1^{(1)}, g_2^{(1)}, \dots, g_{n_1}^{(1)}$ do tipo gaussiana) e finaliza na camada neural de saída (neurônios com funções de ativação $g_1^{(2)}, \dots, g_{n_2}^{(2)}$ lineares). A figura 13 ilustra a topologia deste tipo de rede.

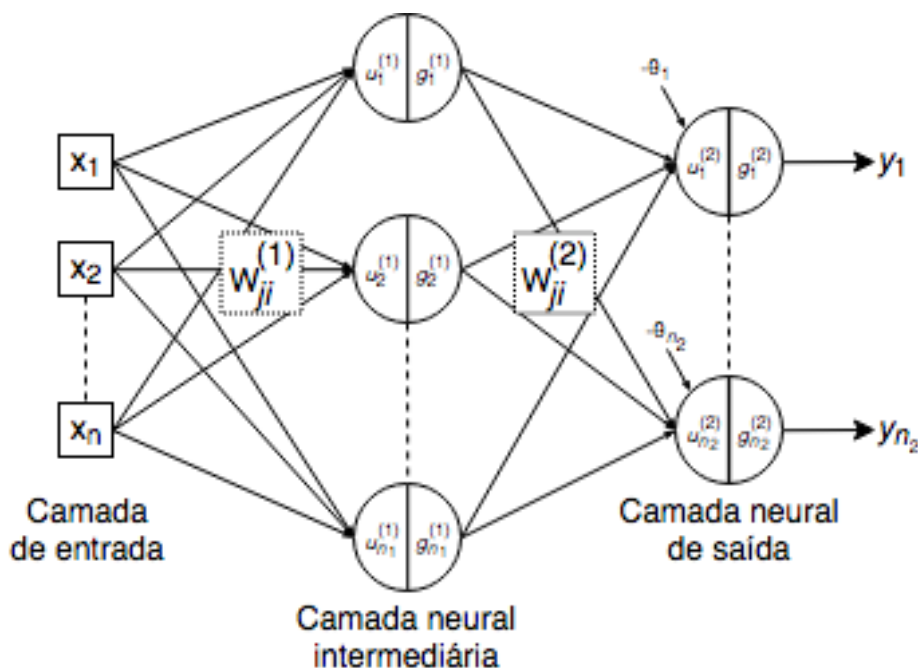


Figura 13: Rede de função de base radial.

A estratégia de treinamento utilizada para o ajuste dos pesos de suas duas únicas camadas neurais consiste de dois estágios:

1. Ajustes dos pesos dos neurônios da camada intermediária usando um método de aprendizagem não-supervisionado (ou seja, que só utiliza as características dos dados de entrada) para posicionar as funções de base radial.
2. Ajustes dos pesos dos neurônios da camada de saída, utilizando um critério de aprendizagem baseado em retropropagação do erro, frequentemente referido como regra delta generalizada (RUMELHART; HINTON; WILLIAMS, 1985).

Além da estratégia de treinamento recém citada, cujos detalhes são apresentados ainda nesta seção, foi proposta na literatura uma versão totalmente supervisionada deste tipo de rede, nomeada GRBF (do inglês *Generalized Radial Basis Functions*) (POGGIO; GIROSI, 1989) (WETTSCHERECK; DIETTERICH, 1992).

Etapa 1: treinamento não-supervisionado Como dito anteriormente, este estágio consiste em ajustar os pesos dos neurônios da camada intermediária, cujas funções de ativação são do tipo gaussiana (representadas pela equação (2.20)). Nelas, o parâmetro c define o centro da função gaussiana e σ^2 denota a sua variância, que indica a distância de um potencial de ativação u em relação ao centro c .

$$g(u) = e^{\left(\frac{-(u-c)^2}{2\sigma^2}\right)} \quad (2.20)$$

Portanto, os parâmetros a serem ajustados neste estágio são justamente o centro c e a variância σ_2 das funções gaussianas. No contexto da rede neural, “o centro c está diretamente associado aos seus próprios pesos, sendo que a entrada $u_j^{(1)}$ de cada um será o próprio vetor de entrada \vec{x} , que representa os n sinais externos que serão aplicados na rede” (SILVA; SPATTI; FLAUZINO, 2010). Consequentemente, a saída de cada neurônio j da camada intermediária é expressa pela equação (2.21).

$$g_j^{(1)}(u_j^{(1)}) = g_j^{(1)}(\mathbf{x}) = e^{\left(-\frac{\sum_{i=1}^n (x_i - w_{ji}^{(1)})^2}{2\sigma_j^2}\right)}, \text{ onde } j = 1, \dots, n_1 \quad (2.21)$$

Segundo Powell (1987), “quanto mais próxima uma determinada amostra estiver do centro de uma gaussiana, mais significativa será o valor produzido pelo campo receptivo radial da função de ativação”. Portanto, deve-se posicionar os centros das gaussianas da forma mais apropriada possível; isto pode ser feito com técnicas de agrupamento (*clustering*). Um destes algoritmos, chamado *k-means* (MACQUEEN, 1967), visa posicionar os centros de k gaussianas em regiões onde os padrões de entrada tenderão a se agrupar. No contexto da rede neural, “o valor do parâmetro k é igual ao número de neurônios da camada intermediária, já que a função de ativação de cada um é gaussiana” (SILVA; SPATTI; FLAUZINO, 2010).

Dado um conjunto de k gaussianas com centros em $\vec{W}_{1,j}^{(1)}, \vec{W}_{2,j}^{(1)}, \dots, \vec{W}_{k,j}^{(1)}$, o algoritmo atribui cada amostra de treinamento x ao *cluster* $\Omega^{(i)}$ mais próximo de acordo com a medida da distância Euclidiana. Em seguida, os centroides dos *clusters* são movidos para posição média das amostras que foram atraídas a ele. Este processo é repetido até que não haja mais movimentações a serem feitas.

Finalmente, os campos receptivos produzidos pelos pesos dos neurônios associados aos *clusters* $\vec{\Omega}$ poderão ser obtidos calculando-se a variância das funções de ativação por meio da distância quadrática média, dada pela equação (2.22).

$$\sigma_j^2 = \frac{1}{|\Omega^{(j)}|} \sum_{x^{(k)} \in \Omega^{(j)}} \sum_{i=1}^n (x_i^{(k)} - W_{ji}^{(1)})^2 \quad (2.22)$$

A figura 14 ilustra as iterações inicial, final e os raios de abrangência dos campos receptivos formados ao redor dos vetores de pesos ajustados pelo algoritmo *k-means*.

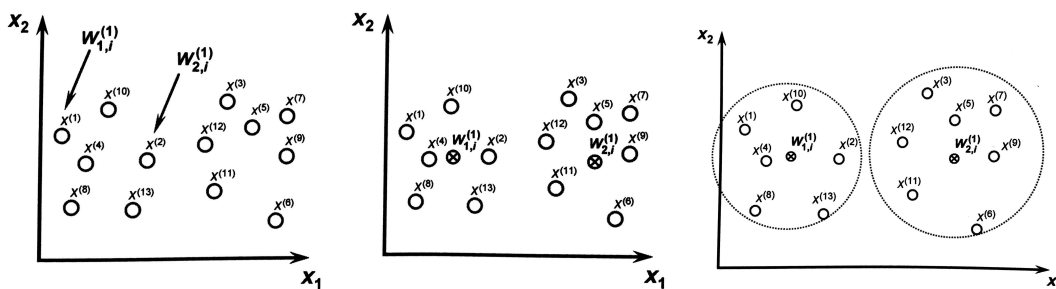


Figura 14: Iterações inicial, final e campos receptivos de uma rede RBF.

Etapa 2: treinamento supervisionado Como já foi dito, este estágio consiste em ajustar os pesos dos neurônios da camada de saída. Neste caso, diferentemente do estágio anterior, “o conjunto de treinamento será constituído por amostras rotuladas, em que as entradas serão as respostas das funções de ativação gaussianas dos neurônios da camada intermediária frente às respectivas amostras de treinamento” (SILVA; SPATTI; FLAUZINO, 2010), ou seja:

$$u_j^{(2)} = \sum_{i=1}^{n_1} W_{ji}^{(2)} \cdot g_i^{(1)}(u_i^{(1)}) - \theta_j, \text{ sendo } j = 1, \dots, n_2 \text{ e:} \quad (2.23)$$

- $W_{ji}^{(2)}$ os pesos dos neurônios da camada de saída;
- θ_j os limiares dos neurônios da camada de saída;
- $g_i^{(1)}(u_i^{(1)})$ as saídas da camada intermediária, calculadas usando os valores de $W_{ji}^{(1)}$ e $\sigma_j^{(2)}$ obtidos no primeiro estágio de treinamento.

Basicamente, os valores dos pesos são ajustados em direção oposta ao gradiente, a fim de minimizar o erro, como pode ser visto na equação (2.24):

$$\Delta W_{ji}^{(2)} = \eta \cdot \frac{\partial E}{\partial W_{ji}^{(2)}} = \eta \cdot (d_j - g_j^{(2)}) \cdot g' (u_j^{(2)}) = \eta \cdot (d_j - g (u_j^{(2)})) \cdot g' (u_j^{(2)}) \quad (2.24)$$

Na qual:

- $\Delta W_{ji}^{(2)}$ é o ajuste que será realizado nos pesos sinápticos da camada de saída;
- η é a taxa de aprendizado da rede, cujo valor situa-se entre 0 e 1;
- d_j é o valor esperado da j -ésima característica da amostra;
- $g_j^{(2)}$ é a saída do j -ésimo neurônio da camada de saída;
- $u_j^{(2)}$ é a entrada ponderada em relação ao j -ésimo neurônio de saída.

Finalmente, adotando-se uma função de ativação linear na camada de saída, seus neurônios realizarão uma combinação linear das gaussianas produzidas pelos neurônios da camada anterior. Portanto, a partir da equação (2.23) e com auxílio da figura 13, as respostas produzidas pelos neurônios de saída serão dadas por:

$$y_j = g_j^{(2)}(u_j^{(2)}) = u_j^{(2)}, \text{ sendo } j = 1, \dots, n_2 \quad (2.25)$$

Desta forma, os estágios de treinamento ajustam os parâmetros da rede RBF.

2.6.3 Árvores de Decisão

Há diversas maneiras de representar classificadores e talvez a árvore de decisão seja a mais comum. Originalmente, fora estudada nos campos de teoria de decisão e estatística, mas sua efetividade em outras disciplinas, tais como aprendizado de máquina e reconhecimento de padrões, vem se tornando cada vez mais notável.

Árvores de decisão são também implementadas em várias aplicações do mundo real, e em Quinlan (1986) são mencionadas as três a seguir:

- Diagnóstico de condição médica a partir de um ou mais sintomas;
- Determinação do valor (na teoria dos jogos) de uma posição de xadrez;
- Opinar, a partir de observações atmosféricas, sobre a possibilidade de chover.

Uma árvore de decisão é um classificador expresso como uma partição recursiva do espaço de instâncias. Consiste em nós que formam uma árvore enraizada (ou seja, há um nó "raiz" que não possui arestas de entrada). Um nó com arestas de saída é chamado de interno ou nó de teste. Todos os outros nós são chamados de folhas (também conhecidos como nós terminais ou de decisão).

Em árvores de decisão, os nós com arestas de saída (chamados de internos ou de teste) dividem o espaço de instâncias em dois ou mais subespaços com base nos valores dos atributos de entrada (ROKACH; MAIMON, 2005):

- No caso de atributos discretos/categóricos, cada teste considera um único atributo, particionando o espaço de acordo com o valor do atributo.
- No caso de atributos contínuos/numéricos, a condição refere-se a um intervalo.

Os nós do tipo folha (também chamados de terminais ou de decisão) representam as possíveis saídas que uma instância pode assumir. A classificação de uma instância dá-se pela travessia da árvore, começando pela raiz e terminando em uma folha, de acordo com a saída dos testes ao longo do caminho.

A figura 15 ilustra uma árvore de decisão que, baseada em atributos de salário, tempo de deslocamento e oferta de café grátis, decide se uma proposta de emprego

será aceita. Nela, nós de teste são representados por losangos rotulados com o atributo que testa, folhas por elipses e as saídas como rótulos das conexões.

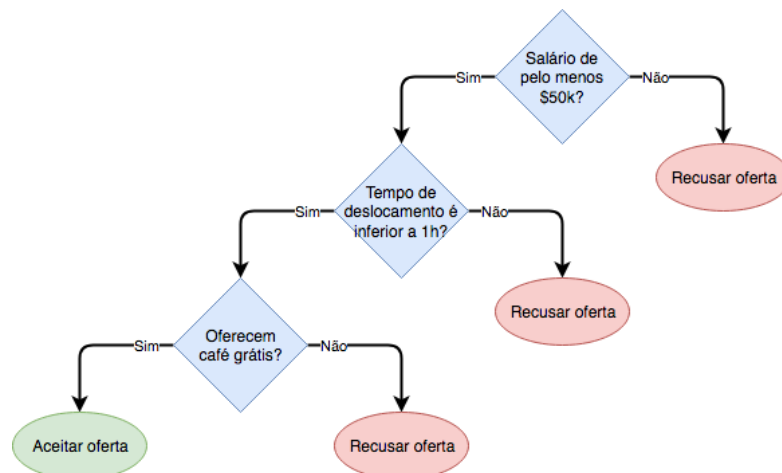


Figura 15: Árvore de decisão apresentando resposta à oferta de emprego.

Dispondo de um classificador deste tipo, um recrutador é capaz de prever a resposta de um candidato em particular e entender as características comportamentais de uma população inteira de candidatos.

O processo de criação de árvores de decisão a partir de um conjunto de dados de treinamento é chamado de indução. Diversas técnicas com esta finalidade foram propostas, tais como os algoritmos CART (BREIMAN et al., 1984), ID3 (QUINLAN, 1986), C4.5 (QUINLAN, 1993) e ADTrees (FREUND; MASON, 1999).

Um aspecto importante sobre árvores de decisão é sua complexidade estrutural. De acordo com (BREIMAN et al., 1984), “tal característica (normalmente medida em termos do número total de nós, folhas, quantidade de atributos usados e profundidade total) tem um efeito crucial na precisão de classificação da árvore”.

Há diversas técnicas que visam simplificar árvores excessivamente complexas. Os chamados algoritmos de poda visam aumentar tanto a compreensibilidade de uma árvore de decisão para um usuário humano quanto sua capacidade de generalização. Em Mingers (1989b), observou-se que, para diversos conjuntos de dados, a poda aumentou a precisão de classificação de 20% a 25%.

2.6.3.1 Reduced-Error Pruning Tree (REPT)

Como já foi apresentado na seção 2.6.3, a complexidade de uma árvore de decisão exerce forte influência no seu desempenho de classificação, e a forma mais comum de simplificação (visando aumentar a capacidade de generalização) é a poda, ou seja, remoção de subárvores que tornam sua estrutura desnecessariamente complexa (vide exemplo na figura 16).

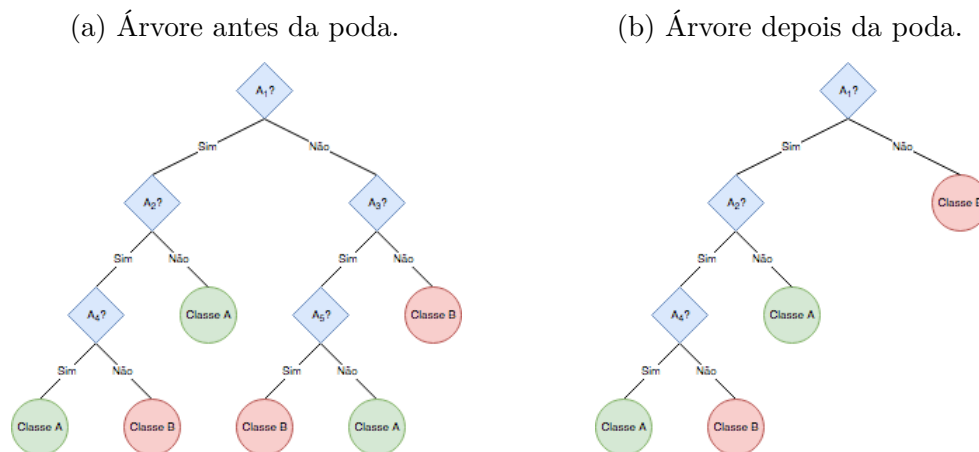


Figura 16: Exemplo de poda de árvore de decisão.

Diversos métodos de poda foram propostos na literatura. Nos trabalhos de Esposito, Malerba e Semeraro (1993) e Mingers (1989a), há revisões e estudos comparativos de vários deles, incluindo *Error-Complexity Pruning*, *Minimum-Error Pruning*, *Reduced-Error Pruning* e *Pessimistic Error Pruning*.

Neste trabalho, foi adotado o método *Reduced-Error Pruning* (QUINLAN, 1987), devido à sua simplicidade conceitual, velocidade de execução e eficácia na simplificação das árvores. Este tipo de poda funciona da seguinte forma:

```

V := conjunto de validação
Para cada nó N na árvore, faça:
  P(V) := precisão de classificação do conjunto de validação antes da poda em N
  Podar N, transformando-o em uma folha da classe mais comum em N
  P(V)' := precisão de classificação do conjunto de validação após a poda em N
  Se P(V)' < P(V), então desfazer a remoção da subárvore
  Caso contrário, então a poda de N foi realizada com sucesso

```

2.6.3.2 Boosted Reduced-Error Pruning Tree (AB-M1)

Boosting (em português “impulsionamento”) é um método para melhorar o desempenho de qualquer algoritmo de aprendizado. Em teoria, ele pode ser usado para reduzir significativamente o erro de qualquer algoritmo de aprendizado “fraco”, ou seja, que gere consistentemente classificadores que sejam apenas um pouco melhores do que adivinhações aleatórias. (SCHAPIRE, 1990) (FREUND, 1995)

Em Freund e Schapire (1996), são propostas duas versões do algoritmo: AdaBoost.M1, para problemas de classificação binária, e AdaBoost.M2, para problemas com mais de duas classes. Como este trabalho aborda justamente uma classificação binária, ou seja, atribuir um rótulo *ham* ou *spam* a um *e-mail*, o foco das explicações a seguir será na primeira versão do algoritmo.

O método, que está ilustrado na figura 17, funciona executando repetidamente um algoritmo de aprendizado “fraco” em diferentes distribuições sobre os dados de treinamento e, em seguida, combinando os classificadores h_1, \dots, h_n produzidos pelo algoritmo com seus respectivos pesos $\alpha_1, \dots, \alpha_n$ formando assim um único classificador composto H .

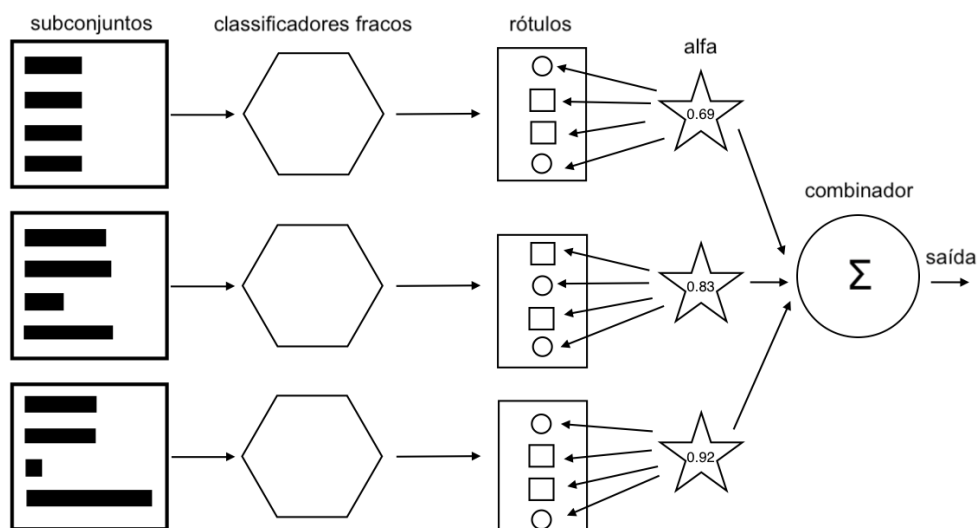


Figura 17: Ilustração do processo de funcionamento do AdaBoost.

Inicialmente, treina-se o primeiro classificador (h_1) atribuindo probabilidades iguais aos exemplos de um subconjunto aleatório de treinamento. Em seguida, é computado o peso α_1 deste classificador usando $t = 1$ na equação (2.26), na qual ϵ_t representa sua taxa de erro sobre o subconjunto.

$$\alpha_t = \frac{1}{2} \ln \left(\frac{1 - \epsilon_t}{\epsilon_t} \right) \quad (2.26)$$

Analisando a equação (2.26), notam-se os seguintes comportamentos de α para alguns valores de ϵ (MCCORMICK, 2013):

- α cresce positiva e exponencialmente à medida que ϵ se aproxima de 0.
- α é zero se ϵ for 0,5, pois um classificador com precisão de 50% não é melhor do que adivinhação aleatória e, portanto, deve ser ignorado.
- α cresce negativa e exponencialmente quando ϵ se aproxima de 1.

Após calcular α_1 , os pesos dos exemplos de treinamento (representados pelo vetor \vec{D}) são atualizados usando a equação (2.27), na qual i é o índice de cada amostra.

$$D_{t+1}(i) = \frac{D_t(i)e^{(-\alpha_t y_i h_t(x_i))}}{\sum_i D_t(i)} \quad (2.27)$$

Cada valor D_i representa a probabilidade de o i -ésimo exemplo ser incluído no conjunto de treinamento. Entretanto, para que \vec{D} represente uma distribuição, é necessário que $\sum_i D_i = 1$. Isto é feito via normalização, ou seja, dividindo cada D_i pelo somatório de todos os pesos, como mostra o denominador da equação (2.27).

Depois de treinar todos os T classificadores “fracos”, a saída do classificador “forte” final é dada pela equação (2.28):

$$H(x) = \text{sgn} \left(\sum_{t=1}^T \alpha_t h_t(x) \right) \quad (2.28)$$

Neste trabalho, foi adotado como classificador “fraco” o *Reduced-Error Pruning Tree* (REPT) já explicado na seção 2.6.3.1, o que tornou possível avaliar se o AdaBoost.M1 foi capaz de prover uma melhora na precisão de classificação.

2.6.4 Máquina de Vetores de Suporte

Originalmente propostas por Vapnik e Lerner (1963), posteriormente aprimoradas por Boser, Guyon e Vapnik (1992) e sendo a formulação atualmente tomada como padrão a estabelecida por Cortes e Vapnik (1995), as máquinas de vetores de suporte consistem em conjuntos de modelos supervisionados de aprendizado e algoritmos associados que analisam dados usados para classificação e regressão.

No contexto da classificação, dado um conjunto de exemplos de treinamento rotulados, um algoritmo de treinamento de SVM constrói um modelo capaz de atribuir classes a novos exemplos, com a propriedade de separar categorias distintas por uma lacuna o mais larga possível (vide figura 18), no intuito de minimizar o erro de generalização do classificador.

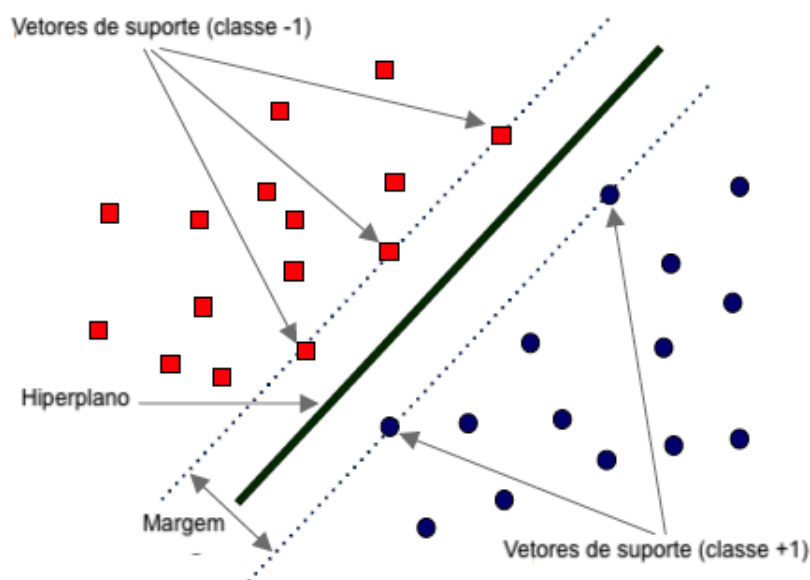


Figura 18: Hiperplano, vetores de suporte e margem de separação.

Na figura 18 é possível notar que há amostras de classes distintas em cada extremo da margem de separação. Cada amostra (\vec{x}, y) é uma tupla composta por um ponto n -dimensional \vec{x} e uma saída associada y . Justamente por “apoiarem” a margem de separação, estes dados situados nas fronteiras são chamados de vetores de suporte, o que dá nome também à técnica SVM.

Além de classificação linear, as SVMs podem realizar classificação não-linear de forma eficiente usando o chamado *kernel trick*. Esta técnica mapeia o espaço original de entradas (em que as classes não são linearmente separáveis) em espaços de características de igual ou mais alta dimensionalidade (vide figura 19), nos quais a separabilidade linear é possível. Para efetuar tal transformação, é usada uma função de *kernel*.

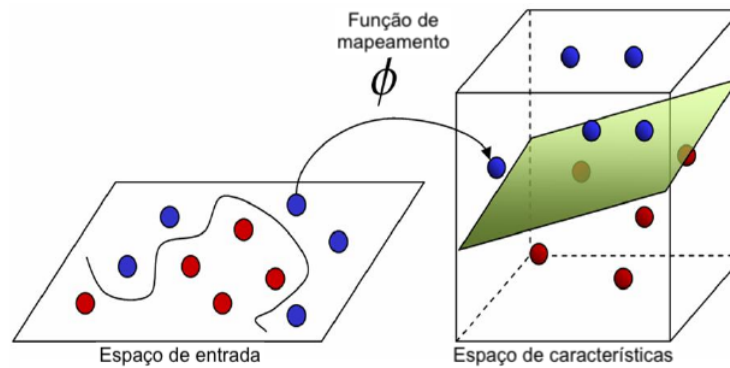
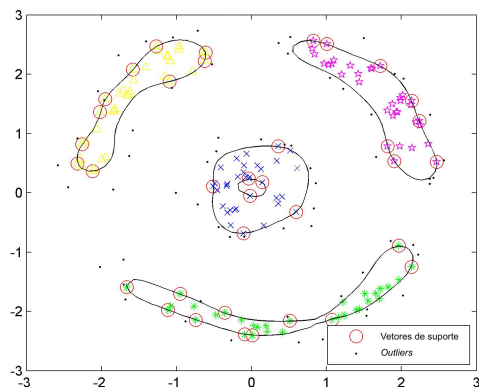


Figura 19: Aumento de dimensionalidade por meio do *kernel trick*.

As SVMs também conseguem encontrar agrupamentos de dados quando as amostras de treinamento não são rotuladas. Para isto, podem ser adotadas técnicas baseadas em grafos completos (BEN-HUR et al., 2001) ou de proximidade (YANG; ESTIVILL-CASTRO; CHALUP, 2002) (vide figuras 20a e 20b, respectivamente).

(a) Agrupamento com grafo completo.



(b) Agrupamento com grafo de proximidade.

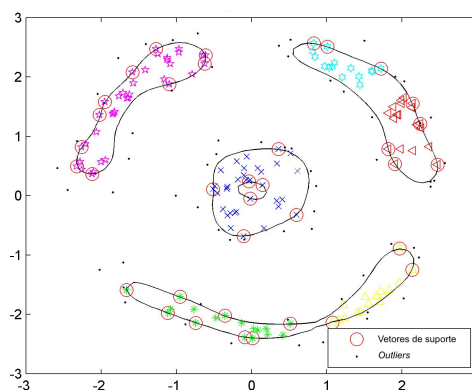


Figura 20: Agrupamentos de dados usando máquinas de vetores de suporte.

2.6.4.1 Linear Support Vector Machine (L-SVM)

Como dito na seção 2.6.4, o conjunto de treinamento utilizado pelas SVMs em cenários de classificação são tuplas (\vec{x}, y) compostas por um ponto n -dimensional \vec{x} e uma saída associada y . Neste sentido, a construção do modelo consiste em encontrar um hiperplano capaz de separar as amostras de treinamento de forma que haja a maior distância possível entre as categorias do problema.

Quando os dados de treinamento são linearmente separáveis, é possível encontrar um hiperplano de “margem rígida” que atenda ao requisito de separabilidade máxima desde que não haja erros e/ou *outliers* nos dados. Entretanto, caso os dados não sejam linearmente separáveis e/ou deseja-se permitir certa tolerância de erros de classificação (visando evitar *overfitting*), deve-se então encontrar um hiperplano de “margem suave”, que não se ajustará totalmente e apresentará capacidade de generalização potencialmente superior.

(a) SVM de margem rígida.

(b) SVM de margem suave.

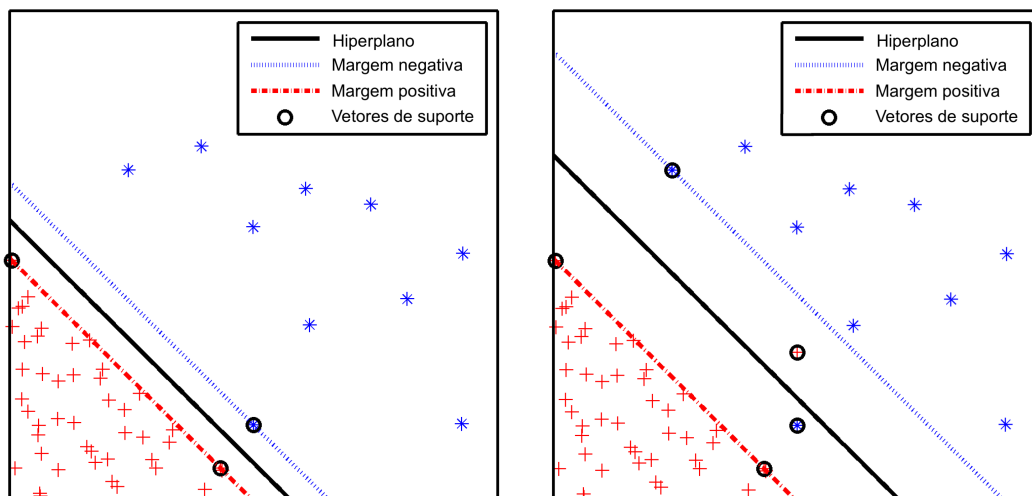


Figura 21: Representação gráfica das margens de separação rígida e suave.

Nas figuras 21a e 21b estão ilustradas SVMs de margem rígida e suave, respectivamente. Na primeira não há pontos com classificação incorreta; entretanto, a margem de separabilidade é estreita. Já na segunda, há uma amostra de cada classe classificada incorretamente; porém, a margem de separabilidade é larga.

A “margem” é a região delimitada por um par de hiperplanos paralelos que separam as duas classes de dados com a maior distância possível, e o “hiperplano de margem máxima” é aquele que fica no meio do caminho entre eles.

Para encontrar tal hiperplano em casos linearmente separáveis, deve-se resolver o seguinte problema de otimização: “dado o conjunto de treinamento $\{(\vec{x}_i, d_i)\}_{i=1}^N$, encontre os valores ótimos do vetor peso \vec{w} e *bias* b de modo que satisfaçam as restrições $d_i (\vec{w}^T \vec{x}_i + b \geq 1)$ para $i = 1, 2, \dots, N$ e o vetor peso \vec{w} minimize a função de custo $\Phi(\vec{w}) = \frac{1}{2} \vec{w}^T \vec{w}$ ” (HAYKIN, 2001).

Em casos não linearmente separáveis, são introduzidas as chamadas variáveis soltas ξ_i , que medem o desvio de um ponto de dado da condição ideal de separabilidade de padrões. Desta forma, obtém-se o seguinte problema de otimização: “dado o conjunto de treinamento $\{(\vec{x}_i, d_i)\}_{i=1}^N$, encontre os valores ótimos do vetor peso \vec{w} e *bias* b de modo que satisfaçam às restrições $d_i (\vec{w}^T \vec{x}_i + b \geq 1 - \xi_i)$ para $i = 1, 2, \dots, N$ e $\xi_i \geq 0$ para todo i e de modo que o vetor peso \vec{w} e as variáveis soltas ξ_i minimizem a função de custo $\Phi(\vec{w}, \xi) = \frac{1}{2} \vec{w}^T \vec{w} + C \sum_{i=1}^N \xi_i$ em que C é um parâmetro positivo especificado pelo usuário” (HAYKIN, 2001).

O parâmetro C indica o quanto deseja-se evitar erros de classificação. Portanto, valores grandes implicam na escolha de um hiperplano de margem menor, enquanto valores pequenos resultarão em hiperplanos com margem maior. A figura 22 ilustra a influência deste parâmetro na fronteira de decisão de uma SVM.

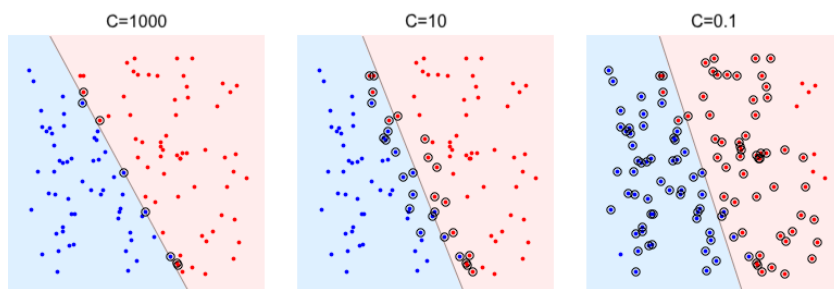


Figura 22: Ilustração da influência do parâmetro C na margem de uma SVM.

Tanto no caso linear quanto no não-linear, os \vec{w} e b que solucionam os problemas de otimização determinam o classificador $\vec{x} \mapsto \text{sgn}(\vec{w} \cdot \vec{x} - b)$.

2.6.4.2 Non-Linear Support Vector Machine (NL-SVM)

Na seção 2.6.4.1, foram apontados os problemas de otimização cujas soluções resultam em SVMs de margem rígida e suave. Como dito, o primeiro tipo destina-se a conjuntos de dados linearmente separáveis (sobre os quais não são tolerados erros de classificação), enquanto o segundo é normalmente utilizado em conjuntos de dados não-linearmente separáveis (sobre os quais podem haver erro de classificação de modo a preservar a capacidade de generalização do classificador).

Em ambos os casos, é importante ressaltar que o classificador produzido é de natureza linear, por trabalhar com hiperplanos construídos no espaço original de características e por utilizar como função de *kernel* o produto interno. A ideia de criar classificadores não-lineares foi proposta em Boser, Guyon e Vapnik (1992), cuja abordagem consiste na aplicação da *kernel trick* (AIZERMAN; BRAVERMAN; ROZONOÉR, 1964) aos hiperplanos de margem máxima.

Uma SVM não-linear funciona de forma similar a uma linear; entretanto, os produtos internos são substituídos por uma função de *kernel* não-linear (vide alguns exemplos na tabela 4 (HAYKIN, 2001)). Isto permite ao algoritmo ajustar o hiperplano de margem máxima em um espaço de características transformado de igual ou maior dimensionalidade e linearmente separável. Somente a máquina de vetores de suporte de *kernel* polinomial (P-SVM) foi usada nos experimentos.

Tabela 4: Exemplos de funções de *kernel*.

Tipo	Função
Linear	$K(\vec{x}, \vec{x}_i) = \vec{x}^T \vec{x}_i$
Polinomial	$K(\vec{x}, \vec{x}_i) = (\vec{x}^T \vec{x}_i)^d$
Gaussiana	$K(\vec{x}, \vec{x}_i) = \exp\left(-\frac{1}{2\sigma^2} \ \vec{x} - \vec{x}_i\ ^2\right)$
Sigmoide	$K(\vec{x}, \vec{x}_i) = \tanh(\eta \vec{x} \vec{x}_i + \theta)$

Outras abordagens que utilizam a *kernel trick* incluem a categorização de textos com base em subsequências de caracteres (LODHI et al., 2002) e o uso de *kernels* compostos na classificação de imagens espectrais (CAMPS-VALLS et al., 2006).

3 Revisão Bibliográfica

O estudo e desenvolvimento de sistemas anti-*spam* são assuntos recorrentes em trabalhos da área de Ciência da Computação. Pesquisadores vêm abordando o tema das mais diversas formas, propondo soluções que utilizam desde técnicas estatísticas e redes neurais artificiais à análise da reputação dos remetentes e detecção de padrões e estilos de escrita frequentes em mensagens indesejadas.

O objetivo do presente trabalho é associar técnicas de seleção de características, redução de dimensionalidade e modelos de aprendizado de máquina de modo a construir um classificador de *e-mails* automático, robusto e adaptável. Portanto, neste capítulo são apresentados e discutidos (em ordem cronológica de data de publicação) outros trabalhos e pesquisas envolvendo classificação de *e-mails*.

Drucker et al. (1999) comparam as SVMs a três outras técnicas (Boosting, Ripper e Rocchio) aplicadas à classificação de *e-mails*. Os testes foram conduzidos em dois conjuntos de dados oriundos da companhia americana de telecomunicações AT&T, submetidos a múltiplos métodos de seleção de características (binário ou baseados em frequências dos termos) e a uma filtragem de palavras vazias. Além de resultados de precisão em torno de 98%, os autores destacaram características importantes da SVM: a melhora de velocidade de treinamento ao usar características binárias¹, a sensibilidade ao parâmetro C^2 e a capacidade de lidar com grandes

¹Neste caso, a maioria das características é nula, resultando em vetores esparsos cujo produto interno pode ser otimizado com rotinas não-multiplicativas.

²Parâmetro de regularização que controla o balanço entre alcançar um erro baixo nos dados de treinamento e minimizar a norma dos pesos.

quantidades de características (em vez de subconjuntos).

Utilizando o classificador probabilístico *Naïve Bayes* aliado a técnicas linguísticas de lematização³ e remoção de palavras vazias⁴, Androutsopoulos et al. (2000a) obtiveram taxas de precisão de classificação superiores a 95% no *corpus* Ling Spam (ANDROUTSOPOULOS et al., 2000b), com um custo computacional relativamente baixo. No trabalho, o autor observa que, apesar de a adoção das técnicas linguísticas ter melhorado sensivelmente a precisão na classificação, elas introduzem novas possibilidades na análise das mensagens. Por fim, foram realizados experimentos com diferentes custos atribuídos a falsos-positivos e falsos-negativos. Este aspecto experimental levou à conclusão de que este modelo é inviável nas situações em que *e-mails* são descartados quando identificados como *spam*.

Meyer e Whateley (2004) avaliaram o desempenho de um classificador *Naïve Bayes* usado em conjunto com esquemas de combinação (probabilidades χ^2 a fim de pontuar mensagens como *ham*, *spam* ou incertas), *n*-gramas⁵ (uni ou bi-gramas) e *tokenização* (símbolos de cabeçalhos, corpo ou “sintéticos” — características indiretas dos *e-mails*, como ausência de remetente, assunto etc.). O treinamento foi feito sobre cinco *corpora* — quatro SpamBayes e o Spam Assassin (The Apache Software Foundation, 2005), utilizando um regime de treinamento configurável (que indica ao algoritmo “quando” a aprendizagem deve ser feita: com todas as mensagens, somente com mensagens dentro de uma faixa de pontuação etc.) e validação cruzada. A quantidade de mensagens com necessidade de classificação manual após o treinamento foi de apenas pouco mais de 1%, comprovando a eficácia da combinação de técnicas escolhida.

Berger e Merkl (2005) realizaram um estudo comparativo das técnicas de aprendizado supervisionado SVM e não-supervisionado SOM (do inglês *Self-Organizing Map*) aplicadas à categorização de *e-mail*. Nos experimentos, as características

³Representação de palavras em sua forma básica, desconsiderando mudanças gramaticais como o tempo e a pluralidade.

⁴Palavras que são filtradas antes ou depois do processamento de dados de linguagem natural (texto); geralmente referem-se às palavras mais comuns em um idioma.

⁵Subseqüência de *n* elementos de uma dada seqüência.

dos *e-mails* (oriundos de um conjunto pessoal de mensagens) foram representadas tanto como palavras quanto como *n*-gramas, sob a justificativa destes apresentarem maior robustez em relação a erros de ortografia e independerem de contexto. Nos resultados experimentais, ambas as técnicas apresentaram precisão superior a 90%, tendo a SOM resultados levemente inferiores à SVM, superando as expectativas dos pesquisadores devido à natureza não-supervisionada de seu treinamento.

Carpinteiro et al. (2006) propuseram o pré-processamento de *e-mails* visando simplificar as mensagens e melhorar o levantamento de suas características. Analisando diversas partes dos *e-mails* (cabeçalho, corpo — tanto texto quanto HTML⁶ — e tipos dos anexos) e aplicando diferentes métodos de seleção de características, foi treinada uma rede neural do tipo MLP (do inglês *Multilayer Perceptron*). Os experimentos mostraram precisão de classificação próxima a 100% na base Spam Assassin (The Apache Software Foundation, 2005). Para trabalhos futuros, são sugeridos outros modelos neurais, formas de pré-processamento de mensagens e métodos de seleção de características.

Sirisanyalak e Sornil (2007) propuseram uma solução para classificação de *e-mails* baseada em um sistema imunológico artificial. Simulando a resposta imunológica biológica, o sistema deve ser capaz tanto de gerar e clonar anticorpos⁷ eficazes quanto de detectar e eliminar anticorpos ineficazes (ou seja, aqueles que atacam células do próprio organismo). Submetendo o *corpus* Spam Assassin (The Apache Software Foundation, 2005) a uma análise sintática de MIME⁸, remoção de palavras vazias e lematização, um conjunto de detectores⁹ é gerado, aprimorado e sintetizado por meio de algoritmos genéticos. Finalmente, é realizada a extração de quatro características específicas¹⁰ e a base é submetida a avaliações previamente

⁶Linguagem de marcação utilizada na construção de páginas na *Web*.

⁷Unidades dedicadas ao reconhecimento de um ou mais antígenos (neste caso, características de *spams*) por meio de propriedades particulares a estes.

⁸Do inglês *Multipurpose Internet Mail Extensions*; padrão da internet para o formato das mensagens de correio eletrônico.

⁹Neste caso, consistem em expressões regulares capazes de detectar palavras e suas variações causadas por técnicas de ofuscação, por exemplo.

¹⁰Número de termos na mensagem, de vezes que a mensagem casa com o conjunto de detectores,

submetidas a métodos de escolha de parâmetros e realizadas por uma função de regressão logística¹¹, obtendo resultados de precisão superiores a 98%.

Zhang et al. (2009) propuseram um sistema de classificação de *e-mails* baseado em reputação, no qual um servidor armazena grupos de endereços IP (em uma estrutura de dados do tipo árvore) e suas respectivas pontuações. Compostas a partir do histórico de envio de *e-mails* destes grupos de usuários e realimentadas tanto por outros sistemas anti-*spam* quanto pelos destinatários, as notas indicam o quanto um remetente é confiável. Utilizando um *corpus* composto por quase 3 milhões de *e-mails* oriundos de um servidor universitário, os experimentos mostraram que o sistema proposto atingiu mais de 95% de acurácia, precisão e revocação.

Pérez-Díaz et al. (2012) propuseram o uso da técnica de conjuntos aproximados (PAWLAK, 1991) na categorização de *e-mails* aliada a heurísticas¹² para endereçar o problema de mensagens com classificação duvidosa. Após submeter a base Spam Assassin (The Apache Software Foundation, 2005) a uma remoção de palavras vazias e utilizando representações binárias ou em termos de frequência, os autores executam a técnica de conjuntos aproximados em regime de validação cruzada e comparam os resultados de classificação aos das técnicas AdaBoost, Flexible Bayes, Naïve Bayes e SVM, usando diferentes quantidades de características para cada. O método obteve taxas de escore F_1 de 98%, superando os demais modelos adotados na comparação. Por fim, os autores destacam a importância de regenerar periodicamente o conjunto de regras derivado pela técnica devido ao *concept drift*¹³ e alertaram para o grande tempo de treinamento, sugerindo a adoção de métodos de redução do espaço de características.

de detectores que casam com a mensagem e um fator de tendência de *spam*.

¹¹Técnica estatística que tem como objetivo produzir, a partir de um conjunto de observações, um modelo que permita a predição de valores tomados por uma variável categórica, frequentemente binária, a partir de uma série de variáveis explicativas contínuas e/ou binárias.

¹²MFD (do inglês *Most Frequent Decision*, LNO (do inglês *Largest Number of Objects*) e LTS (do inglês *Largest Total Strength*).

¹³Noção de que propriedades estatísticas da variável alvo, das quais o modelo está tentando prever, mudam ao longo do tempo de maneira imprevisível, tornando as previsões precisas à medida que o tempo passa.

No trabalho de Yevseyeva et al. (2013), três algoritmos evolutivos são aplicados à otimização de filtros anti-*spam* (sendo um deles de único objetivo — Grindstone4SPAM (MÉNDEZ et al., 2012) — e os outros dois, multiobjetivos — NSGA-II (do inglês *Nondominated Sorting Genetic Algorithm II*) (DEB et al., 2002) e SPEA2 (do inglês *Strength Pareto Evolutionary Algorithm 2*) (ZITZLER; LAUMANN; THIELE, 2001)) e comparados em termos de métricas específicas aos domínios de algoritmos estocásticos — EAF (do inglês *empirical attainment function* (LÓPEZ-IBÁÑEZ; PAQUETE; STÜTZLE, 2010)) — e de pesquisa de sistemas anti-*spam* — precisão, revocação, escore F_1 etc. Foram elaboradas duas estratégias de avaliação para abordar estes domínios, respectivamente: traçado de gráficos dos *fronts* de Pareto¹⁴ (a fim de facilitar a visualização da eficiência de cada método evolutivo) e o treinamento e avaliação da influência da otimização em um regime de validação cruzada. Por fim, os resultados apontaram que o NSGA-II obteve o melhor desempenho, e os autores destacam que a variedade de soluções obtidas pelos algoritmos multiobjetivos não foi totalmente explorada no trabalho justamente pelas restrições impostas pelo arcabouço necessário à comparação de algoritmos de tipos diferentes.

Barigou, Beldjilali e Atmani (2014) propuseram uma melhoria ao algoritmo KNN (do inglês *K-Nearest Neighbors*), introduzindo autômatos celulares para reduzir a quantidade de documentos selecionados ao classificar uma amostra. Tal alteração conferiu uma redução de uso de memória¹⁵ e aumento do desempenho¹⁶ do classificador que, quando exercitado com a base Ling Spam (ANDROUTSOPOULOS et al., 2000b), apresentou precisão superior a 98%. Para trabalhos futuros, os autores sugerem avaliar a mesma técnica com a base Spam Assassin (The Apache Software Foundation, 2005) e estabelecer mais critérios de avaliação.

Youn (2014) propôs um filtro de *e-mails* de dois níveis baseado em ontologia¹⁷. A partir de um determinado conjunto de dados, as características são selecionadas

¹⁴Conjunto de soluções não-nomeadas, sendo escolhido como ótimo, se nenhum objetivo pode ser melhorado sem sacrificar pelo menos um outro objetivo.

¹⁵Devido à estrutura celular em que os dados do conjunto treinamento são organizados.

¹⁶Devido à redução do nível de ruído dos subconjuntos e à resistência a desbalanceamento.

¹⁷Representações formais de um conjunto de conceitos dentro de um domínio e suas relações.

pelo método TF-IDF (do inglês *Term Frequency–Inverse Document Frequency*)¹⁸ e uma árvore de decisão C4.5 (QUINLAN, 1993) é induzida para posteriormente ser convertida para o formato RDF (do inglês *Resource Description Framework*) (CYGANIAK; WOOD; LANTHALER, 2014) e usada para produzir uma ontologia global (primeiro nível). As mensagens do conjunto de testes são objeto de consulta para esta ontologia, que fornece uma classificação baseada nos atributos das mensagens e em seu próprio modelo, que é periodicamente adaptado e pode ser usado para conceber uma ontologia personalizada (segundo nível). Submetido a um conjunto de dados particular, o método apresentou precisões em torno de 91% e 95% usando os filtros global e personalizado, respectivamente.

Visando melhorar o desempenho do algoritmo de seleção negativa, Idris et al. (2015) propuseram um modelo híbrido denominado NSA-PSO (do inglês *Negative Selection Algorithm — Particle Swarm Optimization*), substituindo o gerador de detectores aleatórios pelo algoritmo de otimização por enxame de partículas (J; C, 1995), usando como função de *fitness* a LOF (do inglês *Local Outlier Factor*) (BREUNIG et al., 2000). Os autores também propuseram um modelo de arquitetura modular para o sistema anti-*spam*, facilitando o caminho para uma implantação real. Por fim, experimentando diferentes configurações de limiar e quantidade de detectores, o modelo foi exercitado com o *corpus* Spambase (HOPKINS et al., 1999) e apresentou taxa de precisão de 77% e escore F_1 de 83% (contra respectivos 69% e 36% do algoritmo original), comprovando uma melhora de desempenho.

Kaya e Ertuğrul (2016) propuseram um novo método de seleção de características baseado na “probabilidade do uso de caracteres de ordens similares com respeito aos seus valores UTF-8 (do inglês *8-bit Unicode Transformation Format*)¹⁹”, determinados pelo operador de baixo nível *shifted 1D-LBP* (do inglês *One-*

¹⁸Medida estatística que tem o intuito de indicar a importância de uma palavra de um documento em relação a uma coleção de documentos.

¹⁹Tipo de codificação binária de comprimento variável capaz de representar qualquer caractere universal padrão do *Unicode* e *ASCII* (do inglês *American Standard Code for Information Interchange*).

Dimension Linear Binary Pattern)²⁰. Nos experimentos, o método é aplicado aos conjuntos Ling Spam (ANDROUTSOPOULOS et al., 2000b), Spam Assassin (The Apache Software Foundation, 2005) e TREC (CORMACK; LYNAM, 2006) e, em seguida, alguns modelos de aprendizado de máquina²¹ são treinados com validação cruzada e avaliados em termos de diversas métricas²². Os resultados de acurácia de classificação se mostraram promissores, atingindo aproximadamente 92%, 93% e 95% para os conjuntos de dados mencionados anteriormente, respectivamente.

Shams e Mercer (2016) propuseram a extração e o uso de atributos de estilometria²³ para treinar classificadores de *e-mails*. Como exemplos destes atributos, pode-se citar a quantidade de erros ortográficos e gramaticais, indicadores de facilidade de leitura (índice Gunning fog), SMOG (do inglês *Simple Measure of Gobbledygook*), FRES (do inglês *Flesch Reading Ease Score*), Forcast e legibilidade de Flesch-Kincaid), quantidades de palavras simples (com até duas sílabas) e complexas (com três ou mais sílabas) e tamanho médio de *e-mail* e de palavras. Aliados a atributos clássicos como IG (do inglês *Information Gain*) e TF-IDF, o método foi testado com os classificadores NB, RF e SVM e os *ensembles*²⁴ Bagging e Adaboost.M1 sobre os *corpora* CSDMC2010 (TAO, 2010), Spam Assassin (The Apache Software Foundation, 2005), Ling Spam (ANDROUTSOPOULOS et al., 2000b) e Enron-Spam (COHEN, 2008). Os métodos Bagging e Adaboost.M1 obtiveram os melhores resultados, e concluiu-se que a técnica é relevante na detecção de *spam* em *corpora* personalizados (ou seja, naqueles em que a coleta de *e-mails* não é aleatória), mas limitada em *corpora* não-personalizados, devido à ausência de padrões únicos de escrita neste tipo de conjunto de dados.

²⁰Para uma dada mensagem, o operador é definido como um conjunto de comparações binárias entre o valor central e os valores vizinhos, resultando em um valor inteiro que representa um padrão de informação estrutural.

²¹FLDA (do inglês *Fisher Linear Discriminant Analysis*), NB, BN (do inglês *BayesNet*), FT (do inglês *Functional Tree*), RT (do inglês *Random Tree*) e RF (do inglês *Random Forest*).

²²Precisão, revocação, acurácia, MCC (do inglês *Matthews Correlation Coefficient*), escore F₁ e ROC (do inglês *Receiver Operating Characteristic*).

²³Aplicação do estudo do estilo linguístico, geralmente à linguagem escrita, frequentemente usado para atribuir autoria a documentos anônimos ou contestados.

²⁴Uso de vários algoritmos de aprendizagem para obter melhor desempenho preditivo do que poderia ser obtido de qualquer um dos algoritmos de aprendizado constituintes sozinho.

Yang et al. (2016) propuseram um filtro anti-*spam* baseado no algoritmo IB (do inglês *Information Bottleneck*) destinado a conjuntos de treinamento pequenos, sob a justificativa de que a frequente mudança de temas dos *e-mails* reduz a disponibilidade imediata de amostras com contextos atuais. Basicamente, o algoritmo extrai amostras altamente significativas do conjunto, a fim de construir centroides usados para classificar mensagens em *ham* ou *spam* por meio de uma função de similaridade (neste caso, a divergência de Jensen–Shannon (LIN, 1991)). Submetendo os *corpora* Ling Spam (ANDROUTSOPOULOS et al., 2000b), Spambase (HOPKINS et al., 1999), PU3 (ANDROUTSOPOULOS; PALIOURAS; MICHELAKIS, 2003) e TREC (CORMACK; LYNAM, 2007) a uma extração de características baseadas em IG, o classificador apresentou resultados (de precisão, revocação e escore F_1) comparáveis aos competidores (SVM, NB e AdaBoost) quando trabalhando com conjuntos de treinamento grandes. Com a redução da quantidade de amostras, o algoritmo teve menor deterioração de desempenho, confirmando sua capacidade de predição mesmo com conjuntos pequenos.

Kumaresan, Saravanakumar e Balamurugan (2017) propuseram um *framework* para classificação de *e-mails* baseado em SVM. Combinando os *corpora* Ling Spam (ANDROUTSOPOULOS et al., 2000b) e Spam Archive (DREDZE; GEVARYAHU; ELIAS-BACHRACH, 2007a) (DREDZE; GEVARYAHU; ELIAS-BACHRACH, 2007b), o conjunto de dados utilizado é composto por mensagens com texto e imagens, cujas características textuais são selecionadas pelo método TF (do inglês *Term Frequency*) e as visuais pelos métodos correlograma²⁵ e *wavelet moment*²⁶. Após a extração, o espaço de características é reduzido por uma versão modificada do algoritmo de busca *Cuckoo search* com voos de Lévy²⁷ (YANG; DEB, 2010) e os vetores são submetidos à classificação realizada por uma SVM de *kernel* híbrido²⁸. Após os experimentos, o método obteve taxa precisão de classificação superior a

²⁵Gráfico utilizado em séries temporais para traçar autocorrelações, ou seja, a correlações entre uma série e ela mesma defasada.

²⁶Técnica usada para medir a regularidade local de um sinal.

²⁷Sucessão de passos aleatórios cujos comprimentos seguem uma distribuição de probabilidade de cauda pesada (exemplos: distribuições de Pareto, Lévy, Cauchy, Burr e *t* de Student).

²⁸Utiliza uma combinação de dois ou mais *kernels*, como linear, polinomial e quadrático.

97%, superando os 94% obtidos pelos demais modelos adotados na comparação e evidenciando o aumento de desempenho proporcionado pela técnica proposta.

Os trabalhos revisados nesta seção propuseram diversas abordagens e modelos para categorização de texto e detecção de *spam*. Os modelos foram avaliados sobre diversos *corpora* e frequentemente obtiveram acerto de classificação superior a 90%. Entretanto, a redução de dimensionalidade do espaço de características das mensagens não foi explorada (no máximo, fora citada por (PÉREZ-DÍAZ et al., 2012) como sugestão para continuação de sua pesquisa) e, portanto, o modelo proposto difere-se dos demais neste aspecto.

4 Metodologia

4.1 Modelo proposto

A fim de avaliar o desempenho dos diferentes algoritmos de aprendizado aliados às técnicas de seleção de característica, redução de atributos e balanceamento de classes explicadas no capítulo 2, foi elaborado um modelo ilustrado no diagrama de blocos da figura 23, cujos elementos são explicados individualmente na sequência.

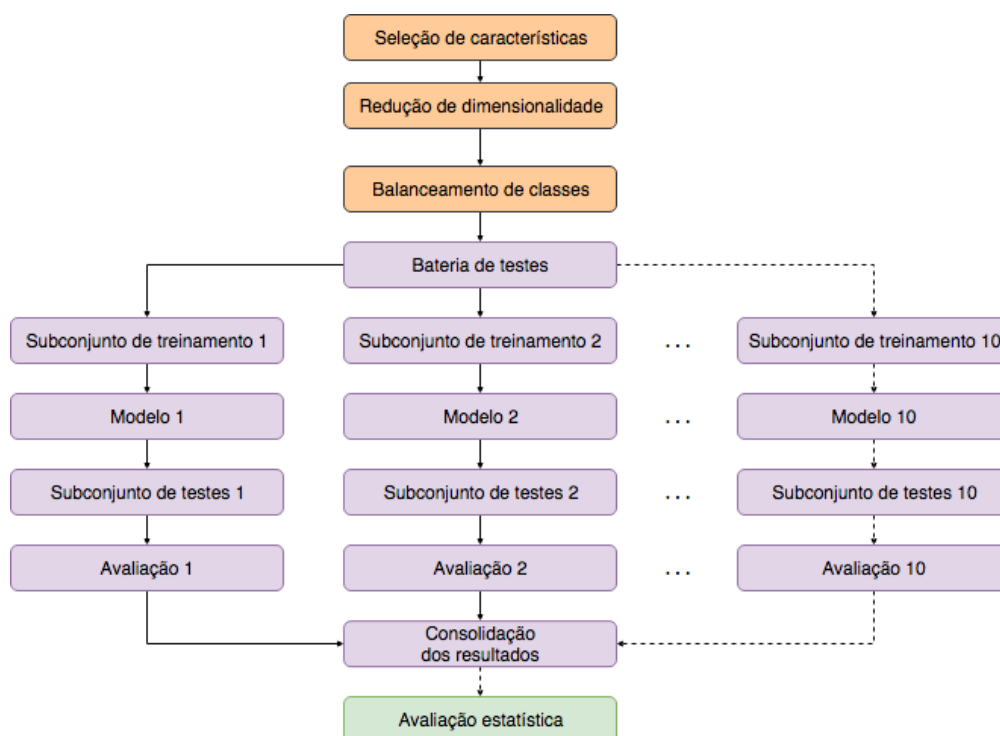


Figura 23: Diagrama de blocos da metodologia proposta.

- **Seleção de características:** Etapa que analisa os *e-mails* dos conjuntos de dados em sua forma textual e extrai suas características de acordo com os métodos descritos em 2.3, gerando uma representação numérica de cada mensagem, que será utilizada pelos modelos de aprendizado de máquina.
- **Redução de dimensionalidade:** Usando o método explicado na seção 2.4, é efetuada uma redução de dimensionalidade no conjunto de dados gerado na etapa anterior, preservando os atributos mais significativos de cada amostra e descartando os menos relevantes perante a classe de cada instância.
- **Balanceamento de classes:** Os conjuntos de dados já reduzidos são balanceados, a fim de equilibrar a contribuição de cada classe no treinamento do modelo. Esta técnica é descrita mais detalhadamente na seção 2.5.
- **Bateria de testes:** Nesta etapa, o conjunto de dados é embaralhado e dividido ao meio¹, e as metades são destinadas às etapas de treinamento e classificação (aqui, os padrões zerados² mencionados na tabela 6 são incluídos no subconjunto). Em seguida, cada modelo é treinado e avaliado 10 vezes.
- **Subconjunto de treinamento:** Consiste na primeira metade do conjunto de dados reduzido e embaralhado.
- **Modelo:** É o artefato resultante do treinamento do modelo de aprendizado. Pode ser persistido no sistema de arquivos e recuperado para uso posterior.
- **Subconjunto de testes:** Consiste na segunda metade do conjunto de dados reduzido e embaralhado, somada aos padrões zerados.
- **Avaliação:** O modelo é avaliado por meio da classificação do subconjunto de testes, e o resultado de cada métrica (explicadas na seção 4.4) é calculado.
- **Consolidação:** As métricas calculadas na etapa anterior são consolidadas para que as médias e intervalos de confiança sejam calculados.
- **Avaliação estatística:** Os resultados dos 10 experimentos são consolidados para que a média e o intervalo de confiança sejam calculados e apresentados.

¹Conjuntos de treinamento maiores que os de teste representam cenários irrealistas, pois a quantidade de dados disponível para pesquisa e aquisição de conhecimento é consideravelmente inferior àquela processada diariamente por grandes servidores de *e-mail* em ambientes reais.

²Mensagens cuja representação vetorial só possui valores nulos (ou seja, iguais a zero).

4.2 Ferramentas utilizadas

O projeto WEKA (do inglês *Waikato Environment for Knowledge Analysis*) (HALL et al., 2009) visa fornecer uma coleção abrangente de algoritmos de aprendizado de máquina e ferramentas de pré-processamento de dados para pesquisadores e profissionais. Ele permite que os usuários experimentem rapidamente e comparem diferentes métodos de aprendizagem de máquina em novos conjuntos de dados. Sua arquitetura modular e extensível permite que processos sofisticados de mineração de dados sejam construídos a partir da ampla coleção de algoritmos e ferramentas de aprendizagem disponíveis. A extensão do *kit* de ferramentas é fácil graças a uma API (do inglês *Application Programming Interface*) simples, mecanismos de *plugin* e instalações que automatizam a integração de novos algoritmos de aprendizagem com as interfaces gráficas de usuário do WEKA (vide figura 24).



Figura 24: Tela principal do WEKA.

O projeto inclui algoritmos para regressão, classificação, agrupamento, mineração de regras de associação e seleção de atributo. A exploração preliminar de dados é bem atendida por telas de visualização de dados e muitas ferramentas de pré-processamento. Estas, quando combinados com avaliação estatística de esquemas de aprendizagem e visualização dos resultados de aprendizagem, suporta modelos de processo de mineração de dados.

O WEKA conta também com um gerenciador de pacotes (vide figura 25), com o qual é possível pesquisar e baixar diversas extensões feitas pela comunidade.

Dentre elas encontram-se algoritmos de pré-processamento, seleção de características, classificação, regressão etc. Também é possível instalar pacotes não-oficiais, que muitas vezes ainda estão sendo desenvolvidos e, por isso, ainda não foram disponibilizados oficialmente.

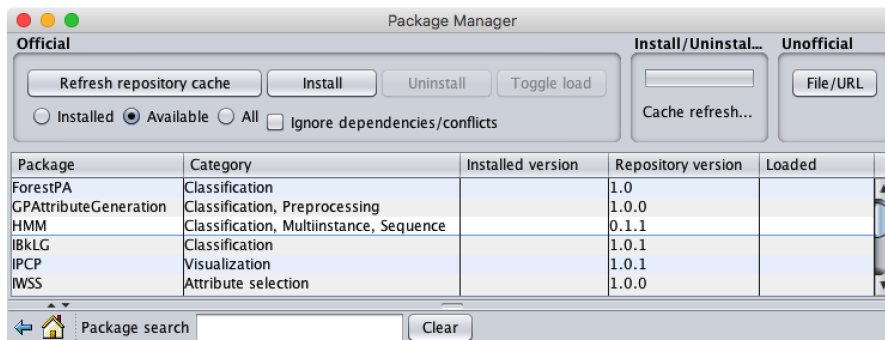


Figura 25: Tela do gerenciador de pacotes do WEKA.

Todos os oito métodos de classificação propostos neste trabalho estão disponíveis no WEKA. Nos experimentos com máquinas de vetores de suporte, as célebres bibliotecas LIBLINEAR (FAN et al., 2008) e LIBSVM (CHANG; LIN, 2011), totalmente suportadas pelo WEKA, poderiam ter sido usadas para treinar as máquinas de vetores de suporte linear e polinomial, respectivamente.

Na página da LIBSVM, porém, é possível encontrar uma versão alternativa da LIBLINEAR, chamada LIBLINEAR-poly2, cujo código “é mais eficiente do que a LIBSVM em ambas as fases de treinamento e teste para SVM com *kernels* polinomiais de grau 2”, como afirmam Chang et al. (2010) na documentação da biblioteca. Neste mesmo documento, há um exemplo mostrando que a diferença de precisão de classificação entre elas foi irrisória (da ordem de 0,001%) e que o tempo de execução foi cerca de sete vezes menor na segunda biblioteca.

Portanto, a fim de manter a comparação entre SVMs consistente em termos de plataforma/linguagem e, ao mesmo tempo, tirar proveito das vantagens da LIBLINEAR-poly2, decidiu-se por adotar as bibliotecas LIBLINEAR e LIBLINEAR-poly2 para treinar as SVMs linear e polinomial, respectivamente.

4.3 Conjuntos de dados utilizados

O treinamento dos algoritmos de aprendizado de máquina exige conjuntos de dados que sejam suficientemente representativos para o problema proposto. Os estudos avaliativos dos classificadores utilizaram três bases de dados públicas e uma privada, compostas por *e-mails* reais. A descrição dos *corpora* é dada a seguir.

4.3.1 Ling Spam

Compilado e disponibilizado em domínio público por Androutsopoulos et al. (2000a), o *corpus* Ling Spam (ANDROUTSOPOULOS et al., 2000b) é composto por 2893 mensagens de *e-mail* e apresenta as seguintes características:

- 2.412 *hams* de uma lista de discussão moderada sobre a profissão e ciência de linguística, obtidas por meio do *download* aleatório de resumos dos arquivos, separação de mensagens e remoção de texto adicionado pelo servidor da lista.
- 481 *spams* recebidos pelo primeiro autor. Anexos, *tags* HTML e mensagens duplicadas recebidas no mesmo dia não foram incluídos.
- É disponibilizado em quatro versões:
 - *bare*: lematização desativada e remoção de palavras vazias desativada.
 - *lemm*: lematização ativada e remoção de palavras vazias desativada.
 - *lemm_stop*: lematização ativada e remoção de palavras vazias ativada.
 - *stop*: lematização desativada e remoção de palavras vazias ativada.

Embora as mensagens de *spam* e *ham* sejam extraídas de fontes diferentes e algumas informações (ex: anexos e *tags* HTML) tenham sido removidas, este conjunto apresentou uma avaliação mais realista do que qualquer outro disponível quando do seu lançamento, sendo usado em diversos estudos, como Barigou, Beldjilali e Atmani (2014), Kaya e Ertuğrul (2016), Shams e Mercer (2016), Yang et al. (2016) e Kumaresan, Saravanakumar e Balamurugan (2017).

Neste trabalho, as quatro versões foram consolidadas em um único *corpus*.

4.3.2 Spam Assassin

O *corpus* Spam Assassin (The Apache Software Foundation, 2005) é composto por 6.047 mensagens de *e-mail* reais, com aproximadamente 31% de *spam*. As mensagens provêm de várias fontes, incluindo fóruns públicos e mensagens enviadas de e para os desenvolvedores do Spam Assassin. Todos os cabeçalhos foram reproduzidos em sua totalidade. Há ofuscação de endereços³ e os *hostnames* foram trocados por `spamassassin.taint.org` em algumas situações. Entretanto, na maioria dos casos, os cabeçalhos aparecem da mesma forma que foram recebidos. O *corpus* é composto por cinco partes, descritas na tabela 5.

Tabela 5: Subconjuntos do *corpus* Spam Assassin.

Subconjunto	<i>Hams</i>	<i>Spams</i>	Descrição
<i>spam</i>	-	500	<i>Spams</i> recebidos de fontes reais
<i>easy_ham</i>	2.500	-	Mensagens de fácil identificação, por serem completamente diferentes de spams
<i>hard_ham</i>	250	-	Mensagens de difícil identificação, por se parecerem em muitos aspectos com o spam típico
<i>easy_ham_2</i>	1.400	-	Adições recentes ao <i>corpus</i>
<i>spam_2</i>	-	1.397	

O *corpus* tem sido usado tanto para o desenvolvimento do *software* Apache Spam Assassin (FOUNDATION, 2015) e outros filtros de *spam* quanto como alvo de avaliação em publicações como Meyer e Whateley (2004), Carpinteiro et al. (2006), Sirisanyalak e Sornil (2007), Pérez-Díaz et al. (2012), Barigou, Beldjilali e Atmani (2014), Kaya e Ertuğrul (2016) e Shams e Mercer (2016).

Neste trabalho, as cinco partes deste *corpus* foram consolidadas, gerando um conjunto com mais mensagens e maior variabilidade.

³Técnica que visa esconder (ou dificultar) a identificação de *e-mail* por *bots*, mantendo-os ao mesmo tempo visíveis para humanos.

4.3.3 TREC

O TREC Spam Track (CORMACK; LYNAM, 2005) (CORMACK; LYNAM, 2006) (CORMACK; LYNAM, 2007) foi a maior e mais realista avaliação de laboratório até o momento de sua disponibilização. Nos três anos que funcionou, dez conjuntos de teste com um total de 721.461 mensagens foram utilizados para testar filtros enviados por 35 participantes. Cada participante apresentou até quatro filtros, que foram testados com os *e-mails* dos *corpora*. (CORMACK et al., 2008)

Dentre os diversos *corpora* disponíveis, três foram utilizados neste trabalho:

- *TREC 2005* (CORMACK; LYNAM, 2005): é composto pelo *e-mail* armazenado de 150 executivos da Enron⁴, coletados e lançados no domínio público como resultado da investigação federal dos EUA sobre seu colapso. Essas mensagens foram rotuladas e somadas a mensagens de *spam* oriundas de uma fonte pública, alteradas de forma a parecer terem sido entregues à Enron durante o mesmo período. Esta base estabeleceu que os *corpora* privados e públicos poderiam produzir resultados razoavelmente consistentes.
- *TREC 2006* (CORMACK; LYNAM, 2006): compreende dois *corpora* privados e dois públicos (*trec06p* - em inglês, extraído da Web e acrescido de *spam* público, e *trec06c* - em chinês, obtido a partir de uma lista de *e-mails* e acrescido de *spams* capturados no mesmo site por um *honeypot*⁵).
- *TREC 2007* (CORMACK; LYNAM, 2007): compreende um *corpus* privado e um público (*trec07p*), sendo o último composto por mensagens entregues a um servidor particular durante três meses.

Neste trabalho, os *corpora* públicos de 2005, 2006 (somente a versão em inglês) e 2007 foram consolidados, gerando um conjunto de origem mista e com uma quantidade expressiva de mensagens.

⁴Companhia de energia americana que decretou falência em 2001 e foi alvo de investigação federal após ter sido descoberto um escândalo fiscal da ordem de bilhões de dólares.

⁵Campo invisível inserido em formulários de páginas da Web, que são preenchidos somente por *bots* e servem como evidência da natureza não humana e maliciosa destes mecanismos.

4.3.4 Unifei 2017

A base de dados Unifei foi coletada durante o segundo semestre de 2016 na Universidade Federal de Itajubá (Unifei) pelo Grupo de Pesquisas em Engenharia de Sistemas e de Computação (GPESC). Consiste em um universo de *e-mails* que representa a realidade da universidade, com mensagens recebidas pelos professores, STAs e alunos da universidade.

Contendo 862.229 *e-mails*, sendo 353.152 *hams* e 509.077 *spams*, é um *corpus* consideravelmente maior que os outros já apresentados. Entretanto, é importante mencionar que os rótulos dos *e-mails* não foram atribuídos manualmente, mas sim de forma automática pelo anti-*spam* proprietário CanIt-PRO da Roaring Penguin.

Cada mensagem foi submetida aos seguintes pré-processamentos:

- Substituição de *tags* HTML (e seus atributos), símbolos monetários, *links*, números, palavras pequenas⁶ e grandes⁷ por *tokens*.
- Remoção de pontuação, substituição de caracteres com diacríticos⁸ e/ou maiúsculos por seus equivalentes minúsculos simples e correção de codificação.

Depois de pré-processar as mensagens, foram executados os métodos de seleção de características (CHI2, FD e MI) de forma a obter representações vetoriais (numéricas) de cada mensagem.

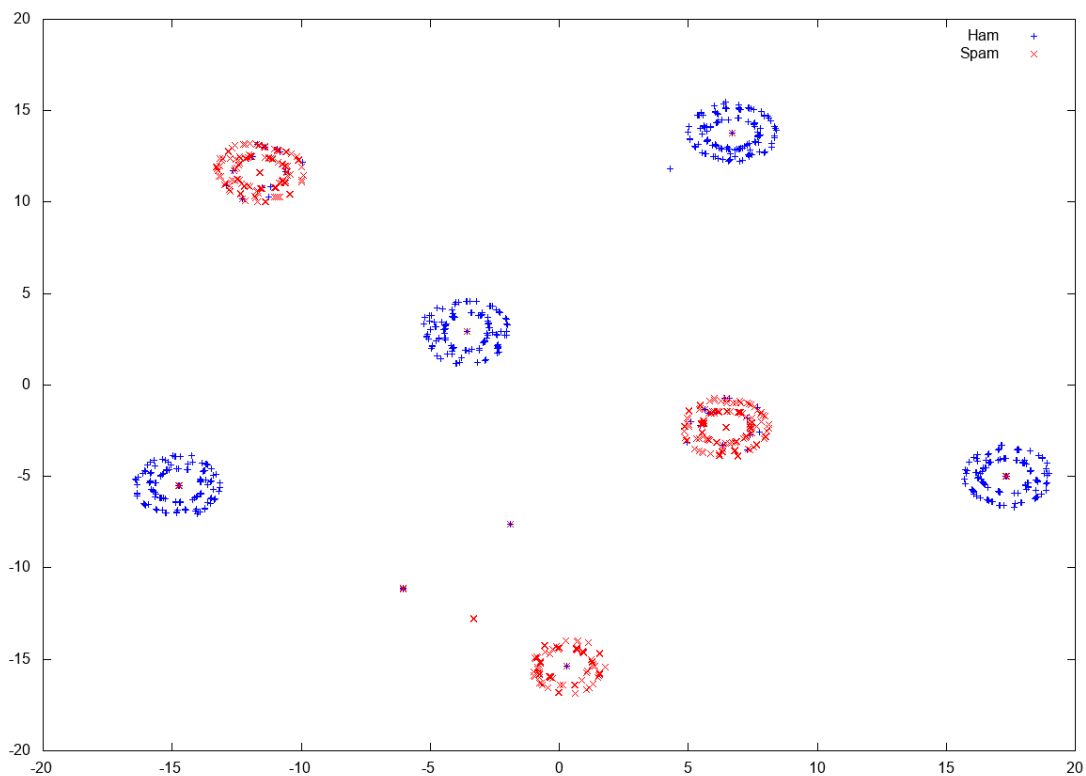
Por fim, utilizando a técnica t-SNE (do inglês *t-Distributed Stochastic Neighbor Embedding*) (MAATEN; HINTON, 2008), foram geradas visualizações bidimensionais do conjunto para cada método de seleção de características, a fim de detectar possíveis inconsistências nos dados. Os resultados apresentados nas figuras 26, 27 e 28, indicam a existência de amostras ambíguas (ou seja, mensagens de classes diferentes mas que ocupam a mesma posição no gráfico). Tal problema motivou a criação do conjunto Unifei 2018, que será apresentado a seguir, na seção 4.3.5.

⁶Cuja quantidade de caracteres é inferior a um limite pré-estabelecido.

⁷Cuja quantidade de caracteres é superior a um limite pré-estabelecido.

⁸Sinais gráficos que alteram a realização fonética (som) de uma letra.

(a) Resultado do t-SNE com vetores de 8 características.



(b) Resultado do t-SNE com vetores de 1024 características.

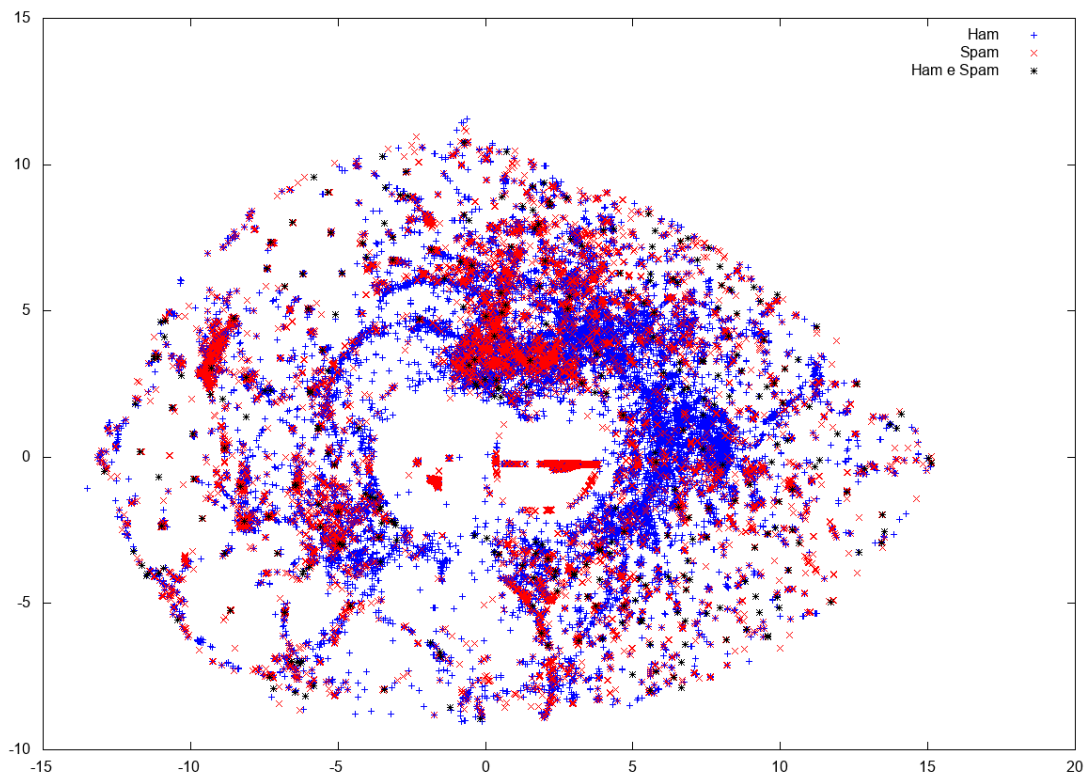
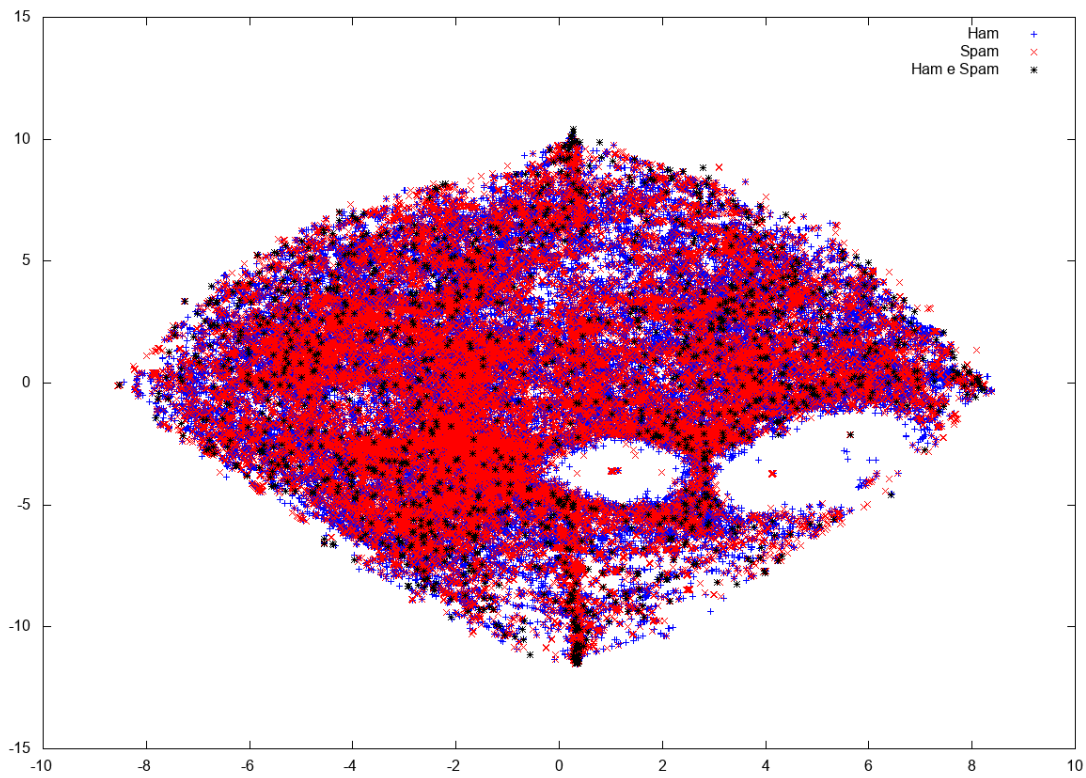


Figura 26: Visualização 2D da base Unifei 2017 (método CHI2).

(a) Resultado do t-SNE com vetores de 8 características.



(b) Resultado do t-SNE com vetores de 1024 características.

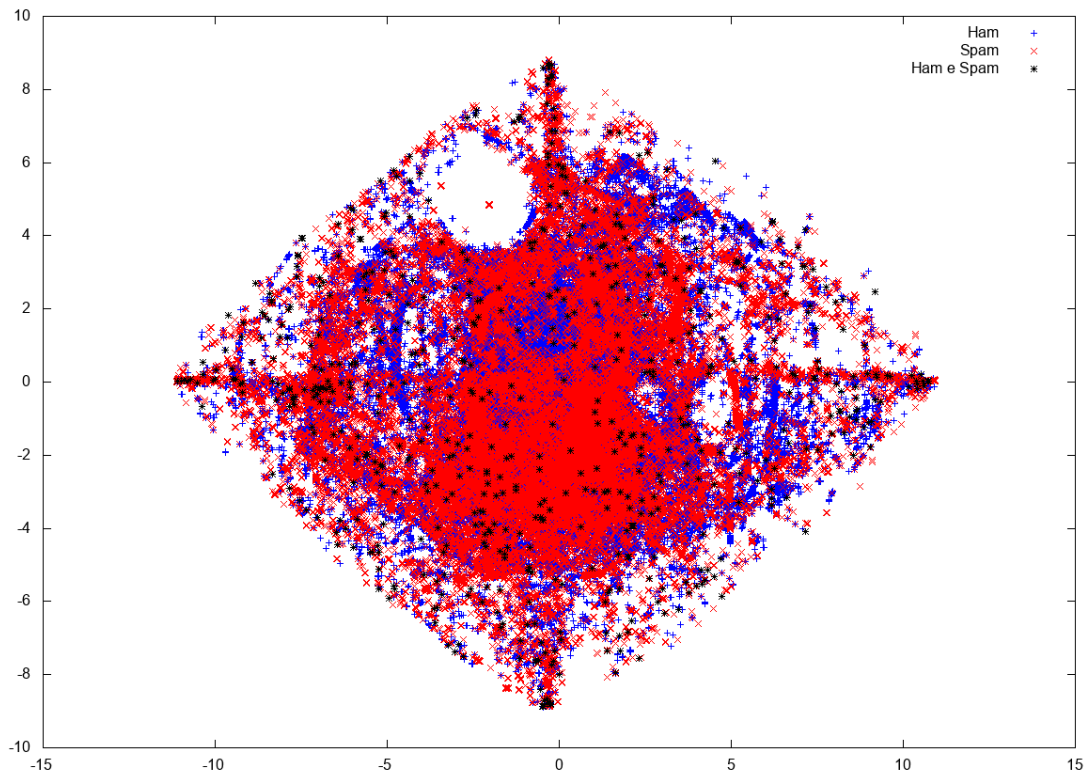
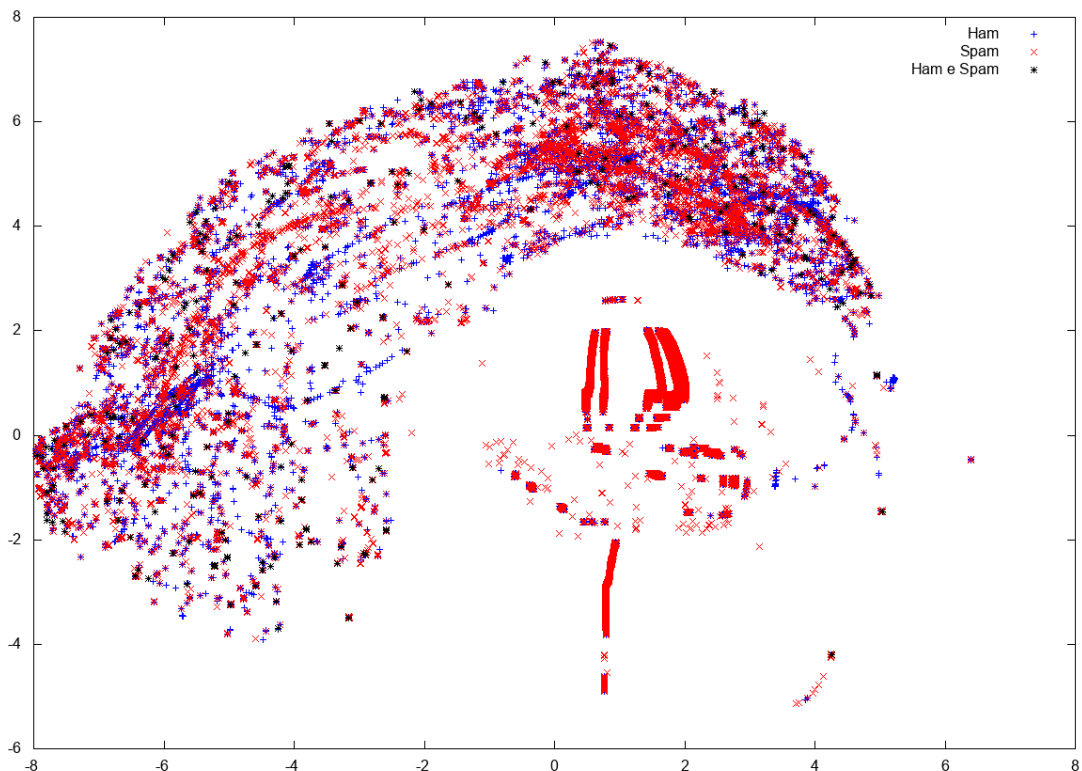


Figura 27: Visualização 2D da base Unifei 2017 (método FD).

(a) Resultado do t-SNE com vetores de 8 características.



(b) Resultado do t-SNE com vetores de 1024 características.

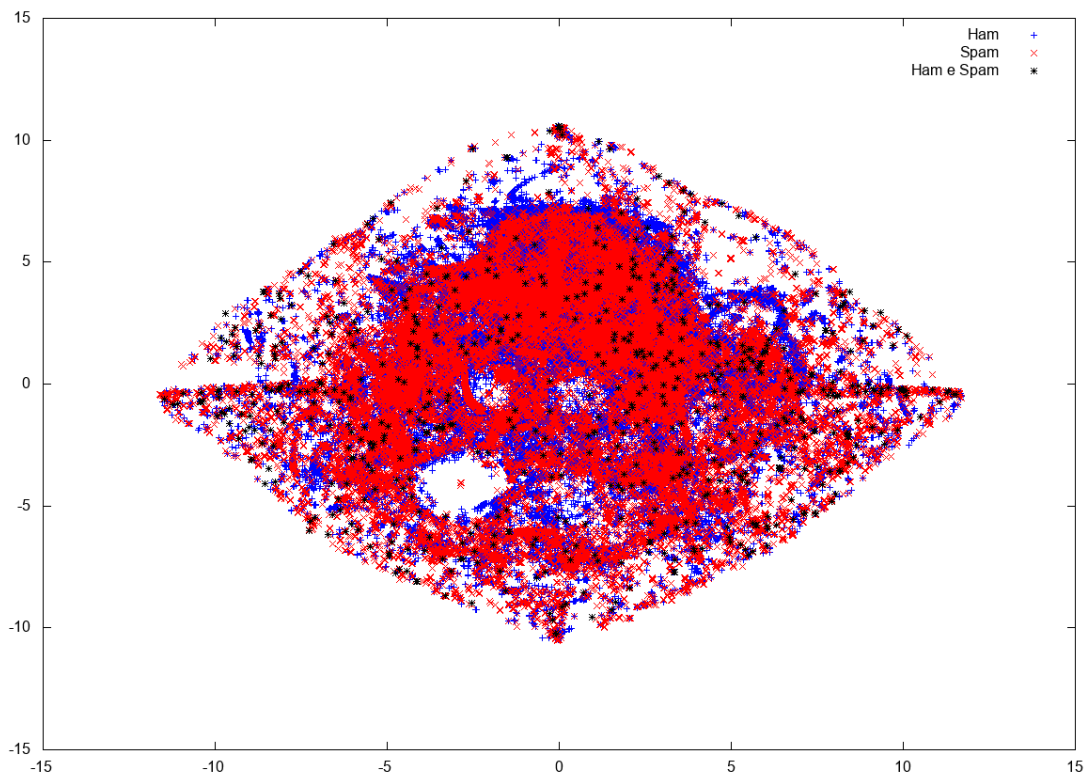


Figura 28: Visualização 2D da base Unifei 2017 (método MI).

4.3.5 Unifei 2018

Considerando que alguns *e-mails* (potencialmente de classes distintas) diferem-se totalmente de todos os outros do conjunto e, portanto, não apresentam qualquer relação de informação mútua nem possuem as palavras mais frequentes do conjunto, as representações vetoriais destas mensagens tornam-se vazias. Neste sentido, os dois problemas a seguir foram apontados na base Unifei 2017:

1. Existência de *e-mails* de classes diferentes com representações vetoriais iguais;
2. Presença de vetores nulos (ou seja, $\vec{x} = [x_1, x_2, \dots, x_n] = [0, 0, \dots, 0]$).

Visando endereçar as consequências do primeiro problema, foi realizado um procedimento de rerrotulagem de mensagens/vetores. Considerando que haja H *hams* e S *spams* com representação vetorial igual, o procedimento resume-se a:

- Se $H > S$, então os S *spams* são rerrotulados como *hams*.
- Se $S > H$, então os H *hams* são rerrotulados como *spams*.

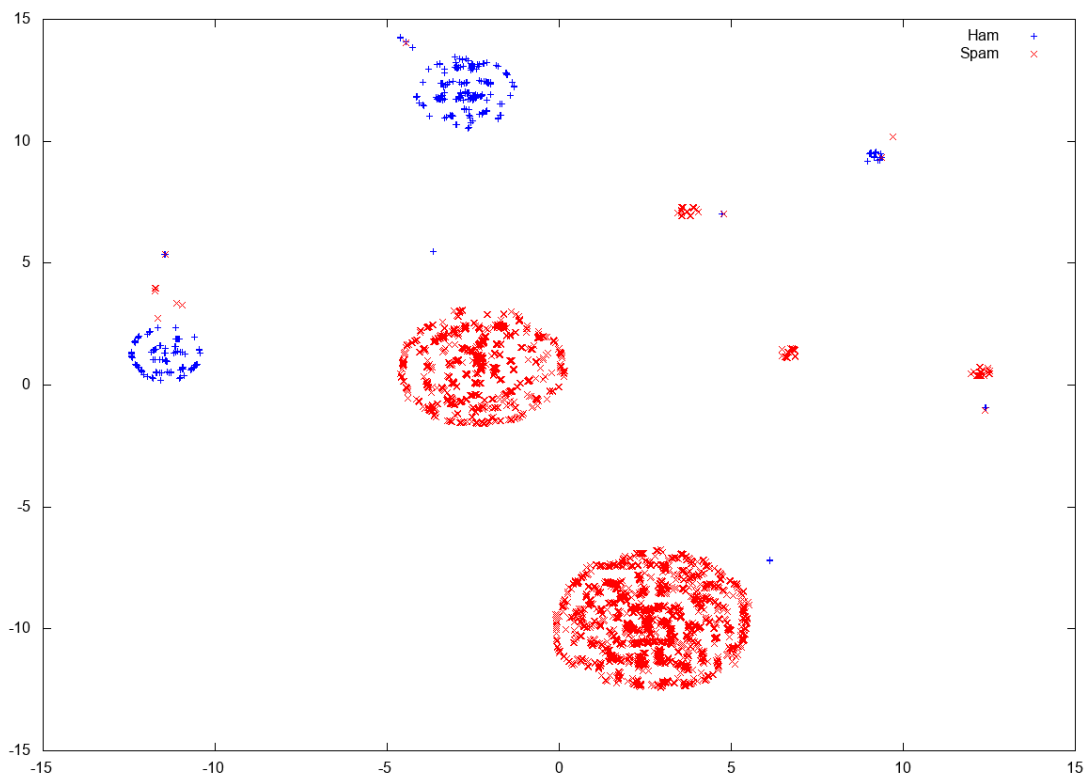
Para contornar o segundo problema, os vetores nulos foram incluídos apenas na etapa de testes (classificação), ou seja, eles não foram apresentados aos métodos de aprendizado de máquina durante a etapa de treinamento (geração de modelo).

Portanto, a base Unifei 2018 consiste em uma versão corrigida e melhorada da base Unifei 2017, obtida pela execução dos seguintes passos:

1. Uso de todos os *e-mails* da base Unifei 2017;
2. Detecção de **mensagens** duplicadas (de classes diferentes);
3. Rerrotulagem das **mensagens** duplicadas identificadas na etapa anterior;
4. Execução dos métodos de seleção de características para obtenção dos vetores (semelhante à Unifei 2017);
5. Detecção dos **vetores** duplicados (de classes diferentes);
6. Rerrotulagem dos **vetores** duplicados identificadas na etapa anterior.

Os gráficos das figuras 29, 30 e 31 comprovam a ausência de amostras ambíguas (ou seja, mensagens de classes diferentes que ocupam a mesma posição no gráfico).

(a) Resultado do t-SNE com vetores de 8 características.



(b) Resultado do t-SNE com vetores de 1024 características.

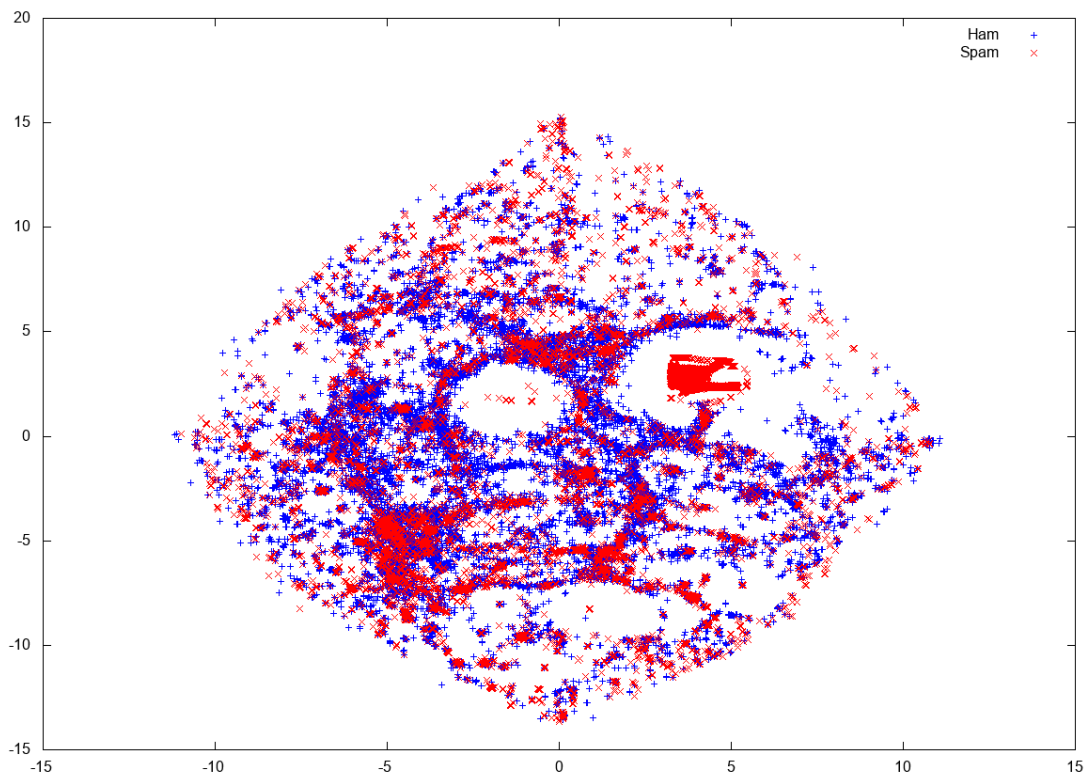
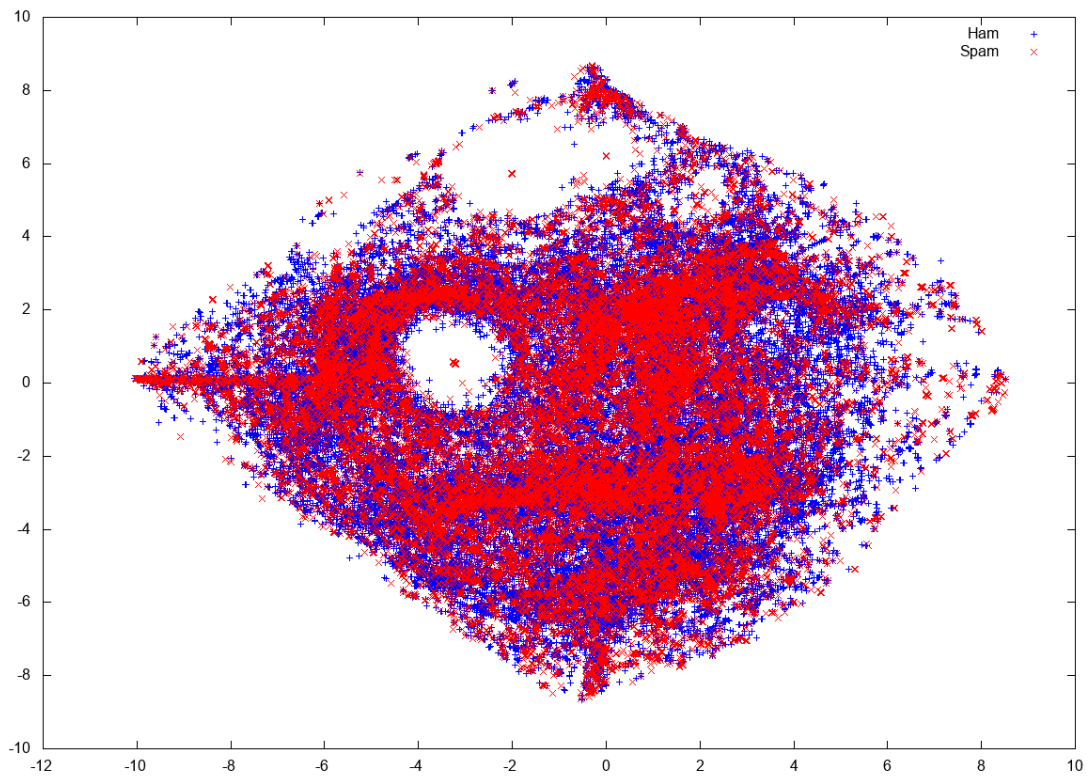


Figura 29: Visualização 2D da base Unifei 2018 (método CHI2).

(a) Resultado do t-SNE com vetores de 8 características.



(b) Resultado do t-SNE com vetores de 1024 características.

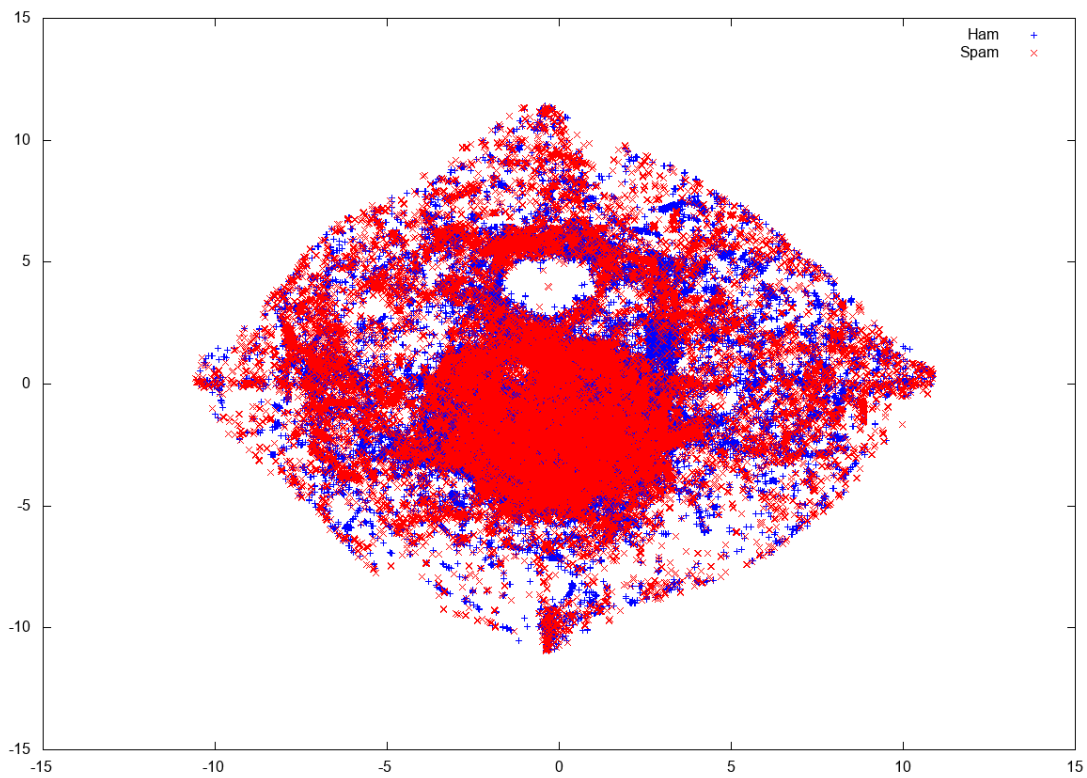
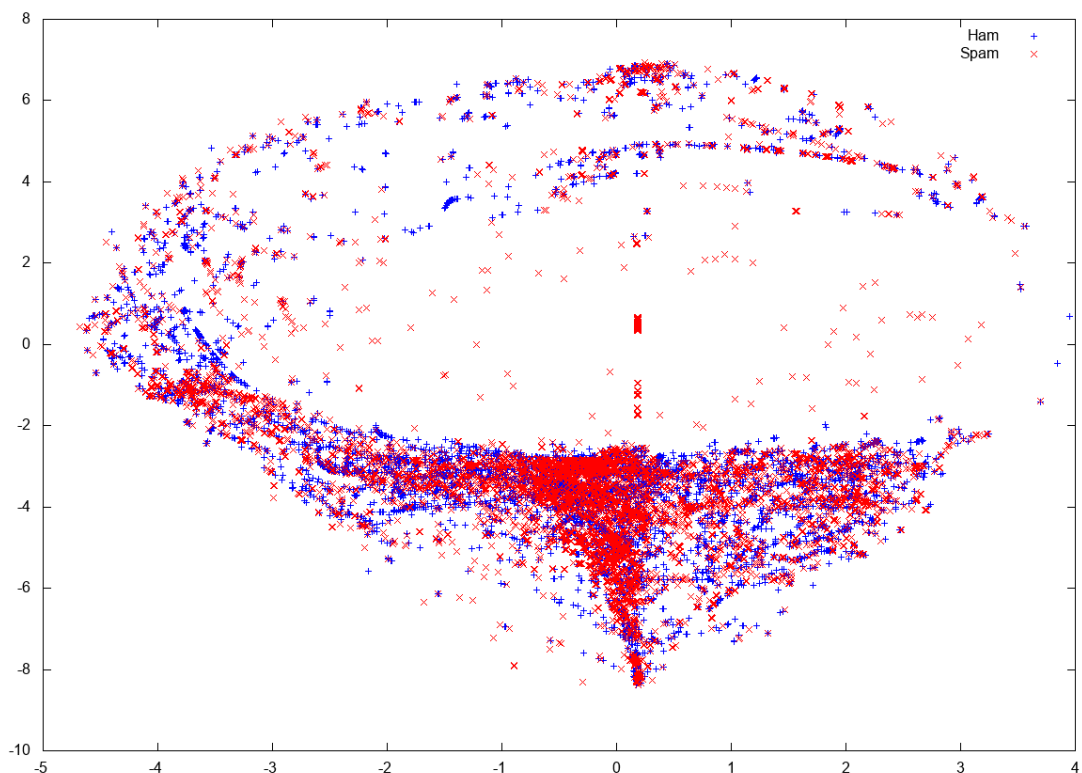


Figura 30: Visualização 2D da base Unifei 2018 (método FD).

(a) Resultado do t-SNE com vetores de 8 características.



(b) Resultado do t-SNE com vetores de 1024 características.

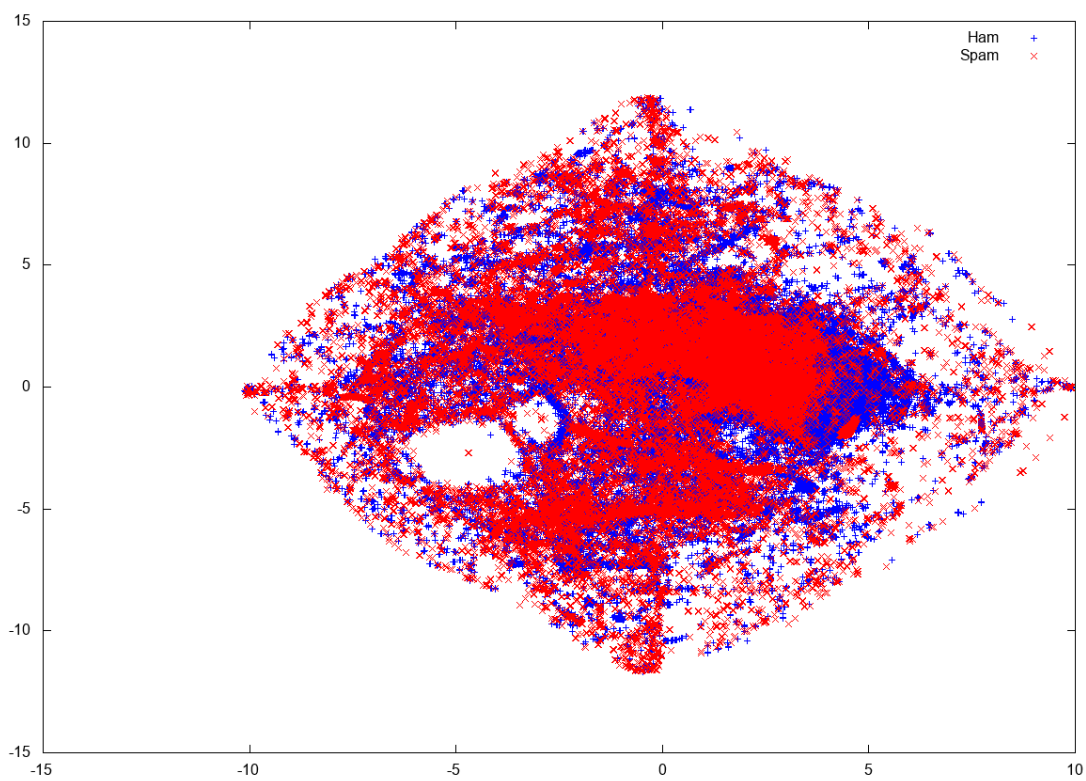


Figura 31: Visualização 2D da base Unifei 2018 (método MI).

4.4 Métricas de desempenho

Trabalhos que envolvam comparações de qualquer natureza requerem a definição de qual é o sujeito de avaliação, bem como sob quais óticas este será avaliado.

No contexto de sistemas anti-*spam*, é desejável que estes sistema não classifiquem mensagens legítimas incorretamente, caso contrário o usuário poderá perder conteúdo importante de forma irreversível, dependendo da política de tratamento de *spam* do sistema. Quando o sistema falha neste tipo de classificação, diz-se que foi obtido um resultado falso positivo.

Além disso, também deseja-se que estes sejam capazes de identificar e filtrar a maior quantidade possível de mensagens indesejadas, protegendo o destinatário de conteúdo potencialmente malicioso ou que implique em perda financeira e/ou de produtividade. Quando o sistema falha neste tipo de classificação, diz-se que foi obtido um resultado falso negativo.

Tendo estas definições em mente, são propostas as métricas de precisão e revocação para avaliação dos modelos de aprendizado de máquina, que correspondem às ausências de resultados falsos-positivos e falsos-negativos, respectivamente. Nos apêndices, os resultados são dados em termos de *ham* e *spam* separadamente.

Adicionalmente, o escore F_1 (explicado em detalhes na seção 4.4.3) faz uma combinação da precisão e revocação, apresentando-se como uma métrica capaz de avaliar um modelo em termos de ambos os tipos de erro e servindo como medida geral da eficácia de um sistema anti-*spam*, já que os resultados são dados em termos de *ham* e *spam* simultaneamente.

Por fim, são apresentadas métricas relacionadas aos tempos de treinamento e execução. A primeira mostra o tempo necessário para construir um modelo preditivo a partir das amostras de treinamento, e a segunda consiste no intervalo necessário para classificar todo o conjunto de testes usando o modelo em questão.

Vale destacar que os resultados são dados em termos de média e intervalo de confiança (com $C = 0,95$) (NEYMAN, 1937) considerando todas as dez execuções.

Para as equações (4.1), (4.2), (4.3), (4.4), (4.5), (4.6), (4.7), (4.8) e (4.9), considere que:

- N_{HAM} é a quantidade total de *hams* no conjunto de testes;
- N_{SPAM} é a quantidade total de *spams* no conjunto de testes;
- $n_{H \rightarrow H}$ é a quantidade de *hams* classificados corretamente;
- $n_{H \rightarrow S}$ é a quantidade de *hams* classificados como *spam*;
- $n_{S \rightarrow S}$ é a quantidade de *spams* classificados corretamente;
- $n_{S \rightarrow H}$ é a quantidade de *spams* classificados como *ham*.

4.4.1 Precisão

Pode ser vista como uma medida de exatidão, pois denota a ausência de falsos positivos. É calculada tanto para *ham* quanto para *spam* pelas equações (4.1) e (4.2), respectivamente. A precisão geral é dada pela equação (4.3).

$$P(HAM) = \frac{n_{H \rightarrow H}}{n_{H \rightarrow H} + n_{S \rightarrow H}} \quad (4.1) \quad P(SPAM) = \frac{n_{S \rightarrow S}}{n_{S \rightarrow S} + n_{H \rightarrow S}} \quad (4.2)$$

$$P(GERAL) = \frac{N_{HAM} * P(HAM) + N_{SPAM} * P(SPAM)}{N_{HAM} + N_{SPAM}} \quad (4.3)$$

4.4.2 Revocação

Pode ser vista como uma medida de completude, pois denota a ausência de falsos negativos. É calculada tanto para *ham* quanto para *spam* pelas equações (4.4) e (4.5), respectivamente. A revocação geral é dada pela equação (4.6).

$$R(HAM) = \frac{n_{H \rightarrow H}}{n_{H \rightarrow H} + n_{H \rightarrow S}} \quad (4.4) \quad R(SPAM) = \frac{n_{S \rightarrow S}}{n_{S \rightarrow S} + n_{S \rightarrow H}} \quad (4.5)$$

$$R(GERAL) = \frac{N_{HAM} * R(HAM) + N_{SPAM} * R(SPAM)}{N_{HAM} + N_{SPAM}} \quad (4.6)$$

4.4.3 Escore F_1

Na análise estatística da classificação binária, o escore F_1 é uma medida que considera tanto a precisão como a revocação de um teste para calcular sua pontuação. É dada pela média harmônica entre estes valores (vide equações (4.7), (4.8) e (4.9)) e atinge seu melhor e pior valor em 1 e 0, respectivamente.

$$F_1(HAM) = 2 \cdot \frac{1}{\frac{1}{P(HAM)} + \frac{1}{R(HAM)}} = 2 \cdot \frac{P(HAM) \cdot R(HAM)}{P(HAM) + R(HAM)} \quad (4.7)$$

$$F_1(SPAM) = 2 \cdot \frac{1}{\frac{1}{P(SPAM)} + \frac{1}{R(SPAM)}} = 2 \cdot \frac{P(SPAM) \cdot R(SPAM)}{P(SPAM) + R(SPAM)} \quad (4.8)$$

$$F_1(GERAL) = \frac{N_{HAM} * F_1(HAM) + N_{SPAM} * F_1(SPAM)}{N_{HAM} + N_{SPAM}} \quad (4.9)$$

4.4.4 Tempo de Treinamento

É definido como o tempo que o algoritmo de classificação levou para ser treinado e gerar um modelo aplicável ao conjunto de testes. É dado por:

$$t_{treinamento} = t'_T - t_T \quad (4.10)$$

Em que t'_T e t_T marcam, respectivamente, os instantes inicial e final do treinamento.

4.4.5 Tempo de Classificação

É definido como o tempo que o modelo levou para classificar todas as amostras do conjunto de treinamento. É dado por:

$$t_{classificação} = t'_C - t_C \quad (4.11)$$

Em que t'_C e t_C marcam, respectivamente, os instantes inicial e final da classificação.

5 Experimentos e Resultados

Visando facilitar a análise e entendimento dos resultados, as conclusões foram agrupadas em três seções que representam cada aspecto macroscópico deste trabalho:

- Seleção de características e padrões zerados (seção 5.1)
 - métodos de seleção de características e quantidade de padrões zerados;
 - quantidade de características e quantidade de padrões zerados;
 - métodos de seleção de características e desempenho (em termos de escore F_1 médio) em cada base de dados.
- Redução de dimensionalidade (seção 5.2)
 - métodos de seleção de características e quantidades (absoluta e relativa) de características após a redução de dimensionalidade;
 - proporcionalidade entre as quantidades de características antes e após a redução de dimensionalidade, bem como as anomalias identificadas.
- Treinamento e classificação (seção 5.3)
 - comparação dos escores F_1 , tempo de treinamento e de classificação obtidos pelos métodos de classificação no *ranking* de cada base de dados;
 - métodos de seleção de características e ocorrências no *ranking*;
 - quantidade original de características e ocorrências no *ranking*.

Vale ressaltar que os resultados experimentais, usados para *renderizar* as visualizações e tirar as conclusões, foram incluídos na íntegra nos apêndices. Além disso, detalhes sobre as configurações do *hardware* e as versões específicas dos *softwares* utilizados também foram disponibilizados.

5.1 Seleção de características e padrões zerados

A tabela 6 e as figuras 32, 33, 34, 35, 36 e 37 mostram, para cada método de seleção de características, os percentuais de padrões zerados em cada conjunto de dados.

Tabela 6: Quantidades de padrões zerados nos conjuntos de dados — NF : número de características; H_0 : percentual de *hams* zerados; S_0 : percentual de *spams* zerados.

Método	NF	Ling Spam		Spam Assassin		TREC		Unifei 2017		Unifei 2018	
		$H_0(\%)$	$S_0(\%)$	$H_0(\%)$	$S_0(\%)$	$H_0(\%)$	$S_0(\%)$	$H_0(\%)$	$S_0(\%)$	$H_0(\%)$	$S_0(\%)$
CHI2	8	97.8	99.3	4.2	11.7	99.6	98.9	99.7	99.9	98.8	98.6
	16	95.8	99.3	3.9	9.9	99.1	97.5	98.6	98.9	98.0	97.6
	32	91.9	99.3	3.5	9.1	98.9	97.2	98.4	98.8	96.0	96.2
	64	85.2	99.3	2.6	8.5	98.2	96.6	98.3	98.8	93.6	94.6
	128	74.2	99.3	2.0	6.5	97.2	94.6	97.8	98.7	89.9	88.7
	256	59.7	99.3	0.6	0.6	94.1	91.9	95.5	97.9	85.9	86.9
	512	34.2	61.9	0.6	0.6	82.2	78.4	86.3	88.2	69.5	60.2
	1024	34.2	49.1	0.5	0.5	57.5	45.1	42.2	47.3	41.2	36.2
FD	8	0.0	0.0	0.0	0.7	4.4	2.4	0.8	2.1	0.2	2.5
	16	0.0	0.0	0.0	0.2	4.3	2.4	0.8	2.1	0.1	2.5
	32	0.0	0.0	0.0	0.2	4.2	2.4	0.8	2.1	0.1	2.5
	64	0.0	0.0	0.0	0.2	4.1	2.4	0.8	2.1	0.1	2.5
	128	0.0	0.0	0.0	0.2	4.0	2.4	0.8	2.0	0.1	2.5
	256	0.0	0.0	0.0	0.2	4.0	2.4	0.8	2.0	0.1	2.5
	512	0.0	0.0	0.0	0.2	3.9	2.3	0.8	2.0	0.1	2.5
	1024	0.0	0.0	0.0	0.2	3.9	2.3	0.7	2.0	0.1	2.5
MI	8	0.0	0.0	0.0	1.0	7.8	3.0	15.1	14.9	14.9	15.0
	16	0.0	0.0	0.0	0.9	5.9	3.0	4.7	6.7	4.0	7.2
	32	0.0	0.0	0.0	0.2	4.4	2.8	4.7	6.7	4.0	7.2
	64	0.0	0.0	0.0	0.2	4.3	2.4	0.8	2.1	0.1	2.5
	128	0.0	0.0	0.0	0.2	4.3	2.4	0.8	2.1	0.1	2.5
	256	0.0	0.0	0.0	0.2	4.2	2.4	0.8	2.0	0.1	2.5
	512	0.0	0.0	0.0	0.2	4.2	2.4	0.8	2.0	0.1	2.5
	1024	0.0	0.0	0.0	0.2	3.9	2.3	0.8	2.0	0.1	2.5

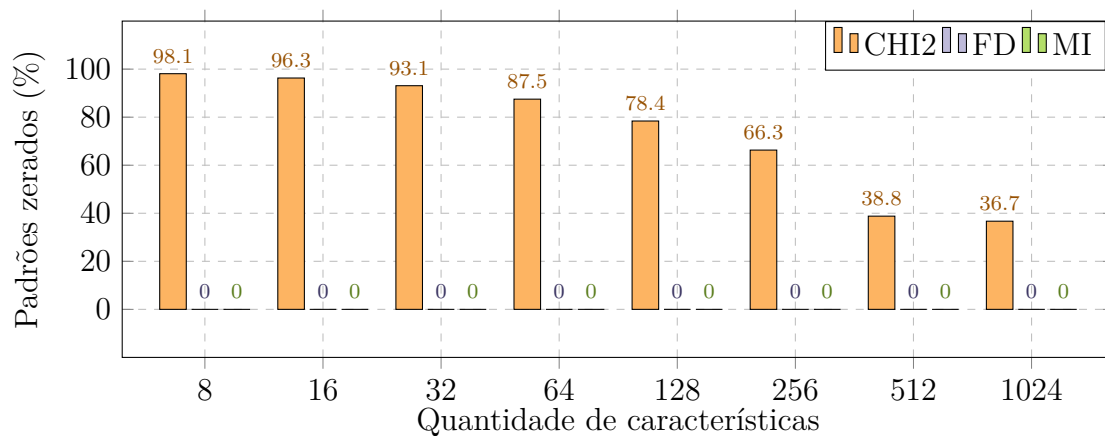


Figura 32: Padrões zerados na base Ling Spam.

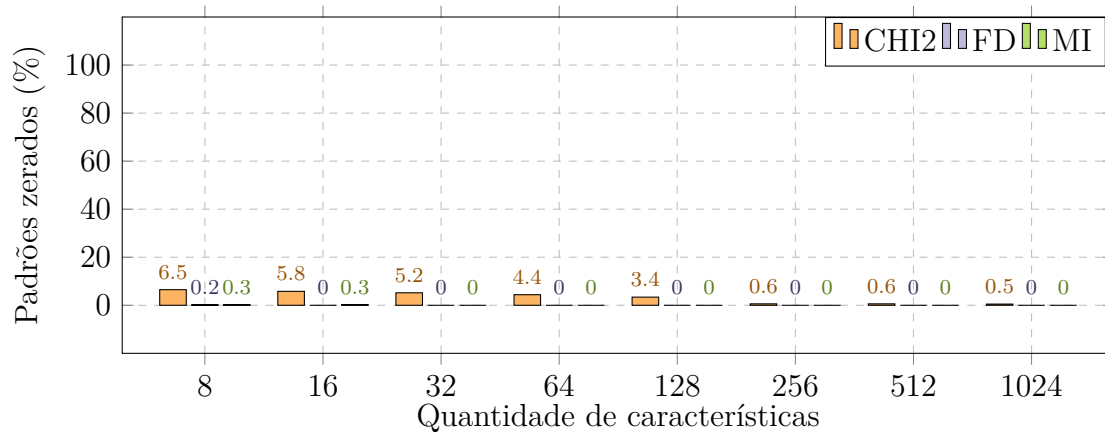


Figura 33: Padrões zerados na base Spam Assassin.

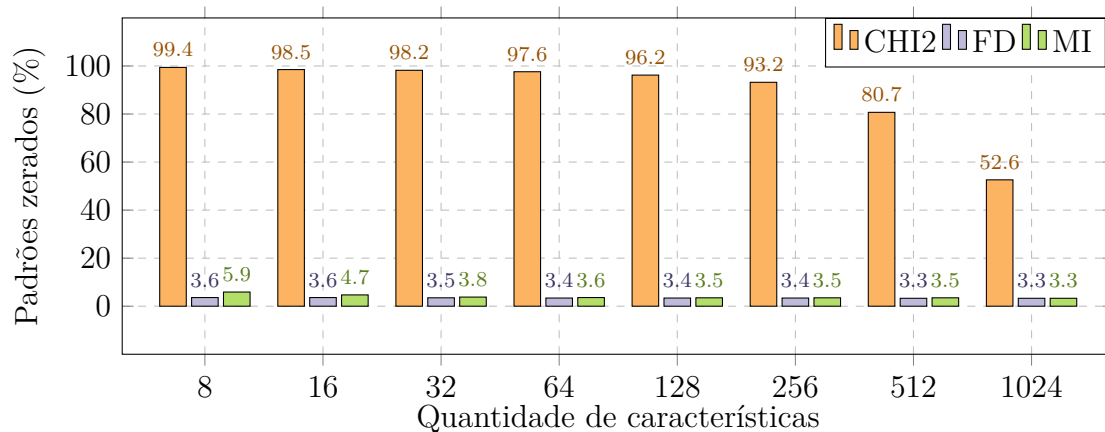


Figura 34: Padrões zerados na base TREC.

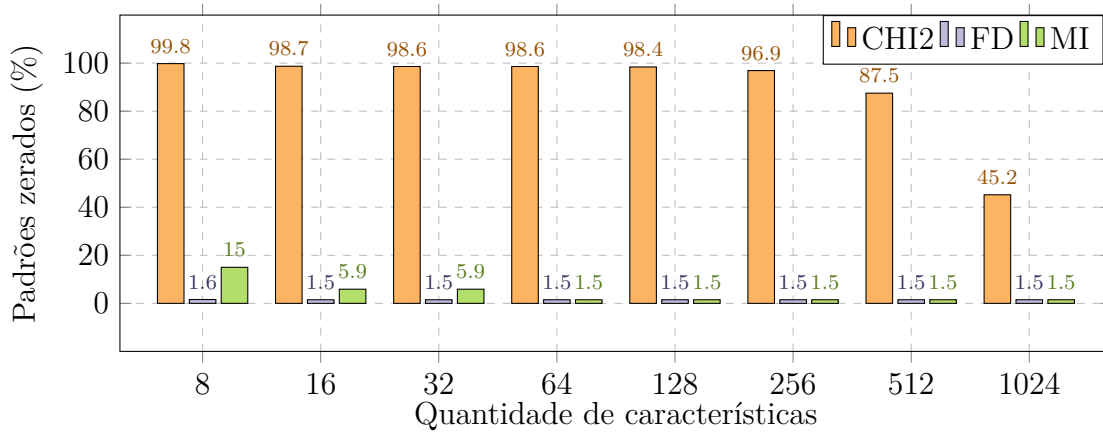


Figura 35: Padrões zerados na base Unifei 2017.

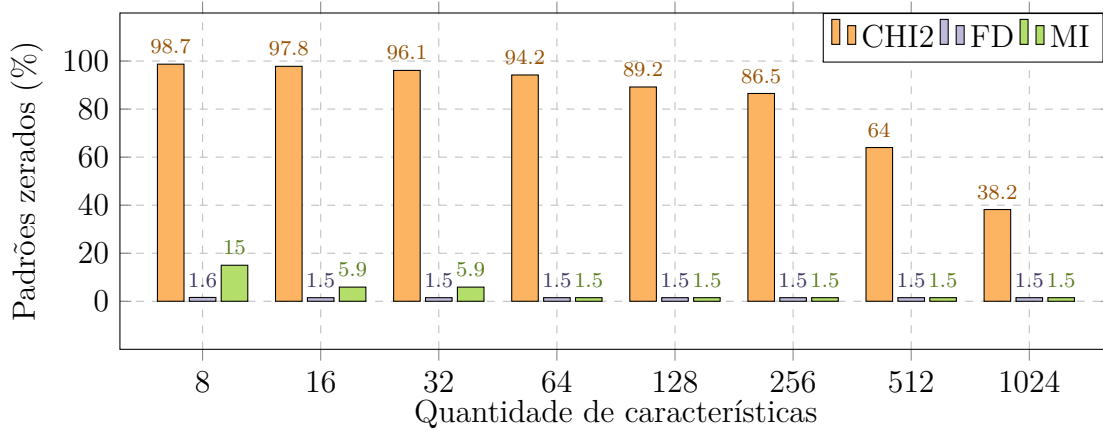


Figura 36: Padrões zerados na base Unifei 2018.

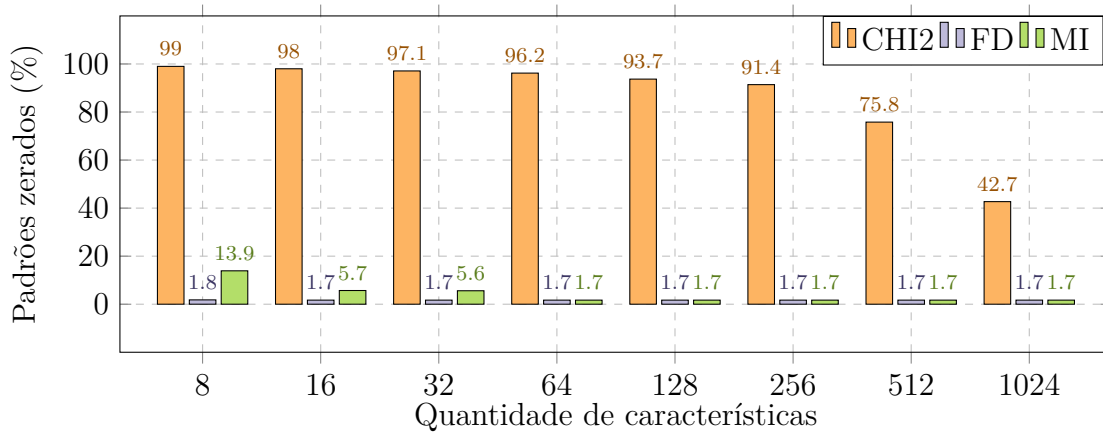


Figura 37: Padrões zerados em todas as bases.

Para cada conjunto de dados, os resultados de todos os modelos de aprendizado de máquina foram agrupados em faixas de escore F_1 médio — “ruim” ($<80\%$), “médio” (entre 80% e 90%) ou “bom” ($>90\%$). Em seguida, foram contabilizadas as ocorrências dos métodos de seleção de características em cada faixa de escore F_1 . Estes resultados estão representados nas figuras 38, 39, 40, 41, 42 e 43.

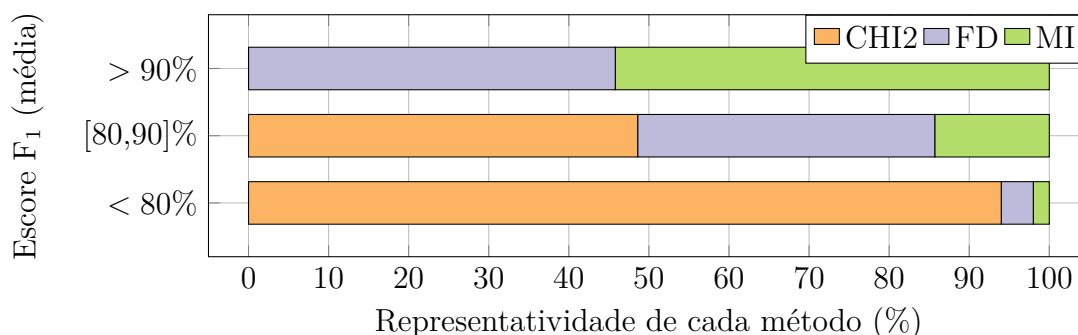


Figura 38: Escore F_1 vs métodos de seleção de características (Ling Spam).

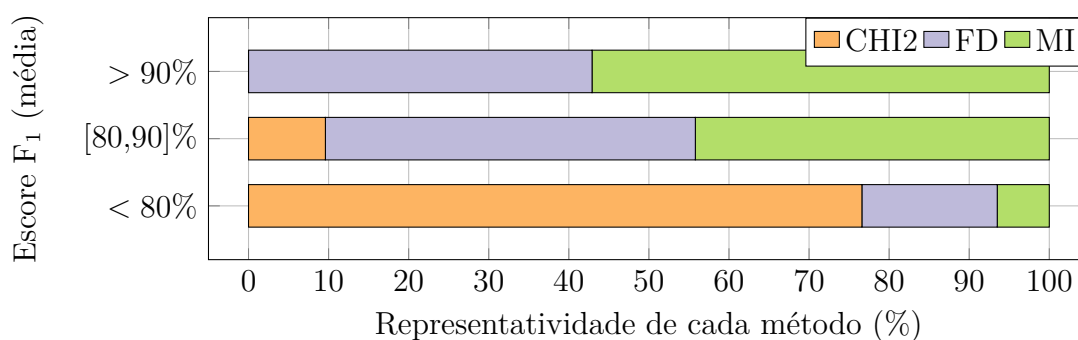


Figura 39: Escore F_1 vs métodos de seleção de características (Spam Assassin).

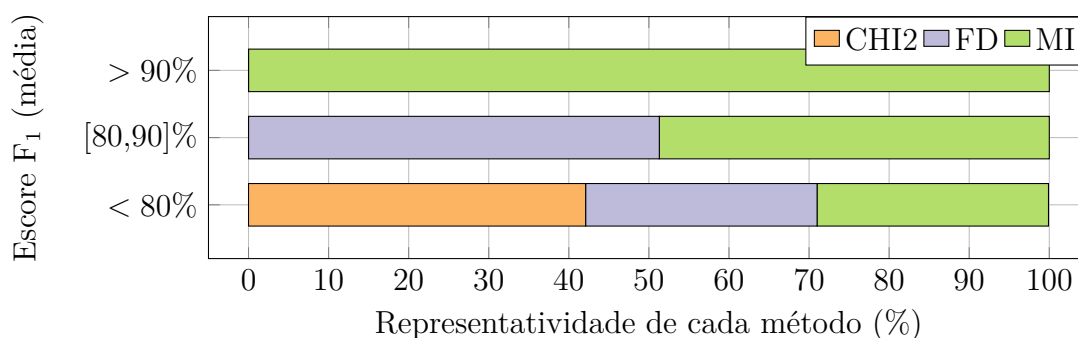


Figura 40: Escore F_1 vs métodos de seleção de características (TREC).

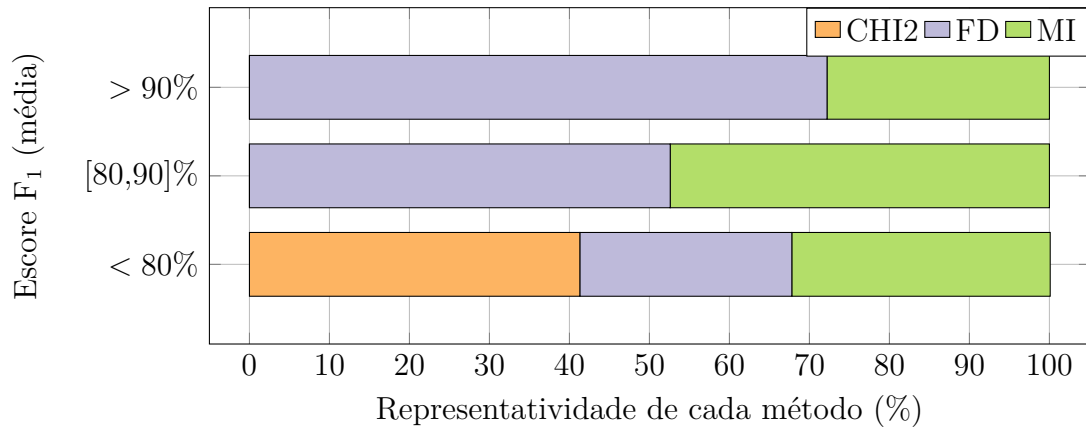


Figura 41: Escore F_1 vs métodos de seleção de características (Unifei 2017).

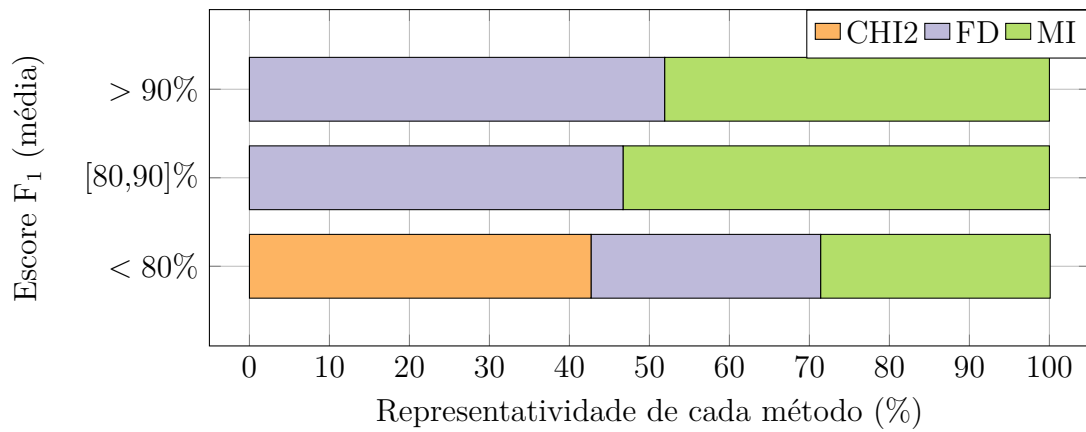


Figura 42: Escore F_1 vs métodos de seleção de características (Unifei 2018).

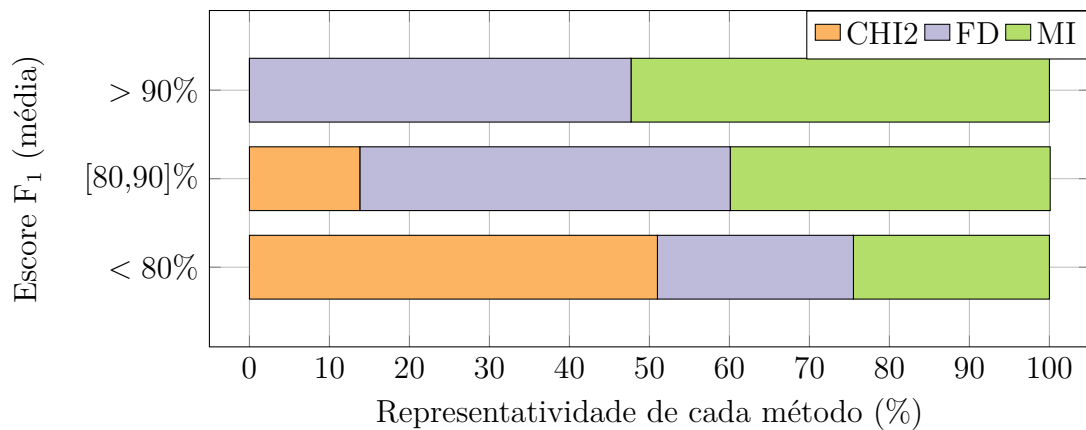


Figura 43: Escore F_1 vs métodos de seleção de características em todas as bases.

Diante dos resultados obtidos, podem ser tiradas as seguintes conclusões:

- Em todos os conjuntos de dados, o método CHI2 foi o que apresentou maior quantidade de padrões zerados. Exceto na base Spam Assassin, esta quantia chegou a representar quase a totalidade (mais de 99%) das mensagens;
- Em todos os conjuntos, o método FD foi o que apresentou menos padrões zerados, tendo o maior percentual (3,6%) ocorrido na base TREC. No caso do MI, o maior percentual (15%) ocorreu nas bases Unifei 2017 e 2018;
- Para todos os métodos de seleção de características e em todos os conjuntos de dados, a quantidade de padrões zerados mostrou-se inversamente proporcional ao número de características das amostras;
- A maioria (51%) dos resultados ruins de escore F_1 médio foi obtida com o método CHI2 (especialmente em Ling Spam e Spam Assassin), seguidos por FD (24,5%) e MI (24,5%);
- A maioria (46,3%) dos resultados médios de escore F_1 médio foi obtida com o método FD (especialmente em TREC e Unifei 2017), seguido por MI (40%) e CHI2 (13,7%);
- A maioria (52,3%) dos resultados bons de escore F_1 médio foi obtida com o método MI (especialmente em Ling Spam, Spam Assassin e TREC), seguido por FD (47,7%). Não houve ocorrências com o método CHI2;
- Nas bases Ling Spam e Spam Assassin, o método MI apresentou os melhores resultados de escore F_1 médio na maioria (54,2% e 57,1%, respectivamente) dos casos;
- Na base TREC, o método MI obteve os melhores resultados de escore F_1 médio em 100% dos casos;
- Nas bases Unifei 2017 e 2018, o método FD apresentou os melhores resultados de escore F_1 médio na maioria (72,2% e 51,9%, respectivamente) dos casos.

5.2 Redução de dimensionalidade

As quantidades de características — NF e NF' , antes e após a redução de dimensionalidade, respectivamente — em cada base de *e-mails*, para cada método de seleção, podem ser encontradas nas tabelas 7, 8, 9, 10 e 11 e nas figuras 44, 45, 46, 47 e 48.

Tabela 7: Redução de dimensionalidade na base Ling Spam.

NF	$NF'(CHI2)$	$NF'(FD)$	$NF'(MI)$
8	1 (12,5%)	4 (50%)	6 (75%)
16	1 (6,3%)	4 (25%)	12 (75%)
32	1 (3,1%)	10 (31,3%)	19 (59,4%)
64	3 (4,7%)	19 (29,7%)	34 (53,1%)
128	4 (3,1%)	38 (29,7%)	52 (40,6%)
256	14 (5,5%)	58 (22,7%)	68 (26,6%)
512	14 (2,7%)	103 (20,1%)	139 (27,1%)
1024	242 (23,6%)	121 (11,8%)	157 (15,3%)

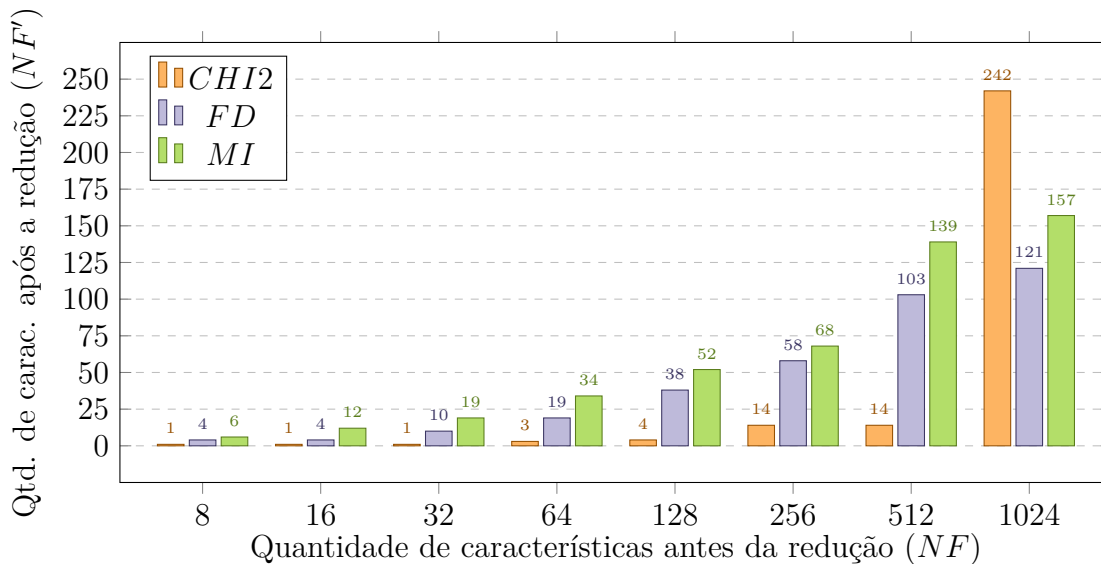


Figura 44: Redução de dimensionalidade na base Ling Spam.

Tabela 8: Redução de dimensionalidade na base Spam Assassin.

NF	$NF'(CHI2)$	$NF'(FD)$	$NF'(MI)$
8	4 (50%)	5 (62,5%)	3 (37,5%)
16	7 (43,8%)	6 (37,5%)	6 (37,5%)
32	14 (43,8%)	9 (28,1%)	18 (56,3%)
64	20 (31,3%)	16 (25%)	33 (51,6%)
128	25 (19,5%)	26 (20,3%)	58 (45,3%)
256	29 (11,3%)	53 (20,7%)	85 (33,2%)
512	7 (1,4%)	67 (13,1%)	96 (18,8%)
1024	20 (2%)	58 (5,7%)	84 (8,2%)

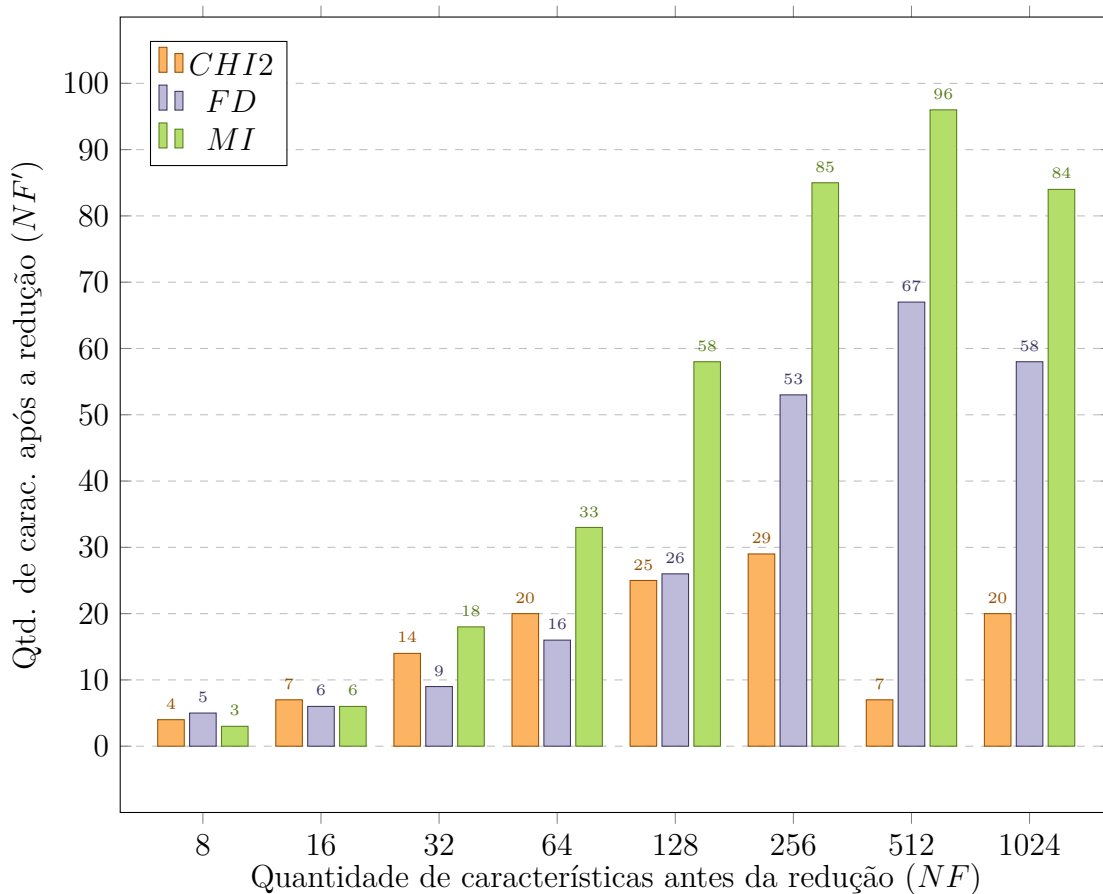


Figura 45: Redução de dimensionalidade na base Spam Assassin.

Tabela 9: Redução de dimensionalidade na base TREC.

NF	$NF'(CHI2)$	$NF'(FD)$	$NF'(MI)$
8	4 (50%)	6 (75%)	5 (62,5%)
16	11 (68,8%)	7 (43,8%)	6 (37,5%)
32	10 (31,3%)	14 (43,8%)	9 (28,1%)
64	12 (18,8%)	7 (10,9%)	16 (25%)
128	45 (35,2%)	22 (17,2%)	26 (20,3%)
256	68 (26,6%)	24 (9,4%)	37 (14,5%)
512	133 (26%)	51 (10%)	66 (12,9%)
1024	47 (4,6%)	49 (4,8%)	204 (19,9%)

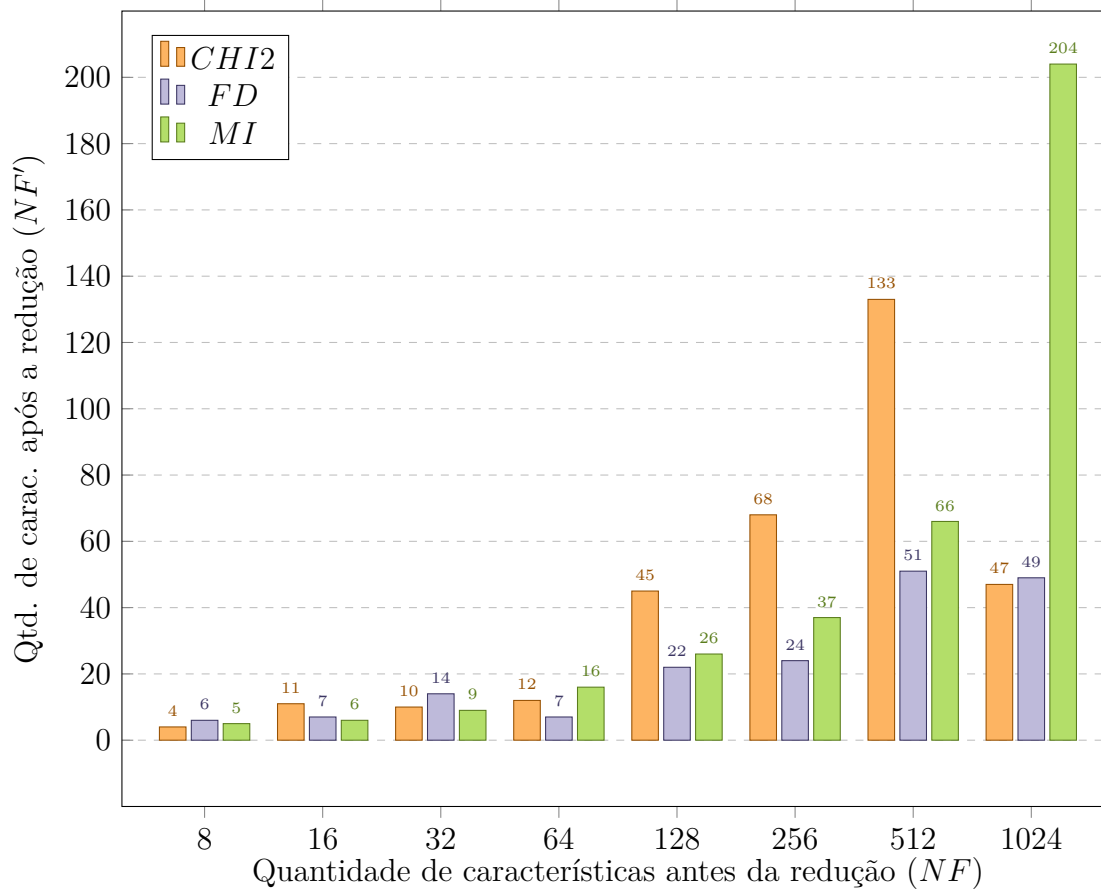


Figura 46: Redução de dimensionalidade na base TREC.

Tabela 10: Redução de dimensionalidade na base Unifei 2017.

NF	$NF'(CHI2)$	$NF'(FD)$	$NF'(MI)$
8	8 (100,0%)	5 (62,5%)	4 (50,0%)
16	10 (62,5%)	3 (18,8%)	7 (43,8%)
32	20 (62,5%)	6 (18,8%)	5 (15,6%)
64	48 (75,0%)	22 (34,4%)	9 (14,1%)
128	86 (67,2%)	65 (50,8%)	18 (14,1%)
256	34 (13,3%)	51 (19,9%)	31 (12,1%)
512	45 (8,8%)	101 (19,7%)	125 (24,4%)
1024	17 (1,7%)	105 (10,3%)	208 (20,3%)

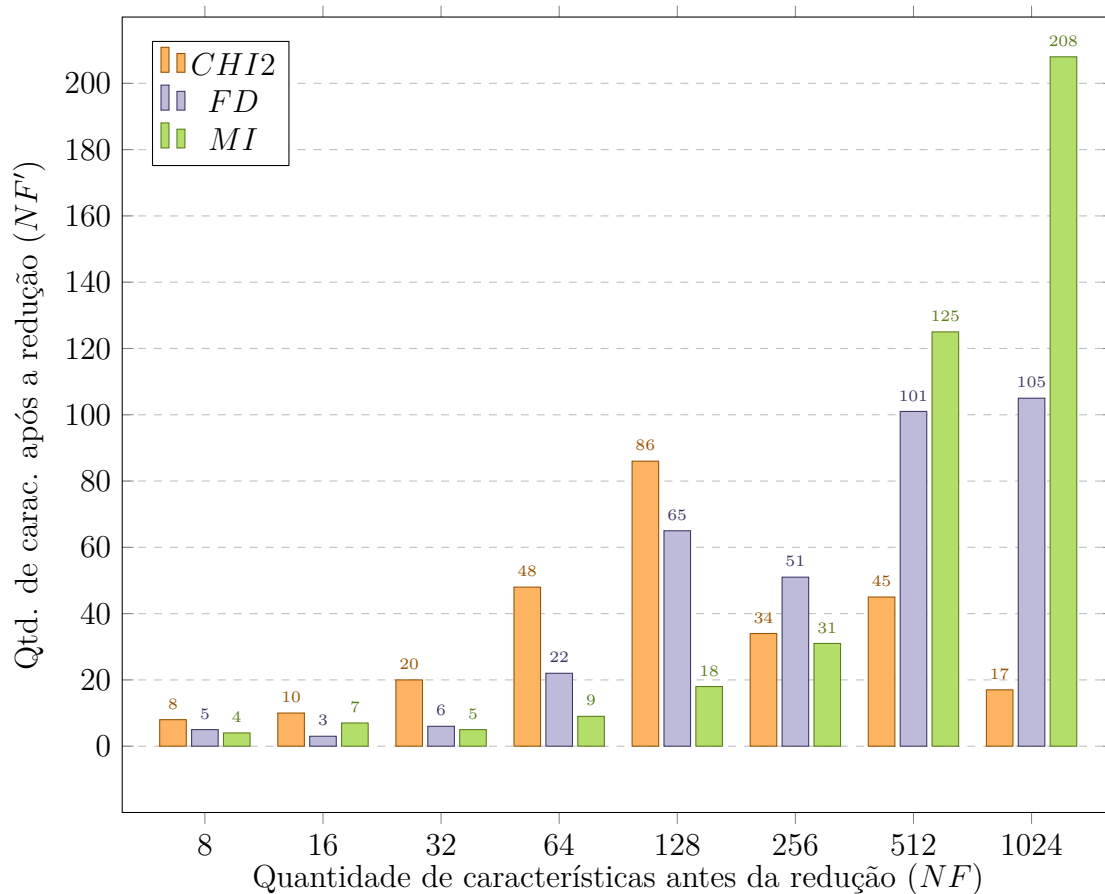


Figura 47: Redução de dimensionalidade na base Unifei 2017.

Tabela 11: Redução de dimensionalidade na base Unifei 2018.

NF	$NF'(CHI2)$	$NF'(FD)$	$NF'(MI)$
8	4 (50%)	5 (62,5%)	3 (37,5%)
16	16 (100%)	3 (18,8%)	3 (18,8%)
32	28 (87,5%)	8 (25%)	7 (21,9%)
64	42 (65,6%)	11 (17,2%)	10 (15,6%)
128	69 (53,9%)	24 (18,8%)	30 (23,4%)
256	116 (45,3%)	60 (23,4%)	49 (19,1%)
512	121 (23,6%)	79 (15,4%)	118 (23%)
1024	152 (14,8%)	77 (7,5%)	155 (15,1%)

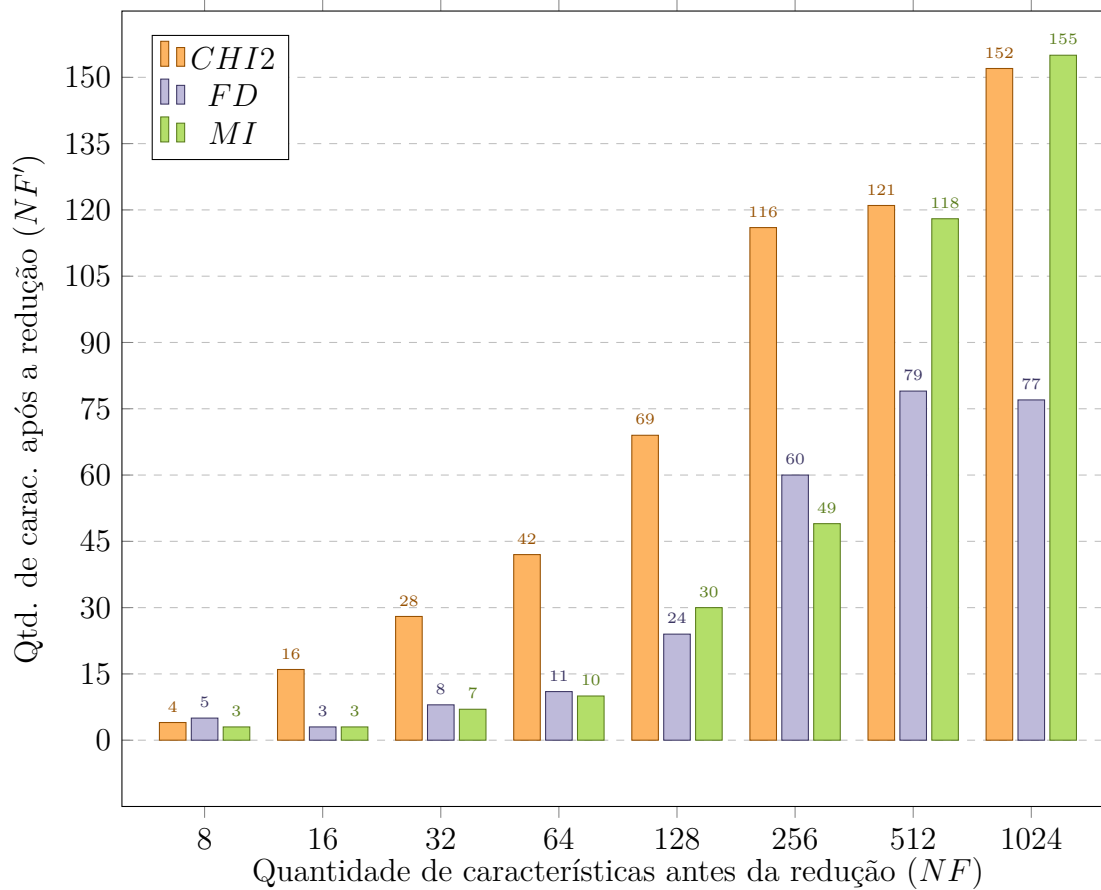


Figura 48: Redução de dimensionalidade na base Unifei 2018.

Diante dos resultados obtidos, podem ser tiradas as seguintes conclusões:

- Nas bases Ling Spam e Spam Assassin, o método CHI2 foi o que mais teve características descartadas pelo algoritmo de redução de dimensionalidade;
- Nas bases TREC, Unifei 2017 e Unifei 2018, o método FD foi o que mais teve características descartadas pelo algoritmo de redução de dimensionalidade;
- Na maioria dos casos, conforme a quantidade de características antes da redução aumentava, a quantidade após a redução também aumentava. Entretanto, algumas anomalias foram detectadas (especialmente quando a dimensionalidade original era de 1024 características):
 - CHI2: bases Spam Assassin (512/7 e 1024/20), TREC (32/10 e 1024/47) e Unifei 2017 (256/34, 512/45 e 1024/17);
 - FD: bases Spam Assassin (1024/58), TREC (64/7 e 1024/49), Unifei 2017 (256/51) e Unifei 2018 (16/3 e 1024/77);
 - MI: bases Spam Assassin (1024/84) e Unifei 2017 (32/5).
- A maior redução percentual para o método CHI2 foi de 98,6% — de 512 (100%) para 7 (1,4%) características —, na base Spam Assassin e as menores foram nas bases Unifei 2017 e 2018, nas quais não houve redução quando a dimensionalidade original era de 8 e 16 características;
- A maior redução percentual para o método FD foi de 95,2% — de 1024 (100%) para 49 (4,8%) características — e a menor foi de 25% — 8 (100%) para 6 (75%) características —, ambas na base TREC;
- A maior redução percentual para o método MI foi de 91,8 % — de 1024 (100%) para 84 (8,2%) características —, na base Spam Assassin e a menor foi de 25 % — de 8 e 16 (100%) para 6 e 12 (75%) características, respectivamente —, na base Ling Spam.

5.3 Treinamento e classificação

Conforme descrito no capítulo 4, cada modelo de aprendizado de máquina é treinado e avaliado dez vezes com cada conjunto de dados e, em seguida, são calculadas as métricas, e seus valores são representados em termos de média e intervalo de confiança. Assim, um experimento foi realizado para cada método e, como já dito, os resultados individuais completos estão disponíveis nos apêndices.

Os melhores resultados, ordenados por escore F_1 médio, para cada base de *e-mails*, são exibidos nas tabelas 12 a 16 e nas figuras 49 a 53. Nas tabelas, *FS* é o método de seleção de características e *NF* e *NF'* são, respectivamente, as quantidades de características antes e após a redução de dimensionalidade.

Além do *ranking* dos modelos de aprendizado de máquina apresentados nas figuras já mencionadas, resultados complementares podem ser encontrados nas figuras 54 e 55.

O gráfico da figura 54 mostra o quanto cada método de seleção de características esteve presente nos melhores resultados, considerando todos os conjuntos de dados. Esta visualização foi obtida contabilizando quantas vezes cada método (CHI2, FD e MI) apareceu nos *rankings* das tabelas 12, 13, 14, 15 e 16.

De forma análoga, o gráfico da figura 55 mostra o quanto cada quantidade de características antes da redução esteve presente nos melhores resultados, considerando todos os conjuntos de dados. Esta visualização foi obtida contabilizando quantas vezes cada quantidade (8, 16, ..., 1024) apareceu nos *rankings* das tabelas 12, 13, 14, 15 e 16.

Tabela 12: *Ranking* dos classificadores na base Ling Spam.

FS	NF	NF'	Método	Escore F_1	Tempo de treinamento	Tempo de classificação
FD	1024	121	AB-M1	99,23 \pm 0,07	00:00:05.321 \pm 00:00:00.249	00:00:00.018 \pm 00:00:00.001
FD	512	103	AODE	98,55 \pm 0,05	00:00:00.236 \pm 00:00:00.000	00:00:00.418 \pm 00:00:00.000
MI	512	139	SLP	98,42 \pm 0,21	00:00:06.196 \pm 00:00:00.017	00:00:00.011 \pm 00:00:00.000
FD	1024	121	REPT	98,26 \pm 0,14	00:00:00.308 \pm 00:00:00.010	00:00:00.004 \pm 00:00:00.000
FD	1024	121	P-SVM	98,20 \pm 0,15	00:00:03.431 \pm 00:00:00.352	00:00:00.388 \pm 00:00:00.002
FD	1024	121	L-SVM	98,16 \pm 0,08	00:00:00.250 \pm 00:00:00.005	00:00:00.145 \pm 00:00:00.002
MI	512	139	RBF	97,01 \pm 0,32	00:00:18.229 \pm 00:00:00.918	00:00:05.242 \pm 00:00:00.021
MI	64	34	NB	95,20 \pm 0,50	00:00:00.022 \pm 00:00:00.000	00:00:00.098 \pm 00:00:00.001

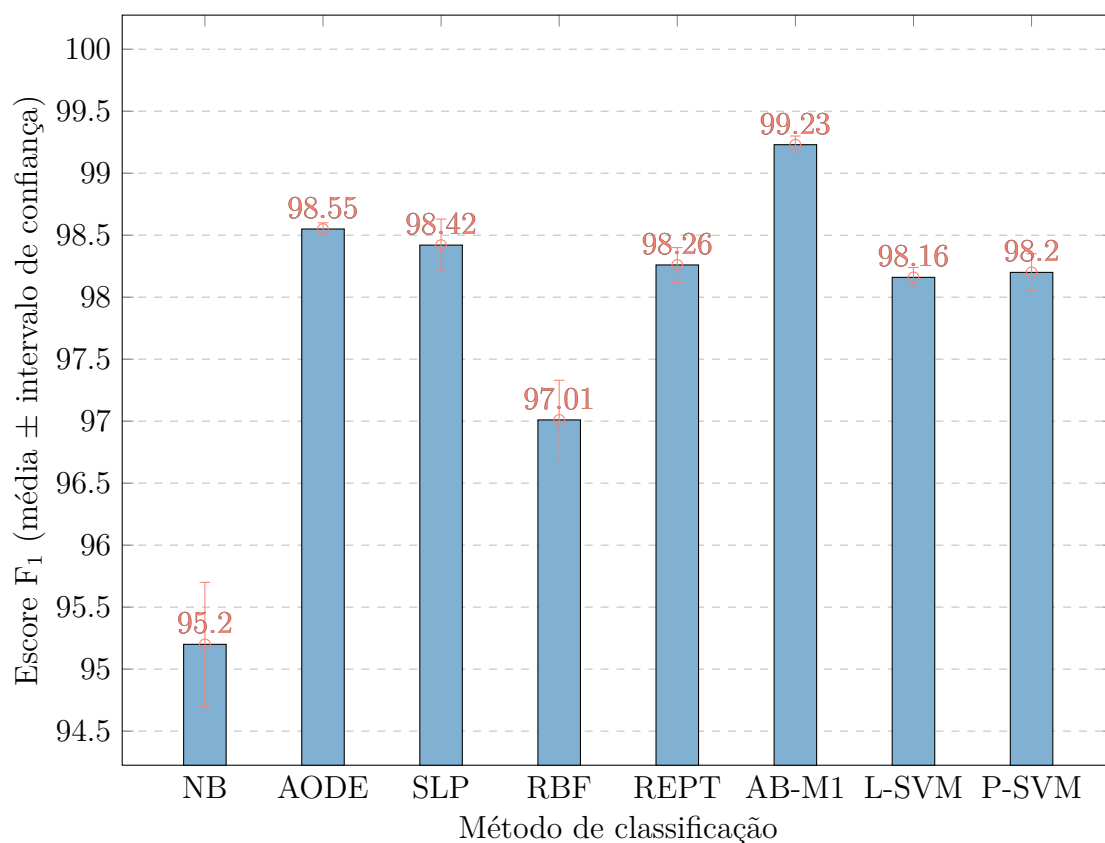
Figura 49: *Ranking* dos classificadores na base Ling Spam.

Tabela 13: *Ranking* dos classificadores na base Spam Assassin.

FS	NF	NF'	Método	Escore F_1	Tempo de treinamento	Tempo de classificação
MI	128	58	AB-M1	$96,97 \pm 0,29$	$00:00:00.955 \pm 00:00:00.034$	$00:00:00.007 \pm 00:00:00.000$
MI	256	85	P-SVM	$95,48 \pm 0,27$	$00:00:01.835 \pm 00:00:00.091$	$00:00:00.135 \pm 00:00:00.001$
MI	128	58	SLP	$95,22 \pm 0,20$	$00:00:01.104 \pm 00:00:00.029$	$00:00:00.002 \pm 00:00:00.000$
MI	128	58	L-SVM	$95,04 \pm 0,28$	$00:00:00.087 \pm 00:00:00.001$	$00:00:00.071 \pm 00:00:00.003$
MI	128	58	AODE	$94,70 \pm 0,22$	$00:00:00.037 \pm 00:00:00.000$	$00:00:00.066 \pm 00:00:00.000$
MI	128	58	REPT	$94,04 \pm 0,35$	$00:00:00.064 \pm 00:00:00.003$	$00:00:00.001 \pm 00:00:00.000$
FD	256	53	RBF	$93,38 \pm 0,47$	$00:00:02.782 \pm 00:00:00.221$	$00:00:00.729 \pm 00:00:00.002$
FD	512	67	NB	$89,74 \pm 0,37$	$00:00:00.017 \pm 00:00:00.000$	$00:00:00.078 \pm 00:00:00.002$

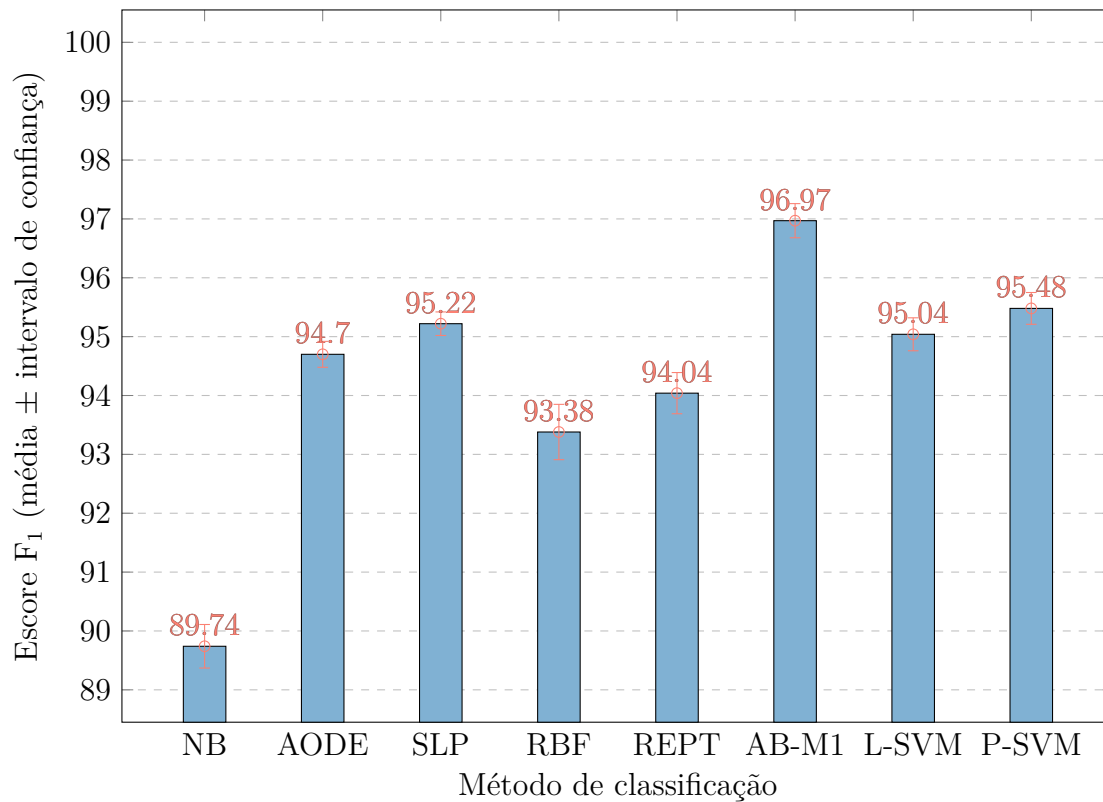
Figura 50: *Ranking* dos classificadores na base Spam Assassin.

Tabela 14: *Ranking* dos classificadores na base TREC.

FS	NF	NF'	Método	Escore F_1	Tempo de treinamento	Tempo de classificação
MI	1024	204	AB-MI	$90,63 \pm 0,07$	$00:19:21.268 \pm 00:00:29.870$	$00:00:01.610 \pm 00:00:00.104$
MI	1024	204	REPT	$89,05 \pm 0,06$	$00:01:30.832 \pm 00:00:01.972$	$00:00:00.136 \pm 00:00:00.003$
MI	1024	204	AODE	$87,67 \pm 0,07$	$00:00:13.285 \pm 00:00:00.050$	$00:00:30.387 \pm 00:00:00.568$
MI	1024	204	P-SVM	$86,80 \pm 0,11$	$00:08:49.209 \pm 00:01:02.402$	$00:00:09.555 \pm 00:00:00.008$
MI	1024	204	L-SVM	$82,75 \pm 2,73$	$00:00:03.994 \pm 00:00:00.113$	$00:00:01.943 \pm 00:00:00.013$
MI	1024	204	SLP	$81,67 \pm 3,73$	$00:01:48.535 \pm 00:00:00.208$	$00:00:00.201 \pm 00:00:00.000$
MI	128	26	RBF	$77,20 \pm 1,42$	$00:01:54.845 \pm 00:00:24.283$	$00:00:08.976 \pm 00:00:00.026$
FD	256	24	NB	$72,08 \pm 0,70$	$00:00:00.312 \pm 00:00:00.000$	$00:00:00.813 \pm 00:00:00.008$

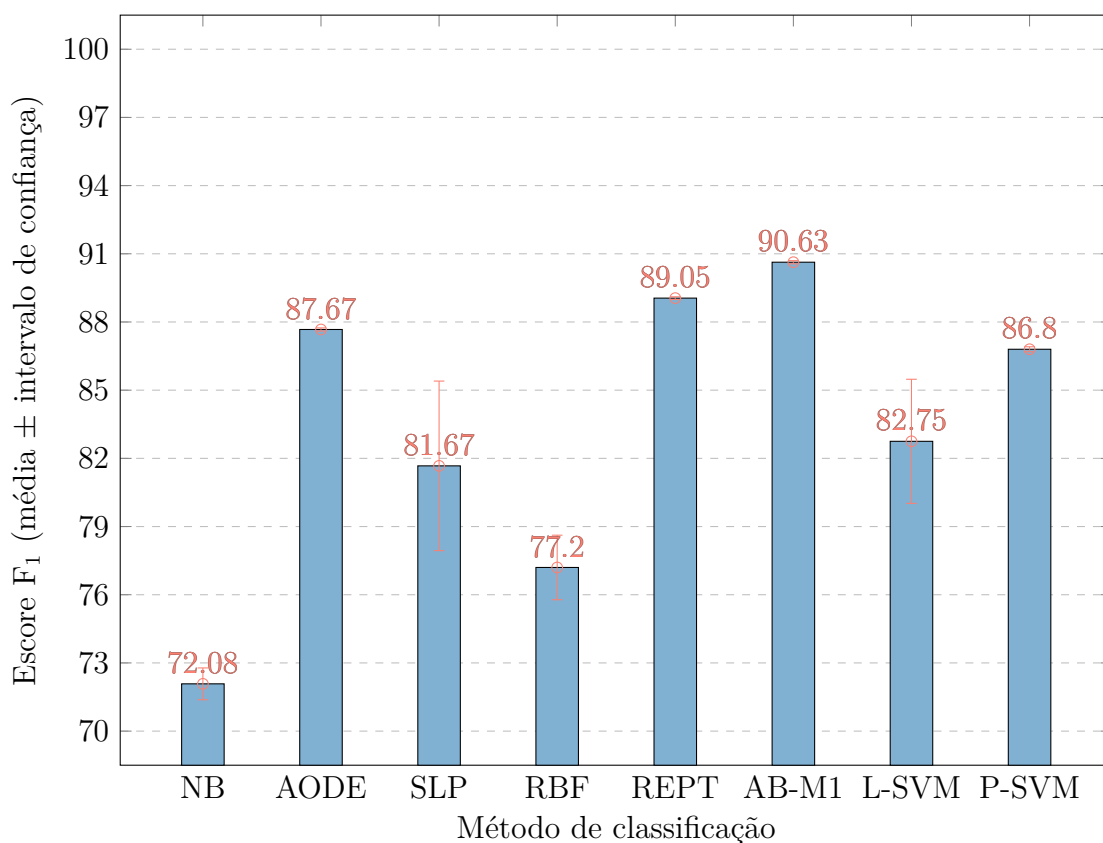
Figura 51: *Ranking* dos classificadores na base TREC.

Tabela 15: *Ranking* dos classificadores na base Unifei 2017.

FS	NF	NF'	Método	Escore F_1	Tempo de treinamento	Tempo de classificação
FD	128	65	AB-M1	$93,96 \pm 0,02$	$00:17:11.884 \pm 00:00:38.035$	$00:00:03.276 \pm 00:00:00.121$
FD	128	65	REPT	$93,27 \pm 0,02$	$00:01:18.676 \pm 00:00:01.901$	$00:00:00.356 \pm 00:00:00.002$
FD	128	65	AODE	$92,37 \pm 0,06$	$00:00:29.522 \pm 00:00:00.336$	$00:00:17.162 \pm 00:00:00.255$
MI	1024	208	P-SVM	$83,00 \pm 0,09$	$01:32:39.551 \pm 00:28:47.995$	$00:00:40.832 \pm 00:00:00.019$
MI	1024	208	L-SVM	$80,91 \pm 0,41$	$00:00:22.785 \pm 00:00:01.653$	$00:00:08.090 \pm 00:00:00.028$
MI	1024	208	SLP	$78,78 \pm 1,32$	$00:07:43.886 \pm 00:00:00.809$	$00:00:00.839 \pm 00:00:00.011$
MI	512	125	RBF	$77,52 \pm 0,23$	$00:24:12.225 \pm 00:03:07.822$	$00:04:07.769 \pm 00:00:00.561$
MI	512	125	NB	$75,44 \pm 0,17$	$00:00:07.338 \pm 00:00:00.064$	$00:00:16.914 \pm 00:00:00.034$

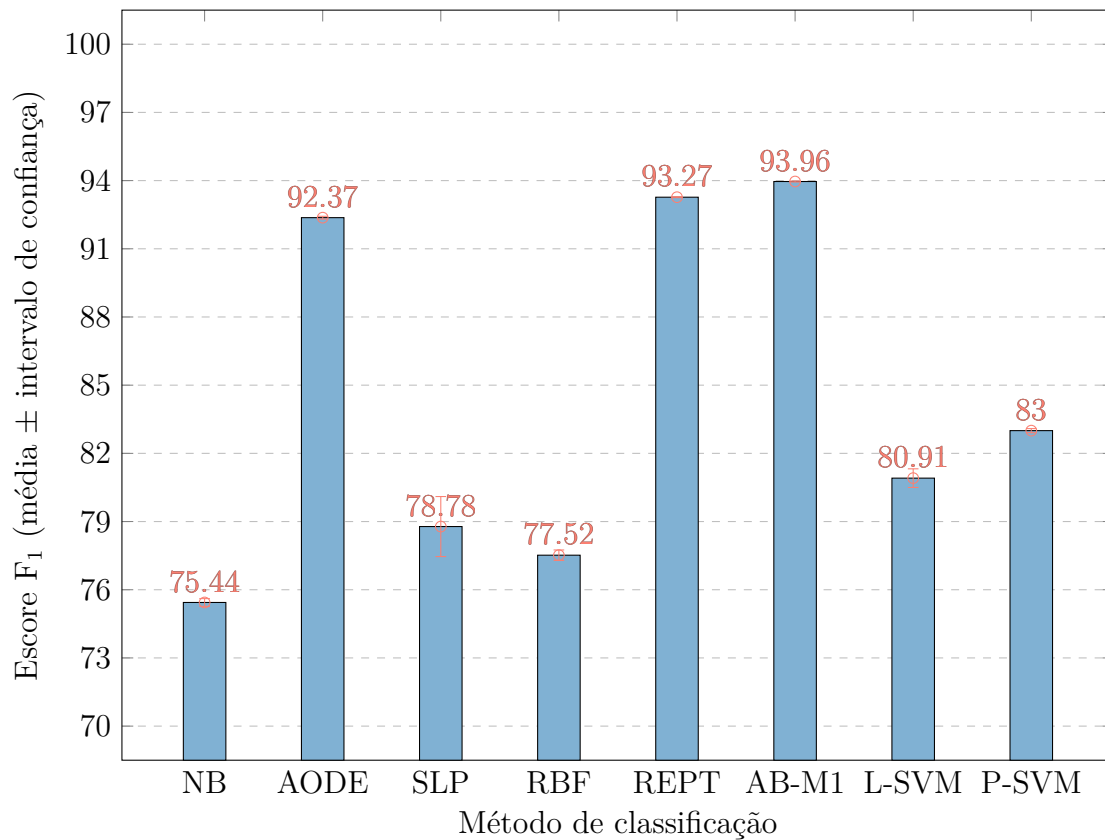
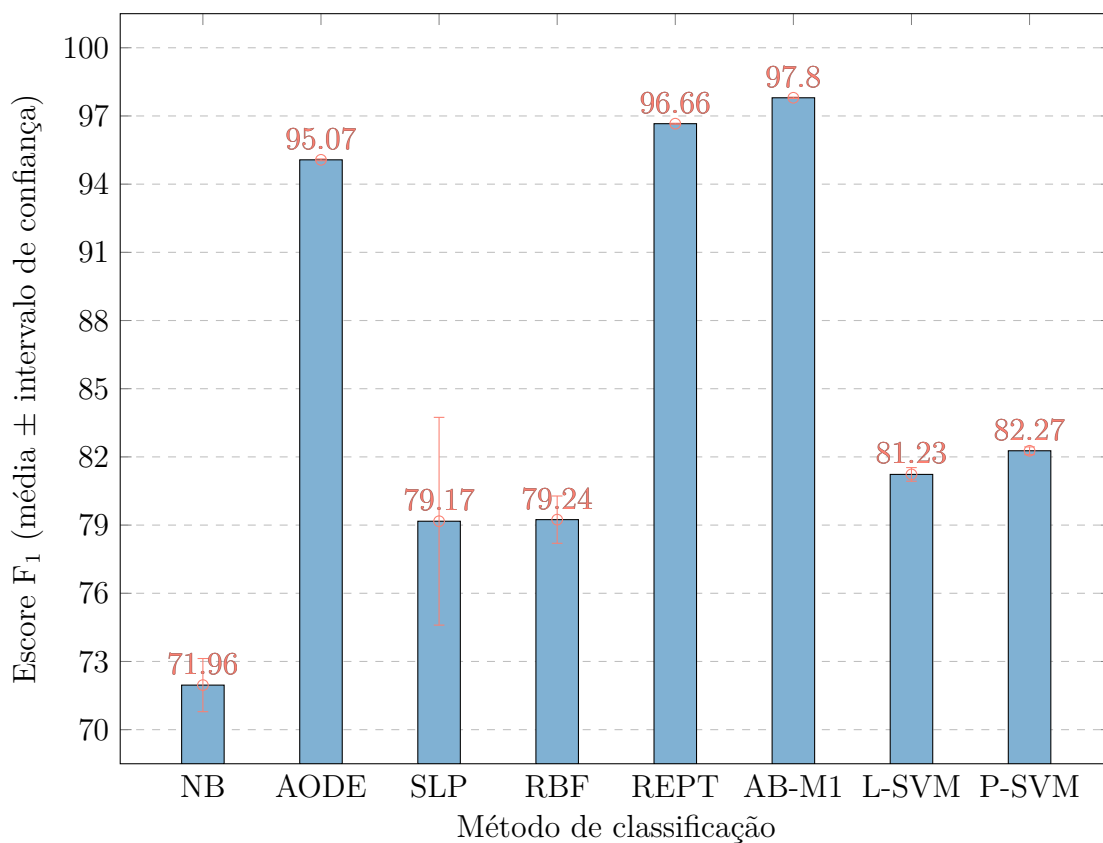
Figura 52: *Ranking* dos classificadores na base Unifei 2017.

Tabela 16: *Ranking* dos classificadores na base Unifei 2018.

FS	NF	NF'	Método	Escore F_1	Tempo de treinamento	Tempo de classificação
FD	8	5	AB-M1	$97,80 \pm 0,03$	$00:00:39.308 \pm 00:00:01.229$	$00:00:02.491 \pm 00:00:00.066$
FD	8	5	REPT	$96,66 \pm 0,02$	$00:00:04.153 \pm 00:00:00.014$	$00:00:00.293 \pm 00:00:00.003$
FD	8	5	AODE	$95,07 \pm 0,04$	$00:00:02.406 \pm 00:00:00.028$	$00:00:00.747 \pm 00:00:00.040$
MI	1024	155	P-SVM	$82,27 \pm 0,17$	$00:58:43.293 \pm 00:12:58.240$	$00:00:25.196 \pm 00:00:00.013$
MI	1024	155	L-SVM	$81,23 \pm 0,30$	$00:00:14.037 \pm 00:00:00.911$	$00:00:06.054 \pm 00:00:00.018$
MI	32	7	RBF	$79,24 \pm 1,04$	$00:01:46.418 \pm 00:00:03.602$	$00:00:07.948 \pm 00:00:00.026$
MI	512	118	SLP	$79,17 \pm 4,57$	$00:04:52.090 \pm 00:00:00.531$	$00:00:00.540 \pm 00:00:00.012$
FD	512	79	NB	$71,96 \pm 1,17$	$00:00:04.986 \pm 00:00:00.061$	$00:00:10.410 \pm 00:00:00.082$

Figura 53: *Ranking* dos classificadores na base Unifei 2018.

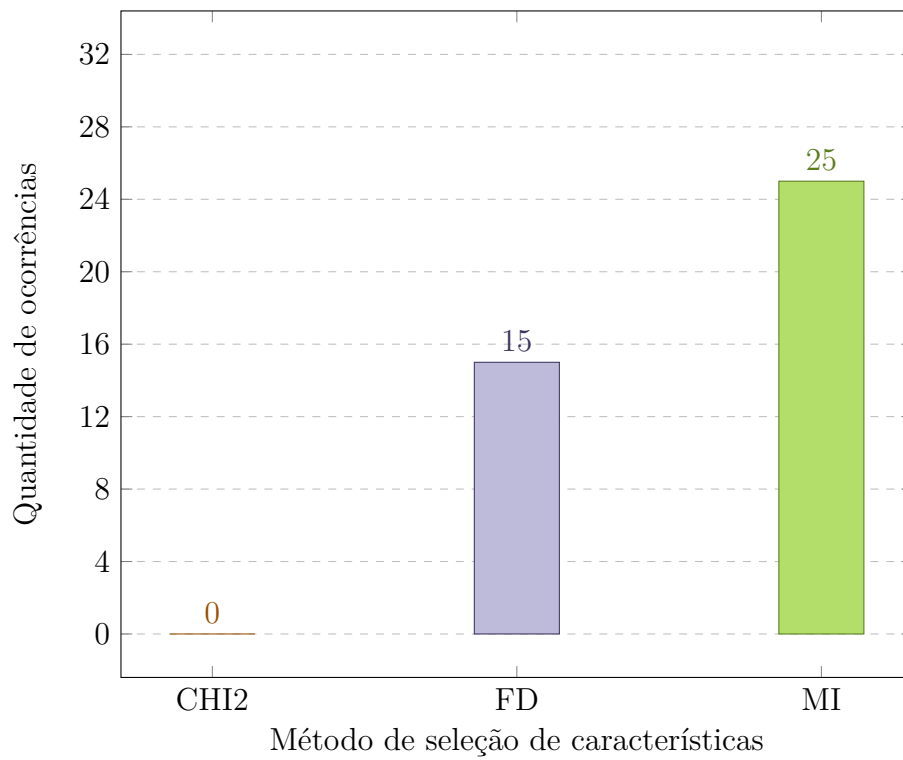


Figura 54: Distribuição dos métodos de seleção de características no *ranking*.

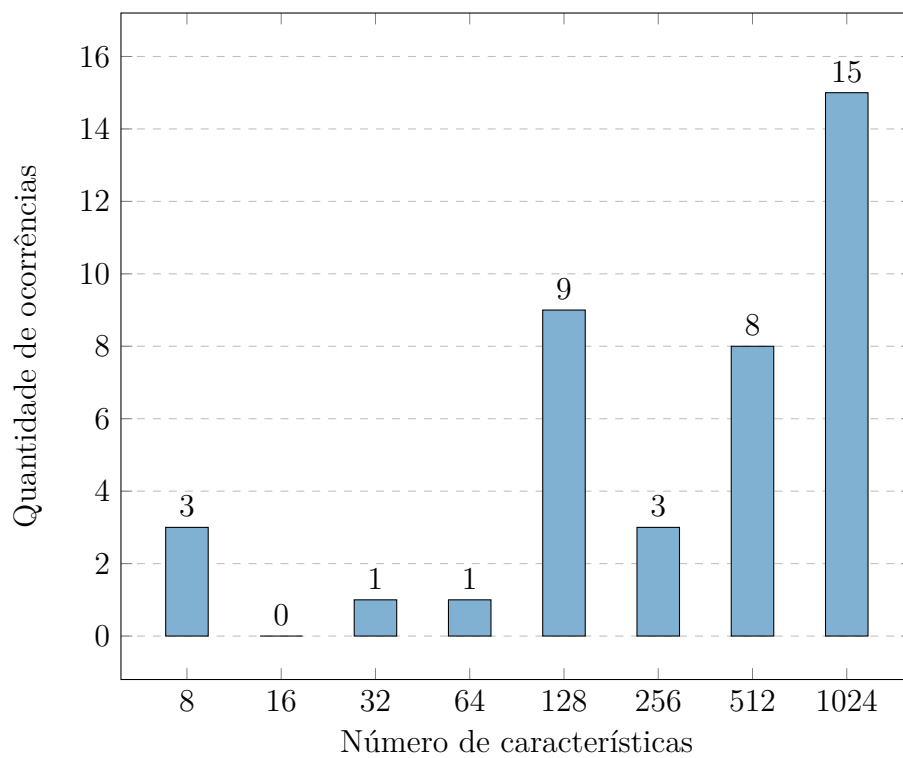


Figura 55: Distribuição das quantidades de características no *ranking*.

Diante destes resultados e levando-se em conta empates técnicos devido à sobreposição de intervalos de confiança, podem ser tiradas as seguintes conclusões:

- Escore F_1
 - O classificador baseado em árvore de decisão e *boosting* (AB-M1) ocupou o primeiro lugar em todos os conjuntos de dados. Portanto, superou a abordagem sem *boosting* (REPT) em todos os cenários, que ocupou os quarto e quinto lugares nas bases menores (Ling Spam e Spam Assassin), e terceiro lugar nas demais (TREC, Unifei 2017 e 2018).
 - O classificador probabilístico não-tão-ingênuo (AODE) superou o ingênuo (NB) em todos os conjuntos de dados. Nas bases menores (Ling Spam, Spam Assassin) e na média (TREC), o AODE ficou em segundo e quinto lugar, respectivamente, e nas demais (Unifei 2017 e 2018), ocupou sempre o terceiro lugar. Em contrapartida, o NB ficou em último colocado em todos os conjuntos de dados.
 - O classificador neural de camada única (SLP) superou a rede de função de base radial (RBF) em todas as bases com exceção da Unifei 2018, ocupando o terceiro lugar nas bases Ling Spam e Spam Assassin, sexto lugar nas bases TREC e Unifei 2017 e sétimo lugar na Unifei 2018. Em contrapartida, o RBF ocupou o sétimo lugar em todas as bases, com exceção da Unifei 2018, onde ocupou a sexta posição.
 - Em todos os conjuntos de dados, a máquina de vetores de suporte de *kernel* polinomial (P-SVM) apresentou resultado superior à de *kernel* linear (L-SVM). A versão polinomial ficou em quinto e segundo lugar nas bases Ling Spam e Spam Assassin, respectivamente, e em quarto lugar nas bases TREC, Unifei 2017 e 2018. Já a versão linear ocupou a sexta posição na base Ling Spam, quarta posição na Spam Assassin e quinto lugar nas demais.

- Tempo de treinamento
 - Nos conjuntos de dados menores (Ling Spam e Spam Assassin), o método RBF apresentou o maior tempo de treinamento. Na base média (TREC), essa posição foi ocupada pela árvore com *boosting* (AB-M1) e, nas bases maiores (Unifei 2017 e 2018), pela SVM de núcleo polinomial.
 - O método NB apresentou o menor tempo de treinamento em todos os conjuntos, com exceção do Unifei 2018, onde o treinamento mais rápido foi da sua contraparte não-tão-ingênua (AODE).
- Tempo de classificação
 - O método RBF apresentou o maior tempo de classificação nos conjuntos Ling Spam, Spam Assassin e Unifei 2017. Nas bases TREC e Unifei 2018, os métodos AODE e P-SVM foram os mais demorados, respectivamente.
 - Em todos os conjuntos de dados, a árvore de decisão sem *boosting* (REPT) apresentou o menor tempo de classificação.
- Seleção de características e redução de dimensionalidade
 - O método de seleção de características MI foi o que teve maior percentual de ocorrências no *ranking* (62,5%), seguido pelo FD (37,5%). Não houve registros nesse sentido para método CHI2.
 - A dimensionalidade original de 1024 características foi a que teve maior percentual de ocorrências no *ranking* (37,5%), seguida por 128 (22,5%) e 512 (20%). As demais quantidades somam 20% das ocorrências, e não houve ocorrências para a configuração com 16 características no *ranking*.

6 Conclusão

6.1 Discussão dos resultados e considerações finais

Nos dias de hoje, com o advento das redes de dados móveis, um número cada vez maior de pessoas possui acesso à Internet de qualquer lugar. Seja no trabalho, em casa ou até mesmo no trânsito, é possível receber e enviar *e-mails* a qualquer momento. Isto implica em um aumento na quantidade de usuários de *e-mail* e, conseqüentemente, na quantidade de trocas de mensagens. Infelizmente, não são somente a quantidade de usuários e a facilidade de acesso que têm aumentado. Desde 2005, pelo menos metade do tráfego de *e-mails* é composto por *spam*.

Qualquer *e-mail* não-solicitado é considerado *spam*. Geralmente, este tipo de mensagem é enviada a um grande número de pessoas para promover um produto ou serviço. Isso inclui anúncios, esquemas de pirâmide, doações, correntes de mensagens, *e-mail* político, assessoria no mercado de ações e avisos únicos.

Este tráfego massivo de mensagens indesejadas provoca conseqüências negativas, tais como o uso excessivo da largura de banda, impacto ambiental, perda de produtividade, exposição de conteúdo impróprio a públicos sensíveis, prejuízos financeiros e legais e redução da efetividade de propagandas legítimas.

Uma técnica largamente empregada por servidores de *e-mail*, a fim de lidar com o problema de *spam*, é através de *listas negras*. Uma lista negra de *e-mail* é um banco de dados em tempo real que determina se um endereço IP está enviando *e-mails* que podem ser considerados *spam*. Este tipo de sistema apresenta diversas

desvantagens, tais como código fechado, parcialidade questionável, limitações devido aos limites dos protocolos IPv4/IPv6 e lentidão na classificação.

Uma outra abordagem ao problema de classificação de *e-mails* consiste em usar os dados contidos na própria mensagem (isto é, remetente, destinatário, assunto, cabeçalho, corpo, anexos etc.) e um modelo de aprendizado de máquina a fim de construir um sistema de decisão automático, inteligente e imparcial. Estes devem ser capazes de identificar padrões e de se ajustarem às características comuns a ambos os tipos de *e-mail*, resultando em classificações cada vez mais precisas, com o decorrer das épocas de treinamento.

Portanto, o objetivo geral deste trabalho é aplicar modelos de aprendizado de máquina à classificação de *e-mails*, a fim de obter modelos utilizáveis em sistemas anti-*spam* do mundo real. Para isto, foi necessário o estudo de métodos de seleção de características e de redução de dimensionalidade, bem como uma revisão teórica de diversos algoritmos usados na categorização de texto. Em seguida, foi elaborado um *framework* de *software* científico genérico, capaz de treinar diversos modelos de aprendizado de máquina e salvá-los para posterior reuso.

Após extensa simulação sobre três conjuntos de dados públicos (Ling Spam, Spam Assassin e TREC) e dois privados (Unifei 2017 e Unifei 2018), os resultados experimentais de cada combinação de modelo de aprendizado de máquina e método de seleção de características foram coletados, consolidados e avaliados em termos de métricas bem difundidas na comunidade científica, tais como precisão, revocação e escore F_1 , dados em termos de média e intervalo de confiança entre as repetições.

Em termos de seleção de características e padrões zerados, constatou-se que os métodos CHI2 e FD foram os que apresentaram maior e menor quantidade de padrões zerados, respectivamente. Além disso, em todos os conjuntos de dados e para todos os métodos de seleção de características, observou-se que a quantidade de padrões zerados foi inversamente proporcional ao número de características das amostras. Por fim, concluiu-se que a grande maioria dos resultados ruins foi obtida com o método CHI2 e a maioria dos resultados bons foi obtida com o método MI.

Em termos de redução de dimensionalidade, os métodos CHI2 e FD foram o que tiveram as maiores quantidades de características descartadas nas bases menores (Ling Spam e Spam Assassin) e maiores (TREC, Unifei 2017 e Unifei 2018), respectivamente. Além disso, observou-se que a redução de até 512 características foi crescente em praticamente todos os casos, sofrendo anomalias especialmente quando a dimensionalidade original era de 1024 características.

Em termos de desempenho na classificação, um método baseado em árvores de decisão “impulsionadas” (AB-M1) mostrou excelentes taxas de acerto em todos os conjuntos de dados. Além disso, foi mostrado que um classificador baseado em redes neurais simples (SLP) foi capaz que apresentar resultados superiores a um modelo mais elaborado (RBF). Neste sentido, também foi possível observar que a suposição de independência em métodos probabilísticos (NB) resultou em desempenho inferior à suposição de dependência restrita (AODE). Finalmente, notou-se que a máquina de vetores de suporte de núcleo polinomial (P-SVM) apresentou resultado superior à de núcleo linear (L-SVM) em todos os casos.

Em termos dos tempos de processamento, enquanto os métodos probabilísticos (ingênuo — NB — e de dependência restrita — AODE) foram treinados mais rapidamente em todas as bases, o modelo de árvore de decisão “impulsionada” (AB-M1) e a máquina de vetores de núcleo polinomial apresentaram os maiores tempos de treinamento. Para classificar as amostras, entretanto, o método que se mostrou mais rápido em todos os casos foi a árvore de decisão simples (REPT). Finalmente, os que mais demoraram para efetuar a classificação variaram, sendo a rede neural de função de base radial (RBF) a mais frequente.

Os resultados experimentais confirmam a hipótese de que um modelo de aprendizado de máquina (especialmente aqueles compostos por árvores de decisão “impulsionadas”), aliado ao método de seleção de características baseado em informação mútua e à técnica de redução de dimensionalidade baseada em correlação pode ser aplicado com sucesso à categorização automática de *e-mails*, protegendo os usuários de conteúdo impróprio e/ou malicioso com alto grau de confiabilidade.

6.2 Sugestões para novos trabalhos

Tanto na etapa de pesquisas e revisões — bibliográficas e de fundamentação teórica — preliminares, quanto durante escrita deste documento, foram identificados tópicos e pontos de possíveis melhora e/ou extensão de pesquisa. Os mais notáveis, identificados como sugestões para novos trabalhos, foram discriminados na listagem a seguir, bem como referências relevantes a cada item.

- Uso de outros métodos de seleção de características. Exemplos: *Gini Index* (GINI, 1912), BNS (do inglês *Bi-Normal Separation*), IG e *Odds Ratio* (FORMAN; GEORGE, 2003).
- Exploração do espaço e otimização de parâmetros de configuração de cada modelo de aprendizado de máquina. (CHAPELLE et al., 2002)
- Uso de técnicas de neuroevolução. Exemplos: GNARL (do inglês *GeNeralized Acquisition of Recurrent Links*) (ANGELINE; SAUNDERS; POLLACK, 1994), NEAT (do inglês *NeuroEvolution of Augmenting Topologies*) (STANLEY; MIIKKULAINEN, 2002) e EANT2 (do inglês *Evolutionary Acquisition of Neural Topologies, Version 2*) (SIEBEL; SOMMER, 2007).
- Uso de bibliotecas de aprendizado de máquina e/ou plataformas computacionais com suporte a paralelismo. Exemplo: Nvidia CUDA (do inglês *Compute Unified Device Architecture*) (KIRK, 2007), que utiliza GPUs (do inglês *Graphics Processing Unit*) dedicadas para acelerar o processamento de problemas com grande capacidade de paralelização.
- Uso de técnicas de *deep learning* (DENG; YU, 2014).

7 Anexos

O código do *framework* desenvolvido¹ e os conjuntos de dados utilizados nos experimentos² foram disponibilizados sob a licença GNU GENERAL PUBLIC LICENSE versão 3, de 29 de junho de 2007³.

This program is free software: you can redistribute it and/or modify it under the terms of the GNU General Public License as published by the Free Software Foundation, either version 3 of the License, or (at your option) any later version.

This program is distributed in the hope that it will be useful, but WITHOUT ANY WARRANTY; without even the implied warranty of MERCHANTABILITY or FITNESS FOR A PARTICULAR PURPOSE. See the GNU General Public License for more details.

You should have received a copy of the GNU General Public License along with this program. If not, see <<http://www.gnu.org/licenses/>>.

¹Disponível em <<https://github.com/marcelovca90/anti-spam-weka-cli>>

²Disponível em <<https://github.com/marcelovca90/anti-spam-weka-data>>

³Disponível em <<https://www.gnu.org/licenses/gpl-3.0.html>>

8 Apêndices

A tabela 17 apresenta os parâmetros de configuração de cada modelo de aprendizado de máquina utilizado no WEKA.

As tabelas 18, 19, 20, 21, 22, 23, 24, e 25 mostram os resultados completos das execuções dos classificadores NB, AODE, SLP, RBF, REPT, AB-M1, L-SVM e P-SVM, respectivamente.

As tabelas 26 e 27 mostram os detalhes de *hardware* (processador, memória e sistema operacional) e de *software* (Java, WEKA e bibliotecas/dependências), respectivamente, utilizados na execução dos experimentos.

Tabela 17: Parâmetros de configuração usados no WEKA.

Método	Parâmetro	Descrição	Valor ^{a b c}
NB	N/D	N/D	N/D
AODE	N/D	N/D	N/D
SLP	-I	quantidade de iterações de treinamento	500
	-L	taxa de aprendizado	0.1
	-B	valor do limiar de ativação	1.0
	-R	semente do gerador de números aleatórios	1
	-M	função de taxa de aprendizado	<i>LD</i>
RBF	-B	número de agrupamentos a serem gerados	$(\log_2 NF)^2$
	-S	semente do gerador de números aleatórios	1
	-R	valor de pico (<i>ridge</i>) para a regressão	10^{-8}
	-M	limite de iterações da regressão logística	-1
	-W	desvio padrão mínimo para os agrupamentos	0.1
REPT	-M	peso total mínimo de amostras em uma folha	2
	-V	menor proporção da variância de todos os dados que deve estar presente em um nó para realizar divisão em árvores de regressão	0.001
	-N	quantidade de frações dos dados de treinamento usadas na poda da árvore de decisão	3
	-S	semente do gerador de números aleatórios	1
	-L	profundidade máxima da árvore	-1
AB-M1	-P	limiar de peso para poda	100
	-S	semente do gerador de números aleatórios	1
	-I	quantidade de iterações de <i>boosting</i>	10
	-W	classificador sobre o qual será feito o <i>boosting</i>	<i>REPT</i>
L-SVM	-s	tipo do algoritmo solucionador	2
	-c	parâmetro de regularização	128.0
P-SVM	-s	tipo do algoritmo solucionador	2
	-c	parâmetro de regularização	128.0
	-g	gamma da função de mapeamento	1
	-r	coef0 da função de mapeamento	1

^a*NF* é a quantidade original de características (ou seja, antes da redução de dimensionalidade)^b*LD* é uma função que decrementa a taxa de aprendizado com o decorrer do treinamento^c*REPT* é o classificador *Reduced-Error Pruning Tree*, explicado com detalhes em 2.6.3.1

Tabela 26: Detalhes do *hardware* utilizado nos experimentos.

Processador	Marca	Intel
	Modelo	Intel(R) Core(TM) i5-4570
	Frequência mínima	3,20 GHz
	Frequência máxima	3,60 GHz
	Cache	6 MB
	Soquete	FCLGA1150
Memória	Tamanho	32 GB
	Tipo	DDR3-1333
Sistema Operacional	Versão	18.3
	Codinome	sylvia
	Edição	Cinnamon 64-bit
	Descrição	Linux Mint 18.3 Sylvia

Tabela 27: Detalhes do *software* utilizado nos experimentos.

Java	openjdk version "1.8.0_162"
	OpenJDK Runtime Environment (build 1.8.0_162-8u162-b12-0ubuntu0.16.04.2-b12)
	OpenJDK 64-Bit Server VM (build 25.162-b12, mixed mode)
libatlas3-base	3.10.2-9
libopenblas-base	0.2.18-ubuntu1
WEKA	Stable 3.8.1
LIBLINEAR	2.20
LIBLINEAR-poly2	2.01

Referências

- AGGARWAL, C. C.; ZHAI, C. A survey of text classification algorithms. In: *Mining text data*. Boston, MA: Springer, 2012. p. 163–222.
- AIZERMAN, M. A.; BRAVERMAN, E. M.; ROZONOÉR, L. I. Theoretical foundations of the potential function method in pattern recognition learning. *Automation and Remote Control*, v. 25, p. 821–837, 1964.
- ALI, A. Comparison study between IPV4 & IPV6. *International Journal of Computer Science Issues*, Citeseer, v. 9, n. 3, p. 314–317, 2012. ISSN 16940784.
- ANDROUTSOPOULOS, I. et al. An evaluation of naive bayesian anti-spam filtering. *arXiv preprint cs/0006013*, Barcelona, p. 9–17, 2000.
- ANDROUTSOPOULOS, I. et al. *Ling-Spam Corpus*. 2000. Disponível em: https://labs-repos.iit.demokritos.gr/skel/i-config/downloads/lingspam{_}public.tar.
- ANDROUTSOPOULOS, I.; PALIOURAS, G.; MICHELAKIS, E. *PU1, PU2, PU3, and PUA corpora*. "DEMOKRITOS", National Center for Scientific Research, 2003. Disponível em: <https://labs-repos.iit.demokritos.gr/skel/i-config/downloads/>.
- ANGELINE, P. J.; SAUNDERS, G. M.; POLLACK, J. B. An evolutionary algorithm that constructs recurrent neural networks. *IEEE transactions on neural networks / a publication of the IEEE Neural Networks Council*, v. 5, n. 1, p. 54–65, 1994. ISSN 1045-9227.
- ARAGÃO, M. V. et al. Factorial design analysis applied to the performance of SMS anti-spam filtering systems. *Expert Systems with Applications*, v. 64, p. 589–604, 2016. ISSN 09574174.
- BARIGOU, F.; BELDJILALI, B.; ATMANI, B. Using cellular automata for improving knn based spam filtering. *International Arab Journal of Information Technology (IAJIT)*, v. 11, n. 4, p. 345–353, 2014.
- BEN-HUR, A. et al. Support vector clustering. *Journal of machine learning research*, v. 2, n. Dec, p. 125–137, 2001.

- BERGER, H.; MERKL, D. A comparison of support vector machines and self-organizing maps for e-mail categorization. In: *Proceedings of AusDM 2005, the 4th Australasian Data Mining Conference*. Sydney, Australia: [s.n.], 2005. p. 189–204. ISBN 1863657169.
- BISHOP, C. M. *Pattern Recognition and Machine Learning*. 1. ed. Secaucus, NJ, USA: Springer-Verlag New York, Inc., 2006. ISBN 978-0-387-31073-2.
- BOSER, B. E.; GUYON, I. M.; VAPNIK, V. N. A training algorithm for optimal margin classifiers. In: ACM. *Proceedings of the 5th ACM Workshop on Computational Learning Theory*. Pittsburgh, PA, USA, 1992. p. 144–152.
- BREIMAN, L. et al. *Classification and regression trees*. 1. ed. Boca Raton, FL, USA: Chapman & Hall/CRC, 1984. 358 p. ISSN 19424787. ISBN 978-0412048418.
- BREUNIG, M. M. et al. LOF: Identifying Density-Based Local Outliers. In: ACM. *Proceedings of the 2000 ACM SIGMOD International Conference on Management of Data*. New York, NY, USA, 2000. v. 29, n. 2, p. 1–12. ISBN 1581132182. ISSN 01635808.
- BROOMHEAD, D. S.; LOWE, D. *Radial basis functions, multi-variable functional interpolation and adaptive networks*. Malvern, Worcs, UK, 1988.
- CAMPS-VALLS, G. et al. Composite kernels for hyperspectral image classification. *IEEE Geoscience and Remote Sensing Letters*, IEEE, v. 3, n. 1, p. 93–97, 2006. ISSN 1545598X.
- CARPINTEIRO, O. A. S. et al. A neural model in anti-spam systems. In: SPRINGER. *International Conference on Artificial Neural Networks*. Berlin, Heidelberg, Germany, 2006. p. 847–855.
- CHANG, C.-C.; LIN, C.-J. LIBSVM: A Library for Support Vector Machines. *ACM Transactions on Intelligent Systems and Technology*, Acm, v. 2, n. 3, p. 1–27, 2011. ISSN 21576904.
- CHANG, Y.-W. et al. Training and Testing Low-degree Polynomial Data Mappings via Linear SVM. *Journal of Machine Learning Research*, v. 11, n. Apr, p. 1471–1490, 2010. ISSN 1532-4435.
- CHAPELLE, O. et al. Choosing multiple parameters for support vector machines. *Machine learning*, Springer, v. 46, n. 1, p. 131–159, 2002.
- COHEN, W. W. Enron Email Dataset. 2008. Disponível em: <<https://www.cs.cmu.edu/~wcohen/>>.

- CORMACK, G. V.; LYNAM, T. R. *2005 TREC Public Spam Corpus*. 2005. Disponível em: <<https://plg.uwaterloo.ca/~gvcormac/treccorpu>>.
- CORMACK, G. V.; LYNAM, T. R. *2006 TREC Public Spam Corpora*. 2006. Disponível em: <<https://plg.uwaterloo.ca/~gvcormac/treccorpus0>>.
- CORMACK, G. V.; LYNAM, T. R. *2007 TREC Public Spam Corpus*. 2007. Disponível em: <<http://plg.uwaterloo.ca/~gvcormac/treccorpus0>>.
- CORMACK, G. V.; OTHERS. Email spam filtering: A systematic review. *Foundations and Trends in Information Retrieval*, Now Publishers, Inc., v. 1, n. 4, p. 335–455, 2008. ISSN 1554-0669.
- CORTES, C.; VAPNIK, V. Support-vector networks. *Machine learning*, Springer, v. 20, n. 3, p. 273–297, 1995. ISSN 15730565.
- COVER, T. M.; THOMAS, J. A. *Elements of information theory*. 2. ed. New York: John Wiley & Sons, 1991. 563 p.
- COVÕES, T. F.; HRUSCHKA, E. R. Towards improving cluster-based feature selection with a simplified silhouette filter. *Information Sciences*, Elsevier, v. 181, n. 18, p. 3766–3782, 2011. ISSN 00200255.
- CYGANIAK, R.; WOOD, D.; LANTHALER, M. RDF 1.1 Concepts and Abstract Syntax. *W3C Recommendation 25 February 2014*, World Wide Web Consortium, n. February, p. 263–270, 2014. ISSN 13514180.
- DEB, K. et al. A Fast and Elitist Multiobjective Genetic Algorithm: NSGA-II. *IEEE Transactions on Evolutionary Computing*, IEEE, v. 6, n. 2, p. 182–197, 2002. ISSN 1089778X.
- DEERWESTER, S. et al. Indexing by Latent Semantic Analysis. *Journal of the American Society for Information Science*, American Documentation Institute, v. 41, n. 6, p. 391–407, 1990. ISSN 10974571.
- DENG, L.; YU, D. Deep Learning: Methods and Applications. *Foundations and Trends in Signal Processing*, Now Publishers, Inc., v. 7, n. 3-4, p. 197–387, 2014. ISSN 1932-8346.
- DONOHO, D. L. High-Dimensional Data Analysis: The Curses and Blessings of Dimensionality. *American Math. Society Lecture-Math Challenges of the 21st Century*, Citeseer, v. 1, p. 1–33, 2000.
- DREDZE, M.; GEVARYAHU, R.; ELIAS-BACHRACH, A. Image Spam Dataset. 2007. Disponível em: <http://www.cs.jhu.edu/~mdredze/datasets/image{_}>.

- DREDZE, M.; GEVARYAHU, R.; ELIAS-BACHRACH, A. Learning Fast Classifiers for Image Spam. In: *Conference on Email and Anti-Spam (CEAS) 2007*. Berlin, Germany: [s.n.], 2007.
- DRUCKER, H. et al. Support vector machines for spam categorization. *IEEE Transactions on Neural networks*, IEEE, v. 10, n. 5, p. 1048–1054, 1999.
- ESPOSITO, F.; MALERBA, D.; SEMERARO, G. Decision Tree Pruning as a Search in the State Space. In: SPRINGER. *Machine Learning: ECML-93, European Conference on Machine Learning, Proceedings*. Berlin, Heidelberg, Germany, 1993. v. 667, p. 165–184. ISBN 9783540566021. ISSN 16113349.
- FAN, R.-E. et al. LIBLINEAR: A Library for Large Linear Classification. *Journal of Machine Learning Research*, v. 9, n. Aug, p. 1871–1874, 2008. ISSN 15324435.
- FERREIRA, I. C. Estudo e Desenvolvimento de um Sistema Anti-Spam. 2018.
- FORMAN; GEORGE. An extensive empirical study of feature selection metrics for text classification. *The Journal of Machine Learning Research*, v. 3, n. Mar, p. 1289–1305, 2003. ISSN 1532-4435.
- FOUNDATION, T. A. S. *Apache SpamAssassin Project*. 2015. Disponível em: <<http://spamassassin.apache.org/>>.
- FREUND, Y. Boosting a weak learning algorithm by majority. *Information and computation*, Elsevier, v. 121, n. 2, p. 256–285, 1995.
- FREUND, Y.; MASON, L. The alternating decision tree learning algorithm. In: *Proceeding of the Sixteenth International Conference on Machine Learning*. San Francisco, CA, USA: [s.n.], 1999. v. 99, p. 10. ISBN 9780874216561. ISSN 14602431.
- FREUND, Y.; SCHAPIRE, R. E. Experiments with a New Boosting Algorithm. In: *Proceedings of the 13th International Conference on Machine Learning*. Bari, Italy: [s.n.], 1996. v. 96, p. 148–156.
- FRIEDMAN, J. H. Greedy function approximation: A gradient boosting machine. *The Annals of Statistics*, JSTOR, v. 29, n. 5, p. 1189–1232, 2001.
- F.R.S., K. P. On lines and planes of closest fit to systems of points in space. *The London, Edinburgh, and Dublin Philosophical Magazine and Journal of Science*, Taylor & Francis, v. 2, n. 11, p. 559–572, 1901.
- GINI, C. Variabilità e mutabilità. *Reprinted in Memorie di metodologica statistica (Ed. Pizetti E, Salvemini, T)*. Rome: Libreria Eredi Virgilio Veschi, 1912.

- GUYON, I. et al. Gene selection for cancer classification using support vector machines. *Machine Learning*, Springer, v. 46, n. 1-3, p. 389–422, 2002. ISSN 08856125.
- HALL, M. et al. The WEKA data mining software: an update. *ACM SIGKDD explorations newsletter*, ACM, v. 11, n. 1, p. 10–18, 2009.
- HALL, M. A. Correlation-based feature selection of discrete and numeric class machine learning. University of Waikato, Department of Computer Science, 2000.
- HAYKIN, S. S. *Redes neurais: princípios e prática*. 2. ed. Porto Alegre, RS, Brazil: Bookman Companhia Editora, 2001. 900 p. ISBN 8573077182.
- HEBB, D. O. *The Organization of Behavior; A Neuropsychological Theory*. New York, New York, USA: Wiley & Sons, 1949. 378 p.
- HINTON, G. E.; ROWEIS, S. T. Stochastic Neighbor Embedding. In: *Advances in Neural Information Processing Systems 15*. Vancouver, Canada: [s.n.], 2002. p. 833–840. ISBN 0262025507. ISSN 10495258.
- HOFMANN, T. Probabilistic Latent Semantic Indexing. In: ACM. *Sigir'99*. New York, NY, USA, 1999. v. 51, n. 2, p. 50–57. ISBN 1581130961. ISSN 15206882.
- HOPKINS, M. et al. *Spam base dataset*. 1999. Disponível em: <<https://archive.ics.uci.edu/ml/datasets/spambase>>.
- HUGHES, G. E. On the Mean Accuracy of Statistical Pattern Recognition Accuracy. *IEEE Transactions on Information Theory*, IEEE, IT-14, n. 1, p. 55–63, 1968.
- IDRIS, I. et al. A combined negative selection algorithm–particle swarm optimization for an email spam detection system. *Engineering Applications of Artificial Intelligence*, Elsevier, v. 39, p. 33–44, 2015.
- J, K.; C, E. R. Particle swarm optimization. In: *Proceedings of IEEE International Conference on Neural Networks*. Perth, WA, Australia: Springer, 1995. p. 1942–1948. ISBN 9781612840529.
- JIMENEZ, F. et al. An evolutionary algorithm for constrained multi-objective optimization. In: IEEE. *Proceedings of the 2002 Congress on Evolutionary Computation*. Honolulu, HI, USA, 2002. v. 2, p. 1133–1138.
- JUNG, J.; SIT, E. An empirical study of spam traffic and the use of DNS black lists. In: ACM. *Proceedings of the 4th ACM SIGCOMM Conference on Internet Measurement*. Taormina, Sicily, Italy, 2004. p. 370–375. ISBN 1-58113-821-0.

KAYA, Y.; ERTUĞRUL, Ö. F. A novel approach for spam email detection based on shifted binary patterns. *Security and Communication Networks*, Wiley Online Library, v. 9, n. 10, p. 1216–1225, 2016. ISSN 19390122.

KIRK, D. NVIDIA CUDA Software and GPU Parallel Computing Architecture. In: *Proceedings of the 6th International Symposium on Memory Management*. Montreal, QC, Canada: [s.n.], 2007. v. 7, p. 103–104. ISBN 9781595938930.

KOHONEN, T. Learning Vector Quantization. In: *Self-Organizing Maps*. Berlin, Heidelberg, Germany: Springer, 1995. p. 175–189. ISBN 3642976123.

KUMARESAN, T.; SARAVANAKUMAR, S.; BALAMURUGAN, R. Visual and textual features based email spam classification using S-Cuckoo search and hybrid kernel support vector machine. *Cluster Computing*, Springer, p. 1–14, 2017.

LIN, J. Divergence measures based on the Shannon entropy. *IEEE Transactions on Information theory*, IEEE, v. 37, n. 1, p. 145–151, 1991.

LIU, H.; SETIONO, R. A probabilistic approach to feature selection - a filter solution. In: UNIVERSITA' DI BARI, UNIVERSITA' DI TORINO. *Proc 13th International Conference on Machine Learning*. Bari, Italy, 1996. v. 96, p. 319–327.

LODHI, H. et al. Text Classification using String Kernels. *Journal of Machine Learning Research*, v. 2, n. Feb, p. 419–444, 2002. ISSN 0003-6951.

LÓPEZ-IBÁÑEZ, M.; PAQUETE, L.; STÜTZLE, T. Exploratory analysis of stochastic local search algorithms in biobjective optimization. In: *Experimental Methods for the Analysis of Optimization Algorithms*. Berlin, Heidelberg, Germany: Springer, 2010. p. 209–222. ISBN 978-3-642-02537-2.

MAATEN, L. van der; HINTON, G. Visualizing data using t-SNE. *Journal of machine learning research*, v. 9, n. Nov, p. 2579–2605, 2008.

MACQUEEN, J. Some methods for classification and analysis of multivariate observations. In: STATISTICAL LABORATORY OF THE UNIVERSITY OF CALIFORNIA, BERKELEY. *Proceedings of the 5th Berkeley Symposium on Math., Stat. and Prob.* Oakland, CA, USA: University of California Press, 1967. v. 1, n. 14, p. 281–296.

McAfee Inc.; ICF International. The Carbon Footprint of Email Spam Report. *McAfee Resource Center*, p. 1–28, 2009.

MCCALLUM, A.; NIGAM, K. A Comparison of Event Models for Naive Bayes Text Classification. In: JUST RESEARCH AND CARNEGIE MELLON

UNIVERSITY SCHOOL OF COMPUTER SCIENCE. *AAAI/ICML-98 Workshop on Learning for Text Categorization*. Madison, WI, USA, 1998. v. 752, n. 1, p. 41–48. ISBN 0897915240. ISSN 0343-6993.

MCCORMICK, C. *AdaBoost Tutorial*. 2013. 1–9 p. Disponible em: <<http://mccormickml.com/2013/12/13/adaboost-tutorial/>>.

MCCULLOCH, W. S.; PITTS, W. A logical calculus of the ideas immanent in nervous activity. *The bulletin of mathematical biophysics*, Springer, v. 5, n. 4, p. 115–133, 1943. ISSN 0007-4985.

MCGONAGLE, J. *Naive Bayes Classifier*. 2018. 1–4 p. Disponible em: <<https://brilliant.org/wiki/naive-bayes-classifier/>>.

MÉNDEZ, J. R. et al. Grindstone4Spam: An optimization toolkit for boosting e-mail classification. *Journal of Systems and Software*, Elsevier, v. 85, n. 12, p. 2909–2920, 2012. ISSN 01641212.

MEYER, T. A.; WHATELEY, B. SpamBayes: Effective open-source, Bayesian based, email classification system. In: *Conference on Email and Anti-Spam (CEAS) 2004*. Mountain View, CA, USA: [s.n.], 2004.

MINGERS, J. An empirical comparison of pruning methods for decision tree induction. *Machine learning*, Springer, v. 4, n. 2, p. 227–243, 1989.

MINGERS, J. An empirical comparison of selection measures for decision-tree induction. *Machine learning*, Springer, v. 3, n. 4, p. 319–342, 1989.

NEYMAN, J. Outline of a Theory of Statistical Estimation Based on the Classical Theory of Probability. *Philosophical Transactions of the Royal Society A: Mathematical, Physical and Engineering Sciences*, The Royal Society, v. 236, n. 767, p. 333–380, 1937.

PAWLAK, Z. *Rough Sets: Theoretical Aspects of Reasoning About Data*. Dordrecht, Netherlands: Springer Science & Business Media, 1991. v. 9.

PÉREZ-DÍAZ, N. et al. Rough sets for spam filtering: Selecting appropriate decision rules for boundary e-mail classification. *Applied Soft Computing Journal*, Elsevier, v. 12, n. 11, p. 3671–3682, 2012. ISSN 15684946.

POGGIO, T.; GIROSI, F. *A theory of networks for approximation and learning*. Cambridge, MA, USA, 1989. 1–65 p.

POWELL, M. J. D. Radial basis functions for multivariable interpolation: a review. *Algorithms for approximation*, Clarendon Press, p. 143–167, 1987.

- POWERS, D. Evaluation: From Precision, Recall and F-Measure To Roc, Informedness, Markedness & Correlation. *Journal of Machine Learning Technologies*, v. 2, n. 1, p. 37–63, 2011. ISSN 2229-3981.
- QUINLAN, J. R. Induction of Decision Trees. *Machine Learning*, Springer, v. 1, n. 1, p. 81–106, 1986. ISSN 1573-0565.
- QUINLAN, J. R. Simplifying decision trees. *International journal of man-machine studies*, Elsevier, v. 27, n. 3, p. 221–234, 1987.
- QUINLAN, J. R. *C 4.5: Programs for machine learning*. San Mateo, CA, USA: Morgan Kaufmann Publishers, 1993. 1–302 p. ISBN 1558602380.
- RADICATI, S. *Email statistics report, 2018-2022*. Palo Alto, CA, USA, 2018. 1–3 p. Disponível em: <<https://www.radicati.com/?p=15185>>.
- ROKACH, L.; MAIMON, O. Top-down induction of decision trees classifiers-a survey. *IEEE Transactions on Systems, Man, and Cybernetics, Part C (Applications and Reviews)*, IEEE, v. 35, n. 4, p. 476–487, 2005. ISSN 1094-6977.
- RUMELHART, D. E.; HINTON, G. E.; WILLIAMS, R. J. *Learning internal representations by error propagation*. La Jolla, CA, USA, 1985. 318–362 p.
- RUSSELL, S.; NORVIG, P. *Inteligencia Artificial*. Rio de Janeiro, RJ, Brazil: Elsevier Editora, 2013. 1–1016 p. ISBN 978-8535237016.
- SAHAMI, M. Learning Limited Dependence Bayesian Classifiers. In: *Second International Conference on Knowledge Discovery and Data Mining*. Portland, OR, USA: [s.n.], 1996. v. 96, n. 1, p. 91–94. ISBN 1-57735-004-9.
- SAMUEL, A. L. Some studies in machine learning using the game of checkers. *IBM Journal of research and development*, IBM, v. 3, n. 3, p. 210–229, 1959.
- SCHAPIRE, R. E. The strength of weak learnability. *Machine learning*, Springer, v. 5, n. 2, p. 197–227, 1990.
- SHAMS, R.; MERCER, R. E. Supervised classification of spam emails with natural language stylometry. *Neural Computing and Applications*, Springer, v. 27, n. 8, p. 2315–2331, 2016.
- SIEBEL, N. N. T.; SOMMER, G. Evolutionary reinforcement learning of artificial neural networks. *International Journal of Hybrid Intelligent Systems*, IOS Press, v. 4, n. 3, p. 171–183, 2007. ISSN 1448-5869.

SILVA, I. N.; SPATTI, D. H.; FLAUZINO, R. A. Redes Neurais Artificiais Para Engenharia e Ciências Aplicadas. *São Paulo: Artliber*, p. 31–55, 2010.

SIRISANYALAK, B.; SORNIL, O. An artificial immunity-based spam detection system. In: IEEE. *2007 IEEE Congress on Evolutionary Computation, CEC 2007*. Singapore, 2007. p. 3392–3398. ISBN 1424413400.

SNYDER, J. *Comparing Industry-Leading Anti-Spam Services: Head-to-Head Testing Results*. Tucson, AZ, USA, 2015. 1–6 p.

SPERBER, D. et al. Epistemic vigilance. *Mind and Language*, Wiley Online Library, v. 25, n. 4, p. 359–393, 2010. ISSN 02681064.

SPRUYT, V. *The Curse of Dimensionality in classification*. 2014. 1–16 p. Disponível em: <<http://www.visiondummy.com/2014/04/curse-dimensionality-affect-classification/>>.

STANLEY, K. O.; MIIKKULAINEN, R. Evolving neural networks through augmenting topologies. *Evolutionary computation*, MIT Press, v. 10, n. 2, p. 99–127, 2002. ISSN 1063-6560.

Symantec Corporation. *Internet Security Threat Report 8*. Cupertino, CA, USA, 2005. 1–106 p. Disponível em: <<https://www.symantec.com/content/dam/symantec/docs/security-center/archives/istr-05-sept-en.pdf>>.

Symantec Corporation. *Internet Security Threat Report 17*. Mountain View, CA, USA, 2012. v. 17, 1–52 p. Disponível em: <<https://www.symantec.com/content/dam/symantec/docs/security-center/archives/istr-12-april-volume-17-en.pdf>>.

Symantec Corporation. *Internet Security Threat Report 23*. Mountain View, CA, USA, 2018. v. 23, 1–89 p. Disponível em: <<https://www.symantec.com/security-center/threat-report>>.

TAO, B. *CSDMC2010 SPAM corpus*. International Conference on Neural Information Processing, 2010. Disponível em: <<http://csmining.org/index.php/spam-email-datasets-.html>>.

TENENBAUM, J. B.; SILVA, V. de; LANGFORD, J. C. A Global Geometric Framework for Nonlinear Dimensionality Reduction. *Science*, American Association for the Advancement of Science, v. 290, n. 5500, p. 2319–2323, 2000.

The Apache Software Foundation. *SpamAssassin Public Mail Corpus*. 2005. Disponível em: <<https://spamassassin.apache.org/publiccorpus/>>.

The Spamhaus Project. *The Spamhaus Block List*. 2018. 1–1 p. Disponível em: <<https://www.spamhaus.org/sbl/>>.

The Spamhaus Project. *Understanding DNSBL Filtering*. 2018. 1–2 p. Disponível em: <https://www.spamhaus.org/whitepapers/dnsbl{_}functi>.

VAPNIK, V. N.; LERNER, A. Pattern recognition using generalized portrait method. *Automation and Remote Control*, v. 24, p. 774–780, 1963.

WEBB, G. I.; BOUGHTON, J. R.; WANG, Z. Not so naive Bayes: aggregating one-dependence estimators. *Machine learning*, Springer, v. 58, n. 1, p. 5–24, 2005.

WETTSCHERECK, D.; DIETTERICH, T. Improving the Performance of Radial Basis Function Networks by Learning Center Locations. In: *Advances in Neural Information Processing Systems*. San Francisco, CA, USA: Morgan Kaufmann Publishers, 1992. v. 4, p. 1133–1140.

WHITE, B. W.; ROSENBLATT, F. Principles of Neurodynamics: Perceptrons and the Theory of Brain Mechanisms. *The American Journal of Psychology*, v. 76, n. 4, p. 705, 1963. ISSN 00029556.

WIDROW, B.; HOFF, M. E. *Adaptive switching circuits*. Stanford, CA, USA, 1960.

YANG, C. et al. An Anti-Spam Filter Based on One-Class IB Method in Small Training Sets. *International Arab Journal of Information Technology (IAJIT)*, v. 13, n. 6, p. 677–685, 2016. ISSN 1683-3198.

YANG, J.; ESTIVILL-CASTRO, V.; CHALUP, S. K. Support vector clustering through proximity graph modelling. In: IEEE. *Proceedings of the 9th International Conference on Neural Information Processing*. Singapore, 2002. v. 2, p. 898–903. ISBN 9810475241.

YANG, X.-S.; DEB, S. Cuckoo Search via Levy Flights. In: IEEE. *Nature & Biologically Inspired Computing, 2009. NaBIC 2009. World Congress on*. Coimbatore, India, 2010. p. 210–214. ISBN 978-1-4244-5053-4. ISSN 0941-0643.

YANG, Y.; JOACHIMS, T. Text categorization. *Scholarpedia*, v. 3, n. 5, p. 4242, 2008. ISSN 1941-6016. Disponível em: <http://www.scholarpedia.org/article/Text{_}categorizat>.

YANG, Y.; PEDERSEN, J. A Comparative Study on Feature Selection in Text Categorization. In: *Proceedings of ICML-97, 14th International Conference on Machine Learning*. Lille, France: [s.n.], 1997. v. 97, p. 412–420.

YEVSEYEVA, I. et al. Optimising anti-spam filters with evolutionary algorithms. *Expert Systems with Applications*, Elsevier, v. 40, n. 10, p. 4010–4021, 2013. ISSN 09574174.

YOUN, S. SPONGY (SPam ONtology): email classification using two-level dynamic ontology. *The Scientific World Journal*, Hindawi Publishing Corporation, v. 2014, 2014. ISSN 1537744X.

ZHANG, H. et al. IPGroupRep: A novel reputation based system for anti-spam. *UIC-ATC 2009 - Symposia and Workshops on Ubiquitous, Autonomic and Trusted Computing*, p. 513–518, 2009.

ZITZLER, E.; LAUMANN, M.; THIELE, L. SPEA2: Improving the Strength Pareto Evolutionary Algorithm. *Evolutionary Methods for Design Optimization and Control with Applications to Industrial Problems*, p. 95–100, 2001. ISSN 03772217.