

JACKSON MALLMANN

**PRODUÇÃO DE PROVAS DIGITAIS A PARTIR
DE RASTREAMENTO EM
RELACIONAMENTOS POR *e-MAILS***

Dissertação apresentada ao Programa de Pós-Graduação em Informática da Pontifícia Universidade Católica do Paraná como requisito parcial para obtenção do título de Mestre em Informática.

CURITIBA

2011

JACKSON MALLMANN

**PRODUÇÃO DE PROVAS DIGITAIS A PARTIR
DE RASTREAMENTO EM
RELACIONAMENTOS POR *e-MAILS***

Dissertação apresentada ao Programa de Pós-Graduação em Informática da Pontifícia Universidade Católica do Paraná como requisito parcial para obtenção do título de Mestre em Informática.

Área de Concentração: *Ciência da Computação*

Linha de Pesquisa: *Computação Forense e Biometria*

Orientadora: Prof. Dra. Cinthia O. de A. Freitas

Co-orientador: Prof. Dr. Altair Olivo Santin

CURITIBA

2011

Mallmann, Jackson

Produção de Provas Digitais a partir de Rastreamento em Relacionamento por e-mails. Curitiba, 2011. 116p.

Dissertação – Pontifícia Universidade Católica do Paraná. Programa de Pós-Graduação em Informática.

1. Rastreamento 2. *e-mail* 3. Extração de Características 4. Agrupamento 5. Classificação 6. Nexo Causal. I. Pontifícia Universidade Católica do Paraná. Centro de Ciências Exatas e de Tecnologia. Programa de Pós-Graduação em Informática II-t

O Deus Pai, pois dele emana toda a sabedoria.

Agradecimentos

A meus Pais, por terem me ensinado o real significado da dignidade.

A Françoise e Wagner, pela compreensão, motivação e auxílio para realização do curso de Mestrado.

Minha orientadora, Professora Doutora Cinthia e meu co-orientador, Professor Doutor Altair, que em nenhum momento mediram esforços para poderem dar alicerces acadêmicos durante esta jornada.

A todos os Professores do Programa de Mestrado, pelo excelente trabalho dedicado na construção de meu conhecimento científico.

A todos os colegas, em especial a Rodolfo Botto e Lucas Galete, pelo suporte no entendimento teórico para questões ligadas a estudos realizados ao longo do curso.

A Adriano Nunes, pelo auxílio no entendimento de seu trabalho dedicado na investigação de crimes de assédio moral em *e-mails*.

E a todos que aqui não estão descritos. Aqueles que me incentivaram e, principalmente, respeitaram meus inúmeros momentos de silêncio que este estudo impôs.

Sumário

Agradecimentos	v
Sumário	vi
Lista de Figuras	ix
Lista de Tabelas	x
Lista de Símbolos	xi
Lista de Abreviaturas	xii
Resumo	xiii
Abstract	xiv
Capítulo 1	
Introdução	16
1.1. Escopo	17
1.2. Objetivo Geral	17
1.3. Objetivos Específicos	18
1.4. Contribuições	18
1.5. Motivação	18
1.6. Estrutura do Trabalho	19
Capítulo 2	
Fundamentação Teórica	20
2.1. <i>E-mail</i> : O que é?	21
2.2. <i>Softwares</i> de <i>e-mail</i> : Clientes e <i>Webmail</i>	27
2.3. Crimes Virtuais e Evidências Digitais	30
2.4. Grafos Direcionados	34
2.5. Técnicas de Processamento de Textos	35
2.5.1. Pré-processamento Textual e Extração de Características de Documentos	35
2.5.2. Representação dos Atributos Textuais	36
2.5.3. Agrupamento	37
2.5.4. Classificação	42

Capítulo 3

Trabalhos Relacionados	49
3.1. Processamento Automático de <i>e-mails</i>	49
3.2. Visualização de Conversações	56
3.3. Outros Trabalhos	57
3.4. Considerações Finais	60

Capítulo 4

Mecanismo para Rastreamento de Relacionamentos em <i>e-mails</i>	62
4.1. Visão Geral	62
4.2. Mecanismo	63
4.2.1. FASE I	63
4.2.2. FASE II	64
4.2.3. FASE III	67
4.2.4. FASE IV	68
4.3. Bases Experimentais	72
4.3.1. Dicionário Criminoso	73
4.3.2. Base de <i>e-mails</i> de Treinamento	73
4.3.3. Base de <i>e-mails</i> de Teste	75
4.4. Resultados Proporcionados	79

Capítulo 5

Resultados Experimentais	80
5.1. <i>Reader</i> – Implementação do Mecanismo para Rastreamento de Relacionamentos em <i>e-mails</i>	80
5.1.1. Leitura dos <i>e-mails</i> e processo de “imagem”	80
5.1.2. Pré-processamento	81
5.1.3. Crime Investigado	81
5.1.4. Extração e Geração dos Arquivos de Atributos	81
5.1.5. Algoritmos de Aprendizagem de Máquina	82
5.1.6. Exposição dos Resultados	83

5.1.7. Restrições do Mecanismo	85
5.2. Resultados Experimentais	85
5.2.1. Agrupamento de Conversações	86
5.2.2. Classificação das Conversações	88
5.2.3. Apresentação dos Resultados Experimentais	90
5.2.4. Análise dos Resultados	95
Capítulo 6	
Conclusão	98
Referências Bibliográficas	101
Apêndice A	
Apresentação do <i>Framework Reader</i>	109
A.1. Tela Principal	109
A.2. Processo de “Imagem” Forense e Pré-processamento	110
A.3. Gerar Arquivos de Atributos	110
A.4. Agrupamento de Conversações	111
A.5. Classificação das Conversações	111
A.6. Resultados	112

Lista de Figuras

Figura 2.1	Partes de um <i>e-mail</i> : Cabeçalho e Corpo	24
Figura 2.2	Exemplo de Grafo G	34
Figura 2.3	Objetos Agrupados [KUM06]	38
Figura 2.4	Classificação dos Tipos de Agrupamento	41
Figura 2.5	Modelo de Classificação [TAM05]	43
Figura 2.6	Exemplo do Classificador DT [TAM05]	47
Figura 3.1	Modelo Proposto [NAG10]	50
Figura 3.2	<i>Software BuzzTrack</i> [CSE06]	52
Figura 3.3	Modelo Proposto [IQB10]	55
Figura 3.4	Visualização de uma Conversação [PUP10]	57
Figura 3.5	Mecanismo para Detecção de Assédio Moral em <i>e-mails</i> [NUN09]	58
Figura 3.6	Vista Parcial do <i>Software</i> FTK versão 1.81.6 [DAT10]	59
Figura 4.1	FASE I do Mecanismo	64
Figura 4.2	FASE II do Mecanismo	65
Figura 4.3	Arquivo de Atributos – FASE II	66
Figura 4.4	FASE III do Mecanismo	67
Figura 4.5	FASE IV do Mecanismo	69
Figura 4.6	Resultado do Agrupamento e Classificação das Conversações	69
Figura 4.7	Grafos Gerados	71
Figura 4.8	Fragmento do Laudo Pericial Exemplo	72
Figura 5.1	Vista Parcial do Código para Geração da Matriz de Adjacências	84
Figura 5.2	Grafo Não Criminoso – Base Teste – Exemplo 1	92
Figura 5.3	Grafo Não Criminoso – Base Teste – Exemplo 2	93
Figura 5.4	Grafo Criminoso – Base Teste – Exemplo 1	94
Figura 5.5	Fragmento da Lista de <i>e-mails</i> Criminosos Detectados	94
Figura 5.6	Fragmento da Lista de Palavras Criminosas Utilizadas na Conversação 29	95
Figura A.1	Tela Principal – <i>Reader</i>	109
Figura A.2	Grafo Criminoso	115

Lista de Tabelas

Tabela 2.1	Representação da Matriz de Adjacências	43
Tabela 2.2	Representação da Matriz de Confusão	52
Tabela 4.1	Base de <i>e-mails</i> de Treinamento	74
Tabela 4.2	Base de <i>e-mails</i> de Teste	78
Tabela 5.1	Agrupamento de Conversações - Agrupamento	87
Tabela 5.2	Agrupamento de Conversações - Classificação	88
Tabela 5.3	Classificação de Conversações	89
Tabela 5.4	Comparação entre <i>Reader</i> e Processamento Manual	97
Tabela A.1	Processo de “Imagem” e Pré-processamento	110
Tabela A.2	Geração dos Arquivos de Atributos (FASE II e FASE III)	111
Tabela A.3	Agrupamento de Conversações	111
Tabela A.4	Classificação de Conversações	112
Tabela A.5	Laudo Pericial Exemplo	112
Tabela A.6	Fragmento com as Palavras Criminosas	112

Lista de Símbolos

\vec{D}	<i>E-mail</i> representado por seus atributos (vetor)
n	Quantidade de atributos
$tf-idf$	Peso para cada atributo (técnica tf-idf)
$Freq_{atributo.texto}$	Proporção de vezes que um atributo ocorre em um texto do <i>e-mail</i>
$DocFreq_{atributo}$	Número de <i>e-mails</i> em que um atributo é detectado
Q	Quantidade de <i>e-mails</i>
x_i	Atributo
$S_{X,Y}$	Distância entre 2 <i>e-mails</i> – Distância <i>Euclidiana</i>
X, Y	<i>E-mails</i> 1 e 2
W_{Xk}	Atributo do <i>e-mail</i> X
W_{Yk}	Atributo do <i>e-mail</i> Y
Coc	Intersecção entre os atributos de 2 <i>e-mails</i> comparados
Cit	União entre os atributos de 2 <i>e-mails</i> comparados
C	Classes
C	Conjunto de classes
X	Quantidade de classes
T	Conjunto de <i>e-mails</i> (treinamento) representados por seus atributos
f_{11}	Quantidade de registros (<i>e-mails</i>) classificados corretamente na Classe 1
f_{10}	Quantidade de registros (<i>e-mails</i>) classificados corretamente na Classe 0
f_{01}	Quantidade de registros (<i>e-mails</i>) classificados incorretamente na Classe 1
f_{00}	Quantidade de registros (<i>e-mails</i>) classificados incorretamente na Classe 0
$P(. .)$	Probabilidade condicional conforme um elemento
$P(.)$	Probabilidade condicional de um elemento
V_{NB}	Probabilidade da escolha de uma classe a um <i>e-mail</i>
π	A representação $(. .)$ generaliza elementos que possam aparecer de diferentes maneiras no texto do trabalho. Exemplo: $P(\vec{d}_j c_i)$.

Lista de Abreviaturas

Agrupamento	<i>Cluster</i>
API	<i>Application Programming Interface</i>
BD	Banco de Dados
<i>Browse</i>	Navegador
Cabeçalho	<i>Header</i>
Corpo	<i>Body</i>
DT	<i>Decision Trees</i>
FTP	<i>File transfer protocol</i>
HTML	<i>HyperText Markup Language</i>
HTTP	<i>Hypertext Transfer Protocol</i>
IP	<i>Internet Protocol</i>
LSI	<i>Latent Semantic Indexing</i>
MIME	<i>Multipurpose Internet Mail Extensions</i>
MP	Medida Provisória
NB	<i>Naïve Bayes</i>
OST	<i>Offline Folder File</i>
PL	Projeto de Lei
PST	<i>Personal Storage Table</i>
RFC	<i>Request for Comments</i>
SVM	<i>Support Vector Machines</i>
Tf-idf	<i>Term frequency-inverse document frequency</i>
URL	<i>Uniform Resource Locator</i>

Resumo

No uso da tecnologia *e-mail* podem ser cometidos crimes virtuais, e em caso de investigação judicial do crime, é necessária a produção de provas digitais. Aplicam-se métodos científicos para comprovação de autoria, e assim, formalização do nexos causal. Apresenta-se neste trabalho um mecanismo para produção de provas digitais a partir de rastreamento em relacionamentos existentes em bases de *e-mail* auxiliando na formação do nexos causal. Investigam-se os relacionamentos entre o usuário proprietário de uma base de *e-mails* e seus contatos a partir do agrupamento de conversações (*e-mails* que compartilhem do mesmo conjunto textual em seu corpo) criminosas. A análise criminosa faz a comparação por similaridade de palavras pertencentes a uma conversação e as de um dicionário criminoso de palavras. Assim, são agrupadas e classificadas as conversações de uma base de *e-mails*, que posteriormente são visualizadas graficamente (grafos direcionados). No trabalho, estudam-se os métodos de classificação (*K-means* no uso de diferentes funções - *Euclidiana*, *Manhattan*, *Cossine* e *Jaccard*, *Naïve Bayes*, *Árvore de Decisão* e *Máquinas de Vetor de Suporte*). Os resultados experimentais para uma base de 570 *e-mails* com 33 conversações submetidas à classificação de crime de assédio moral atingiram taxas de acertos próximas de 99%, demonstrando ser a pesquisa promissora no que se refere ao rastreamento em relacionamentos por *e-mails* para a produção de provas digitais.

Palavras-Chave: Rastreamento; *e-mail*; Extração de Características; Agrupamento; Classificação; Nexos Causal.

Abstract

The e-mail may be involved in a cybercrime and it is necessary find the digital forensic evidences. Methods and techniques are applied to provide the proofs of authorship, and thus formalizing the causal nexus. This work presents a mechanism to produce digital evidences based on relationships existing in e-mails supporting the causal nexus. The mechanism investigates the relationships between the user-owner of an e-mails database and its contacts based on classification of the criminal conversation (e-mails that share the same textual set in its body). The criminal analysis makes the comparison by similarity of words from a conversation and a dictionary of criminal words. Therefore, are identified and classified the conversations existing on a set of e-mails, which are then displayed graphically (directed graph). For this purpose, this work presents a study on classification methods (K-means in use of different functions - Euclidian, Manhattan, Cossine and Jaccard, Naïve Bayes, Decision Tree and Support Vector Machine). The experimental results obtained for a base of 570 e-mails with 33 conversations submitted to the classification bullying crime at work environment achieved rates of 98% proving to be promising to provide digital evidences from tracking of relationship in a e-mails database.

Keywords: Tracking; e-mail; Extraction of Features; Cluster; Classification; Causal Nexus.

Capítulo 1

Introdução

A *Internet* a partir da década de 1990 se tornou popular, deixando de pertencer a grupos específicos, passando a oferecer um vasto número de serviços virtuais aos usuários da sociedade. Citam-se alguns exemplos de serviços disponibilizados: *home banking*, comunicação instantânea, *e-mails*, leilões virtuais, compras *on-line* (lojas virtuais), contratos virtuais, documentos eletrônicos e *sites* de relacionamento.

Não resta dúvida dos benefícios e da massiva utilização dos serviços oferecidos através da tecnologia *Internet*. Entretanto, as informações empregadas por serviços virtuais ficam expostas para todos os usuários também conectados na *Internet*. De acordo com PINHEIRO (2009a) “usuários mal intencionados, confiando na falsa sensação de anonimato, praticam condutas criminosas através da *Internet*” [PIN09a].

Condutas criminosas denominada de crimes no cyber espaço, ou ainda, os *cibercrimes* [PIN09a] [BRO06] [SUR05] [CIA04]. Cita-se como exemplo, a propagação de crimes de pedofilia na *Internet* [CHO06], e casos de assédio sexual [SIP99].

Como estatística informa-se a pesquisa realizada pela instituição ONG *SaferNet* Brasil¹ – associação civil de direito privado, a qual é responsável pela Central Nacional de Denúncias de Crimes Cibernéticos, sendo que esta aponta os seguintes crimes para denúncia: apologia e incitação a crimes contra a vida, homofobia, intolerância religiosa, maus tratos contra animais, neonazismo, pornografia infantil, racismo e xenofobia.

Ao final de todos os meses a *SaferNet* apresenta relatório dos *cibercrimes* denunciados. Durante o mês de novembro de 2010, foram registradas 3.015 denúncias de

¹ ONG *SaferNet* - Denúncias de *cibercrimes*, disponível em <http://www.safernet.org.br/site/indicadores>, acesso em 18 de janeiro de 2011.

crimes ligados ao *site* de relacionamentos *Orkut*². Também foram registradas 3.743 denúncias que não estavam ligados ao domínio *Orkut*.

Para que o infrator do *cibercrime* sofra punição em decorrência de sua prática perante uma legislação, o *cibercrime* deve ser registrado aos órgãos competentes. Nesta situação, serão coletadas evidências digitais do *cibercrime* cometido, e realizada a produção de provas digitais.

Na produção de provas digitais de crime registrado em *e-mail*, como casos de troca de mensagens eletrônicas com comentários de forma negativa, sobre raças ou religiões [STJ08], fazem-se necessário a existência de mecanismos que facilitem este tipo de trabalho, principalmente ao saber do alto volume de *e-mails* que um usuário pode ter.

1.1. Escopo

A análise de evidências digitais estuda os elementos coletados na cena de um *cibercrime* para comprovação de que realmente estejam ligados ao autor-infrator [OLI02] [PHI09]. Mediante a análise do caso concreto, comprova-se a existência de relação entre a conduta do infrator e o crime cometido, tendo-se o nexa causal [DEL00] [JES02].

Este trabalho apresenta um mecanismo que auxilia os peritos da área de Computação Forense na produção e análise de evidências digitais a partir de rastreamento em *e-mails*, permitindo assim a comprovação do nexa causal.

Analisa-se conversações entre contatos de um usuário existentes em uma base de *e-mails* não cifrada/criptografada e armazenada preferencialmente no servidor de *e-mail* por questões de garantias jurídicas, mas em situações em que isto não ocorre, pode-se considerar a base de *e-mails* armazenada em pasta/diretório designado pelo usuário. Deste modo, o mecanismo apresentado neste trabalho possibilita a análise de conversações entre diferentes endereços de *e-mail* do remetente e destinatário envolvidos em uma conversação.

1.2. Objetivo Geral

Este trabalho apresenta um mecanismo para produção de provas digitais que possam ser encontradas na base de *e-mails* de um usuário. Rastreiam-se evidências a partir de conversações supostamente criminosas existentes entre o proprietário da base de *e-mails* e

² Site do *Orkut*, disponível em <http://www.orkut.com>, acesso em 18 de janeiro de 2011.

seus contatos. Objetiva-se auxiliar na comprovação do nexo causal que provará o relacionamento entre o infrator de um crime e seu delito.

1.3. Objetivos Específicos

Tendo-se por plano de fundo a determinação do nexo causal, apresentam-se os objetivos específicos para realização deste estudo. Entre eles: levantamento bibliográfico, entendimento de métodos de agrupamento e classificação, pré-processamento textual, formatação de *e-mails*, grafos direcionados, crimes virtuais e conceitos jurídicos, tais como nexo causal.

1.4. Contribuições

Com o desenvolvimento do mecanismo apresentado tem-se um método semi-automático de identificação de *cibercrimes* através do rastreamento de trocas de mensagens de *e-mails*. O rastreamento permite a construção dos relacionamentos criminosos que servirão de base para comprovação do nexo causal.

O estudo permite a aplicação de técnicas e métodos em situações de assédio moral [NUN09]. Além disto, no emprego de técnicas e métodos já utilizados por outros autores, contribui-se na geração do mecanismo para forense digital, haja vista que técnicas e métodos estudados em trabalhos como [NAG10] [IQB10] [DAT10] [NUN09] [BAL08a] [BAL08b] [TAM08] [CSE06] [DRE06] [VIE06] [WAT03] proporcionam a investigação de evidências digitais.

Uma contribuição relevante é o fato do mecanismo ter sido aplicado em um contexto real relacionado com assédio moral. Isto foi possível por meio da associação deste trabalho com os desenvolvimentos realizados por Nunes et al. (2009) [NUN09]. O mecanismo foi publicado no X Simpósio Brasileiro de Segurança da Informação e de Sistemas Computacionais - SBSeg 2010 [MAL10].

1.5. Motivação

As motivações em realizar este trabalho estão na inexistência de mecanismos que realizem rastreamento de conversações praticadas por *e-mail*, ou mesmo na existência de poucos trabalhos científicos que utilizem métodos que possam ser utilizados pela investigação forense. Além disto, observa-se um número crescente de crimes virtuais [STJ08].

1.6. Estrutura do Trabalho

O Capítulo 2 apresenta a fundamentação teórica estudada para realização do trabalho de Mestrado, dentre eles, conceitos técnicos aplicados no Capítulo 4 e 5: técnicas de processamento de textos, grafos, métodos de agrupamento e classificação. No Capítulo 3, os chamados trabalhos relacionados. No Capítulo 4 apresenta-se o mecanismo para produção de provas digitais a partir de conversações encontradas em *e-mails*, assim como as bases de *e-mails* e dicionário utilizado durante os experimentos. No Capítulo 5 relatam-se os testes experimentais realizados pela aplicação do mecanismo. Descrições do *framework* implementado, os resultados, assim como as restrições (limitações) do *framework*. E no Capítulo 6 são apresentadas as conclusões do trabalho.

Capítulo 2

Fundamentação Teórica

Para o *site* de referência em *Internet*³, *e-mail* é a ferramenta de comunicação mais utilizada no uso dos serviços virtuais proporcionados pelo uso da *Internet*. E também, em estatísticas expostas no *site Internet World Stats*⁴ o uso da *Internet* no mundo, considerando uma população mundial equivalente a 6.845.609.960 de habitantes, registrou até o final do mês de junho de 2010, a existência de 1.966.514.816 utilizadores de serviços eletrônicos providos pela *Internet*. E ainda, comparando os números apresentados em 2010 com a pesquisa realizada no ano de 2000, em que se registrou 360.985.492 usuários de serviços *Internet*, totaliza-se um aumento de 544,76%.

A grande maioria dos usuários utiliza o computador para ler, pesquisar, e se corresponder via *e-mail*. Como comprovação, apresenta-se a pesquisa realizada no Brasil pelo Centro de Estudos sobre as Tecnologias da Informação e da Comunicação - CETIC⁵, a qual concluiu que no ano de 2007, dentre 5.823 usuários da *Internet*, 78% utilizavam os serviços da *Internet* para comunicação via *e-mail*. No ano de 2008, foram entrevistados 8.207 usuários, e destes 77% utilizavam o *e-mail* como serviço de comunicação. E ainda, no ano de 2009, dentre 9.747 entrevistados, a pesquisa revela que 85% utilizavam *e-mails*.

Assim, no trabalho aqui realizado, é proposto e implementado um mecanismo para produção de provas digitais em *e-mails*. O mecanismo faz uso de diferentes tecnologias. Relata-se neste Capítulo uma revisão da tecnologia de *e-mail* - definição, formatação,

³ *The World's First Book Published On The Web* (2000), disponível em <http://www.livinginternet.com/e/e.htm>, acesso em 18 de janeiro de 2011.

⁴ *Internet Usage Statistics*, disponível em <http://www.internetworldstats.com/stats.htm>, acesso em 18 de janeiro de 2011.

⁵ Comitê Gestor da Internet Brasil, Centro de Estudos sobre Tecnologias da Informação e da Comunicação, disponível em <http://www.cetic.br>, acesso em 18 de janeiro de 2011.

funcionamento, modos de acesso aos *e-mails* na utilização de *softwares* cliente ou um navegador (*browser*) e armazenamento de *e-mails*. Também se estudam os crimes virtuais, as evidências digitais e as tipificações para o *cibercrime* frente à atual legislação Brasileira. Por fim, este Capítulo, fundamenta os conceitos de grafos direcionados, métodos de agrupamento e classificação, os quais foram aplicados ao longo do desenvolvimento deste trabalho.

2.1. *E-mail*: O que é?

O *e-mail* é uma das aplicações de rede mais antigas [PET04], sendo um método para enviar e receber mensagens eletrônicas mediante utilização de sistemas eletrônicos de comunicação [LIM07]. Desta forma, pessoas em diferentes pontos geográficos podem se comunicar de forma fácil, rápida, e a baixo custo [KUR10].

E-mail tornou-se a forma mais comum para transferência de informações de seus usuários [PHI09]. E ainda, pesquisa realizada pelo grupo *Radicatti* no ano de 2009⁶, registrou aproximadamente 1.4 bilhões de usuários de *e-mails* e o número de 247 milhões de mensagens enviadas por dia. Baseado nessas informações, o grupo *Radicatti* estimou a existência de 1.9 bilhões de usuários de *e-mail* até o ano de 2013. E também, projetou que em 2013, o total de *e-mails* enviados por dia será de aproximadamente 507 bilhões.

WHITTAKER et al. (1996) estudou a evolução da tecnologia *e-mail*. Na realização daquele trabalho já era possível enviar e receber mensagens que continham documentos incorporados em sua estrutura [WHI96]. Citam-se como exemplos, os arquivos anexados, tais como figuras, vídeos, músicas e binários. Assim, esta comunicação eletrônica incorporou outras finalidades, como por exemplo, a transferência de arquivos.

Cabe, especificar os pré-requisitos para funcionamento da transmissão de mensagens de *e-mail* [SOA08]:

- Um usuário remetente e *software* de *e-mail* instalado em dispositivo eletrônico remetente;
- Um ou mais usuários destinatários e *software* de *e-mail* instalado em dispositivo eletrônico destinatário;
- Meio digital para condução da mensagem eletrônica entre remetente e destinatário(s) (*Internet*);

⁶ Pesquisa disponível em <http://www.radicati.com/?p=3237>, acesso em 18 de janeiro de 2011.

- Existência de uma conta de *e-mail* para o remetente;
- Existência de uma conta de *e-mail* para cada destinatário.

Agregado a uma conta de *e-mail* tem-se um endereço virtual. Sua formação é dividida pelo nome do usuário e um domínio da *Internet*, separados pelo símbolo @. Exemplo disso é o endereço user1@exemplo.com.br. O *login* de acesso deste endereço é user1, sendo único dentro do domínio exemplo, o qual foi registrado como pertencente ao país Brasil, visto a indicação de “br” representando o domínio.

Uma conta de *e-mail* pode ser solicitada/criada em servidores de *e-mail* disponibilizados gratuitamente ou não. Algumas empresas privadas que disponibilizam este tipo de serviço gratuitamente são a *Yahoo*⁷ e *Gmail*⁸.

Através do endereço de *e-mail* é possível identificar informações de um usuário proprietário [PIN09a]. Caso o endereço pertença a um servidor particular, os dados do usuário não podem ser omitidos com facilidade, como por exemplo, contas de *e-mail* em uma corporação. Entretanto em um servidor público, os dados do usuário podem ser omitidos com maior facilidade, entretanto existem formas de rastreamento.

Em contra partida, o armazenamento de *e-mails* no recebimento de mensagens eletrônicas, tanto no remetente como no destinatário, acarretam num acúmulo de *e-mails* que acabam sendo armazenados no computador ou servidor de usuários [WHI96] [FIS06].

WHITTAKER et al. (1996) analisou e reportou o comportamento de 20 usuários da ferramenta *e-mail* [WHI96]. Após dez anos FISHER et al. (2006) refez a análise, entretanto utilizando 600 participantes, chegando aos mesmos resultados [FIS06]. É importante relatar que após 10 anos da análise, apresentou-se um aumento em média de 10 vezes no tamanho da base de *e-mails* armazenados em dispositivos eletrônicos.

Evidenciou-se a necessidade de se gerenciar estas mensagens e então funcionalidades foram adaptadas aos *softwares* de *e-mail*. Citam-se algumas tarefas que os indivíduos passaram a realizar em seus *softwares* de *e-mail*, conforme [WHI96]:

- Entrega de documentos;
- Arquivamento de documentos;

⁷ Site da empresa *Yahoo*, disponível em <http://www.yahoo.com>, acesso em 21 de janeiro de 2011.

⁸ Site da empresa *Gmail*, disponível em <http://www.gmail.com>, acesso em 21 de janeiro de 2011.

- Delegação de tarefas de trabalho;
- Acompanhamento de tarefas;
- Armazenamento de nomes pessoais e endereços (contatos);
- Envio de lembretes;
- Solicitação de assistência;
- Agendamento de consultas;
- Manipulação de consultas de suporte técnico.

Por outro lado diversos serviços de comunicação são disponibilizados na *Internet*, como programas de comunicação instantânea, *e-mails*, fóruns, *sites* [PIN09a]. Dentre estas, a ferramenta *e-mail* pode ser utilizada em diversos ambientes, sejam estes corporativos ou não.

Embora pessoas possam se relacionar e transmitir tipos de informações variadas através da tecnologia *e-mail*, confidenciais ou não, de forma rápida, ainda existe a falta de controle sobre o conteúdo transmitido [WHI96]. Para tal, DUANE et al. (2004) realizou estudo para implantação de políticas e monitoramento do uso de mensagens eletrônicas [DUA04]. Neste estudo o autor descobriu que uma grande quantidade de mensagens de *e-mail* é a causa de efeitos negativos aos usuários desta ferramenta, causando perda de produtividade. Cita-se como exemplo a utilização desta tecnologia para finalidades ilícitas, como o encaminhamento de mensagens ofensivas (assédio moral).

Verificou-se anteriormente que as informações transmitidas em uma mensagem eletrônica normalmente ficam armazenadas no computador do usuário, ou então no servidor de *e-mail*. Através destas informações é possível analisar e averiguar o histórico de conversações de usuários [VIE06]. Desta forma, através da utilização das conversações entre um usuário e seus contatos em *e-mails* armazenados em dispositivos eletrônicos é possível demonstrar a existência de provas digitais que comprovem o cometimento de um crime executado por *e-mail* [PIN09b].

As conversações realizadas por *e-mails* envolvem uma sequência de *e-mails* entre contatos, fazendo com que compartilhem de um mesmo contexto (conjunto textual). Como exemplo, cita-se um *e-mail* enviado de um contato remetente ao contato destinatário, e depois respondido. Caso o *e-mail* respondido mantenha em seu corpo a cópia do conjunto textual composto pelo primeiro *e-mail* (enviado), proporciona-se a existência de uma conversação envolvendo dois *e-mails* (*e-mail* enviado e *e-mail* respondido). Estes 2 *e-mails* estarão

armazenados na base de *e-mails* do contato remetente e do destinatário, visto que ambos receberam e enviaram 1 *e-mail*. E ainda, uma conversação pode envolver n contatos, haja vista que um *e-mail* pode ser enviado para diferentes contatos, como será visto nas explicações sobre os campos do cabeçalho de um *e-mail*.

Sendo possível a análise das informações transmitidas em formato de *e-mails*, estuda-se na sequência a formatação do *e-mail*. O *e-mail* é composto por uma série de caracteres, separados em dois grupos: cabeçalho e corpo. Exemplo de formatação pode ser visualizado em formato texto na Figura 2.1.

As especificações da formatação de um *e-mail* texto podem ser obtidas na RFC (*Request for Comments*) 5322⁹. De acordo com a RFC 5322, quando um *e-mail* é transferido de um usuário a outro(s), informações são adicionadas nos *e-mails*. Por exemplo, informações são incorporadas na parte inicial da mensagem, no cabeçalho.



Figura 2.1 – Partes de um *e-mail*: Cabeçalho e Corpo

⁹ *Request For Comments: Internet Message Format*, disponível em <http://www.ietf.org/rfc/rfc5322.txt>, acesso em 18 de janeiro de 2011.

O cabeçalho de um *e-mail* é composto por várias linhas, sendo que cada linha possui um campo e uma informação, tendo o sinal de dois pontos separando o campo e a informação. Na sequência são descritos alguns campos, os quais foram identificados na Figura 2.1.

A maioria dos campos pertencentes ao cabeçalho possui padronização. Embora haja essa padronização, cada *software* de *e-mail* também pode fazer uso de campos adicionais com sua própria formatação. Na sequência, listam-se as informações essenciais para identificação de um *e-mail*:

- Assunto da mensagem (*Subject*);
- Endereço de *e-mail* daquele que enviou a mensagem (*From* ou *Sender*);
- Endereço de *e-mail* dos destinatários (*To*);
- Quando um *e-mail* é enviado com cópia (*CC*);
- Data e horário de envio da mensagem (*Date*).

Cada *e-mail* possui identificadores. No campo identificador da mensagem (*Message-ID*) existe um número acompanhado do nome da máquina de onde partiu a mensagem. Em outro campo, é feita a referência à mensagem que originou a atual (*In-Reply-To*).

Além de possuir identificadores, um *e-mail* pode fazer referência às mensagens que pertencem às conversações anteriores (*References*). O campo *References* é útil quando uma sequência de mensagens é respondida. Entretanto, quando *e-mails* são encaminhados (*Forward*), o campo *References* apenas detalha o *e-mail* que originou esta mensagem.

Sabe-se ainda que um *e-mail* pode ser enviado para um ou mais usuários, podendo-se utilizar vários campos existentes no cabeçalho de uma mensagem, a saber: o campo Para (*To*), Com Cópia (*CC*), ou ainda, o campo Cópia Carbono ou Cópia Oculta (*CO*). Estes campos são oferecidos durante a composição de uma mensagem eletrônica.

De forma contrária, quando um *e-mail* é recebido por um usuário, é possível identificar através do cabeçalho, para quem este *e-mail* foi replicado. Entretanto, não é possível visualizar através do cabeçalho, aquele usuário que recebeu a mensagem de forma oculta.

Outro campo que permite o rastreamento de um *e-mail* é o campo *Received*. O campo *Received* descreve a transferência da atual mensagem ao longo de seu percurso na rede

mundial de computadores. Citam-se como exemplo, os servidores que realizaram tratamento da mensagem.

Além de campos específicos do cabeçalho descritos pela RFC 5322, *softwares* de correio eletrônico podem fazer uso de outros. Os campos desenvolvidos por empresas proprietárias visam à manipulação de determinadas informações do cabeçalho. Campos iniciados pelo caractere X, como por exemplo, *X-Originating-IP*, *X-Apparently-To*, *X-YMailISG* e *X-BeenThere* representados na Figura 2.1. No campo *X-Originating-IP* é identificado o endereço IP (*Internet Protocol*) do equipamento utilizado pelo remetente da mensagem.

Ainda explicando a formatação do cabeçalho, verifica-se a existência de campos pertencentes ao MIME (*Multipurpose Internet Mail Extensions*). O MIME inclui uma estrutura para descrever no cabeçalho do *e-mail* o conteúdo que o corpo do *e-mail* transporta. Cita-se como exemplo o campo *Content-Type*. Este campo identifica o tipo e formato do conteúdo do corpo. Informações detalhadas sobre os formatos MIME podem ser encontradas nas RFC 2045 até RFC 2049¹⁰.

Após as informações do cabeçalho, uma linha em branco realiza a separação desse com o corpo do *e-mail*. No corpo do *e-mail* encontra-se a mensagem eletrônica transportada, assim como o conteúdo incorporado (anexos), sendo que a mensagem e os possíveis anexos pertencentes a um *e-mail* são delimitados pelo código gerado em *Content-Type*. Assim, cada anexo será delimitado por um código gerado.

Durante a explicação sobre a formatação de uma mensagem de *e-mail* recebido focou-se principalmente nos campos que possibilitam o rastreamento deste *e-mail*. Analisam-se agora os campos de um *e-mail* encaminhado ou respondido. Uma mensagem encaminhada ou respondida faz com que o destinatário desta mensagem envie uma mensagem com um novo cabeçalho. Ademais, o histórico de mensagens incorporadas neste *e-mail* perde muitas informações de seus cabeçalhos.

Na nova mensagem, mantêm-se apenas os campos *From*, *To*, *Subject* e *Date*. Tais informações permanecem no corpo da mensagem recebida e não no cabeçalho. Interessante salientar que os usuários podem alterar tais informações.

A RFC 5322 aconselha desenvolvedores de *softwares* de manipulação de *e-mail* que sempre armazenem todas as informações do cabeçalho de uma mensagem. Verifica-se neste

¹⁰ *Request For Comments* disponíveis em <http://www.ietf.org/>, acesso em 21 de janeiro de 2011.

conselho a possibilidade do rastreamento de uma mensagem eletrônica ao longo de todo o seu percurso. Inclui-se aqui a lista de todos os servidores de *e-mail* que realizaram tratamento na mensagem. Entretanto, na prática, este conselho não tem sido implementado pelos desenvolvedores de *softwares* gerenciadores de *e-mail*.

Com a utilização de informações existentes no cabeçalho de *e-mails* é possível realizar o rastreamento de *e-mails*. Mas estas informações podem ser alteradas/omitidas por usuários mal intencionados ao enviarem *e-mails*. E informações contidas no corpo de *e-mails*, uma vez que um *e-mail* é respondido ou encaminhado, recebe, sempre, uma cópia do conteúdo textual do *e-mail* que o gerou, caso este conteúdo não seja excluído pelo usuário que irá responder um *e-mail*. Assim, uma tentativa de realizar o rastreamento de *e-mails* para produção de provas digitais é a investigação pelo conjunto textual existente no corpo e campo assunto dos *e-mails* pertencentes a uma conversação.

Com base na definição de *e-mail* apresentada e detalhada pode-se seguir explicando sobre esta ferramenta de comunicação via *e-mail*. Mediante conhecimento do próximo tópico, haverá um melhor entendimento na ferramenta *e-mail*.

2.2. Softwares de e-mail: Clientes e Webmail

Usuários de *e-mail* possuem duas maneiras para gerenciamento de suas mensagens eletrônicas. Com auxílio de *software* cliente instalado na máquina do usuário ou por intermédio de *Webmail* (via um *browser*). Indiferente da maneira escolhida, usuários são responsáveis pelo conteúdo de seus *e-mails*.

Os *softwares* de *e-mail* clientes são programas instalados em computadores pessoais e não são necessariamente vinculados a uma conta de *e-mail* específica. Adentra-se na possibilidade de poderem se comunicar com o *software* instalado no servidor de qualquer tipo de conta. De forma contrária, o acesso via *Webmail*, é disponibilizado pela empresa proprietária do servidor em que o usuário possui uma conta de *e-mail*. Um *software* é instalado pela empresa proprietária do domínio. Esta instalação é efetuada em um servidor, e para que o usuário de *e-mail* possa utilizar este *software*, precisar acessá-lo via um *browser*.

O número de *softwares* clientes e *Webmails* existentes no mercado é grande. Diferenciam-se por aspectos de funcionalidades que cada um apresenta, podendo-se citar: facilidade de utilização, armazenamento dos *e-mails* do usuário, gerenciamento de contatos,

calendário, existência de filtro para arquivamento e maneira como são visualizados os *e-mails* armazenados por este *software*.

Pesquisa realizada e divulgada pela empresa privada *CampaignMonitor* [CAM09a], lista os *softwares* clientes de *e-mail* mais populares utilizados por usuários de *e-mail*. A análise refere-se aos meses entre Janeiro e Junho de 2009, e por meio de um *software* próprio, realizou o envio de *e-mail* para seus clientes contendo um anexo, no caso uma figura/imagem. Quando os pesquisados realizavam a leitura do *e-mail* e visualizavam a figura/imagem, o *software* realizava o aviso para a empresa *CampaignMonitor*. Assim, registrou-se 300 milhões de *e-mails* recebidos. Entretanto a empresa não divulgou o número de usuários utilizados no teste. Os resultados da pesquisa apontaram os seguintes *softwares* cliente de *e-mail* como os mais utilizados:

- *Outlook* 2000, 2003 e *Express* da empresa *Microsoft Windows*¹¹: 32,08%;
- *Yahoo Mail*¹²: 15,65%;
- *Hotmail*¹³: 5,35%;
- *Outlook* 2007 (*Microsoft Windows*): 7,55%;
- *Mozilla Thunderbird*¹⁴ versão 2: 1,12%.

Nesta pesquisa, usuários que não abriram o anexo que estava no *e-mail* enviado pela empresa *CampaignMonitor* não participaram da pesquisa. A cada visualização de um mesmo anexo foi contabilizado o tipo de *software* utilizado pelo usuário. Considerando que alguns *softwares* existentes no mercado são configurados para abertura automática de imagens em *e-mails* recebidos, evidencia-se que grupos de usuários foram favorecidos no resultado da pesquisa. De forma contrária, *softwares* que visualizam *e-mails* recebidos em formato texto não participaram da pesquisa.

Demonstra-se também a pesquisa sobre a utilização de *e-mail* apresentada na dissertação de CSELLE (2006). Computou-se valores por GROUP (2006) entre os anos de 2001 e 2006 na utilização de 8.900 *e-mails* [GRO06]. Os 3 *softwares* para estações clientes

¹¹ Site da empresa *Microsoft Windows* disponível em <http://www.microsoft.com>, acesso em 21 de janeiro de 2011.

¹² Acesso ao *Webmail* do *Yahoo*, disponível em <http://mail.yahoo.com.br>, acesso em 21 de janeiro de 2011.

¹³ Acesso ao *Webmail* do *Hotmail*, disponível em <http://www.hotmail.com>, acesso em 21 de janeiro de 2011.

¹⁴ Site da comunidade *Mozilla*, disponível em <http://www.mozilla.org>, acesso em 21 de janeiro de 2011.

mais utilizados são listados abaixo, sendo que não foi divulgada a técnica utilizada nesta pesquisa [CSE06]:

- *Microsoft Outlook Express* com 16,2%;
- *Mozilla Thunderbird* com 13,9%;
- *Microsoft Outlook* com 10,6%;
- *Microsoft Outlook Web Access* com 7,4%;
- *Gmail* com 5,9%;
- *Yahoo Mail* com 1,6%.

Em outra pesquisa, foram entrevistados por *e-mail* 484 participantes que trabalhavam em diferentes companhias dos Estados Unidos [DAB05]. Deste modo, obteve-se a lista dos *softwares* mais utilizados, a saber:

- *Microsoft Outlook*: 76%;
- *Lotus Notes*¹⁵: 7%;
- *Novell GroupWise*¹⁶: 6%;
- *Mozilla Thunderbird*: 2%;
- Outros *softwares* totalizaram 10%.

E, constata-se mediante pesquisa divulgada e realizada pela empresa *Net Application*¹⁷, em que a grande maioria de usuários de computador utiliza o Sistema Operacional da empresa *Microsoft Windows*. Entende-se desta forma haver um aumento na probabilidade da ocorrência de delitos que possam ser deflagrados neste tipo de Sistema Operacional. Sendo que, dentre os *softwares* de correio eletrônico, positiva-se o *software Windows Live Mail*¹⁸, conforme reportagem¹⁹, em função da análise positiva de seu antecessor, o *software Outlook Express*, que não esta mais sendo comercializado.

¹⁵ *Lotus Software*, disponível em <http://www.lotus.com>, acesso em 21 de janeiro de 2011.

¹⁶ Site da empresa *Novell*, disponível em <http://www.novell.com>, acesso em 21 de janeiro de 2011.

¹⁷ Pesquisa disponível em <http://marketshare.hitslink.com/operating-system-market-share.aspx?qprid=8>, acesso em 23 de janeiro de 2011.

¹⁸ *Download* disponível em <http://windows.microsoft.com/pt-BR/windows/downloads?T1=downloadsWinLive>, acesso em 23 de janeiro de 2011.

¹⁹ Reportagem disponível em <http://windows.microsoft.com/pt-BR/windows-vista/So-long-Outlook-Express-Introducing-Windows-Live-Mail>, acesso em 23 de janeiro de 2011.

Em resumo, esta Subseção apresentou os tipos de *software* para que um usuário possa gerenciar seus *e-mails*. A filosofia de funcionamento entre os tipos de *software* cliente e *browser*, difere. No primeiro, *e-mails* são gerenciados em um computador pessoal, e no segundo, a partir de um servidor. Assim, apresentou-se e justificou-se o *software* cliente de *e-mails* e o Sistema Operacional escolhidos para realização de experimentos que são relatados no Capítulo 4.

Visto isso, explanam-se no próximo tópico questões relativas à criminalidade virtual. Apresentam-se a definição de crime virtual, evidência digital e sua coleta, legislação, *e-mail* como prova jurídica, ato de periciar, além de seu agente, o profissional perito.

2.3. Crimes Virtuais e Evidências Digitais

Com o advento dos serviços virtuais para a sociedade, houve a inserção do cometimento de novos tipos de crimes [BRO06]. Além dos crimes convencionais, a sociedade passou a enfrentar crimes cometidos por meio eletrônico. Crimes chamados de *cibercrimes* [PIN09a] [BRO06] [CIA04] [STU99].

Crimes que frequentemente podem ser tipificados como crimes tradicionais, embora sejam executados com maior rapidez, em número grande de vítimas, e com impacto entre jurisdições, em relação aos crimes convencionais. Isto devido ao fato dos crimes virtuais não obrigarem vítima e criminoso a estarem em um mesmo ambiente físico.

A quantidade de ocorrências envolvendo crimes virtuais vem aumentando, visto que pesquisa do Supremo Tribunal de Justiça – STJ mostra que no ano de 2008 houveram 17.000 decisões judiciais no Brasil no tocante a criminalidade virtual. No ano de 2002, eram apenas 400 decisões [STJ08].

A lista de crimes cometidos eletronicamente é grande. De acordo com o Supremo Tribunal Federal de Justiça, STJ (2008) e ainda com o trabalho de BROADHURST (2006) são citados exemplos de crimes virtuais. Entre outros, fraude, assédio sexual e moral, racismo, ameaças e pedofilia [PHI09] [STJ08] [BRO06]. E ainda, PINHEIRO (2009) e WENDT (2010) afirmam ser a fraude e a pedofilia os crimes virtuais que mais ocorrem no Brasil atualmente [WEN10] [PIN09a].

No Brasil não existe uma legislação específica para tipificação de crimes virtuais [PIN09a] [NOG09] [STJ08], ao contrário de países que possuem uma legislação mais preparada, como o Canadá [FIL10] e a Argentina [MES08]. Sendo consenso entre diversos

autores que quando crimes virtuais ocorrem no Brasil, acabam sendo tipificados através do Código Penal, Civil e legislações específicas na esfera brasileira. Com base em STJ (2008) afirma-se que em 95% dos casos, os crimes envolvendo informática são tipificados mediante a utilização do Código Penal, sendo que apenas 5% dos tipos de crimes por computador não podem ser tipificados [STJ08]. E ainda, para STJ (2008) afirma que mesmo na ausência de leis específicas para crimes virtuais no Brasil, mediante o conhecimento de aplicadores da lei, esses crimes podem ser puníveis [STJ08].

Existem projetos de lei no Brasil que visam à regulamentação de novas leis para tipificação de crimes virtuais [NOG09] [PIN09b] [STJ08] [SOA08]. Entre eles, o mais comentado é o Projeto de Lei 84/99 (PL) que foi aprovado em primeira instância pela Câmara dos Deputados, revisado pelo Senado, e que agora aguarda segunda aprovação da Câmara [DEP99]. O PL 84/99 visa enquadrar juridicamente 13 tipos de crimes cometidos no universo da informática. E ainda, para os crimes que já podem ser tipificados através do Código Penal, propõem-se um aumento da penalidade [SEN08].

Crimes virtuais podem ocorrer mediante o uso da comunicação via *e-mail*. O STJ (2008) lista vários crimes que podem ser caracterizados em *e-mails* [STJ08]. Como exemplo cita-se a realização da troca de mensagens eletrônicas com comentários de forma negativa, sobre raças ou religiões, ou ainda, enviar e/ou trocar fotos de crianças nuas.

Havendo uma denúncia formal sobre prática de crime virtual, os órgãos competentes o tipificam conforme a legislação onde ocorreu a prática. Em conjunto, profissionais (peritos) recolhem as evidências deste crime. As evidências digitais são resquícios dos crimes virtuais cometidos, e podem ser encontradas em diferentes formatos, como por exemplo, arquivos de *log*, banco de dados, ou mesmo em repositórios de *e-mail*. E ainda, os peritos realizam um processo de análise (perícia) nas evidências recolhidas, para que seja comprovada a ligação do crime cometido com o acusado [NOR98].

Conforme definição de PAULO (2004), o ato de periciar envolve o exame técnico de pessoa nomeada pelo juiz de Direito, por indicação ou consentimento das partes, para averiguar uma coisa ou um fato objeto de litígio e a ele relacionado [PAU04].

Os resultados (laudo pericial) da análise de uma perícia podem ser utilizados durante um julgamento, haja vista o artigo 436 do Código Processual Civil²⁰:

²⁰ Lei nº 5.869, de 11 de janeiro de 1973, disponível em <http://www.planalto.gov.br/CCIVIL/Leis/L5869.htm>, acesso em 21 de janeiro de 2011.

Art. 436. O juiz não está adstrito ao laudo pericial, podendo formar a sua convicção com outros elementos ou fatos provados nos autos.

No Direito Brasileiro, evidências digitais de um crime encontradas em uma base de *e-mails* por um perito, poderão ser utilizadas como prova válida em um processo Judicial [SOA08]. Ademais, uma mensagem eletrônica enviada por um usuário trafega por diversos equipamentos eletrônicos até sua chegada ao dispositivo destino. É possível que a mensagem possa ser interceptada, excluída, copiada, ou informações nela contidas sejam alteradas em algum desses dispositivos [SOA08]. Soares (2008) cita algumas situações:

- Uma conta de *e-mail* pode ser administrada por um usuário ilegal, desde que este possua informações de acesso a esta conta, como por exemplo, a senha;
- Envio de mensagens utilizando um endereço de *e-mail* adulterado. Muitos servidores de *e-mail* não realizam autenticação para envio de *e-mails*. Com isso, um usuário mal intencionado pode encaminhar *e-mails* com maior facilidade.

Ao analisar estas duas situações apresentadas, conclui-se que a troca de mensagens eletrônicas é uma forma frágil de comunicação. Esta comunicação pode possuir problemas de segurança em seu contexto. De acordo com KUROSE (2010) soluções para tornar a transmissão de *e-mails* mais segura são a criptografia e a certificação digital, a saber [KUR10]:

- Na criptografia uma mensagem é codificada pelo remetente, sendo decodificada pelo usuário destinatário. Objetiva-se a autenticidade e a veracidade [STA08];
- Na certificação digital existe a participação de uma terceira entidade, que assegura a procedência e a integridade de um *e-mail* [STA08].

Diante dessas informações é entendível o questionamento da aceitação do *e-mail* como prova válida perante a legislação brasileira. Entretanto quando esta tecnologia de comunicação é utilizada respeitando-se determinados critérios, torna-se regulamentada a comprovação da veracidade de um *e-mail*. Os critérios devem obedecer a aspectos de

segurança, regulamentados pela Medida Provisória (MP) 2.200-2/01²¹, com a qual foi instituída a infra-estrutura de chaves públicas Brasileira (ICP-Brasil).

Na investigação de crimes virtuais, a área de Computação Forense utiliza técnicas científicas ou tecnológicas para tratamento das evidências digitais recolhidas após um crime [SHI08] [PHI09]. PHILIPP et al. (2009) cita que peritos seguem procedimentos durante todas as etapas deste longo processo, os quais envolvem o tratamento das evidências digitais. Cita-se como exemplo que peritos devem tomar precauções para garantir a integridade das evidências digitais coletadas durante o procedimento pericial. Assim evita-se qualquer tipo de alegação contra o trabalho resultante [PHI09].

Cita OLIVEIRA (2002), na ausência de um método especial, o laudo pericial pode ser contestado na ação ajuizada, podendo fazer com que o processo seja prejudicado ou anulado [OLI02]. Entre alguns métodos, cita-se a necessidade da realização de uma cópia exata da evidência digital recolhida na cena de um crime virtual, processo chamado de “imagem” [SHI08] [PHI09]. Esta cópia será usada durante a análise da evidência. Neste procedimento objetiva-se a manutenção da integridade das evidências digitais coletadas.

Como constatação, REITH (2002) confirma que na investigação de evidências digitais do cometimento de crimes, os procedimentos para análise forense podem não ser totalmente consistentes e nem padronizados. Ainda, o mesmo autor, afirma que muitos métodos têm sido desenvolvidos, porém com enfoque muito técnico [REI02].

Com isto, durante a investigação de evidências, pode-se acarretar muitos casos de falha nos procedimentos, como por exemplo, esquecimento da etapa de isolamento da área antes de ser realizada a coleta de evidências [REI02].

Fez-se desta forma um estudo do *e-mail*, seu funcionamento, formatação, características, e por fim, a criminalidade virtual, focando nos crimes cometidos por *e-mail*. Na próxima Seção relata-se o estudo de grafos direcionados, método utilizado pelo mecanismo na apresentação das provas digitais.

²¹ MP nº 2.200-2, disponível em <http://www6.senado.gov.br/legislacao/ListaPublicacoes.action?id=233404>, acesso em 21 de janeiro de 2011.

2.4. Grafos Direcionados

Define-se grafos como o conjunto de pontos que podem estar interligados [PRE00]. Entre suas divisões encontram-se os grafos direcionados, objeto de aplicação para o mecanismo utilizado para produção de provas digitais, e que se explana neste tópico.

Grafos direcionados, dirigidos ou também conhecidos como dígrafos são representados por $G = (V, E)$. Um conjunto de vértices V e arestas E . Um grafo é formado por vértices interligados por arestas [PRE00].

Apresentam-se algumas das características deste tipo de grafo e que podem ser visualizadas no grafo G da Figura 2.2:

- Conjunto de vértices e arestas não vazios;
- Arestas possuem direção;
- Possibilidade da existência de interligação de vértices de duas formas: até duas arestas em direções contrárias;
- É possível que uma aresta realize a ligação de um vértice a ele próprio;
- Vértices ou arestas podem receber um peso.

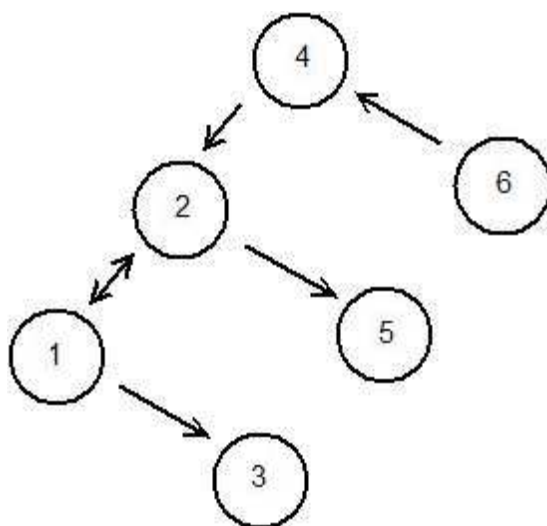


Figura 2.2 – Exemplo de Grafo G

Um grafo é representado pelo conjunto de vértices V e um conjunto de arestas E . O grafo G da Figura 2.2 é representado pelo conjunto $V = \{1,2,3,4,5,6\}$, sendo o conjunto de arestas $E = \{ (1,3), (1,2), (2,5), (4,2), (6,4), (2,1) \}$. Também se pode representar o grafo

através de uma matriz de adjacências, onde a indicação do valor 0 implica na não ocorrência de aresta entre dois vértices, e 1 o contrário. Apresenta-se na Tabela 2.1 o resultado final da matriz de adjacência referente à Figura 2.2.

Tabela 2.1 – Representação da Matriz de Adjacências

	1	2	3	4	5	6
1	0	1	1	0	0	0
2	1	0	0	0	1	0
3	0	0	0	0	0	0
4	0	1	0	0	0	0
5	0	0	0	0	0	0
6	0	0	0	1	0	0

Vértices do grafo podem corresponder aos usuários participantes de conversações criminosas investigada em *e-mails*, e as arestas, a quantidade de *e-mails* trocados entre os usuários. Assim, na próxima Seção apresentam-se métodos para normalizar as palavras existentes em *e-mails* e assim permitir o agrupamento de conversações criminosas encontradas em uma base de *e-mails*.

2.5. Técnicas de Processamento de Textos

A seguir serão abordados brevemente aspectos teóricos sobre o pré-processamento textual, extração e representação de características textuais, métodos de agrupamento e de classificação aplicáveis em *e-mails*.

2.5.1. Pré-processamento Textual e Extração de Características de Documentos

No texto localizado no corpo e cabeçalho de *e-mails*, identificam-se, entre outros, diversas palavras. Por meio do procedimento de normalização, o qual integra o pré-processamento do conjunto textual de *e-mails*, visa-se facilitar a localização de palavras que caracterizem nos *e-mails* algum tipo de ação criminosa [TAM08]. Algumas técnicas utilizadas neste processo são [NET00]:

- *Case Foldering*: Todas as palavras que representam o texto de um *e-mail* são convertidas para sua forma minúscula;

- *Stop-Words*: Retiram-se do texto do *e-mail* palavras de pouco valor semântico, como por exemplo, artigos e preposições;
- *N-grams*: Cada palavra do texto do *e-mail* é dividida em pequenas partes. Cita-se como exemplo a palavra “azul”, que ao se utilizar a técnica de *tri-grams*, têm-se a possibilidade da criação de termos com até 3 letras oriundas desta palavra, gerando-se todas as possibilidades de subsequências a partir da palavra em questão: “_az”, “azu”, “zul” e “ul_”.

Assim, as palavras mais importantes de um *e-mail* são identificadas, sendo estas palavras consideradas as características (atributos) dos *e-mails*. Cada *e-mail* é formado por um conjunto de atributos. Pode-se realizar uma representação de cada *e-mail* por seus atributos [SEB02]. Mediante análise de um *e-mail* representado pelo vetor D (Equação 2.1) verifica-se que cada atributo é representado por x_i , definindo, assim, n atributos, sendo que $i = 1, 2, \dots, n$. Desta forma, n equivale a quantidade total de atributos escolhidos para representação dos *e-mails*. A quantidade de atributos para representação de uma base de *e-mails* deverá ser o mesmo, fazendo com que cada *e-mail* de uma base tenha a mesma quantidade de atributos.

$$\vec{D} = \{x_1, x_2, x_3, \dots, x_n\} \quad (2.1)$$

Pode-se assim representar os *e-mails* de uma base. Na próxima Subseção são apresentados os métodos para execução da representação dos *e-mails* utilizando seus atributos.

2.5.2. Representação dos Atributos Textuais

Feito o pré-processamento do texto dos *e-mails* a ser investigado e definido o conjunto de atributos relevantes para sua caracterização, pode-se obter a representação dos atributos em vetores, matrizes, sequências (*codebooks*), contagem e verificação (*assertions*) e, finalmente, grafos [TRI96]. Dentre as técnicas utilizadas para representação dos atributos dos *e-mails*, destacam-se: *tf-idf* (*term frequency–inverse document frequency*) [TAM08] [HUA08] [CSE06] [VIE06] e *bag of words* [BAL08a] [BAL08b]. Resumidamente, tais técnicas, referem-se à:

- Tf-idf: Técnica utilizada como medida estatística para avaliar o quanto é importante os atributos do texto de um *e-mail*, atribuindo um peso a cada atributo. Esta medida na verdade é a representação da proporção de vezes que um atributo ocorre em um texto do *e-mail* em relação ao número de *e-mails* em que o mesmo atributo foi localizado no conjunto textual (base de *e-mails*) [SAL88], tal qual definido pela Equação 2.2:

$$tf - idf = \frac{Freq_{atributo.texto}}{DocFreq_{atributo}} \quad (2.2)$$

- *Bag of words*: Nesta técnica, os atributos dos *e-mails* são representados como uma coleção de palavras. Assim, representam-se todos os atributos de um *e-mail* em um único vetor. O vetor apresenta as quantidades de repetições dos atributos existentes no conjunto textual formado pelos *e-mails*.

A seguir apresenta-se o método de agrupamento, o qual quando aplicado em textos de *e-mails*, tem por objetivo agrupar conjuntos de *e-mails* que possuam atributos semelhantes, ou seja, pertencentes a uma mesma conversação.

2.5.3. Agrupamento

No método de agrupamento objetos são agrupados de acordo com seus atributos. Na Figura 2.3 apresenta-se o resultado da aplicação de um método de agrupamento num conjunto de 20 objetos. Estes objetos foram agrupados em, um (a), dois (b), quatro (c) e seis (d) grupos [KUM06]. Assim, quanto maior a proximidade entre objetos, maior a facilidade com que pertençam a um mesmo grupo. A proximidade entre objetos é dada pela escolha de atributos que os representem. Objetos que possuam características similares podem pertencer ao mesmo grupo. No exemplo da Figura 2.3, para facilitar a visualização dos agrupamentos resultantes, os objetos de um mesmo grupo são representados pelo mesmo símbolo. E ainda, para identificação do melhor grupo formado, existem técnicas, tais como a taxa de categorização correta [KUM06].

Define-se como a aplicação de métodos de agrupamento a exploração de semelhanças, reunindo D objetos (*e-mails*) padrões. Faz-se com que cada grupo possua *e-mails* mais semelhantes do que outros grupos [JAI99], sendo considerada uma abordagem para

classificação não supervisionada, pois não existe o treinamento para a identificação dos agrupamentos resultantes [JAI99].

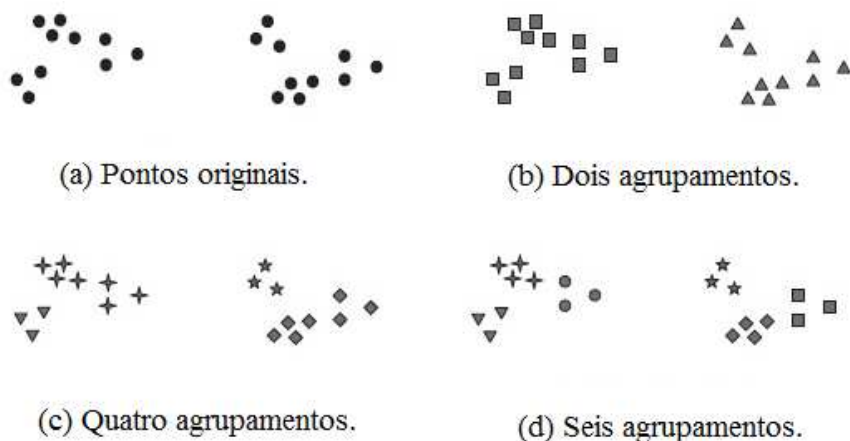


Figura 2.3 – Objetos Agrupados [KUM06]

JAIN et al. (1999) descreve os passos referentes a tarefa de um agrupamento [JAI99], como sendo: Representação, Definição de uma Medida de Proximidade, Agrupamento, Abstração dos Dados e Validação. Assim, explica o autor cada uma das tarefas:

- Representação: Inicialmente é realizada a preparação dos padrões. Definição dos padrões que representem um conjunto de *e-mails*, por exemplo. Cita-se a extração de características;
- Definição de um Padrão de Proximidade: Função que define a distância entre dois padrões;
- Agrupamento: O algoritmo de agrupamento escolhido pelo usuário agrupa os *e-mails* que possuam características similares;
- Abstração dos Dados: Descrição da formação de cada grupo;
- Validação: Métodos para validação dos grupos gerados.

O agrupamento, terceiro passo entre as tarefas apresentadas, é realizado por um algoritmo específico de agrupamento, o qual se utiliza de uma métrica [JAI99]. Esta métrica pode ser estabelecida pelo usuário durante o segundo passo das tarefas, ou seja, a definição de um padrão de proximidade. Através de cálculos de similaridade, o algoritmo analisa o quão

similares são dois *e-mails*, para então realizar a reunião de *e-mails*. Entre os métodos existentes cita-se:

- Distância *Euclidiana*: Corresponde à distância entre 2 *e-mails* [JAI99]. A distância *Euclidiana* é utilizada para avaliar a proximidade entre *e-mails* que aqui são representados por X e Y . X e Y possuem 2 atributos cada. X_1 representa um atributo que está no vetor de atributos do *e-mail* X , e Y_1 representa outro atributo que esta na mesma sequência, entretanto, do vetor de atributos do *e-mail* Y . Da mesma forma ocorre para X_2 e Y_2 . Visualiza-se na Equação 2.3 a distância *Euclidiana*;

$$Euclidiana = ((|X_1 - X_2| + |Y_1 - Y_2|)^2)^{\frac{1}{2}} \quad (2.3)$$

- Distância de *Manhattan*: Calcula a distância absoluta entre atributos de 2 *e-mails* [YAT99]. Na Equação 2.4 observa-se o cálculo de *Manhattan* entre 2 *e-mails* (X e Y), sendo a descrição de X_1 , X_2 , Y_1 e Y_2 realizada na explicação da Equação 2.3.

$$Manhattan = |X_1 - X_2| + |Y_1 - Y_2| \quad (2.4)$$

- Similaridade do *Cossine*: Função utilizada para calcular o grau de similaridade entre o texto de dois *e-mails* que estejam representados por vetores [YAT99]. O valor deste cálculo representa o ângulo entre os vetores (similaridade do *cossine*), sendo uma das medidas mais populares utilizadas para cálculo da similaridade entre documentos textos [YAT99]. Na Equação 2.5 apresenta-se a fórmula para cálculo da similaridade do *cossine* entre dois *e-mails* representados pelos vetores t_a e t_b ;

$$Cossine = \frac{\vec{t}_a \cdot \vec{t}_b}{|\vec{t}_a| \cdot |\vec{t}_b|} \quad (2.5)$$

O resultado da similaridade entre os dois *e-mails* não pode ser negativo, e deve estar entre os valores [0,1]. Quanto mais próximo do valor 1, mais idêntico são os dois *e-mails* comparados;

- Similaridade de *Jaccard*: Utilizado para calcular a similaridade entre 2 *e-mails* representados na utilização de vetores (X e Y) [HAM89]. Na Equação 2.6 observa-se o cálculo da similaridade de *Jaccard*, onde o resultado da similaridade entre X e Y é expresso por $S_{(X,Y)}$, sendo $coc(X,Y)$ a intersecção entre os atributos de 2 vetores comparados, e divididos pela união desses 2 vetores, dado por $(cit(X) + cit(Y) - coc(X,Y))$. O resultado desta similaridade demonstra a proximidade ou não entre dois *e-mails*, sendo que o resultado será 0 (*e-mails* idênticos) ou 1 (*e-mails* diferentes).

$$Jaccard = \frac{coc(X,Y)}{cit(X) + cit(Y) - coc(X,Y)} \quad (2.6)$$

Os algoritmos utilizados com frequência para formação de agrupamentos dividem-se em técnicas Hierárquicas e Particional [JAI99], como apresentado na Figura 2.4. Os algoritmos Hierárquicos podem fazer com que grupos tenham subgrupos. Dois grupos reunidos produzem um grupo do nível superior. Algoritmos Hierárquicos dividem-se em Aglomerativo e Divisivo. No Aglomerativo, cada *e-mail* é representado por um grupo. A cada passo, unem-se os pares de grupos mais próximos, até que todos os *e-mails* sejam reunidos em um grupo [JAI99]. Este processo é finalizado quando existir apenas 1 ou n grupos.

Contrário ao Aglomerativo, os métodos Divisivos iniciam o agrupamento com apenas um grupo que contém todos os *e-mails*. Entretanto grupos se originam, e com isso o grupo originário passa a ter menos *e-mails*. Este processo continua até o momento em que cada grupo tenha um *e-mail* [JAI99]. Entretanto este processo deverá ser auxiliado, já que existe a necessidade de informar ao algoritmo de agrupamento, qual o grupo a sofrer nova divisão.

No método Particional deve ser informado o número exato de agrupamentos resultantes, sendo que os grupos não se sobrepõem. Os *e-mails* são agrupados em um dos grupos informados [JAI99].

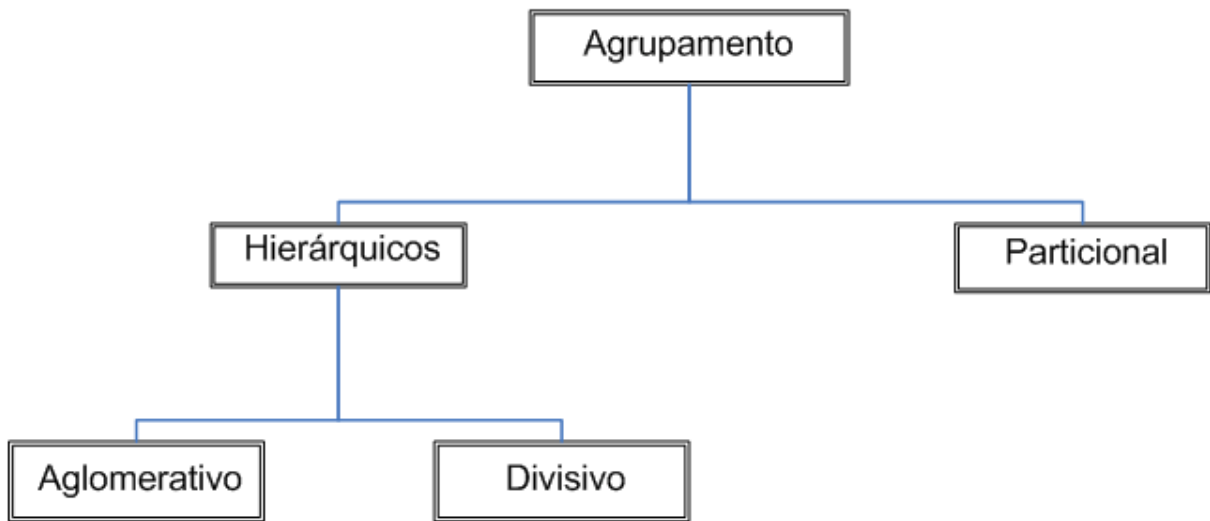


Figura 2.4 - Classificação dos Tipos de Agrupamento

Existem vários algoritmos para realização de agrupamentos. Dente eles, o mais comumente utilizado é o algoritmo *K-means* [JAI99]. Considerado um algoritmo Particional.

O algoritmo *K-means* realiza a divisão do conjunto de *e-mails* em *K*-grupos. O valor de *K* deve ser informado na inicialização do algoritmo, por isso é um método supervisionado. Centróides são formados em um hiperplano, e a partir da verificação da similaridade (Medida de Similaridade) de cada *e-mail* a ser designado como um centróide pode-se obter os conjuntos de *e-mails* (grupos).

Os centróides são atualizados do início ao fim do processo de agrupamento, e a medida de similaridade utilizada normalmente é a *Euclidiana* [JAI99], já apresentada anteriormente. JAIN et al. (1999) afirmam ainda que a proximidade entre *e-mail* e centróide agrupa os *e-mails* similares.

No quarto passo das tarefas, é realizada a Abstração dos *e-mails* agrupados. Os resultados devem ser interpretados, ou seja, realiza-se a identificação dos *e-mails* em suas respectivas alocações. Desta forma, cada *e-mail* terá um rótulo, o qual informa a qual grupo pertence. Esta identificação pode ser fornecida pelo algoritmo responsável no processo de agrupamento. E no quinto passo, procede-se a Validação ou Análise dos Grupos formados, sendo que existem vários métodos para validação de um agrupamento. Uma maneira de se realizar tal tipo de validação é verificar a taxa de acertos quando concluído o agrupamento dos *e-mails* conhecidos.

Pode-se também utilizar métodos estatísticos a fim de analisar a qualidade dos grupos formados. Cita JAIN et al. (1999) que dentre as técnicas utilizadas, destaca-se a matriz de similaridade [JAI99]. Esta matriz apresenta valores da similaridade entre dois *e-mails*. Após o preenchimento da matriz, os valores são comparados com uma matriz ideal para o agrupamento realizado.

A seguir apresenta-se a classificação de *e-mails* através de aprendizagem supervisionada, de modo complementar as técnicas de agrupamento. Isto visto que o objetivo final do mecanismo proposto é identificar se as conversações agrupadas possuem indícios de ações criminosas.

2.5.4. Classificação

A classificação consiste em atribuir objetos a classes [TAM05]. Existem várias aplicações que utilizam este conjunto de técnicas consideradas de aprendizagem supervisionada. Na classificação, atribui-se uma classe c a um objeto respectivamente. As classes c pertencem ao conjunto C , composto por k classes, conforme Equação 2.7:

$$C = \{c_1, c_2, \dots, c_k\} \quad (2.7)$$

Para atribuição de um objeto em uma dada classe, é necessária a existência de objetos treinados, os quais podem ser representados pelo conjunto T . No caso do mecanismo, os *e-mails* treinados representados por seus atributos (D), em uma quantidade q , vide Equação 2.8:

$$T = \{D_1, D_2, D_3, \dots, D_q\} \quad (2.8)$$

O classificador treinado pode realizar a classificação de novos objetos automaticamente. Um processo de classificação é apresentado na Figura 2.5, onde os objetos pertencentes ao conjunto de teste (Base de Teste) são classificados mediante comparação aos objetos pertencentes ao conjunto de treinamento (Base de Treinamento). O conjunto de treinamento são objetos que possuem uma rotulação pré-determinada (*No* ou *Yes*), a qual foi atribuída por um usuário, e através do Algoritmo de Aprendizagem (AA) foi gerado o Modelo.

Este treinamento é realizado pelo AA mediante a Indução de uma quantidade de amostras. Tendo-se o Modelo, possibilita-se ao classificador aplicar o Modelo em informações não treinadas, ou seja, para o exemplo da Figura 2.5, as informações de teste (Base de Teste). Pode-se assim a partir deste momento, classificar cada novo registro em uma das classes propostas de forma automática.

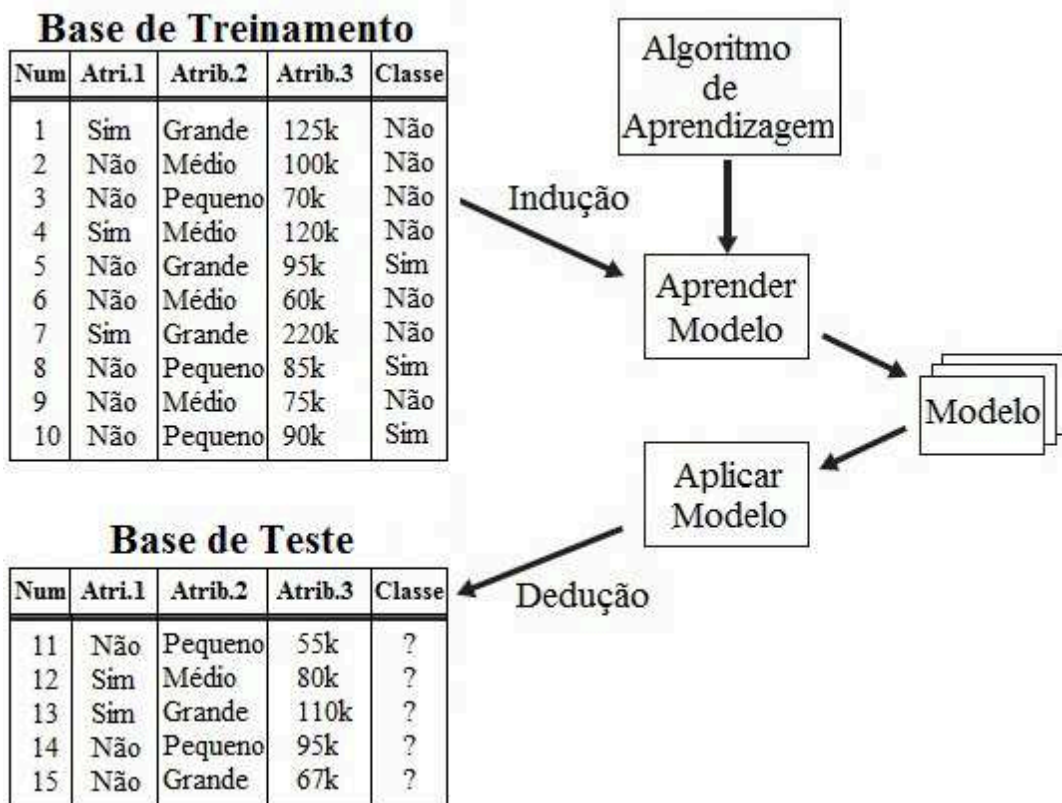


Figura 2.5 – Modelo de Classificação [TAM05]

No caso de *e-mails*, um classificador poderá automaticamente associar *e-mails* em classes, por exemplo. Considerando que um *e-mail* possa ser definido como pertencente a uma classe (*single-label*) ou quando se classifica em mais de uma classe (*multi-label*).

Conforme TAM et al. (2005), existem várias técnicas de classificação. Cada uma utiliza o algoritmo de aprendizagem que melhor se adapta as necessidades do problema [TAM05].

Os algoritmos de aprendizagem precisam gerar modelos que consigam classificar com a maior exatidão possível um conjunto de *e-mails* desconhecidos em classes pré-definidas. TAM et al. (2005) cita como exemplo as técnicas de *Naïve Bayes* (NB), *Árvore de Decisão*

(*Decision Tree* - DT) e Máquinas de Vetor de Suporte (SVM - *Support Vector Machine*) [TAM05]. Tais algoritmos serão detalhados na sequência, já que são os algoritmos mais utilizados e estudados na classificação automática de documentos texto, conforme a literatura consultada durante este trabalho.

Para avaliação do desempenho do classificador, podem-se utilizar diferentes métricas, pois necessita-se verificar a precisão da classificação resultante. Nesta métrica, a avaliação do desempenho de um modelo de classificação baseia-se na taxa de categorização correta e incorreta do método. Esses valores são anotados e apresentados em uma Matriz de Confusão (*Confusion Matrix*). Assim é possível realizar a comparação entre diferentes classificadores na solução de um mesmo problema.

Cita-se como exemplo a Tabela 2.2, onde é apresentado a Matriz de Confusão para um problema de classificação em dois grupos: Classe 1 e Classe 0. Na coluna Classe 1 é apresentado na primeira linha, a quantidade de registros classificados corretamente para a Classe 1, e na segunda linha, o número de registros da Classe 0 classificados incorretamente como pertencentes a Classe 1. Da mesma forma ocorre na segunda coluna, onde é demonstrado os registros corretamente e incorretamente classificados como pertencentes a Classe 0.

Tabela 2.2 – Representação da Matriz de Confusão

Classe 1	Classe 0
f_{11}	f_{10}
f_{01}	f_{00}

Desta forma, pode-se exemplificar o preenchimento da Matriz de Confusão mediante a classificação de 20 *e-mails* em duas classes: Classe 1 e Classe 0. Imaginando-se a existência de 12 *e-mails* pertencentes a Classe 0 e outros 8, como Classe 1. No processo de classificação, o método classificador consegue identificar 1 *e-mail* (f_{10}) como não sendo da Classe 0, e 2 *e-mails* (f_{01}) como não sendo da Classe 1. Assim, resulta-se em 11 *e-mails* (f_{11}) e 6 *e-mails* (f_{00}).

Através dos valores de acertos e erros armazenados na Matriz de Confusão, torna-se possível a utilização de métricas para avaliação do modelo de classificação utilizado. E ainda, a fim de sintetizar e facilitar a comparação entre o resultado da aplicação de diferentes

modelos aplicados calcula-se a Exatidão e a Taxa de Erro da aplicação de um modelo de Classificação, conforme Equação 2.9 e 2.10 respectivamente:

$$\text{Exatidão} = \frac{\text{Número de Classificações Corretas}}{\text{Número Total de Classificações}} = \frac{f_{11} + f_{00}}{f_{11} + f_{10} + f_{01} + f_{00}} \quad (2.9)$$

$$\text{Taxa de Erro} = \frac{\text{Número de Classificações Incorretas}}{\text{Número Total de Classificações}} = \frac{f_{10} + f_{01}}{f_{11} + f_{10} + f_{01} + f_{00}} \quad (2.10)$$

A seguir explicam-se os métodos de classificação NB, SVM e DT aplicáveis na produção de provas digitais em *e-mails*:

- **NB:** Classificador probabilístico baseado no teorema de *Bayes* [GOD65] [SEB02]. Determina a classe de maior probabilidade para cada novo *e-mail* apresentado ao classificador, baseando-se na probabilidade da existência de determinados atributos encontrados no *e-mail*. Na Equação 2.11 apresenta-se a fórmula do Teorema de *Bayes*.

$$P(c_i | \vec{d}_j) = \frac{P(c_i) P(\vec{d}_j | c_i)}{P(\vec{d}_j)} \quad (2.11)$$

O valor resultante da equação será a probabilidade de um *e-mail* representado por um vetor \vec{d}_j pertença a uma classe c_i . Sendo:

- $P(c_i)$ a probabilidade que um dado *e-mail* pertença a uma classe;
- $P(\vec{d}_j | c_i)$ a probabilidade de um *e-mail* dada a classe;
- E por fim, $P(\vec{d}_j)$, a probabilidade da escolha aleatória de um *e-mail* que tenha o vetor \vec{d}_j como sua representação.

Para classificar um novo *e-mail*, o algoritmo determina a classe mais provável, conforme os atributos que descrevem esse *e-mail*. Para realizar o cálculo da classe de maior probabilidade, o classificador NB utiliza-se da Equação 2.12:

$$V_{NB} = \arg \max_{v_j \in V} P(v_j) \prod_i P(a_i | v_j) \quad (2.12)$$

O resultado da Equação 2.12 define a maior probabilidade para cada *e-mail* dada por $P(v_j)$. Sendo a probabilidade de um *e-mail* pertencer a uma classe, mediante o produtório de $P(a_i|v_j)$. Por fim, tem-se a probabilidade deste *e-mail* pertencer à determinada classe;

- SVM: Algoritmo classificador [VAP63]. Pode utilizar conceitos de separação de dados lineares (single-label) ou não lineares (multi-label) [LOR07]. Baseia-se em encontrar um hiperplano ótimo que divida linearmente dois conjuntos resultantes possibilitando a melhor classificação, e por isso, SVM é considerado um classificador estatístico [LOR07].

Existe uma distância (margem) entre o hiperplano e determinados documentos de cada classe. Os documentos mais próximos do hiperplano são considerados os vetores de suporte, e determinam a atual localização do hiperplano.

E ainda, LORENA et al. (2007) cita utilizações de SVM em aplicações bem sucedidas, dentre elas, a categorização de textos [LOR07]. Da mesma forma RODRIGUES (2009), onde cita que diversas aplicações para classificação de texto que utilizam SVM podem ser encontradas atualmente, como por exemplo, a classificação de *e-mails*;

- **DT:** Técnica de classificação simples [TAM05], baseado em método estatístico que utiliza treinamento supervisionado para classificação e previsão de dados. Várias possibilidades são organizadas na forma de uma árvore de decisão hierarquicamente estruturada. Visualiza-se na Figura 2.6 um exemplo de DT.

Uma DT é composta por [TAM05]:

- Um nodo principal (Temperatura do Corpo): nodo que não possui entradas;

- Nodos internos (Gerar): tem exatamente uma entrada e uma ou mais saídas;
- Nodos terminais (Mamífero e Não Mamífero): em qual classe pertence determinada instância.

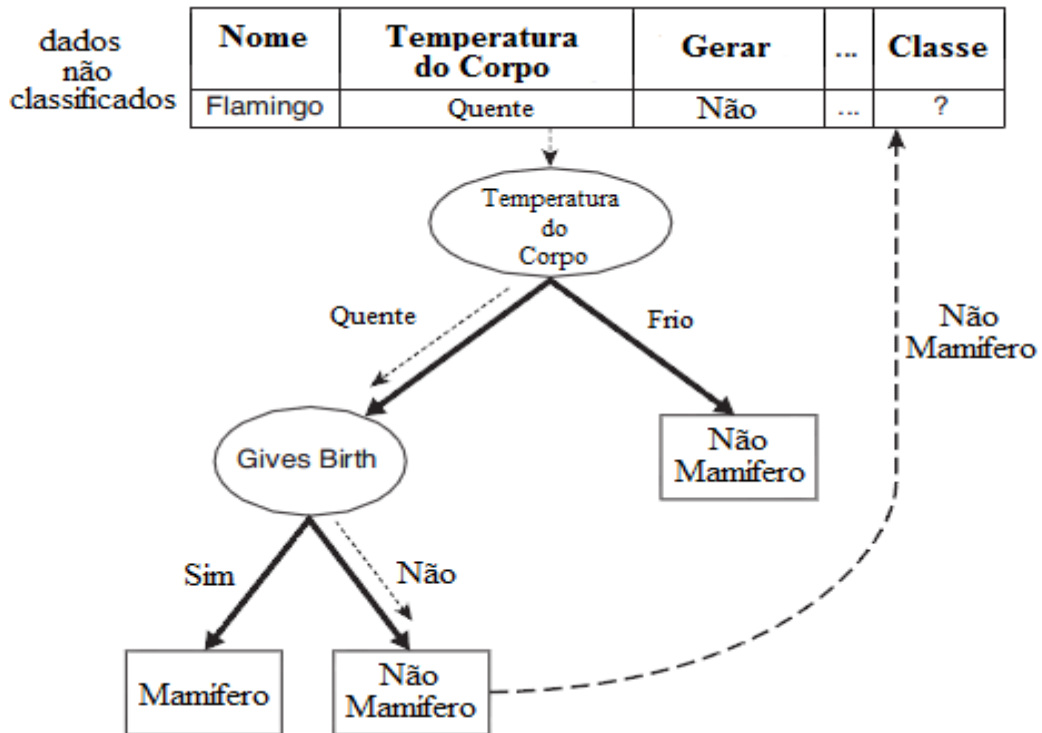


Figura 2.6 – Exemplo do Classificador DT [TAM05]

Uma DT classifica os dados a partir do nodo principal em direção a nodos internos. O nodo principal é a característica que distingue classes. Cada nodo interno especifica determinado teste, podendo direcionar a outro nodo interno, ou então, a um resultado classificatório (nodo terminal). Conforme TAM et al. (2005) várias são as possibilidades para a representação de um problema em formato de DT [TAM05] e por isso algumas DT são mais precisas que outras.

A classificação DT consiste em um conjunto de regras que são aplicadas de maneira sequencial, finalizando com uma decisão, sendo assim, uma DT pode ser representada por um conjunto de regras “se-então”. E também, o classificador DT pode ser utilizado na classificação de *e-mails*, como por exemplo, a existência de

palavras criminosas no corpo de um *e-mail*, fazendo com que o mesmo se classifique em criminoso, havendo um prévio treinamento para esta classificação.

Neste Capítulo foram expostos estudos para realização do trabalho de Mestrado. Entendimento da ferramenta de comunicação *e-mail*, sua definição e formato técnico são explorados. Criminalidade resultante da utilização indevida dos serviços virtuais, dentre eles, no *e-mail* e o estudo de técnicas de processamento de textos, grafos, agrupamento e classificadores. No próximo Capítulo, relatam-se estudos encontrados na literatura técnica, os quais permitem compreender trabalhos relacionados ao tema aqui estudado.

Capítulo 3

Trabalhos Relacionados

O estudo dos Trabalhos Relacionados permitiu a análise de vários trabalhos que visam automatizar a investigação criminal em ferramenta *e-mail*. Apresentam-se na sequência, trabalhos envolvendo o processamento automático de *e-mails*, a visualização do histórico de conversações, estudos de métodos para investigação criminal em *e-mails*, assim como as considerações finais relacionadas a todos os trabalhos estudados e apresentados neste Capítulo.

3.1. Processamento Automático de *e-mails*

No estudo de [NAG10] é proposto um modelo que utiliza o algoritmo de agrupamento *k-means* (similaridade do *Cossine*). O modelo calcula a similaridade entre *e-mails*, considerando como atributos o texto do cabeçalho (campos assunto, destinatário – *From* e anexos) e o texto do corpo dos *e-mails*. Comparam-se os atributos dos *e-mails* individualmente, como pode ser observado na Figura 3.1, onde se apresenta o modelo proposto pelos autores [NAG10]. Os *e-mails* similares são agrupados através de várias comparações, e a comparação finaliza quando todos os *e-mails* similares estiverem agrupados.

Este modelo [NAG10] foi testado para agrupamento de 310 *e-mails* existentes na base da *Enron Corpus* através da implementação de um *software* em linguagem *Java*. Assim, obteve uma precisão de 78% na utilização da técnica *10-fold Cross Validation* [RAB93].

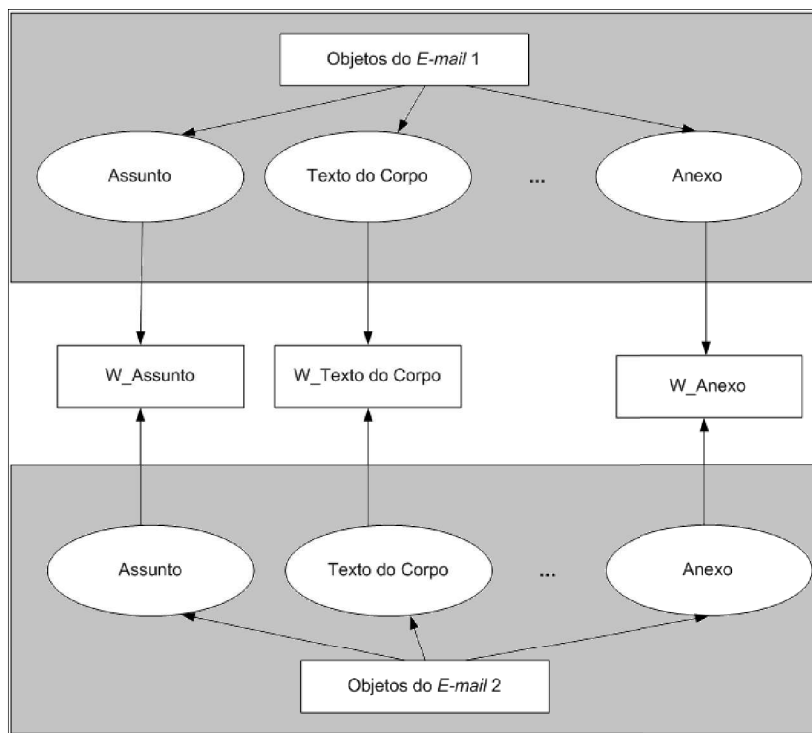


Figura 3.1 – Modelo Proposto [NAG10]

Citam os autores em suas conclusões [NAG10], que o modelo proposto pode ser utilizado para classificação dos *e-mails* de um usuário em pastas, por exemplo.

Ademais, na sequência, estudos realizados por CSELLE (2006) apresentam categorias para armazenamento e classificação automática de *e-mails*:

- Agrupamento: Pode utilizar sofisticados métodos de agrupamento mediante a utilização de vetores representados pela técnica *tf-idf*, o qual descreve a significância de palavras encontradas no texto de um *e-mail*. Os métodos são utilizados para agrupamento de mensagens por tópicos. A vantagem da utilização de agrupamento é que usuários não precisam criar pastas durante o processo. Todo o processo de classificação e armazenamento da mensagem é automático;
- Extração de conteúdo específico dos *e-mails*, como a detecção da assinatura do usuário, ou outra informação;
- Classificação por priorização. Um método manual é utilizado após a primeira triagem (recebimento dos *e-mails* pelo usuário);

- Criação de modelos através do histórico de conversações entre usuários. Realiza-se a análise de *e-mails* enviados e recebidos;
- Utilização de padrões temporais. Um artefato é alimentado por informações dos *e-mails* recebidos pelo usuário, prevendo o possível comportamento de um usuário.

E ainda, CSELLE (2006) expõe um grupo com três classificadores automáticos [CSE06]:

- Classificadores *Bayesianos* (NB): Algoritmos que utilizam um conjunto de palavras provenientes de um documento, como por exemplo, um dicionário. E baseado em regras estipuladas pelo desenvolvedor, o algoritmo analisa probabilidades baseado no encontro das palavras existentes em um dicionário e as palavras de *e-mails*;
- Classificadores *Ripper* (baseado em regras): Algoritmos que criam regras baseado na aprendizagem de máquina. Tem por base a utilização de informações recebidas por uma heurística;
- Classificadores fundamentados no Remetente: Algoritmos classificam mensagens baseando-se apenas no remetente desta.

No trabalho realizado por [CSE06] foi aplicado algoritmos de classificação e agrupamento. Este trabalho foi implementado em um *software* denominado de *BuzzTrack*, o qual pode ser visto na Figura 3.2. Autor utilizou códigos de programação *JavaScript* e *Python*. *BuzzTrack* realiza a análise de *e-mails* provenientes do *software Thunderbird* 1.5, sendo que este *software* organiza uma base de *e-mails* por assuntos.

Para testes, [CSE06] objetivou organizar 3 bases de *e-mails*. O autor testou 1.346 *e-mails* pré-organizados em 76 assuntos. Os resultados experimentais apresentaram que: o classificador NB obteve uma taxa de acertos de 55,1%; o classificador baseado em regras, 52,0%; e o classificador fundamentado no remetente da mensagem, 47,4%. Além destes classificadores, os autores testaram as bases com um algoritmo de agrupamento (*Single Link Cluster* - tf-idf), alcançando uma taxa média de acertos de 82%.

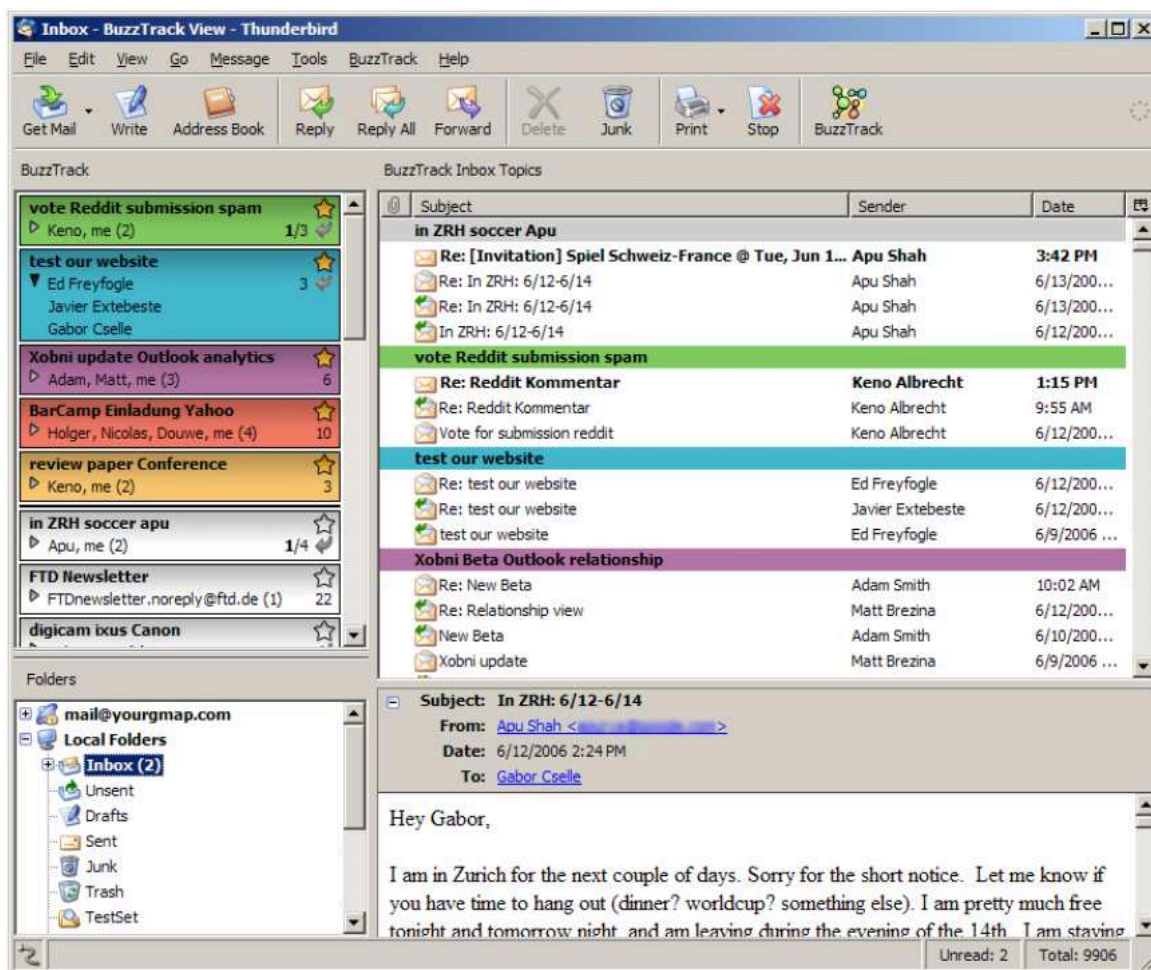


Figura 3.2 – Software BuzzTrack [CSE06]

Existem outros trabalhos para classificação automática de *e-mails*, os quais utilizam uma ampla variedade de métodos. Em BALAMURUGAN et al. (2008a), comparam a aprendizagem do comportamento dos classificadores para detecção de ameaças por *e-mail*, utilizando os métodos de classificação [BAL08a]:

- DT: Algoritmo mais popular de aprendizagem indutivo;
- SVM: Algoritmo considerado o melhor método classificatório com desempenho em categorização de texto. Esta técnica é popular para aprendizagem de máquina;
- NB: Classificador frequentemente utilizado para classificação de *e-mails*. Simples mas altamente efetivo no que se refere ao modelo de aprendizagem.

Para [BAL08a] os resultados dos testes utilizando as palavras do corpo dos *e-mails* mostraram que o algoritmo DT apresentou 97,4%, o algoritmo SVM 95,6% e o NB 95,4% de eficácia na classificação de *e-mail*, concluindo que o classificador DT é uma abordagem promissora para detecção automática em *e-mail*; tendo melhor desempenho que o SVM e sendo menos complexo que o classificador NB. Deste modo, o algoritmo DT torna-se mais atrativo e apresenta maior facilidade para manipulação, além de ser executado eficientemente em grandes bases de *e-mails* e número de recursos.

Em um segundo trabalho, BALAMURUGAN et al. (2008b), os autores propõem o algoritmo *Ad Infinitum*, melhoria do algoritmo DT, dando continuidade ao trabalho desenvolvido (detecção de ameaças em *e-mails*) [BAL08b]. Os autores consideraram este algoritmo com um novo método para classificação de texto, informando que este é mais rápido e fácil de ajustar, e pode controlar um grande grupo de características.

E ainda, [BAL08b] realizaram comparações de *Ad Infinitum* com os algoritmos DT, SVM e NB. Como resultado, relatou-se o tempo desses algoritmos para classificação: NB foi o mais rápido, entretanto o menos eficiente. Em ordem crescente, o mais rápido foi: NB, *Ad Infinitum*, DT e SVM [BAL08b].

Em outro trabalho, na utilização de um classificador *Bayesiano*, DREDZE et al. (2006) apresentam a classificação de *e-mails* por atividades. Os métodos básicos foram o uso do classificador *Bayesiano* e mensagens de resposta para *threads* a fim de determinar em qual atividade deve ser aderido um *e-mail* [DRE06].

Para testes, utilizaram 149 *e-mails* pré-definidos em 27 classes para treinamento em dois classificadores distintos. Estas classes foram utilizadas na classificação automática de uma base composta por 1.146 *e-mails* [DRE06].

Assim, [DRE06] aplicaram no primeiro classificador como métrica de similaridade os contatos (destinatário e remetente) identificados no cabeçalho de *e-mails*. E no outro classificador, foi identificada a similaridade entre o conteúdo do corpo de um *e-mail* teste com os *e-mails* de treinamento, usando uma variação do algoritmo LSI (*Latent Semantic Indexing*) para medir a similaridade. Para ambos os classificadores projetados, quanto maior a similaridade, maior a probabilidade de um *e-mail* que está sendo testado ser classificado corretamente em uma classe.

DREDZE et al. (2006) também realizou a comparação com o classificador NB. Os resultados experimentais demonstraram que os classificadores propostos alcançaram uma taxa

de acerto de 94% dos *e-mails* testados. Entretanto o algoritmo NB não obteve bons resultados, visto que os testes demonstraram que tal algoritmo obteve uma taxa de erros de 50% na classificação dos *e-mails* [DRE06].

Autores como TAM et al. (2008) propuseram a organização de *e-mails* [TAM08]. O processo foi separado em quatro partes, a saber:

- Redução da dimensão do texto;
- Representação das mensagens de forma a que estas possam ser processadas automaticamente;
- Atribuição de um peso às palavras que compõe cada *e-mail* (através da verificação da frequência das palavras, como por exemplo através do uso da técnica *tf-idf*, ou a aplicação da técnica denominada de *Singular Value Decomposition*);
- Agrupamento dos *e-mails* por determinada semelhança (utilização de técnica de aprendizagem não supervisionada, nomeadamente o algoritmo de agrupamento *K-Means*).

A verificação da verdadeira autoria de um *e-mail* sendo investigado pela computação forense foi estudada pelos autores [IQB10]. Neste trabalho foram utilizadas duas bases de *e-mail* para a identificação [IQB10], como pode ser visto na Figura 3.3, onde se apresenta o modelo proposto pelos autores.

Uma possuía *e-mails* do suspeito (*E-mails of Suspect S*), e a outra, um conjunto muito grande de *e-mails* pertencentes a diferentes indivíduos (*Large E-mail Dataset U*). As duas bases são submetidas a técnicas de pré-processamento e extração de características (*Preprocessing and Features Extraction*), onde a base suspeita é transformada em arquivos de atributos de treinamento, e a segunda base, arquivos de atributos de teste [IQB10].

Após, estes dois arquivos são submetidos a métodos de classificação (*Classification Model*) e de regressão (*Regression Model*). O modelo dos autores [IQB10] possibilita na sequência, a comparação de um *e-mail* anônimo (*Anonymous E-mail*) com o modelo gerado (*Validated Class Model* ou *Validated Reg. Model*).

O sistema apresentou a taxa média de erro de 17%. Foram aplicados vários métodos de classificação e de regressão. Nos métodos de classificação, foram aplicados: *Adaboost.M1*,

Discriminative Multinomial Naïve Bayes e Bayesian Network, na base de e-mails da Enron Corpus.

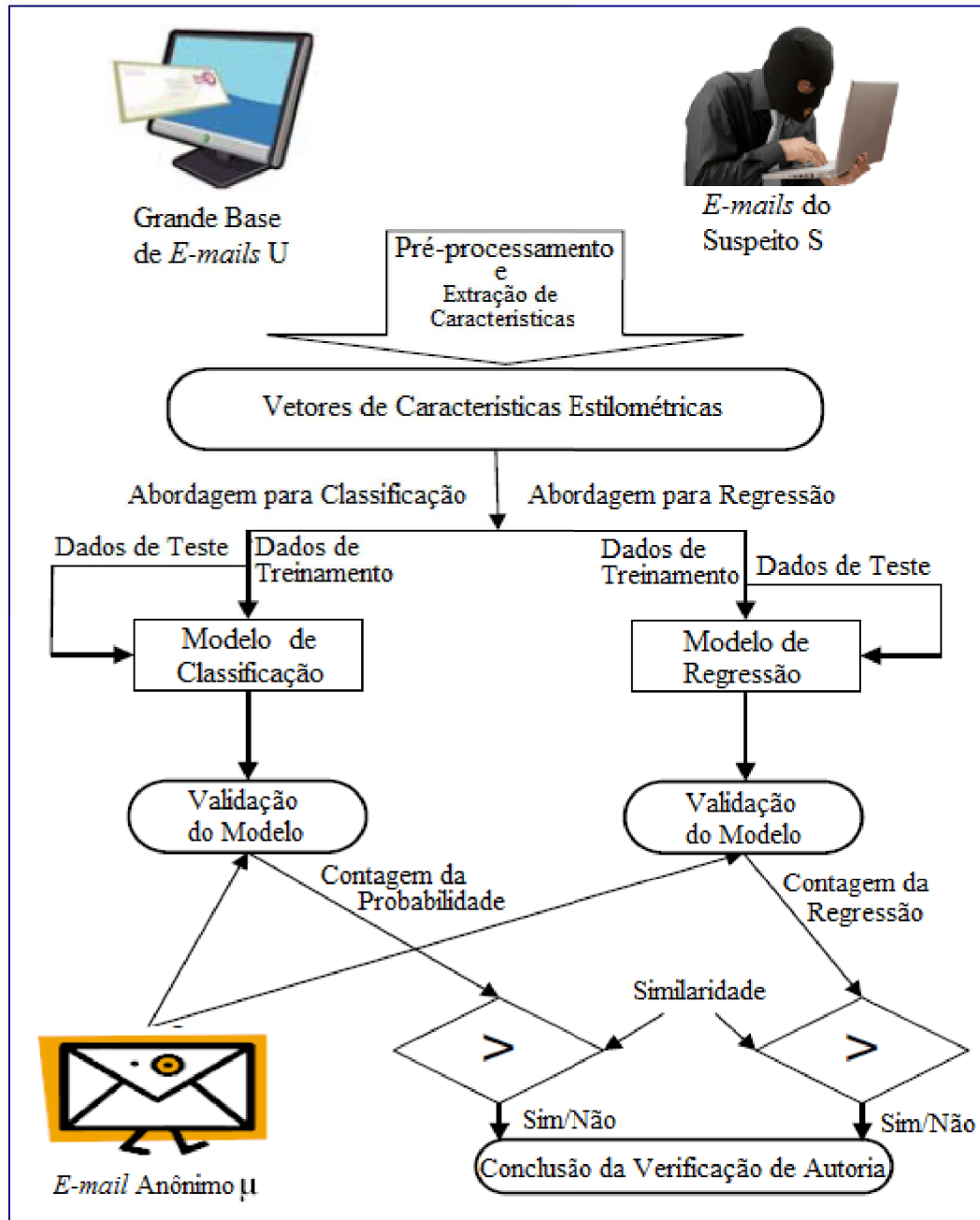


Figura 3.3 – Modelo Proposto [IQB10]

Na próxima Seção relatam-se trabalhos envolvendo a visualização do histórico de conversações identificadas, entre outros, em base de *e-mails*.

3.2. Visualização de Conversações

Em outro foco, VIÉGAS et al. (2006) realizam a visualização do histórico de conversações entre o usuário proprietário de uma base de *e-mails* e um de seus contatos [VIE06]. O contato deve ser escolhido pelo usuário proprietário da base de *e-mails* sendo analisada.

Para visualização do histórico entre o usuário e contato escolhido, realiza-se a contagem das palavras utilizadas na conversação (conjunto de *e-mails* - histórico) escolhida mediante aplicação da técnica tf-idf [VIE06].

VIÉGAS et al. (2006) implementaram o aplicativo *Themail*²², que depende do conteúdo dos *e-mails* (texto do corpo), em lugar do conteúdo do cabeçalho dos *e-mails* [VIE06]. Assim contabilizam-se e se visualizam as palavras mais utilizadas em um relacionamento escolhido.

O aplicativo *Themail* realiza a construção de uma apresentação visual de interações ao longo do tempo. Esta ferramenta pode apresentar os resultados de seu algoritmo, tanto mensalmente, como anualmente.

Desta forma, o aplicativo desenvolvido por [VIE06], possibilita que um usuário detecte em sua base de *e-mails*, as principais palavras detectadas entre as conversações com seus contatos, assim como quais as diferenças existentes entre duas conversações.

De forma análoga, autores PUPYREV et al. (2010) analisam conversações através de visualização gráfica [PUP10]. Consideram conversações aquelas realizadas em tempo real e armazenadas em dispositivo eletrônico, como *chats* ou bases de *e-mail*.

A visualização gráfica tem o objetivo de proporcionar entendimento da evolução de uma conversação [PUP10]. Para identificação de conversações, autores baseiam-se na análise das propriedades estruturais de rede do método analisado. No caso de *e-mails*, os campos do cabeçalho, e em *chats*, o usuário remetente e o destinatário de uma mensagem.

Na Figura 3.4, apresentam-se conversações detectadas em uma base de *e-mails* praticada ao longo de 7 dias pelo usuário proprietário denominado de leonardo.pacheco. Na Figura 3.4, optou-se em visualizar a conversação realizada no dia 13 de Novembro de 2000, apresentando a existência de relacionamento entre o usuário leonardo.pacheco e 12 contatos com quem trocou *e-mails* neste dia.

²² Aplicativo disponível em <http://alumni.media.mit.edu/~fviegas/projects/themal/study/index.htm>, acesso em 29 de janeiro de 2011.

Autores [PUP10] analisaram duas bases: 131.879 mensagens de *chat*, base chamada de VNC e pertencente a uma corporação Russa, e também, analisaram 873.963 *e-mails* da base *Enron Corpus*. Por fim, autores concluem que através do mecanismo estudado, proporcionam o entendimento da evolução das relações sociais entre contatos, além da determinação de padrões, como a frequência do envio ou recebimento de *e-mails* ou mensagens.

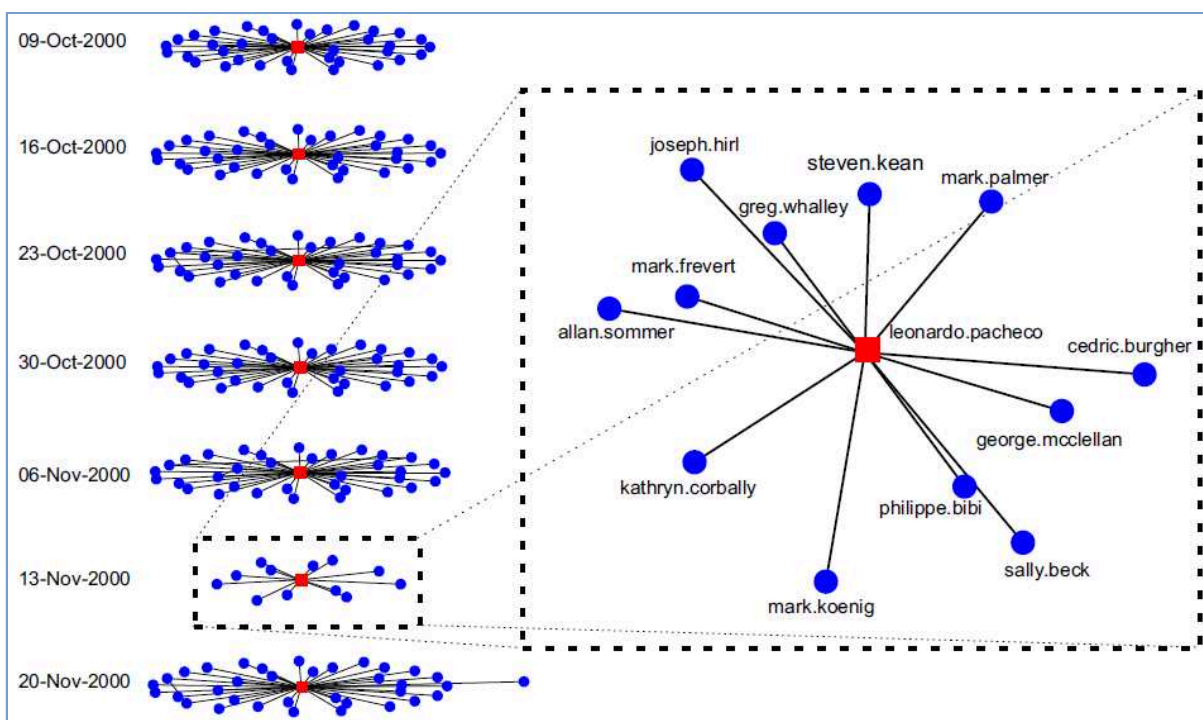


Figura 3.4 – Visualização de uma Conversação [PUP10]

Na próxima Seção relatam-se trabalhos que adicionam propriedades na investigação em *e-mails*. Na primeira, autores preparam um conjunto de palavras (dicionário) e uma base de *e-mails* para investigação de crime de assédio moral registrado em *e-mails*, e no segundo trabalho, apresenta-se uma ferramenta forense comercial utilizada atualmente na investigação de bases de *e-mail*.

3.3. Outros Trabalhos

No estudo de [NUN09], palavras que caracterizam assédio moral no ambiente de trabalho foram coletadas de diversas maneiras, passando inclusive por pré-processamento. No

pré-processamento aplicou-se a técnica de *3-grams* sobre duas bases de dados contendo 25 *e-mails* e 512 palavras de um dicionário de assédio moral.

Assim, na Figura 3.5 apresenta-se o mecanismo utilizado para detecção de assédio moral em *e-mails* [NUN09]. Inicialmente é fornecida ao mecanismo a Entrada de informações: uma coleção de *e-mails*, um arquivo texto que represente um *e-mail*, ou mesmo uma frase. A Entrada (textos) é convertida em vetores de palavras através do processo de Manipulação, após aplicação de pré-processamento na Entrada. Aplica-se na próxima fase o algoritmo *N-gram* no vetor de palavras originado na fase anterior, assim como no dicionário de palavras com registro de crime de assédio moral já formalizado. Assim, na fase de Análise realiza-se a comparação entre os *N-grams* gerados a partir das palavras registradas no vetor e os *N-grams* do dicionário de palavras. Por fim, na fase de Resultados pode-se verificar a existência de crime de assédio moral na Entrada informada.

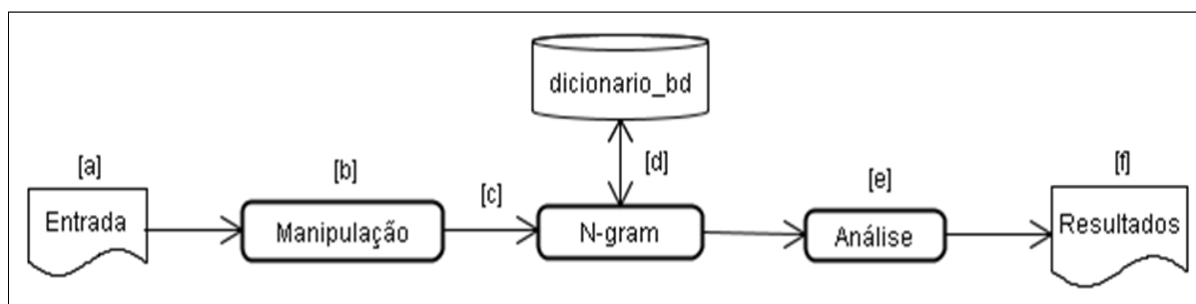


Figura 3.5 – Mecanismo para Detecção de Assédio Moral em *e-mails* [NUN09]

Para efeito de testes, foram comparados todos os *3-grams* de uma palavra do *e-mail* com os *3-grams* de todas as palavras do dicionário, individualmente. Na classificação dos 25 *e-mails*, por similaridade, utilizando o dicionário de assédio moral, os autores conseguiram uma taxa de acerto de 90,91%, com uma taxa de falso positivo igual a 9%. E ainda, atualmente já existe a formalização de 40 *e-mails* e o dicionário é composto por 539 palavras.

Além de trabalhos científicos, verificou-se a existência de *softwares* comerciais para análise forense. Dentre eles, os desenvolvidos pela empresa *Access Data Forensic Toolkit* (FTK) [DAT10]. O *software* FTK pode ser visualizado na Figura 3.6, onde uma base composta por 90 *e-mails* esta em análise, sendo que um dos *e-mails* da base é selecionado, e seu conteúdo apresentado (visualizado).

FTK permite a realização de cópia (“imagem”), análise, recuperação e a formalização das análises realizadas (construção de relatório - laudo pericial) de evidências digitais. Para análise de evidências digitais em formato *e-mail*, FTK é a ferramenta comercial mais rápida [PHI09], sendo que ela proporciona a leitura de diferentes formatos de arquivos de *e-mail* (ex.: PST – *Personal Storage Table* e OST – *Offline Folder File*) e a visualização de características de cada *e-mail* (cita-se a visualização do *e-mail* em formato hexadecimal, texto, ou mesmo seus parâmetros de configuração). Ademais, FTK possibilita a procura por palavras/expressões no conjunto textual formado pelos *e-mails*.

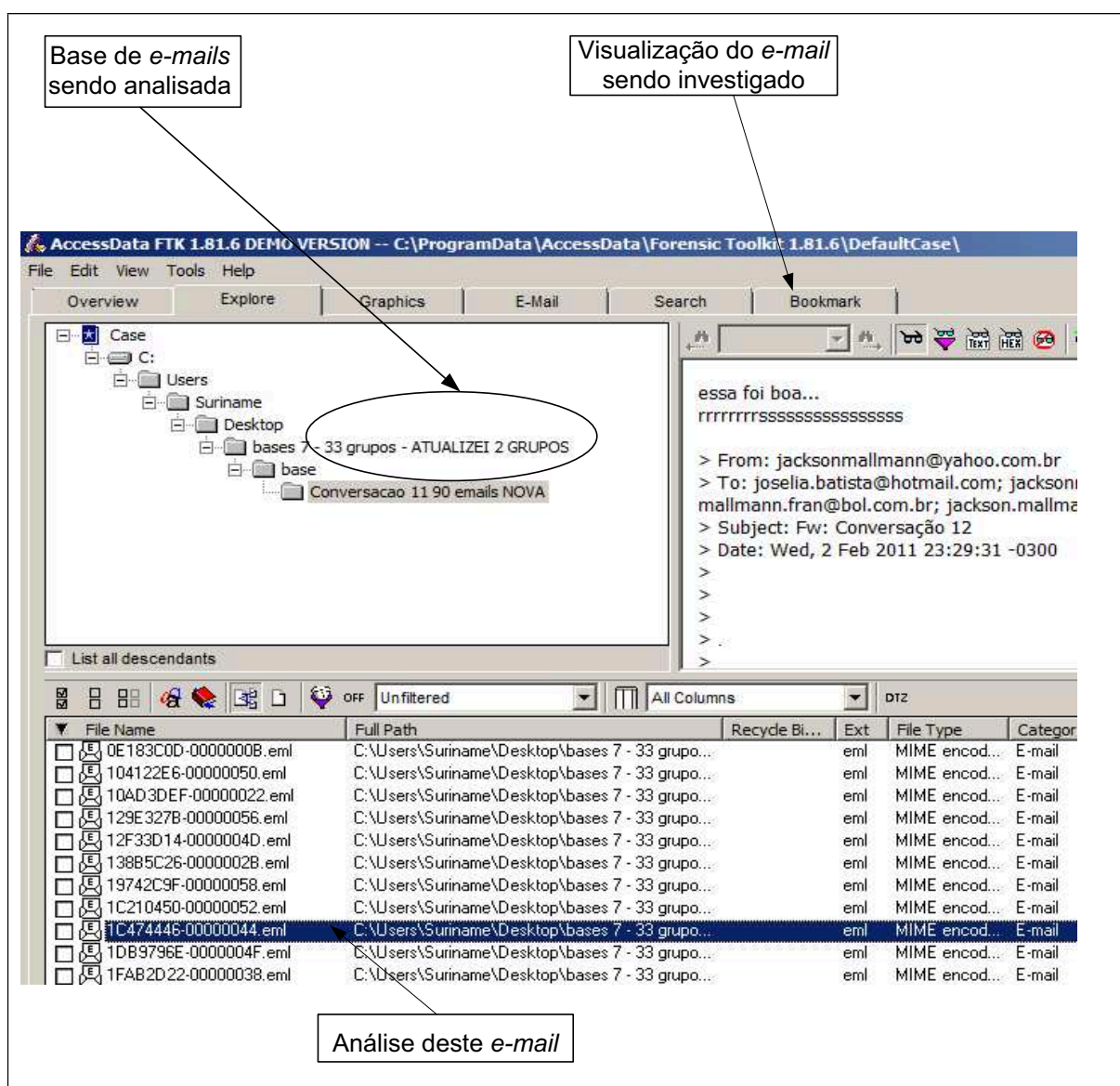


Figura 3.6 – Vista Parcial do *Software* FTK versão 1.81.6 [DAT10]

Realizado este estudo, apresenta-se na próxima Seção considerações do estudo de todos os Trabalhos Relacionados apresentados neste Capítulo.

3.4. Considerações Finais

Os trabalhos descritos neste Capítulo utilizam alguns dos métodos que serviram de base na formulação e na implementação do mecanismo de rastreamento de relacionamentos em *e-mails*, visando na produção de provas digitais. Sendo assim, relata-se que nos trabalhos de [CSE06] e [DRE06] foi utilizado o classificador NB, sendo que este classificador se demonstrou inadequado para realizar agrupamento/classificação de *e-mails*.

Porém, quando se trata de somente duas classes o desempenho do algoritmo NB é bom, o que pode ser constatado no trabalho de [BAL08a] e [BAL08b], no qual foi proposto uma classificação de *e-mails* em duas classes: ameaçador e não ameaçador, sendo que esses autores apresentam comparações entre os classificadores DT, SVM e NB com duas bases de *e-mails*, com 2.099 e 3.893 *e-mails*, respectivamente.

E ainda, nos trabalhos de [TAM08] verifica-se a possibilidade do uso de método de agrupamento para identificação de conversações existentes em base de *e-mails*, e no trabalho de [IQB10] verifica-se o uso de métodos para classificação de *e-mails*.

Nos trabalhos de [VIE06] e [PUP10], usam-se recursos tecnológicos para visualização do histórico, ou conversações, em *e-mails* identificados em uma base de *e-mails*. A idéia destes 2 trabalhos pode ser utilizada pelo mecanismo para produção de provas digitais apresentado neste trabalho de Mestrado para visualização gráfica de conversações criminosas ou não que tenham sido identificadas.

Ademais, nos trabalhos apresentados por [VIE06] e [PUP10], apenas se faz uso dos contatos que enviam ou recebem mensagens (*e-mails* ou *chat*) para a construção de gráficos, ou mesmo a visualização do histórico (*e-mails*) de uma conversação. No mecanismo deste estudo, pode-se utilizar a construção de gráficos, entretanto investigando o conteúdo textual existente em conversações.

E embora haja conflitos entre os resultados apresentados, como por exemplo, resultados do classificador NB entre diferentes autores como [DRE06] e [BAL08a], é visível a utilização dos conceitos relatados neste Capítulo para a produção de provas digitais em *e-mails*.

No estudo do *software* FTK, foi identificado que ele não classifica ou identifica conversações em *e-mails*, e depende da habilidade do usuário (perito) na utilização do *software*. Quanto maior o conhecimento do perito para identificação de um *e-mail* criminoso, mais eficiente se torna a utilização de FTK.

Assim, mediante a aplicação de métodos de agrupamento ou de classificação [NAG10], realizar a identificação de conversações, e na aplicação de métodos de classificação, a possibilidade de caracterização criminal das conversações, sendo estas conversações posteriormente representadas graficamente, similar a [PUP10]. Para caracterização criminal, é notória a possibilidade do uso do dicionário criminoso formalizado por [NUN09].

Enfim, os trabalhos estudados e apresentados neste Capítulo aplicam vários métodos. Na união de alguns destes métodos, viabiliza-se a aplicação do mecanismo de produção de provas digitais. Deste modo, o mecanismo que se apresenta no próximo Capítulo utiliza alguns dos conceitos aqui apresentados e descritos.

Capítulo 4

Mecanismo para Rastreamento de Relacionamentos em *e-mails*

Apresenta-se neste Capítulo especificações do mecanismo. Definição do problema a ser resolvido, ambiente em que ocorre o problema, bases de *e-mails* utilizadas nos experimentos, funcionamento e técnicas empregadas pelo mecanismo.

4.1. Visão Geral

Durante um processo judicial que envolva crime virtual, pode ser solicitado a um perito que as provas digitais sejam analisadas para comprovação do crime [NOR98]. Na análise das provas, pode-se determinar o nexos causal que faça a ligação entre um crime cometido e o acusado deste crime. Através do nexos causal, proporciona-se uma adequada interpretação daqueles que são responsáveis pela tomada de decisões referentes ao processo judicial [DEL00][JES02].

E ainda, em bases de *e-mails* podem ser encontradas provas digitais do cometimento de diferentes crimes virtuais [STJ08]. Sendo que, como já mencionado anteriormente, o repositório de *e-mails* pode estar armazenada em um computador pessoal ou mesmo em um servidor, sendo que para procedimentos periciais o conjunto de *e-mails* pode ser coletado para realização de análise pericial.

Faz-se assim a necessidade da existência de mecanismos que auxiliem os peritos para análise de *e-mails*. Principalmente ao saber do grande volume de *e-mails* que os usuários armazenam em seus computadores, tornando o tempo de análise uma questão importante.

Além disto, deve-se considerar o trabalho exaustivo e susceptível a erros que é para o perito a análise manual de bases de *e-mails*.

Assim, proporciona-se neste trabalho um mecanismo que realiza o rastreamento em informações encontradas em contextos criminosos registrados na base de *e-mails* de um usuário, visando na produção de provas digitais em *e-mails*. O contexto é formado pelo conjunto de palavras contidas nas mensagens trocadas entre contatos envolvidos em uma conversação por *e-mail*. Uma conversação envolve *e-mails* que possuam características em comum, neste caso, o mesmo contexto.

Cita-se como exemplo, a existência de um grupo de *e-mails* que tenham sido trocados entre dois ou mais contatos. Analisando-se o conteúdo destes *e-mails*, verifica-se que os *e-mails* possuem contextos em comum, ou seja, palavras ou expressões em comum. E na existência de um contexto criminoso, pode-se realizar a produção de provas digitais, assim como o rastreamento do contexto.

Mediante aplicação de métodos de agrupamento e classificação em conversações, visa-se a produção de provas digitais e seu rastreamento. Estes métodos são aplicados em ocasiões oportunas, conforme se descreve na Seção 4.2.

4.2. Mecanismo

O mecanismo realiza a classificação de conversações em uma base de *e-mails*, buscando palavras que definam contextos de mensagens criminosas. Utiliza-se uma grande quantidade de informações que identificam *e-mails* (discussão realizada no Capítulo 2), considerando o conteúdo textual formado pelo corpo e campo assunto (cabeçalho) dos *e-mails*. Portanto, o mecanismo funciona em quatro FASES. Nas próximas Subseções explicam-se estas FASES.

4.2.1. FASE I

Na FASE I, ilustrada na Figura 4.1, procede-se a leitura digital e cópia do conjunto de *e-mails* (procedimento de “imagem” forense) que se pretende investigar, pois as evidências digitais (neste caso, os *e-mails*) não podem sofrer modificações durante um processo de coleta/análise pericial [PHI09] [SHI08] [REI02].

Considerando que os *e-mails* suspeitos tenham sido copiados integralmente, e por isso, a partir deste momento o mecanismo utiliza apenas o conjunto formado pelos *e-mails* copiados.

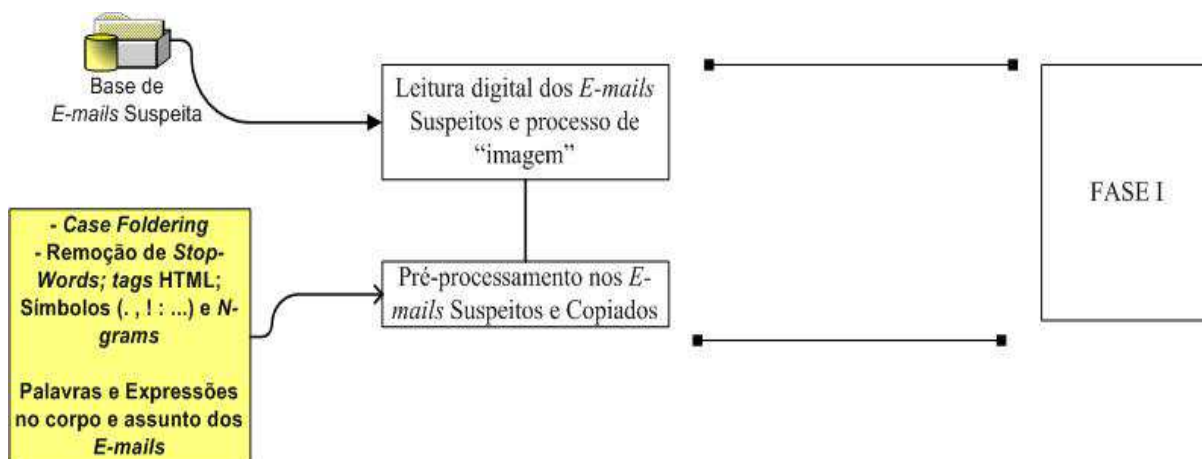


Figura 4.1 – FASE I do Mecanismo

Na FASE I, realiza-se um pré-processamento no conteúdo textual dos *e-mails*, retirando-se elementos existentes no corpo e campo assunto: exclusão de códigos HTML, URLs, símbolos e números. Aplicam-se também as técnicas de *case foldering*, *stop-words* e *n-grams*. Resultando assim, apenas as palavras e termos essenciais para facilitar na determinação das conversações criminosas existentes em uma base de *e-mails*.

Ressalta-se que aplicação do mecanismo não existe modificação nos *e-mails* originais de uma base de *e-mails*, mesmo pelo fato do mecanismo ter como pré-requisito a prática da “imagem” (cópia autêntica) da base de *e-mails* a ser periciada – este procedimento é classicamente executado em processos de investigação computacional.

4.2.2. FASE II

Na FASE II, o mecanismo agrupa os *e-mails* pertencentes a uma mesma conversação, ou seja, realiza a extração de conversações da base de *e-mails* sendo analisada. Representa-se a FASE II na Figura 4.2. Desta forma, *e-mails* que possuam palavras ou expressões repetidas são agrupados mediante a utilização de métodos de agrupamento ou de classificação.

Aplica-se esta FASE visto a possibilidade de usuários nem sempre utilizarem o mesmo endereço de *e-mail* durante uma conversação, e assim é possível descobrir todos os usuários (contatos) envolvidos em uma conversação.

Para realização da FASE II, o mecanismo realiza a extração de atributos a partir dos *e-mails* pré-processados na FASE I. Os atributos para esta etapa do processo podem constar de palavras ou expressões repetidas no corpo dos *e-mails*, assim como de palavras ou expressões encontradas no campo assunto do cabeçalho dos *e-mails*, formalizando desta forma várias possibilidades para análise dos resultados provenientes do método de agrupamento. A escolha do atributo é determinada pelo usuário utilizar do mecanismo.

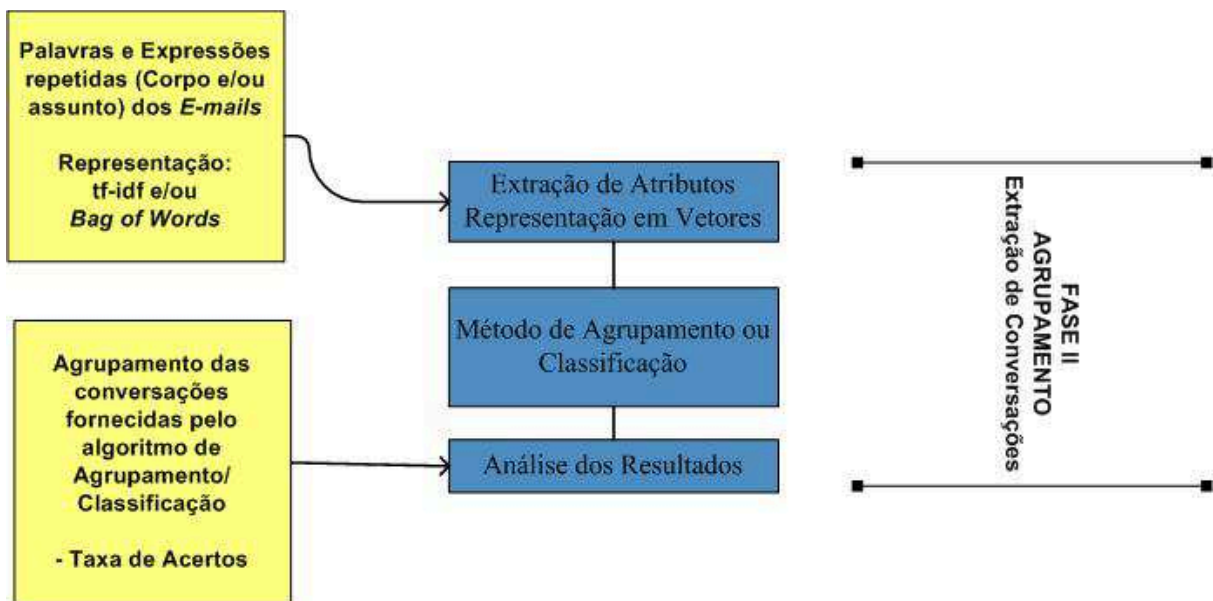


Figura 4.2 – FASE II do Mecanismo

Os atributos citados são representados em arquivos de atributos, que são formalizados em formato digital com extensão ARFF (*Attribute-Relation File Format*). Apresenta-se na Figura 4.3 um arquivo de atributos tendo as 22 palavras pertencentes a um conjunto de 8 *e-mails* como sendo os atributos. O arquivo é dividido em duas partes. Na primeira, são listados os atributos escolhidos a partir dos *e-mails*, assim como seu tipo.

Na segunda etapa do arquivo tem-se a contagem do número de ocorrências para cada atributo existente em cada *e-mail* pertencente a um conjunto de *e-mails*. E na última linha que descreve os atributos apresenta-se a linha “@attribute class {CLASSE1, CLASSE2, CLASSE3}”. No exemplo exposto na Figura 4.4, a última linha é utilizada para representar as conversações de uma base de *e-mails*. Optou-se em atribuir cada *e-mail* a uma conversação. Em uma análise envolvendo 8 *e-mails* por exemplo, destinou-se os *e-mails* em 3

conversações (CLASSES), assim como poderia ser criado 8 CLASSES. Uma CLASSE para cada *e-mail*.

Atributos (22)	<pre> @attribute palavras_ola_ : real @attribute palavras_user_ : real @attribute palavras_tudo_ : real @attribute palavras_bem_ : real @attribute palavras_ab_ : real @attribute palavras_voce_ : real @attribute palavras_parece-me_ : real @attribute palavras_que_ : real @attribute palavras_abobado_ : real @attribute palavras_agilize_ : real @attribute palavras_desculpe_ : real @attribute palavras_por_ : real @attribute palavras_ser_ : real @attribute palavras_desculpado_ : real @attribute palavras_mas_ : real @attribute palavras_marcado_ : real @attribute palavras_horas_ : real @attribute palavras_dia_ : real @attribute palavras_ok_ : real @attribute palavras_nao_ : real @attribute palavras_esqueca_ : real @attribute palavras_documento_ : real @attribute class {CLASSE1,CLASSE2,CLASSE3} @data </pre>
Contagem dos Atributos. Para uma base contendo 8 e-mails.	<pre> 2,9,3,3,2,0,CLASSE1 2,9,3,3,2,4,0,CLASSE1 0,9,0,0,0,4,3,3,5,4,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,CLASSE2 0,9,0,0,0,4,3,3,5,4,2,2,2,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,CLASSE2 0,9,0,0,0,4,3,3,5,4,2,2,2,1,1,0,0,0,0,0,0,0,0,0,0,0,0,0,CLASSE2 0,3,3,3,0,0,0,0,CLASSE3 0,3,3,3,2,0,0,0,CLASSE3 0,3,3,3,2,1,1,1,CLASSE3 </pre>

Figura 4.3 – Arquivo de Atributos – FASE II

Para representação dos atributos na FASE II, utiliza-se mais de uma técnica (*bag of words* ou *tf-idf*), possibilitando a formação dos arquivos de atributos, e assim, proporcionando o encontro de mais de um resultado e permanecendo com aquele que melhor obtiver resultados durante o processo de experimentos realizados durante a FASE II.

No total, o mecanismo propicia a formação de 9 arquivos, considerando-se os atributos: palavras do corpo (*bag of words* e *tf-idf*); expressões do corpo (*bag of words* e *tf-idf*); palavras do campo assunto (*bag of words* e *tf-idf*); expressões do campo assunto (*bag of words* e *tf-idf*) e palavras do corpo e palavras do campo assunto (*bag of words*).

Desta forma, os métodos (agrupamento e classificação) aplicáveis na FASE II analisarão cada arquivo de atributos, analisando se determinado *e-mail* pertence à

conversao predestinada. Utilizam-se diferentes mtodos, buscando-se o encontro do mtodo e arquivo de atributo que propicie o melhor resultado para utilizao na aplicao do mecanismo. No caso dos mtodos de agrupamento, avalia-se a matriz de similaridade, assim como a taxa de erros. E no caso de classificadores, a matriz de confuso e os valores de exatido e taxa de erro. A avaliao  realizada por intermdio do *software WEKA*²³.

4.2.3. FASE III

Aps ter sido realizada a extrao de caractersticas, na FASE III, usando-se mtodos de classificao. Classificam-se os *e-mails* das conversaes, visualizando-se a sequncia das atividades da FASE III apresentados na Figura 4.4. Para isso, todos os *e-mails* das conversaes so classificados com crime, ou no. No caso deste trabalho de Mestrado, crimes de assdio moral, ou no.

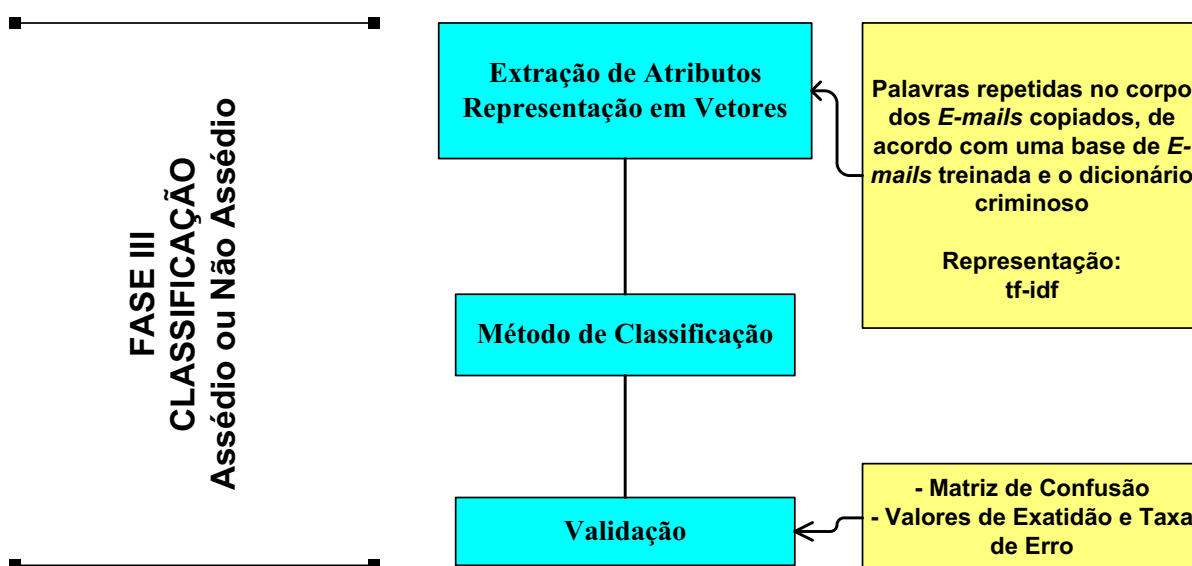


Figura 4.4 – FASE III do Mecanismo

Nesta FASE geram-se outros arquivos de atributos. Para isto, utilizam-se *e-mails* pertencentes à base de *e-mails* pré-processada (FASE I), de uma base de *e-mails* chamada de treinamento e do auxílio de um dicionrio criminoso. Desta forma, o mecanismo propicia a gerao de 4 arquivos de atributos teste e 4 arquivos de atributos treinamento.

²³ Weka. *Data Mining Software in Java*, disponvel em <http://www.cs.waikato.ac.nz/ml/weka/>, acesso em 23 de janeiro de 2011.

O mecanismo propõe a geração de arquivos de atributos para a FASE III no uso da técnica de *n-grams* e representação via tf-idf. Assim, realiza-se a representação em arquivo digital com extensão ARFF. Comparam-se os termos resultantes da aplicação da técnica *n-grams* nas palavras do corpo dos *e-mails* (base teste e base treinamento), com a aplicação da técnica de *n-grams* nas palavras do dicionário criminoso, utilizando similaridade de 67% ou 100%. E representam-se os vetores formados no uso da técnica tf-idf.

Para tal, geram-se 4 arquivos de teste: 3-grams 67% e 100% e 4-grams 67% e 100%, e 4 arquivos de treinamento (3-grams 67% e 100%) e (4-grams 67% e 100%).

Os arquivos de atributos de treinamento são gerados a partir de uma base de *e-mails* de treinamento, onde existem *e-mails* considerados criminosos, e *e-mails* não criminosos.

Posteriormente, os arquivos (teste e treinamento) são submetidos aos métodos de classificação (SVM – *kernels Polynomial, Radial e Sigmoid*, NB e DT) no uso do *software* WEKA. Para validação dos métodos de classificação empregados na FASE III, realizam-se a construção da matriz de confusão, assim como a obtenção dos valores para exatidão e taxa de erro.

4.2.4. FASE IV

Na FASE IV apresentam-se os resultados alcançados pelo mecanismo, conforme apresentado na Figura 4.5. O mecanismo efetua a representação dos resultados mediante a construção de grafos dirigidos, podendo assim formalizar um laudo pericial, onde constará onexo causal. Facilita-se desta forma o entendimento e rastreamento de cada grupo criminoso (conversação criminosa) encontrado. E ainda, mediante análise dos resultados proporcionados pelo mecanismo, é possível proporcionar as respostas apresentadas na Seção 4.3. (Resultados Proporcionados).

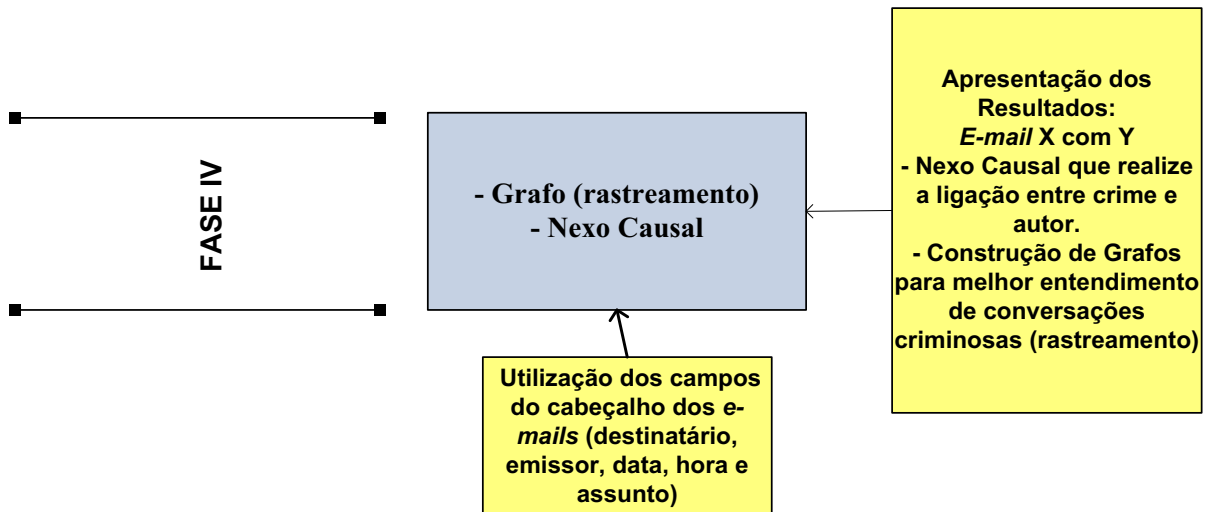


Figura 4.5 – FASE IV do Mecanismo

Identificados os resultados da FASE II e FASE III, obtém-se a informação de quais métodos o mecanismo irá utilizar durante os experimentos. Portanto, uma planilha é formulada com os resultados proporcionados pelos métodos (um para FASE II e um para FASE III). Apresenta-se na Figura 4.6 exemplo desta planilha, cujos resultados são da aplicação do mecanismo em uma base contendo 8 *e-mails*. A planilha está dividida em 5 campos, sendo: *email* (usuário remetente), *grupo* (conversação), *criminoso* (classificação criminosa do *e-mail*), *receiver* (contato que recebeu o *e-mail*), e *idEmail* (número de identificação do *e-mail*).

FASE II		FASE III			
Query Editor					
Table					Database Connection
results					mail
email	grupo	criminoso	receiver	idEmail	
user02@dominio.com	2	1	user01@dominio.com	1	
user01@dominio.com	0	1	user02@dominio.com	2	
user02@dominio.com	1	0	user03@dominio.com	3	
user03@dominio.com	1	1	user02@dominio.com	4	
user02@dominio.com	1	0	user03@dominio.com	5	
user02@dominio.com	2	1	user04@dominio.com	6	
user04@dominio.com	2	1	user02@dominio.com	7	
user02@dominio.com	2	1	user04@dominio.com	8	

Figura 4.6 – Resultado do Agrupamento e Classificação das Conversações

Na planilha apresentada (Figura 4.6), imaginando que o método tenha utilizado o classificador NB para execução da FASE II e o método SVM *linear* para a FASE III, na construção desta planilha, pode-se afirmar que **NB** relacionou todos os *e-mails* em respectivas conversações (FASE II). Como exemplo, cita-se o *e-mail* 1 (*idEmail*). Este *e-mail* foi agrupado por **NB** como pertencente a conversação 2. E o método **SVM *linear***, classificou o *e-mail* 1 (*idEmail*) como 1 (não criminoso).

A planilha serve de apoio para a construção de uma matriz de adjacências. Para construção da planilha da Figura 4.6, o mecanismo verifica os endereços de *e-mails* encontrados nos campos *To* e *From* existentes no cabeçalho de cada *e-mail* pertencente à conversação. Um *e-mail* enviado para 2 contatos será listado por 2 vezes nesta planilha. Assim, analisando os resultados proporcionados por esta planilha, expõem-se a quantidade de *e-mails* existentes em cada uma das conversações existentes em uma base de *e-mails*, o qual será fornecido para matriz de adjacências.

Na matriz de adjacências têm-se resumidamente as quantidades de *e-mails* enviados e recebidos entre remetentes e destinatários, possibilitando a construção de grafos dirigidos para cada conversação (FASE IV). Para construção dos grafos, utilizam-se as informações da matriz de adjacências. Assim, cada grafo possuirá todos os contatos encontrados em uma conversação agrupada, e em cada nodo de um grafo (contato), a soma da quantidade de *e-mails* enviados e recebidos por contato, a quantidade de *e-mails* enviados e aqueles que foram recebidos, seguindo a nomenclatura: TOTAL [ENVIADOS, RECEBIDOS].

No exemplo da base de *e-mails* utilizada para construção da planilha exposta na Figura 4.6, foram agrupadas 3 conversações, e por isso, a FASE IV do mecanismo proporciona a criação de 3 grafos. No caso da conversação 01 (grupo 1 na Figura 4.6), *User01* possui um total de 1 *e-mail*, sendo que enviou 1 *e-mail* e recebeu 0. *User02*, possui um total de 1 *e-mail*, sendo que enviou 0 *e-mails*, e recebeu 1 *e-mail*.

Por padronização, o grafo de cada conversação é colorido em **azul** na **não** existência de *e-mail* criminoso no relacionamento representado, e em cor **vermelha** do contrário. Na Figura 4.7 é exposto exemplo da construção dos grafos elaborados pela aplicação do mecanismo na base de *e-mails* representada na planilha da Figura 4.6. Desta forma, por meio da exibição do grafo, um perito pode mais facilmente analisar a rede de contatos que

envolvem uma conversação, assim como realizar a rastreabilidade que envolve uma conversação.

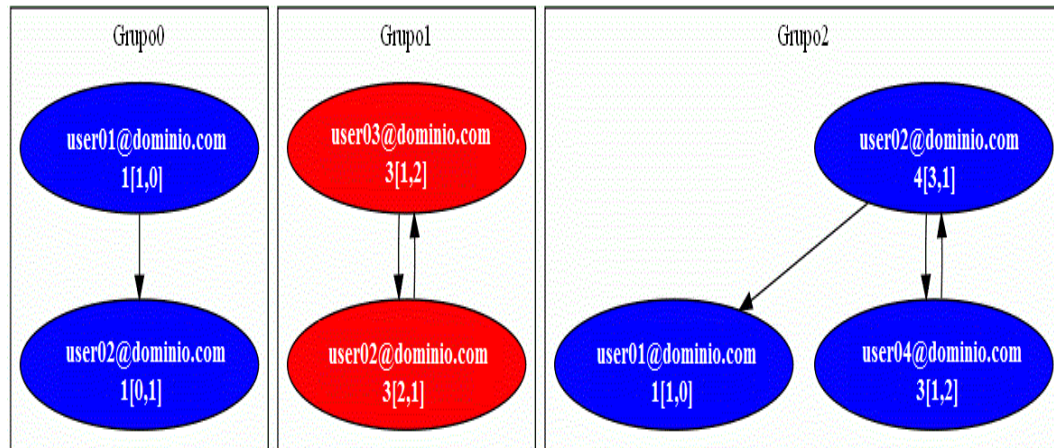


Figura 4.7 – Grafos Gerados

Pela execução da FASE IV descobrem-se todos os endereços de *e-mail* pertencentes a uma conversação investigada. Contatos que utilizam diferentes endereços de *e-mail* para se comunicarem com o proprietário da respectiva base, por exemplo.

Finalmente, na Figura 4.8 apresenta-se um exemplo de laudo pericial da análise de conversações encontradas na investigação de uma base de *e-mails* de exemplo. Neste laudo pericial de exemplo estão contidas as informações relevantes aos profissionais que recebem o resultado de toda a análise proveniente do mecanismo, como por exemplo, a data e hora em que foi realizada a análise de uma base, a lista de *e-mails* das conversações criminosas e métodos empregados.

Na Figura 4.8 constam alguns campos em destaque. Data em que foi gerado o relatório e aplicado o mecanismo. Quantidade de *e-mails* analisados. Método utilizado para Agrupamento e Classificação das conversações. *E-mails* Criminosos – esta sendo listado o usuário remetente e destinatário de cada *e-mail* criminoso, a data e horário, assim como o conteúdo do campo assunto.

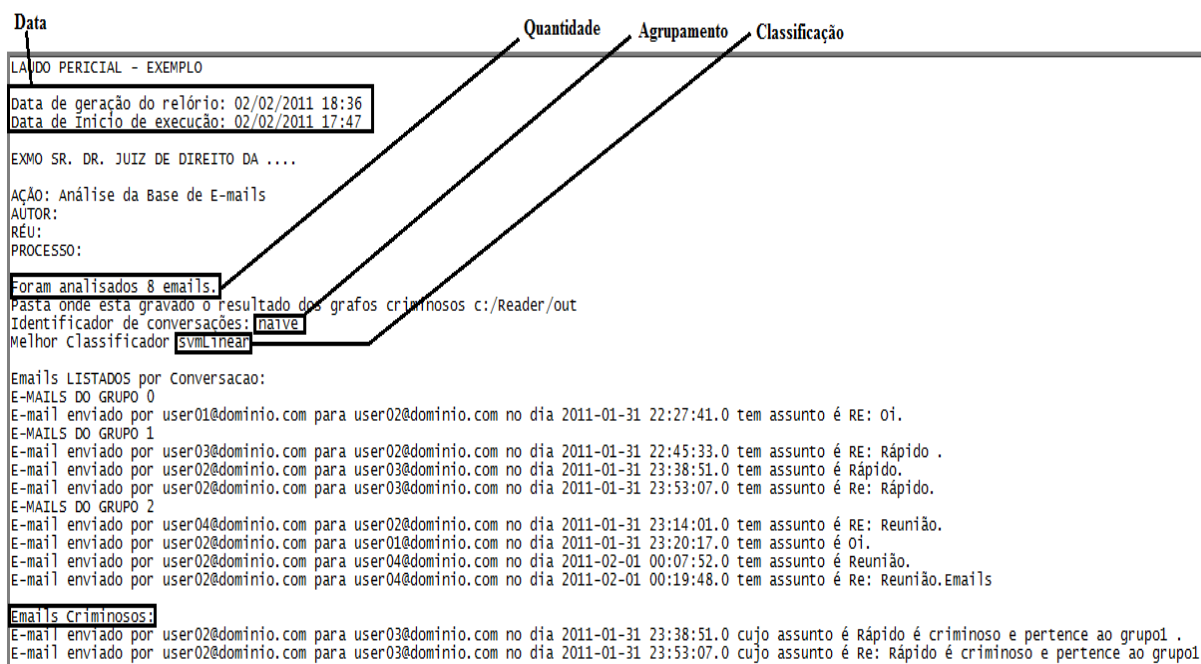


Figura 4.8 – Fragmento do Laudo Pericial Exemplo

O mecanismo faz uso de vários métodos de agrupamento e classificação, assim como de diversas técnicas. Para seu funcionamento, realiza-se primeiramente o agrupamento de conversações, para então concretizar a classificação das conversações como criminosas, ou não. Desta forma, torna-se viável a rastreabilidade de conversações criminosas, já que as conversações são expostas através de grafos. E com a aplicação do mecanismo, busca-se o fornecimento de máxima precisão para produção das provas digitais. Sendo assim, na próxima Seção listam-se os resultados proporcionados pela aplicação do mecanismo.

4.3. Bases Experimentais

Para realização dos experimentos, foram elaboradas duas bases de *e-mails* em idioma português (Brasil) e que seguiram as características daquelas apresentadas e estudadas na Seção Trabalhos Relacionados. Para elaboração das bases de *e-mails*, necessitou-se de um dicionário criminoso. Detalhes deste dicionário são descritos na Subseção 4.3.1. E nas Subseções 4.3.2 e 4.3.3 apresentam-se as bases de *e-mails* elaboradas.

4.3.1. Dicionário Criminoso

No caso de crime registrado em *e-mails*, aplica-se o mecanismo no uso de um dicionário. O dicionário deverá possuir palavras que caracterizem o crime que se queira investigar.

Entre os possíveis crimes, verifica-se um aumento no número de processos judiciais no Brasil envolvendo crime de assédio moral, conforme estatísticas para o crescimento do número de ações trabalhistas²⁴. O crime de assédio moral pode ser cometido na utilização da ferramenta *e-mail*, já que o *e-mail* tornou-se a ferramenta de comunicação que possibilita a rápida comunicação entre contatos, seja dentro ou fora de uma empresa, por exemplo.

VILLATORE et al. (2009) estudaram as definições e conceitos jurídicos para o tema de crime de assédio moral frente ao entendimento de diversos autores [VIL09]. O crime de assédio moral ocorre a partir de um ambiente de trabalho, tendo o crime de assédio moral responsabilidades trabalhista, civil e criminal para o agressor [VIL09].

Visto a importância deste tipo de crime, aplica-se experimentalmente o mecanismo para produção de provas digitais em *e-mails* contendo crimes de assédio moral. Para isso, neste trabalho utilizou-se como fonte de pesquisa os resultados do estudado de [NUN09], e que foram descritos na Seção 3.3. Um dicionário e uma base de *e-mails* contendo resquícios de crimes. O dicionário possui atualmente 539 palavras, e a base, 40 *e-mails*.

As palavras do dicionário foram utilizadas para geração de parte da base de *e-mails* de treinamento, sendo que na Subseção 4.3.2, explicam-se os procedimentos para elaboração de 1.876 *e-mails*. Além da base de Treinamento, as palavras do dicionário são utilizadas como atributos para a representação dos *e-mails* de Teste para execução da FASE III do mecanismo. Explica-se esta representação na Seção 4.3.3.

E ainda, a base de *e-mails* do trabalho [NUN09] foi dividida na razão 80:20. A base de Treinamento recebeu um total de 32 *e-mails* (80%), e a base de Teste, um total de 8 *e-mails* (20%). Nas próximas Subseções apresentam-se especificações da base treinamento e teste.

4.3.2. Base de *e-mails* de Treinamento

A base de *e-mails* de Treinamento foi formalizada como referência para classificação do crime de assédio moral [NUN09] existente nos *e-mails* da base de Teste. Optou-se em investigar criminalmente as palavras encontradas no corpo dos *e-mails*, e em função disso, a

²⁴ Pesquisa disponível em http://www.protecao.com.br/site/content/noticias/noticia_detalhe.php?id=Jyy5AJjj, acesso em 24 de janeiro de 2011.

base de Treinamento sofreu análise forense apenas nas palavras do corpo de seus *e-mails*. Esta análise forense ocorreu haja vista que palavras criminosas normalmente não são expressas no campo assunto de *e-mails* trocados entre contatos.

Apresentam-se na Tabela 4.1 estatísticas da base de Treinamento, que é composta por 3.816 *e-mails*, tendo uma média de 501,28 palavras. A base de Treinamento está dividida na razão 50:50. Assim, 1.908 *e-mails* referem-se à classe Não Crime de Assédio Moral, e outros 1.908 na classe Crime de Assédio Moral (1.876 + 32).

Originaram-se 1.876 *e-mails* característicos de crime mediante a combinação das 539 palavras do dicionário de assédio moral [NUN09] somado a expressões particulares. Inicialmente originaram-se 539 *e-mails* contendo as 539 palavras do dicionário, sendo adicionada em cada *e-mail* a expressão “você é um”. Logo após, combinaram-se a primeira palavra do dicionário com as demais, sendo que se originaram outros 539 *e-mails* com as seguintes características: a expressão “você é um”, somando a primeira palavra do dicionário, somando a expressão “e um”, e por fim, somando-se a segunda palavra do dicionário. E ainda, outros 798 *e-mails* foram originados utilizando-se do mesmo raciocínio, totalizando-se 1.876 *e-mails*, e com uma média de 7,82 palavras existentes no corpo de cada um desses *e-mails*.

Tabela 4.1 – Base de *e-mails* de Treinamento

Classe – Crime de Assédio Moral	Quantidade de <i>E-mails</i>	Média de Palavras - Corpo
Crime	1.876	7,82
	32	12,46
Não Crime	1.908	481
	3.816	501

Ainda na formalização dos *e-mails* da classe Crime, utilizou-se 80% dos *e-mails* pertencentes à base de *e-mails* do trabalho de [NUN09]. Ou seja, 32 *e-mails* criminosos foram adicionados a base de Treinamento, sendo que este segundo conjunto da classe Crime possui uma média de 12,46 palavras no corpo de seus *e-mails*.

Na formação da base de Treinamento, adicionou-se 1.908 *e-mails* não criminosos. Os *e-mails* não criminosos foram selecionados manualmente, sendo que estes *e-mails* contêm propagandas, mensagens cotidianas, frases e textos, entre outros. Todos os *e-mails* foram lidos para certificação da não existência de indícios de palavras ou expressões que estejam

relacionadas ao contexto de assédio moral ou outro tipo de crime, e possuem uma média de 481 palavras em seu corpo.

Visto que a base de Treinamento é utilizada como referência para classificação de crime de assédio moral nos *e-mails* de Teste, indica-se que na aplicação do método de classificação são utilizadas como referência as palavras do dicionário criminoso [NUN09]. E na sequência explica-se a base de *e-mails* de Teste.

4.3.3. Base de *e-mails* de Teste

A base de Teste foi coletada a partir dos *e-mails* do autor desta Dissertação e mapeada em grupos (conversações). Nas conversações se objetivou aplicar métodos durante os experimentos para agrupamento, classificação e posteriormente representação gráfica. Desta forma, apresentam-se nesta Subseção, estatísticas e explicações da formulação da base de Teste, sendo que na Tabela 4.1 são relatadas informações desta base.

Totaliza-se a seleção de 570 *e-mails* adquiridos entre os *e-mails* existentes na caixa de entrada do *software* de *e-mail* do autor. Todos os *e-mails* selecionados sofreram modificações nos campos de seus cabeçalhos. A saber, existe um total de 50 contatos envolvidos nesta base de *e-mails*, e os endereços do remetente e destinatário desses *e-mails* foram renomeados, proporcionando sigilo a todos os contatos participantes das conversações.

Na sequência identificaram-se as conversações existentes dentre a base de *e-mails* de Teste. Para identificação das conversações, individualmente, foi realizada a leitura de cada *e-mail* da base de Teste pelo autor, e assim rotulando os *e-mails* em conversações. Fez-se a leitura do conjunto textual pertencente ao corpo de cada *e-mail* e dos endereços eletrônicos existentes no cabeçalho de cada *e-mail* também (campos *To* e *From*). Assim, os *e-mails* da base de Teste que tiveram as mesmas palavras em seu corpo, foram rotulados pertencentes à mesma conversação. Assim, obteve-se a quantidade de contatos envolvidos em cada conversação, vide Quantidade de Contatos – Tabela 4.2.

Classificadas e quantificadas as conversações, informa-se que dentre os 570 *e-mails* da base de Teste, existem um total de 33 conversações. Cada conversação envolve um número diferente de *e-mails* (Quantidade de *E-mails* – Tabela 4.2), assim como um número diferente de contatos (Quantidade de Contatos – Tabela 4.2). Observa-se os valores expressos na Tabela 4.2, um número mínimo de *e-mails* nas conversações *D* e *JP* com 2 *e-mails* cada, e um número máximo na conversação *GPF* com 90 *e-mails*.

Fortalecendo as estatísticas da base de Teste, totalizou-se a quantidade de palavras encontradas em cada conversação. Esta totalização foi realizada considerando a soma das palavras do corpo dos *e-mails* envolvendo as conversações (Total de Palavras - Corpo – Tabela 4.2) e das palavras do campo assunto dos respectivos *e-mails* (Total de Palavras – Assunto – Tabela 4.2). Foi detectado no corpo dos *e-mails* o equivalente a 1.477.375 palavras e no campo assunto, um total de 1.714 palavras. Em ambos os caso, foram somados apenas as palavras representativas do corpo ou do campo assunto dos *e-mails*, ou seja, a contagem apresentada foi realizada após a aplicação da FASE I do mecanismo na base de Teste (pré-processamento), conforme explicado na Subseção 4.2.1.

Assim, tem-se o número médio de palavras utilizado em cada conversação, tanto para as palavras existentes no corpo (Média de Palavras - Corpo – Tabela 4.2) quanto para as palavras existentes no campo assunto (Média de Palavras - Assunto – Tabela 4.2) dos *e-mails*. Obteve-se uma média equivalente a 1.702,85 palavras encontradas no corpo, e 2,81 palavras para o campo assunto.

Ademais, calculou-se o tamanho médio das palavras de cada conversação. O tamanho médio refere-se à soma do número de caracteres das palavras do corpo ou campo assunto dos *e-mails* de uma conversação, dividido pela quantidade de *e-mails* envolvidos nesta conversação. Obteve-se um tamanho médio equivalente a 7,35 caracteres (Tamanho Médio das Palavras – Corpo – Tabela 4.2) no uso das palavras do corpo, e 6,09 caracteres utilizando palavras do campo assunto (Tamanho Médio das Palavras – Assunto – Tabela 4.2) dos *e-mails* de Teste. Esta informação é importante, pois ao se aplicar o método de *N-grams*, o tamanho das palavras influencia o resultado da obtenção das sub-cadeias de caracteres.

As conversações apresentadas na Tabela 4.2 diferem-se pela quantidade de *e-mails* e contatos envolvidos, assim como na quantidade média de palavras existentes no corpo e campo assunto dos *e-mails* pertencentes a cada conversação.

Entre as conversações, a de número 29 possui 30 *e-mails*. Os *e-mails* foram trocados entre 7 contatos, sendo que em 90% desses *e-mails*, observa-se que cada *e-mail* da conversação foi enviado para mais de um contato. Os *e-mails* possuem frases retiradas de 8 *e-mails* (20%) pertencentes a base de *e-mails* de assédio moral elaborada por [NUN09], podendo-se citar como exemplos: “Essa **negrinha** tentou me derrubar, mas não conseguiu”, “não agir com espírito de cavalo selvagem ou até mesmo com espírito de **porco**”, “o próximo que escrever **besteira** será demitido”, “você é um **abobado**”.

Desta forma, todos os *e-mails* da conversação 29 possuem palavras criminosas, conforme o dicionário de assédio moral. E analisando as conversações da base de Teste, verificou-se que os *e-mails* de uma conversação nem sempre tinham a mesma quantidade de palavras no campo assunto. Cita-se como exemplo as conversações que tiveram os *e-mails* com número médio de palavras existentes no campo assunto (Média de Palavras – Assunto – Tabela 4.2) com valor fracionado, como por exemplo, as conversações 0, 5, 9 e 16. Isto identifica a probabilidade dos *e-mails* dessas conversações possuírem modificações nas palavras existentes no campo assunto, evidenciando que uma conversação pode possuir *e-mails* com diferentes palavras no campo assunto. Entretanto, o que não ocorre com as palavras existentes no corpo dos *e-mails*. Para que exista uma conversação, seus *e-mails* devem compartilhar do mesmo conjunto de palavras.

Sendo assim, apresentou-se nesta Subseção características da base de *e-mails* de Teste. Nesta base, somente os *e-mails* da conversação 29 possuem características criminosas de assédio moral, ao contrário de todas as outras conversações.

Sendo assim, finaliza-se nesta Subseção a apresentação do ambiente experimental. Na próxima Seção apresenta-se a implementação do mecanismo, para então serem disponibilizados os resultados alcançados pela aplicação do mecanismo na produção de provas digitais.

Tabela 4.2 – Base de e-mails de Teste

Conversaço	Descrição	Quantidade de Contatos	Quantidade de e-mails	Total de Palavras Corpo	Média de Palavras Corpo	Total de Palavras Assunto	Média de Palavras Assunto	Tamanho Médio das Palavras Corpo	Tamanho Médio das Palavras Assunto
0	C	2,00	11,00	23.147,00	2.104,27	34,00	3,09	7,370	7,58
1	D	2,00	2,00	2.772,00	1.386,00	6,00	3,00	12,92	4,00
2	JP	2,00	2,00	1.267,00	633,50	10,00	5,00	7,51	3,57
3	M	2,00	6,00	5.648,00	941,33	6,00	1,00	6,39	2,00
4	P	2,00	3,00	4.338,00	1.446,00	3,00	1,00	7,78	3,00
5	B	2,00	12,00	32.585,00	2.715,40	40,00	3,33	6,13	4,28
6	HY1	2,00	8,00	6.178,00	772,25	8,00	1,00	7,86	3,00
7	L	2,00	5,00	10.326,00	2.055,20	5,00	1,00	7,83	4,00
8	A	2,00	7,00	8.770,00	1.252,85	7,00	1,00	7,07	5,00
9	CI	3,00	15,00	103.321,00	6.880,00	43,00	2,86	7,83	6,06
10	GPF	16,00	90,00	105.002,00	1.160,68	90,00	1,00	5,62	11,00
11	GP	7,00	18,00	20.997,00	1166,50	18,00	1,00	7,12	11,00
12	F1	2,00	4,00	5.117,00	1.279,25	16,00	4,00	7,83	6,75
13	F2	2,00	6,00	5.504,00	917,33	42,00	7,00	5,72	3,71
14	GI	3,00	4,00	3.472,00	868,00	4,00	1,00	9,89	7,00
15	LA	2,00	3,00	1.341,00	447,00	6,00	2,00	5,72	5,50
16	RE	3,00	6,00	2.911,00	485,16	14,00	2,33	7,22	7,22
17	SB	4,00	29,00	51.635,00	1.780,52	29,00	1,00	5,53	12,00
18	OT	5,00	5,00	2.434,00	486,80	40,00	8,00	4,57	4,37
19	NA	6,00	5,00	5.986,00	1.197,20	33,00	6,60	6,90	6,44
20	HY2	2,00	54,00	249.163,00	4.614,13	54,00	1,00	8,10	11,00
21	HY3	3,00	45,00	90.707,00	2.015,70	135,00	3,00	8,53	2,66
22	SM	2,00	4,00	1.464,00	366,00	9,00	2,25	7,42	4,75
23	BB	2,00	11,00	6.827,00	620,64	11,00	1,00	7,93	11,00
24	J1	2,00	61,00	443.993,00	7.278,57	122,00	2,00	8,17	6,50
25	GL	2,00	7,00	2.928,00	418,28	21,00	3,00	7,92	6,00
26	GE	2,00	13,00	12.882,00	990,92	13,00	1,00	7,08	7,00
27	CA	2,00	27,00	102.838,00	3.808,80	27,00	1,00	11,50	9,00
28	J2	2,00	30,00	39.205,00	1.306,80	60,00	2,00	4,44	6,50
29	CRIMINOSOS	7,00	20,00	32.642,00	1.632,10	60,00	3,00	7,31	7,33
30	C31	2,00	17,00	17.421,00	1.024,76	34,00	2,00	7,45	6,50
31	C32	13,00	21,00	22.079,00	1.051,38	42,00	2,00	5,30	6,50
32	C33	8,00	19,00	12.086,00	636,10	56,00	2,95	6,60	5,58
Total	33	-	50,00	1.436.986,00	-	1.098,00	-	-	-
Média	-	-	-	-	1.689,07	-	2,50	7,35	5,96

4.4. Resultados Proporcionados

Mediante análise dos resultados (FASE IV - Figura 4.1) respondem-se questionamentos da existência de relacionamentos entre usuários envolvidos em uma conversação e que foram representados por grafos. Citam-se as respostas:

- Existência de relacionamento entre contatos sendo investigados (grafo): Análise das conversações em que os contatos desta conversação utilizaram diferentes endereços de *e-mail* durante este relacionamento (comunicação), por exemplo;
- Quantidade de *e-mails* envolvidos em uma conversação: *E-mails* que foram enviados e respondidos pelo proprietário da base de *e-mails*;
- Horário do envio e recebimento de cada *e-mail*: Principalmente tempo entre um *e-mail* enviado e outro recebido;
- Número de contatos desconhecidos encontrados em uma conversação: Exemplo: contatos que apenas enviam *e-mails* ao proprietário da base sendo investigada, como por exemplo, *e-mails* de propaganda;
- Dentre os *e-mails* que compõem uma conversação analisada, em quais destes existem diálogos criminosos? Será listado o horário e data de envio ou recebimento para cada um;
- Quais as palavras criminosas frequentemente utilizadas em cada conversação agrupada.

Neste Capítulo apresentou-se o mecanismo para produção de provas digitais em *e-mails* envolvendo o rastreamento de relacionamentos. Seu funcionamento e suas vantagens foram explorados. No próximo Capítulo relata-se a implementação do mecanismo, assim como os detalhes técnicos e resultados dos experimentos realizados em uma base com 570 *e-mails*.

Capítulo 5

Resultados Experimentais

Neste Capítulo são apresentados os resultados experimentais da aplicação do mecanismo. Descreve-se a implementação do mecanismo em um *framework*, assim como os resultados dos experimentos para o agrupamento, a classificação, e exposição gráfica das conversações existentes na base de *e-mails* Teste.

5.1. *Reader* – Implementação do Mecanismo para Rastreamento de Relacionamentos em *e-mails*

Obtiveram-se resultados da aplicação do mecanismo mediante a implementação de um *framework* desenvolvido em linguagem *Java* [DEI05], sendo o *framework* denominado de *Reader*. *Reader* faz uso de várias APIs (*Application Programming Interface*) para seu funcionamento, sendo que os resultados experimentais apresentados neste Capítulo são expostos com maiores detalhes no Apêndice A. Sendo assim, relatam-se nas próximas Subseções as etapas tecnológicas desenvolvidas na elaboração de *Reader*.

5.1.1. Leitura dos *e-mails* e processo de “imagem”

Reader analisa evidências digitais em formato *e-mail* que tenham sido armazenados por usuários utilizadores do Sistema Operacional *Windows*, já que este é o Sistema Operacional mais utilizado por usuários de computador. *Reader* analisa a base de *e-mails* colhida pelo perito em procedimentos de produção de provas, sendo que a base deverá estar no mesmo formato digital que aquela proveniente do *software* cliente de *e-mail Windows Live Mail* (utilizando a extensão *eml*). As justificativas para escolha do tipo de Sistema

Operacional e do *software* cliente de *e-mail* foram apresentadas na Subseção 2.2. A base de *e-mails* que se deseja investigar deverá ser informada ao *framework Reader*.

Para realização da FASE I (Figura 4.1) - leitura e processo de “imagem” de uma base de *e-mails*, *Reader* faz uso da API *JavaMail*²⁵. No uso desta API, os *e-mails* são lidos digitalmente e seus campos copiados (corpo, assunto, data/hora, remetente e destinatários) para o banco de dados (BD) *mysql*²⁶.

5.1.2. Pré-processamento

No corpo e campo assunto dos *e-mails* lidos digitalmente, *Reader* aplica as técnicas de pré-processamento textual (FASE I – Figura 4.1): *case foldering*, *stop words*, remoção de códigos HTML, URLs, símbolos e números.

Assim, **para efeitos de padronização, na continuação da análise forense realizada por *Reader*, usam-se sempre os campos pré-processados nesta Subseção.** Desta forma, nas próximas Subseções, todos os *e-mails* são transformados em arquivos de atributos, assim como se define um tipo de crime a ser investigado por *Reader*.

5.1.3. Crime investigado

Reader pode ser utilizado na investigação de qualquer crime que contenha um dicionário formalizado. Para obtenção dos resultados experimentais, definiu-se a utilização do dicionário de crimes de assédio moral [NUN09]. Este dicionário foi detalhado no Capítulo 4.

5.1.4. Extração e Geração dos Arquivos de Atributos

Reader origina arquivos digitais específicos para a FASE II e outros para a FASE III do mecanismo.

Na geração dos arquivos utilizados na FASE II – Figura 4.2, *Reader* considera as palavras ou expressões dos *e-mails* lidos digitalmente como sendo atributos, gerando assim, um total 9 arquivos de atributos para utilização na FASE II.

E na FASE III – Figura 4.4, *Reader* considera como atributos as palavras do dicionário de [NUN09], gerando quatro arquivos de atributos, onde é realizada a comparação das

²⁵ API *JavaMail* 1.4.4, disponível em <http://www.oracle.com/technetwork/java/index-138643.html>, acesso em 23 de janeiro de 2011.

²⁶ Disponível em <http://www.mysql.com/>, acesso em 04 de fevereiro de 2011.

palavras do dicionário com as palavras encontradas no corpo dos *e-mails* sendo testados, e mais 4 arquivos formados a partir de uma base de treinamento.

Assim, *Reader* possibilita a construção de diferentes arquivos de atributos, ou seja, para a FASE II (agrupamento das conversações), 9 arquivos e na FASE III (classificação das conversações), um total de 8 arquivos. Todos os arquivos de atributos são posteriormente submetidos à análise de algoritmos de aprendizagem (métodos de agrupamento/classificação), descritos na próxima Subseção.

5.1.5. Algoritmos de Aprendizagem de Máquina

Descreveu-se anteriormente a extração e geração de atributos em arquivos. Esses arquivos são aplicados por diversos métodos de agrupamento/classificação para agrupamento de conversações, e posteriormente, a classificação das conversações. Nesta Subseção apresentam-se os métodos de agrupamento/classificação implementados no *framework Reader*.

Reader faz uso de algoritmos disponíveis no *software Weka*. O *software Weka* possui vários algoritmos de Aprendizagem de Máquina disponibilizados em formato de APIs.

Na FASE II, o *framework* utiliza vários métodos de agrupamento/classificação para o agrupamento de conversações a partir dos 9 arquivos de atributos gerados a partir da base de *e-mails* sendo investigada. Citam-se os métodos de agrupamento: *K-means* (*Euclidiana*, *Manhattan*, *Cossine* e *Jaccard*), e os métodos de classificação: SVM (*Polynomial*, *Radial* e *Sigmoid*), DT e NB.

Aplicaram-se diversos métodos na análise do encontro das melhores taxas de acertos para o agrupamento de conversações, e assim, proporcionou-se a determinação do melhor método, assim como do melhor arquivo de atributos contendo os *e-mails* representados. Os resultados das taxas de acertos são proporcionados pela execução do algoritmo de agrupamento/classificação.

Na FASE III do mecanismo, utilizam-se os métodos de classificação: SVM (*Polynomial*, *Radial*, *Sigmoid* e *Linear*), DT e NB para classificação das conversações agrupadas. No uso dos métodos, os 4 arquivos de atributos da base de *e-mails* investigada são testados, respectivamente, aos 4 arquivos de atributos da base de Treinamento. De forma similar a FASE II, os algoritmos se responsabilizam em informar a taxa de acertos, assim como qual a classificação de um *e-mail* testado. Os algoritmos utilizados na FASE III são

similares aos utilizados na FASE II, entretanto fazem uso da comparação entre o arquivo de atributos da base investigada e do arquivo de atributos de Treinamento.

Nesta Subseção apresentaram-se vários métodos que foram implementados no *framework Reader* para análise de arquivo de atributos. Os arquivos de atributos e os métodos de agrupamento/classificação foram detalhados no Capítulo 4. E salienta-se que os métodos escolhidos para implementação no *framework* são os mais utilizados na categorização de textos, já que proporcionam os resultados com as maiores taxas de acertos, conforme exposto no Capítulo 2. E na próxima Subseção apresentam-se os resultados proporcionados pela aplicação do *framework Reader* na base de Teste e de Treinamento.

5.1.6. Exposição dos Resultados

Realizada o agrupamento e classificação de conversações existentes em uma base de *e-mails*, *Reader* expõe os resultados proporcionados pelo mecanismo (FASE IV – Figura 4.5). As formas de apresentação dos resultados são explicados nesta Subseção.

Os resultados facilitam a construção do nexos causal, fazendo com que um perito consiga com maior agilidade e precisão comprovar a existência de relação entre uma conduta e o suposto infrator sendo investigado. Dentre os resultados proporcionados pelo *framework*, apresentam-se:

- Agrupada as conversações criminosas e não criminosas de uma base de *e-mails*, constroem-se grafos. No caso, *Reader* utiliza a API do *software Graphviz*²⁷, para auxílio na visualização gráfica de informações. *Reader* representa graficamente cada conversação utilizando os campos destinatário e remetente existente em cada *e-mail* que compõem a conversação. Para isso, primeiramente *Reader* faz a construção de uma matriz de adjacências, onde é verificado o relacionamento entre os contatos encontrados na conversação, assim como a quantidade de *e-mails* enviados e recebidos por contatos dentro de cada conversação. Um pedaço do código para geração da matriz de adjacências de *Reader* pode ser visto na Figura 5.1, onde se utiliza as informações da planilha apresentada na Figura 4.6. Visualizam-se todos os contatos envolvidos dentro de uma conversação, desde que estes contatos estejam listados nos campos do cabeçalho denominados destinatário ou remetente, sendo que os grafos são

²⁷ API *Graphviz*, disponível em <http://www.graphviz.org>, acesso em 23 de janeiro de 2011.

padronizados em cor azul ou vermelha. Após a API *Graphviz* realiza a leitura da matriz de adjacência e realiza a construção dos grafos.

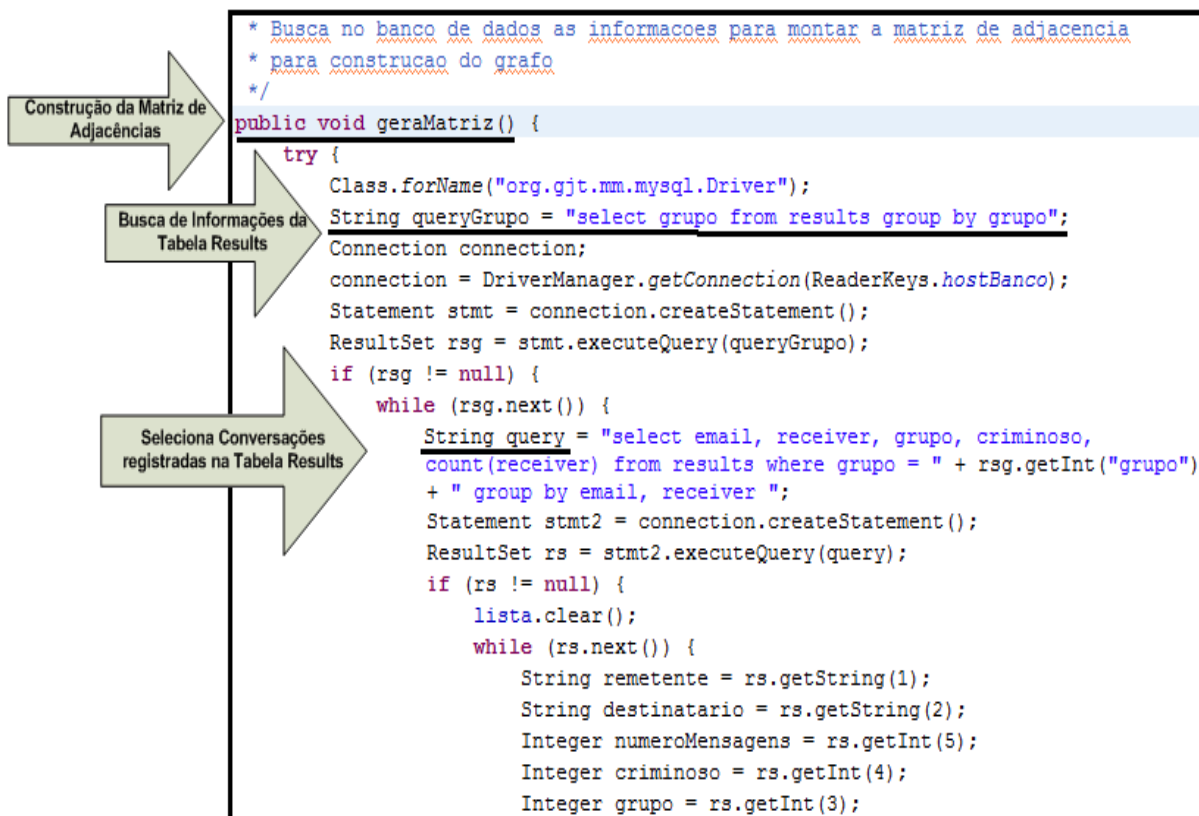


Figura 5.1 – Vista Parcial do Código para Geração da Matriz de Adjacências

- Cada conversação agrupada pode ter seus *e-mails* analisados individualmente, os quais são listados em ordem cronológica. *Reader* faz isso mediante a análise dos campos data/hora, assunto e corpo de cada *e-mail* pré-processado e copiado. Tem-se o horário do envio e recebimento de cada *e-mail*, assim como o tempo entre um *e-mail* enviado e outro recebido;
- Na análise de uma conversação criminosa, proporciona-se a visualização dos *e-mails* que possuam diálogos criminosos, ou não;
- Além de informar quais são os *e-mails* criminosos agrupados em uma conversação, possibilita-se também a detecção de quais as palavras/expressões criminosas frequentemente utilizadas em uma conversação investigada.

Baseado nos resultados expostos se facilita a construção do laudo pericial que servirá de base para aqueles que darão continuidade no processo Judicial. *Reader* proporciona um modelo, onde se relatam resultados obtidos pela aplicação do mecanismo, o qual pode ser utilizado como base pelo perito da área de Computação Forense. Todavia, este laudo pericial não é utilizado em sua essência pelo mecanismo, por questões de padronização.

Apresentado os resultados provenientes da aplicação do framework,

5.1.7. Restrições do Mecanismo

Apresentam-se algumas restrições para o funcionamento do mecanismo, tornando-o dependente do usuário utilizador, no caso, o perito. Citam-se:

- O mecanismo suporta a investigação de qualquer crime que contenha um dicionário criminoso formalizado;
- Necessita-se da formulação de uma base de *e-mails* de Treinamento. É necessário que o utilizador do mecanismo “adapte” esta base, onde deverá existir uma quantidade de *e-mails* relacionados ao crime que se investiga, assim como *e-mails* que não estejam vinculados ao crime que se investiga;
- Alguns métodos de agrupamento/classificação que podem ser aplicados na FASE II ou FASE III podem precisar sofrer ajustes em seus parâmetros por parte do perito, como por exemplo, no uso do método *K-means*. Para funcionamento do método *K-means* para agrupamento de conversações, o perito deverá especificar o número de conversações existentes em uma base de *e-mails* sendo investigada;

Na próxima Seção apresentam-se Resultados dos experimentos realizados na base de *e-mails* de Teste a partir da base de *e-mails* de Treinamento, cuja formação foi referenciada no dicionário e base de *e-mails* com crime de assédio moral [NUN09]. Expõem-se os resultados e comentários referentes à aplicação do *framework Reader*.

5.2. Resultados Experimentais

Nesta Seção informam-se os resultados e comentários obtidos durante a realização de experimentos. Aplicou-se *Reader* na investigação de crime de assédio moral da base de *e-mails* de Teste (Subseção 4.3.3).

5.2.1. Agrupamento de Conversações

Nesta Subseção apresentam-se resultados do agrupamento de conversações no uso de métodos de agrupamento e de classificação. Os métodos foram aplicados na base de *e-mails* de Teste. A base foi lida digitalmente e pré-processada por *Reader*, para então realizar-se a extração de atributos e representação em vetores (FASE II). Assim, originaram-se 9 arquivos de atributos diferentes. Os arquivos foram detalhados na Subseção 4.2.2.

Apresentam-se na Tabela 5.1 as taxas de acerto (percentagem) para agrupamento de conversações no uso dos métodos de agrupamento *K-means* usando diferentes funções de distâncias. São elas: *Euclidiana*, *Manhattan*, *Cossine* e *Jaccard*. Mediante aplicação dos métodos de agrupamento, experimentou-se o agrupamento de cada um dos 9 arquivos (base Teste).

Os resultados do agrupamento alcançaram taxa de 81,23% no uso do arquivo de atributos contendo as palavras do campo assunto submetido ao método *K-means Manhattan*. Ademais, para os arquivos que fizeram uso das palavras ou expressões do corpo dos *e-mails*, o resultado mais alto alcançado foi de 43,34% (*K-means Euclidiana*), sendo que este valor aumentou a partir dos testes realizados com o arquivo palavras do corpo e assunto dos *e-mails*. Através da junção de palavras (corpo e assunto), obteve-se uma taxa de 50,35% (*K-means Euclidiana*). E ainda, lembra-se que os resultados proporcionados pela utilização dos arquivos de atributos com palavras/expressões envolvendo o campo assunto foram superiores em relação ao demais (Corpo), haja vista que 24 conversações mapeadas possuíam as mesmas palavras no campo assunto de seus *e-mails*, o que não ocorre na prática. E também, os resultados proporcionados pelo método *K-means Jaccard* são frutos de seu cálculo, já que *Jaccard* faz apenas o agrupamento de *e-mails* similares, ou não.

Analisados os 9 arquivos de atributos pelos métodos de agrupamento, estes foram submetidos aos métodos de classificação SVM com *kernels* (*Polynomial*, *Radial* e *Sigmoid*), NB e DT. Expõem-se na Tabela 5.2 os resultados das taxas de acertos (percentagem) referentes ao agrupamento de conversações no uso dos métodos de classificação.

Considerou-se o arquivo de atributo testado como seu arquivo de atributo treinamento também. Assim, o classificador agrupou conversações baseado no arquivo de treinamento, que neste caso, foi o próprio arquivo testado.

Tabela 5.1 – Agrupamento de Conversações - Agrupamento

Arquivo de Atributos	<i>K-means Euclidiana</i>	<i>K-means Cossine</i>	<i>K-means Manhattan</i>	<i>K-means Jaccard</i>
Palavras Corpo – tf-idf	41,58	44,22	42,28	15,79
Expressões Corpo – tf-idf	34,04	54,04	35,62	15,79
Palavras Assunto – tf-idf	79,30	57,37	81,23	60,00
Expressões Assunto – tf-idf	58,77	47,37	59,13	57,00
Palavras Corpo – <i>bag of words</i>	43,34	28,57	37,20	15,79
Expressões Corpo – <i>bag of words</i>	35,04	54,04	35,62	15,79
Palavras Assunto – <i>bag of words</i>	78,94	53,51	79,30	77,00
Expressões Assunto – <i>bag of words</i>	58,77	47,37	59,13	57,67
Palavras Corpo e Assunto – <i>bag of words</i>	50,35	42,11	47,54	49,84

Dentre os métodos de classificação aplicados no agrupamento de conversações, o método de classificação NB foi o método que alcançou as melhores taxas. Nos arquivos de atributos que fizeram uso das palavras ou expressões do corpo dos *e-mails*, utilizando a representação via *bag of words*, o classificador NB obteve, respectivamente, 98,07% e 97,90% de acertos. E no uso dos arquivos de atributos contendo as palavras e expressões do campo assunto dos *e-mails*, o classificador NB obteve taxas de acertos na ordem de 97,00 a 98,00%, tanto na representação via tf-idf como *bag of words*. E no arquivo com a junção das palavras do corpo e do campo assunto, usando representação via *bag of words*, NB obteve 97,90% de acertos.

O classificador SVM (*kernel Polynomial*) obteve bom resultado na classificação para agrupamento de conversações quando utilizado o arquivo de atributos com as palavras do corpo dos *e-mails* representados via *bag of words*, tendo obtido 97,55%, e ainda, o mesmo classificador teve taxa de acertos equivalente a 97,20% no uso do arquivo contendo as palavras do corpo e do campo assunto dos *e-mails* (representação *bag of words*).

O uso do arquivo de atributos com as palavras encontradas no campo assunto dos *e-mails* proporcionou boas taxas de acertos, tanto nos resultados apresentados na Tabela 5.1 como na Tabela 5.2. Entretanto, as conversações de uma base de *e-mails* não precisam utilizar o mesmo conjunto de palavras no campo assunto dos *e-mails*. *E-mails* de uma conversação podem ter as palavras/expressões do campo assunto modificados.

Nos experimentos para agrupamento de conversações de *e-mail* utilizando *K-means Cossine*, verificou-se que ao contrário do agrupamento de textos, este não é um método apropriado para agrupamento de conversações em base de *e-mails*. Textos podem conter uma

grande quantidade de palavras, o que não ocorre necessariamente com *e-mails*. Desta forma, contrariando os bons resultados apresentados por [NAG10].

Tabela 5.2 – Agrupamento de Conversações – Classificação

Arquivo de Atributos	SVM <i>Polynomial</i>	SVM <i>Radial</i>	SVM <i>Sigmoid</i>	DT	NB
Palavras Corpo – tf-idf	15,44	15,78	15,78	71,41	97,72
Expressões Corpo – tf-idf	42,46	49,65	42,80	67,90	97,90
Palavras Assunto – tf-idf	15,44	15,78	15,78	54,22	98,94
Expressões Assunto – tf-idf	15,44	15,78	15,79	43,68	78,60
Palavras Corpo – <i>bag of words</i>	97,55	27,36	15,78	72,63	98,07
Expressões Corpo – <i>bag of words</i>	48,25	50,35	42,80	67,90	97,90
Palavras Assunto – <i>bag of words</i>	93,16	94,73	75,96	54,22	98,94
Expressões Assunto – <i>bag of words</i>	42,29	26,49	15,78	43,86	78,60
Palavras Corpo e Assunto – <i>bag of words</i>	97,20	35,96	15,78	71,23	97,90

Ademais, os resultados apresentados para o agrupamento de conversações, fazem corroborar com os métodos utilizados por [CSE06] [BAL08a] [BAL08b] e [DRE06]. Embora estes autores não tenham trabalhado com o agrupamento de conversações, utilizaram o método de classificação NB para classificação de *e-mails*.

Por fim, o método de classificação NB foi o melhor método em termos de taxas de acerto para o agrupamento de conversações. Na próxima Subseção, os *e-mails* da base de Teste sofrem uma nova classificação, entretanto para verificar a existência, ou não, de crime de assédio moral.

5.2.2. Classificação das Conversações

Nesta Subseção, a base de *e-mails* de Teste é classificada em crime de assédio moral, ou não. Para isso, *Reader* realiza a extração de atributos da base de *e-mails* Teste e de uma base de *e-mails* de Treinamento, e os representam em vetores (FASE III). Estes vetores, também chamados de arquivos de atributos, foram apresentados na Subseção 4.2.3. Obtiveram-se naquela Seção, 8 arquivos. Sendo que 4 arquivos são obtidos a partir da base de *e-mails* Teste, e outros 4 a partir da base de *e-mails* de Treinamento.

Os arquivos de atributos da base Teste foram submetidos aos classificadores. Os classificadores, baseados nos padrões existentes nos arquivos de atributos da base de

Treinamento, realizaram a classificação de cada *e-mail* da base de Teste nas classes crime ou não crime de assédio moral.

Na Tabela 5.3 apresentam-se as taxas de acertos alcançadas na classificação dos *e-mails* da base de Teste. Nesta FASE, o mecanismo avalia a taxa de classificação correta dos métodos de classificação: SVM com *kernels* (*Linear*, *Polynomial*, *Radial* e *Sigmoid*), DT e NB.

Tabela 5.3 – Classificação de Conversações

Arquivo de Atributos	SVM <i>Linear</i>	SVM <i>Polynomial</i>	SVM <i>Radial</i>	SVM <i>Sigmoid</i>	DT	NB
3-grams 67%	72,28	69,82	57,01	80,00	99,64	99,64
4-grams 100%	72,28	70,00	65,43	95,61	99,82	99,82
3-grams 100%	94,73	76,14	23,68	82,98	94,38	85,08
4-grams 67%	94,91	77,54	39,12	90,35	94,91	95,08

Analisando os resultados da Tabela 5.3, verifica-se que os classificadores NB e DT obtiveram resultados similares na classificação dos *e-mails* da base de Teste no uso dos arquivos de atributos 4-grams – 100% e 3-grams – 67%. O melhor resultado foi obtido aplicando-se o classificador NB, o qual obteve 99,82% de taxa de acerto no uso do arquivo de atributos formalizado pela técnica 4-grams 100%. E o classificador SVM com *kernel Linear* também se mostrou como uma boa opção para a classificação dos *e-mails* de Teste em duas classes, principalmente nos uso dos arquivos de atributos onde foram aplicadas as técnicas de 3-grams - 100% e 4-grams - 67%. Embora os resultados obtidos pelo classificador SVM (*kernel Linear*) possam ser melhorados, principalmente a partir de um melhor treinamento na base de Treinamento.

Os resultados apresentados na Tabela 5.3 fazem concordar com os resultados expostos na Seção Trabalhos Relacionados (Capítulo 3), onde se utilizaram classificadores para categorização de *e-mails* [CSE06] [BAL08a] [BAL08b] [DRE06]. E baseando-se nos resultados expostos nesta Subseção, apresenta-se a opção em utilizar o classificador NB para efetivação das fases do mecanismo, já que NB apresentou o melhor resultado na classificação dos *e-mails* de Teste (99,82%), no uso do arquivo de atributos gerado pela técnica de 4-grams 100%.

Assim, os *e-mails* de Teste foram classificados em crime e não crime. Unindo os resultados do agrupamento de conversações com a classificação de *e-mails* realizada nesta

Subseção, resulta-se na classificação das conversações, as quais serão visualizadas mediante aplicação de grafos direcionados (próxima Subseção). Mediante esta visualização, proporciona-se o rastreamento de conversações existentes em bases de *e-mail*.

5.2.3. Apresentação dos Resultados Experimentais

No uso dos resultados proporcionados pelo agrupamento e classificação das conversações, *Reader* constrói grafos para cada conversação agrupada a partir da matriz de adjacências. Para uma base de *e-mails* investigada com n conversações, *Reader* propicia a construção de n grafos.

Para formação da matriz de adjacências, utilizam-se os contatos encontrados nos campos *To* e *From* existentes no cabeçalho dos *e-mails* da base Teste.

Na Figura 5.2 visualiza-se um dos grafos elaborados por *Reader* após o agrupamento e classificação de conversações na base de Teste. A conversação se refere ao grupo (conversação) 31, a qual possui um total de 13 contatos, sendo que os usuários *User08* e *User11*, *User08* e *User29*, *User08* e *User09* e *User08* e *User18* trocaram *e-mails* durante esta conversação, ao contrário dos outros usuários, que apenas receberam *e-mails*. Nenhum destes contatos enviou ou recebeu *e-mail* criminoso, caracterizando a conversação como não criminosa. Em função disto, o grafo foi colorido em azul.

Apresenta-se na Figura 5.3 o grafo gerado após análise da conversação 19 da base de Teste, sendo que esta conversação apresenta um total de 6 contatos. Observa-se que o contato *User1* enviou *e-mails* para 3 contatos, sendo que nenhum destes contatos respondeu seus *e-mails*. E os contatos *User04* e *User02* deram continuidade nesta conversação não criminosa.

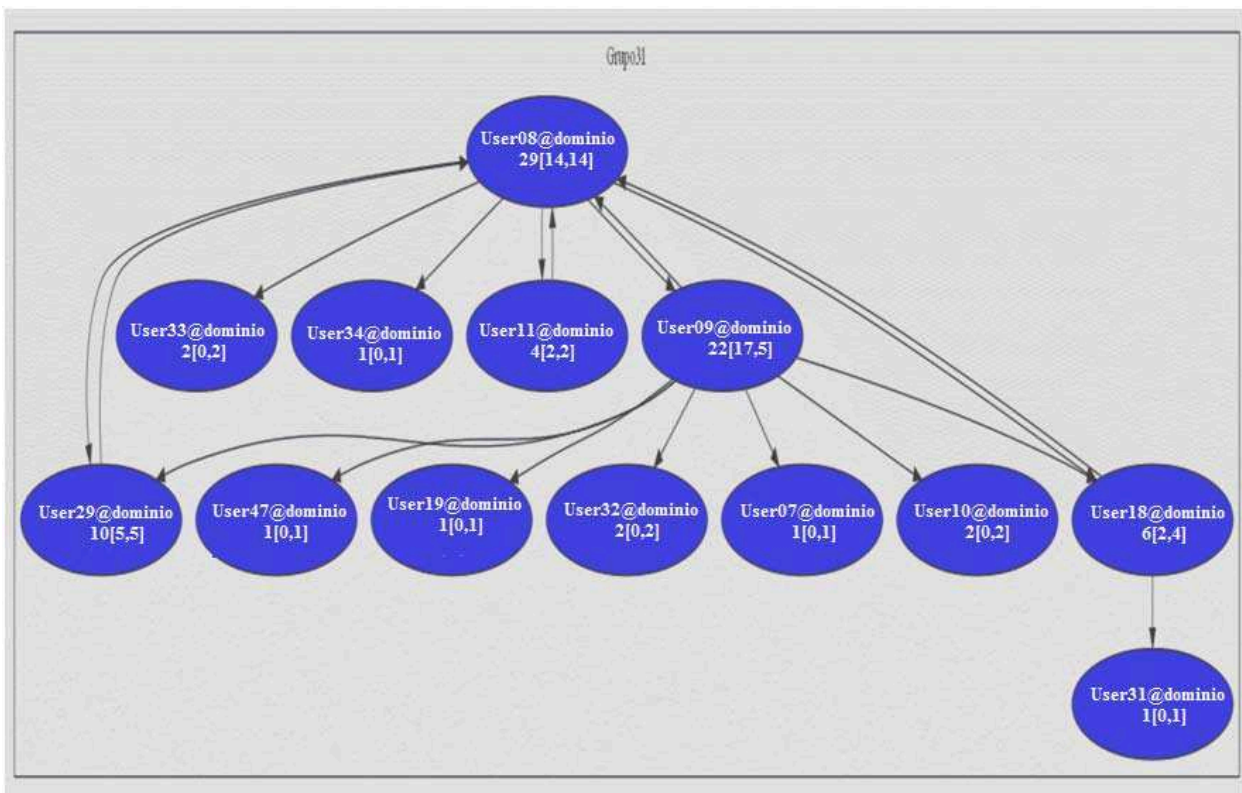


Figura 5.2 – Gráfico Não Criminoso – Base Teste – Exemplo 1

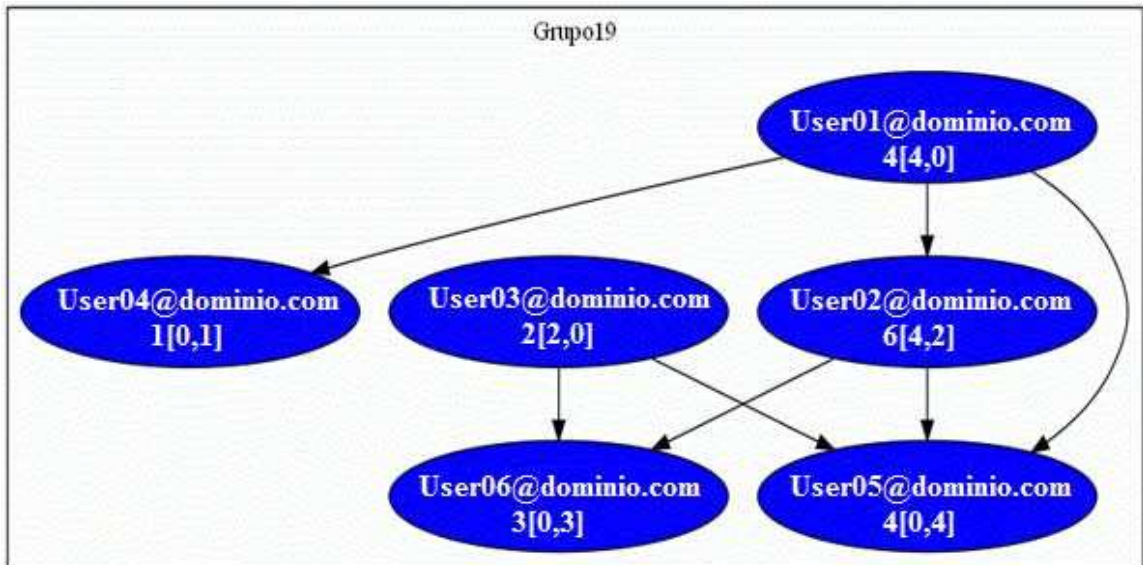


Figura 5.3 – Grafo Não Criminoso – Base Teste – Exemplo 2

Uma conversa o criminoso (grafo em cor vermelha)   apresentada na Figura 5.4, representando a conversa o 29 da base de Teste. Verifica-se que os contatos User29 e User08 trocaram *e-mails*, entretanto os demais contatos apenas receberam *e-mails*.

Atrav s da constru o de grafos a partir do relacionamento em conversa es por *e-mails*, existe a facilidade de entendimento e rastreamento de conversa es e contatos. Um *e-mail* recebido, por exemplo, pode ter sido enviado para um ou mais contatos. Na Figura 5.4, o contato User29 enviou um *e-mail* para o contato User36, User08 e para User35, ao mesmo tempo.

Uma conversa o criminoso pode ter apenas um *e-mail* criminoso. Em fun o disto, todo o grafo ser  colorido em vermelho. Para facilitar a continua o da an lise forense de um perito, o mecanismo prop e que os *e-mails* criminosos sejam listados. Assim, um perito que realize a an lise de evid ncias digitais consegue com facilidade detectar qual ou quais s o os *e-mails* criminosos detectados.

Na Figura 5.5 visualizam-se fragmento da listagem dos *e-mails* criminosos detectados na base de Teste. Eles pertencem a conversa o 29. Para facilitar a an lise, Reader lista os *e-mails* por contatos. Como exemplo cita-se o caso do contato A enviar um *e-mail* criminoso para B e para C. Por isso, s o listados 2 *e-mails* criminosos, e n o apenas um.

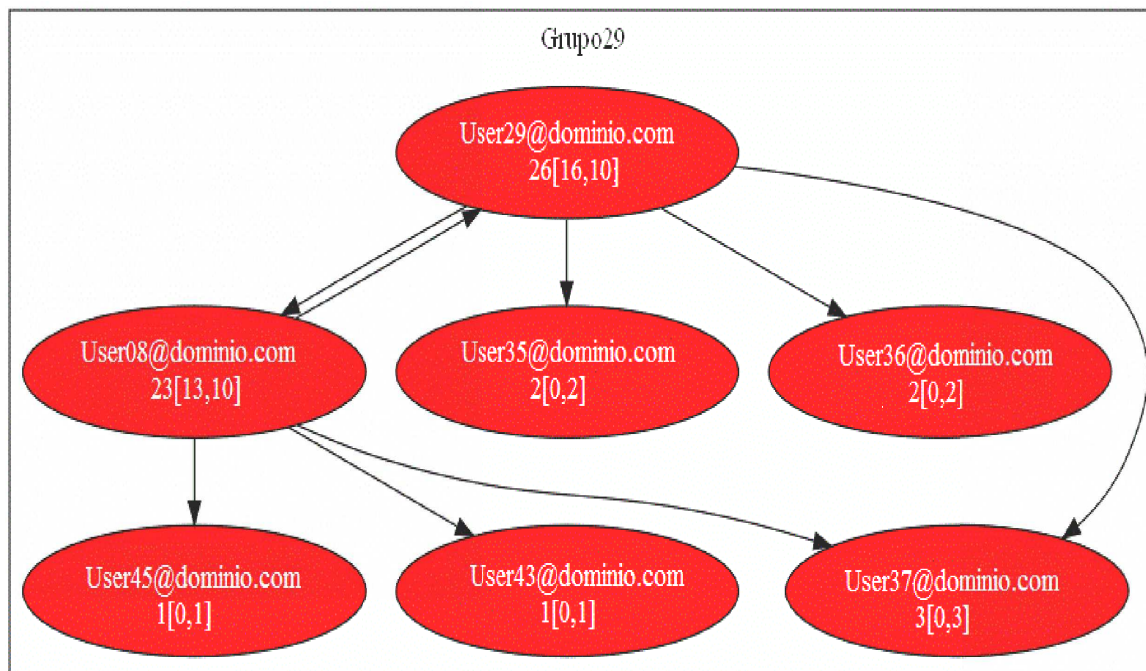


Figura 5.4 – Grafo Criminoso – Base Teste – Exemplo 1

Na prática, observa-se que na conversação 29, foram detectados por *Reader*, 19 *e-mails* criminosos. Entretanto o mecanismo propicia a listagem de 28 *e-mails* criminosos, já que um *e-mail* pode ser enviado para mais de um contato.

Lista de <i>e-mails</i> Criminosos	Usuário Remetente	Usuário Destinatário	Hora/Data	Assunto do <i>e-mail</i>
↓	↓	↓	↓	↓
Email enviado por User29@dominio.com para User08@dominio.com no dia 2010-12-03 01:08:24 cujo assunto é RE: Conversação	User29@dominio.com	User08@dominio.com	2010-12-03 01:08:24	RE: Conversação
Email enviado por User08@dominio.com para User29@dominio.com no dia 2010-12-04 00:26:38 cujo assunto é RE: Conversação	User08@dominio.com	User29@dominio.com	2010-12-04 00:26:38	RE: Conversação
Email enviado por User29@dominio.com para User08@dominio.com no dia 2010-12-04 00:02:20 cujo assunto é RE: Conversação	User29@dominio.com	User08@dominio.com	2010-12-04 00:02:20	RE: Conversação
Email enviado por User08@dominio.com para User29@dominio.com no dia 2010-12-02 11:31:21 cujo assunto é Conversação CRI	User08@dominio.com	User29@dominio.com	2010-12-02 11:31:21	Conversação CRI
Email enviado por User08@dominio.com para User29@dominio.com no dia 2010-12-04 01:34:11 cujo assunto é Re: Conversação	User08@dominio.com	User29@dominio.com	2010-12-04 01:34:11	Re: Conversação
Email enviado por User29@dominio.com para User08@dominio.com no dia 2010-12-03 20:33:19 cujo assunto é RE: Conversação	User29@dominio.com	User08@dominio.com	2010-12-03 20:33:19	RE: Conversação
Email enviado por User08@dominio.com para User29@dominio.com no dia 2010-12-04 01:21:56 cujo assunto é Re: Conversação	User08@dominio.com	User29@dominio.com	2010-12-04 01:21:56	Re: Conversação
Email enviado por User29@dominio.com para User08@dominio.com no dia 2010-12-04 00:58:17 cujo assunto é RE: Conversação	User29@dominio.com	User08@dominio.com	2010-12-04 00:58:17	RE: Conversação
Email enviado por User08@dominio.com para User29@dominio.com no dia 2010-12-03 02:01:26 cujo assunto é Re: Conversação	User08@dominio.com	User29@dominio.com	2010-12-03 02:01:26	Re: Conversação
Email enviado por User08@dominio.com para User29@dominio.com no dia 2010-12-02 17:59:30 cujo assunto é Re: Conversação	User08@dominio.com	User29@dominio.com	2010-12-02 17:59:30	Re: Conversação
Email enviado por User08@dominio.com para User29@dominio.com no dia 2010-12-02 15:23:11 cujo assunto é Re: Conversação	User08@dominio.com	User29@dominio.com	2010-12-02 15:23:11	Re: Conversação
Email enviado por User29@dominio.com para User08@dominio.com no dia 2010-12-04 00:08:04 cujo assunto é RE: Conversação	User29@dominio.com	User08@dominio.com	2010-12-04 00:08:04	RE: Conversação
Email enviado por User08@dominio.com para User29@dominio.com no dia 2010-12-04 01:03:49 cujo assunto é Re: Conversação	User08@dominio.com	User29@dominio.com	2010-12-04 01:03:49	Re: Conversação
Email enviado por User08@dominio.com para User43@dominio.com no dia 2010-12-04 01:03:49 cujo assunto é Re: Conversação	User08@dominio.com	User43@dominio.com	2010-12-04 01:03:49	Re: Conversação
Email enviado por User29@dominio.com para User08@dominio.com no dia 2010-12-03 23:51:30 cujo assunto é RE: Conversação	User29@dominio.com	User08@dominio.com	2010-12-03 23:51:30	RE: Conversação
Email enviado por User29@dominio.com para User08@dominio.com no dia 2010-12-02 14:26:03 cujo assunto é RE: Conversação	User29@dominio.com	User08@dominio.com	2010-12-02 14:26:03	RE: Conversação
Email enviado por User29@dominio.com para User35@dominio.com no dia 2010-12-02 14:26:03 cujo assunto é RE: Conversação	User29@dominio.com	User35@dominio.com	2010-12-02 14:26:03	RE: Conversação

Figura 5.5 – Fragmento da Lista de *e-mails* Criminosos Detectados

Para melhor entendimento das conversações criminosas, *Reader* proporciona listar as palavras criminosas mais utilizadas dentro de uma conversação criminosa. Na Figura 5.6 visualiza-se a detecção parcial realizada para a conversação 29 da base de Teste. Para cada *e-mail* criminoso detectado, apresenta-se na Figura 5.6 o número de ocorrência de cada palavra criminosa detectada, conforme o dicionário de crimes de assédio moral.

Nesta Subsecção foram apresentados resultados que auxiliam um perito na constituição do nexos causal. Entendimento e rastreamento de uma conversação. Na próxima Subsecção apresentam-se análises comparativas do uso de *Reader* e o processo de análise manual, além de análises críticas da aplicação do mecanismo.

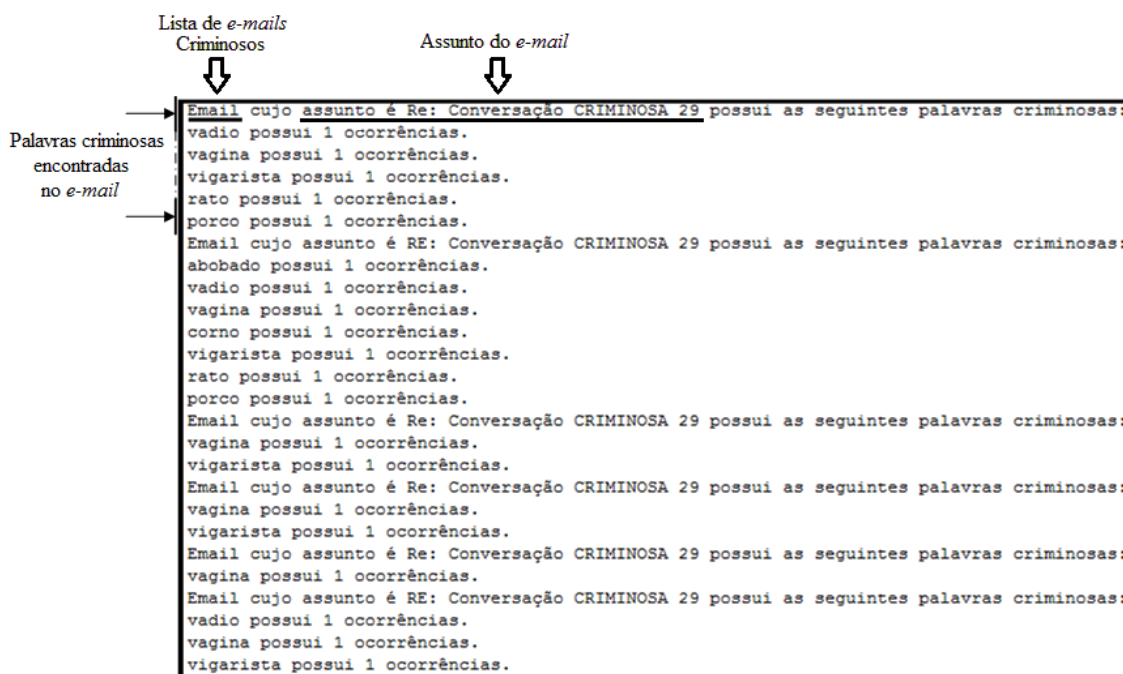


Figura 5.6 – Fragmento da Lista de Palavras Criminosas Utilizadas na Conversação 29

5.2.4. Análise dos Resultados

Os resultados expostos na Subsecção 5.2.1, demonstram ter sido NB o método de classificação que obteve as melhores taxas para agrupamento de conversações no uso da base de Teste. Mediante análise, observa-se que os métodos de agrupamento aplicados para agrupamento de conversações não obtiveram bom resultado, ao contrário de seus usos na classificação textual, como relatado nos Trabalhos Relacionados (Capítulo 3).

E na classificação criminal (Subsecção 5.2.2) dos *e-mails* da base de Teste, obtiveram as melhores taxas de acertos de classificação em duas classes (crime e não crime), os métodos

NB, DT e SVM com *kernel Linear*. Os resultados fazem concordar com Trabalhos Relacionados apresentados (Capítulo 3).

Tendo-se os resultados aplicados por *Reader*, comparam-se neste momento os tempos de análise da base de *e-mails* de Teste na utilização de *Reader* e através da análise manual que um perito gastaria na produção de provas digitais. Os tempos registrados estão expostos na Tabela 5.4. Para os testes utilizando *Reader*, utilizou-se um computador provido de processador Intel (R) Core (TM) Duo CPU 2.20 GHz, com 4 G de RAM.

Necessitou-se de aproximadamente 2 horas, 07 minutos e 25 segundos, com média de 13,41 segundos para processamento de cada *e-mail* da base de Teste para a análise realizada por *Reader*. Considerou-se a utilização dos métodos com a melhor obtenção de resultados (Subseção 5.2.1 e 5.2.2). Citam-se:

- FASE I: Leitura digital, processo de “imagem” e pré-processamento da base de Teste;
- FASE II: Formação do arquivo de atributos referente à base de Teste (Palavras Corpo – *bag of words*) e agrupamento das conversações mediante aplicação do método NB;
- FASE III: Formação do arquivo de atributos referente à base de Teste aplicando a técnica de *4-grams* - 100%, e após, a classificação da base de Teste tendo-se o arquivo de atributos similar da base de Treinamento já formado. Nesta análise foi utilizado o método de classificação DT;
- FASE IV: Apresentação dos grafos com todas as conversações agrupadas (criminosas e não criminosas), além de todas as respostas apresentadas na Subseção 5.3.3.

Na análise manual da base de Teste, foi necessário em média, 60 segundos para realizar o mesmo processo (leitura, agrupamento das conversações, classificação dos *e-mails* de cada conversação e análise das conversações resultantes), totalizando em 9 horas e 50 minutos.

Comparando os testes entre as duas formas de análise da base de Teste, constata-se que o mecanismo implementado em *Reader* reduziu o tempo de análise em aproximadamente 447% em relação ao tempo que um perito utilizaria no mesmo tipo de análise. Observa-se que a diminuição de tempo proporcionado pelo mecanismo é de grande importância para peritos, haja vista que bases de *e-mails* podem ter uma grande quantidade de *e-mails*, chegando a ordem de GBytes. Além do que, na análise manual de uma grande base de *e-mail* possibilita o

cometimento de erros por parte dos peritos, como por exemplo, na análise de uma amostra de uma base suspeita onde não contenham os *e-mails* criminosos.

Além do menor tempo proporcionado pela aplicação do mecanismo, expõem-se a facilidade para rastreabilidade de conversações. Mediante análise de um grafo, tornam-se simples o entendimento de conversações criminosas agrupadas, já que o mecanismo proporciona a visualização das conversações. Já na análise manual, isto dependerá da experiência daquele que fará a investigação criminal de uma base de *e-mails*.

Tabela 5.4 – Comparação entre *Reader* e Processamento Manual

FASE	Tempo <i>Reader</i> (hh:mm:ss)	Tempo Manual
I	00:04:11	09:50:00
II	00:02:48	
III	01:55:30	
IV	00:04:36	
TOTAL	02:07:25	
MÉDIA	13,41 segundos	60,00 segundos

Desta forma, este Capítulo apresentou os resultados obtidos nos experimentos realizados a partir da investigação na base de *e-mails* de Teste, composta por 570 *e-mails* divididos em 33 conversações. Nesta base investigou-se o agrupamento e a classificação das conversações (crime ou não crime), baseando-se no dicionário de crimes com assédio moral. Os resultados do agrupamento e classificação foram expostos e representados no uso de grafos, além da apresentação de outros resultados. Além do mais, fez-se uma análise entre os resultados proporcionados pela aplicação do mecanismo e de peritos, para análise da base de *e-mails* de Teste. No próximo Capítulo apresentam-se as conclusões, assim como a indicação de trabalhos futuros que podem ser originados a partir deste estudo de Mestrado.

Capítulo 6

Conclusão

Atualmente, verifica-se na atual sociedade que a utilização da ferramenta eletrônica *e-mail* e o cometimento de cibercrimes esta em aumento. Crimes cometidos no uso de *e-mail* também precisam ser denunciados aos órgãos competentes para que infratores possam ser penalizados. Para isso, os órgãos investigam os resquícios dos crimes cometidos e formalizam o nexu causal (vínculo entre a conduta ilícita e o dano).

A investigação criminal em *e-mails* é um processo árduo, visto a grande quantidade de mensagens eletrônicas que usuários possuem. No trabalho de investigação, peritos da Computação Forense precisam encontrar com rapidez e precisão as provas que evidenciem o crime cometido. Visto isto, verifica-se a necessidade da existência de mecanismos que auxiliem a formalização do nexu causal a partir de repositórios de *e-mail*.

Assim, apresentou-se neste trabalho de Mestrado um mecanismo semi-automático para produção de provas digitais a partir do rastreamento em relacionamentos existentes em base de *e-mails*, investigando conversações trocadas entre o usuário proprietário e seus contatos.

As conversações são agrupadas pela existência da troca de mensagens e que compartilhem de um mesmo conjunto textual no corpo dos *e-mails*. O agrupamento das conversações baseia-se no conjunto textual encontrado no corpo de *e-mails*. Desta forma, conversações são agrupadas pelo mecanismo mesmo quando usuários se utilizam de diferentes endereços de *e-mails*. Agrupada as conversações, estas são classificadas.

A classificação se baseia no uso de um dicionário criminal de palavras e de uma base de Treinamento que contenham *e-mails* com registros do crime que esta sendo investigado,

assim como do contrário. Ao se obter a classificação dos *e-mails* das conversações em criminoso ou não, têm-se os resultados comprobatórios das evidências digitais.

Os resultados do agrupamento e da classificação consistem, dentre eles, da representação gráfica. Ainda, cita-se que para exposição dos resultados da aplicação do mecanismo, utilizam-se alguns campos existentes no cabeçalho dos *e-mails*.

Para aplicação do mecanismo fizeram-se vários estudos. Técnicas de pré-processamento textual (*case foldering*, *stop-word* e *n-grams*) e da extração e representação de características (*tf-idf* e *bag of words*), de métodos de agrupamento/classificação (*K-means - Euclidiana*, *Manhattan*, *Cossine* e *Jaccard*, NB, AD e SVM) e de técnicas para representação gráfica das conversações (grafos direcionados).

Para validação do mecanismo, implementou-se um *framework* desenvolvido em linguagem *Java*. Este *framework* se chama *Reader*, o qual analisa base de *e-mails* provenientes do *software Windows Live Mail*, Sistema Operacional *Windows*. E para obtenção de resultados experimentais, originou-se uma base composta por 570 *e-mails* divididos em 33 conversações.

Nos experimentos objetivou-se o agrupamento de conversações com indícios de crime de assédio moral. Para tal, aplicou-se o dicionário e a base de *e-mails* contendo palavras caracterizadoras de crime de assédio moral [NUN09]. Ainda, formalizou-se uma base de Treinamento com 3.816 *e-mails*, divididos na razão 50:50, ou seja, 50% dos *e-mails* continham evidências de crime de assédio moral, e os outros 50%, não.

Desta forma, comprovou-se a eficiência do mecanismo, já que os resultados atingiram taxas de acertos próximas de 99% na utilização do classificador NB. Os resultados experimentais demonstram ser possível obter evidências digitais quanto à autoria e/ou quanto à materialidade de possível delito. De tal modo, as análises realizadas podem ser formalizadas pelos membros da ciência jurídica de modo a caracterizar o nexos causal.

Aplicando-se o mecanismo, obtiveram-se resultados conclusivos para formalização do nexos causal. Além do mais, durante os estudos de Mestrado, analisaram-se ferramentas computacionais para análise de *e-mails*. Constatou-se a falta de dispositivos eletrônicos para este tipo de análise no mercado forense e que utilizem métodos científicos investigatórios, além das técnicas triviais, com a análise textual individual de *e-mails*.

Também, verificou-se o tempo para análise do mecanismo em relação ao processamento manual de uma mesma base de *e-mails* durante os experimentos. Os resultados

demonstraram que o mecanismo desempenhou a tarefa de processamento da base de *e-mails* utilizada num tempo inferior a 447% aquele proporcionado pela análise manual que um perito realizaria.

No mecanismo apresentado neste trabalho, cabe a indicação de trabalhos futuros. Dentre eles, a formulação de novas práticas computacionais no *framework* desenvolvido. Apresentação dos grafos de forma diferenciada, podendo o usuário do *framework* interagir mais rapidamente após a construção dos grafos, como por exemplo, clicando no grafo de uma conversa e visualizando o seu conteúdo.

Até se podem utilizar outros campos do cabeçalho dos *e-mails* que o *framework* analisa, como por exemplo, *e-mails* enviados com cópia oculta. Além do mais, aperfeiçoar o mecanismo através da junção de funcionalidades entre os diferentes métodos de classificação estudados, e também testar outros métodos para aplicação do mecanismo no agrupamento e classificação das conversações, possibilitando que os resultados apresentados pelo mecanismo possam ser aperfeiçoados.

Algumas das funcionalidades citadas como trabalhos futuros já estão em desenvolvimento, proporcionando a existência de novos conteúdos para a publicação de novos artigos.

Referências Bibliográficas

- [BAL08a] BALAMURUGAN, S.A.A.; RANJARAM, R. *Learning to classify threatening e-mail*. Int. J. Artificial Intelligence and Soft Computing, Vol. 1, No. 1, 2008, p.39-51.
- [BAL08b] BALAMURUGAN, S.A.A.; RANJARAM, R.; MUTHUPANDIAN, M.; ATHIAPPAN, G. *Automatic mining of Threatening e-mail using Ad Infinitum algorithm*. Int. J. of Information Technology, Vol. 14, No. 2, 2008, p.81-108.
- [BRO06] BROADHURST, R. *Developments in the global law enforcement of cyber-crime*. Policing: An Int. Journal of Police Strategies & Management. Vol. 29, No. 3, 2006, p.408-433.
- [CAM09a] CampaignMonitor. *Most popular email clients in 2009*. Disponível em <<http://www.campaignmonitor.com/stats/email-clients/>>. Acesso em 22 de dezembro de 2010.
- [CHO06] CHOPRA, M.; MARTIN, M. V.; RUEDA, L.; HUNG, P. C. K. *Toward new Paradigms to Combating Internet Child Pornography*. Canadian Conference on Electrical and Computer Engineering, CCECE, Ottawa, Maio 2006, p.1012-1015.
- [CIA04] CIA, S. Ó. *An Extended Model of Cybercrime Investigations*. Int. Journal of Digital Evidence, 2004, Vol. 3, Issue 1.
- [CSE06] CSELLE, G. *Organizing email*. Master's thesis, ETH Zurich, 2006.

- [DAB05] DABBISH, L. A.; KRAUT, R. E.; FUSSEL, S.; KIESLER, S. *Understanding Email Use: Predicting Action on a Message*. Proc. of the SIGCHI conference on Human factors in computing systems, USA, Abril 2005, p.190-197.
- [DAT10] DATA, A. *User Guide Forensic Toolkit. Find, Organize, & Analyse Computer Evidence*. Disponível em <<http://accessdata.com/>>. Acesso em 21 de janeiro de 2011.
- [DEI05] DEITEL, H. M. *Java: como Programar?* Editora Pearson Prentice Hall. 6ª. Edição, São Paulo. 2005.
- [DEL00] DELMANTO, C. *Código penal comentado*. Editora Renovar. 5ª. Edição atual. e ampl., Rio de Janeiro. 2000.
- [DEP99] DEPUTADOS, D. C. *Projeto de Lei 84/99*. Disponível em <http://imagem.camara.gov.br/dc_20.asp?selCodColecaoCsv=D&Datain=11/5/1999&txpagina=19975&altura=700&largura=800>. Acesso em 21 de janeiro de 2011.
- [DRE06] DREDZE, M.; LAU, T.; Kushmerick, N. *Automatically Classifying Emails into Activities*. Proc. of the 11th International Conference on Intelligent User Interfaces, Sydney, Janeiro, 2006, p.70-77.
- [DUA04] DUANE, A.; FINNEGAN, P. *Managing Email Usage: A cross Case Analysis of Experiences with Electronic Monitoring and Control*. Proc. of the 6th International Conference on Electronic Commerce [Table of Contents](#), 2004, p.229-238.
- [FIL10] FILHO, J. M. A. *Legislação e as Unidades Especializadas em cibercrime no Canadá*. Disponível em <<http://mariano.delegadodepolicia.com/o-cibercrime-no-canada-e-%E2%80%9Chigh-tech-crime-unit%E2%80%9D/>>. Acesso em 03 de março de 2010.
- [FIS06] FISHER, D.; BRUSH, A. J.; GLEAVE, E.; SMITH, M. A. *Revisiting Whittaker & Sidner's "email overload" ten years later*. Proc. of the 2006 20th Anniversary

Conference On Computer Supported Cooperative Work: CSCW'06, Banff, Alberta, Canada, Novembro 2006, p.309-312.

[GOD65] GOOD, I. J. *The Estimation of Probabilities: An Essay on Modern Bayesian Methods*. M.I.T. Press, 1965.

[GRO06] GROUP, R. *Microsoft Exchange market share statistics*. Disponível em <<http://www.radicati.com>>. Acesso em 21 de janeiro de 2011.

[HAM89] HAMERS, L.; HEMERYCK, Y.; Herweyers, G.; JANSSEN, M; KETERS, H.; ROUSSEAU, R.; VANHOUTTE, A. *Similarity measures in scientometric research: the Jaccard index versus Salton's cosine formula*. Information Processing and Management, Vol. 25, N. 3, 1989, p. 315-318.

[HUA08] HUANG, A. *Similarity Measures for Text Document Clustering*. Proc. of the Sixth New Zealand Computer Science Research Student Conference NZCSRSC2008, Christchurch, New Zealand, 2008, p.49-56.

[IQB10] IQBAL, F.; KHAN, L. A.; FUNG, B. C. M.; DEBBABI, M. *E-mail Authorship Verification for Forensic Investigation*. SAC'10 March 22-26, 2010, Sierre, Switzerland, p.1591-1598.

[JAI99] JAIN, A. K.; MURTY, M. N.; FLYNN, P. J. *Data Clustering: A Review*. ACM Computing Surveys, Vol. 31, No. 3, September 1999.

[JES02] JESUS, D. E. *Direito Penal*. Editora Saraiva. São Paulo. 2002.

[KUM06] KUMAR, V., TAN, P., e STEINBACH, M. *Cluster Analysis: Basic Concept and Algorithms*. Capítulo 8, p. 487-586.

[KUR10] KUROSE, J. F. *Redes de Computadores e a Internet: Uma Abordagem Top-Down*. Editora Addison Wesley. 5ª. Edição, São Paulo, 2010.

- [LIM07] LIMEIRA, T. M. V. *E-Marketing*. Editora Saraiva. 2ª. Edição, São Paulo, 2007.
- [LOR07] LORENA, A. C.; CARVALHO, A. C. P. L. F. *Uma introdução às Support Vector Machines (In Portuguese)*. Revista de Informática Teórica e Aplicada, v. 14, p. 43, 2007.
- [MAL10] MALLMANN, J.; FREITAS, C. O. A.; SANTIN, A. *Produção de Provas Digitais a partir de Rastreamento em Relacionamentos por e-mails*. In: X Simpósio Brasileiro de Segurança da Informação e de Sistemas Computacionais - SBSeg 2010, 2010, Fortaleza. X Simpósio Brasileiro de Segurança da Informação e de Sistemas Computacionais - SBSeg 2010. Porto Alegre : SBC - Sociedade Brasileira de Computação, 2010.
- [MES08] MESQUITA, J. *Crimes na Internet com leis mais rígidas na Argentina*. Disponível em <http://www.leicordem.com.br/crimes-na-internet-com-leis-mais-rigida-na-argentina.html>>. Acesso em 21 de janeiro de 2011.
- [NAG10] NAGWANI, N. K.; BHANSALI, A. *An Object Oriented Email Clustering Model Using Weighted Similarities between Emails Attributes*. Int. Journal of Research and Reviews in Computer Science (IJRRCS), 2010, V. 1, N. 2, p.1-6.
- [NET00] NETO, J. L.; SANTOS, A. D.; KAESTNER, C. A. A.; FREITAS, A. A. *Document Clustering and Text Summarization*. 2000. Artigo disponível em <http://citeseerx.ist.psu.edu/viewdoc/summary?doi=10.1.1.43.4634>>. Acesso em 21 de abril de 2010.
- [NOG09] NOGUEIRA, S. D. *Crimes de Informática*. Editora BH. 2ª. Edição, São Paulo, 2009.
- [NOR98] NORONHA, E. M. *Curso de Direito Processual Penal*. Editora Saraiva. 26ª. Edição, São Paulo, 1998.

- [NUN09] NUNES, A. V.; FREITAS, C. O. A.; PARAÍSO, E. C. *Detecção de Assédio Moral em e-mails*. In: I Student Workshop on Information and Human Language Technology, 2009, São Carlos. I Student Workshop on Information and Human Language Technology - 7th Brazilian Symposium in Information and Human Language Technology. Porto Alegre : SBC, 2009. v. 1. p.01-05.
- [OLI02] OLIVEIRA, F. S.; GUIMARÃES, C.; GEUS, P. L. G. *Resposta a Incidentes para Ambientes Corporativos Baseados em Windows*. Disponível em: <http://www.guiatecnico.com.br/PericiaForense/>. Acesso em 19 de janeiro de 2010.
- [PAU04] PAULO, A. D. *Pequeno Dicionário Jurídico*. Editora DP&A, 2º Edição, Rio de Janeiro, 2004, p. 261.
- [PET04] PETERSON, L. L.; DAVIE, B. S. *Redes de Computadores: uma abordagem de sistemas*. Editora Elsevier. 3ª Edição, Rio de Janeiro, 2004.
- [PHI09] PHILIPP, A.; COWEN, D.; DAVIS, C.; *Hacking Exposed Computer Forensics: Computer Forensics Secrets & Solutions*. McGraw-Hill (Ed.), 2ª Edição, 2009.
- [PIN09a] PINHEIRO, P. P. *Direito Digital*. Editora Saraiva. 3ª Edição, São Paulo, 2009.
- [PIN09b] PINHEIRO, P. P.; MORAES, C. S. *Direito Digit@l no dia-a-dia*. Áudio Livro. Editora Saraiva, São Paulo, 2009.
- [PRE00] PREISS, B. R. *Estrutura de Dados e Algoritmos: Padrões de Projetos orientados a objeto com Java*. Editora Campus, Rio de Janeiro, 2000.
- [PUP10] PUPYREV, S.; TIKHONOV, A. *Analyzing Conversations with Dynamic Graph Visualization*. Intelligent Systems Design and Applications (ISDA), 10th International Conference on, 2010, p. 748-753.

- [RAB93] RABINGER, L.; JUANG, B.H. *Fundamentals of speech recognition*. Prentice Hall Inc., London, UK, 1993, p.506.
- [REI02] REITH, M.; CARR, C.; GUNSCH, G. *An Examination of Digital Forensic Models*. Int. Journal of Digital Evidence, 2002, Vol. 1, Issue 3.
- [SAL88] SALTON, G.; BUCKLEY, C. *Term Weighting Approaches in Automatic Text Retrieval*. Information Processing and Management, 1988, 24, 5, p.513-523.
- [SEB02] SEBASTIANI, F. *Machine Learning in Automated Text Categorization*. ACM Computing Surveys, Vol. 34. No. 1, Março 2002, p.1-47.
- [SEN08] SENADO, A. *Senado aprova projeto que pune crimes praticados pela internet*. Disponível em http://oglobo.globo.com/pais/mat/2008/07/10/senado_aprova_projeto_que_pune_crimes_praticados_pela_internet-547186663.asp. Acesso em: 28 de fevereiro de 2010.
- [SHI08] SHINDER, D. L.; CROSS, MICHAEL. *Scene of the Cybercrime*. Editora Syngress. 2ª. Edição, UK, 2008.
- [SIP99] SIPIOR, J. C.; WARD, B. T. *The Dark Side of Employee Email*. Communications of the ACM, Vol. 42, N°. 7, Julho 1999, p.88-95.
- [SOA08] SOARES, C. P. *O E-mail como prova no Direito Brasileiro*. Monografia, Brasil, 2008, 2ª Edição, Rio de Janeiro, 1995.
- [STA08] STALLINGS, W. *Criptografia e Segurança de Redes*. Editora Pearson Prentice Hall. 4ª Edição, São Paulo, 2008.
- [STJ08] JUSTIÇA, S. T. *Justiça usa Código Penal para combater crime virtual*. Disponível em

<http://www.stj.gov.br/portal_stj/publicacao/engine.wsp?tmp.area=398&tmp.texto=90108>. Acesso em 02 de março de 2010.

- [STU99] STUMVOLL, V. P.; QUINTELA, V.; DOREA, L. E. *Criminalística*. Editora Sagra Luzzatto, Porto Alegre, 1999.
- [SUR05] SURENDRAN, A. C.; PLATT, J. C.; RENSHAW, E. *Automatic Discovery of Personal Topics to Organize Email*. Proc. Second Conference on Email and Anti-Spam (CEAS 2005).
- [TAM05] TAM, P-N.; STEINBACH, M.; KUMAR, V. *Introduction to Data Mining*. Addison-Wesley. 2005.
- [TAM08] TAM, T.; LOURENÇO, A. *Estudo Exploratório para a Organização Automática de Mensagens de Correio Electrónico*. JETC'08, ISEL, Lisboa, Novembro 2008.
- [TRI96] TRIER, O.; JAIN, A. K.; TAXT, T. *Feature extraction methods for character recognition*. Pattern Recognition, V. 29, N. 4, 1996, p.641–662.
- [VAP63] VAPNIK, V.; e LERNER , A. *Pattern recognition using generalized portrait method*. Automation and Remote Control, 24, p.774–780, 1963.
- [VIE06] VIÉGAS, F. B.; GOLDBER, S.; DONATH, J. *Visualizing Email Content: Portraying Relationships from Conversational Histories*. Proc. of the SIGCHI conference on Human Factors in computing systems. Visualization 2, Canada, Abril, 2006, p.979–988.
- [VIL09] VILLATORE, M. A.; FREITAS, C. O. A. *Palavras e expressões dicionarizadas podem coibir, no trabalho, provável assédio moral nas mensagens eletrônicas*. In: XVIII Encontro Nacional do CONPEDI, 2009, Maringá. Anais do XVIII Encontro Nacional do CONPEDI. Porto Alegre : CONPEDI, 2009. v. 1. p. 2615-2629.

- [WEN10] WENDT, E. *Reportagem Ameaça invisível na rede*. Disponível em <<http://www.emersonwendt.com.br/search/label/Polícia%20Civil>>. Acesso em 03 de março de 2010.
- [WHI96] WHITTAKER, S.; SIDNER, C. *Email overload: exploring personal information management of email*. In Proc. of CHI 1996, ACM Press (1996), p.276-283.
- [YAT99] YATES, R. B.; NETO, R. B. *Modern Information Retrieval*. Addison-Wesley (Ed.), New York, 1999.

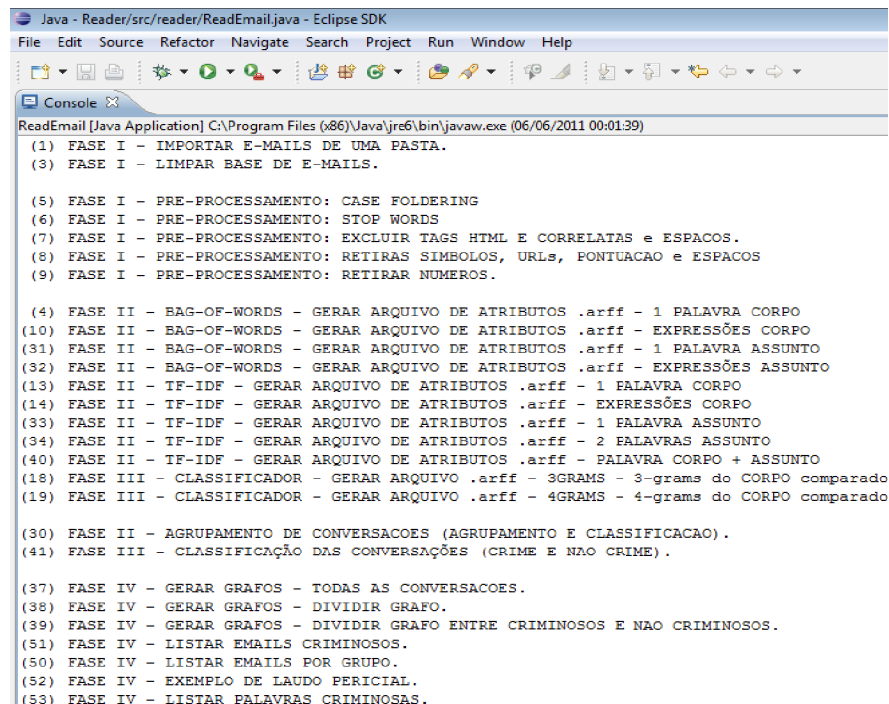
Apêndice A

Apresentação do *Framework Reader*

Este apêndice demonstra a utilização do *framework Reader* e os resultados experimentais apresentados no Capítulo 5 onde a base de *e-mails* Teste e a base de *e-mails* Treinamento foram empregadas.

A.1. Tela Principal

Reader foi implementado tendo sua visualização em modo texto. Na Figura A.1 apresenta-se a tela principal de *Reader*.



```
Java - Reader/src/reader/ReadEmail.java - Eclipse SDK
File Edit Source Refactor Navigate Search Project Run Window Help
Console
ReadEmail [Java Application] C:\Program Files (x86)\Java\jre6\bin\javaw.exe (06/06/2011 00:01:39)
(1) FASE I - IMPORTAR E-MAILS DE UMA PASTA.
(3) FASE I - LIMPAR BASE DE E-MAILS.

(5) FASE I - PRE-PROCESSAMENTO: CASE FOLDING
(6) FASE I - PRE-PROCESSAMENTO: STOP WORDS
(7) FASE I - PRE-PROCESSAMENTO: EXCLUIR TAGS HTML E CORRELATAS e ESPACOS.
(8) FASE I - PRE-PROCESSAMENTO: RETIRAR SIMBOLOS, URLs, PONTUACAO e ESPACOS
(9) FASE I - PRE-PROCESSAMENTO: RETIRAR NUMEROS.

(4) FASE II - BAG-OF-WORDS - GERAR ARQUIVO DE ATRIBUTOS .arff - 1 PALAVRA CORPO
(10) FASE II - BAG-OF-WORDS - GERAR ARQUIVO DE ATRIBUTOS .arff - EXPRESSÕES CORPO
(31) FASE II - BAG-OF-WORDS - GERAR ARQUIVO DE ATRIBUTOS .arff - 1 PALAVRA ASSUNTO
(32) FASE II - BAG-OF-WORDS - GERAR ARQUIVO DE ATRIBUTOS .arff - EXPRESSÕES ASSUNTO
(13) FASE II - TF-IDF - GERAR ARQUIVO DE ATRIBUTOS .arff - 1 PALAVRA CORPO
(14) FASE II - TF-IDF - GERAR ARQUIVO DE ATRIBUTOS .arff - EXPRESSÕES CORPO
(33) FASE II - TF-IDF - GERAR ARQUIVO DE ATRIBUTOS .arff - 1 PALAVRA ASSUNTO
(34) FASE II - TF-IDF - GERAR ARQUIVO DE ATRIBUTOS .arff - 2 PALAVRAS ASSUNTO
(40) FASE II - TF-IDF - GERAR ARQUIVO DE ATRIBUTOS .arff - PALAVRA CORPO + ASSUNTO
(18) FASE III - CLASSIFICADOR - GERAR ARQUIVO .arff - 3GRAMS - 3-grams do CORPO comparado
(19) FASE III - CLASSIFICADOR - GERAR ARQUIVO .arff - 4GRAMS - 4-grams do CORPO comparado

(30) FASE II - AGRUPAMENTO DE CONVERSACOES (AGRUPAMENTO E CLASSIFICACAO).
(41) FASE III - CLASSIFICACAO DAS CONVERSACOES (CRIME E NAO CRIME).

(37) FASE IV - GERAR GRAFOS - TODAS AS CONVERSACOES.
(38) FASE IV - GERAR GRAFOS - DIVIDIR GRAFO.
(39) FASE IV - GERAR GRAFOS - DIVIDIR GRAFO ENTRE CRIMINOSOS E NAO CRIMINOSOS.
(51) FASE IV - LISTAR EMAILS CRIMINOSOS.
(50) FASE IV - LISTAR EMAILS POR GRUPO.
(52) FASE IV - EXEMPLO DE LAUDO PERICIAL.
(53) FASE IV - LISTAR PALAVRAS CRIMINOSAS.
```

Figura A.1 – Tela Principal – *Reader*

A tela principal de *Reader* possui 28 opções organizadas de acordo com as FASES para funcionamento do mecanismo.

A.2. Processo de “Imagem” Forense e Pré-processamento

A opção 3 limpa o BD de antigas análises forenses. A opção 1 realiza o processo de “imagem” – cópia forense dos *e-mails* suspeitos que serão investigados e que deve ser indicado pelo utilizador de *Reader*. E as opções 5, 6, 7, 8 e 9 realizam o pré-processamento nos *e-mails* (corpo e campo assunto) copiados. Na Tabela A.1 todas as opções que tratam esta Seção são executadas.

Tabela A.1 – Processo de “Imagem” Forense e Pré-processamento

3	Limpendo o BD
	0 email(s) adicional(is) removido(s)
	0 associacoes Receiver removida(s)
	0 associacoes Sender removida(s)
	0 Pessoa(s) removida(s)
	0 anexo(s) removido(s)
	0 email(s) removido(s)
	0 resultado(s) removido(s)
1	DIGITE O CAMINHO DO REPOSITARIO DE E-MAILS:
	E:\Mestrado\02 Mestrado PPGIA - PUC - 2009\Trabalhos Linha de Pesquisa\Dissertacao Jackson\bases 7 - 33 grupos\base
	IMPORTANDO...
	E-MAILS COPIADOS!
5	
6	
7	
8	
9	
	PRÉ-PROCESSAMENTO REALIZADO NA BASE DE E-MAILS ARMAZENADA NO BD.

A.3. Gerar Arquivos de Atributos

As opções 4, 10, 13, 14, 31, 32, 34 e 40 geram os arquivos de atributos para ser utilizado no agrupamento de conversações. E as opções 18 e 19 geram arquivos de atributos para ser utilizado para classificação de conversações. Na Tabela A.2, o arquivo de atributos que utiliza as palavras do corpo de *e-mails* como atributos (opção 4), e representados pela técnica *bag of words*, e o arquivo de atributos para classificação dos *e-mails* da Base teste são gerados (opção 19).

Tabela A.2 – Geração dos Arquivos de Atributos (FASE II e FASE III)

<p>4</p> <p>Buscando e-mails do banco de dados As palavras do corpo dos e-mails são os atributos - bag of words Gerando o ARFF MUDAR O NOME DO ARQUIVO CONCLUÍDO NOVO ARQUIVO - C:/reader/OPCAO.arff</p> <p>19</p> <p>Gerando arquivo de atributos .ARFF - teste - 4-grams nas palavras do CORPO comparado com as palavras do DICIONARIO.TXT. Similaridade? 100% MUDAR O NOME DO ARQUIVO CONCLUÍDO NOVO ARQUIVO - C:/reader/baseteste.arff</p>
--

A.4. Agrupamento das Conversações

A opção 30 realiza o agrupamento de conversações no uso de um dos arquivos de atributos gerados durante a FASE II (opção 4, 10, 13, 14, 31, 32, 34 ou 40). Na Tabela A.3 visualiza-se a execução da opção 30 na utilização do arquivo de atributos gerado pela opção 4.

Tabela A.3 – Agrupamento de Conversações

<p>30</p> <p>QUAL O ARQUIVO (FASE II) DESEJA USAR? OPCAO4</p> <p>RELATORIO - Agrupamento de Conversações - TAXA DE ACERTOS SVM Polynomial: 2.456140350877193 % SVM Radial: 27,36 % SVM Sigmoid: 15,78 % Euclidian: 56.6667 % Manhattan: 62.807 % Cossine: 71,43% Jaccard: 84.2105 % Decision Tree J48: 27.36842105263158 % NaiveBayes: 1.9298245614035088 % Relatorio - TAXA DE ERROS Melhor resultado: 1.9298245614035088 - naive</p>
--

A.5. Classificação das Conversações

Tendo-se gerado o arquivo de atributo (opção 19) para a base de *e-mails* Teste e imaginando-se a existência de um arquivo de atributos gerado da mesma forma para a base Treinamento, realiza-se a classificação das conversações mediante a opção 41, como pode ser visto na Tabela A.4.

Tabela A.4 – Classificação das Conversações

<p>41</p> <p>Identificando conversações criminosas - TAXA DE ERROS QUAL O ARQUIVO DESEJA USAR? TESTE: 4-GRAMS SIMILARIDADE 100% TREINAMENTO: 4-GRAMS SIMILARIDADE 100% RELATORIO - TAXA DE ACERTOS</p> <p>Decision Tree results : Correct 99.82456140350877 Naive results : Correct 99.82456140350877 SVM Radial results : Correct 65.43859649122807 SVM Polynomial results : Correct 70.00000000000000 SVM Linear results : Correct 72.28070175438596 SVM Sigmoid results : Correct 95.61403508771930</p> <p>Melhor Classificador : naive</p>
--

A.6. Resultados

Agrupada e classificadas as conversações, podem ser consultados resultados da aplicação do mecanismo mediante as opções 37, 38, 39, 50, 51, 52 e 53. Na Tabela A.5 apresenta-se o laudo pericial exemplo (opção 52).

Tabela A.5 – Laudo Pericial Exemplo

<p>52</p> <p>EXEMPLO DE LAUDO PERICIAL. LAUDO PERICIAL - EXEMPLO</p> <p>Data de geração do relório: 03/02/2011 04:00 Data de Inicio de execução: 03/02/2011 05:30</p> <p>EXMO SR. DR. JUIZ DE DIREITO DA</p> <p>AÇÃO: Análise da Base de E-mails AUTOR: RÉU: PROCESSO:</p> <p>Foram analisados 570 emails.</p> <p>Pasta onde esta gravado o resultado dos grafos criminosos c:/Reader/out</p> <p>Agrupamento de conversações: naive Melhor Classificador Identificador de conversações: naive</p> <p>Emails LISTADOS por Conversação: E-MAILS DO GRUPO 0 E-MAILS DO GRUPO 1 E-MAILS DO GRUPO 2 E-MAILS DO GRUPO 3 E-MAILS DO GRUPO 4</p> <p>E-mails Criminosos: E-mail enviado por User08@dominio.com para User29@dominio.com no dia 2010-12-02 11:31:21.0 cujo assunto é Conversação CRIMINOSA 29 é criminoso e pertence ao grupo29 . E-mail enviado por User29@dominio.com para User08@dominio.com no dia 2010-12-02 13:36:08.0 cujo assunto é RE: Conversação CRIMINOSA 29 é criminoso e pertence ao grupo29 .</p>
--

E-mail enviado por User29@dominio.com para User36@dominio.com no dia 2010-12-02 13:36:08.0 cujo assunto é RE: Conversação CRIMINOSA 29 é criminoso e pertence ao grupo29 .

E-mail enviado por User29@dominio.com para User35@dominio.com no dia 2010-12-02 13:36:08.0 cujo assunto é RE: Conversação CRIMINOSA 29 é criminoso e pertence ao grupo29 .

E-mail enviado por User29@dominio.com para User37@dominio.com no dia 2010-12-02 13:36:08.0 cujo assunto é RE: Conversação CRIMINOSA 29 é criminoso e pertence ao grupo29 .

E-mail enviado por User29@dominio.com para User08@dominio.com no dia 2010-12-02 14:26:03.0 cujo assunto é RE: Conversação CRIMINOSA 29 é criminoso e pertence ao grupo29 .

E-mail enviado por User29@dominio.com para User36@dominio.com no dia 2010-12-02 14:26:03.0 cujo assunto é RE: Conversação CRIMINOSA 29 é criminoso e pertence ao grupo29 .

E-mail enviado por User29@dominio.com para User35@dominio.com no dia 2010-12-02 14:26:03.0 cujo assunto é RE: Conversação CRIMINOSA 29 é criminoso e pertence ao grupo29 .

E-mail enviado por User29@dominio.com para User37@dominio.com no dia 2010-12-02 14:26:03.0 cujo assunto é RE: Conversação CRIMINOSA 29 é criminoso e pertence ao grupo29 .

E-mail enviado por User08@dominio.com para User29@dominio.com no dia 2010-12-02 15:23:11.0 cujo assunto é Re: Conversação CRIMINOSA 29 é criminoso e pertence ao grupo29 .

E-mail enviado por User29@dominio.com para User08@dominio.com no dia 2010-12-02 17:06:57.0 cujo assunto é RE: Res: Conversaçã9 CRIMINOSA 30 é criminoso e pertence ao grupo29 .

E-mail enviado por User08@dominio.com para User29@dominio.com no dia 2010-12-02 17:59:30.0 cujo assunto é Re: Conversação CRIMINOSA 29 é criminoso e pertence ao grupo29 .

E-mail enviado por User29@dominio.com para User08@dominio.com no dia 2010-12-03 01:08:24.0 cujo assunto é RE: Conversação CRIMINOSA 29 é criminoso e pertence ao grupo29 .

E-mail enviado por User08@dominio.com para User29@dominio.com no dia 2010-12-03 02:01:26.0 cujo assunto é Re: Conversação CRIMINOSA 29 é criminoso e pertence ao grupo29 .

E-mail enviado por User29@dominio.com para User08@dominio.com no dia 2010-12-03 20:33:19.0 cujo assunto é RE: Conversação CRIMINOSA 29 é criminoso e pertence ao grupo29 .

E-mail enviado por User08@dominio.com para User29@dominio.com no dia 2010-12-03 20:38:24.0 cujo assunto é Re: Conversação CRIMINOSA 29 é criminoso e pertence ao grupo29 .

E-mail enviado por User29@dominio.com para User08@dominio.com no dia 2010-12-03 23:51:30.0 cujo assunto é RE: Conversação CRIMINOSA 29 é criminoso e pertence ao grupo29 .

E-mail enviado por User29@dominio.com para User08@dominio.com no dia 2010-12-04 00:02:20.0 cujo assunto é RE: Conversação CRIMINOSA 29 é criminoso e pertence ao grupo29 .

E-mail enviado por User29@dominio.com para User08@dominio.com no dia 2010-12-04 00:08:04.0 cujo assunto é RE: Conversação CRIMINOSA 29 é criminoso e pertence ao grupo29 .

E-mail enviado por User08@dominio.com para User29@dominio.com no dia 2010-12-04 00:26:38.0 cujo assunto é Re: Conversação CRIMINOSA 29 é criminoso e pertence ao grupo29 .

E-mail enviado por User29@dominio.com para User08@dominio.com no dia 2010-12-04 00:58:17.0 cujo assunto é RE: Conversação CRIMINOSA 29 é criminoso e pertence ao grupo29 .

E-mail enviado por User08@dominio.com para User29@dominio.com no dia 2010-12-04 01:01:16.0 cujo assunto é Re: Conversação CRIMINOSA 29 é criminoso e pertence ao grupo29 .

E-mail enviado por User08@dominio.com para User37@dominio.com no dia 2010-12-04 01:01:16.0 cujo assunto é Re: Conversação CRIMINOSA 29 é criminoso e pertence ao grupo29 .

E-mail enviado por User08@dominio.com para User43@dominio.com no dia 2010-12-04 01:01:16.0 cujo assunto é Re: Conversação CRIMINOSA 29 é criminoso e pertence ao grupo29 .

E-mail enviado por User08@dominio.com para User29@dominio.com no dia 2010-12-04 01:03:49.0 cujo assunto é Re: Conversação CRIMINOSA 29 é criminoso e pertence ao grupo29 .

E-mail enviado por User08@dominio.com para User45@dominio.com no dia 2010-12-04 01:03:49.0 cujo assunto é Re: Conversação CRIMINOSA 29 é criminoso e pertence ao grupo29 .

E-mail enviado por User08@dominio.com para User29@dominio.com no dia 2010-12-04 01:21:56.0 cujo assunto é Re: Conversação CRIMINOSA 29 é criminoso e pertence ao grupo29 .

E-mail enviado por User08@dominio.com para User29@dominio.com no dia 2010-12-04 01:34:11.0 cujo assunto é Re: Conversação CRIMINOSA 29 é criminoso e pertence ao grupo29 .

Na execução da opção 53, palavras criminosas detectadas no conjunto de *e-mails* criminoso são listadas, conforme fragmento exposto na Tabela A.6.

Tabela A.6 – Fragmento com as Palavras Criminosas

<p>53</p> <p>Email cujo assunto é Re: Conversação CRIMINOSA 29 possui as seguintes palavras criminosas: abobado possui 1 ocorrências. vadio possui 1 ocorrências. vagina possui 1 ocorrências. entendido possui 1 ocorrências. vigarista possui 1 ocorrências. rato possui 1 ocorrências. porco possui 1 ocorrências. doente possui 1 ocorrências. abelhudo possui 1 ocorrências.</p> <p>Email cujo assunto é Re: Conversação CRIMINOSA 29 possui as seguintes palavras criminosas: abobado possui 1 ocorrências. vadio possui 1 ocorrências. safado possui 1 ocorrências. fifi possui 1 ocorrências. vagina possui 1 ocorrências. entendido possui 1 ocorrências. corno possui 1 ocorrências. vigarista possui 1 ocorrências. rato possui 1 ocorrências. porco possui 1 ocorrências. doente possui 1 ocorrências. abelhudo possui 1 ocorrências.</p> <p>Email cujo assunto é Re: Conversação CRIMINOSA 29 possui as seguintes palavras criminosas: abobado possui 1 ocorrências. vadio possui 1 ocorrências. vagina possui 1 ocorrências. entendido possui 1 ocorrências. vigarista possui 1 ocorrências. abelhudo possui 1 ocorrências.</p> <p>Email cujo assunto é Re: Conversação CRIMINOSA 29 possui as seguintes palavras criminosas: abobado possui 1 ocorrências. vadio possui 1 ocorrências. vagina possui 1 ocorrências. entendido possui 1 ocorrências. vigarista possui 1 ocorrências. abelhudo possui 1 ocorrências.</p> <p>Email cujo assunto é Re: Conversação CRIMINOSA 29 possui as seguintes palavras criminosas: abobado possui 1 ocorrências. vadio possui 1 ocorrências. vagina possui 1 ocorrências. entendido possui 1 ocorrências. abelhudo possui 1 ocorrências.</p> <p>Email cujo assunto é Re: Conversação CRIMINOSA 29 possui as seguintes palavras criminosas: abobado possui 1 ocorrências. vadio possui 1 ocorrências. vagina possui 1 ocorrências. entendido possui 1 ocorrências. vigarista possui 1 ocorrências. doente possui 1 ocorrências. abelhudo possui 1 ocorrências. vadio possui 1 ocorrências.</p>
--

Na execução da opção 39, listam-se os grafos criminosos e não criminosos. Na Figura A.2 visualiza-se o grafo criminoso para a análise feita na base de *e-mails* Teste.

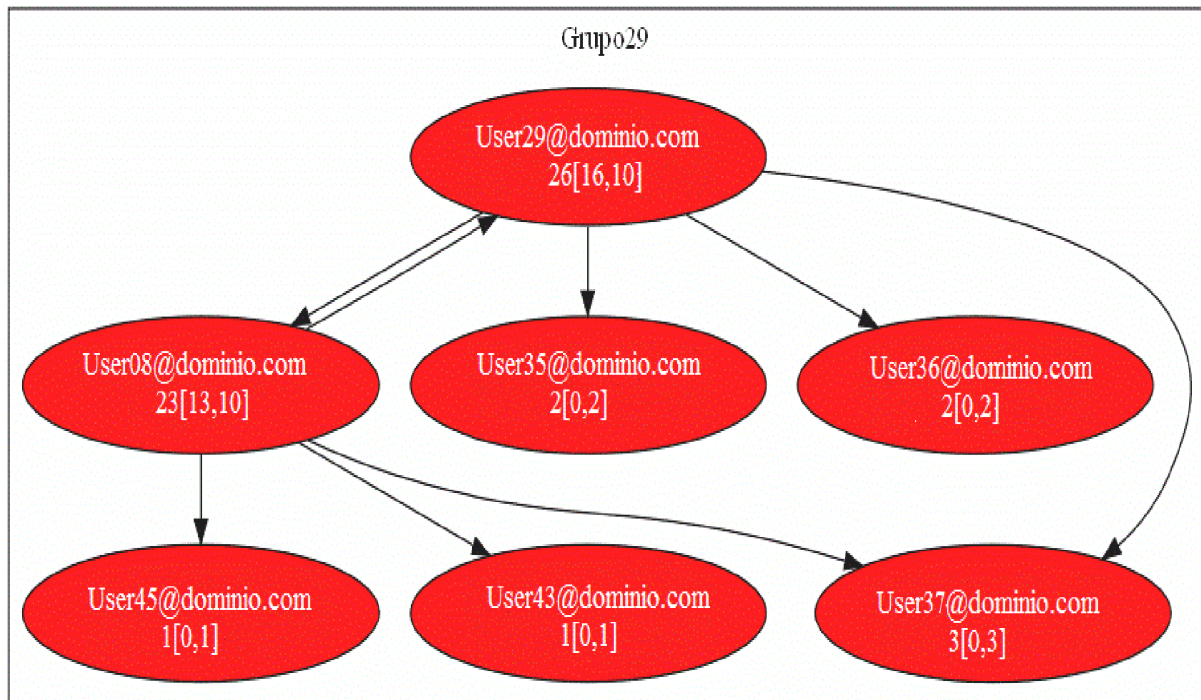


Figura A.2 – Grafo Criminoso