

UNIVERSIDADE DE CAXIAS DO SUL

FILIPE CHAVES DE MACEDO

**ANÁLISE DE SENTIMENTOS EM REDES SOCIAIS : UMA ABORDAGEM
BASEADA EM EMOÇÕES**

Caxias do Sul
2017

Filipe Chaves de Macedo

**ANÁLISE DE SENTIMENTOS EM REDES SOCIAIS : UMA ABORDAGEM
BASEADA EM EMOÇÕES**

Trabalho de Conclusão de Curso apresentado
como parte dos requisitos para a obtenção do
título de Engenheiro de Controle e Automação
da Universidade de Caxias do Sul.

Orientador:
Prof.^a Dra. Carine Geltrudes Webber

Caxias do Sul
2017

Filipe Chaves de Macedo

**ANÁLISE DE SENTIMENTOS EM REDES SOCIAIS : UMA ABORDAGEM
BASEADA EM EMOÇÕES**

Trabalho de Conclusão de Curso apresentado
como parte dos requisitos para a obtenção do
título de Engenheiro de Controle e Automação
da Universidade de Caxias do Sul.

Orientador:
Prof.^a Dra. Carine Geltrudes Webber

Aprovado em ____/____/____

Banca Examinadora

Prof. Dra. Carine Geltrudes Webber(orientadora)
Universidade de Caxias do Sul - UCS

Prof. Dra. Maria de Fatima Webber do Prado Lima
Universidade de Caxias do Sul - UCS

Prof. Dra. Helena Graziottin Ribeiro
Universidade de Caxias do Sul - UCS

Aos meus pais, por tornarem isto possível.

AGRADECIMENTOS

Gostaria de agradecer em especial aos meus pais e minha família, que sempre foram a base e espelho para que eu pudesse seguir e chegar até aqui.

À Mariana Santos, que por vezes me ajudou, aconselhou e dedicou do seu tempo para me ajudar em diversos aspectos relacionados a este trabalho.

À todos os colegas que por vezes me acompanharam nos estudos extraclasse, acabando por ficar sábados inteiros na biblioteca estudando para provas e trabalhos.

À minha orientadora prof^a. Dra. Carine Geltrudes Webber, por ter me dado todas as direções, por me incentivar e por sempre demonstrar disposição com o esclarecimento de dúvidas e questionamentos.

A todos que direta ou indiretamente contribuíram para que este trabalho fosse possível.
Meu mais sincero obrigado!

“Aprender é a única coisa de que a mente nunca se cansa, nunca tem medo e nunca se arrepende.”
Leonardo da Vinci

RESUMO

Nos dias atuais as redes sociais online testemunham um aumento significativo de usuários, com isso acabam por se tornar um alvo em potencial para entidades mal intencionadas que tem por objetivo espalhar dados maliciosos comprometendo a experiência do usuário e o colocando em risco. A presente pesquisa apresenta a importância do cuidado com a segurança nas redes sociais online e quais métodos são utilizados atualmente para evitar as fraudes contra usuários em tais meios sociais.

Com o intuito de entender opiniões de usuários sobre diferentes temas e agregar conhecimento à análise realizada, este trabalho aborda o conceito de emoções e suas características. Com o apoio de um léxico que torna possível a ligação de determinadas palavras com sentimentos expressos por seres humanos, analisou-se postagens que foram classificadas de acordo com o sentimento contido nelas.

A análise de dados textuais encontrados na rede social online foi abordada neste trabalho e para isso foram apresentados meios de coleta e tratamento de dados, classificadores de texto e métodos de avaliação de conhecimento. Posteriormente, com o objetivo de verificar a eficácia dos métodos apresentados foram feitas coletas de postagens na rede social *Facebook*, tentando avaliar o quão efetiva foi a classificação das postagens coletadas.

Palavras-chave: Análise de sentimentos. Redes Sociais. Inteligência Artificial. Mineração de Texto.

ABSTRACT

Nowadays, online social networks have witnessed a significant increase in users, which in turn has become a potential target for malicious entities that aim to spread malicious data by compromising the user experience and putting it at risk. The present research presents the importance of taking care of security in online social networks and what methods are currently used to avoid fraud against users in such social media.

In order to understand user's opinions on different subjects and to aggregate knowledge an analysis carried out deals with the concept of emotions and their characteristics. With the support of a lexicon that makes possible a connection of certain words with feelings expressed by humans as posts is also classified according to the feeling that is contained in them.

The analysis of textual data found in the online social network will be approached and for this will be presented ways of collecting and processing data, text classifiers with unsupervised learning and methods of knowledge evaluation. Subsequently, in order to verify the effectiveness of the presented methods will be made collections of posts in the social network Facebook, trying to evaluate how effective was the classification of the collected posts.

Keywords: Sentiment Analysis. Social network analysis. Artificial intelligence. Text Mining.

LISTA DE FIGURAS

Figura 1 – Falso e-mail de intimação	15
Figura 2 – Links Maliciosos	16
Figura 3 – Processo de Extração de Conhecimento Não Supervisionada	18
Figura 4 – Método de Luhn	20
Figura 5 – Processo de extração de conhecimento nas redes sociais	26
Figura 6 – Tipos de comportamento individual	28
Figura 7 – Comparação da curva ROC(<i>Receiver operating characteristic</i>) para diferentes modelos com dados fora do conjunto de teste	32
Figura 8 – Menu de Buscas	37
Figura 9 – Avisos de alerta Ok	37
Figura 10 – Avisos de Alerta Problema	38
Figura 11 – Nuvem de Palavras	38
Figura 12 – Botão Análise Gráfica	39
Figura 13 – Análise Gráfica de Sentimentos	39
Figura 14 – Botão Exportar Postagens	40
Figura 15 – Exemplo de Postagem Exportada	40
Figura 16 – Fluxograma de funcionamento do software	41
Figura 17 – Mensagem endereço inválido	42
Figura 18 – Exemplos de Nuvens Disponíveis	45
Figura 19 – Postagem classificada como neutra	48

LISTA DE TABELAS

Tabela 1 – Redes Sociais e seu objetivo	14
Tabela 2 – Tabela Documento-Termo	20
Tabela 3 – Tabela de Contingência	23
Tabela 4 – Resultado Alimentos – K-Means	30
Tabela 5 – Resultados Obtidos por (CARVALHO FILHO, 2014)	31
Tabela 6 – Performance dos modelos para dados fora do conjunto de treinamento	32
Tabela 7 – Categorias de emoções	36
Tabela 8 – Tipos de URL páginas x grupos	41
Tabela 9 – Remoção de caracteres especiais	43
Tabela 10 – Remoção de <i>stopwords</i> e caracteres especiais	43
Tabela 11 – Classificação Inicial do Grupo Mackbook	48
Tabela 12 – Classificação Final do Grupo Mackbook	49
Tabela 13 – Classificação do Grupo Previdência Social	49
Tabela 14 – Classificação da Página Outback Brasil	50
Tabela 15 – Taxa média de acerto na classificação de postagens dos três grupos analisados	50
Tabela 16 – Comparativo de trabalhos relacionados	55

SUMÁRIO

1	INTRODUÇÃO	12
1.1	OBJETIVOS	13
1.1.1	Objetivo Geral	13
1.1.2	Objetivos Específicos	13
1.2	LIMITES DO TRABALHO	14
2	EXTRAÇÃO DE CONHECIMENTO NAS REDES SOCIAIS	15
2.1	MEIOS DE FRAUDE NAS REDES SOCIAIS ONLINE	15
2.1.1	Spam	15
2.1.2	Links Falsos	16
2.1.3	Falsificação de Identidade	16
2.1.4	Conteúdo Impróprio	16
2.2	EXTRAÇÃO AUTOMÁTICA DE CONHECIMENTO	17
2.3	EXTRAÇÃO EM TEXTOS	17
2.4	ETAPAS DA EXTRAÇÃO	18
2.4.1	Pré-Processamento dos Documentos	18
2.4.2	Extração de Padrões usando Agrupamento de Textos	20
2.4.2.1	Medida de proximidade	21
2.4.2.2	Estratégia de agrupamento	22
2.4.2.3	Seleção de descritores para agrupamento	23
2.4.3	Avaliação de conhecimento	24
2.4.3.1	Silhueta	24
2.4.3.2	Entropia	25
2.5	EXTRAÇÃO DE CONHECIMENTO NAS REDES SOCIAIS	26
2.5.1	Coleta de Dados	26
2.5.2	Mineração de Texto	27
2.6	ANÁLISE DE COMPORTAMENTO INDIVIDUAL DE USUÁRIOS	27
2.6.1	Método Para Análise de Comportamento	28
2.7	TRABALHOS RELACIONADOS	29
2.7.1	Análise de Dados Para Apoio a Tomada de Decisão Com Ênfase no Marketing Digital	29
2.7.2	Mineração de Textos: Análise de Sentimento Utilizando Tweets Referentes à Copa do Mundo 2014	30
2.7.3	Detectando Clusters de Contas Falsas em Redes Sociais Online	31
2.7.4	Detectando Conteúdo Malicioso no Facebook	32
2.8	CONSIDERAÇÕES FINAIS	34
3	ANÁLISE DE SENTIMENTOS NO FACEBOOK: IMPLEMENTAÇÃO E RESULTADOS	35
3.1	TEORIAS DA EMOÇÃO	35
3.1.1	Léxico de Emoções	35
3.1.2	Palavras relacionadas a segurança e linguagem inapropriada	36
3.2	APRESENTAÇÃO E IMPLEMENTAÇÃO DO SOFTWARE DESENVOLVIDO	36
3.2.1	Apresentação do Software	37
3.2.1.1	Menu de Buscas	37
3.2.1.2	Avisos de Alerta	37

3.2.1.3	Palavras mais frequentes	38
3.2.1.4	Análise gráfica	38
3.2.1.5	Exportar Postagens Analisadas	39
3.2.2	Implementação do Software	40
3.2.2.1	Identificando o grupo ou página de coleta a partir de uma URL	41
3.2.2.2	Pré-Processamento dos Dados Adquiridos	42
3.2.2.3	Classificação das postagens	43
3.2.2.4	Implementação do Gráfico	44
3.2.2.5	Implementação da Nuvem de Palavras	44
3.2.2.6	Implementação da Opção Exportar PDF	45
3.3	DESCRIÇÃO DO EXPERIMENTO	45
3.3.1	Escolha da Rede Social	46
3.3.2	Definição dos Grupos de Coleta	47
3.3.3	Definição de Dados para Coleta	47
3.3.4	Método de Análise dos Dados Coletados	47
3.4	AVALIAÇÃO DE RESULTADOS	48
4	CONSIDERAÇÕES FINAIS	51
4.1	SÍNTESE DO TRABALHO	51
4.2	CONTRIBUIÇÕES DO TRABALHO	51
4.3	TRABALHOS FUTUROS	52
	REFERÊNCIAS	53
	APÊNDICE A – TABELA COMPARATIVA DE TRABALHOS RELACIONADOS	55

1 INTRODUÇÃO

Uma rede social pode ser entendida como um organismo pluricelular de tal forma que seus usuários são as células (FRANÇA et al., 2014). Como ocorre em todos os organismos, existem muitos tipos de célula, algumas perduram por pouco tempo, outras por muito tempo, existem também as boas e as más que podem beneficiar ou prejudicar o organismo e as demais células. De fato, a mesma dinâmica que mantém ou prejudica a vida se repete nas redes sociais em geral. Estudar e avaliar os efeitos desta dinâmica a fim de desenvolver tecnologias de monitoramento e diagnóstico autônomas que garantam a privacidade e segurança dos usuários é um dos papéis da Inteligência Artificial. Além disso, o monitoramento das redes sociais é de grande importância para as organizações em geral. Por exemplo, a análise do relacionamento entre funcionários de uma empresa permite diagnosticar problemas no ambiente de trabalho, já a identificação de notificações na área da saúde possibilita mapear uma propagação de epidemia.

Pessoas precisam pertencer a grupos sociais, é uma necessidade humana. Some-se a este fato a facilidade de se obter acesso à internet nos dias atuais. Como resultado, segundo dados de 2012, cerca de 96% dos usuários da internet utilizam em algum momento uma rede social online (MARQUES; PINTO; ALVAREZ, 2016). Infelizmente, grande parte destes usuários não compreende os riscos de segurança aos quais estão expostos no ambiente social da internet, e estes riscos podem ser relacionados a violação de privacidade do usuário, roubo de informações, falsificação de identidade ou até mesmos assédio sexual. Estes são apenas alguns exemplos, no ano de 2015 cerca de 45% das fraudes em cartão de crédito nos Estados Unidos foram relatadas como roubo de informações pelo meio virtual (HOLMES, 2015) e como nas redes sociais não há filtro de publicações muitas pessoas podem ser enganadas por pessoas que postam links de sites falsos para roubo destas informações. Nas redes sociais os usuários compartilham informações profissionais, pessoais e também comportamentais, mesmo que sem a intenção de fazê-los, além de criarem vínculos de amizade e confiança. A partir disso torna-se possível a análise do seu perfil e o levantamento de traços específicos de sua personalidade, uma vez que as informações inseridas expressam suas opiniões e sentimentos.

Além destes perigos ainda há o risco de que pessoas possam agredir os usuários verbalmente com comentários tóxicos que podem ser desagradáveis quando lidos, este tipo de situação se torna ainda pior quando se trata de jovens que utilizam o meio social online. Os jovens por sua vez, são de modo geral mais ingênuos e suscetíveis aos riscos da internet, por isso, mães e pais ficam preocupados com quem seus filhos estão se relacionando nas redes sociais, uma vez que não podem estar o tempo todo acompanhando suas ações, devido às rotinas de vida corridas. Ainda, pais e mães não querem tirar seus filhos das redes sociais, visto que hoje em dia este acesso ao meio social eletrônico é muito importante para os jovens e não pertencer a uma rede social pode excluir uma criança de relações sociais. Se por um lado não deseja-se excluir as crianças e adolescentes das redes sociais, não se pode ignorar o fato de que existem perigos adjacentes ao seu uso sem monitoramento. Sabendo que nem sempre é possível garantir que o

acesso às redes sociais é feito sob supervisão de um adulto, é necessário criar-se mecanismos que apoiem a navegação segura. Um estudo realizado por (OLIVEIRA; FERREIRA, 2013) avaliou alunos de uma escola pública (9 a 14 anos) usuários do *Facebook* a fim de observar quais os seus conhecimentos a respeito de segurança na rede social. Como resultado descobriram que 80% dos alunos avaliados não tinha conhecimento sobre as opções de segurança (filtros de visualização de fotos, postagens, compartilhamentos e dados pessoais) da rede social. Foi averiguado também que nem os pais passam orientação mínima sobre questões de segurança e privacidade na rede.

Percebe-se que de fato há uma lacuna entre o usuário e a sua segurança. Além da necessidade de orientação para que os jovens se tornem cidadãos virtuais seguros, pode ser necessário o desenvolvimento de software que auxilie neste processo, e é neste escopo que se insere este trabalho, visando a proposta de uma ferramenta que colete e analise relações interpessoais em uma rede social, realizando o monitoramento para identificar e sinalizar usuários potencialmente perigosos.

1.1 OBJETIVOS

1.1.1 Objetivo Geral

O objetivo principal deste trabalho é avaliar as postagens em grupos nas redes sociais analisando o sentimento expresso nas mesmas. A partir daí encontrar uma maneira de classificá-las. Isto permitirá que o usuário a partir de sua rede social possa ter mais clareza acerca dos grupos que vai acessar em sua rede social, para que estes grupos não venham a causar desconforto ou até perigo, este processo de verificação de postagens manualmente se tornaria um processo demorado e desgastante visto que hoje as redes sociais possuem um grande volume de perfis.

1.1.2 Objetivos Específicos

O presente estudo trata da coleta de dados nas redes sociais e como tais informações podem ser utilizadas para a análise de mecanismos de classificação. Um dos grandes desafios da coleta é criar uma definição de quais dados disponíveis nas redes foram utilizados para definir os perfis dos usuários, em paralelo como foram filtrados, uma vez que grande parte dos dados é preenchido de forma manual pelo usuário e cada um pode preencher da maneira escolhida, onde os dados que deveriam representar a mesma coisa podem ser escritos de maneira diferente, após filtrados devem ser classificados para que possam retornar a informação que desejamos, ou seja, a correta classificação do usuário. Após os dados terem sido adquiridos inicia-se a estruturação dos algoritmos. Para esta estruturação foram utilizadas bibliotecas disponíveis já implementadas uma vez que o principal objetivo não é desenvolver o método, mas sim analisar

quais deles melhor se encaixam no tipo de estudo que esta sendo feito e entender o funcionamento do algoritmo para que se possa avaliar o porquê de sua superioridade ou inferioridade em relação a outros métodos.

1.2 LIMITES DO TRABALHO

O estudo apresentado utiliza dados adquiridos da rede social Facebook, devido a características desta rede social que facilitarão a implementação. Sendo que uma delas é a ferramenta Facebook Graph API que é disponibilizada pela plataforma e de livre acesso nos permitindo coletar informações da rede social como postagens, atividades recentes do usuário, dados pessoais, entre outros para que possamos realizar a coleta de informações. O Facebook também se mostra mais atraente em relação a outras redes devido a suas características de proporcionar que você conheça pessoas pelo meio *online*, disponibilizando ao usuário a chance de criar grupos de interesses em comum. Na Tabela 1 pode-se verificar uma breve classificação das redes sociais mais utilizadas hoje em dia.

Tabela 1 – Redes Sociais e seu objetivo

Nome	Objetivo
Facebook	Relações Interpessoais
Twitter	Relações Interpessoais
Instagram	Postagem de fotos
Youtube	Postagem de vídeos
LinkedIn	Relações Profissionais

Fonte: Elaborada pelo autor.

Uma futura análise de resultados deve ser feita a partir da coleta de dados, para isto foram utilizados bancos de dados onde será possível comparar sintaticamente expressões e palavras que possam indicar comentários desagradáveis ou até mesmo atitudes suspeitas da parte de outros usuários. Os dados foram classificados a partir da extração e classificação de sentimentos e opiniões que envolvam particularmente entidades humanas, ou seja, envolvam ou sejam direcionadas aos próprios usuários.

Os métodos de classificação foram utilizados e escolhidos a partir de trabalhos prévios os quais se assemelham com o que está sendo desenvolvido. Também foi utilizado amplo material de pesquisa que auxiliará no desenvolvimento do trabalho.

O material de pesquisa foi escolhido pelo método de revisão sistemática, definindo fontes de pesquisa específicas, palavras chave e também data de publicação dos itens utilizados, uma vez que itens antigos podem estar defasados.

Para a organização dos textos no capítulo dois são apresentados as fundamentações e conceitos estudados, assim como trabalhos relacionados, já no capítulo três é apresentada a proposta de trabalho, materiais e métodos a serem utilizados.

2 EXTRAÇÃO DE CONHECIMENTO NAS REDES SOCIAIS

Neste capítulo é apresentada a fundamentação teórica em que o trabalho se baseia para que se possa entender como funciona e quais métodos são utilizados na extração de conhecimento de redes sociais online a partir de dados textuais.

2.1 MEIOS DE FRAUDE NAS REDES SOCIAIS ONLINE

O mau uso das redes sociais pode deixar o usuário suscetível a uma série de riscos. Nas redes sociais online pode-se encontrar um grande número de perfis de usuários que utilizam do meio para realizarem algum tipo de ato malicioso. Abaixo serão expostos alguns dos métodos que são utilizados e foram citados por (WÜEST, 2010) e como podem prejudicar os membros de diferentes meios sociais online. No site (CAIS, 2017)¹ podemos encontrar uma lista com onze mil e trezentas fraudes cadastradas desde o ano de 2008, onde grande parte delas são divulgadas e disseminadas em redes sociais online.

2.1.1 Spam

O método Spam geralmente é utilizado para enviar e-mails para o usuário, sendo que estes geralmente contém arquivos maliciosos. A consequência gerada pode ser, por exemplo, quando realizado o download para o computador da vítima pode se instalar um vírus que vai acabar roubando senhas, tendo acesso a webcam ou pode até mesmo danificar a máquina do usuário modificando configurações de processamento. Na Figura 1 podemos visualizar o exemplo de um falso e-mail relatado pelo site (CAIS, 2017) que possui um arquivo anexado se fazendo parecer uma intimação do Ministério Público em formato pdf. O arquivo quando clicado realiza o download da suposta intimação e quando esta é aberta acaba transferindo o vírus para o computador do usuário.

Figura 1 – Falso e-mail de intimação



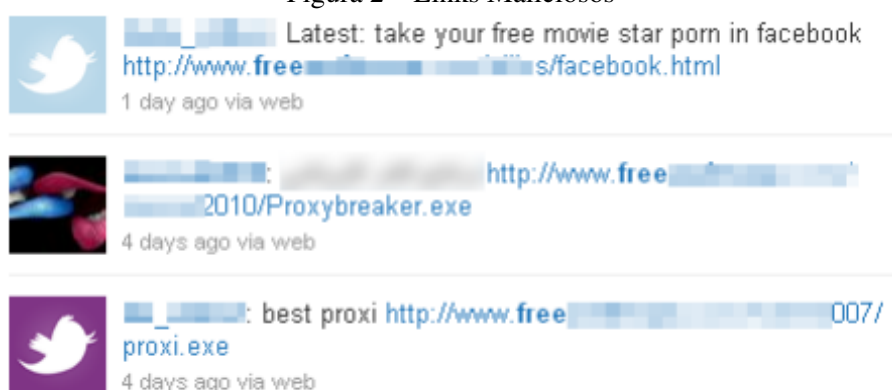
Fonte: (CAIS, 2017)

¹<https://memoria.rnp.br/cais/fraudes.php?>

2.1.2 Links Falsos

Outra maneira de espalhar as fraudes pelas redes sociais são os links contendo arquivos maliciosos. De acordo com (WÜEST, 2010) este tipo de link malicioso geralmente é postado nas redes junto com comentários, vinculados geralmente com assuntos que estão em alta na rede social. Na Figura 2 pode-se visualizar alguns links de exemplo, onde foram utilizados assuntos diversos para auxiliar na divulgação dos links maliciosos.

Figura 2 – Links Maliciosos



Fonte: (WÜEST, 2010)

2.1.3 Falsificação de Identidade

Um tipo de fraude um tanto diferente dos citados anteriormente é a criação de perfis falsos. O fraudadores criam perfis de famosos ou de amigos encontrados na rede social da vítima, pois de acordo com (WÜEST, 2010) relações com pessoas conhecidas facilitam o acesso a informações as quais não seria possível ter acesso. Da mesma maneira funciona com perfis de famosos, onde a informação lá postada ganha credibilidade, por serem de certa forma conhecidos dos usuários.

2.1.4 Conteúdo Impróprio

As redes sociais online não conseguem filtrar cem por cento do conteúdo que está inserido nelas, então muitas vezes o usuário corre o risco de acabar encontrando conteúdos desagradáveis durante a navegação. Em outros casos pode-se encontrar temas que seriam de classificação indicativa adulta e estão sem nenhum tipo de filtro na rede, sendo mostrados para usuários de todas as idades.

2.2 EXTRAÇÃO AUTOMÁTICA DE CONHECIMENTO

A extração automática deve começar pelos dados, mas antes que se entenda o método e quais serão os meios utilizados para adquiri-los, filtrá-los e também classificá-los é preciso que se defina, o que é um dado, ou seja, a matéria prima. De acordo com (AMARAL, 2016) dados são fatos coletados e geralmente armazenados. A informação é o dado analisado e com algum significado, o conhecimento é a informação interpretada, entendida e aplicada para um fim.

Ao exemplo de um carro que possui diversos sensores para que possa informar ao usuário ou ao próprio sistema informações importantes de segurança, como por exemplo, o sensor de velocidade que indica ao motorista o quão rápido o carro está se deslocando pela via. Este tipo sensor emite um sinal de frequência proporcional a velocidade do carro, se o veículo se desloca a baixa velocidade o sensor emite um sinal de baixa frequência, que aumenta proporcionalmente a velocidade.

Os dados podem estar basicamente em dois formatos diferentes, eletrônico ou não. No formato eletrônico podem se dividir entre digital e analógico, já no formato não eletrônico são encontrados em diferentes formas, impresso, sonoro, em pintura ou em inúmeras outras maneiras. Para que estes dados façam sentido precisamos que se encontrem em um formato ou estruturação de dados, por isso após o armazenamento os dados precisam passar por uma etapa de transformação para que se possa realizar sua análise. A etapa de transformação também deve prever a maneira como o dado será apresentado para o usuário.

A extração automática de conhecimento surge devido ao grande número de informações textuais produzidas, as quais seriam quase impossível de ser analisadas em sua totalidade se esta fosse feita manualmente.

2.3 EXTRAÇÃO EM TEXTOS

Neste trabalho foi feita a extração de conhecimento a partir de dados eletrônicos no formato digital. A produção de dados de formato digital pode ser feita de inúmeras maneiras, com sensores, com o teclado do computador enquanto digitamos algum texto e até mesmo com câmeras ao gravar vídeos.

A mineração de texto, uma forma de análise de dados, tem um papel importante na aquisição de conhecimento, visto que hoje cerca de 80% dos dados produzidos se encontra em formato textual (MARCACINI; MOURA; REZENDE, 2011). As análises textuais podem ser diretamente relacionadas com dados não estruturados. Muitos modelos de análise utilizados hoje não são adequados para este tipo de dado, mas sim para dados estruturados os quais podem ser planilhas, tabelas entre outros.

2.4 ETAPAS DA EXTRAÇÃO

O processo de mineração de texto pode ser executado de diversas formas no momento da aquisição de conhecimento, os métodos com estrutura automática recebem grande atenção devido ao seu potencial de não precisar de nenhuma base acerca do que está sendo analisado ou do banco de dados utilizado. De acordo com (MARCACINI; MOURA; REZENDE, 2011) um processo de mineração de textos se divide em três etapas: pré-processamento dos documentos, extração de padrões com agrupamento de textos e avaliação, como pode-se observar na Figura 3.

Figura 3 – Processo de Extração de Conhecimento Não Supervisionada



Fonte: O Autor.

2.4.1 Pré-Processamento dos Documentos

Esta etapa tem como princípio selecionar os dados de texto e apresentá-los de uma forma organizada e padronizada para que se possa extrair o conhecimento do texto. Nesta etapa é definido como os dados são apresentados ao algoritmo, por isso se faz necessário definir os principais termos de interesse e também a documentação de texto, visando manter as características e evitando diminuir a assertividade do resultado.

Quando os dados textuais são adquiridos, eles podem se encontrar em diferentes formatos, visto que hoje possuímos diversas ferramentas de textos que podem criar este tipo de dado. Comumente os textos são convertidos para uma forma limpa e sem formatação.

De acordo com (MARCACINI; MOURA; REZENDE, 2011) pode-se dizer que o grande desafio do processamento de textos é a grande variedade de termos que a língua escrita possui, ou seja, palavras que quando unidas podem ter diferentes significados do que quando comparadas se estiverem sozinhas. Outra característica dos termos é que muitos são desnecessários para a análise e são considerados termos de baixo peso ou *stopwords*.

Com a seleção de termos procura-se obter um aglomerado de termos representativo no contexto da análise, primeiramente deve-se remover as *stopwords*, ou seja, termos de baixo

peso, irrelevantes se não estiverem em conjunto com alguma outra palavra ao exemplo de artigos, advérbios, pronomes e preposições. A junção de várias *stopwords* chama-se *stoplist*, uma vez que este tipo de palavras e termos não tem representação significativa para a análise podemos removê-los sem ter perdas de conteúdo. Após a eliminação das *stopwords* é realizado o *stemming*, ou seja, reduzir palavras flexionadas a sua base.

Existe também outra maneira de realizar este processo para filtragem dos termos é simplesmente selecioná-los pelo número de vezes que aparecem no texto, chama-se *term frequency*. No processo de *term frequency* os termos de maior frequência são de menor importância, pois aparecem em grande parte do texto não trazendo informações relevantes, já os termos que aparecem pouco geralmente também não acrescentam informações importantes.

A definição de valores utilizada para faixa de corte dos termos de baixa e alta frequência de aparição no texto são definidas por métodos estatísticos. A frequência destes termos é regida pelo Método de Luhn, um método que tem por base a Lei de Zipf, conhecida também como Princípio do Menor Esforço. Conforme a Equação 2.1, onde P_n é a frequência que o termo aparece no texto, n é o n -ésimo termo onde a é uma constante de valor muito próximo de 1.

$$P_n \sim \frac{1}{n^a} \quad (2.1)$$

Diante da Equação chegamos à curva de Zipf que mostra a frequência com que aproximadamente os termos aparecem em um texto, a partir daí pode-se associá-la ao método de Luhn que nos mostra o quão significativo estes termos são dentro de um contexto geral. Em (LUHN, 1958) é relatado que um autor repete certas palavras enquanto escreve um texto e que quanto mais vezes esta palavra for encontrada mais significativa ela se torna.

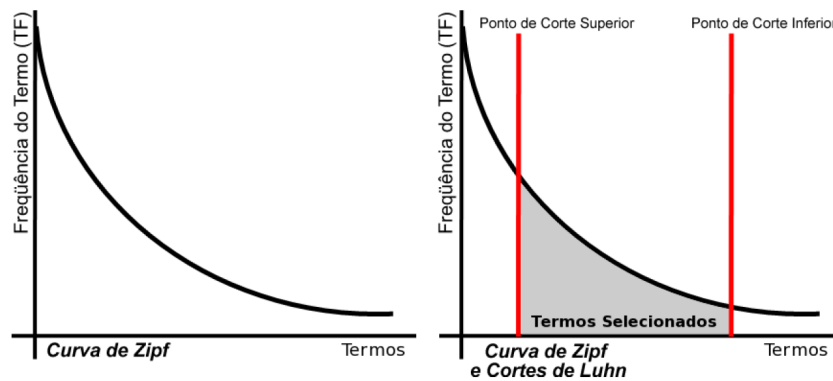
Embora algumas palavras apareçam muitas vezes, estas podem ter a função de ligação, logo são encontradas em grande parte do texto mas não tem significado importante e podem então ser desconsideradas. Deve-se criar um limite inferior para as palavras de menor frequência definindo assim a gama mais útil de palavras a serem consideradas.

Deve-se considerar que para que se estabeleça um local ótimo para as linhas de corte é necessário realizar testes com amostras e textos, para que se possa variar as palavras selecionadas e modificar também os resultados de saída. Esta forma de trabalho é a principal desvantagem deste método, pois a definição do ponto de corte acaba sendo por meio empírico e acaba tornando o ponto de corte algo um tanto subjetivo.

A Figura 4 exemplifica o Método de Luhn e também a relação entre termos e sua frequência de aparição no texto. Devido ao baixo processamento necessário deste método ele é geralmente o escolhido para coleções textuais de grandes tamanhos.

O próximo passo após a seleção de termos significativos é estruturar estes documentos, para que eles possam ser lidos pelo algoritmo que processará os resultados. Para isso o

Figura 4 – Método de Luhn



Fonte: (MARCACINI; MOURA; REZENDE, 2011).

método utilizado é o de espaço-vetorial, no qual cada documento é um vetor com um número determinado de valores, que são a quantidade de vezes em que um termo específico aparece na totalidade dos documentos. Pode-se chamar este método de *bag-of-words*, este por sua vez apresentado na Tabela 2, onde d_i é i -ésimo documento, t_j representa o j -ésimo termo, documentos são representados por vetores $d_i=(a_{i1},a_{i2},a_{i3},\dots,a_{ij})$. O valor a_{ij} geralmente pode ser encontrado pelo número de vezes que está presente em um documento ou o indicador de sua importância no documento.

Tabela 2 – Tabela Documento-Termo

	t_1	t_1	...	t_j
d_1	a_{11}	a_{12}	...	a_{1j}
d_2	a_{21}	a_{22}	...	a_{2j}
\vdots	\vdots	\vdots	...	\vdots
d_I	a_{I1}	a_{I2}	...	a_{Ij}

Fonte: (MARCACINI; MOURA; REZENDE, 2011)

2.4.2 Extração de Padrões usando Agrupamento de Textos

Tendo estruturado os dados na parte de processamento agora precisamos organizar os documentos conforme sua semelhança, ou seja, agrupá-los conforme a similaridade de seus termos, pois no aprendizado não supervisionado não se possui classes pré-determinadas de termos semelhantes para treinamento, o que acontece com algoritmos supervisionados.

Este tipo de agrupamento é baseado em três etapas principais: a medida de proximidade, estratégia de agrupamento e a seleção de descritores para o agrupamento. Estas etapas são descritas abaixo.

2.4.2.1 Medida de proximidade

As medidas de proximidade definem de qual maneira será verificada a similaridade entre textos ou documentos. Este método geralmente é definido de acordo com o tipo de dados que estamos utilizando, assim pode-se definir qual o objetivo destas medidas, calcular similaridade ou diferença entre os dados. Abaixo são apresentados dois métodos aplicados em dados do tipo texto: Cosseno e Jaccard. Para que seja entendido como estes métodos funcionam deve-se considerar dois vetores documento $x_i=(x_{i1},x_{i2},x_{i3},\dots,x_{im})$ e $x_j=(x_{j1},x_{j2},x_{j3},\dots,x_{jm})$. O método do cosseno pode definir a semelhança de acordo com o ângulo que os vetores formam entre si, onde o cosseno é 0 o ângulo entre os documentos é 90° , o que significa que não possuem nenhuma semelhança. Todavia se o cosseno for igual a 1 quer dizer que o documento possui ângulo igual a 0° , logo são muito semelhantes quando não iguais (MARCACINI; MOURA; REZENDE, 2011).

$$\text{cosseno}(x_i, x_j) = \frac{\sum_{l=1}^m x_{il}x_{jl}}{\sqrt{\sum_{l=1}^m x_{il}^2} \sqrt{\sum_{l=1}^m x_{jl}^2}} \quad (2.2)$$

É importante salientar que o método acima considera o peso dos termos dentro dos vetores, ou seja, quanto mais um termo x aparecer, maior será sua influência para o resultado de similaridade.

Em alguns tipos de vetores os dados estão em forma binária, ou seja, os vetores documento-termo consideram apenas se o termo aparece no respectivo documento ou não, sendo assim não importa quantas vezes o termo aparece ele terá o mesmo peso que os demais. Neste tipo de análise é utilizada a medida de Jaccard, como dito anteriormente x_i e x_j são dois vetores termo-documento e a partir disso podemos considerar:

termos_{11} : termos presentes nos dois vetores.

termos_{01} : termos presentes apenas no vetor x_j .

termos_{10} : termos presentes apenas no vetor x_i .

$$\text{jaccard}(x_i, x_j) = \frac{f_{11}}{f_{10} + f_{01} + f_{11}} \quad (2.3)$$

A medida Jaccard fica entre o intervalo $[0,1]$ pode-se dizer que quanto mais próximo a 1 maior a similaridade dos vetores e quanto mais próximo a 0 maior a diferença entre ambos (MARCACINI; MOURA; REZENDE, 2011).

Vistos os métodos acima não se pode definir qual é a maneira mais eficaz de se realizar uma medida de proximidade sem realizar a análise prévia dos vetores que se possui, uma vez que são métodos utilizados para diferentes tipos de estruturação de dados de entrada.

2.4.2.2 Estratégia de agrupamento

Após apresentados os métodos que podem ser utilizados para a medida de proximidade, devemos definir como será feito o agrupamento dos vetores com semelhança, para isso definimos o método de agrupamento. O método de agrupamento pode basicamente ser dividido entre Particional e Hierárquico. No Particional os vetores são apenas divididos em grupos simples, ou seja, apenas uma camada, e no método Hierárquico podemos encontrar diversos sub-níveis. Os algoritmos podem definir que os vetores pertençam a mais de um grupo, para que isso não aconteça deve-se utilizar estratégias de agrupamento sem sobreposição.

Resumidamente as técnicas de agrupamento procuram realizar o agrupamento de vetores da melhor forma possível (MARCACINI; MOURA; REZENDE, 2011). O *k-means* é um dos modos mais utilizados para que se realize o agrupamento de vetores quando se trata de agrupamento particional, onde o objetivo é dividir os vetores em k grupos. Abaixo são apresentadas algumas características acerca do *k-means*, explicando esta técnica conforme (STEINBACH et al., 2000).

O *k-means* necessita de um vetor médio, o qual se denomina o centroide, e é calculado a partir dos demais vetores que se possui no grupo. Abaixo temos a equação do centroide C para um grupo de vetores de tamanho G , onde x é um vetor documento-termo.

$$C = \frac{1}{|G|} \sum_{x \in G} x \quad (2.4)$$

Com o *K-means* se obtêm a distribuição média de vetores no grupo. Vale lembrar que o *k-means* só é válido a partir do momento em que se consegue calcular o centroide. Para o algoritmo *k-means* são dadas as entradas $X=(x_1, x_2, x_3, \dots, x_m)$ que é um conjunto de vetores documento-termo e k que é o número de grupos referente ao qual queremos que os vetores se dividam. A partir daí o algoritmo deve selecionar k vetores aleatórios como centroides iniciais e repetir a operação de calcular a dissimilaridade de cada vetor para cada centroide e inserir o documento no grupo a qual ele mais se assemelha, e então deve-se calcular a nova centroide deste grupo tendo em vista os vetores que foram adicionados a ele.

Conforme (STEINBACH et al., 2000) o critério de parada deste algoritmo pode ser escolhido de duas maneiras, pelo número de iterações ou quando não ocorrem mais modificações entre os vetores, ou seja, os vetores não trocam mais de grupos. Este método pode ser considerado o de diminuição de erro, onde ao minimizar o erro tem-se o melhor balanceamento entre grupos e conseqüentemente a melhor separação possível.

O algoritmo *k-means* pode possuir algumas desvantagens, ou restrições. Uma delas é necessitar que sejam informados inicialmente o número de grupos a serem utilizados e também pode ter respostas diferentes considerando que os vetores iniciais de centroides são aleatórios e podem influenciar no resultado.

2.4.2.3 Seleção de descritores para agrupamento

Com os grupos de documentos determinados precisa-se definir como apresentar estes dados, de forma que o algoritmo que irá ler os grupos de vetores saiba o que cada grupo significa e o usuário também possa saber se estes dados ficaram agrupados corretamente ou não.

Alguns métodos para a seleção de descritores utilizam o centroide, uma vez que este mantém a característica central de todos os grupos, logo seria o representante mais apropriado para o grupo. Outra estratégia é tomar como representante do grupo o termo que mais aparece em tal grupo, mas conforme (MARCACINI; MOURA; REZENDE, 2011) os resultados obtidos não são totalmente satisfatórios.

Em (MANNING et al., 2008) são expostas as técnicas de aprendizado de máquina que podem ser utilizadas para a seleção de descritores de agrupamento. Com estas técnicas é possível que se obtenha uma classificação dos termos que podem representar o grupo de melhor forma.

Com um grupo G de centroide C , os termos que estão em C se denominam t . Tem-se como objetivo criar uma listagem ordenada dos termos e selecionar os melhores como descritores de agrupamento daquele grupo. Diferentes critérios podem ser definidos para a construção do ranking e estes podem derivar da Tabela 3.

Tabela 3 – Tabela de Contingência

Q(t) / G	Relevante	Não Relevante
Relevante	acertos	ruído
Não Relevante	perda	rejeitos

Fonte: O Autor.

Em cada termo t é executada uma busca de documentos que o contenham, com o conjunto de documentos $Q(t)$ levantados e também o grupo de documentos G é preenchida a tabela acima, onde:

acertos: número de documentos listados por $Q(t)$ que pertencem ao grupo G ;

perda: número de documentos G que não foram listados por $Q(t)$;

ruído: número de documentos listados por $Q(t)$ que não pertencem ao grupo G ;

rejeitos: número de documentos que não estão em G que não foram listados por $Q(t)$;

Com estes números pode-se definir o critério que será utilizado para definir o quão significativo um termo é para representar determinado grupo. Uma maneira é utilizar o método F-Measure, mas para isso precisa-se que sejam calculadas previamente o valor da *Precisão(t)* e *Revocação(t)* conforme as fórmulas:

$$Precisao(t) = \frac{acertos}{acertos + ruido} \quad (2.5)$$

$$Revocacao(t) = \frac{acertos}{acertos + perda} \quad (2.6)$$

Tendo obtido os dois valores pode-se aplicar a equação F-Measure (2.7), esta função varia no intervalo de [0,1]. Quanto mais perto de 1 o valor da função F-Measure melhor o termo discriminará o grupo de vetores. A função é aplicada a todos os termos e a partir daí pode-se criar o ranking e então os n valores mais próximos a 1 serão os termos descritores do grupo G e vice-versa.

$$F - Measure = \frac{2 * Precisao(t) * Revocacao(t)}{Precisao(t) + Revocacao(t)} \quad (2.7)$$

2.4.3 Avaliação de conhecimento

Para a avaliação de conhecimento foi descrita a forma objetiva, ou seja, por meio de índices estáticos, sendo que estes indicam a qualidade obtida nos resultados.

A qualidade da organização dos documentos está diretamente relacionada com a qualidade do agrupamento na extração de padrões. Assim, a validação do conhecimento extraído é realizada por meio de índices utilizados na análise de agrupamentos.(STEINBACH et al., 2000)

De maneira geral pode-se dizer que existem três critérios principais a serem considerados: critérios internos, relativos e externos.

O critério interno avalia as posições dos objetos e se as mesmas correspondem à matriz de proximidade. O critério relativo compara todos os agrupamentos para verificar se os dados estão no grupo mais adequado, em quanto isso o critério externo contam com alguma informação que é fornecida por algum meio externo ao algoritmo, normalmente sobre a estrutura de agrupamento desejada. Pode ser realizada também a comparação dos dados agrupados automaticamente com um pequeno grupo de dados classificados manualmente.

Pode-se citar três índices de validação distintos, os quais podem nos indicar o quão bom é o agrupamento de acordo com diferentes maneiras de avaliação.

2.4.3.1 Silhueta

De acordo com (KAUFMAN; ROUSSEUW, 2009) a Silhueta pode ser considerada um critério relativo, que serve para que se possa avaliar agrupamentos. Considerando que $x_j=(x1,x2,x3,...,xj)$ é um conjunto de vetores documento termo e $s(x_j)$ o valor da silhueta, obtemos a Equação 2.8.

$$s(x_j) = \frac{b(x_j) - a(x_j)}{\max(a(x_j), b(x_j))} \quad (2.8)$$

Na equação deve-se considerar que $a(x_j)$ é a diferença média entre x_j e os demais documentos que estão no mesmo grupo e $b(x_j)$ a diferença média entre x_j e os documentos do agrupamento vizinho. O denominador da equação $\max(a(x_j), b(x_j))$ nada mais é do que um fator para normalização.

O valor final desta expressão se aplicada corretamente deve estar no intervalo de $[-1, 1]$. Os valores abaixo de zero mostram que o documento pode estar no grupo escolhido erroneamente e os valores positivos representam que o documento se encontra satisfatoriamente alocado no grupo que está.

Após realizada a análise dos documentos termo, deve-se avaliar o valor global de todas as silhuetas. Isso é feito para que se possa comparar diferentes agrupamentos e algoritmos.

$$S_p = \frac{1}{N} \sum_{j=1}^N s(x_j) \quad (2.9)$$

Uma vez que a equação global da silhueta S_p depende dos valores de silhueta de cada documento $s(x_j)$ e que, como dito anteriormente, quanto maior for este valor melhor alocado esta o documento, entende-se que S_p (de acordo com a Equação 2.9) deve também ser o mais alto possível. O que irá por consequência aumentar o valor de $a(x_j)$ e diminuir o valor de $b(x_j)$ na Equação 2.8.

2.4.3.2 Entropia

Outro método para avaliação de conhecimento pode ser a Entropia, que de acordo com (STEINBACH et al., 2000) a melhor entropia ocorre quando os dados estão perfeitamente alocados nos melhores grupos para estes. A Entropia pode ser considerada um critério externo, pois precisa de alguma base dos grupos de documentos. Em cada um dos grupos deve-se calcular a probabilidade de um documento do grupo j pertencer a uma classe i que é a informação externa fornecida ao algoritmo, conforme a Equação 2.10.

$$P_{ij} = \frac{|i \cap j|}{|j|} \quad (2.10)$$

Com isso aplica-se a fórmula padrão da entropia conforme a Equação 2.12, onde o somatório é feito a partir de todas as classes.

$$E_j = - \sum_i P_{ij} \log(P_{ij}) \quad (2.11)$$

A entropia total de um aglomerado de grupos é calculada a partir da soma de todas as entropias da coleção, ponderadas pelo tamanho de cada um dos grupos referentes, de acordo com a Equação 2.12 onde n_j é o tamanho do agrupamento j , m o número de grupos e n que

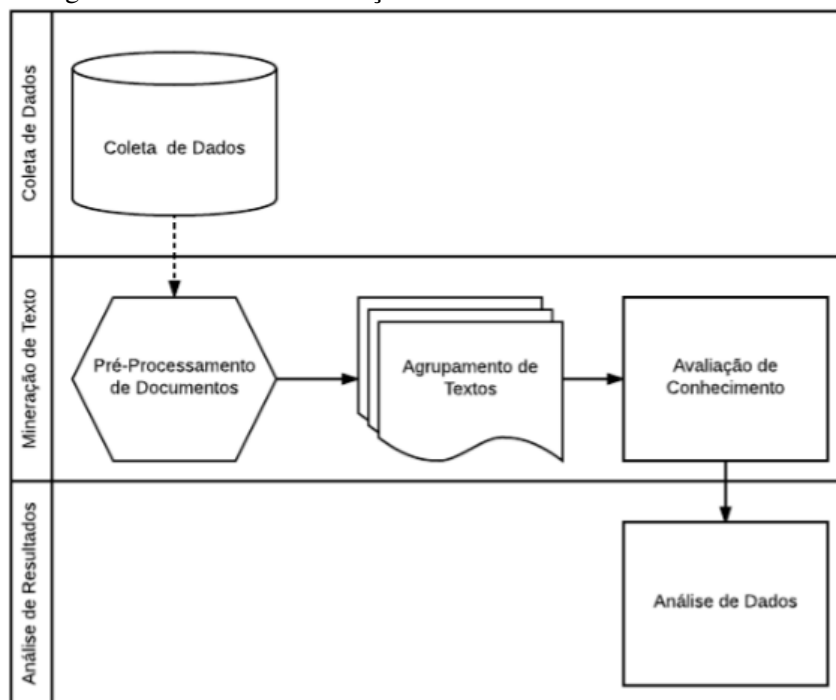
representa o número de documentos de toda a base textual.

$$Entropia = \sum_{j=1}^m \frac{n_j E_j}{n} \quad (2.12)$$

2.5 EXTRAÇÃO DE CONHECIMENTO NAS REDES SOCIAIS

Entendido como se executa a extração de conhecimento em texto deve-se avaliar como é feita a coleta de dados textuais em redes sociais de maneira geral, para que estes posteriormente possam ter conhecimento extraído das informações que apresentam. Para isso podemos utilizar o processo KDT (*Knowledge Discovery in Texts*), que segundo (BARION; LAGO, 2015) trata do processo de extração de informações que de alguma maneira sejam importantes em textos, que como já dito anteriormente são dados não estruturados. Este processo segue o fluxo da Figura 5.

Figura 5 – Processo de extração de conhecimento nas redes sociais



Fonte: O Autor.

2.5.1 Coleta de Dados

Na coleta de dados é onde acontece a aquisição dos dados textuais a serem utilizados. De acordo com (ARANHA; VELLASCO, 2007) a coleta de dados tem o objetivo de que se possa formar uma base de textos para serem utilizados no trabalho. Um grande desafio na coleta de dados é adequar o formato dos dados para que se possa utilizar no processo de mineração de

texto.

Tendo em vista o resultado a que se quer chegar deve-se definir quais textos são utilizados, pois isso torna a análise mais precisa e eleva o grau de confiabilidade do resultado (CECI, 2010). Esta coleta pode ser realizada por programas que se denominam *web crawlers*. Estes programas que criam uma conexão com determinado endereço eletrônico e extraem de lá a informação desejada. Estes tipos de programas utilizados no contexto de redes sociais geralmente se valem das APIs disponibilizadas pelas plataformas sociais.

2.5.2 Mineração de Texto

Em (TAN et al., 1999) o autor descreve a mineração de textos como sendo a descoberta de conhecimento a partir de textos. Diz também que este processo refere-se a extração de padrões de interesse e não-triviais.

De acordo com (GOMES; CASTRO NETO; HENRIQUES, 2013) a grande massa de dados textuais não estruturada em sua grande parte, não é passível de ser utilizada diretamente para extração de conhecimento por computadores. Uma vez que as ferramentas computacionais tratam esse tipo de dado somente como uma sequência de caracteres, então deve-se tratar este tipo de dados com a mineração de texto, que conforme a Figura 5 possui três etapas principais:

- **Pré-Processamento de Documentos:** Consiste na preparação dos dados, onde se utilizam os procedimentos de remoção de *stop words*, *stemming* e *term frequency*.
- **Agrupamento de Textos:** Tem por objetivo agrupar os documentos de texto por sua semelhança, utilizando medidas de proximidade e algoritmos de agrupamento.
- **Avaliação de Conhecimento:** Indica a qualidade do agrupamento realizado, para esta avaliação pode-se utilizar os métodos de silhueta ou entropia.

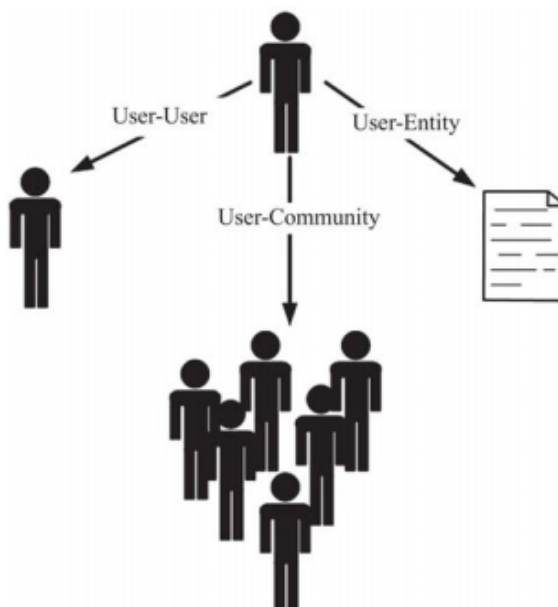
Estas etapas foram descritas na seção 2.4 deste trabalho.

2.6 ANÁLISE DE COMPORTAMENTO INDIVIDUAL DE USUÁRIOS

Para que se possa analisar o comportamento de usuários é necessário que se entenda o que os motiva a se juntar a determinada rede social, pois cada usuário tem um comportamento específico, apresenta diferentes emoções e também interesses. De acordo com (ZAFARANI; ABBASI; LIU, 2014) o comportamento individual pode ser classificado em três tipos distintos que são mostrados na Figura 6 .

Comportamento *Usuário-Usuário* está ligado a relações entre usuários singulares, ou seja, de um único indivíduo para outro seja com comentários curtindo fotos ou qualquer outro tipo de interação entre singulares. *Usuário-Comunidade* é quando o usuário fala para algum grupo ou comunidade expressando sua opinião para outros usuários que se encontram na rede

Figura 6 – Tipos de comportamento individual



Fonte: (ZAFARANI; ABBASI; LIU, 2014)

social. *Usuário-Entidade* é onde um usuário se relaciona com entidades específicas que não sejam grupos de pessoas ou usuários singulares. Deve-se perceber que todos estes tipos de comportamento partem de um único indivíduo, isso o caracteriza como comportamento individual.

Conforme (ZAFARANI; ABBASI; LIU, 2014) a análise do comportamento individual nos ajuda a entender como determinados fatores afetam os comportamentos de determinados indivíduos e para que se possa realizar esta análise precisamos entender e conseguir relacionar padrões de comportamento com determinadas características esperadas na rede social, a partir disso verificar se esta característica é encontrada no comportamento do indivíduo.

2.6.1 Método Para Análise de Comportamento

Um método citado por (ZAFARANI; ABBASI; LIU, 2014) pode ser utilizado para que se analise o comportamento de usuários que estão na rede social e para isso deve-se atender aos requisitos exemplificados a seguir:

- Comportamento deve ser possível de observar: Para isso precisamos ter acesso aos dados do usuário a ser analisado, isso significa ter acesso a alguma informação que explicita sua opinião ou algo semelhante.
- Características Relevantes: Deve-se definir quais são as características que são interessantes para a análise, ou seja, são necessários fatos que relevem algo para que posteriormente possa ser realizada uma avaliação clara de ações do usuário analisado.
- Associação entre Características e Comportamento: Este passo tem como objetivo ligar

as características escolhidas no passo anterior a tipos específicos de comportamentos. Visando ligar as características levantadas com um comportamento.

- Avaliação da Estratégia: Esta etapa avalia se todos os dados de características são realmente ligados a elas ou se podem ter sido confundidos com algum fator externo que não representa o que se procura.

2.7 TRABALHOS RELACIONADOS

Nesta seção serão apresentados trabalhos relacionados ao tema, para que se possa analisar resultados obtidos por outros autores. O primeiro apresenta um estudo de caso que analisa dados disponíveis em redes sociais relacionados a alimentação, para que possa transformá-los em informações que auxiliem a direcionar o foco de empresas de marketing digital. O segundo trabalho que será apresentado trata da análise de sentimento utilizando tweets referentes à copa do mundo 2014. Ambos tratam de mineração de textos e extração de dados de redes sociais. O terceiro e o quarto trabalho vão apresentar casos relacionados a segurança, que é o assunto a ser abordado nesse trabalho. De tal forma poderemos avaliar as implementações em casos distintos verificando sua aplicação para diferentes estudos de caso.

2.7.1 Análise de Dados Para Apoio a Tomada de Decisão Com Ênfase no Marketing Digital

O trabalho apresentado por (LONGHI, 2017) abordou conceitos de *marketing* digital bem como as estratégias utilizadas para isso e sua importância. Tratou sobre a extração de conhecimento em textos e falou sobre quais foram os métodos e técnicas utilizados para a mineração de dados textuais. Abordou diferentes maneiras de coletar dados em redes sociais e tratou de um estudo de caso referente a alimentos saudáveis com o intuito de auxiliar uma empresa do ramo a tomar decisões referentes a marketing digital.

Primeiramente (LONGHI, 2017) levantou um grupo de palavras chave as quais foram apresentadas pela empresa como mais procuradas por seus clientes e a partir daí foram utilizadas como base na pesquisa por *tweets* de usuários. Para a busca dos tweets foi utilizado uma aplicação JAVA que por meio de conexão HTTP realiza a busca no Twitter em determinadas regiões definidas pelo usuário. A ferramenta utilizada na busca de tweets já realiza um pré-processamento de dados removendo caracteres especiais, vírgulas, pontuações e também passa todas as letras para minúsculas.

Após a etapa de coleta (LONGHI, 2017) realizou a etapa de pré-processamento onde definiu classes de palavras específicas, verificando se as mesmas se encontravam no tweets relacionados. Após este procedimento foi realizado o processo de agrupamento utilizando o algoritmo K-Means e Apriori. Então foram obtidos os agrupamentos da Tabela 4 que indica o que as pessoas procuram em determinado tipo de produto, ao exemplo de bolo, pão e biscoito

que foram agrupados ao grupo sem glúten o que indica que neste nicho de produtos as pessoas buscam alternativas que não contenham glúten.

Tabela 4 – Resultado Alimentos – K-Means

Sem Lactose	Sem Glúten	Diet	Light
Pão, Leite Condensado, Bolos, Doce de Leite, Biscoitos e Creme de Leite.	Bolos, Pão e Biscoitos	Bolo, Pão, Biscoito e Achocolatado	Pão, Bolo, Biscoito, Iogurte, Creme de Leite, Biscoitos e Doce de Leite.

Fonte: (LONGHI, 2017)

Com estes resultados o autor informa a possibilidade de que a empresa possa definir de melhor forma o foco de marketing de seus produtos. Ao fim do trabalho foi realizado um questionário com a empresa estudada e a mesma respondeu que a análise foi esclarecedora, pois de acordo com o informado pela empresa eles não tinham dados concretos a respeito de preferência de clientes.

2.7.2 Mineração de Textos: Análise de Sentimento Utilizando Tweets Referentes à Copa do Mundo 2014

O respectivo trabalho escrito por (CARVALHO FILHO, 2014) teve como objetivo analisar *tweets* referentes a Copa do Mundo de 2014, durante a ocorrência dos jogos e classificá-los de acordo com o sentimento neles expresso. Para isso o autor utilizou o processo de mineração de textos.

O método utilizado por (CARVALHO FILHO, 2014) começa com a coleta de *tweets* relacionados à copa. Para isso ele criou um script em *python* que busca por *tweets* que contenham as *hashtags* informadas pelo usuário e neste caso específico foram informados parâmetros relacionados a Copa do Mundo 2014. O processo de coleta ocorreu paralelamente com o acontecimento dos jogos. Após a coleta de dados o autor iniciou a etapa de pré-processamento onde removeu os caracteres especiais e *stopwords* utilizando a plataforma NLTK (*Natural Language ToolKit*).

Para a etapa de análise de sentimentos o autor utilizou o algoritmo Naive Bayes do software Apache Mahout, onde primeiramente ele separou uma pequena quantidade de *tweets* e classificou os mesmos manualmente para que pudessem servir de modelo de classificação para o software. Posteriormente ainda dividiu este conjunto em dois subconjuntos denominados de treino e teste. A partir do conjunto de treino criou-se um modelo de classificação o qual foi utilizado para classificar os demais *tweets* a respeito dos jogos.

Como resultados o autor analisou a precisão obtida em cada uma das categorias definidas como neutro, positivo, negativo e ambíguo. Utilizando os conjuntos de treino, teste e

um terceiro conjunto denominado pelo autor de amostra. Abaixo na Tabela 5 são mostrados os dados coletados no trabalho.

A partir desta tabela pode-se verificar que o algoritmo utilizado encontrou uma dificuldade maior para classificar tweets ambíguos, o que pode ser justificado conforme descrito pelo autor pela aparição de termos positivos e também negativos neste tipo de classificação, o que dificulta o trabalho do algoritmo.

Tabela 5 – Resultados Obtidos por (CARVALHO FILHO, 2014)

CONJUNTO	POSITIVO	NEGATIVO	NEUTRO	AMBÍGUO	ACURÁCIA
Treino	94%	93%	84%	75%	88,91%
Teste	82%	84%	62%	40%	74,80%
Amostra	98%	76%	75%	36%	85,57%

Fonte: O Autor

2.7.3 Detectando Clusters de Contas Falsas em Redes Sociais Online

No trabalho (XIAO; FREEMAN; HWA, 2015) os autores relatam que um dos meios mais utilizados por usuários mal intencionados enviarem spam, cometerem crimes e diferentes tipos de fraudes nas redes sociais é a utilização de contas *fake*², devido a facilidade de se criar um grande número de contas falsas. A partir daí os autores ressaltam a importância de se detectar e agir nestes tipos de conta o mais rápido possível. Isso mantém a confiança dos usuários na rede social.

Para os experimentos (XIAO; FREEMAN; HWA, 2015) utilizou uma base de dados com 275.000 contas do *LinkedIn*. Estas contas foram registradas ao longo de um período de 6 meses e 55% destas foram classificadas pela equipe de segurança da rede social como falsa ou contas de spam. A partir disso foram utilizadas para treinar o algoritmo de classificação. Foram utilizados para implementação os algoritmos floresta aleatória, SVM(*Support Vector Machine*) e regressão logística.

Os dados utilizados para a classificação das contas foram o IP(*Internet Protocol*)³ dos usuários e a data em que a conta foi criada. Após as contas terem sido classificadas os autores utilizaram uma métrica de avaliação chamada AUC(*Area Under the ROC Curve*), para que pudessem comparar o funcionamento dos algoritmos e definir qual se saiu melhor ou pior na classificação. Este método utiliza a área abaixo de uma curva. Área esta que é dividida em duas partes, falsos positivos e verdadeiros positivos. O método segue o parâmetro que quanto maior for a área abaixo da curva melhor é o classificador, pois isso indica o quão longe ele está de classificar errado uma das contas. A Tabela 6 demonstra os resultados obtidos.

²Contas falsas que se fazem passar por outras pessoas.

³Identificação de um dispositivo conectado a uma rede privada ou pública, cada dispositivo tem um IP único.

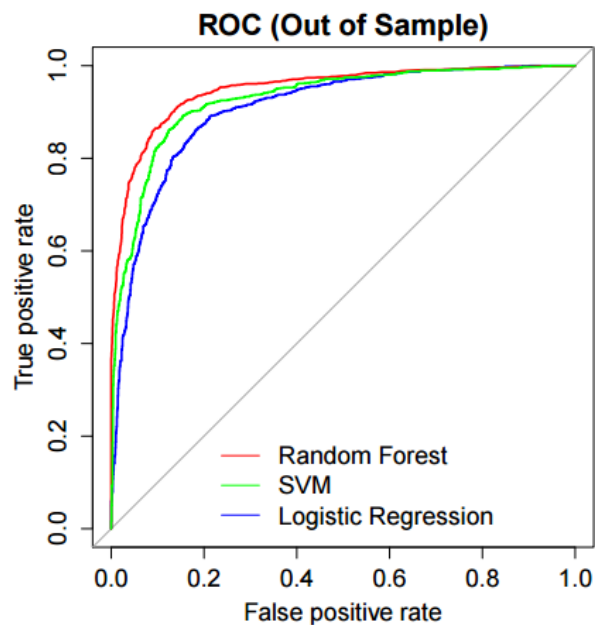
Tabela 6 – Performance dos modelos para dados fora do conjunto de treinamento

Algoritmo	AUC
Floresta Aleatória	0,954
Regressão Logística	0,917
SVM	0,922

Fonte: (XIAO; FREEMAN; HWA, 2015)

Na Figura 7, pode-se encontrar a curva ROC respectiva aos dados da Tabela 6 que compara os três tipos de algoritmos utilizados, onde (XIAO; FREEMAN; HWA, 2015) encontraram os melhores resultados com o algoritmo floresta aleatória.

Figura 7 – Comparação da curva ROC (*Receiver operating characteristic*) para diferentes modelos com dados fora do conjunto de teste



Fonte: (XIAO; FREEMAN; HWA, 2015)

Após a execução dos testes (XIAO; FREEMAN; HWA, 2015) rodaram o algoritmo de floresta aleatória com dados atuais e diários no *LinkedIn*, onde foram encontradas cerca de 15000 contas falsas em um período de 15 dias. Ao fim do artigo os autores sugerem para trabalhos futuros utilizar um algoritmo diferente e falam sobre o modelo *K-Means*.

2.7.4 Detectando Conteúdo Malicioso no Facebook

No trabalho apresentado por (DEWAN; KUMARAGURU, 2015) é relatado que os eventos de grande porte são impulsionadores para que as entidades mal intencionadas possam disseminar suas falcatruas. Os autores de (DEWAN; KUMARAGURU, 2015) citam o *Facebook* como sendo a rede social online como a favorita pelas entidades mal intencionadas, uma

vez que a rede social é a maior e mais ativa rede em atualidade. No artigo (DEWAN; KUMARAGURU, 2015) os autores estimam que os *spammers*⁴ tenham lucrado cerca de US \$200 milhões apenas publicando links que seriam clicados gerando ganhos em anúncios de publicidade em um período de cerca de um ano. O artigo cita que existem inúmeros tipos de ataques maliciosos na rede social *online*. Estes podem ser propagandas contendo vírus, perfis falsos entre outros, mas o artigo focou em apenas identificar um tipo de fraude que são as URLs⁵ com conteúdo malicioso e também criar um método automatizado para que isso funcione em tempo real.

No artigo (DEWAN; KUMARAGURU, 2015) os autores relatam que os métodos existentes se revelam ineficientes nos casos de acontecimentos em tempo real, pois estes identificam URLs maliciosas a partir de bancos de dados de listas negras. Quando se trata de novos acontecimentos as URLs criadas ainda não estão incluídas em tais bancos de dados.

A partir disso os autores analisaram as URLs contidas em seu banco de dados de itens maliciosos, levando em conta três características em especial: A primeira foi o conteúdo textual ligado as URLs, a segunda foi o gerador da postagem e a terceira são as plataformas que criaram as URLs.

Na análise textual, os autores encontraram características que estão ligadas a exploração da curiosidade dos usuários, ou seja, não mais com promoções que prometem itens gratuitos ou preços milagrosos. A característica mais encontrada nas mensagens que continham itens maliciosos foram conteúdo adulto ou nudez estando em um total de 52% das mensagens, em seguida mensagens com conteúdo de reivindicações e conteúdo de ódio com um total de 45,2%.

No segundo quesito os autores avaliaram dois tipos distintos de geradores páginas e usuários, classificando previamente as postagens maliciosas verificaram que apenas 29% destas postagens partiam de páginas e não de perfis de usuários, logo focaram em postagens criadas por perfis de usuários.

Como o terceiro quesito é a plataforma de postagem foi necessário avaliar quais seriam os possíveis meios de postagem e definiram os tipos web, mobile e outros. E o que mais postou mensagens maliciosas foi a categoria outros, onde podem estar os algoritmos de postagens automáticas responsáveis por disseminar as postagens.

Para a classificação os algoritmos utilizados foram Naive Bayes, Árvore de decisão, AdaBoost. O melhor resultado encontrado foi no algoritmo de Árvore de Decisão com uma acurácia de 86,9%.

Ao fim do trabalho os autores implementaram um *plugin* para o navegador que avisa aos usuários postagens que contém links maliciosos.

⁴Criadores e disseminadores de spam.

⁵Um URL se refere ao endereço de rede no qual se encontra algum recurso informático, como por exemplo um arquivo de computador ou um dispositivo periférico.

2.8 CONSIDERAÇÕES FINAIS

Neste capítulo foram apresentados os métodos de extração de conhecimento em textos, bem como, suas etapas e os algoritmos geralmente utilizados para que este processo ocorra da maneira mais eficaz e rápida. Também foram expostos requisitos para a coleta dos dados textuais em redes sociais *online* que utilizam softwares denominados *web crawlers*. Outro tópico importante abordado neste capítulo foi a limpeza destes dados, ou seja, o pré-processamento dos dados para que se possa utilizá-los na análise.

Foram apresentados trabalhos relacionados com o tema abordado e também trabalhos que trataram de temas diferentes mas que utilizaram meios de implementação similares. Pode-se assim verificar que independentemente do tema a ser tratado, o método é funcional e pode ser aplicado, uma vez que foram avaliados os resultados obtidos e as taxas de acerto foram relativamente altas. A partir da tabela comparativa dos trabalhos apresentada no Apêndice A, pode-se identificar que algoritmos tradicionais de classificação e agrupamento de dados são comumente utilizados (árvore de decisão e k-means, entre outros).

Outro ponto importante a se salientar foi o entendimento adquirido sobre as redes sociais *online*, como os dados são apresentados, o que realmente significam e como podemos filtrá-los e entendê-los para que não sejam somente um aglomerado de letras e palavras. A maneira como os usuários se portam na rede social *online*, o que procuram e como funciona o relacionamento entre indivíduo e comunidade também foram possíveis de identificar após a conclusão deste capítulo. Ainda, foi visto como os fraudantes utilizam este meio social *online* para poder disseminar suas comportamentos ilegais e quais os diferentes métodos utilizados por eles, além de verificar como são criadas as soluções para que se possa combater e prevenir este tipo de situação que se torna tão incomoda e desagradável para o usuário, que tenta utilizar a rede social de maneira a não prejudicar outras pessoas ali inseridas. Um forma de analisar os dados textuais das redes sociais é por meio da análise de sentimentos, que estão explícitos nas postagens. A análise de sentimentos permite que seja reconhecida que toda mensagem reflete uma emoção. Emoções positivas refletem opiniões favoráveis, indicando que o conteúdo de um perfil foi benéfico, útil, vantajoso ou proveitoso. Já emoções negativas refletem conteúdos inapropriados, nocivo, prejudicial ou perigoso. Trabalhos relacionados, como o de (CARVALHO FILHO, 2014), indicam que existem uma lacuna de sistemas nesta área a serem integrados em aplicações mais abrangentes.

3 ANÁLISE DE SENTIMENTOS NO FACEBOOK: IMPLEMENTAÇÃO E RESULTADOS

A análise de sentimentos é uma subárea que pertence a mineração de textos. Manipular grandes volumes de textos que podem conter opiniões e percepções úteis se torna inviável em um curto espaço de tempo. Além disso, esta tarefa pode ser bem cansativa para um ser humano, que apesar da sua grande capacidade de compreensão, pode facilmente se distrair, não prestando atenção permanente nos detalhes. Neste cenário meios automáticos podem apoiar o ser humano, realizando ao menos uma parte da análise textual, produzindo sumários e antecipando percepções positivas ou negativas.

A fim de proceder com a análise de sentimentos em um conjunto textual oriundo do Facebook desenvolveu-se uma ferramenta de software apresentada neste capítulo. O software realiza a coleta de dados em grupos ou páginas na rede social *Facebook* e os examina. O objetivo principal do software consiste em analisar as postagens dos usuários de grupos alvo e classificar o conteúdo expresso em termos de emoções. Também são avaliadas postagens que possuam linguagem inapropriada, termos que indiquem fraudes bancárias, a fim de criar um alerta ao usuário para que ele tenha conhecimento do risco apresentado na página do grupo.

3.1 TEORIAS DA EMOÇÃO

Com o objetivo de analisar e classificar as postagens é necessário que se entenda os métodos para análise de emoções. Também é preciso definir quais são as palavras que indicam perigos a segurança e também linguagem inapropriada. Nesta seção serão demonstrados os métodos utilizados para que essa análise pudesse ser feita.

3.1.1 Léxico de Emoções

Para que se possa entender e utilizar a análise de emoções é necessário estudar os termos os quais realmente as representam. Para isso é necessário que se colete as palavras relacionadas a emoções, posteriormente deve-se introduzi-las em contextos que representem personalidade, estado físico e cognitivo. Posteriormente deve-se avaliar as relações afetivas utilizando técnicas de análise fatorial ou escala multidimensional. No caso do trabalho apresentado foi utilizado um léxico que já havia sido desenvolvido a partir de tais métodos. O léxico *WordNet-Affect* é uma extensão do léxico *WordNet* (STRAPPARAVA; VALITUTTI et al., 2004), o *WordNet-Affect* se torna uma base de palavras para que se possa detectar emoções. Uma base que conta com cerca de cinco mil termos. Estas emoções estão divididas em seis grupos distintos alegria, desgosto, medo, raiva, surpresa e tristeza. Estes grupos serão utilizados para classificar as emoções junto ao software implementado. Na Tabela 7 estão alguns exemplos de palavras e suas respectivas emoções.

Tabela 7 – Categorias de emoções

Categoria	Palavras
Alegria	animação, ânimo, comédia
Desgosto	maldoso, chatear, abominável
Medo	desespero, escuridão, friamente
Raiva	assassinar, bravejar, detestar
Surpreza	chocante, espanto, impressionado
Tristeza	cabisbaixo, chorão, culpa

Fonte: O autor.

3.1.2 Palavras relacionadas a segurança e linguagem inapropriada

Além da análise de sentimentos foi implementada uma busca por palavras que possam estar relacionadas a riscos de segurança e também por palavras de linguagem inapropriada. Neste trabalho para a identificação das palavras relacionadas à segurança foram utilizadas palavras apresentadas em (HEPP, 2010). O autor realiza uma busca por identificação de ataques de *phishing* através de mineração textual, com isso ele também identificou algumas palavras que acabaram se relacionando com os ataques e aparecendo frequentemente e com isso criou um corpus onde listou essas palavras. Tendo em vista essa pesquisa e sabendo quais foram as palavras utilizadas foi possível implementar no software uma busca por palavras relacionadas a fraudes de segurança. As palavras que foram levantadas pelo autor foram conta, acesso, banco, crédito, clique, identidade, inconveniência, informação, limitado, registro, minutos, senha, recentemente, risco, social, segurança, serviço, suspenso e cartão.

Já no caso das palavras de linguagem inapropriada foi feito um levantamento de quais são as palavras que podem ser inapropriadas para quem estiver lendo, visto que não há filtro de idade nos grupos e que qualquer um pode acessar esse conteúdo podemos ter crianças visualizando essas postagens.

3.2 APRESENTAÇÃO E IMPLEMENTAÇÃO DO SOFTWARE DESENVOLVIDO

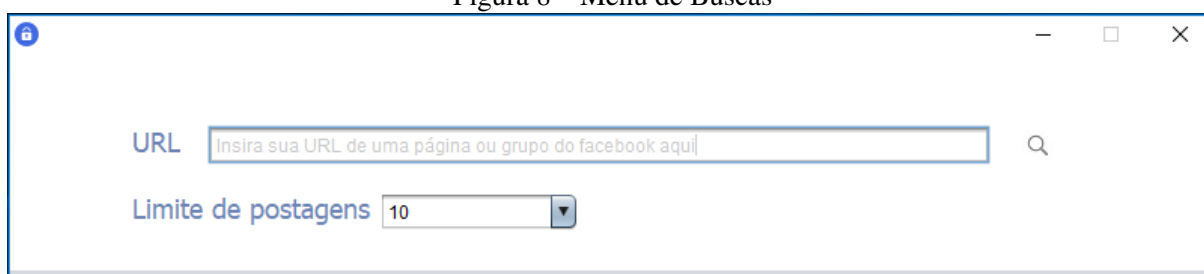
Para que a coleta de dados e análise pudesse ser realizada por usuários sem experiência em programação foi desenvolvido um software de interface relativamente simples responsável por coletar informações e apresentá-las ao usuário. Junto com a apresentação dos dados o software ainda realiza uma análise nas postagens para que o usuário possa interpretar e tirar suas conclusões. Análise esta que conta com uma nuvem de palavras, gráfico de sentimentos encontrados nas postagens, aviso que informa a aparição de palavras de linguagem inapropriada e também palavras relacionadas a fraudes bancárias. Nesta seção serão explicadas essas funcionalidades.

3.2.1 Apresentação do Software

3.2.1.1 Menu de Buscas

O software conta com um menu de busca, onde se encontram a barra de endereços o botão pesquisar e também um limitador de número de postagens coletadas conforme a Figura 8. Inicialmente o usuário deve informar o endereço da página ou grupo no campo URL, posteriormente informar o limite de postagens para que o software busque somente o número de postagens determinado e então clicar no botão de pesquisa para que seja iniciada a busca. Caso o usuário informe um endereço que não seja de uma página ou grupo do *Facebook* o software retorna uma mensagem informando que o endereço inserido é inválido.

Figura 8 – Menu de Buscas

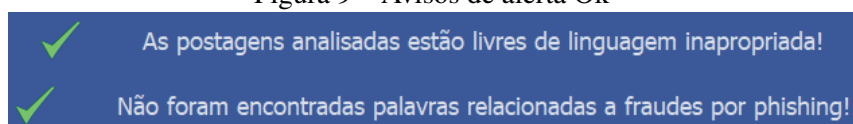


Fonte: O Autor.

3.2.1.2 Avisos de Alerta

Após realizada a busca o software vai apresentar duas mensagens de aviso ao usuário. A primeira é relacionada a linguagem inapropriada, ou seja, informa se nas postagens analisadas foram encontradas palavras de linguagem inapropriada ou não. Já a segunda mensagem informa se foram encontradas palavras que se relacionam a roubo de dados utilizando o método de *phishing*. Na Figura 9 pode-se visualizar os avisos quando não foram encontradas nenhuma palavra relacionada relacionada aos riscos analisados, nesse caso o software não encontrou palavras iguais a que estão em sua base de dados.

Figura 9 – Avisos de alerta Ok



Fonte: O Autor.

Caso o software encontre alguma palavra que se encaixe no escopo analisado ele irá informar ao usuário alterando o descritivo da mensagem e também a imagem antes do texto, para sinalizar o alerta de segurança. Podemos visualizar essas mudanças na Figura 10.

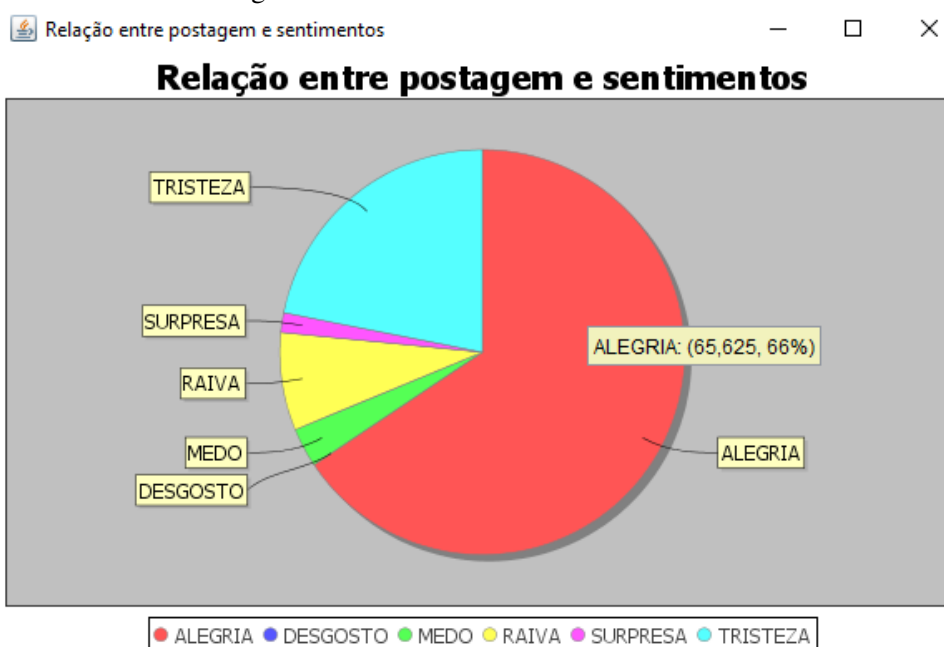
Figura 12 – Botão Análise Gráfica



Fonte: O Autor.

que contém um gráfico circular, este gráfico apresenta a porcentagem de palavras relacionadas a cada sentimento que foram encontradas nas postagens analisadas. Na Figura 13 podemos visualizar a janela do gráfico.

Figura 13 – Análise Gráfica de Sentimentos



Fonte: O Autor.

3.2.1.5 Exportar Postagens Analisadas

A opção de exportar as postagens analisadas permite ao usuário que ele possa ver as postagens que foram analisadas e quem foram os usuários responsáveis por redigir o conteúdo que nelas se encontram. O software permite também que o usuário veja como essas postagens foram classificadas e caso queira ir mais a fundo na pesquisa o software ainda exporta o link da rede social para que se possa encontrar essa postagem no *Facebook*. Para isso é necessário clicar no botão da Figura 14.

Com isso o software vai gerar um arquivo em extensão PDF o qual será exibido na tela do usuário conforme o exemplo da Figura 15, onde se tem uma das postagens analisada com os dados relacionados a ela e também sua classificação.

Figura 14 – Botão Exportar Postagens



Fonte: O Autor.

Figura 15 – Exemplo de Postagem Exportada

```

-----
Usuário:
Globo Esporte
Postagem:
Duelo reúne times que tentam fugir do Z-4
Classificação da postagem:
MEDO,
Link da postagem:
http://fb.com/22344895408\_10155920690745409
-----
  
```

Fonte: O Autor.

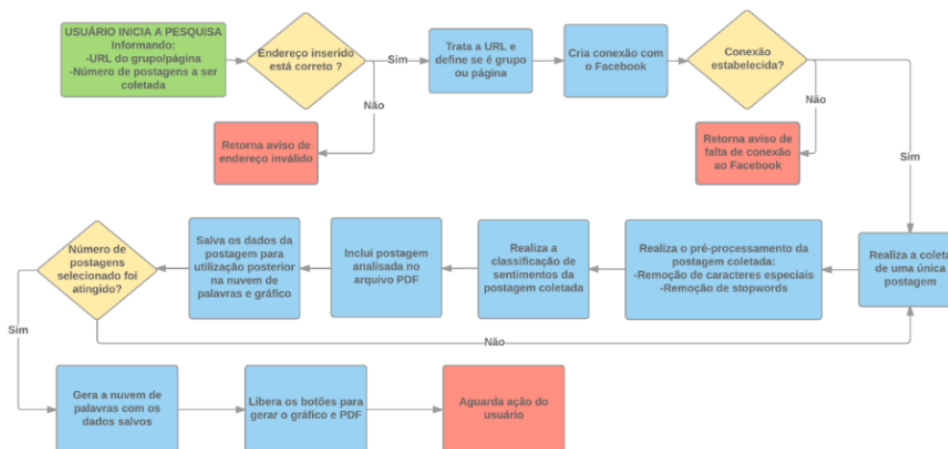
3.2.2 Implementação do Software

Para que fosse possível a coleta de dados foi implementado um *web crawler*, que como dito no capítulo anterior é um software responsável por adquirir informações de sites da internet. Neste trabalho mais especificamente foi desenvolvido um script em Java que realizou a aquisição de dados contando com o auxílio da ferramenta *Facebook Graph API Explorer* a qual é responsável somente por fornecer o código de autorização para coleta, código esse que é conhecido como *Token*. Para que se possa utilizar esta ferramenta o *Facebook* precisa fornecer uma autorização de desenvolvedor que pode ser solicitada no site descrevendo o tipo de trabalho que se pretende realizar.

Ainda no script Java foi utilizada a biblioteca *Restfb* que é a responsável pela criação da conexão com a rede, ou seja, vai criar a conexão de dados entre o script desenvolvido e o *Facebook*. Para que o script retorne com os dados corretos é necessário que seja informado o número de identificação do grupo no Facebook e também os tipos de dados a serem coletados, no caso deste estudo as postagens do grupo determinado e os usuários responsáveis por elas.

Na Figura 16 pode-se visualizar o fluxograma de funcionamento do software que foi elaborado. Inicialmente o usuário deve entrar com o endereço de busca a quantidade de postagens para que seja realizada a coleta da primeira postagem e posterior classificação, ao finalizar esse fluxo o software mantém disponível todos os dados para visualização do usuário.

Figura 16 – Fluxograma de funcionamento do software



Fonte: O Autor.

3.2.2.1 Identificando o grupo ou página de coleta a partir de uma URL

Para que seja possível direcionar a coleta de dados a um determinado grupo ou página que se encontra no *Facebook* é necessário fornecer a identificação do endereço em que se procura realizar a coleta. Essa identificação está contida na URL encontrada na aba do navegador, no caso de páginas a informação se encontra da maneira exposta na linha 1 da Tabela 8. A parte que identifica o conteúdo que estamos buscando é a "paginadeinteresse", ou seja a partir dessa parte da URL é que o banco de dados do *Facebook* vai distinguir a página em que desejamos realizar a pesquisa.

A busca em grupos é bastante parecida, a não ser pela diferenciação que possui na URL onde se tem antes da parte de interesse a palavra *groups* conforme mostrado na linha 2 da Tabela 8, esta diferença é responsável por fazer a distinção de endereços de páginas ou grupos na rede social.

Tabela 8 – Tipos de URL páginas x grupos

Tipo	URL
Página	https://www.facebook.com/paginadeinteresse/
Grupo	https://www.facebook.com/groups/grupodeinteresse/

Fonte: O Autor.

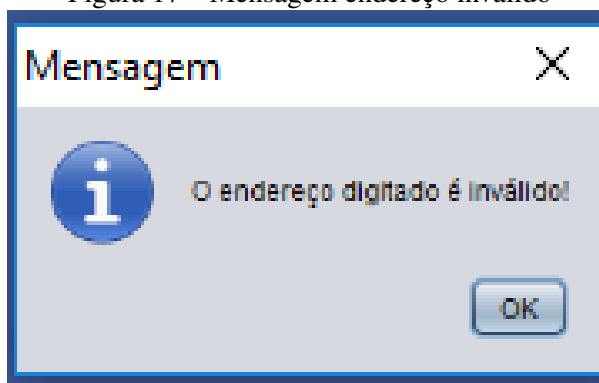
Após a análise de URL e suas diferenças foi visto que é necessário que o script em JAVA consiga distingui-las, então para isso foi implementada uma parte no programa que identifica se a URL informada pelo usuário contém a palavra *groups*, o que sinaliza que será necessário tratar o endereço fornecido de maneira diferenciada.

Após a distinção entre URLs é feita a extração da parte de interesse, para isso precisamos pegar somente a parte que contém o destino que esta se buscando, no caso de páginas deve-se selecionar somente os caracteres que estejam após "https://www.facebook.com/" e

antes de "/", quando se trata de endereço de grupos deve-se coletar o que se encontra após "https://www.facebook.com/groups/" e antes de "/". Para isso foi utilizado um script que localiza determinados caracteres em uma *string* e onde estão localizados na mesma, para que possa salvar especificamente os que são necessários.

Juntamente com isso é feita a verificação onde caso não sejam encontrados os padrões necessários o programa informa ao usuário que o endereço digitado não é válido, conforme a Figura 17.

Figura 17 – Mensagem endereço inválido



Fonte: O Autor.

Finalizando a parte de distinção entre endereços e tendo realizado a separação do endereço do grupo e página do resto da *URL* é informado ao servidor do *Facebook* quais dados estamos procurando dentro do grupo ou página escolhida, pois pode-se solicitar diferentes tipos de dados como por exemplo, nome do grupo ou página, quem são os usuários proprietários ou mesmo quando foi a última vez em que o grupo ou página foram atualizados. No caso da coleta em questão informaremos ao servidor que estamos buscando informações no *feed* que é como o facebook nomeia a página onde se encontram as postagens, posteriormente informaremos que estamos buscando as informações *from* e *message*, que são respectivamente as solicitações feitas para que o servidor nos retorne, quem foi o usuário responsável pela postagem e qual foi a postagem escrita por ele.

3.2.2.2 Pré-Processamento dos Dados Adquiridos

A etapa de pré-processamento é de extrema importância para que se chegue em um bom resultado, mas é de maior importância para a otimização dos algoritmos posteriores que serão utilizados para extração de conhecimento dos dados textuais. Nesta etapa foi incluído no software de coleta a parte de pré-processamento. Foram removidos os caracteres especiais, que dificultam a análise de comparação termo a termo. Muitas pessoas nas redes sociais não se importam em acentuar as palavras escritas corretamente. Então após essa etapa de pré-processamento será obtida uma igualdade no que se refere a acentos e pontuações. É importante salientar que os caracteres especiais não foram removidos, mas sim substituídos por caracteres

de uso normal, conforme a Tabela 9.

Tabela 9 – Remoção de caracteres especiais

ANTES	DEPOIS
econômico	economico
informações	informacoes

Fonte: O Autor

Após finalizada a remoção de caracteres especiais foi realizada a segunda etapa do pré-processamento, a remoção das *stopwords*, que são as palavras que na contextualização da postagem não tem uma representação significativa no contexto da análise. Estas *stopwords* são geralmente palavras de ligação. Um conjunto de *stopwords* é considerado uma *stoplist*. Esta *stoplist* não tem uma definição fechadas de palavras que constam em seu interior, constam como padrão da língua portuguesa as mais comuns ao exemplo de a, e, com, porém entre outras, mas as palavras que não constam e possivelmente podem não ter significado de modo geral devem ser definidas incluídas na *stoplist*. Na Tabela 10 podemos encontrar um exemplo de postagem coletada onde foram removidas as *stopwords* e também os caracteres especiais.

Tabela 10 – Remoção de *stopwords* e caracteres especiais

ETAPA	RESULTADO ADQUIRIDO
ANTES	Chegando agora, uma boa noite e um forte abraço para vocês família.
DEPOIS	Chegando agora boa noite forte abraco voces familia

Fonte: O Autor

Pode-se perceber que no final do pré-processamento tivemos uma redução de doze para apenas oito termos na frase. Esta redução de termos não prejudicou o entendimento da frase, ou seja, ainda é possível identificar se o sentido que a frase esta sendo exposta é negativo ou positivo. Outra avaliação que devemos fazer é em relação aos caracteres especiais onde nem sempre o usuário que esta postando a mensagem irá pontuar ou acentuar a frase corretamente. Por isso a substituição dos caracteres especiais nos dados coletados se torna tão importante.

3.2.2.3 Classificação das postagens

A classificação das postagens é simples e feita após o pré-processamento, as postagens são analisadas termo a termo. A comparação é realizada a partir de uma base de dados criada com o léxico *WordNet-Affect* já apresentado anteriormente. O software seleciona cada postagem após a remoção de *stopwords* e remoção de caracteres e então envia esses dados para uma

classe secundária. Classe essa que faz a análise e retorna com um contador que é o número de palavras relacionadas a sentimentos que foram encontradas na postagem este contador é utilizado para alimentar todas as outras funcionalidades do software. O software faz essa busca na base de dados de palavras para cada classe de sentimento. Cabe salientar que a postagem ao ser analisada pode ter mais do que um sentimento expresso na mesma e então o software indica que todos estes sentimentos estão expressos na postagem analisada. Após a busca de sentimentos o software ainda realiza a busca por palavras relacionadas a problemas de roubo de dados com atividades relacionadas a *phishing* e também linguagem inapropriada.

3.2.2.4 Implementação do Gráfico

Para a implementação do gráfico onde são apresentadas a relação de postagens e sentimentos, foi utilizada a biblioteca *JFreeChart*¹ que é de código aberto e permite criar uma grande variedade de gráficos que podem ser interativos ou não, foi escolhida essa biblioteca devido a sua facilidade de uso e também pelo número de possibilidades de gráficos que ela nos permite criar alterando apenas alguns comandos. Para que o gráfico seja criado e apresente os dados que são solicitados é necessário apenas definir alguns parâmetros. Primeiro foi definido o tipo de gráfico a ser utilizado, no caso foi escolhido o gráfico de "pizza" conforme é conhecido o gráfico de setores, este foi escolhido de maneira a apresentar a porcentagem de postagens de cada sentimento que foram encontradas. Após definido o tipo de gráfico devemos informar ao software quais os dados o gráfico irá exibir, conforme falado na seção de classificação o classificador retorna números inteiros que são os responsáveis por alimentar esse gráfico.

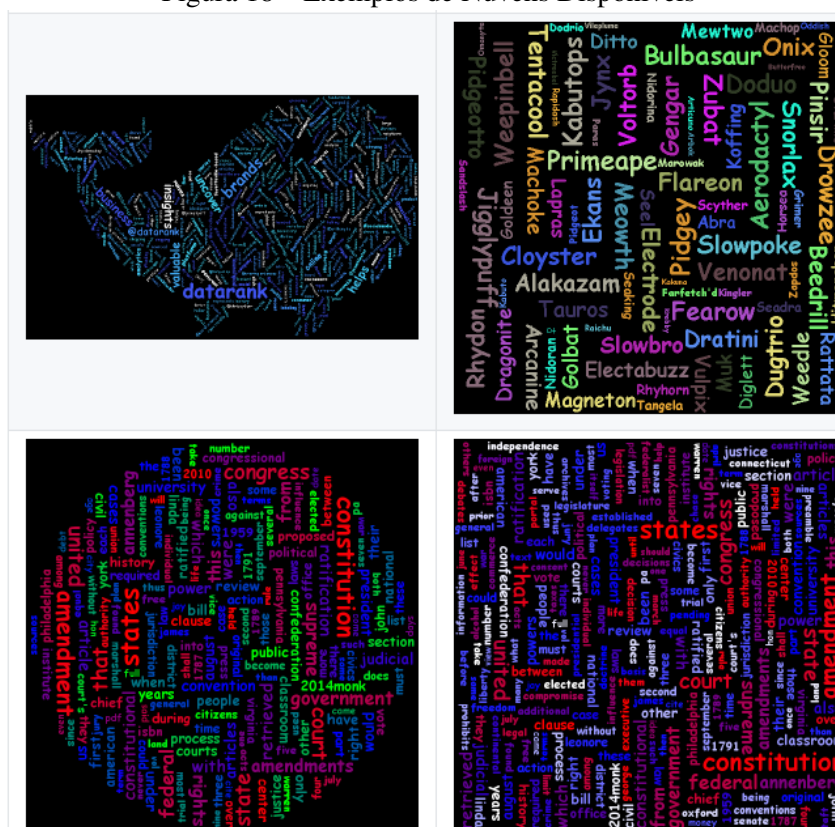
3.2.2.5 Implementação da Nuvem de Palavras

Na implementação da nuvem de palavras foi utilizada a biblioteca Kumo². Esta biblioteca permite a criação de diferentes tipos de nuvens de palavras, podendo variar seu formato cores e até mesmo a fonte em que as letras são escritas, conforme a Figura 18. Para o correto funcionamento dessa biblioteca e assim exibição da nuvem de palavras é necessário parametrizar algumas entradas, deve-se informar o tamanho da nuvem desejada, qual é o seu formato desejado e informar quais são os textos a serem analisados, no caso do software implementado todas as postagens coletadas são armazenadas em um arquivo em formato de texto durante a coleta das informações, ao final da coleta esse arquivo é chamado pelo método *FrequencyAnalyzer* que está contido dentro da biblioteca Kumo, esse método é responsável por analisar a frequência em que as palavras aparecem e que devem aparecer na nuvem de palavras. A partir disso é gerada uma figura com as palavras que foram encontradas mais vezes, diferenciando a sua quantidade de aparição pelo tamanho da fonte.

¹<http://www.jfree.org/jfreechart/>

²<https://github.com/kennycason/kumo>

Figura 18 – Exemplos de Nuvens Disponíveis



Fonte: <https://github.com/kennycason/kumo>

3.2.2.6 Implementação da Opção Exportar PDF

A opção exportar em PDF se fez necessária para que o usuário pudesse visualizar o conteúdo que foi analisado pelo software. Então a partir dessa necessidade foi necessária a utilização da biblioteca *Itex*³. Essa biblioteca nos permite jogar dados textuais e também imagens para dentro de um arquivo com extensão PDF. Foi escolhida esta extensão de arquivo para realizar a exportação das postagens devido a facilidade de implementação e também pela compatibilidade quanto a abertura do arquivo, pois hoje em dia grande parte dos computadores possuem um visualizador de arquivos PDF. Para a implementação no software foi criada uma classe chamada *gerarpdf*, essa classe é responsável por receber os dados textuais que serão incluídos no arquivo PDF, dentro do da estrutura principal do programa foi criado um método que invoca essa classe e envia os dados a ela.

3.3 DESCRIÇÃO DO EXPERIMENTO

Para que se possa definir o método e como aplicá-lo primeiramente deve-se estudar o cenário que se está trabalhando. Quando este estudo não é previamente realizado alguns proble-

³<https://itextpdf.com/>

mas podem surgir posteriormente, tais como a escolha errada do método de análise de dados, tornando a análise muito mais complexa ou inviabilizando-a. Duas definições importantes neste trabalho foram: a escolha da rede social utilizada e quais os dados coletados.

3.3.1 Escolha da Rede Social

Com a vasta gama de redes sociais online encontradas na internet, deve-se avaliar o tipo de dados que se procura e também se esta rede concede a permissão para que se tenha acesso a este tipo de dados. Uma vez que a análise é voltada para dados textuais já se consegue aplicar o primeiro filtro, redes sociais voltadas a imagens e vídeo já podem ser desconsideradas. Entre elas o *Instagram*⁴, *Snapchat*⁵, *Pinterest*⁶ entre outras que tem seu foco voltado a este tipo de mídia.

A próxima qualificação para que uma rede social possa ser escolhida para análise deve ser a disponibilidade de coleta de seus dados, uma vez que sem a autorização da rede social se torna uma tarefa muito mais árdua de se realizar de forma automática, para não se dizer impossível. Entre as redes sociais mais conhecidas duas delas trabalham com dados textuais e disponibilizam estes dados para coleta. A primeira é o *Twitter*, que possui a ferramenta *Twitter API*⁷ a qual permite que sejam coletados dados e também realizadas postagens na página do usuário do *Twitter*. No entanto a rede social *Facebook* também possui uma ferramenta com características semelhantes que permitem a coleta de dados em sua rede social ela se chama *Graph API Explorer*⁸.

De tal maneira qualquer uma das duas redes sociais poderia ser utilizada no desenvolvimento do trabalho, mas devido a algumas características específicas o *Facebook* se mostra mais interessante para a análise, pois a interação entre usuários é menos limitada, uma vez que o *Twitter* só permite que você poste mensagens na sua página ou na página de algum outro usuário o *Facebook* permite que o usuário entre em grupos e interaja diretamente com diferentes tipos de pessoa em uma comunidade sem mesmo conhecê-las. Outra vantagem encontrada no *Facebook* é o seu número de usuários que de acordo com os dados divulgados pelo próprio site conta com 1,8 bilhões de usuários inscritos em quanto o *Twitter* em sua última divulgação revelou ter somente 310 milhões de usuários, ou seja, menos de um quinto do número de usuários inscritos no concorrente *Facebook*.

⁴<https://www.instagram.com>

⁵<https://www.snapchat.com>

⁶<https://www.pinterest.com>

⁷<https://dev.twitter.com/overview/api> - Ferramenta de coleta e inserção de dados no Twitter

⁸<https://developers.facebook.com/tools/explorer/> - Ferramenta de coleta e inserção de dados para o Facebook

3.3.2 Definição dos Grupos de Coleta

Neste ponto foram definidas as páginas e grupos que teriam os dados coletados para análise. Foi decidido por analisar três páginas de conteúdo distinto, buscando avaliar as postagens e levantar opiniões sobre determinados assuntos. A primeira é de um serviço público onde as pessoas podem realizar postagens expressando seus sentimentos a cerca disso. A segunda trata de um grupo de debates em geral acerca de um equipamento eletrônico com a intenção de saber a aceitação e retorno dos compradores. Já o terceiro grupo analisado foi o de uma franquia de restaurantes, onde a ideia foi visualizar a opinião sobre atendimento, comida e demais opiniões em torno da franquia.

3.3.3 Definição de Dados para Coleta

Uma vez que se busca analisar as ações de um usuário avaliando seu comportamento com outros usuários na rede social deve-se lembrar que estamos falando de um tipo de comportamento individual. Conforme apresentado na seção 2.4, um comportamento se classifica como *Usuário-Comunidade* quando o usuário analisado está publicando algo voltado ao grupo que está inserido. Logo, deve-se buscar coletar dados que possam representar este tipo de interação.

Avaliando a rede social do Facebook e os dados disponibilizados por ela, foi concluído que os melhores representantes do tipo de dados que se está procurando seriam as postagem em grupos voltados a algum tipo de assunto em específico. Isso permite representar a relação de um usuário com uma comunidade e, caso seja necessário, será possível também incluir traços específicos da pesquisa.

3.3.4 Método de Análise dos Dados Coletados

Possuindo as postagens classificadas pelo software a partir do léxico é necessário que seja feita uma análise para constatar o quão precisa é a classificação realizada. Então foram coletadas amostras de diferentes tamanhos em todos os grupos analisados para que manualmente fosse realizada uma análise levantando o número de postagens classificadas corretamente ou não. O número de amostras para avaliação foi definido em cinquenta, cem e quinhentas postagens. Uma vez que se aumentarmos muito o número de postagens podemos não conseguir grupos que tenham toda essa quantidade disponível. Cabe salientar que entre as postagens analisadas nem todas são classificadas com algum tipo de sentimento, pois muitas postagens não expressam nenhum dos tipos de sentimentos classificados ou mesmo não expressam nenhum tipo de sentimento conforme demonstra a Figura 19.

Figura 19 – Postagem classificada como neutra

Usuário:
Henrique Vasconcelos

Postagem:
MacBook Pro 13" 2012 i7 8Gb 2.9Ghz 750hd TurboBoost 3.1Ghz
FRETE GRÁTIS PARA TODO BRASIL
Máquina em perfeito estado.
Acompanha carregador original 60W.

Classificação da postagem:
POSTAGEM NEUTRA!

Link da postagem:
http://fb.com/858463467574409_1557101574377258

Fonte: O Autor

3.4 AVALIAÇÃO DE RESULTADOS

A avaliação de resultados foi feita a partir da análise realizada manualmente para avaliar a classificação da coleta realizada com quinhentas postagens. Na Tabela 11 pode-se visualizar os resultados da classificação das quinhentas postagens do primeiro grupo. Na classificação manual foi constatado que alguns sentimentos estavam sendo classificados incorretamente. Visando melhorar os resultados, algumas modificações no léxico utilizado foram necessárias.

Tabela 11 – Classificação Inicial do Grupo Mackbook

Classificação	ALEGRIA	MEDO	TRISTEZA	DESGOSTO	RAIVA	SURPRESA
Correta	1	4	14	2	5	4
Errada	22	2	8	1	0	2

Fonte: O Autor.

O sentimento de alegria acabou tendo o pior resultado quando avaliado manualmente. A partir disso foi necessário avaliar quais foram as palavras que se encontravam no léxico que geraram essa classificação errônea. Para isso foi feita uma alteração no código do software. Inicialmente rastreou-se cada postagem com o sentimento de alegria por meio da palavra contida no léxico que tornava a classificação desta postagem como alegria. As palavras encontradas que mais se relacionavam a classificações incorretas quanto ao sentimento de alegria foram *bom*, *melhor*, *favor* e *excelente*, aparecendo mais de três vezes em classificações erradas. Estas palavras foram desassociadas à alegria no léxico.

No decorrer da análise manual foram visualizadas algumas expressões de uso comum na internet que também ajudam a identificar postagens que expressão alegria, tais como risadas (e.g., *hehe*, *kkk*, *haha*, entre outras). Essas expressões foram incluídas no léxico com a intenção de melhorar a classificação do sentimento de alegria.

Outro sentimento que teve um número grande de classificações incorretas foi o da tristeza. Tomando por base o mesmo método utilizado para aprimorar a detecção do sentimento de alegria, foi feito com que o software exibisse as palavras que eram responsáveis por classificar a postagem como relacionada ao sentimento de tristeza. Com isso foi visto que a palavra *pena*, contida no léxico do sentimento tristeza, estava gerando grande parte da classificação errada. Por exemplo, essa palavra é bastante utilizada para realizar perguntas na expressão "vale a pena?". Ela não está relacionada com o sentimento de *sentir pena* de algo, nem com o sentimento de tristeza.

Após realizar-se estas alterações, o número de postagens classificadas corretamente aumentou tanto para o sentimento de alegria quanto para o sentimento de tristeza. Contudo, o sentimento de alegria continua com uma taxa de acerto inferior a cinquenta por cento, conforme pode-se visualizar na Tabela 12.

Tabela 12 – Classificação Final do Grupo Mackbook

Classificação	ALEGRIA	MEDO	TRISTEZA	DESGOSTO	RAIVA	SURPRESA
Correta	6	8	20	3	5	4
Errada	7	2	4	1	0	2

Fonte: O Autor.

O segundo grupo analisado foi o da previdência social. Assim como o primeiro grupo este também teve suas peculiaridades, onde a que se destacou foi o aparecimento da palavra *benefício*. Esta palavra está associada normalmente com felicidade, mas no grupo analisado essa palavra tem o significado de um auxílio recebido do governo por motivo de doença ou aposentadoria. Após a remoção dessa palavra do léxico foi obtido o resultado da Tabela 13. Pode-se visualizar um grande número de postagens com sentimento tristeza, desgosto e raiva que foram classificadas corretamente. A critério de análise isso pode demonstrar a quem está realizando a busca qual é a opinião das pessoas que escreveram as postagens. As postagens evidenciam uma grande insatisfação das pessoas com o serviços prestados. Grande parte das postagens analisadas nesse grupo falavam sobre a má qualidade do serviço, a demora para o agendamento de perícias por parte da previdência social, a falta de descaso com o contribuinte, problemas no site da previdência entre outros problemas que ficam evidenciados quando feita a análise de sentimentos.

Tabela 13 – Classificação do Grupo Previdência Social

Classificação	ALEGRIA	MEDO	TRISTEZA	DESGOSTO	RAIVA	SURPRESA
Correta	5	5	30	22	19	1
Errada	5	0	5	2	3	0

Fonte: O Autor.

Para a fase final de análise de postagens foi realizada a busca na página Outback Brasil, que é uma franquia de restaurantes espalhada por todo o Brasil. Diferente dos outros grupos não foram identificadas palavras a serem removidas do léxico. Nessa página foi localizado um

número maior de postagens com classificação positiva de acordo com a Tabela 14 quando comparado aos outros grupos. O sentimento de tristeza também teve um número alto de postagens o que demonstra a insatisfação de um grupo de clientes com a franquia. Essa insatisfação fica ainda mais clara quando lidas as postagens. Alguns clientes chegaram a relatar que não pretendem mais voltar a restaurantes que estejam ligados a tal rede. Contudo, o grande número de postagens relacionadas ao sentimento de surpresa são de cunho positivo, ou seja, muitos clientes ficaram surpreendidos com a qualidade da comida e do atendimento e tiveram suas expectativas superadas.

Tabela 14 – Classificação da Página Outback Brasil

Classificação	ALEGRIA	MEDO	TRISTEZA	DESGOSTO	RAIVA	SURPRESA
Correta	41	12	31	15	13	17
Errada	19	1	10	0	2	3

Fonte: O Autor.

Ao fim das buscas e classificação das três páginas foi realizada a comparação da classificação dos três grupos avaliados e exatificada as suas respectivas taxas de acerto. Após foi feita a média dos destas taxas de acerto, juntamente com as médias foi calculado o desvio padrão populacional entre as quantidades de cinquenta, cem e quinhentas postagens. Estas taxas médias de acerto e desvio são demonstradas na Tabela15.

Tabela 15 – Taxa média de acerto na classificação de postagens dos três grupos analisados

Classificação	ALEGRIA	MEDO	TRISTEZA	DESGOSTO	RAIVA	SURPRESA
Taxa de acerto	54,82%	90,77%	81,55%	88,89%	91,01%	75,83%
Desvio padrão(σ)	16,06	2,62	5,20	4,61	4,23	9,87

Fonte: O Autor.

Pode-se visualizar que quatro dos seis sentimentos classificados tiveram taxa de acerto a cima de oitenta por cento, um com taxa a cima de setenta e outro com taxa a cima de cinquenta por cento. Essa variação de percentuais pode estar relacionada às palavras que foram utilizadas em cada léxico. O sentimento de alegria possuía quase o dobro de palavras quando comparado aos demais e isso pode ter influenciado na baixa taxa de acerto, uma vez que as postagens eram classificadas como alegria por conterem palavras pertencentes ao léxico, mas que por sua vez eram postagens que não expressavam nenhum sinônimo de alegria conforme foi percebido na classificação manual. Por sua vez o desvio padrão reflete essa má classificação uma vez que devido ao seu valor consideravelmente alto apresenta que a classificação tem uma grande dispersão em torno da média populacional.

4 CONSIDERAÇÕES FINAIS

Com o intuito de poder demonstrar todo o conteúdo que foi trabalhado e apresentar resumidamente o conhecimento adquirido com a presente pesquisa no capítulo quatro são apresentadas algumas considerações que visam falar sobre pontos importantes visualizados no decorrer do desenvolvimento.

4.1 SÍNTESE DO TRABALHO

Neste trabalho foram apresentados métodos de coleta de dados em redes sociais juntamente com técnicas de extração de conhecimento a partir de dados textuais as quais auxiliaram no desenvolvimento da pesquisa e posteriormente da definição de qual método seria utilizado para que a união de todo esse conhecimento permitisse chegar nos resultados obtidos.

Posteriormente foi realizada a construção de um software que visa auxiliar o usuário a realizar pesquisas por conteúdo em postagens da rede social *Facebook*. No software ainda foi incluída a função de analisar os sentimentos contidos nas postagens e a função de identificar palavras de linguagem inapropriada e relacionadas a *phishing*.

Uma etapa de suma importância para a implementação do software foi a definição dos dados a serem analisados. A escolha da rede social foi um dos principais fatores que permitiu a implementação do software criado para coleta e análise das postagens. Posteriormente a definição de quais dados coletar dentro da rede permitiu que fosse identificado o tipos de relacionamento entre usuários que era buscado.

Os métodos de extração de conhecimento em texto incluindo o pré-processamento possibilitaram que as postagens coletadas fossem minimizadas o que permitiu que não fosse perdido conteúdo de análise, permitindo assim que o processamento de tarefas de análise fosse agilizado diminuindo consideravelmente o tempo de processamento para grandes quantidades de postagens. A análise de sentimentos utilizando um léxico foi outro ponto a ser exaltado, uma vez que a classificação foi baseada quase que totalmente em sentimentos extraídos do próprio léxico e a partir disso foram obtidos resultados de certa forma satisfatórios.

4.2 CONTRIBUIÇÕES DO TRABALHO

A pesquisa apresentada buscou esclarecer diversos aspectos de análise de conteúdo de texto em redes sociais e também utilizou de métodos de coleta de busca e coleta de dados que não são comumente explanados nos demais trabalhos em que foram realizadas pesquisas. O software construído buscou agregar conteúdo a pesquisa e também tornar intuitiva a busca e análise por parte de um usuário final, que por sua vez não tem conhecimentos de programação e não conseguiria desenvolver a busca e análise de postagens. A linguagem escolhida para a implementação do software foi de excepcional ajuda, uma vez que há muito conteúdo na

internet e diferentes tipos de bibliotecas que auxiliam na elaboração de diferentes aspectos no software.

Por fim o resultado obtido foi de certa forma satisfatório. Visto que a utilização dos temas abordados, são eles análise textual e de sentimentos, compreensão e aplicação de linguagem de programação, criação de um software que é intuitivo ao usuário e apresenta maneiras simplificadas de análise foi implementado e se encontra pronto para uso.

4.3 TRABALHOS FUTUROS

Para uma possível continuação do trabalho da pesquisa e desenvolvimento apresentado pode-se pensar na inclusão de um método complementar de classificação textual, talvez um corpus que utilize expressões de internet e não somente analise por termos. Outra melhoria que pode ser feita é a implementação do software criado em alguma linguagem web, para que futuramente possa ser alocado juntamente a rede social estudada e a partir daí realizar análises online a partir de contas do Facebook. Pode-se também pensar em aprimoramento de testes onde possa ser realizada a comparação com outros classificadores para que sirvam como parâmetro.

Todavia pode-se pensar em uma continuação do trabalho em que sejam utilizados os métodos apresentados no capítulo dois, que são classificação a partir do método *K-means*, também análise de classificação a partir de métodos como curva ROC e Entropia os quais não foram utilizados na implementação deste trabalho.

REFERÊNCIAS

- AMARAL, F. **Introdução à ciência de dados: mineração de dados e big data.** [S.l.]: Alta Books Editora, 2016.
- ARANHA, C. N.; VELLASCO, M. Uma abordagem de pré-processamento automático para mineração de textos em português: sob o enfoque da inteligência computacional. **Pontifícia Universidade Católica do Rio de Janeiro, Rio de Janeiro, RJ**, [S.l.], 2007.
- BARION, E. C. N.; LAGO, D. Mineração de textos. **Revista de Ciências Exatas e Tecnologia**, [S.l.], v. 3, n. 3, p. 123–140, 2015.
- CAIS. **Fraudes identificadas e divulgadas pelo cais.** Disponível em: <<https://memoria.rnp.br/cais/fraudes.php?>>. Acesso em: 02 junho 2017.
- CARVALHO FILHO, J. A. Mineração de textos: análise de sentimento utilizando tweets referentes à copa do mundo 2014. , [S.l.], 2014.
- CECI, F. **Um modelo semiautomático para a construção e manutenção de ontologias a partir de bases de documentos não estruturados.** 2010. 131 f. 2010. Tese (Doutorado em Ciência da Computação) — Dissertação (Mestrado)-Universidade Federal de Santa Catarina, Florianópolis, 2010.
- DEWAN, P.; KUMARAGURU, P. Detecting malicious content on facebook. **arXiv preprint arXiv:1501.00802**, [S.l.], 2015.
- FRANÇA, T. C. et al. Big social data: princípios sobre coleta, tratamento e análise de dados sociais. **XXIX Simpósio Brasileiro de Banco de Dados–SBBD**, [S.l.], v. 14, 2014.
- GOMES, H.; CASTRO NETO, M. de; HENRIQUES, R. Text mining: sentiment analysis on news classification. In: **INFORMATION SYSTEMS AND TECHNOLOGIES (CISTI), 2013 8TH IBERIAN CONFERENCE ON**, 2013. **Anais...** [S.l.: s.n.], 2013. p. 1–6.
- HEPP, F. S. Aprendizagem automática de phishing por mineração de dados. , [S.l.], 2010.
- HOLMES, T. E. **Credit card fraud and id theft statistics.** Disponível em: <<http://www.creditcards.com/credit-card-news/credit-card-security-id-theft-fraud-statistics-1276.php>>. Acesso em: 11 março 2017.
- KAUFMAN, L.; ROUSSEEUW, P. J. **Finding groups in data: an introduction to cluster analysis.** [S.l.]: John Wiley & Sons, 2009. v. 344.
- LONGHI, C. Análise de dados para apoio a tomada de decisão com ênfase no marketing digital. , [S.l.], 2017.
- LUHN, H. P. The automatic creation of literature abstracts. **IBM Journal of research and development**, [S.l.], v. 2, n. 2, p. 159–165, 1958.
- MANNING, C. D. et al. **Introduction to information retrieval.** [S.l.]: Cambridge university press Cambridge, 2008. v. 1, n. 1.

MARCACINI, R. M.; MOURA, M. F.; REZENDE, S. O. O uso da mineração de textos para extração e organização não supervisionada de conhecimento. **Revista de Sistemas de Informação da FSMA, Macaé**, [S.l.], v. 1, n. 7, p. 7–21, 2011.

MARQUES, T. P.; PINTO, A. M.; ALVAREZ, M.-J. Psychometric study of the scale of evaluation of risks and opportunities for young facebook users. **REVISTA IBEROAMERICANA DE DIAGNOSTICO Y EVALUACION-E AVALIACAO PSICOLOGICA**, [S.l.], v. 1, n. 41, p. 145–158, 2016.

OLIVEIRA, D. R. de; FERREIRA, D. F. M. Quem mexeu no meu face? uso e percepções de segurança no facebook, por crianças e adolescentes. **ABCiber Associação Brasileira de Pesquisadores em Cibercultura**, [S.l.], 2013.

STEINBACH, M. et al. A comparison of document clustering techniques. In: KDD WORKSHOP ON TEXT MINING, 2000. **Anais...** [S.l.: s.n.], 2000. v. 400, n. 1, p. 525–526.

STRAPPARAVA, C.; VALITUTTI, A. et al. Wordnet affect: an affective extension of wordnet. In: LREC, 2004. **Anais...** [S.l.: s.n.], 2004. v. 4, p. 1083–1086.

TAN, A.-H. et al. Text mining: the state of the art and the challenges. In: PAKDD 1999 WORKSHOP ON KNOWLEDGE DISCOVERY FROM ADVANCED DATABASES, 1999. **Proceedings...** [S.l.: s.n.], 1999. v. 8, p. 65–70.

WÜEST, C. The risks of social networking. **Symantec Corporation**, [S.l.], 2010.

XIAO, C.; FREEMAN, D. M.; HWA, T. Detecting clusters of fake accounts in online social networks. In: ACM WORKSHOP ON ARTIFICIAL INTELLIGENCE AND SECURITY, 8., 2015. **Proceedings...** [S.l.: s.n.], 2015. p. 91–101.

ZAFARANI, R.; ABBASI, M. A.; LIU, H. **Social media mining: an introduction**. [S.l.]: Cambridge University Press, 2014.

APÊNDICE A – TABELA COMPARATIVA DE TRABALHOS RELACIONADOS

Tabela 16 – Comparativo de trabalhos relacionados

Autor	Origem dos Dados	Extração	Algoritmos Classificadores	Objetivos	Resultados	Ano
Camila Longhi	Twitter, Facebook	Aplicação Java criada por Hofman, Netlytic	K-Means, Apriori	Relacionar grupos de alimentos com determinadas características de procura de usuários para que pudesse direcionar o marketing de uma empresa. Ao exemplo de bolo, pão e biscoito que se relacionaram com pesquisas da palavra sem glúten.	Dados coletados pela aplicação java puderam ser utilizados para análise no software Weka, os dados do Netlytics só puderam ser utilizados no próprio aplicativo. O algoritmo K-means se mostrou mais eficaz uma vez que conseguiu identificar relações entre grupos específicos e produtos que eram procurados dentro deste grupo. Já o algoritmo apriori não fez agrupamentos tão significantes.	2016
José Carvalho	Twitter	Script Phyton	Naive Bayes	Avaliar o sentimento nos tweets referente a copa do mundo de 2014 durante a ocorrência dos jogos da seleção brasileira de futebol. Com o objetivo de classifica-los como positivos, negativos, ambíguo e neutro.	O modelo de classificação apresentado retornou resultados mais satisfatórios nas categorias positiva onde teve 82% de precisão e negativa com 84%, já nas categorias neutra e ambígua não se mostrou tão eficiente obtendo apenas 69% e 40% respectivamente.	2014
Xiao, Cao and Freeman, David Mandell and Hwa, Theodore	LinkedIn	Dados fornecidos pela equipe do LinkedIn	Árvore de Decisão	Encontrar perfis falsos criados para disseminar spam pela rede social.	Os autores utilizaram de uma base de dados para treinamento de seu algoritmo. Utilizaram como dados de entrada o IP e a data em que as contas foram criadas. Conseguiram atingir um nível alto devido a supostamente os criadores das contas falsas utilizarem o mesmo IP para criarem várias contas em tempos próximos.	2015
Prateek Dewan, Ponnurangam Kumaraguru	Facebook	Facebook GRAPH API	Árvore de Decisão	Avaliar URLs que contenham conteúdos maliciosos em tempo real	O autor tratou do problema relacionado a postagens de links maliciosos em postagens relacionadas a acontecimentos que estão em alta no momento, ou seja, grandes acontecimentos como show, catástrofes entre outros. Utilizou do texto e dos dados contidos nas URLs. Após testado o algoritmo ainda criaram uma aplicação para navegadores a qual identifica URLs maliciosas e avisa ao usuário.	2015

Fonte: O Autor.

