

# **Análise de Sentimentos/Mineração de Opinião: Uma Revisão Bibliográfica**

Bárbara Maurosso Amaral<sup>1</sup>; Eliana Marinete dos Santos Silva<sup>2</sup>; Alex Marino Gonçalves de Almeida<sup>3</sup>

## **Resumo**

Este trabalho descreve uma análise de mineração de opinião realizada através de alguns trabalhos da área. O mesmo se trata de um referencial bibliográfico onde serão reunidas as técnicas e ferramentas utilizadas e os resultados obtidos, além de apresentar algumas abordagens sobre o tema em questão. A mineração de opinião/análise de sentimento, é uma área de estudos que busca analisar qual é a opinião ou sentimento sobre um determinado assunto. Essa área está se expandindo cada vez mais, principalmente pelo fato das pessoas expressarem suas opiniões através de plataformas como, redes sociais, fóruns, sites, entre outros. Com isso surge então a necessidade de aplicar técnicas e ferramentas que consigam coletar e identificar o conteúdo de opinião, e determinar o sentimento, percepção ou comportamento das pessoas em relação a um determinado alvo. Portanto este trabalho além de realizar essa análise proporciona uma direção para futuras pesquisas.

**Palavras-chaves:** Análise de sentimento, técnicas, ferramentas.

## **Abstract**

This paper describes an opinion mining analysis through some work in the area. The same it is a bibliographic reference where the techniques and tools used and the results obtained will be gathered, and present some approaches to the issue at hand. Mining opinion / sentiment analysis, is a field of study that seeks to analyze what is the opinion or feeling about a particular subject. This area is expanding more and more, mainly because of the people to express their opinions through platforms such as, social networks, forums, websites, among others. With that then arises the need to apply techniques and tools that are able to collect and identify the contents of opinion, and determine the feeling, perception or behavior of people in relation to a particular target. Therefore this work and perform this analysis provides a direction for future research

**Keywords:** Sentiment analysis, techniques, tools.

## **Introdução**

O constante crescimento do uso da tecnologia potencializa o uso de mídias a um nível em que o público em geral produz e consome grandes quantidades de informações. (SOUZA, 2012) diz que a esse novo tipo de mídia, em que muitos produzem para muitos, chama-se Mídia Social. As mídias sociais apresentam-se de várias formas: listas de discussão, blogosfera, sites de redes sociais (como Facebook, Twitter, Orkut, etc.), entre outros. Junto a esses fatores

---

1 Tecnóloga em Segurança da Informação pela Faculdade de Tecnologia de Ourinhos-FATEC. E-mail: barbaramamaral@gmail.

2 Tecnóloga em Segurança da Informação pela Faculdade de Tecnologia de Ourinhos-FATEC. E-mail: elianamari27@gmail.com.

3 Mestre em Ciência da Computação pela Universidade Estadual de Londrina-UEL, professor da Faculdade de Tecnologia de Ourinhos-FATEC. E-mail: alex.marino@fatecourinhos.edu.br.

surtem também uma enorme quantidade de informações relacionadas às pessoas como forma de expor suas opiniões. Com esse avanço também foram desenvolvidos métodos (ferramentas) para analisar essa grande quantidade de informações conhecidas como análise de sentimento/mineração de opinião. Com tantas mudanças visíveis, torna-se claro, que a maneira de se obter dados se expandiu imensamente. Análise de Sentimento (AS) visa identificar as opiniões postadas por usuários na internet, uma vez que a quantidade de comentários na web é muito grande, se faz necessário e muito útil à sumarização de vários comentários disponíveis em um simples resultado. Isso poupa tempo e ajuda os usuários a tomarem decisões com maior convicção (SANTOS, 2014) além disso, a análise de sentimento avalia, classifica as opiniões como positivas, negativas ou neutras e disponibiliza o resultado para o usuário (LIU, 2010). A análise de sentimento é um problema que está sendo cada vez mais desafiado, além de apresentar vários subproblemas, exige técnicas cada vez melhores que ensinam programas de computador a identificar opiniões e sentimentos sobre entidades de forma eficaz (RAVI; RAVI, 2015). Muitas empresas têm dado importância nesta forma de analisar essas opiniões, não só pelo fato de serem mais econômicas, mas também por atingir uma grande demanda de resultados sobre o nível de aceitação de um produto, serviço ou sobre a própria imagem da empresa, ou seja, são levantados os resultados sem se quer precisar de uma entrevista, facilitando na coleta de opiniões sobre determinado assunto. Com base nesses conceitos fica aparentemente claro que a AS é uma área que está sendo descoberta e muito pesquisada. Esse tema auxilia empresas que tem sua marca como principal ativo, a obterem a percepção das pessoas sobre seu produto ou serviço. A análise é trabalhosa, porém contribui claramente ao panorama competitivo das empresas, pois nela são contidos vários métodos com o propósito de obter resultados de forma eficaz. Portanto, neste trabalho será realizada uma pesquisa bibliográfica, consolidando os trabalhos mais recentes em mineração de opinião, onde serão levantadas as técnicas, ferramentas utilizadas para análise de sentimento, tendo como objetivo avaliar, expor os resultados obtidos pelos autores, mostrando que é possível e viável utilizar-se dessa técnica não só em ambiente corporativo, mas em ambiente de pesquisa, com o efeito de obter proveito das informações coletadas. Este referencial bibliográfico, denota à análise de sentimentos como uma técnica que possibilita alcançar grandes resultados, ou seja, com todo seu processo de análise, ela pode contribuir de forma evidente aos leitores que interessem pela a área. Este trabalho está dividido nas seguintes etapas: seção 1, revisão da literatura abrange a teoria sobre análise de sentimento, como *Big data*, análise inicial, pré-processamento, diferentes

níveis de análise, diferentes tipos de opiniões, sumarização e dificuldades na mineração de dados. seção 2, trabalhos relacionados abrange a parte mais prática, pois mostra as técnicas e ferramentas que alguns autores utilizaram, seção 3, metodologia utilizada, seção 4, resultados desses trabalhos e, por fim, seção 5, considerações finais.

## **1 Revisão da literatura**

A análise de sentimento ou mineração de opinião são abordagens que se consistem em identificar um estudo de opiniões, sentimentos, emoções e atitudes, onde é possível detectar, extrair, classificar opiniões, sentimentos e atitudes sobre diversos assuntos. Por ser um tema de pesquisa popular, porém complexo, pois além de abranger um tema de pesquisa Processamento de Linguagem Natural (PLN), abrange outros subproblemas como veremos mais tarde. Anteriormente ao ano 2000 era um tema pouco pesquisado, pelo fato de haver poucos textos disponíveis em formas digitais. Mas a partir do ano 2000 essa área tem se expandido rapidamente até os dias de hoje e ainda existe muito a ser pesquisado, ou seja, é uma área com grande potencial a ser explorados (LIU, 2012).

### **1.1 Big data**

*Big Data*, esse termo é aplicado em grandes quantidades de dados disponíveis na web. É um processo utilizado para coletar, analisar e organizar diferentes tipos de dados, esses dados podem ser estruturados ou não e com variações de tamanhos. Embora muitas pessoas ainda não tenham conhecimento sobre o termo "*big data*", a prática de manipulação de grandes quantidades de informações para subsequente análise é bem antiga. Devido a grande massa de dados hoje disponíveis, torna-se imprescindível a utilização de técnicas, para a realização e organização de todos esses dados, e é por esse motivo que as empresas veem se interessando por este conceito, pois traz consigo benefícios como, velocidade e eficiência, além de ser uma técnica onde os dados são aproveitados e usados para identificar informações úteis, visando novas oportunidades com ideias futuras e decisões estratégicas de negócio. É de conhecimento geral que, cada vez mais o volume de dados irá aumentar, *Big Data* promete grandes avanços no varejo, serviços financeiros, saúde, comunicações, manufatura e outras indústrias. Para aproveitar esses avanços, as organizações precisam de soluções de análise que possam processar e armazenar dados de forma contínua e muito mais rápido do que os sistemas de dados transacionais tradicionais (ORACLE, 2016).

## 1.2 Análise inicial

Primeiramente é necessário realizar a coleta dos dados, essa coleta de dados servirá como base textual, auxiliando no processo de mineração de dados. Os dados coletados podem ser obtidos de diversas fontes, esses dados serão armazenados e posteriormente serão analisados.

## 1.3 Pré-Processamento

O pré-processamento é uma etapa fundamental para mineração de dados, principalmente quando os dados não são estruturados, como os textos. Uma grande parcela dos textos encontrados nas redes sociais não são escritos de maneira correta, por isso se faz necessária a etapa de pré-processamento. Nesta etapa de pré-processamento é feita uma checagem nos dados, onde serão utilizados algoritmos para a padronização das palavras. Nessa checagem é feita uma limpeza, na qual, é eliminado qualquer tipo de acentuação, pontuação, caractere especial, números. São eliminados também, *stopwords* que são palavras que ocorrem no texto com certa frequência, mas que não influenciam no entendimento do mesmo, os erros ortográficos, *stemming* que são palavras no plural e no gerúndio, por exemplo. Essa etapa de reestruturação das palavras pode ser corrigida de maneira simples, utilizando um dicionário. A remoção de tais palavras, diminui a dimensionalidade do problema, reduzindo o tempo de resposta dos algoritmos de classificação e aumentando sua performance (SANTOS, 2014). Toda essa etapa de checagem tem como finalidade evitar o crescimento excessivo do grafo, por esse motivo só são adicionadas ao grafo os substantivos e os adjetivos, ou seja, só são adicionadas palavras em sua forma radical.

## 1.4 Diferentes níveis de análise

Existem diversos problemas de pesquisa que abrangem a análise de sentimento. Por ser uma área que envolve diversas abordagens, alguns desses problemas se baseiam no nível de granularidade da pesquisa factual. A sistemática proposta por (LIU, 2012) as dividem em três níveis: nível de documento, de sentença e de entidade e aspecto.

**Nível de documento:** O nível de documento tem como objetivo classificar se um documento de opinião expressa um sentimento positivo ou negativo. Esse nível de análise conceitua que o documento como um todo se refere a apenas uma entidade (um único produto, por exemplo). Dessa forma, só expressa opiniões para uma única entidade e não é válido para opiniões comparativas.

**Nível de sentença:** O nível de sentença analisa separadamente as sentenças de um documento e as classifica de forma individual, como positivo, negativo ou neutro. Esse nível está associado como classificação de subjetividade, que diferenciam frases que apresentam fatos concretos de frases que apresentam concepções subjetivas e opiniões. Por mais específico, ainda apresentam problemas. Os dois níveis citados acima não são suficientes para identificar o que as pessoas gostam ou não, isto é, esses níveis não conseguem detectar a real opinião das pessoas.

**Nível de entidade e aspecto:** Desse modo, o nível de análise entidade e aspecto realiza uma análise de maneira mais detalhada. Esse nível se difere dos outros, pois ao invés de basear em documentos, parágrafos ou sentenças, ele foca em direção a real opinião, ou seja, se baseia no sentido em que a opinião é formada por um sentimento positivo ou negativo. Pode-se considerar que a opinião é composta de um alvo e de um sentimento. Nesse nível, é necessário identificar a entidade alvo e todos os aspectos dela que estão sendo avaliados.

(PAK; PAROUBEK, 2010) utilizou um classificador *Naive-Bayes* para categorizar *tweets* em positivo ou negativo, com base em N-gramas e na classificação gramatical de partes do texto. Já NASCIMENTO (2011) fez uma comparação de três diferentes classificadores para verificar qual deles se adequa melhor. Eles optaram por utilizar os modelos N-gramas (esses modelos verificam a probabilidade de um agrupamento de palavras que aparecem em uma determinada sequência).

SILVA (2010) apresenta em seu trabalho a comparação de dois sistemas de classificação (*HowGood* e *BestChoice*) utilizando uma base de termos conhecida como *SentiWordNet*. Souza, (2011) utiliza a ferramenta *SentiWordNet* para classificar opiniões de usuários no Twitter. O trabalho foi desenvolvido na linguagem de programação Java com a utilização da *SentiWordNet* 3.0.

*SentiWordNet* é uma ferramenta léxica para a mineração de opinião. Essa ferramenta usa como base de dados o *WordNet*, um banco de dados que contém palavras da língua inglesa, os chamados *synset*. *SentiWordNet* atribui a cada *synset* da *WordNet* três classificações a respeito de sentimento: negatividade, positividade, objetividade (SOUZA, 2011).

### 1.5 Diferentes tipos de opiniões

Existem tipos diferentes de opiniões, mas a que foi abordada até agora é chamada de opinião regular. Um outro tipo de opinião chamada de opinião comparativa será posteriormente

discutido neste trabalho. É possível classificar as opiniões, de acordo como elas são apresentadas no texto, além disso essas opiniões podem ser classificadas como opiniões explícitas ou implícitas (LIU, 2012).

**OPINIÕES REGULARES E COMPARATIVAS.** Segundo (LIU, 2012) opinião regular muitas vezes é referida como uma opinião da literatura e se divide em dois subtipos:

- (i) Opinião direta: Ela se refere diretamente a uma entidade ou aspecto entidade.
- (ii) Opinião indireta: Ela expressa indiretamente uma opinião sobre uma entidade com base em seus efeitos sobre algumas entidades. Como exemplo, “Após a injeção da droga, minhas articulações estavam piores”, descreve um efeito indesejável da droga sobre “minhas articulações”, o que apresenta indiretamente uma opinião ou sentimento negativo à droga. No caso a entidade é a droga e o aspecto é o efeito sobre as articulações.

Já a comparativa apresenta uma relação de diferenças e semelhanças entre duas ou mais entidades ou baseando-se em alguns aspectos da entidade. Como exemplo, “O sabor da coca é melhor do que o sabor da Pepsi”. Existem muitos tipos de opiniões comparativas, que serão discutidas futuramente neste trabalho.

**OPINIÕES EXPLÍCITAS E IMPLÍCITAS.** De acordo com (LIU, 2012) opinião explícita é de fácil entendimento, ou seja, fica claro qual é a intenção do texto. Ela é uma declaração subjetiva, que dá uma opinião regular ou comparativa, por exemplo, "Eu adoro Coca"e "Coca é melhor do que Pepsi". Opinião implícita é uma indicação objetiva que implica uma opinião regular ou comparativa. Ela expressa algo desejável ou indesejável, como, por exemplo, “Comprei um carro semana passada e ele já quebrou” e “A vida útil de bateria do telefone da Samsung é maior do que o Iphone”.

## 1.6 Subjetividade e emoção

Há dois conceitos importantes que estão estreitamente relacionadas com sentimento e opinião, ou seja, a subjetividade e a emoção (LIU, 2012). Para o autor uma sentença objetiva traz informações sobre algo concreto, já uma sentença subjetiva denota sentimentos pessoais, opiniões, alegações, desejos, crenças, suspeitas e especulações. Dessa forma entende-se que são abordagens opostas.

Há grande confusão entre essas abordagens, pelo fato de terem uma relação, mas não serem equivalentes. Para definir se uma sentença é objetiva ou subjetiva, existe uma função chamada classificação de subjetividade. Essa função será mais abordada nas próximas seções.

De acordo com (PINTO, 2001), emoção pode ser definida como um sentimento subjetivo e com medidas mais complexas, pois abrange diversas categorias, cuja classificação não é

conclusiva, além de não se mensurar e não se concretizar. É uma reação complexa desencadeada por um estímulo ou pensamento e envolve reações orgânicas e sensações pessoais. Existem algumas formas que expressão a emoção, como exemplo, alteração de frequência cardíaca, pressão arterial, expressões faciais, postura, etc. A emoção é uma experiência subjetiva que envolve a pessoa toda, a mente e o corpo.

Esse tipo de sentimento tem sido estudado em vários campos, como, psicologia, filosofia e sociologia, e segundo cientistas as emoções são classificadas em categorias. De acordo (PARROTT, 2001), as pessoas têm seis emoções primárias, que são elas:

- Amor;
- Alegria;
- Surpresa;
- Raiva;
- Tristeza e
- Medo.

Essas emoções primárias são subdivididas em outras emoções, ou seja, as emoções podem estar vinculadas entre si, dependendo de sua intensidade, várias emoções podem estar interligadas. A forma como se expressa um sentimento ou opinião estão relacionados também com a intensidade das emoções.

Muito tem se abordado sobre particularidades de uma análise de sentimento, que na maioria das vezes são as avaliações. Segundo pesquisas, o comportamento do consumidor diante das avaliações pode ser classificados em dois tipos: avaliações racionais e avaliações de emocionais. (CHAUDHURI, 2006):

- (i) Avaliação racional: Essas avaliações são de pensamento racional, ou seja, se baseiam em fatos concretos, sem expor o sentimento sobre algo. Por exemplo, "Esse carro é novo".
- (ii) Avaliação Emocional: Estas avaliações são de respostas e são baseadas no sentimento e sua intensidade. Por exemplo, "Eu amo Pizza".

Para fazer uso destes dois tipos de avaliações na prática, nós podemos projetar 5 classificações de sentimento, negativo emocional (-2), negativo racional (-1), neutro (0), positivo racional (+1), e positivo emocional (+ 2). Na prática, neutro muitas vezes significa qualquer opinião ou sentimento expresso.

## 1.7 Sumarização

A sumarização é responsável por exibir ao usuário, de forma simples e clara, os resultados da classificação das opiniões. Subdivididas em duas etapas:

(ii) Sumarização de documento simples: Unir os textos que relatam opiniões parecidas ou que relatam um mesmo aspecto, para mostrar pontos positivos e negativos de cada aspecto (SILVA, 2010).

(ii) Sumarização de Multi-documento: Está relacionada ao resultado de vários objetos analisados (SILVA, 2010). Além de mostrar os pontos positivos e negativos, também é feita a comparação de cada aspecto.

No trabalho desenvolvido por (HU; LIU, 2004) é apresentado um sistema de sumarização automática de revisões de produtos e um mecanismo de mineração de características dos produtos, bem como a identificação das sentenças com opinião positiva e sentenças com opinião negativa.

## 1.8 Dificuldades da mineração de dados

Tendo a web como fonte de informações e comunicação para grande parte das pessoas, surge então à necessidade de saber qual é o parecer dessas pessoas sobre determinado assunto. Essas opiniões expostas por essas pessoas na web podem ser classificadas como positiva, negativa ou neutra. Devido a grande quantidade de opiniões expostas pelas pessoas, se torna indispensável à sumarização dos comentários obtidos, para que possam ser analisados de forma mais rápida. Textos que expõe opiniões são diferentes de outros tipos de textos, pois passam algum tipo de informação de forma mais direta. Uma das maiores dificuldades encontradas na mineração de dados é identificar muitas vezes a subjetividade. Muitas frases em um texto não possuem ou não expressam nenhuma opinião. No entanto, essa identificação de opiniões se torna algo muito complexo para uma máquina.

Podemos encontrar dificuldades em alguns textos, pois em geral, os textos não são estruturados, não tem uma formalidade constituída por sujeito e predicado. A maioria dos textos retirados da Internet apresentam erros gramaticais, sarcasmo, duplo sentido, gírias, abreviações, etc (SANTOS, 2014).

Outras dificuldades podem ser encontradas nas ferramentas utilizadas para realizar a mineração de dados. Por mais que a ferramenta apresente uma interface amigável, é necessário ter um suporte técnico da empresa que fornece a aplicação, não só para suprir as dificuldades que podem surgir na utilização, mas para orientar sobre as particularidades de tal aplicação

(DIAS, 2008). É imprescindível que uma ferramenta utilizada para mineração de dados seja manuseada por alguém que tenha certo conhecimento, conheça o fundamento de sua utilização para a organização, sobre as técnicas que serão utilizadas para a mineração de dados, além disso, o usuário deverá saber como analisar os resultados da coleta, para posteriormente utilizá-los de maneira correta.

## **2 Trabalhos relacionados**

Como citado anteriormente, existem várias técnicas para realizar uma análise de sentimento, alguns autores utilizam software já existente outros criam algoritmo na tentativa de aprimorar os resultados. Devido ao contexto geral da área de Análise de Sentimento, analisamos alguns trabalhos e separamos algumas técnicas, ferramentas e métodos utilizados pelos autores em seus trabalhos.

Devido a grande utilização das plataformas web, como meio de comunicação e exposição de opiniões, pode-se dizer que através dessas ferramentas podemos obter um leque de informações sobre um determinado assunto. As pessoas estão utilizando esses meios de comunicação para expor suas opiniões de maneira geral e por esse motivo surgem algumas preocupações para as organizações. Esse é um tema complexo, que por um lado tem grandes vantagens e, por outro, algumas desvantagens, não só pelo fato de ser um tarefa difícil identificar o que as pessoas escrevem, mas também por ser um método onde as informações se propagam rapidamente permitindo aos usuários tomarem conhecimento sobre a opinião e avaliação de um determinado assunto. (BECKER; TUMITAN, 2013)

Já que a maioria das informações publicadas nas redes sociais são em formas de textos, torna-se algo muito demorado se as pessoas fossem analisar cada uma dessas informações. É por esse e por outros motivos que a análise de sentimentos/mineração de opinião é uma técnica que melhora e agiliza esse processo. O trabalho Análise de Sentimentos em Mídias Sociais de (TORTELLA; COELLO, 2015) relata um estudo feito no período de agosto de 2014 a julho de 2015. De início o mesmo tinha o intuito de levantar trabalhos sobre a área e explorar os tópicos mais relevantes sobre o tema em questão. Mas por fim foi realizado um estudo mais detalhado, utilizando o sistema de análise de sentimentos Piegas, que foi desenvolvido utilizando a linguagem Ruby on Rails e utiliza um classificador Naive Bayes, podendo assim classificar as mensagens publicadas no Twitter em português, de acordo com os sentimentos contidos nessas mensagens. O sistema Piegas foi desenvolvido buscando facilitar a construção de aplicações Web. O programa utiliza-se de uma Application Programming (API) Interface do Twitter, onde

são buscadas as palavras digitadas pelo usuário. As palavras encontradas são filtradas, hiperlinks são descartados, pois podem ser spams. Os tuítes resultantes da filtragem são submetidos a análise feita pelo classificador Naive Bayes, que por sua vez, antes de realizar a classificação ele passa por uma etapa onde são reunidos tuítes classificados manualmente, ou seja, pelo fato do algoritmo Naive Bayes utilizar aprendizado supervisionado. Basicamente ao longo deste trabalho foram realizadas pesquisas sobre bibliografias a respeito do tema, além de implementar o sistema Piegas em conjunto com o Naive Bayes.

(FILHO, 2014) mostra como o processo de mineração de textos foi usado para coletar, estruturar o texto extraído do Twitter e como criar um modelo de classificação de texto que permita mapear a opinião dos usuários na rede social Twitter sobre Copa do Mundo da FIFA Brasil 2014 utilizando uma implementação do algoritmo de classificação Naive Bayes do Apache Mahout, o algoritmo analisa os *tweets* com base nas palavras e hashtags podendo ser classificados como positivos e negativos. O algoritmo Naive Bayes é um simples classificador probabilístico baseado na aplicação do teorema de Bayes, frequentemente utilizado como base na classificação de textos por ser rápido e fácil de implementar. O algoritmo bayesiano utilizado na etapa de Mineração de Textos, uma implementação pertencente ao Apache Mahout<sup>4</sup>. O Mahout é uma biblioteca de aprendizagem de máquina de código aberto do Apache. (OWEN SEAN; ANIL, 2011)

(SANTOS, 2014) também utiliza, em seu trabalho, o o algoritmo mahout juntamente com Hadoop<sup>5</sup>. O Hadoop é um framework de código aberto para processamento em larga escala de grandes quantidades de dados em clusters de computadores de baixo de custo. O projeto é desenvolvido e mantido pela fundação Apache (APACHE, 2016). O objetivo de sua pesquisa foi desenvolver uma aplicação de mineração de opinião capaz de avaliar a eficácia do pré-processamento de textos em uma base de textos do português formada por gírias, abreviações e termos comumente utilizados na Internet. Foram extraídos 759 mil comentários em português da loja de aplicativos Google Play. Antes de utilizar os algoritmos citados, ele utilizou a ferramenta Weka<sup>6</sup> que não foi capaz de trabalhar com modelos de mineração de dados que envolvessem a quantidade total de documentos extraídos da Google Play. Possibilitou a análise completa do corpus de forma à validar as hipóteses feitas na primeira análise e também de forma à iniciar estudos e experimentos em uma recente área de pesquisa: *Big Data*.

---

4 <https://mahout.apache.org/users/classification/bayesian.html>.

5 <http://hadoop.apache.org/>.

6 <http://www.cs.waikato.ac.nz/ml/weka/>.

(SILVA; LIMA; BARROS, 2012) utilizaram um protótipo, o Sentiment Analysis using Pairs (SAPair) um sistema completo para Análise de Sentimentos em nível de característica, codificado em Java. O trabalho tem como foco a AS no nível de característica, contando com duas etapas centrais: Extração de características, que busca identificar as características sobre as quais o texto trata; e Classificação, que atribui uma polaridade (positiva ou negativa) a cada característica com base nas palavras opinativas relacionadas a ela. A etapa de Extração utilizou um corpus de 500 comentários em inglês sobre aparelhos celulares, e a etapa de classificação foi testada com um corpus de 5.500 opiniões no mesmo domínio e os resultados foram comparados com os sistemas SentiworldNet e Turney.

(SANTOS; BECKER; MOREIRA, 2014) O propósito deste trabalho é coletar e analisar as opiniões dos textos, independentemente do idioma em que os textos estão escritos. No mesmo apresenta uma abordagem distinta de outros trabalhos que focam apenas na classificação da polaridade do sentimento, destacando como foco principal a análise das emoções que existem nos textos.

Com isso, este trabalho classifica as emoções baseando-se em dicionário NRC (*word-emotion association*) e tradução automática. Além disso, este estudo procura analisar se é mais viável traduzir o texto das revisões ou das palavras do dicionário, avaliando também, se a utilização de um lematizador melhora os resultados, que nada mais é do que uma técnica utilizada por buscador de palavras em sites, ou seja, quando é feita uma busca por uma determinada palavra em um site, ele retorna como resultado as variações da palavra.

O NRC *Emoticon Lexicon* (MOHAMMAD; TURNEY, 2010) e um método léxico que classifica textos em 8 categorias afetivas, definidas por (PLUTCHIK, 1980), são elas: alegria, tristeza, raiva, medo, confiança, desgosto, antecipação e surpresa. O método foi desenvolvido utilizando uma base de dados chamada pelos autores de EmoLex. O EmoLex consiste de palavras rotuladas no serviço Amazon Mechanical Turk service<sup>7</sup> (GOMES, 2013) utiliza o software SAS 5 (SAS Information Retrieval Studio versão 1.3 e o SAS Web Crawler versão 2.1) objetivo do trabalho consiste na construção de um modelo capaz de avaliar a polaridade (positiva, negativa ou neutra) de títulos de notícias de economia, disponíveis em endereços de RSS Feeds. Ele utiliza<sup>8</sup> tipos de modelos, o modelo estatístico, Modelo baseado em regras (BeR) e o modelo híbrido (estatístico + BeR).

---

7 <https://www.mturk.com/mturk/welcome>.

8 <http://www.sas.com/ptbr/home.html>.

A análise de sentimentos vem se tornando um tema muito abordado, não só pelo fato de grande parte das pessoas utilizarem as redes sociais, mas sim pela maioria das pessoas exporem de maneira pública suas opiniões sobre determinado assunto, objeto ou serviço. (SILVA R. L. A. ROCHA, 2016) busca aprofundar a pesquisa sobre análise de sentimentos, e propor um modelo prático da Roda das Emoções de Plutick juntamente com a utilização de Árvores de Decisão Adaptativa para classificação automática de sentimentos, de menções extraídas de textos da rede social Twitter. No mesmo foi realizada a classificação dos sentimentos onde foram selecionadas as partes mais importantes do objeto em questão, desconsiderando os elementos irrelevantes. Para isso foi utilizado a técnica Roda de Emoções de Plutchik onde foram descritos recursos computacionais que classificam os sentimentos e as emoções. No trabalho são propostos os modelos de ações adaptativas, sendo o primeiro o MFA – Modelo de Filtragem Adaptativa e o MCS – Modelo de Classificação de Sentimentos.

(SANTOS et al., 2010) destina-se em realizar uma análise aos usuários do twitter, tentando mostrar que as redes sociais podem ser um contato direto com os clientes de empresas. No mesmo foi utilizado como método *Support Vector Machine* (SVM), ou seja, método supervisionado de aprendizagem de máquina, onde busca o quanto os usuários expressam suas opiniões sobre o Windows 7, produto da Microsoft. No trabalho foi realizada a classificação das mensagens postadas. As mensagens eram neutras e outras continham opiniões relacionadas ao produto em questão. Para isso foi utilizada a seguinte implementação SVMlight<sup>9</sup>. Para as coletas das mensagens foi utilizada a Twitter *Streaming* API, mas como a Twitter *Streaming* API não possibilita a exata filtragem por palavras-chave com espaços, além de não filtrar por idiomas, o resultado da coleta acabou sendo insatisfatório. Posteriormente foi utilizada *Application Programming Interface* (API) em tradução para o português "Interface de Programação de Aplicativos" do site Lang ID<sup>10</sup> que se baseia na Google Ajax API, que tem como finalidade identificar o idioma de cada mensagem. Depois de definida a base de treinamento foram removidas pontuações, links de páginas e *stopwords* que são palavras que ocorrem com frequência nas mensagens. Após todas as procedimentos e avaliações, obteve-se como alguns resultados uma distinção eficiente entre mensagens neutras e mensagens com sentimento.

Em seu trabalho, (MOREIRA et al., 2016), teve como objetivo principal a realização de uma análise dos dados extraídos de uma plataforma de rede social online. Os dados extraídos

---

9 <http://svmlight.joachims>.

10 <http://langid.net>.

foram analisados de acordo com duas diferentes abordagens: automática – ferramentas de análise de sentimentos (SentiStrength<sup>11</sup>, LIWC<sup>12</sup> 9, SenticNet<sup>13</sup>) e manual – análise por um dos pesquisadores. Nestas abordagens, os comentários foram analisados de acordo com sua polaridade e sentimentos expressados. A partir destes resultados, foi estabelecida uma análise comparativa com os resultados da análise manual dos autores das mensagens com o objetivo de identificar a quantidade de acertos de cada abordagem, bem como sugerir possíveis motivos para as falhas de avaliação.

A análise de sentimentos é uma área que vem crescendo constantemente devido a grande parte das pessoas utilizarem a internet para exporem o que pensam ou sentem. Essa é uma área dentro de Processamento de Linguagem Natural que está sendo pesquisada e que busca identificar qual o sentimento e opiniões das partes subjetivas de um texto. (BHADANE; DALAL; DOSHI, 2015) analisa as técnicas utilizadas para classificar esses elementos em um texto, e verifica se a opinião geral do texto é positiva ou negativa. Também são abordados dois pontos interessantes, classificação de aspecto seguido de classificação de polaridade. As duas principais áreas de investigação na classificação de sentimento são abordagens de aprendizagem lexicais e de máquina e o software LibSVM<sup>14</sup>. Primeiramente para inserir uma abordagem lexical é necessário um dicionário para armazenar os valores que correspondem a polaridade. Para seguir com essa abordagem, foi realizado um cálculo, ou seja, cada palavra do texto é analisada desde que esteja no dicionário, assim ela obtém uma pontuação de polaridade. Dessa forma se um léxico representa a uma palavra marcada como positiva no dicionário, acumula-se a pontuação de sua polaridade. Portanto se a polaridade como um todo no texto é positiva, o texto predomina-se como positivo, se não, predomina-se como negativo. Nessa abordagem são apresentadas alguns de suas variações que são: Abordagem de linha de base, Decorrentes, Parte da fala Tagging, WordNet, N-grams, Regras de conjunção, Pares de Palavras, Métodos de negação. Já na abordagem de aprendizagem de máquina, alguns parâmetros são selecionados e um conjunto com etiquetas é usado para organizar um modelo. Esse modelo, pode classificar esse corpus. Geralmente, uma variedade de unigramas ou N-gramas são escolhidas como parâmetros de classificação. Nessa abordagem é essencial que seja realizada a seleção de características para um bom resultado da classificação. Alguns métodos para

---

11 <http://sentistrength.wlv.ac.uk>.

12 <http://liwc.wengine>.

13 <http://liwc.wengine>.

14 <https://www.csie.ntu.edu.tw/~cjlin/libsvm/>.

realizar a classificação aplicado em aprendizagem de máquina são *Support Vector Machine*, Naive Bayes.

(HUTTO; GILBERT, 2014) apresenta a ferramenta VADER (*Valence Aware Dictionary for sEntiment Reasoning*), um modelo simples baseado em regras para utilização geral em análise de sentimento. Em seu trabalho foi realizada a comparação da sua eficácia com outros softwares: LIWC, ANEW, o *General Inquirer*, SentiWordNet, e técnicas orientadas de Naive Bayes, Máxima Entropia, e *Support Vector Machine* (SVM). Usando o modelo baseado em regras parcimoniosas para avaliar o sentimento de *tweets*. VADER é uma ferramenta utilizada na análise de sentimento validada por humanos e foi desenvolvida com foco em avaliar mensagens no contexto do Twitter e outras mídias sociais.

### 3 Materiais e métodos

Neste trabalho utilizamos o método de pesquisa bibliográfica. A pesquisa bibliográfica procura explicar e discutir um tema com base em referências teóricas publicadas em livros, revistas, periódicos e outros. Busca também, conhecer e analisar conteúdos científicos sobre determinado tema. (MARTINS, 2001).

Para atingirmos o objetivo fizemos um levantamento bibliográfico de alguns materiais como, artigos, trabalhos, dissertações, entre outros, nacionais e internacionais, a partir do ano 2010 sobre o tema análise de sentimento/mineração de opinião. Desses artigos selecionados foram feitos resumos e a partir desses resumos serão discutidos os tipos de sistemas utilizados, técnicas e o resultado final obtido pelos autores.

Os documentos bibliográficos foram retirados dos sites ACM Digital Library<sup>15</sup>, ScienceDirect<sup>16</sup> e Google Acadêmico<sup>17</sup>.

### 4 Análise e discussão dos resultados

A Tabela 1 foi elaborada para melhor compreensão dos resultados, onde pode-se visualizar autor, técnica utilizada, software e o resultado obtido.

---

15 <http://dl.acm.org/>.

16 <http://WWW.sciencedirect.com/science/journals>.

17 <https://scholar.google.com.br/>.

**Tabela 1** - Tabela de resultados

AUTORES	TÉCNICAS	SISTEMAS	RESULTADOS
TORTELLA; COELLO, 2015	Léxico	Piegas	Mediano
FILHO, 2014	PNL	Apache Mahout	Positivo
SANTOS, 2014	SVM e StopWords	Hadoop	Positivo
SILVA; LIMA; BARROS, 2012	Léxico	SAPair	Positivo
SANTOS; BECKER; MOREIRA, 2014	Léxico	NRC Emoticon Lexicon	Positivo
GOMES, 2013	Modelo BeR	SAS	Positivo
SILVA; ROCHA, 2016	Roda de emoções	MFA e MCS	Positivo
SANTOS et al., 2010	SVM	SVMLight e Lang Id	Positivo
MOREIRA et al., 2016	Léxico	SentiStrength, LIWC e SenticNet	Mediano
BHADANE; DALAL; DOSHI, 2015	Léxico e SVM	LibSVM	Positivo
HUTTO; GILBERT, 2014	Léxico	VADER	Positivo

Fonte: Elaborada pelos autores.

(TORTELLA; COELLO, 2015): O trabalho foi considerado mediano pelo autor, pois não houve tempo para a implementação do classificador SVM para determinar o sentimento dos tweets, e comparar o seu desempenho com o do classificador NB. Não foi comparado um sistema Piegas com o Naive Bayes e obteve um bom desempenho.

(FILHO, 2014): O resultado deste trabalho foi positivo. Foi coletado mais de 2 milhões de *tweets* distintos. Após a coleta, foram removidos *tweets* não relevantes para fase de mineração, em seguida foi gerado um modelo utilizando conjunto de treino com 2634 *tweets* que apresentou uma ótima taxa de precisão 88,91% sendo positivo (94%), negativo (93%) e neutra (84%).

(SANTOS, 2014): Após aplicadas as fases de pré-processamento e construídos os modelos de classificação, foram testados os softwares e o resultado mostra que a ferramenta Weka não suportou a grande quantidade de comentários, por isso foi utilizado o Apache Mahout e Hadoop tendo uma precisão de 81,29%. Para o autor, em vista de outros trabalhos, obteve-se um resultado positivo.

(SILVA; LIMA; BARROS, 2012): Como resultado final do experimento, obtiveram um resultado positivo com precisão global, são eles: SentiWordNet 77%, Turney 88%, e SAPair 90%. O SAPair mostrou-se superior aos outros métodos. Depois calculando-se a precisão da classificação das opiniões no nível de documento, SAPair calculando-se a precisão de 83,12%.

(BECKER; TUMITAN, 2013): Esse estudo de caso, obteve um resultado positivo, a tradução do texto da revisão é uma abordagem que produz resultados melhores do que a tradução do dicionário usado. O uso de lemas não produz melhorias estatísticas nos resultados.

Alguns erros de tradução fizeram com que os classificadores não identificassem corretamente algumas emoções e os resultados mostraram que é possível, necessitando de mais experimentos incluindo outras técnicas para melhorá-los.

(GOMES, 2013): O modelo estatístico conseguiu classificar corretamente cerca de 96% negativo e 66% positivo. Modelo baseado em regras (BeR) consegue classificar corretamente cerca de 91% das notícias positivas, 85% das notícias negativas e apenas 38% das neutras. O modelo híbrido (estatístico + BeR) ofereceu os melhores resultados, classificação das negativas, com 94% de acerto, nas positivas, com cerca de 80%. Tendo um resultado positivo.

(SILVA R. L. A. ROCHA, 2016): O uso da camada de persistência pode se mostrar eficaz no MFA, comunicando-se diretamente com a Máscara que por sua vez, retroalimenta o sistema e as gravações nas tabelas de decisões. Tendo os sentimentos já separados pelo MFA, e trabalhando em conjunto com o MCS pode-se perceber a eficiência dos modelos em funcionamento paralelo. O uso de recursos adaptativos no MCS pode se mostrar eficaz e coerente para classificação de sentimentos fundamentada nas oito emoções primárias da Roda das Emoções.

(SANTOS et al., 2010): Analisando os dados coletados, não foi possível realizar a classificação entre positivos, negativos, após a separação das mensagens neutras. Entretanto, o propósito deste trabalho era somente mensurar o quanto os usuários expressam suas opiniões sobre o produto em questão, utilizando a rede Twitter foi obtido esse respaldo apenas com a eliminação das mensagens neutras. Resultado do teste obteve uma precisão de 80%, considerado positivo.

(MOREIRA et al., 2016): A comparação realizada indicou diferenças entre os resultados obtidos e o real, enquanto as ferramentas utilizadas classificaram mais de 40% dos comentários como neutros, a análise dos autores das mensagens indicou que 71% dos comentários eram positivos. Na classificação foram identificadas divergências na categorização da intensidade dos sentimentos por cada método. Nenhuma abordagem, incluindo a realizada pelo pesquisador, foi considerada suficientemente eficiente, uma vez que o maior nível de acurácia obtido foi inferior a 70%.

(BHADANE; DALAL; DOSHI, 2015): Baseando se nos conceitos, a implementação dessas técnicas para a classificação de aspecto e reconheci associados a léxicos específicos, pode-se obter resultados satisfatórios, com certa de 78% de precisão, ou seja, com efeitos bem-sucedidos no desempenho de suas funções.

(HUTTO; GILBERT, 2014): VADER mostrou-se um software de rápida precisão, o léxico e as regras utilizadas são diretamente acessíveis, não requer um conjunto extenso de dados e mostrou ótimos resultados superando a avaliação individual humana. Ele foi comparado com outros softwares de classificação por léxico. Tendo uma precisão de 0,96, mostrando resultado bastante satisfatório, considerado pelo autor, um grande avanço na ciência da computação que em geral é mais favorável do que qualquer um dos softwares.

Análise de sentimento é uma grande área que está se expandindo, e para o seu crescimento se faz necessário a criação de novas técnicas e softwares que possam auxiliar em pesquisas cada vez mais satisfatórias. Notamos que a maioria dos trabalhos pesquisado focam em mídias sociais e não é por menos, pois atualmente, os usuários criaram um hábito de postarem suas opiniões referente a diversos assuntos. Na nossa pesquisa vimos que existem diversas técnicas e softwares que podem ser utilizados. Os trabalhos coletados mostram que alguns autores optaram por verificar o desempenho de um software, fazer uma comparação entre os sistemas disponíveis, criar o próprio algoritmo e outros utilizaram mais de uma técnica em busca de melhores resultados. Nota-se também que alguns utilizam a mesma técnica, mas com sistemas diferentes, o que dificulta fazer uma comparação direta, mas possibilitou ampliar a visão a diversos métodos utilizados. De uma forma geral, os trabalhos, obtiveram resultados positivos e sempre deixando sugestões para que futuros pesquisadores melhorarem o que eles começaram.

## **5 Considerações finais**

Baseando-se nestes conceitos, este trabalho foi elaborado com perspectiva a visualizar os resultados obtidos dos trabalhos coletados, ou seja, expor as técnicas e ferramentas utilizadas, servindo como um referencial bibliográfico.

Através desses trabalhos coletados, pode-se analisar tanto a teoria quanto a prática na área de análise de sentimento/mineração de opinião. E com estes resultados, pode-se reunir os dados e apresentá-los em uma tabela, servindo como norteamento a quem se interesse pela área.

Portanto conclui-se que com o grande crescimento de conteúdos na web fica mais viável utilizar a análise de sentimento para realizar uma pesquisa de opiniões sem precisar, necessariamente, que as pessoas respondam um questionário formal. Mas essa facilidade traz algumas dificuldades em analisar esta grande quantidade de dados, e é por esse motivo que o aperfeiçoamento e o desenvolvimento de novas técnicas se fazem necessários.

## Referências

- APACHE. **Hadoop**. 2016. <http://hadoop.apache.org/>. Acessado em 11 de outubro de 2016.
- BECKER, K.; TUMITAN, D. Introdução à mineração de opiniões: Conceitos, aplicações e desafios. **Simpósio Brasileiro de Banco de Dados**, 2013.
- BHADANE, C.; DALAL, H.; DOSHI, H. Sentiment analysis: Measuring opinions. **Procedia Computer Science**, Elsevier, v. 45, p. 808–814, 2015.
- CHAUDHURI, A. **Emotion and reason in consumer behavior**. [S.l.]: Routledge, 2006.
- DIAS, M. M. Parâmetros na escolha de técnicas e ferramentas de mineração de dados. **Acta Scientiarum. Technology**, v. 24, p. 1715–1725, 2008.
- FILHO, J. A. C. **MINERAÇÃO DE TEXTOS: ANÁLISE DE SENTIMENTO UTILIZANDO TWEETS REFERENTES À COPA DO MUNDO 2014**. Tese (Doutorado) — Trabalho de Conclusão de Curso, Universidade Federal do Ceará, 2014.
- GOMES, H. J. C. **Text Mining: análise de sentimentos na classificação de notícias**. Tese (Doutorado), 2013.
- HU, M.; LIU, B. Mining and summarizing customer reviews. In: **ACM. Proceedings of the tenth ACM SIGKDD international conference on Knowledge discovery and data mining**. [S.l.], 2004. p. 168–177.
- HUTTO, C. J.; GILBERT, E. V. A parsimonious rule-based model for sentiment analysis of social media text. In: **Eighth International AAI Conference on Weblogs and Social Media**. [S.l.: s.n.], 2014.
- LIU, B. **Sentiment analysis and subjectivity**. **Handbook of natural language processing**, v. 2, p. 627–666, 2010.
- LIU, B. Sentiment analysis and opinion mining. **Synthesis lectures on human language technologies, Morgan & Claypool Publishers**, v. 5, n. 1, p. 1–167, 2012.
- MARTINS, G.A; PINTO. R. **Manual para elaboração de trabalhos acadêmicos**. São Paulo: Atlas, 2001.
- MOHAMMAD, S. M.; TURNEY, P. D. Emotions evoked by common words and phrases: Using mechanical turk to create an emotion lexicon. In: **Association for Computational Linguistics. Proceedings of the NAACL HLT 2010 workshop on computational approaches to analysis and generation of emotion in text**. [S.l.], 2010. p. 26–34.
- MOREIRA, V. S. et al. Análise de sentimentos: Comparando o uso de ferramentas e a análise humana. **XII Brazilian Symposium on Information Systems**. Florianópolis, 17-20 mai. 2016.
- NASCIMENTO, P. **Análise de sentimento de tweet com foco em notícias**. Programa de Engenharia de Sistemas e Computação Instituto Alberto Luiz Coimbra de Pós-Graduação e Pesquisa de Engenharia Universidade Federal do Rio de Janeiro, 2011.

ORACLE. **Big Data**. 2016. <https://www.oracle.com/big-data/index.html>. Acessado em 03 de outubro de 2016.

OWEN, S.; ANIL, R.; DUNNING, T.; Friedman, E. **Mahout in Action**. Connecticut: Manning Publications Co, 2011.

PAK, A.; PAROUBEK, P. Twitter as a corpus for sentiment analysis and opinion mining. In: **LREc**. 2010.

PARROTT, W. G. **Emotions in Social Psychology: Essential Readings**. [S.l.]: Psychology Press, 2001.

PINTO, A. C. **Psicologia Geral**. [S.l.]: Universidade Aberta, 2001. 344 p.

PLUTCHIK, R. A general psychoevolutionary theory of emotion. **Theories of emotion**, Academic Press New York, v. 1, 1980, p. 3–31.

RAVI, K.; RAVI, V. A survey on opinion mining and sentiment analysis: tasks, approaches and applications. **Knowledge-Based Systems**, Elsevier, v. 89, 2015, p. 14–46.

SANTOS, A. G. L.; BECKER, K.; MOREIRA, V. Um estudo de caso de mineração de emoções em textos multilíngues. Instituto de Informática, Universidade Federal do Rio Grande do Sul-UFRGS, 2014. Disponível em: <<http://www.inf.ufrgs.br/~aglsantos/publicacoes/AlineLermen-BraSNAM2014.pdf>>. Acesso em: 11 ago. 2017.

SANTOS, F. L. d. **Mineração de opinião em textos opinativos utilizando algoritmos de classificação**. Monografia (Graduação) - Trabalho de Conclusão de Curso, Universidade Brasília, 2013. Ceará, 2014.

SANTOS, L. M. et al. Twitter, análise de sentimento e desenvolvimento de produtos: Quanto os usuários estão expressando suas opiniões? **Revista PRISMA.COM**, n. 13, 2010. p. 1-1. Disponível em: <http://revistas.ua.pt/index.php/prismacom/article/view/790/722>>. Acesso em: 11 ago. 2017.

SENTWORDNET. <http://sentiwordnet.isti.cnr.it/>. Acessado em 04/05/2016. Citado na página 5.

SILVA, N. G. R. **BestChoice**: Classificação de Sentimento em Ferramentas de Expressão de Opinião. Tese (Doutorado) — Tese de graduação, Universidade Federal de Pernambuco, Recife, 2010.

SILVA, N. R.; LIMA, D.; BARROS, F. Sapair: Um processo de análise de sentimento no nível de característica. In: **IV International Workshop on Web and Text Intelligence (WTI-2012)**. 2012. Disponível em: <<http://www.labic.icmc.usp.br/wti2012/artigos/105283.pdf>>. Acesso em: 11 ago. 2017.

SILVA R. L. A. ROCHA, J. J. N. A. M. Análise semântica de sentimentos utilizando árvores de decisão adaptativas. **X Workshop de Tecnologia Adaptativa-WTA**, p. 37-45.

SOUZA, L. V. **Análise de sentimentos no twitter utilizando sentiwordnet**. Trabalho de graduação – Centro de Informática da Universidade Federal de Pernambuco, 2011.

SOUZA, M. V. d. S. **Mineração de opiniões aplicada a mídias sociais**. Dissertação (Mestrado) – Ciências da Computação, Pontifícia Universidade Católica do Rio Grande do Sul, 2012.

TORTELLA, P. L.; COELLO, J. M. A. **Análise de sentimentos em mídias sociais**. Laboratório de Banco de Dados - Departamento de Ciência da Computação da Universidade Federal de Minas Gerais-UFMG, 2015. Disponível em: <http://www.lbd.dcc.ufmg.br/colecoes/sbsi/2013/0047.pdf>>. Acesso em: 11 mar. 2017.