



ANÁLISE DE SENTIMENTO DE REDES SOCIAIS POR CLASSIFICADORES MULTIMODAIS DE APRENDIZADO DE MÁQUINA

Breno Vieira Arosa

Dissertação de Mestrado apresentada ao Programa de Pós-graduação em Engenharia Elétrica, COPPE, da Universidade Federal do Rio de Janeiro, como parte dos requisitos necessários à obtenção do título de Mestre em Engenharia Elétrica.

Orientadores: Luiz Pereira Calôba
Natanael Nunes de Moura
Junior

Rio de Janeiro
Setembro de 2019

ANÁLISE DE SENTIMENTO DE REDES SOCIAIS POR CLASSIFICADORES
MULTIMODAIS DE APRENDIZADO DE MÁQUINA

Breno Vieira Arosa

DISSERTAÇÃO SUBMETIDA AO CORPO DOCENTE DO INSTITUTO
ALBERTO LUIZ COIMBRA DE PÓS-GRADUAÇÃO E PESQUISA DE
ENGENHARIA (COPPE) DA UNIVERSIDADE FEDERAL DO RIO DE
JANEIRO COMO PARTE DOS REQUISITOS NECESSÁRIOS PARA A
OBTENÇÃO DO GRAU DE MESTRE EM CIÊNCIAS EM ENGENHARIA
ELÉTRICA.

Examinada por:

Prof. Luiz Pereira Calôba, Dr.Ing.

Natanael Nunes de Moura Junior, D.Sc.

Prof. Nome do Terceiro Examinador Sobrenome, D.Sc.

Prof. Nome do Quarto Examinador Sobrenome, Ph.D.

Prof. Nome do Quinto Examinador Sobrenome, Ph.D.

RIO DE JANEIRO, RJ – BRASIL
SETEMBRO DE 2019

Vieira Arosa, Breno

Análise de Sentimento de Redes Sociais por Classificadores Multimodais de Aprendizado de Máquina/Breno Vieira Arosa. – Rio de Janeiro: UFRJ/COPPE, 2019.

IX, 15 p.: il.; 29, 7cm.

Orientadores: Luiz Pereira Calôba

Natanael Nunes de Moura Junior

Dissertação (mestrado) – UFRJ/COPPE/Programa de Engenharia Elétrica, 2019.

Referências Bibliográficas: p. 14 – 15.

1. Análise de sentimento. 2. Processamento de linguagem natural. 3. Redes complexas. 4. Aprendizado de máquina. I. Pereira Calôba, Luiz *et al.* II. Universidade Federal do Rio de Janeiro, COPPE, Programa de Engenharia Elétrica. III. Título.

*A alguém cujo valor é digno
desta dedicatória.*

Resumo da Dissertação apresentada à COPPE/UFRJ como parte dos requisitos necessários para a obtenção do grau de Mestre em Ciências (M.Sc.)

ANÁLISE DE SENTIMENTO DE REDES SOCIAIS POR CLASSIFICADORES MULTIMODAIS DE APRENDIZADO DE MÁQUINA

Breno Vieira Arosa

Setembro/2019

Orientadores: Luiz Pereira Calôba
Natanael Nunes de Moura Junior

Programa: Engenharia Elétrica

Nos últimos anos, as redes sociais se tornaram um dos principais meios de comunicação e com isso, houve um aumento da influência que exercem sobre os usuários. Por esse motivo, as mensagens que trafegam por elas passam a ter importância para as mais diversas finalidades como, por exemplo, a avaliação de produtos e eventos. Dentre as possíveis análises, a mineração de opinião é uma das operações com mais aplicações diretas. Nesse sentido, ferramentas de processamento de linguagem natural e de redes complexas são capazes de auxiliar a geração destas análises. Entretanto, o desempenho dessas técnicas, em geral, depende da existência e do volume de bases de treinamento anotadas manualmente, dificultando assim a utilização das mesmas. O presente trabalho aplica métodos de geração de bases de treinamento automatizadas para contornar esse obstáculo e gerar classificadores de análise de sentimento. São avaliados diferentes modelos de classificação textual, tanto lineares, como Naïve Bayes e SVM, quanto por *Deep Learning*, como redes convolucionais e redes recorrentes. Também são analisadas técnicas de redes complexas para caracterização dos usuários das redes, abordando assim diferentes aspectos das informações fornecidas por estas mídias.

Abstract of Dissertation presented to COPPE/UFRJ as a partial fulfillment of the requirements for the degree of Master of Science (M.Sc.)

SENTIMENT ANALYSIS OF SOCIAL NETWORKS BY MULTIMODAL
MACHINE LEARNING CLASSIFIERS

Breno Vieira Arosa

September/2019

Advisors: Luiz Pereira Calôba

Natanael Nunes de Moura Junior

Department: Electrical Engineering

In this work, we present ...

Sumário

Lista de Figuras	viii
Lista de Tabelas	ix
1 Introdução	1
1.1 Motivação	3
1.2 Objetivo	4
1.3 Organização do Texto	5
2 Análise de Sentimento de Redes Sociais	6
3 Classificadores por Processamento de Linguagem Natural	9
4 Modelos de Redes Complexas	10
5 Método	11
6 Resultados e Discussões	12
7 Conclusões	13
Referências Bibliográficas	14

Lista de Figuras

Lista de Tabelas

Capítulo 1

Introdução

Nas últimas duas décadas, as redes sociais se tornaram um dos principais meios de comunicação. Esse crescimento, em parte, se justifica pela massificação do acesso a internet incluindo dispositivos móveis como *smartphones* e *tablets*. Também alavancado pelos avanços computacionais e pelo desenvolvimento acelerado de novas técnicas e algoritmos, o aprendizado de máquina, em especial o processamento de linguagem natural, tem essas redes como importante objeto de estudo.

Desde a chamada Revolução Digital, observamos um progressivo barateamento e facilitação do uso de dispositivos eletrônicos. À medida que essas tecnologias passaram a ser acessíveis, não apenas para as corporações, mas também para os indivíduos, houve um crescente processo de digitalização de diversos aspectos de nossas vidas. Com a comunicação não foi diferente. O email, por exemplo, substituiu desde os anos 70 operações que até então eram apenas possíveis de forma analógica, como pelo uso de cartas. Nesse contexto, as redes sociais, ou mídias sociais, abordam aspectos diferentes da comunicação, mais dinâmica e informal.

Apesar de já existirem desde os anos 90, é com a virada do milênio que as primeiras grandes mídias sociais online aparecem, como *LinkedIn*, *MySpace* e *Orkut*. Desde então há um aumento anual da quantidade de seus usuários. Atualmente, estima-se que 3,5 bilhões de pessoas, ou 45% da população mundial utilize pelo menos uma rede social. Este número torna-se ainda mais interessante quando considerado que 4,4 bilhões de pessoas têm acesso à internet. Portanto, quase 80% dos internautas estão em alguma das mídias sociais. No Brasil esses números se acentuam ainda mais; 70% da população tem acesso à internet e 66% utiliza as redes sociais [1].

Além da alta penetração na sociedade, devido à disponibilidade proporcionada pelos dispositivos móveis, os usuários consomem boa parte de seu tempo nessas redes. No mundo, gasta-se em média 2 horas e 16 minutos por dia. Novamente esse número é ainda superior no Brasil, onde a média é de 3 horas e 34 minutos, sendo o segundo país no mundo a usar por mais tempo as redes, ficando apenas atrás das Filipinas.

Essa forte presença fez com que as mídias sociais não impactassem apenas as comunicações. Hoje em dia esses meios também são comumente utilizados para busca de relacionamentos, compartilhamento de notícias, divulgação de serviços, atendimento ao público, entre outros. As informações que trafegam nas redes exercem grande influência na formação de opinião das pessoas, seja ela em relação a um produto, a um evento ou até mesmo temas políticos, como pôde-se observar nas eleições pelo mundo nos últimos anos.

Portanto, a análise dessas informações, presentes nas redes, é importante para as mais diversas aplicações. Contudo, essa grande quantidade de usuários também se reflete no número de dados provindos das mídias sociais. Dentre as estatísticas de uso do ano de 2018 fornecidas pelas próprias redes sociais, tem-se que, diariamente, 300 milhões de fotos são publicadas no *Facebook*, 5 bilhões de vídeos são vistos no *YouTube*, 43 bilhões de mensagens são enviadas no *WhatsApp* e 100 milhões de usuários interagem pelo *Twitter*.

O massivo volume de dados inviabiliza que essas análises sejam feitas manualmente, tornando-se necessário o desenvolvimento de ferramentas capazes de automatizar esse processo. Entram aí as técnicas desenvolvidas pelo campo do Processamento de Linguagem Natural, do inglês *Natural Language Processing* (NLP). Foi a partir do anos 50 que esse termo passou a aparecer como um ramo da Inteligência Artificial. Devido à sua complexidade, a linguagem natural acabou inclusive se tornando um critério de inteligência, como proposto no teste de Turing [2].

O NLP é uma ampla área de pesquisa, abrangendo diferentes estágios da língua, desde os níveis de abstração mais baixos, como o estudo da fonologia e da sintaxe, quanto os maiores, que lidam com a semântica de determinado conteúdo. Neste ramo, busca-se desenvolver métodos capazes de auxiliar e/ou automatizar tarefas como: o reconhecimento de fala, a análise sintática de frases, a extração de entidades, a segmentação por tópicos, entre outros.

Esse conjunto de ferramentas são essenciais para tratar o volume de textos gerado pelas redes sociais. Com o avanço de técnicas de *Deep Learning*, nos últimos anos pôde-se observar um avanço significativo de desempenho, que permitiu até que tarefas de grande dificuldade, como a automatização de traduções e de aplicações de atendimento ao cliente por conversação, sejam possíveis.

Entretanto, apesar do grande potencial dessas técnicas, as mídias sociais apresentam características que diferenciam seu conteúdo dos tipos textuais que tradicionalmente se analisam em NLP, como artigos e textos jornalísticos. Isso dificulta o processamento dessa informação. Em geral, as redes sociais se apresentam como meio de conversação ágil, logo, as mensagens que circulam por elas costumam ser extremamente curtas, com amplo uso de abreviações. Por seu caráter informal, observa-se uma alta taxa de erros gramaticais e uma grande utilização de *emoti-*

*cons*¹. Além disso, o fato de serem meios de comunicação globais também ressalta a presença de estrangeirismos². Ademais a dinamicidade das redes sociais faz com que a evolução de sentido das palavras seja acelerada. Esses elementos trazem a necessidade de adaptação ou desenvolvimento de novas técnicas para se reproduzir o sucesso obtido pelas técnicas de processamento de texto em documentos com escrita mais formal e estruturada.

Porém, o principal fator que distingue as informações de redes de outros meios é a forte interligação entre diferentes tipos de mídia, como textos, imagens, áudios, fotos e vídeos. Além de metadados importantes, como localização, data e horário, uma propriedade importante das mídias sociais são os atributos referentes às redes de usuário. Exemplos desses dados são o número de amigos de um usuário da rede, e o número de re-compartilhamentos de uma mensagem. Logo, apesar da capacidade das ferramentas de NLP, existe um conjunto de informações que essas técnicas desconsideram, abrindo espaço para que as abordagens sejam multimodais, ou seja, que tratem de diversas destas propriedades.

De certa maneira, as interações entre usuários são o cerne das mídias sociais. Logo, técnicas que se dispõem a analisar esse tipo de informação também são de grande relevância. Nesse quesito, o campo das Redes Complexas, ou Ciência de Redes, é responsável por estudar os algoritmos e comportamentos observados em grafos que representam sistemas reais, como no caso das redes sociais. Assim como o aprendizado de máquina, essa esfera do conhecimento também apresenta um grande crescimento nos últimos anos, fornecendo um novo leque de tecnologias, de forma a descobrir-se aplicabilidades até então inexploradas. Dentre suas aplicações, que possuem importância para o estudos das mídias sociais, podemos ressaltar, por exemplo, a detecção de comunidades, identificação de principais influenciadores, modelagem de propagação de informação.

Finalmente, estes métodos são meios poderosos de análises de redes sociais, principalmente quando aplicados em complemento à informação textual. A inclusão desses dois últimos elementos na análise se faz necessária porque uma mesma mensagem pode ter conotações diferentes quando escrita por usuários de comunidades com ideias distintas. Por isso, um estudo que também considere esses diferentes tipos de elementos se faz importante.

1.1 Motivação

Dados são considerados um dos bens mais valiosos da atualidade, de forma mesmo a serem chamados de “o novo petróleo”. Isso porque, assim como o óleo, os dados

¹Sequência de caracteres ou pequena imagem que transmite um estado emotivo

²Uso de palavras ou expressões estrangeiras

são preciosos e precisam ser refinados para terem utilidade. Um dos aspectos dessa transformação pode ser observado na mudança cultural de organizações e empresas que passam a tomar decisões baseadas em dados e métricas coletadas.

A busca por informação de qualidade sobre um serviço ou produto sempre foi importante para consumidores. Quando não havia as tecnologias que usamos hoje essas pesquisas eram feitas majoritariamente no boca-a-boca ou a partir de revistas especializadas. Com a criação da internet e das redes sociais estas passaram também a exercer essa função, com o benefício de se encontrar opiniões de forma espontânea e em grande quantidade. As mídias sociais se tornaram um dos principais meios de compartilhamento dessa informação. As empresas, por sua vez, têm a oportunidade de utilizar as opiniões que trafegam nas redes para identificar falhas em suas mercadorias, melhorar sua segmentação, planejar novos produtos, entre outras atividades. Com o fácil acesso a coleta desses dados, as ferramentas capazes de extrair o sentimento dessas mensagens tornam-se fundamentais para viabilizar esse procedimento na escala em que ocorrem.

Apesar das dificuldades inerentes a classificação de mensagens de redes sociais técnicas de aprendizado de máquina, sobretudo *Deep Learning*, e de Redes Complexas apresentam êxito em várias tarefas realizadas sobre elas. Entretanto, o sucesso desses modelos, em geral, depende da quantidade de dados anotados disponíveis para treinamento. Como o processo de anotação é manual esse passa a ser o gargalo da construção de classificadores de sentimento.

Esse empecilho se torna ainda mais notável quando consideramos aplicações que requerem a análise de múltiplas línguas ou que tenham foco em um tema específico, necessitando criação de bases de dados próprias para cada caso de uso. Esses fatos motivam a elaboração de métodos que sejam independentes de bases de treinamento.

1.2 Objetivo

Esse projeto visa desenvolver um método capaz de formar classificadores de análise de sentimento sem a necessidade de bases de dados de treinamento. Essas análises serão feitas sobre dados de mídia sociais e serão explorados atributos tanto textuais quando de redes de usuários. A principal meta desse trabalho é viabilizar o emprego de modelos complexos e que apresentem melhores desempenhos, sem o custo proveniente da anotação de dados.

Para análise das mensagens será avaliado o impacto da utilização classificadores de *Deep Learning* em comparação a métodos lineares tradicionalmente aplicados em NLP. Diferentes arquiteturas de redes de aprendizado profundo serão experimentadas, como redes convolucionais e redes recorrentes. Além disso, as estratégias de representação de palavras também serão variadas. O processo será feito de maneira

semi-supervisionada com supervisão distante para anotação automática dos dados.

Técnicas de Redes Complexas serão aplicadas para caracterização de autores das mensagens. Modelos como Node2Vec e redes convolucionais de difusão, ambos também baseados em aprendizado de máquina, serão comparados. Neste caso, além de avaliar os modelos entre si, será analisado se adicionar informação do usuários quando aplicada em conjunto com o classificador textual decorre em alguma alteração de performance do sistema.

Concluindo, há um amplo conjunto de estudo aplicando de processamento de linguagem natural em redes sociais. Apesar de menor, a ciência de redes também têm um grande reportório de pesquisa sobre esse meio de comunicação. Este trabalho visa preencher a lacuna de sistemas que não necessitam de investimento em anotação de dados e que abrangem a multimodalidade da informação.

1.3 Organização do Texto

Esse documento é organizado da seguinte maneira:

- O Capítulo 2 apresenta o problema da análise de sentimento aplicada em mídias sociais e seus desafios. Esse Capítulo contém uma breve revisão bibliográfica de classificadores de análise de sentimento.
- No Capítulo 3 as técnicas de processamento de linguagem natural são apresentadas. O Capítulo descreve tanto os métodos lineares tradicionalmente aplicados a textos quanto os de *Deep Learning*. São caracterizadas também as diferentes formas de representação numérica de palavras.
- As ferramentas de ciência de redes são apresentadas no Capítulo 4. Nesse Capítulo são mostrados as diferentes técnicas aplicadas a modelagem de usuários de mídias sociais.
- O Capítulo 5 descreve o método proposto para desenvolvimento dos classificadores. São apresentados as etapas de formação de bases de dados, de anotação automática da mesma e de classificação.
- Os resultados dos obtidos pelos experimentos propostos são mostrados no Capítulo 6.
- Por fim, o Capítulo 7 avalia os resultados obtidos, apresenta as conclusões e enumera possíveis desdobramentos do trabalho realizado.

Capítulo 2

Análise de Sentimento de Redes Sociais

Análise de Sentimento é o campo de estudos que analisa a opinião, sentimento, atitude e emoções de pessoas em relação a entidades. Essas entidades podem ser pessoas, eventos, produtos, tópicos, entre outros. Na literatura também se encontram os seguintes nomes relacionados a esse ramo: *mineração de opinião*, *extração de opinião*, *mineração de sentimento*, *análise de subjetividade*, *análise de emoção* e *extração de críticas* [3]. Por ser aplicado na maioria das vezes a textos escritos, este campo é filiado ao processamento de linguagem natural, sendo uma de suas ramificações mais ativas.

Segundo CAMBRIA [4], a Análise de Sentimento é dividida entre duas principais tarefas, a extração de polaridade e a detecção de emoções. Enquanto a primeira foca em discernir conteúdo positivo de negativo, a segunda é responsável por classificar em emoções como: felicidade, medo, raiva, tristeza, entre outros.

A autora LIU [3], por sua vez, ressalta que há diferentes níveis de granularidade que essas para execução desta tarefa, a escolha do nível a ser utilizado depende da finalidade pela qual se aplicará a classificação e será uma das principais características para definir o tipo de técnica empregada. A Análise de Sentimento pode ser realizada a nível de documento, no qual um documento, como uma avaliação de produto, é avaliado como um todo. Nesses casos se assume implicitamente que um documento expressa uma opinião sobre uma única entidade, como o produto em questão no exemplo citado, também fica implícita que o documento expressa um sentimento único sobre a entidade. Por serem textos mais extensos, logo com mais informação, a acurácia de modelos nesses casos em geral é mais alta, portanto, até técnicas mais simples como as baseadas em dicionário como será apresentado em (TODO: referenciar exemplo com lexicon (subcapítulo seguinte 2.1)) podem apresentar resultados suficientemente bons. TABOADA *et al.* [5] exemplifica a extração de opinião de documentos a partir de técnicas de dicionário aplicadas em diferen-

tes bases de dados de avaliações de produtos. Por sua vez, DAS e CHEN [6] cacula predição de valor de ações a partir de análise de sentimento das mensagens presentes em um fórum online de investidores.

Visto que as limitações decorrentes de se classificar documentos por inteiro reduzem o escopo de aplicações, pode-se recorrer ao nível seguinte de granularidade, a classificação de sentimento de sentenças. A principal diferença entre essa abordagem e a anterior é a quantidade de informação disponível dado que uma sentença é composta, geralmente, por poucas palavras. Por outro lado, a premissa de sentimento único no conteúdo de uma frase é mais coerente com a realidade comparando-se com a classificação documento como um todo, sendo uma aproximação boa o suficiente para uma nova gama de casos de uso. Entretanto, LIU [3] ressalta que para o caso de sentenças é importante que a classificação de polaridade leve em consideração o sentimento neutro. Isso se torna relevante pois até mesmo dentro de documentos opinativos, como avaliações de produto, há sentenças puramente objetivas, que não expressarão polaridade sobre uma entidade. RILOFF *et al.* [7] apresentam o primeiro trabalho especificamente voltado para classificação de subjetividade de documentos. Devido a quantidade limitada de informação presentes em uma sentença, classificadores baseados em dicionários e em técnicas de aprendizado de máquina por modelos lineares, apresentam indicadores piores na execução da tarefa. São nesses cenário que nos últimos anos os modelos não lineares começaram a sobressair, como apresentam SOCHER *et al.* [8] e SOCHER *et al.* [9] que aplicam diferentes técnicas baseadas em *Deep Learning* para classificação de sentenças.

Por fim, apresenta-se a análise de sentimento de características. Uma entidade pode ser composta de diversos atributos. Ao falar sobre um filme pode se avaliar diferentes aspectos dele, como o roteiro, os atores, os personagens, etc. Uma crítica a esse determinado filme é composta de sentimentos distintos para cada um desses atributos, e o objetivo da análise de sentimento de características é identificar a polaridade de uma mensagem em relação aos atributos presentes. Para realizar essa análise são necessários elementos novos, como o reconhecimento da entidade citada e quais aspectos dela estão sendo avaliados. Também podem ser relevantes a identificação do autor e do momento do documento analisado. NASUKAWA e YI [10] e SNYDER e BARZILAY [11] são exemplos de trabalhos focados em mineiração de opinião focada em aspectos.

A Análise de Sentimento aplicada a redes sociais, em especial ao Twitter, se assemelha a classificação de sentimento de sentenças. Entretanto o perfil de mensagem que circulam as mídias sociais apresentam peculiaridades quando comparadas a meios convencionais. Por se tratar de um ambiente informal e de comunicação rápida, é comum encontrar erros gramaticais e abreviações. Similarmente os *emojis*, ícones ilustrativos de expressões faciais, também tem ampla adesão por serem

métodos práticos de exprimir sentimentos em poucos caracteres. Por se tratarem de redes globais, é frequente o emprego de palavras ou expressões de outras línguas em uma mesma mensagem. Esses fatores são obstáculos aos classificadores de linguagem natural, dificultando a tarefa de extração de polaridade.

O Twitter é uma rede social baseada em interações por mensagens curtas. hash-tag, retweet

Capítulo 3

Classificadores por Processamento de Linguagem Natural

Capítulo 4

Modelos de Redes Complexas

Capítulo 5

Método

Capítulo 6

Resultados e Discussões

Capítulo 7

Conclusões

Referências Bibliográficas

- [1] SOCIAL, W. A. “Digital in 2019: Global Overview”. <https://wearesocial.com/global-digital-report-2019>. acessado em 3 de Junho de 2019.
- [2] TURING, A. M. “Computing machinery and intelligence”, *Mind*, v. 59, n. 236, pp. 433–460, 1950.
- [3] LIU, B. *Opinions, Sentiment, and Emotion in Text*. Sentiment Analysis: Mining Opinions, Sentiments, and Emotions. Cambridge University Press, 2015. ISBN: 9781107017894. Disponível em: <<https://books.google.com.br/books?id=6IdsCQAAQBAJ>>.
- [4] CAMBRIA, E. “Affective computing and sentiment analysis”, *IEEE Intelligent Systems*, v. 31, n. 2, pp. 102–107, 2016.
- [5] TABOADA, M., BROOKE, J., TOFILOSKI, M., et al. “Lexicon-based methods for sentiment analysis”, *Computational linguistics*, v. 37, n. 2, pp. 267–307, 2011.
- [6] DAS, S. R., CHEN, M. Y. “Yahoo! for Amazon: Sentiment extraction from small talk on the web”, *Management science*, v. 53, n. 9, pp. 1375–1388, 2007.
- [7] RILOFF, E., WIEBE, J., PHILLIPS, W. “Exploiting subjectivity classification to improve information extraction”. In: *AAAI*, pp. 1106–1111, 2005.
- [8] SOCHER, R., PENNINGTON, J., HUANG, E. H., et al. “Semi-supervised recursive autoencoders for predicting sentiment distributions”. In: *Proceedings of the conference on empirical methods in natural language processing*, pp. 151–161. Association for Computational Linguistics, 2011.
- [9] SOCHER, R., PERELYGIN, A., WU, J., et al. “Recursive deep models for semantic compositionality over a sentiment treebank”. In: *Proceedings of the 2013 conference on empirical methods in natural language processing*, pp. 1631–1642, 2013.

- [10] NASUKAWA, T., YI, J. “Sentiment analysis: Capturing favorability using natural language processing”. In: *Proceedings of the 2nd international conference on Knowledge capture*, pp. 70–77. ACM, 2003.
- [11] SNYDER, B., BARZILAY, R. “Multiple aspect ranking using the good grief algorithm”. In: *Human Language Technologies 2007: The Conference of the North American Chapter of the Association for Computational Linguistics; Proceedings of the Main Conference*, pp. 300–307, 2007.