

## **APLICAÇÃO PARA ANÁLISE DE SENTIMENTOS EM TWEETS**

**Alexis Cesar Ruiz de Almeida** ; <https://orcid.org/0000-0002-1320-5254>  
Faculdade de Tecnologia de Sorocaba

**Maria das Graças Junqueira Machado Tomazela** ; <https://orcid.org/0000-0002-5471-2658>  
Faculdade de Tecnologia de Sorocaba



## APLICAÇÃO PARA ANÁLISE DE SENTIMENTOS EM TWEETS

### TWEETS SENTIMENT ANALYSIS APPLICATION

*Alexis Cesar Ruiz de Almeida*

<https://orcid.org/0000-0002-1320-5254>

503148448-99

*Centro Paula Souza – Fatec Sorocaba/SP*

*alexis.almeida@fatec.sp.gov.br*

*Orientador: Profa. Dra. Maria das Graças J. M. Tomazela*

<https://orcid.org/0000-0002-5471-2658>

085107058-28

*Centro Paula Souza – Fatec Sorocaba/SP*

*graca.tomazela@fatec.sp.gov.br*

**RESUMO:** As pesquisas na área de Análise de Sentimentos, embora recentes, vem evoluindo constantemente. A análise de sentimentos tem aplicações em diversas áreas, podendo ser encontrada na literatura em pesquisas sociais, extração de opinião de clientes e identificação de eventos. Assim, este trabalho teve como objetivo explorar os algoritmos utilizados na mineração de opinião para alcançar uma maior acurácia no julgamento dos textos obtidos por meio do Twitter. Desta forma foi desenvolvido um modelo de classificação com três algoritmos de inteligência artificial formando um comitê de classificadores visando a alcançar um aumento na acurácia de predições tornando-o assim um modelo confiável. Foi realizado também o desenvolvimento de um software para mineração de opiniões baseado em palavras-chave que faz uso do modelo desenvolvido previamente e que, por sua vez, pode ser utilizado para identificação de sentimento em relação às diversas organizações, temas, partidos etc., e com isso demonstrar de forma visual a utilização de tal tecnologia para o uso real em benefício de pesquisas sociais, corporativas e de cunho político. Este trabalho utilizou a abordagem experimental e para a realização dessa pesquisa foi primeiramente solicitado o acesso a API do Twitter para fins acadêmicos, seguido da fase de construção da base de dados formada por tweets classificados manualmente pelo autor, pré-processamento (com a utilização de técnicas para a obtenção de um melhor resultado de predição), processamento e avaliação do modelo gerado, com um comitê de classificadores composto pelos algoritmos mais utilizados para este fim: Naïve Bayes, Support Vector Machine e Regressão Logística. Como resultado, obteve-se um modelo de classificação bastante eficaz capaz de determinar os sentimentos do público-alvo em relação a um determinado conteúdo, com uma acurácia de 96,21%.

**ABSTRACT:** Research in the area of Sentiment Analysis, although recent, is constantly evolving. Sentiment analysis has applications in several areas, which can be found in the literature in

social research, customer opinion extraction and event identification. Thus, this work aimed to explore the algorithms that are used in opinion mining to obtain better accuracy in the judgment of texts retrieved from Twitter. In this way, a classification model was developed with three different artificial intelligence algorithms that grouped together form an ensemble that aims to experience an increase in predictions making it reliable. In addition, a keyword-based opinion mining software was developed that implements the model created above and that can be used to identify sentiments in relation to different organizations, topics, political parties, etc. And with that, demonstrate the use of such technology for beneficial purposes such as social, corporate or political research. This work used the experimental approach and to carry it out, access to the Twitter API was first requested, followed by the construction phase of the database formed by tweets manually classified by the author, pre-processing (with techniques for better results in the predictions), processing and evaluation of the generated model, with a classification committee composed of the most used algorithms for this purpose: Naïve Bayes, Support Vector Machine and Logistic Regression. As a result, an effective classification model was obtained, capable of determining the target audience's feelings in relation to certain content, with an accuracy of 96.21%.

**PALAVRAS-CHAVE:** Análise de Sentimentos; Aprendizado de Máquina; Processamento de Linguagem Natural; Classificação, Twitter.

**KEYWORDS:** Sentiment Analysis; Machine Learning; Natural Language Processing; Classification; Twitter.

## 1 INTRODUÇÃO

A Análise de Sentimentos é um campo dentro da área de Processamento de Linguagem Natural que tem como objetivo identificar opiniões, sentimentos e até emoções em informações subjetivas (LIU, 2010). Há muitas aplicações para a técnica de análise de sentimento, principalmente no âmbito organizacional em que as empresas podem obter de forma mais ágil milhares de opiniões de seus clientes a respeito de seus produtos como pode ser visto em Sarlan, Nadam e Basri (2014).

A rede social Twitter é uma ótima escolha para o estudo da análise de sentimentos, pois, além de ser uma rede social popular e com foco em textos, as publicações contam com um limite de 280 caracteres, facilitando assim a análise realizada (FELL e LUKIANOVA, 2019).

Para que o computador determine a classe de um texto, são utilizados algoritmos classificadores. Cada classificador tem seu próprio método de atribuir uma classe a um documento, desde cálculos probabilísticos simples até redes neurais. A capacidade de um algoritmo prever a

que conjunto um dado elemento pertença com base em seus atributos se dá por conta de um treinamento a priori, sendo este denominado um aprendizado supervisionado; também há a possibilidade de um algoritmo identificar a classe de um elemento sem um treinamento prévio, sendo um aprendizado não-supervisionado (Feldman, 2013).

Um dos problemas mais comuns durante a análise é a má escrita dos textos, principalmente no Twitter (MARTÍNEZ-CÁMARA, E. et al., 2012). Os usuários costumam utilizar várias gírias, abreviar palavras, se referir ao mesmo elemento de formas diferentes durante o texto etc. Isso está ligado com a quantidade limitada de caracteres na publicação, o que faz com que os usuários adotem uma escrita mais rápida e informal. Para tratar essa questão, são utilizadas técnicas de processamento de linguagem natural, como o processo de redução de palavras, remoção de palavras vazias, remoção de caracteres especiais, links, *hashtags* etc. Há diversas técnicas a se utilizar para tornar o texto mais fácil de se analisar para o computador, e diversas delas são implementadas, por exemplo, na biblioteca Natural Language Toolkit (NLTK) da linguagem Python.

Desta forma, este trabalho teve como objetivo explorar os algoritmos utilizados na mineração de opinião para alcançar uma maior acurácia no julgamento dos textos obtidos por meio do Twitter.

## 2 METODOLOGIA

Nesta pesquisa, optou-se pela abordagem experimental, que consiste em determinar um objeto de estudo, selecionar as variáveis que seriam capazes de exercer influência sobre o objeto, definir as formas de controle e de observação dos efeitos produzidos no objeto por essas variáveis (GIL, 2007).

Para o desenvolvimento do experimento foram realizadas as seguintes etapas: solicitação de acesso à API do Twitter, coleta e anotação manual de tweets e suas polaridades (totalizando 400 tweets, sendo 200 classificados como negativos e 200 como positivos), fase de pré-processamento (a qual envolveu as atividades de: tokenização, filtragem e padronização, remoção de *stop words* e processo de *stemming*), fase de processamento (treinamento dos classificadores, sendo reservada uma quantidade equivalente a 66,66% da base de dados para treino e 33,33% para teste), construção do comitê com o método de votação (*Voting*) e avaliação do modelo.

## 3 DESENVOLVIMENTO

É considerado como sendo mineração de opinião qualquer estudo computacional que envolva opinião (sentimento, avaliação, ponto de vista, emoção e subjetividade) de forma textual (LIU, 2010). Segundo Feldman (2013) a análise de sentimentos pode ser feita por meio das quatro seguintes

abordagens diferentes: em nível de documento, em nível de sentença, em nível de aspecto e comparativa.

O processamento de linguagem natural (PLN) é uma área que explora como os computadores podem ser usados para compreender e manipular textos ou falas em línguas humanas naturais. (CHOWDHURY, 2005). A Análise de Sentimentos utiliza várias técnicas de PLN para facilitar o processo de identificação das palavras. Algumas dessas técnicas são: *Stemming* (LOVINS, 1968), *Tokenization* (WEBSTER e KIT, 1992), Remoção de *Stopwords* (SILVA e RIBEIRO, 2003).

O aprendizado de máquina é um campo da Inteligência Artificial que tem como objetivo construir modelos matemáticos capazes de prever ou classificar algo com base em uma base de dados de exemplo (ZHANG, 2020). A análise de sentimentos utiliza o aprendizado de máquina para julgar um dado texto pertencente a uma classe ou não. Segundo Feldman (2013) e Zhang (2020), os dois principais meios de aprendizado são:

- Aprendizado supervisionado, no qual é assumido que existe um número finito de classes e os dados que serão disponibilizados para o treinamento do modelo são anotados com sua respectiva classe;
- Aprendizado não supervisionado: em que os dados disponibilizados para treino do modelo não são anotados com suas classes. É tarefa do próprio algoritmo analisar os padrões e utilizá-los para discriminar grupos a partir desses dados.

Além destes, também há o aprendizado por reforço, sendo análogo ao aprendizado animal. Diferente do aprendizado supervisionado, não há um treino sobre o que é bom ou ruim; o agente precisa perceber quando algo bom ocorreu por meio de um feedback. Um feedback positivo é chamado de recompensa e o agente precisa estar programado para identificar essa recompensa e não a tratar apenas como qualquer outra entrada sensória. (RUSSEL e NORVIG, 2013).

#### 4 RESULTADOS OBTIDOS

Durante o desenvolvimento deste experimento, foram testadas combinações diferentes de técnicas para o pré-processamento a fim de obter o maior nível de acurácia. A versão final do comitê tem seu pré-processamento composto pelas técnicas de: remoção de *stopwords*, *stemming*, remoção de caracteres não alfanuméricos; remoção de hiperlinks, *hashtags* e menções.

A acurácia para cada classificador e o comitê (modelo composto pelos 3 classificadores), para uma base de dados contendo 400 tweets classificados como positivos ou negativos foi a seguinte:

- *Support Vector Machine*: 82,19%
- Regressão Logística: 81,28%

- Naïve Bayes: 79,16%
- Comitê: 96,21%

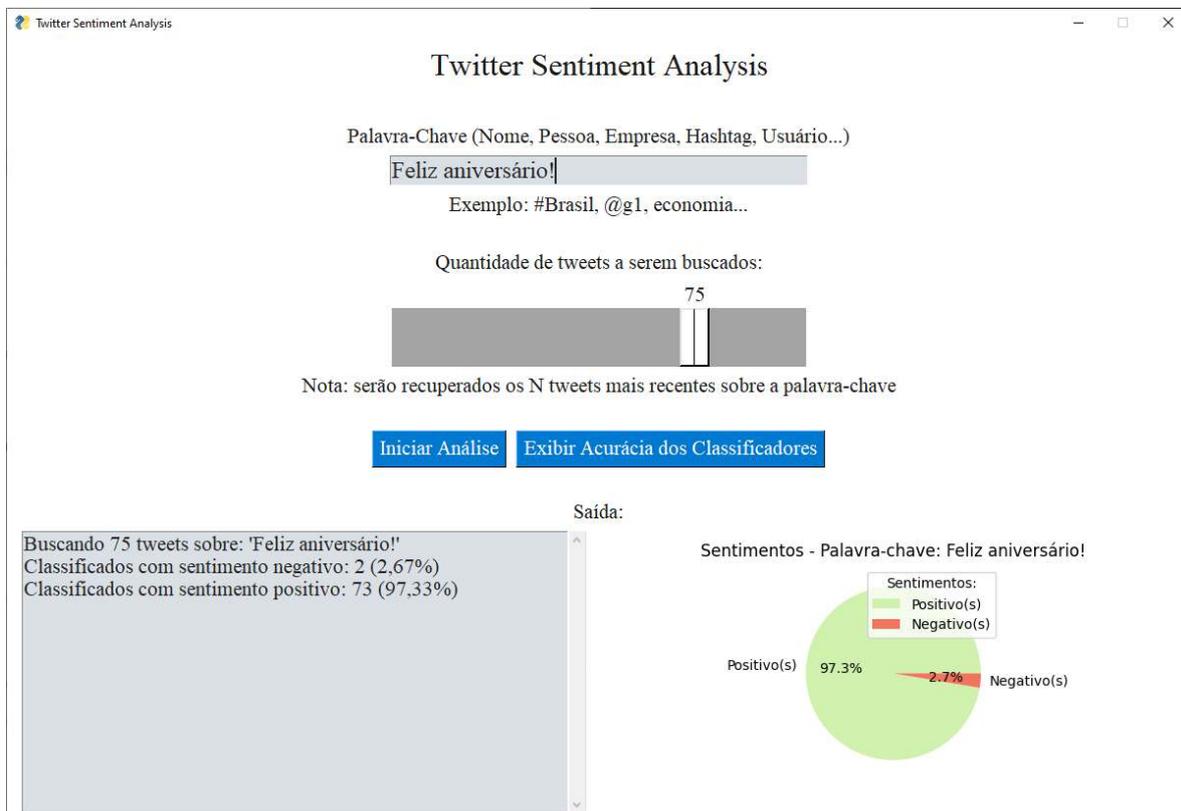
A tabela 1 apresenta os resultados obtidos de cálculos realizados a partir das matrizes de confusão geradas para cada modelo. Tabela 1 – Resultados gerados pelos modelos de classificação

Classificador	Acurácia	Precisão (Positivo)	Precisão (Negativo)	Sensibilidade (Positivo)	Sensibilidade (Negativo)
SVM	82,19%	88,40%	84,12%	85,91%	86,88%
Regressão Logística	81,28%	87,09%	65,71%	69,23%	85,18%
Naïve Bayes	79,16%	79,36%	78,26%	76,92%	80,59%
Comitê	96,21%	96,96%	96,96%	96,96%	96,96%

Fonte: Elaborada pelo autor

Os classificadores foram integrados a uma aplicação gráfica que permite buscar palavras-chaves no Twitter, possibilitando ao usuário escolher um intervalo entre 1 e 100 tweets, ao final exibindo os resultados (junto a um gráfico) referentes a quantidade dos tweets classificados como positivos e a quantidade dos tweets classificados como negativos. A aplicação final é apresentada na figura 7.

Figura 7 – Captura da aplicação final



Fonte: Elaborada pelo autor

Algumas particularidades puderam ser notadas no desenvolvimento do comitê.

Durante a coleta, alguns tweets inseridos na base de dados apresentavam múltipla polaridade da mesma forma. no tópico de Análise a Nível de Aspecto; esses tweets geralmente tratavam de resenhas de produtos ou serviços; alguns tweets desse gênero foram removidos da base de dados para não desequilibrar a polaridade de cada coleção no banco de dados.

Outro ponto foi que, mesmo durante a coleta manual, pôde-se notar que a maioria dos tweets, em diversos assuntos, têm a polaridade negativa predominante. Este fato somado a base de dados desbalanceada proporcionava resultados que praticamente excluía a polaridade positiva. Porém depois do balanceamento da base de dados, as classificações ficaram mais equilibradas. Foi possível perceber também o grande impacto da fase de pré-processamento e suas técnicas, uma vez que os classificadores apresentavam uma acurácia em torno de 60% antes da aplicação das medidas de preparação do texto.

## 5 CONSIDERAÇÕES FINAIS

Foi desenvolvido um modelo de classificação composto por três classificadores, sendo implementado em uma aplicação gráfica. O comitê de classificadores utilizou da abordagem de votação (*Voting*) para decidir a classificação final de um texto. Para o treino e teste, alguns tweets foram coletados e anotados manualmente como sendo positivos ou negativos, sendo gravados posteriormente em um banco de dados as palavras e suas respectivas probabilidades (calculadas pelo algoritmo) de pertencer a cada classe.

Foi possível perceber a grande influência que a fase de pré-processamento exerce na classificação. Os resultados, foram bastante satisfatórios, mostrando que o comitê com uma acurácia de 96,21% superou os todos os três classificadores analisados individualmente.

Algumas das limitações no desenvolvimento deste trabalho foram: Quantidade elevada de tweets neutros, *emojis* e enorme variação de gírias e abreviações.

Apesar das limitações encontradas, os resultados dos experimentos tornam possível concluir que os objetivos deste trabalho foram alcançados.

## REFERÊNCIAS

CHOWDHURY, G. G. **Natural language processing**. Annual Review of Information Science and Technology, 2005.

FELDMAN, R. **Techniques and applications for sentiment analysis**. Communications of the ACM, 2013.

FELL, E; LUKIANOVA, N. **Twitter (Digital Media and Society)**. European Journal of Communcation, 2019.

GIL, A. C. **Métodos e Técnicas de Pesquisa Social**. 6. ed. São Paulo: Atlas, 2008.

HAN J; KAMBER, M; PEI, J. **Data Mining: Concepts and Techniques**. 3. ed. San Francisco: Elsevier, 2011

LIU, B. **Sentiment Analysis and Subjectivity**. Department of Computer Science. University of Illinois. 2. ed. Chicago, 2010.

LOVINS, J. B. **Development of a Stemming Algorithm**. Electronic Systems Laboratory. Massachusetts Institute of Technology. v. 11. Cambridge, 1968.

MARTÍNEZ-CÁMARA, E. et al. **Sentiment analysis in Twitter**. Natural Language Engineering. Cambridge University. 2012.

RUSSEL, N; NORVIG, P. **Inteligência Artificial**. 3. ed. Editora Campus, 2013.

SARLAN, A; NADAM, C; BASRI, S. **Twitter Sentiment Analysis**. Computer Information Science. Universiti Teknologi PETRONAS. Perak, 2014.

SILVA, C; RIBEIRO, B. **The Importance of Stop Word Removal on Recall Values in Text Categorization**. Computer Science. Proceedings of the International Joint Conference on Neural Networks. Coimbra, 2003.

WEBSTER, J; KIT, C. **Tokenization as the initial phase in NLP**. Proceedings of the 14th conference on Computational linguistics. v. 4. City Polytechnic of Hong Kong. Kowloon, 1992.

ZHANG, X.-D. **A Matrix Algebra Approach to Artificial Intelligence**. 1. ed. Springer Singapore, 2020. p. 223.