

Ricardo César Gonçalves Sant'Ana

Moisés Lima Dutra

Guilherme Ataíde Dias

Organizadores

WIDaT 2018

II WORKSHOP DE INFORMAÇÃO,
DADOS E TECNOLOGIA

ANAIS
WIDaT 2018

Organização do WIDaT 2018

- **Organização Geral:**

Guilherme Ataíde Dias (PPGCI-UFPB) - Coordenador geral do evento
Moisés Lima Dutra (PPGCIN-UFSC) - Vice-coordenador

- **Coordenador da Comissão Científica:**

Ricardo César Gonçalves Sant'Ana (PPGCI-UNESP)

- **Comissão científica**

Adilson Luiz Pinto (PPGCIN-UFSC)
Ana Alice Baptista (Universidade do Minho, Portugal)
Ana Carolina Simionato (PPGCI-UFSCar)
Angela Maria Grossi de Carvalho (PPGCI-UNESP)
Bernardina Maria Juvenal Freire de Oliveira (PPGCI-UFPB)
Cristian Berrío-Zapata (PPGCI-UFPA)
Dalton Lopes Martins (FCI-UnB)
Denysson Axel Ribeiro Mota (PPGB-UFCA)
Douglas Dyllon Jeronimo de Macedo (PPGCIN-UFSC)
Ed Porto Bezerra (PPGI-UFPB)
Edgar Bisset Alvarez (PPGCIN-UFSC)
Edna Gusmão de Goés Brennand (MPGOA-UFPB)
Edna Gomes Pinheiro (DCI-UFPB)
Elaine Parra Affonso (FATEC-SP)
Elvis Fusco (UNIVEM-Marília)
Enrique Muriel Torrado (PPGCIN-UFSC)
Evandro de Barros Costa (IC-UFAL)
Fábio Paraguaçu (IC-UFAL)
Fernando de Assis Rodrigues (PPGCI-UNESP)
Gustavo Medeiros de Araújo (PPGCIN-UFSC)
Henry Pôncio Cruz de Oliveira (PPGCI-UFPB)
Joana Coeli Ribeiro Garcia (PPGCI-UFPB)
José Eduardo Santarém Segundo (USP-FFCLRP)
Leonardo Castro Botega (UNIVEM-Marília)
Luana Farias Sales Marques (PPGCI-IBICT-UFRJ)
Marckson Roberto Ferreira de Sousa (PPGCI-UFPB)
Luís Fernando Sayão (CNEN)
Marcelo Morandini (EACH-USP)
Márcio Matias (PPGCIN-UFSC)
Marcos Mucheroni (CBD-USP)
Marynice de Medeiros Matos Autran (PPGCI-UFPB)

Maurício Barcellos Almeida (PPGGOC-UFMG)
Moisés Lima Dutra (PPGCIN-UFSC)
Plácida Leopoldina V. da Costa Santos (PPGCI-UNESP)
Pedro Luiz Pizzigatti Corrêa (POLI-USP)
Renata Baracho (PPGGOC-UFMG)
Ricardo César Gonçalves Sant'Ana (PPGCI-UNESP)
Robson Rodrigues Lemos (UFSC-Araranguá)
Rogério Ramalho (PPGCI-UFSCar)
Ryan Ribeiro de Azevedo (UFRPE-UAG)
Sandra de Albuquerque Siebra (PPGCI-UFPE)
Sandro Rautenberg (DECOMP-UNICENTRO)
Silvana Aparecida Borsetti G. Vidotti (PPGCI-UNESP)
Virginia Bentes Pinto (PPGCI-UFC)
Wagner Junqueira de Araújo (PPGCI-UFPB)
Zaira Regina Zafalon (PPGCI-UFSCar)

- **Coordenador do Cerimonial:**

André Luiz Dias de França (PPGCI-UFPB)

- **Coordenador da Equipe Técnica Local:**

Laerte Pereira da Silva Júnior (CCHLA-UFPB)

- **Equipe Técnica Local:**

Adriana Alves Rodrigues (PPGCI-UFPB)
Antonio Felipe dos Santos (MPGOA-UFPB)
Débora Gomes de Araújo (PPGCI-UFPB)
Pedro Augusto de Lima Barroso (PPGCI-UFPB)
Pollianna Marys de Souza e Silva (PPGCI-UFPB)
Renata Lemos dos Anjos (PPGCI-UFPB)

ANÁLISE DE SENTIMENTOS:

Identificando sentimentos em comentários da Rede Humaniza SUS

SENTIMENT ANALYSIS:

Identifying sentiment in comments on Humaniza SUS network

Eduardo Alves Silva¹, Luis Felipe Rosa de Oliveira⁽²⁾

(1) Universidade Nova de Lisboa - NOVAIMS, Campolide - Lisboa, easilva91@gmail.com

(2) Universidade Federal de Goiás, Goiânia, Goiás, luisfelipeprf@gmail.com

Resumo:

A análise de sentimentos e áreas correlatas ao processamento de linguagem natural tem-se tornado cada vez mais comuns em diversos contextos, seja no empresarial ou acadêmico. Esse tipo de análise facilita a compreensão a respeito de opiniões e sentimentos de diferentes indivíduos em relação a determinado produto ou temática. Ao levar em consideração a crescente difusão das redes sociais e a interação entre usuários nessas redes, a captação de dados para análise de sentimento se tornou mais simples do que era anos atrás, onde frequentemente eram utilizados questionários ou formulários para captar a opinião de uma pessoa a respeito de algo. Com esse intuito e com a alta densidade de informação da rede humaniza SUS (RHS), a primeira rede social brasileira criada a partir de uma política pública, voltada para o tema da humanização da saúde. A rede humaniza SUS torna-se um ambiente de estudo promissor para análise de sentimentos, fazendo-se uso de uma abordagem não supervisionada que visa o uso de um dicionário ou léxico de sentimentos criado com palavras em Português, e com a possibilidade de aplicação utilizando a linguagem de programação python, para apresentar a classificação de sentimentos dos comentários feitos pelos usuários da rede.

Palavras-chave: Rede Humaniza SUS. Análise de Sentimentos. Léxico de Sentimentos.

Abstract:

The sentiment analysis and areas related to the natural language processing has become increasingly common in several contexts, be it business or academic. This type of analysis facilitates the understanding of the opinions and feelings of different individuals regarding a given product or thematic. By taking into account the growing diffusion of social networks and the interaction among users in these networks, capturing data for feeling analysis has become simpler than it was years ago, where questionnaires or forms were often used to capture a person's opinion about something. With this intention and with the high density of information of the humaniza SUS network (RHS), the first Brazilian social network created from a public policy, focused on the humanization of health. The humaniza SUS network becomes a promising study environment for analysis of feelings, making use of an unsupervised approach that aims to use a dictionary or sentiment lexicon created with words in Portuguese, and with the possibility of application using the programming language python, to present the classification of feelings of the comments made by the users of the network.

Keywords: Humaniza SUS Network. Sentiment Analysis. Sentiment Lexicon.

I INTRODUÇÃO

A crescente popularização das redes sociais, blogs e plataformas que incentivam a participação e colaboração de usuários, tornou-se comum permitindo o compartilhamento de mensagens que refletem suas opiniões e por vezes sentimentos a respeito de determinada temática ou produto.

Nesse sentido, a análise de sentimentos, ou mineração de opinião, corresponde ao problema de tentar identificar ou extrair emoções, opiniões ou pontos de vista expressados em um texto (DUARTE, 2013).

Liu (2012) menciona que esse tipo de análise está ligado a linguística e ao processamento de linguagem natural (PLN), gerando grande impacto nessa área do conhecimento, podendo ter grandes reflexos também em estudos sobre ciências políticas, economia, e ciências sociais, uma vez que são afetados pela opinião das pessoas.

O uso da análise de sentimentos é extenso, sendo usado para verificar o quanto um produto é bem aceito, de acordo com as opiniões de usuários, ou para compreender o sentimento gerado por uma publicação no Facebook, ou uma hashtag no Twitter. Essa variabilidade torna o uso da análise de sentimentos e opinião, um amplo campo de estudo.

Para a definição de sentimentos em um texto, normalmente é utilizado um léxico de sentimento, servindo como um dicionário de polaridade (FREITAS; VIEIRA, 2015), que agrega palavras que comumente expressam sentimentos positivos ou negativos (LIU, 2012). Atualmente os léxicos e modelos de análise de sentimento na língua inglesa são bastante comuns e assertivos em relação a análise de opinião e sentimento, no entanto, em Português são conhecidos quatro léxicos: OpLexicon, SentiLex, Brazilian Portuguese Linguistic Inquiry and Word Count (LIWC) e Onto.PT (FREITAS; VIEIRA, 2015).

O presente trabalho, tem por objetivo identificar os sentimentos expressados por comentários, fazendo uso do OpLexicon, da primeira rede social do Brasil criada no âmbito de uma política pública. Trata-se da rede humaniza SUS (RHS)¹ que se encontra em atividade desde o ano de 2008, atualmente com mais de 30 mil usuários por todo o Brasil, onde muitos são responsáveis pelas mais de 14 mil publicações existentes na rede e os quase 40 mil comentários, existentes nessas publicações.

A rede humaniza SUS, traz consigo uma temática bastante específica, a humanização da saúde. Nesse sentido, foi verificado que tipo de sentimentos/opiniões os usuários da rede expressam a partir dos comentários, servindo também como um experimento para o uso de análise de sentimentos em Português.

2 ANÁLISE DE SENTIMENTOS

A análise de sentimento é, talvez, a aplicação mais popular da análise de texto, com um grande número de tutoriais, sites e aplicativos que se concentram em analisar o sentimento de vários recursos textuais, desde pesquisas corporativas até análises de crítica de filmes (SARKAR, 2016).

Em sua utilização a análise de sentimento tem uma série de fatores a serem considerados para que se consiga resultados que possam ser considerados coerentes, passando pela fase de pré-

¹ Site da RHS – www.redehumanizausus.net

processamento do texto, até o uso de abordagens para identificação de sentimentos, potencialmente utilizando métodos supervisionados que envolvem aprendizado de máquina e não supervisionados, que faz uso de bancos de dados, ontologias ou léxicos de sentimento.

Um dos principais objetivos da análise de sentimentos é analisar um determinado gênero de texto para compreender o que ele expressa. Nesse sentido alguns fatores devem ser considerados, como a polaridade do sentimento e sua subjetividade, de acordo com Sarkar (2016) a análise de sentimentos funciona melhor em textos que têm um contexto subjetivo (opiniões) do que em textos objetivos (fatos).

Vale ressaltar que, a simples busca por palavras como “bom” ou “ruim” não é o suficiente para expressar o sentimento em um texto (DOSCIATTI; FERREIRA, 2013).

2.1 Léxico de Sentimento

Um léxico de sentimento pode ser um dicionário, vocabulário ou conjunto de palavras que tem uma polaridade positiva ou negativa atribuída (SARKAR, 2016). Para o presente trabalho foi utilizado um método não supervisionado. Sendo assim utilizado o léxico OpLexicon.

O OpLexicon é constituído de um total de 32.191 itens (24.475 adjetivos e 6.889 verbos, Tabela 1), tendo sua construção feita com base em textos jornalísticos e resenhas de filmes escritas em Português do Brasil, além do uso de tesouros e a tradução do léxico de opinião em inglês (SOUZA et al, 2012).

Tabela 1 – Quantidade de palavras por tipo OpLexicon

Tipo	Quantidade
adj (adjetivo)	24.475
vb (verbo)	6.889
hashtag	471
vb (verbo) det (determinante) n (nome) prp (preposição)	103
vb (verbo) n (nome) prp (preposição)	91
vb (verbo) adj (adjetivo)	74
emoticons	66
vb (verbo) adv (adverbio)	22

Fonte: OpLexicon v3.0

3 PROCEDIMENTOS METODOLÓGICOS

O primeiro passo para se atingir o objetivo proposto foi captar os comentários da RHS em um recorte de 9 anos, a partir do banco de dados foram extraídos comentários de 2008 até o fim de 2017, o que gerou um total de 36.346 comentários.

Após essa etapa, foi iniciado o tratamento do texto dos comentários, utilizando a linguagem de programação python, juntamente com as etapas de sumarização e normalização de texto.

Segundo Sakar (2016) a normalização de texto é um processo de limpeza, normalização e padronização de dados com técnicas de remoção de símbolos e caracteres especiais, remoção de tags HTML (*Hypertext Markup Language*), remoção de stop words (preposições, pronomes, artigos), correção de grafia, *stemming* (reduzir palavras ao seu radical) e lematização (reduzir a flexão das palavras).

Os itens de normalização utilizados neste trabalho foram:

1. Remoção de tags HTML.
2. Remoção de símbolos e caracteres especiais.
3. Remoção de palavras irrelevantes (*Stopwords*).
4. Lematização.

Uma vez que os comentários da RHS seguem um padrão de escrita compreensível e coerente, itens como a correção gráfica não se mostra de extrema importância nessa análise, por outro lado, os comentários tem inúmeras tags HTML que foram tratadas durante o processo (Apêndice A).

Por sua vez, seguindo as outras etapas da normalização, como a remoção de caracteres especiais que removem itens como (“”, “?”, “/”, “;”, “:”), e a remoção de palavras irrelevantes, que é feita a partir de um dicionário de artigos, preposições, pronomes (a, e, é, ali, uns), entre outros, com isso o texto passa a ser transformado novamente, como é demonstrado no Apêndice B.

Vale mencionar que as *stopwords*, podem seguir determinados padrões, nessa análise foi usado um padrão comum da língua portuguesa e foram adicionados alguns termos que se encontravam nos comentários e não agregavam valor para a análise final.

A lematização por sua vez, se trata do da remoção de afixos das palavras, fazendo com que essa palavra volte para sua base raiz (SARKAR, 2016), levando em consideração a parte do discurso em que a palavra se encontra, ou seja, caso a palavra se encontre na posição de verbo dentro de um texto, ela será lematizada de uma forma, caso seja um adjetivo será de outra forma e assim por diante.

Uma vez que para a lematização é necessário identificar a que parte do discurso a palavra pertence e após isso lematizar a palavra, é possível fazer o uso de diferentes bibliotecas e ferramentas do python para essa tarefa, neste trabalho foi utilizada a biblioteca *Spacy*², que trabalha com uma variedade de idiomas incluindo o Português.

A lematização se faz importante pois, como demonstrado o léxico utilizado tem um número relevante de palavras que são representadas como parte do discurso, como adjetivos e verbos, conseguir lematizar o texto impacta no resultado final de identificação do sentimento de cada comentário.

Para identificar os sentimentos, foi utilizada uma abordagem similar ao utilizado para verificação de sentimentos com o léxico de opinião (HU;LIU, 2004), que se trata de um léxico com palavras positivas (1) e negativas (-1), entanto, o OpLexicon traz palavras neutras (0), ao adaptar o script da biblioteca NLTK³ que faz uso do léxico de opinião foi possível identificar o sentimento dos comentários.

2 Biblioteca Spacy - <https://spacy.io/>

3 NLTK (Natural Language Toolkit) é uma biblioteca em Python que reúne funções e *corpus* textuais para tratamento e análise de texto. - <https://www.nltk.org/>

De forma geral, o script acessa cada um dos comentários, separa as palavras, verifica se essa palavra se encontra no dicionário léxico e qual sua polaridade, após esse processo, é feita uma contagem e caso existam mais palavras positivas no comentário, o mesmo é considerado positivo e assim por diante.

4 RESULTADOS

O trabalho com texto gera resultado diversos, que vão além da identificação de sentimentos, dessa forma foi feita uma verificação nos textos e palavras que geraram uma visão descritiva dos dados. Como por exemplo as palavras mais comuns encontradas nesses mais de 36 mil comentários (Apêndice C).

Além das palavras mais comuns é possível visualizar a interação em relação ao número de palavras utilizadas nos comentários e qual seria a média padrão de palavras.

De forma geral, foram encontrados mais de 25 mil comentários que foram classificados como positivos, 4,594 negativos e em torno de 4.692 comentários considerados neutros.

Faz-se importante ressaltar que o resultado final não apresenta indicadores de um modelo de aprendizado de máquina como, accuracy, fl-score, entre outros pelo fato de todo o trabalho ter sido desenvolvido em volta de um léxico de sentimento, o que foi feito em si representa a identificação de sentimentos nos comentários de acordo com palavras que estão presentes no léxico de sentimento, não sua classificação, a partir desse trabalho o treino de um classificador de sentimentos se dá em uma etapa posterior.

4 CONSIDERAÇÕES FINAIS

O estudo de identificação de sentimentos nos comentários da rede humaniza SUS, apresenta um resultado que demonstra estar de acordo com aquilo que a rede se propõe, que seria a colaboração entre os usuários em torno de uma temática específica.

Identificar sentimentos positivos em sua grande maioria pode ser visto como um incentivo para que o trabalho de desenvolvimento da rede tenha continuidade, foi possível perceber durante o estudo que o teor dos comentários em sua maioria, são de mensagens de agradecimento, por divulgação de encontros ou iniciativas relacionadas a humanização da saúde no Brasil, gerando um assim um ambiente de aprendizado e gratidão por todos aqueles que estão envolvidos nesse processo.

Os resultados obtidos em relação a esse estudo demonstram uma abordagem inicial em relação a tudo aquilo que pode ser produzido a partir do uso da análise de sentimentos e processamento de linguagem natural. A rede humaniza sus é um ambiente próspero para a possível criação de um léxico de sentimento relacionado a saúde e a humanização da saúde.

Essa mesma abordagem gera caminhos para que, a partir dos resultados desse estudo, seja criado um classificador de sentimentos utilizando métodos supervisionados e aprendizado de máquina, o que daria um contexto mais abrangente para os resultados, bem como para o uso dos dados da RHS.

Por fim, este trabalho agrega um valor experimental para o desenvolvimento de análises de sentimento utilizando o Português brasileiro, com ênfase em um método não supervisionado e com a possibilidade de construção de dados para sua posterior aplicação em algoritmos e classificadores de aprendizagem de máquina.

REFERÊNCIAS

DOSCIATTI, M. M.; FERREIRA, E. C. L. P. C. Identificando emoções em textos em português do brasil usando máquina de vetores de suporte em solução multiclasse. In: **ENIAC-Encontro Nacional de Inteligência Artificial e Computacional**. Fortaleza, Brasil, 2013.

DUARTE, E. S. **Sentiment analysis on twitter for the portuguese language**. Tese (Doutorado) — Faculdade de Ciências e Tecnologia, 2013.

FREITAS, L. d; VIEIRA, R. Exploring resources for sentiment analysis in portuguese language. In: IEEE. In: **2015 Brazilian Conference on Intelligent Systems**. BRACIS. [S.l.], 2015. p. 152–156.

HU, M; LIU, B. Mining and summarizing customer reviews. **Proceedings of the ACM SIGKDD International Conference on Knowledge Discovery & Data Mining**, Seattle, Washington, p.168-177, 2004

LIU, B. **Sentiment Analysis and Opinion Mining: Synthesis Lectures on Human Language Technologies**. Morgan & Claypool Publishers, 180p., 2012.

SARKAR, D. **Text Analytics with Python: A Practical Real-World Approach to Gaining Actionable Insights from Your Data**. Library of Congress Control Number, Apress, 385p., 2016.

SOUZA, M.; VIEIRA, R.; Busetti, D.; CHISHMAN, R. e ALVES, I. M. Construction of a Portuguese Opinion Lexicon from multiple resources. In: **8th Brazilian Symposium in Information and Human Language Technology**, 2012

APÊNDICE A – TÍTULO

Apêndice A – Exemplo de remoção de tags HTML

```
<P>Excelente post, Mariella!</P>  
<P>Acho que você e o Bruno poderão e deverão ter um papel  
fundamental nesta Rede, orientando-nos em relação a como  
redigir um post</P>  
<P>Precisaremos, certamente, de um "manual de redação"!</P>  
<P>Um abraço,</P>  
<P>Ricardo<BR></P>
```

Excelente post, Mariella! Acho que você e o Bruno poderão e deverão ter um papel fundamental nesta Rede, orientando-nos em relação a como redigir um post. Precisaremos, certamente, de um "manual de redação"! Um abraço, Ricardo

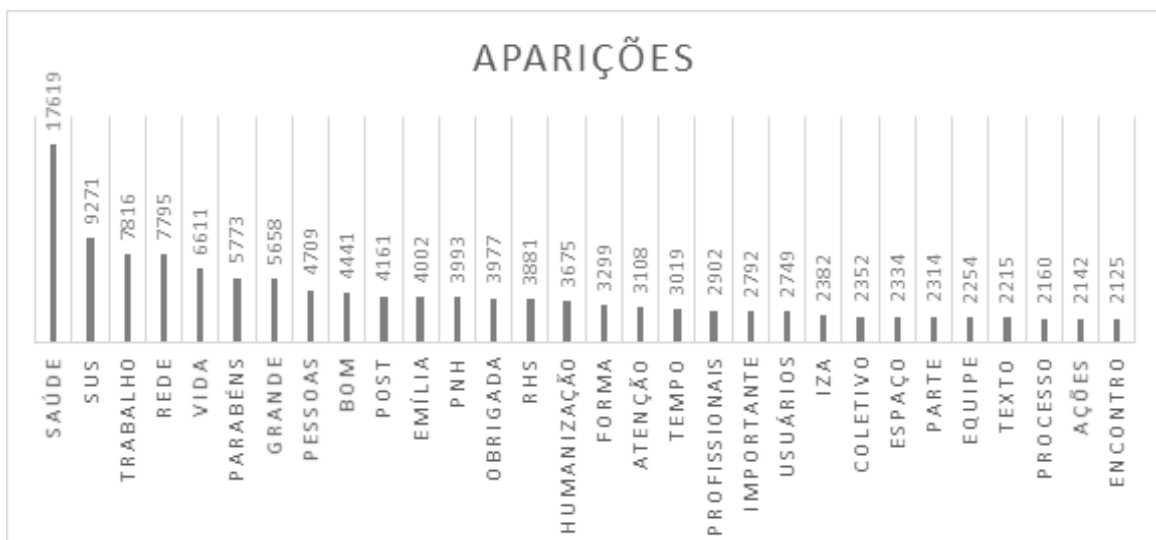
Fonte: Dados da pesquisa, 2018

Apêndice B – Exemplo de remoção de símbolos e caracteres especiais e stopwords

excelente post mariella acho bruno poderão deverão papel fundamental rede orientando-nos redigir post precisaremos certamente manual redação abraço ricardo

Fonte: Dados da pesquisa, 2018

Apêndice C – 30 palavras mais comuns nos comentários



Fonte: Dados da pesquisa, 2018