

IDENTIFICAÇÃO DE PADRÕES DE TWEETS: UM ESTUDO DE CASO SOBRE A COVID-19 NO BRASIL

Bruna Pereira Paz¹, Leila Cristina C. Bergamasco²

^{1,2}Departamento de Ciência da Computação, Centro Universitário FEI

unifbpaz@fei.edu.br leila.cristina@fei.edu.br

Resumo: A partir de março de 2020, a população passou a conviver com a pandemia do Coronavírus (COVID-19) e tem tentado, por meio dos dados gerados encontrar soluções. Um dos aspectos estudados tange à reação social e psicológica da sociedade perante tal fenômeno, que pode ser expressa por meio de redes sociais. Este projeto de pesquisa busca criar a primeira base de dados referente as interações no Twitter sobre a doença do COVID-19 e identificar associações dos conteúdos publicados antes e após o início da vacinação da população brasileira.

1. Introdução

A doença Coronavírus (COVID-19) causada pelo vírus de síndrome severa aguda coronavírus 2 (*severe acute respiratory syndrome coronavirus 2 (SARS-CoV-2)*) foi anunciada oficialmente pela Organização Mundial de Saúde (OMS) em fevereiro de 2020 e em março de 2020, e declarou estado de pandemia devido à doença. Desde então mais de 180 milhões de pessoas foram infectadas e mais de 4 milhões de pessoas vieram à óbito [1]. O Brasil ocupa a 7ª posição entre países com maior mortalidade, com 2542 mortes por milhão de habitante [2].

COVID-19 é a primeira pandemia após o início da transformação digital da sociedade, a partir do século XX, e dessa forma, uma variedade de dados relacionados à doença está sendo disponibilizados ao público e pesquisadores, fomentando diversas aplicações incluindo previsão, planejamento, gestão, tomada de decisão e rastreabilidade da doença [3]. Além de pesquisas focadas em aspectos diretamente associadas à COVID-19, nota-se um interesse sobre as consequências psicológicas e sociais geradas pela doença. Nesse sentido, redes sociais como Facebook, Twitter e Instagram são uma fonte rica de dados. O Twitter possui mais de 350 milhões de usuários ativos e gera por dia mais de 200 bilhões de *tweets* por dia, sendo a 7ª rede social mais utilizada no Brasil.

Observa-se, portanto, um ao alto volume de dados textuais gerados diariamente e oriundo de diversas fontes, incluindo as redes sociais. Sendo necessário o desenvolvimento de soluções computacionais para extrair informações e padrões relevantes de tal massa de dados. Tais técnicas pertencem ao domínio de Mineração de Texto (MT) que, por sua vez, é uma das estratégias de mineração dentro do contexto de Mineração de Dados.

2. Metodologia

Para a construção do projeto serão utilizados conceitos da mineração de Texto, que pode ser definida como sendo a descoberta não-trivial de padrões ocultos ou desconhecidos e potencialmente úteis em textos. Tais padrões podem oferecer informações uteis em contextos

específicos e, diferentemente de outros tipos de dados, os textos muitas vezes são não-estruturados, ambíguos e de difícil processamento [4]. Para trabalhar com os dados serão divididos alguns aspectos da MT em blocos a serem trabalhados durante o projeto (Figura 1).



Figura 1 – Integração da área de MT com diferentes domínios de conhecimento.

Dado que o objetivo do presente projeto de pesquisa é encontrar padrões sobre o comportamento social dos usuários do Twitter referente a pandemia do COVID-19 e, conseqüentemente, fornecer para a comunidade científica uma base de dados inédita, com dados a partir do início da vacinação da população, grande parte dos esforços envolvidos nesse trabalho se darão na etapa de pré-processamento dos *tweets* recuperados. Dessa forma durante o projeto, serão avaliadas as técnicas computacionais pertinentes a cada uma das etapas de pré-processamento como a *tokenização* e remoção de *stopwords*. Adicionalmente o algoritmo Apriori será utilizado para descobrir associações frequentes e válidas no conteúdo analisado.

Será utilizada API do Twitter para recuperar os *tweets* relacionados ao tema do presente projeto de pesquisa. Adicionalmente será utilizada a linguagem Python e suas bibliotecas *numpy*, *nlk* e *apriori* para o desenvolvimento e análise dos resultados.

$$\text{Suporte}(A) = \frac{NA}{NT}$$

$$\text{Confiança}(A \rightarrow B) = \frac{\text{Suporte}(A \cup B)}{\text{Suporte}(A)}$$

Figura 2 – Equações de suporte e confiança.

Para avaliar a qualidade das regras encontradas, serão utilizadas as métricas de suporte e confiança, descritas nas equações acima (Figura 2), respectivamente. O termo *NA* se refere ao número de

