

# Análise de sites disseminadores de *fake news*

Davi P. Guimarães, Guilherme M. Moreira, Matheus E. Fagundes, Nilson M. Lazarin

<sup>1</sup>Bacharelado em Sistemas de Informação – Centro Federal de Educação Tecnológica Celso Suckow da Fonseca (Cefet/RJ) – Nova Friburgo, RJ – Brazil

{davipguimaraes.dev, guilherme.muller.m, matheusefagundes}@gmail.com

nilson.lazarin@cefet-rj.br

**Abstract.** *Fake news contains incorrect or imprecise information that are distributed for some purpose. The impacts resulted by the consumption of fake news can be large, affecting political, social, economic and personal aspects of the population live. Furthermore, there aren't scalable and trusty definitive solutions to identifying fake news. The purpose of this work is to identify characteristics on websites that distribution fake news by analyzing different audit data and content that may be useful by classifiers.*

**Resumo.** *Fake news são notícias que contém informações incorretas ou imprecisas e são disseminadas com algum objetivo. Os impactos causados pelo consumo de notícias falsas podem ser grandes, afetando aspectos políticos, sociais, econômicos e pessoais da vida da população. Apesar disso, não existem soluções definitivas confiáveis e escaláveis de se identificar uma notícias falsas. O objetivo deste artigo é identificar características em sites disseminadores de fake news, através da análise de diferentes dados de auditoria e conteúdo, que possam ser úteis a classificadores de confiabilidade de notícias.*

## 1. Introdução

Caracteriza-se como *fake news* as informações fabricadas que imitam as notícias confiáveis em formato e estrutura mas não possuem o mesmo objetivo nem adotam o mesmo processo editorial [Lazer et al. 2018]. Portais divulgadores de *fake news* não adotam padrões da mídia que garantem a qualidade e confiabilidade da informação e geralmente publicam notícias com algum tipo de objetivo a ser alcançado [Kovic 2018]. Socialmente as *fake news* podem provocar efeitos que vão do inofensivo a extremamente prejudicial. Há notícias falsas que são usadas apenas para captar o clique do leitor e fazê-lo compartilhar para alcançar novos leitores, outras podem ter um propósito mais malicioso como mudar o rumo de uma disputa eleitoral.

O objetivo desse trabalho é identificar características presentes em sites disseminadores de *fake news* que não estejam presentes em sites confiáveis. A identificação desses padrões poderá contribuir com trabalhos de classificação automática. Outra contribuição deste trabalho é a construção e disponibilização de uma base de URLs de páginas contendo *fake news* em português brasileiro. Os critérios adotados para classificar páginas com *fake news*, bem como as técnicas utilizadas para levantar os padrões, são apresentados na seção Metodologia.

## 2. Fundamentação Teórica

Apesar da popularidade do tema, poucos trabalhos científicos que abordam esse assunto foram encontrados. Ao se considerar trabalhos realizados em português e com o foco no ecossistema de notícias brasileiro, a quantidade é ainda menor. Nos artigos encontrados durante a revisão bibliográfica, a abordagem escolhida para a detecção de notícias falsas geralmente utiliza técnicas de análise de conteúdo e texto. Isto acontece pois, segundo [Jaworski et al. 2014], esta é, atualmente, a forma mais eficiente de algoritmos identificarem *fake news*, pois estão menos suscetíveis a manipulação.

É proposto por [Tanaka et al. 2010] um método para avaliar a credibilidade extraindo informações do conteúdo, do suporte social e do autor da página, embora o objetivo do trabalho não foi, necessariamente, detectar fake news. As métricas utilizadas foram: o *Grau Cobertura*, o quanto o conteúdo discorre sobre os tópicos do assunto da página, e o *Grau de Profundidade* do tópico, em que profundidade o assunto da página é tratado pelo texto. É apresentado em [Kakol et al. 2017] um modelo para prever se o conteúdo de uma página web possui credibilidade. A ideia é identificar fatores que possam ser utilizados em futuras automações para classificação da credibilidade. Segundo os pesquisadores, uma das maiores dificuldades em se criar modelos para predição de credibilidade é que, quando um humano faz uma avaliação em site, utiliza muitos critérios subjetivos. A proposta apresentada em [Dong et al. 2015] busca classificar a veracidade do conteúdo de páginas web, baseado na premissa de que uma fonte que tem poucos fatos falsos é considerada confiável. Através de um modelo probabilístico denominado *Knowledge-Based Trust* (KBT), foram analisadas 119 mil páginas web e posteriormente parte dos resultados foi verificada manualmente para avaliar a eficiência da técnica. Foram analisadas manualmente 100 páginas consideradas confiáveis pelo modelo. Destas, 85 foram consideradas confiáveis pela abordagem manual.

Este artigo se difere dos citados em alguns pontos: idioma, objeto de análise e contribuições. No que se refere ao idioma, os trabalhos encontrados lidam apenas com notícias em inglês, não atendendo a necessidade de se analisar notícias e portais em português brasileiro. Já o objeto de análise se difere no fato de que este artigo se propõe a identificar as características dos portais disseminadores de *fake news*, enquanto que os trabalhos citados analisam o conteúdo das notícias. Além disso, outro objetivo deste trabalho é disponibilizar uma base de notícias comprovadamente falsas que sirva de suporte para trabalhos futuros.

## 3. Metodologia de Pesquisa

Para determinar os padrões comuns às páginas de notícias falsas, é necessário antes ter uma lista com URLs com conteúdo previamente caracterizado como falso. Uma vez que este trabalho visa analisar *fake news* em português brasileiro e não foi encontrada nenhuma base de dados pública, o primeiro passo foi construir uma base através do seguinte procedimento: (1) Eleger sites reconhecidos como checadores de notícias ou desmentidores de boatos. Para tal considerou-se a recomendação realizada por [CGI.br 2018] que indica os sites: E-farsas, Boatos.org e Quatro Cantos; (2) Através das postagens nos sites checadores, buscou-se armazenar a URL de onde a *fake news* em questão foi divulgada.

Caso a URL não fosse divulgada no site checador, uma consulta em um buscador era realizada com o texto da *fake news*. Nos links retornados pelo buscador, checou-se o

conteúdo para observar se tratava-se da mesma notícia avaliada pelo site checador e se a fonte encontrada divulgava o conteúdo como verdadeiro; Caso ambas observações fossem verdadeiras a URL era adicionada a base de dados.

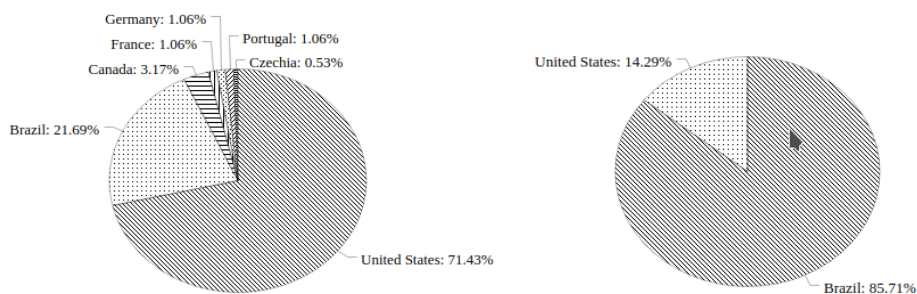
Por fim, foi necessário obter uma base de dados de notícias consideradas verdadeiras. Esta base foi construída através da API RESTful *News API*<sup>1</sup> que permite consultar URL de notícias de várias fontes. Foram usadas 7 fontes<sup>2</sup> de notícia brasileiras disponibilizadas pela API. Essas fontes foram escolhidas por possuírem em suas páginas informações sobre as empresas de comunicação a qual são associadas, além de terem boa circulação.

A segunda etapa consistiu em analisar todas as páginas web apontadas na base de URLs, tanto de notícias comprovadamente falsas, quanto de notícias consideradas verdadeiras. Buscou-se comparar as características dos serviços de hospedagem, sendo elas: Localização geográfica (obtida através do endereço IP); Sistema operacional do servidor; Linguagem de programação utilizada e servidor web. Além disso, analisou-se o conteúdo das páginas em busca do uso excessivo de determinadas palavras, conforme recomendação de [CGI.br 2018] para identificação de boatos.

#### 4. Resultados

A base obtida na fase de busca e catalogação está disponível em *urlfakelist*<sup>3</sup>. Esta base contém 256 URLs de páginas com notícias confirmadas como *fake news*, através dos sites checadores. Durante a fase de análise dos sites, foram observadas algumas características que podem ajudar na distinção de portais de notícias e portais disseminadores de *fake news*, sendo elas:

- **País de hospedagem:** Pode-se observar na figura 1(a) que 78.3% dos disseminadores de *fake news* estão hospedados fora do Brasil, preferencialmente nos Estados Unidos e que 85.7% dos portais de notícias estão hospedados em território nacional, conforme figura 1(b).



(a) Disseminadores de fake news

(b) Portais de Notícias

**Figura 1. País de hospedagem**

<sup>1</sup><https://newsapi.org/>

<sup>2</sup>São elas: [extra.globo.com](http://extra.globo.com), [www.infomoney.com.br](http://www.infomoney.com.br), [oglobo.globo.com](http://oglobo.globo.com), [blogs.oglobo.globo.com](http://blogs.oglobo.globo.com), [g1.globo.com](http://g1.globo.com), [kogut.oglobo.globo.com](http://kogut.oglobo.globo.com) e [globoesporte.globo.com](http://globoesporte.globo.com)

<sup>3</sup><https://github.com/nilsonmori/urlfakelist>

- **Palavras alarmistas:** Verificou-se que sites disseminadores de *fake news* possuem em média 4.5 palavras alarmistas, enquanto sites de notícias tendem a não utilizar essas palavras em suas páginas. As palavras alarmistas utilizadas nesta pesquisa foram *ameaçar, atenção, colabore, compartilhe, contatos, corja, divulgue, enganação, enganado, enganar, espalhem, farsa, grave, gravíssimo, perigo, repassem, sacanagem, urgente e vergonha*.
- **Servidor web:** Sites disseminadores de *fake news* utilizam diversos servidores web, dos quais destacam-se *Google Servlet Engine 28%, nginx 19.2%, cloudflare 15.2%, e Apache 14.8%*. Portais de notícias ocultaram essa informação em 99.4% dos casos.

## 5. Conclusão

Este trabalho disponibilizou uma base com URLs de páginas contendo notícias em português do Brasil classificadas como *fake news* por sites checadores. Algumas características da hospedagem de sites e do conteúdo de páginas foram exploradas na segunda fase. Observou-se que em sites disseminadores de *fake news* a ocorrência de palavras alarmistas é recorrente. Na maioria das vezes esses sites estão hospedados fora do Brasil e utilizam diversas tecnologias de hospedagem.

Utilizando a base disponibilizada, trabalhos futuros podem analisar outras características de notícias falsas apontadas por [CGI.br 2018], tais como: Uso excessivo de determinadas tags HTML; Erros ortográficos; Ausência de data da notícia.

## Referências

- [CGI.br 2018] CGI.br (2018). Cartilha de Segurança para Internet: Fascículo Boatos. <https://cartilha.cert.br/fasciculos/boatos/fasciculo-boatos.pdf>.
- [Dong et al. 2015] Dong, X. L., Gabrilovich, E., Murphy, K., Dang, V., Horn, W., Lugaresi, C., Sun, S., and Zhang, W. (2015). Knowledge-based trust: Estimating the trustworthiness of web sources. *Proc. VLDB Endow.*, 8(9):938–949.
- [Jaworski et al. 2014] Jaworski, W., Rejmund, E., and Wierzbicki, A. (2014). Credibility microscope: relating web page credibility evaluations to their textual content. *IEEE/WIC/ACM International Joint Conferences on Web Intelligence (WI) and Intelligent Agent Technologies (IAT)*.
- [Kakol et al. 2017] Kakol, M., Nielek, R., and Wierzbicki, A. (2017). Understanding and predicting web content credibility using the content credibility corpus. *Information Processing and Management*, 53(5):1043 – 1061.
- [Kovic 2018] Kovic, M. (2018). A typology of fake news.
- [Lazer et al. 2018] Lazer, D. M. J., Baum, M. A., Benkler, Y., Berinsky, A. J., Greenhill, K. M., Menczer, F., Metzger, M. J., Nyhan, B., Pennycook, G., Rothschild, D., Schudson, M., Sloman, S. A., Sunstein, C. R., Thorson, E. A., Watts, D. J., and Zittrain, J. L. (2018). The science of fake news. *Science*, 359(6380):1094–1096.
- [Tanaka et al. 2010] Tanaka, K., Ohshima, H., Jatowt, A., Nakamura, S., Yamamoto, Y., Sumiya, K., Lee, R., Kitayama, D., Yumoto, T., Kawai, Y., Zhang, J., Nakajima, S., and Inagaki, Y. (2010). Evaluating credibility of web information. In *ICUIMC*.