
Análise descritiva das *fake news* da saúde através de mineração de textos no Portal da Saúde¹

Larissa Machado VIEIRA²
Núbia Rosa da SILVA³
Douglas Farias CORDEIRO⁴

RESUMO

As chamadas *fake news* tem causado uma sensação de insegurança informativa, trazendo um cenário de vulnerabilidade em várias esferas de partilha de informações, desde questões políticas, financeiras, e até mesmo no âmbito da saúde. Neste contexto, o Ministério da Saúde desenvolveu uma iniciativa de enfrentamento às *fake news*, chamado “Saúde sem *fake news*”, o qual prevê a manutenção de um portal de divulgação e investigação sobre notícias potencialmente inverídicas sobre saúde, de modo a auxiliar a população na obtenção de informações pautadas em evidências científicas. Diante disso, esse trabalho se propõe a realizar uma análise descritiva acerca do conteúdo das notícias publicadas no Portal da Saúde, através da utilização de soluções baseadas em mineração de textos, no sentido de geração insumos e informações para discussões sobre os impactos das fakes news na área da saúde.

PALAVRAS-CHAVE: *fake news*; saúde; mineração de dados; análise.

1. Introdução

As mudanças comunicacionais que ocorreram ao longo da história da humanidade alteraram as estruturas sociais e criaram novas ambiências, promovendo um cenário onde “o desenvolvimento da mídia vem entrelaçado de modo fundamental com as principais transformações institucionais que modelaram o mundo moderno” (THOMPSON, 1998, p. 9). Nesse ambiente de constante evolução tecnológica, com a popularização de dispositivos de Tecnologia da Informação e do Conhecimento (TICs), o acesso à informação tornou-se algo

¹ Trabalho apresentado na DT 6 – Interfaces Comunicacionais do XXI Congresso de Ciências da Comunicação na Região Centro-Oeste, realizado de 22 a 24 de maio de 2019.

² Mestranda em Comunicação, Universidade Federal de Goiás (UFG), vieira.mlarissa@gmail.com

³ Doutora em Ciência da Computação e Matemática Computacional, USP. Professora adjunta da Unidade Acadêmica Especial de Biotecnologia (UFG), nubia@ufg.br

⁴ Doutor em Ciência da Computação e Matemática Computacional, USP. Professor adjunto da Faculdade de Informação e Comunicação (UFG), cordeiro@ufg.br

mais simplificado, tendo seus usos determinados pelas concepções, costumes e necessidades dos indivíduos que usufruem dos dispositivos (CASTELLS, 2005).

Esses fatores acabam por emergir um cenário onde cada indivíduo, enquanto unidade de um conjunto, passa a ser parte de um espaço de constante partilha, permeado pelo estabelecimento de teias de relações, através das redes sociais, que por sua vez, “[...] alteram os modos de ver e ler, as formas de reunir-se, falar e escrever, de amar e saber-se amado à distância, ou, talvez, imaginá-lo” (CANCLINI, 2008, p. 54). Entretanto, é fundamental destacar que, embora a velocidade e facilidade no acesso à informação possa representar uma série de vantagens sob os mais variados aspectos, também pode gerar problemas característicos e específicos dos ambientes virtuais, entre os quais, se destaca a insegurança informativa (MORETZSOHN, 2017).

Neste ambiente, a propagação de notícias falsas (*fake news*) torna-se um problema crônico, sendo comum sua disseminação em larga escala, espalhadas mundialmente nas redes sociais, implicando em uma problemática alarmante: o declínio da importância da verdade (KAKUTANI, 2018). O autor ainda traz um panorama sobre o grave problema das *fake news*, afirmando que:

[...] esse cenário vem sendo exponencialmente acelerado pelas redes sociais, que conectam usuários que pensam da mesma forma e os abastecem com notícias personalizadas que reforçam suas ideias preconcebidas, permitindo que eles vivam em bolhas, ambientes cada vez mais fechados e sem comunicação com o exterior (KAKUTANI, 2018, p. 16 e 17).

Na atualidade, observa-se que esses conteúdos enganosos, disseminados no mundo virtual e informações alarmantes sem cunho científico comprobatório, têm influenciado determinados grupos em relação a questões de saúde pública, como por exemplo, no que se refere à necessidade de vacinação. Isto gera insegurança e incerteza em alguns países acerca da vacinação dos filhos e da sua própria, e, por conseguinte, aumenta a vulnerabilidade da população a doenças.

No Brasil, durante o terceiro trimestre de 2018, de acordo com PSafe (2018)⁵, no 5º Relatório de Segurança Digital, relativo ao terceiro trimestre de 2018, 46,3% das *fake news* detectadas abordaram o tema política, seguido pelo tema saúde, em segundo lugar, com 41,6% das identificações realizadas. Cabe ressaltar que houve um avanço significativo na detecção de conteúdos falsos relativos à saúde do segundo para o terceiro trimestre do ano a que se refere o relatório. Enquanto no segundo trimestre o tema saúde estava em 4º lugar no nível de identificação, com 19,1% de detecções, no terceiro trimestre ele subiu para 2º lugar (41,6%), sinalizando um aumento de 22,5%. A partir disso, pode-se inferir que houve uma ampliação da propagação de *fake news* na esfera da saúde, e, por conseguinte, sua detecção. O relatório ainda apresenta que as três ferramentas onde as *fake news* foram mais disseminadas, nesta ordem, são o aplicativo de mensagens WhatsApp, os navegadores de internet e o Facebook.

A propagação de *fake news* no âmbito da saúde possui uma série de implicações, as quais acabam por demandar uma grande atenção e relevância na detecção e tratamento dessas notícias falsas, englobando questões que se referem à compreensão dos aspectos sociais, comunicacionais e informacionais envolvidos, os quais permeiam as motivações que levam um indivíduo a disseminar *fake news*. Por outro lado, esses aspectos também estão relacionados a questões puramente de cunho tecnológico, relativas às facilidades que as TICs proporcionam na propagação de *fake news*, ou mesmo na aplicação de soluções automatizadas para detecção ou análise.

Nessa conjuntura, importa ressaltar a imprescindibilidade de que estudos desta natureza aprofundem-se, para que políticas públicas de combate às *fake news* sejam pensadas e construídas, haja vista os riscos decorrentes da disseminação desses conteúdos ardilosos, os quais apontaremos alguns nesta pesquisa.

Neste contexto, no âmbito do presente trabalho, será abordada a aplicação de soluções baseadas em mineração de dados, mais especificamente em mineração de textos, para a realização de análises descritivas sobre *fake news* na área da saúde sobre o conjunto de

⁵ Relatório da segurança digital no Brasil: terceiro trimestre - 2018, disponível em <<https://www.psafe.com/dfndr-lab/pt-br/relatorio-da-seguranca-digital/>>, acessado em 21/03/2019

publicações do portal “Saúde sem Fake News”. O objetivo é, através da aplicação de técnicas específicas, gerar informações que possam servir de insumo no que tange ao processo de investigação do cenário das *fake news* na área da saúde brasileira, através de identificação das relações semânticas entre termos de maior relevância, e identificação das classes dos temas abordados.

2. Fake news

Fake news, cuja tradução significa “notícias falsas”, são conteúdos virais, propositalmente maliciosos e feitos com a intenção de enganar (BAKIR; MCSTAY, 2017, TANDOC; LIM; LING, 2017). Conforme Ferrari (2018, p. 29), as *fake news* podem ser descritas como “uma variedade de desinformações que podem variar entre a correta utilização de dados manipulados, a utilização errada de dados verdadeiros, a incorreta utilização de dados falsos e outras combinações possíveis”.

No contexto da problematização acerca do termo *fake news*, é importante ressaltar a existência de possíveis objeções inerentes à sua tradução literal. Segundo as teorias da comunicação, uma notícia jamais é considerada falsa; ela não representa com total exatidão a realidade, mas é uma construção feita a partir dela (TRAQUINA, 2005). Diante dessa percepção, o fato de uma representação não ser fidedigna em sua totalidade, não implica necessariamente em qualificação de uma notícia como falsa. Por outro lado, a disseminação proposital e maliciosa de informações falsas é compreendida e considerada como notícia falsa.

D’Ancona (2017) apresenta um panorama histórico sobre as *fake news*, apontando dois marcos importantes para as pesquisas desta natureza: o primeiro refere-se às eleições norte-americanas de 2016, que culminaram com a designação de Donald Trump para ocupar o alto do executivo dos Estados Unidos. Outro fator que alavancou as discussões sobre este termo foi o processo de saída da Grã-Bretanha da União Europeia, denominado *Brexit*, onde foi observado um movimento constante de inverdades disseminadas à população por meio de *fake news*.

Ainda segundo D’Ancona (2017), a eleição do presidente Trump e as informações falsas relacionadas ao *Brexit* não são a causa, mas uma consequência preocupante do valor declinante da verdade, visto que inúmeras inverdades contribuíram para que os resultados desses processos fossem alcançados como foram. O autor relata que:

Donald Trump depreciou a suposição de que o líder do mundo livre deve ter ao menos uma familiaridade oblíqua com a verdade: de acordo com o site PolitiFact, que checa informações e é ganhador do Prêmio Pulitzer, 69% das declarações de Trump são “predominantemente falsas”, “falsas” ou “mentirosas”. No Reino Unido, a campanha a favor da saída da União Europeia triunfou com *slogans* que eram comprovadamente não verdadeiros ou enganosos, mas também comprovadamente ressonantes (D’ANCONA, 2018, p. 20).

Nesse âmbito, surge um termo estritamente ligado às *fake news*: a pós-verdade, vocábulo que significa, conforme definição do *Oxford Living Dictionaries*⁶, “circunstâncias em que os fatos objetivos são menos influentes em formar a opinião pública do que apelos à emoção e à crença pessoal”, revelando as *fake news* como um produto da pós-verdade.

Neste cenário, o Ministério da Saúde lançou em 2018 o programa “Saúde sem Fake News”, que tem como objetivo confrontar notícias falsas sobre saúde disseminadas na internet. Foi aberto um canal via WhatsApp para que os internautas enviem aos encarregados do programa as informações que circulam sobre o tema e que causam dúvidas nos usuários. Assim, os jornalistas selecionam os conteúdos recebidos e repassam aos responsáveis técnicos pela apuração das informações. Após a checagem sobre se há comprovação científica acerca dos tópicos de saúde, os jornalistas apresentam no Portal da Saúde a confirmação da veracidade ou não daquele conteúdo, inserindo o selo “ISTO É *FAKE NEWS!*” ou “ESTA NOTÍCIA É VERDADEIRA”, a depender do resultado da verificação.

Gilberto Occhi, ex-ministro da Saúde, considera que a baixa cobertura vacinal contra a gripe, que não atingiu sua meta em 2018, seja em virtude das *fake news* disseminadas em

⁶ Palavra eleita do ano de 2016 pelo *Oxford Living Dictionaries*: pós-verdade. Disponível em: <https://en.oxforddictionaries.com/word-of-the-year/word-of-the-year-2016>. Acesso em 07 jan. 2018.

torno do assunto⁷. Segundo Occhi, a meta era vacinar 90% do público-alvo, porém, o período da campanha terminou com apenas 78% de cobertura vacinal realizada. Importa alertar que não foram apresentados pelo Ministério da Saúde estudos e dados concretos que estabeleçam a relação entre o menor índice de vacinação e a disseminação de notícias falsas, porém, apresenta-se aqui a suspeita que o órgão levantou na véspera do lançamento do programa “Saúde sem Fake News”⁸.

No tocante a questões de saúde pública, existem alguns grupos (como os *anti-vaxxers*, grupos antivacinação) que fazem apologia a existência de uma pseudociência, onde a ciência é vista como um campo conspiratório, em contraposição ao seu caráter investigativo (D’ANCONA, 2017), que é extremamente necessário para o desenvolvimento de pesquisas em torno de doenças que, se prevenidas e controladas, melhorarão demasiadamente a qualidade de vida dos indivíduos.

Ressalta-se que não é o objetivo deste artigo conhecer os critérios que movem os jornalistas em suas rotinas produtivas na escolha dos conteúdos recebidos, ou como é feita a análise dos conteúdos pela equipe técnica e quem a compõe, recortes estes que poderão ser abordados em outra pesquisa. Neste trabalho, serão analisados os conteúdos que compreendem os seis primeiros meses de atuação do programa, quais sejam, de Agosto de 2018 a Fevereiro de 2019.

3. Metodologia

O presente trabalho trata da construção de uma análise exploratória com base na utilização de soluções de mineração de texto sobre conjuntos de dados relacionados a *fake news* no âmbito da saúde. Para tanto, será considerada como amostra de pesquisa um conjunto de conteúdos disponibilizados pela equipe do programa “Saúde sem *Fake News*”, do

⁷Ministro culpa *fake news* por não cumprimento de meta em vacinação contra gripe. Disponível em: <https://jovempan.uol.com.br/programas/jornal-da-manha/ministro-culpa-fake-news-por-nao-cumprimento-de-meta-em-vacinacao-contragripe.html>. Acesso em 29 de março de 2019.

⁸ Saúde sem *fake news* - <http://portalms.saude.gov.br/fakenews>

Ministério da Saúde do Brasil, datados entre Agosto de 2018 e Fevereiro de 2019, os quais contemplam todos os dados disponibilizados no portal do programa.

A metodologia do trabalho é baseada no processo conhecido como KDD (*Knowledge Discovery in Databases* ou Descoberta de Conhecimento nas Bases de Dados), proposto por Fayyad *et al.* (1996). O processo trata de um conjunto de etapas sequenciais, as quais devem ser realizadas e avaliadas, de forma a gerar resultados concisos e úteis aos propósitos demandados, no âmbito da geração de informação e descoberta de conhecimento. As cinco etapas que compõe o KDD são: extração de dados, tratamento de dados, padronização de dados, mineração de dados, avaliação da informação. Esse modelo possui a vantagem de ser suficientemente flexível para ser aplicado nas mais diversas áreas de investigação, como é o caso da comunicação, provendo mecanismos para geração de análises sobre grandes conjuntos de dados, de forma assertiva e inovadora.

A primeira etapa a ser realizada consiste da obtenção dos dados da amostra considerada. Os dados a serem utilizados se encontram disponibilizados através do portal do programa “Saúde sem *Fake News*”. É importante destacar que o portal não possui um mecanismo de extração dos dados, sendo necessária a utilização de uma abordagem alternativa que permita a execução de tal ação diretamente nas páginas onde se encontram os conjuntos de dados textuais. As páginas são construídas utilizando linguagem HTML (do inglês, *Hypertext Markup Language*). Para extração dos dados foi desenvolvida uma solução baseada em *Web Scraping*. De modo geral, *Web Scraping* pode ser descrito como um conjunto de técnicas de extração de dados em sites, onde é utilizada a estrutura sintática dos códigos HTML para detecção, seleção e coleta de dados de interesse. Importa salientar que a extração de dados, no presente caso, consiste de duas fases, na primeira fase são obtidos os endereços para as páginas que contém, individualmente, os conteúdos sobre as *fake news*, os quais são armazenados em um arquivo estruturado no padrão CSV (do inglês, *comma separated values*). Na segunda fase, cada uma das páginas é acessada, extraindo o conteúdo de interesse para um arquivo textual, conforme é apresentado na Figura 1.

A segunda etapa do processo KDD refere-se ao tratamento dos dados. Nesta fase é necessário realizar a remoção de caracteres especiais, os quais podem representar problema

durante a aplicação dos métodos de análise adotados. Diante disso, foi desenvolvida uma rotina automatizada, utilizando a linguagem Python, que realiza a tarefa de forma automatizada.



Figura 1 - Processo de coleta de dados.

Fonte: autores.

A partir da obtenção dos dados textuais referentes a cada uma das páginas com conteúdo sobre *fake news*, é então necessário converter para o padrão processável pela solução abordada. No presente trabalho são utilizadas soluções de análise disponibilizadas através do software Iramuteq⁹, que utiliza formato textual com os dados identificados através de rótulos únicos, identificados por uma sequência de quatro asteriscos, como mostrado em recorte do corpus textual apresentado na Figura 2.

A solução utilizada através do software Iramuteq, implementa alguns tipos de análises baseadas em métodos estatísticos e de mineração de textos. No presente trabalho serão exploradas análises voltadas à geração de contagem de frequência de termos, grafo de similitude, e detecção de similaridade através do método de classificação ALCESTE, proposto por Reinert (1990) . O cálculo da frequência dos termos permite uma visualização dos termos de maior relevância no corpus textual analisado, sendo insumo para análises relacionadas aos temas abordados em cada elemento textual. Por outro lado, a visualização através do grafo de similitude provê informações sobre o relacionamento semântico entre os

⁹ <http://www.iramuteq.org/>

termos. Finalmente, o método ALCESTE permite a identificação das classes presentes no corpus textual. Através de tais resultados, é possível construir um panorama das *fake news* no âmbito da saúde, assim como os assuntos específicos aos quais abordam, conforme será discutido na próxima seção.

```
**** *noticia4
Furar dedos com agulha ajuda a salvar pessoa com AVC - É FAKE NEWS!
Olá! Essa mensagem é falsa! Não compartilhe. Não há nenhuma comprovação científica que furar os dedos de uma pessoa a ajudaria em caso de Acidente Vascular Cerebral. O tratamento do AVC é feito nos Centros de Atendimento de Urgência, que são os estabelecimentos hospitalares que desempenham o papel de referência para atendimento aos pacientes com AVC. Essas unidades de saúde disponibilizam e realizam o procedimento com o uso de trombolítico, conforme Protocolo Clínico e Diretrizes Terapêuticas (PCDT) específico. IMPORTANTE: O AVC é uma doença que é totalmente dependente do tempo. Isso quer dizer que quanto mais rápido for o tratamento, maiores serão as chances de recuperação completa. Desta forma, torna-se primordial a identificação dos sinais e sintomas e o atendimento médico imediato. Para saber mais, acesse: saude.gov.br/avc

**** *noticia5
Repelente de insetos causa reação química - É FAKE NEWS!
Olá! Essa mensagem é falsa, não compartilhe! Conforme esclarecimentos prestados pela fabricante do produto, a partir de testes de eficácia e segurança do produto, o uso do repelente de insetos é seguro e não causa reação química quando aplicado na pele.

**** *noticia6
Paracetamol pode diminuir eficácia de vacinas em crianças - É VERDADE!
Olá, essa mensagem é verdadeira e a informação inclusive consta na publicação "Manual de Normas e Procedimentos para Vacinação". Em estudos observou-se que as crianças que receberam paracetamol profilático apresentaram uma redução nos títulos de anticorpos das vacinas administradas. Por isso, de uma forma geral, não é indicado o uso de paracetamol antes ou imediatamente após a vacinação para não interferir na imunogenicidade da vacina.
```

Figura 2 - Recorte de corpus textual gerado.

Fonte: autores.

4. Resultados e Conclusão

Através da aplicação do método de coleta desenvolvido foi extraído todo o conjunto de publicações do portal “Saúde sem *Fake News*”, resultando em um total de oitenta notícias, categorizadas como elementos de um corpo textual único, o qual serviu de entrada para as rotinas aplicadas através do software Iramuteq, seguindo o processo de tratamento descrito na metodologia deste trabalho. Diante disso, primeiramente foi realizada a contagem de frequência dos termos, através da qual foi gerada uma nuvem de palavras (Figura 3), que revela os termos de maior relevância dentro de todo o conjunto de elementos textuais. Com base nesse resultado, é possível observar alguns termos de maior sensibilidade no sentido de identificação dos possíveis temas abordados, tais como: “câncer”, “vacina”, “alimento”, e “tratamento”. É importante destacar que alguns termos específicos, tais como “saúde”,

quais existe uma proximidade maior entre as classes 3 e 4, e separadamente, entre as classes 2 e 5. Esse cálculo permite ainda identificar as notícias que são mais similares entre si, como pode ser observado na Figura 6.

Finalmente, a Figura 7, apresenta a distribuição de termos para cada uma das classes. Através dos resultados é possível notar que a Classe 1, de forma ligeiramente isolada, engloba notícias que possuem seu conteúdo mais fortemente ligado ao próprio enfoque da notícia. A Classe 2 aborda notícias com tema focado em câncer. A Classe 3 trata de notícias que tem como conteúdo principal a vacinação. A Classe 4 refere-se ao tema medicamento. E por fim, a Classe 5 é relativa a notícias sob o tema alimentação. Desta maneira, tais resultados permitem inferir que os conteúdos abordados no programa “Saúde sem *Fake News*” estão relacionados, principalmente, aos temas: câncer, vacinação, alimentação e medicamentos. A Tabela 1 apresenta os termos de maior relevância para cada uma das classes.

Classe	Termos principais
Classe 1	mensagem, notícia, hospital, boato, texto, informação, alerta, combate.
Classe 2	câncer, tratamento, risco, relação, infecção, mama, evidência, causa, lesão.
Classe 3	vacinação, vacina, criança, recomendação, população, país, dose, reação.
Classe 4	programa, registro, medicamento, imunização, eficácia, qualidade, vigilância.
Classe 5	alimento, alimentação, consumo, atividade, prevenção, dieta, natura, vitamina.

Tabela 1. Termos de maior relevância nas classes detectadas.

Os resultados obtidos através da aplicação das soluções explorados permitem tecer um retrato do panorama sobre as *fake news* da saúde no Brasil. Tais resultados podem servir de insumo para o desenvolvimento de investigações posteriores sob diversas vertentes, focadas, inclusive, nos temas que apresentaram maior destaque, a fim de compreender, por exemplo, as motivações e implicações das *fake news* em cada uma destas áreas.

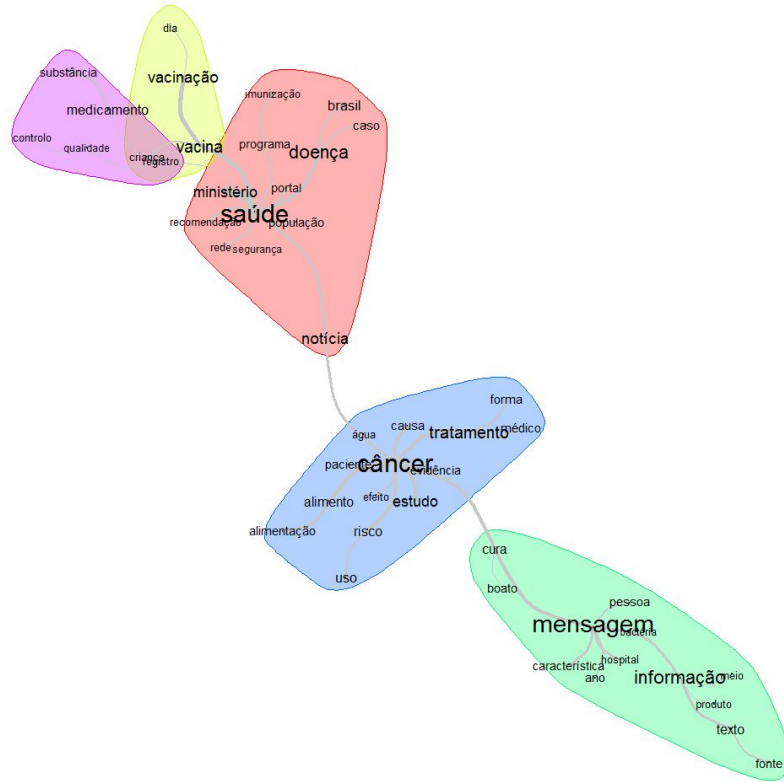


Figura 4. Grafo de similitude.

Fonte: autores.

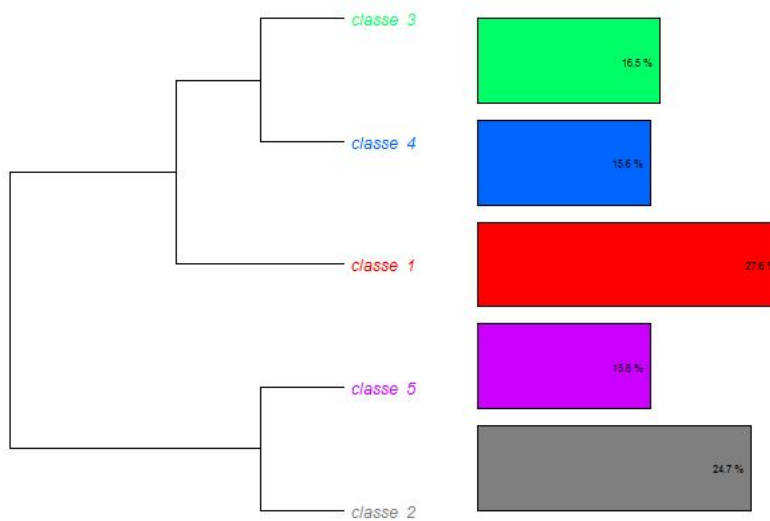


Figura 5. Dendrograma de classes.

Fonte: autores.

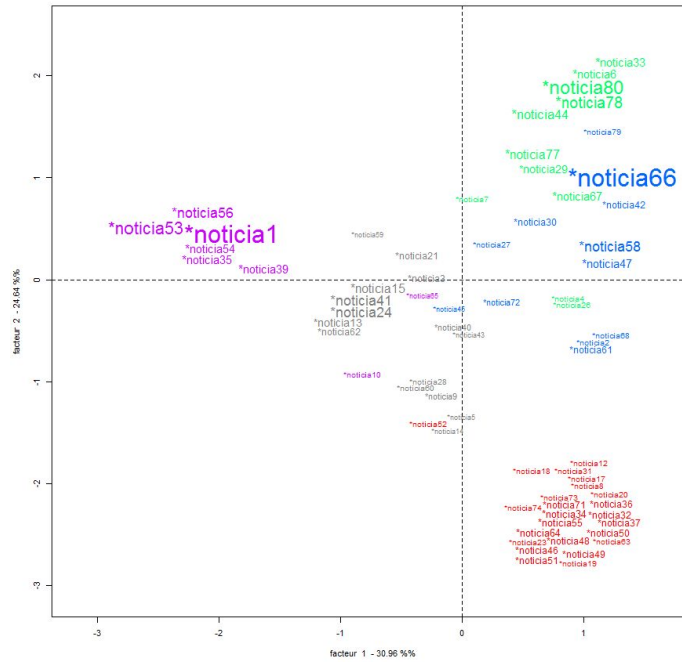


Figura 6. Distribuição dos elementos textuais por similaridade.

Fonte: autores.

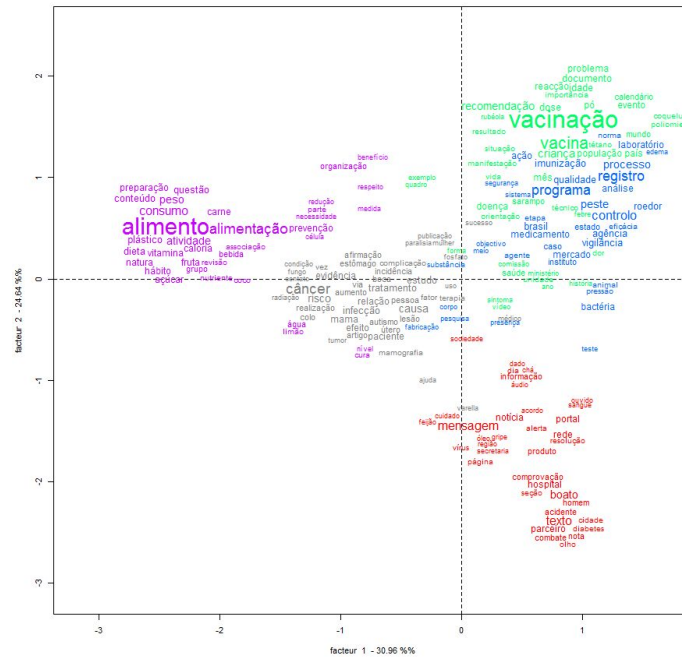


Figura 7. Distribuição de termos com base na ocorrência por classes detectadas.

Fonte: autores

Referências

CAMILO, C.O.; SILVA, J.C.D. **Mineração de dados: conceitos, tarefas, métodos e ferramentas**. Relatório Técnico. Goiânia: Universidade Federal de Goiás, 2009.

CANCLINI, Néstor García. **Leitores, espectadores e internautas**. São Paulo: Iluminuras, 2008.

CASTELLS, M.; CARDOSO, G.(orgs.). **A sociedade em rede: do conhecimento à ação política**. Brasília: Imprensa Nacional, 2005.

D'ANCONA, Matthew. **Pós-verdade: a nova guerra contra os fatos em tempos de fake news**. São Paulo: Faro Editorial, 2018.

FAYYAD, U., Piatetsky-Shapiro, G., and Smyth, P. From Data Mining to Knowledge Discovery in Databases. **American Association for Artificial Intelligence**, v. 17, n. 3, p. 37-54, 1996

FERRARI, Pollyana. **Como sair das bolhas**. São Paulo: EDUC, 2018.

KAKUTANI, M. **A morte da verdade: notas sobre a mentira na era Trump**. Rio de Janeiro: Intrínseca, 2018.

MORETZSOHN, S, D. Uma legião de imbecis: hiperinformação, alienação e o fetichismo da tecnologia libertária. **Liinc em Revista**, v. 13, n. 2, p. 294-306, 2017.

REINERT, M. (1990). ALCESTE, une méthodologie d'analyse des données textuelles et une application: Aurélia de G. de Nerval. **Bulletin de méthodologie sociologique**, (28) 24-54.

TANDOC, Edson C.; LIM, Zheng Wei; LING, Richard. Defining “Fake News”: a typology of scholarly definitions. **Digital Journalism**, v. 6, n. 2, p. 137-153, 2017.

THOMPSON, J. B. **A mídia e a modernidade: uma teoria social da mídia**. 5 ed. Petrópolis: Vozes, 1998.

TRAQUINA, Nelson. **Teorias do Jornalismo I: porque as notícias são como são?** Florianópolis, Insular, 2005.