

Análise de Propagandas Eleitorais Antecipadas no Twitter

Samuel Guimarães¹, Márcio Silva², Josemar Caetano¹, Marcelo Araújo¹, Jonatas Santos¹,
Julio C. S. Reis³, Ana P. C. Silva¹, Fabrício Benevenuto¹, Jussara M. Almeida¹

¹Depto. Ciência da Computação, Universidade Federal de Minas Gerais (UFMG)

²Faculd. Ciência da Computação, Universidade Federal de Mato Grosso do Sul (UFMS)

³Depto. Informática, Universidade Federal de Viçosa (UFV)

{samuelsg, jonatashds}@ufmg.br, marcio@facom.ufms.br, jreis@ufv.br,
{josemarcaetano, marceloaraujo, ana.coutosilva, fabricio,
jussara}@dcc.ufmg.br

Abstract. *Electoral advertising are an essential part of an election. The popularization of online social networks has offered a promising way for candidates to communicate with the electorate at large. In fact, the use of these applications to share electoral ads has already being pointed out, even outside the period allowed by Brazilian law. Yet, fighting this practice is hampered by the lack of a broader knowledge about the characteristics of this type of content, which allows effective detection solutions. This study aims to contribute to this knowledge through a broad characterization of the textual content associated with a set of early electoral advertisements shared on Twitter in pre-election periods associated with recent elections in Brazil (2016, 2018 and 2020). Our main findings are that ads tend to have a negative or neutral sentiment, a certain constant structure and more than half tend to explicitly mention a candidate or party to be chosen or avoided.*

Resumo. *Propagandas eleitorais são parte essencial de uma eleição. A popularização das redes sociais online ofereceu um meio promissor para que candidatos se comuniquem com o eleitorado em larga escala. De fato, já foi apontado o uso destas aplicações para divulgar propagandas eleitorais, inclusive fora do período permitido pela legislação brasileira (propagandas antecipadas). Porém, o combate desta prática é frustrado pela ausência de um conhecimento mais amplo das características deste tipo de conteúdo, permitindo soluções de detecção eficazes. Este estudo visa contribuir para tal conhecimento através da ampla caracterização do conteúdo textual de um conjunto de propagandas eleitorais antecipadas compartilhadas no Twitter em períodos pré-eleitorais associados a eleições recentes no Brasil (2016, 2018 e 2020). Como principal conclusão, observou-se que as propagandas tendem a ter sentimento negativo ou neutro, certa estrutura constante e mais da metade tendem a citar explicitamente um candidato ou partido a ser escolhido ou evitado.*

1. Introdução

Tradicionalmente, as campanhas eleitorais sempre foram realizadas através do horário político gratuito em mídias tradicionais ou presencialmente. Nesse cenário, para garantir o cumprimento da lei eleitoral na utilização de meios de comunicação em massa basta

monitorar o que é veiculado em mídias como rádio e televisão, por exemplo. Porém, a popularização das plataformas de redes sociais online ofereceu um novo e promissor mecanismo de comunicação em larga escala para que candidatos alcancem seu eleitorado. De fato, trabalhos anteriores já mostraram evidências do uso destas plataformas para a disseminação de propagandas políticas (Silva et al. 2021), incluindo compartilhamentos fora do período eleitoral estabelecido pelo Tribunal Superior Eleitoral (TSE), prática vedada pela legislação eleitoral vigente¹ e conhecida como *propaganda eleitoral antecipada*. Entretanto, o combate desta prática é prejudicado pela ausência de ferramentas digitais e métodos para detecção destas propagandas nas plataformas de mídia social, levantando sérias preocupações sobre como mitigar este problema (Sosnovik and Goga 2021). Por outro lado, o projeto de tais métodos precisa ser baseado em um conhecimento amplo e fundamental sobre características marcantes deste tipo de conteúdo, conhecimento este ainda muito limitado na literatura disponível (Silva et al. 2020).

Neste contexto, o objetivo deste estudo é contribuir para tal conhecimento apresentando uma caracterização de possíveis propagandas eleitorais antecipadas identificadas em dados públicos coletados do Twitter, uma das redes sociais mais utilizadas para discussão política online. Busca-se responder à seguinte pergunta: *Quais são as características do conteúdo textual associado às propagandas eleitorais antecipadas compartilhadas no Twitter em diferentes períodos pré-eleitorais?*

Para responder esta pergunta, o estudo realizado cobriu dados (tweets) coletados de três períodos distintos, associados a eleições diferentes no Brasil (2016, 2018 e 2020). Para cada período, foram avaliados vários atributos associados ao *conteúdo textual* de mensagens contendo propagandas eleitorais antecipadas, incluindo atributos associados a propriedades semânticas e sintáticas. O foco no Twitter se deve à grande presença de atores políticos nesta rede, bem como a lacuna de trabalhos que caracterizam este tipo de conteúdo em bases de dados representativas, diferentemente, por exemplo, de dados provenientes do Facebook (Silva et al. 2020).

Nossos resultados mostram que as propagandas antecipadas normalmente possuem um sentimento negativo ou neutro, com o sentimento neutro aumentando ao longo dos anos e possuem um padrão recorrente em relação ao uso de *hashtags* e *links*, além de menções à outros perfis (*arrobas*). As diferenças encontradas são geralmente devido ao contexto eleitoral. Além disso, propagandas fazem mais referências a atributos sintáticos, como pronomes e preposições. Em relação a atributos semânticos, destacamos atributos relacionados a decisões (*votar, eleger*) e a espaço (*cidade, próximo*). Por fim, em torno da metade das propagandas eleitorais referenciam alguma entidade (partidos, locais, pessoas) ou outros perfis no Twitter. Nossos resultados podem direcionar a coleta de dados com base nos atributos textuais encontrados, auxiliar o monitoramento ativo de propagandas, através da detecção de *links* ou *hashtags* de interesse e, por fim, permitir a caracterização textual de campanhas com *bots*.

O restante deste artigo está organizado como segue. A seção 2 apresenta trabalhos relacionados e a seção 3 descreve sucintamente a metodologia adotada para a caracterização. A seção 4 detalha os resultados nas diversas dimensões analisadas, e, por fim, a seção 5 apresenta as conclusões e aponta possíveis direções de trabalhos futuros.

¹Lei Eleitoral nº 13.488, 6 de outubro de 2017. http://www.planalto.gov.br/ccivil_03/_ato2015-2018/2017/lei/L13488.htm

2. Trabalhos Relacionados

Vários trabalhos abordaram a detecção de mensagens com discurso político em redes sociais online, seja usando conjunto de palavras-chave, seja pelo uso de aprendizado de máquina. No primeiro grupo, *hashtags* foram usadas para caracterizar *tweets* com conteúdo político em (Conover et al. 2011), enquanto os autores de (Grimaldi 2019) usaram termos em geral para criação de um léxico visando filtrar *tweets* com informações referentes aos candidatos da eleição espanhola de 2019. Outros trabalhos abordaram a identificação automática de mensagens com conteúdo político, utilizando métodos de aprendizado de máquina (Oliveira et al. 2020; Silva et al. 2020). Em (Oliveira et al. 2020), por exemplo, *tweets* de congressistas foram usados para detectar automaticamente mensagens com conteúdo político, visando estimar o quanto esses perfis comentam sobre política. Já em (Silva et al. 2020), os autores rotularam propagandas da biblioteca pública de anúncios do Facebook para treinar um classificador para detectar propagandas eleitorais. Por fim, em (Sosnovik and Goga 2021) os autores usaram propagandas políticas no Facebook para demonstrar a complexidade da tarefa de identificar este tipo de conteúdo.

Até onde sabemos, a maioria dos trabalhos anteriores analisou mensagens com conteúdo político no geral. Os únicos que abordaram *propagandas eleitorais*, foco do presente estudo, foram (Silva et al. 2020; Sosnovik and Goga 2021). Porém, diferentemente deste estudo, esses trabalhos analisaram propagandas políticas no Facebook ou focaram em conteúdo em inglês. Este estudo dá continuidade a um trabalho anterior (Silva et al. 2021), indo além de análises superficiais de propagandas eleitorais antecipadas. A contribuição inédita é a caracterização ampla de propriedades textuais de propagandas políticas, em português, no Twitter, uma plataforma frequentemente usada para o discurso político.

Outras pesquisas, que tangenciam este trabalho, visam avaliar como discussões em redes sociais podem influenciar a opinião pública. Por exemplo, em (Ribeiro et al. 2019b), os autores analisaram como anúncios russos no Facebook influenciaram a eleição presidencial americana de 2016. Outros esforços abordaram a monitoração de campanhas políticas no Facebook (Kim et al. 2018), a caracterização de interações entre usuários no contexto de movimentos sociais e políticos específicos (Kou et al. 2017; Caetano et al. 2017) e correlações entre atividades online e opiniões ou intenção de voto em candidatos específicos (Maruyama et al. 2014; Ribeiro et al. 2019a). Esses estudos são ortogonais ao trabalho aqui apresentado, já que não abordam propagandas eleitorais especificamente.

3. Metodologia

3.1. Coleta e Tratamento dos Dados

O conjunto de dados analisado compreende uma amostra de *tweets* publicados em períodos associados às três últimas eleições brasileiras em 2016, 2018 e 2020. A coleta de dados foi realizada utilizando a API histórica do Twitter², que fornece acesso a mensagens compartilhadas desde 2006. Como o nosso foco é na análise de propagandas eleitorais antecipadas, os *tweets* abrangem os períodos pré-eleitorais das referidas eleições.

Para buscar *tweets* de interesse, utilizamos um dicionário de palavras-chave, sugeridas por especialistas de instituições fiscalizadoras das eleições brasileiras, que indicam pedidos de voto, implícitos ou explícitos. Exemplos de palavras-chave são “*vote em mim*”,

²<https://developer.twitter.com/en/docs/twitter-api/search-overview>

Tabela 1. Coleções de dados rotulados.

Coleção de Dados	Período	#Propaganda Eleitoral	#Não Propaganda	Total de Mensagens
Twitter 2016	1/1 - 15/08/2016	147 (22,21%)	515 (77,79%)	662
Twitter 2018	1/1 - 15/08/2018	188 (37,38%)	315 (62,62%)	503
Twitter 2020	1/1 - 25/09/2020	786 (45,15%)	955 (54,85%)	1741

“*vote no candidato*” e “*conto com seu voto*”. Também utilizamos uma *blacklist*³ de termos buscando filtrar mensagens sobre votações *online* de programas de TV, como o BBB.

As mensagens textuais coletadas passaram por um processo de rotulação manual de modo a identificar aquelas com conteúdo que pode ser caracterizado como propaganda eleitoral. Especificamente, cada mensagem coletada foi manualmente rotulada (como contendo propaganda eleitoral ou não) por três voluntários, seguindo instruções de especialistas sobre como identificar este tipo de conteúdo. Foram rotuladas como propagandas mensagens contendo pedidos de voto para si mesmo, pedido de voto para outra pessoa, propaganda eleitoral contra um adversário ou propaganda contra um partido rival. Mensagens com falas sobre eleição sem pedido de voto, mensagens com conteúdo sem teor político, ataques pessoais ou elogios a políticos ou partidos, e questões políticas de outros países foram rotuladas como não propaganda. A confiabilidade da rotulação foi avaliada usando o percentual de concordância entre os voluntários e o coeficiente de Cohen’s Kappa (κ) (Landis and Koch 1977). De forma geral, observou-se uma concordância entre pares de voluntários superior a 75%, com κ variando entre 0,48 e 0,58 (dependendo do par de voluntários). A concordância total entre os três voluntários apresentou um Fleiss’ Kappa κ de 0,53. Conforme (Landis and Koch 1977), estes valores de κ sugerem uma concordância “moderada”, o que é aceitável considerando que em muitas mensagens a propaganda eleitoral não aparece de forma explícita, dando margem para múltiplas interpretações, como em “Não VOTE em quem aprovou esse absurdo! LINK”, que cita um político indiretamente. O rótulo final de cada mensagem (propaganda eleitoral ou não propaganda) foi dado pela maioria dos três voluntários. A Tabela 1 apresenta uma breve descrição dos dados coletados⁴, incluindo período de coleta, bem como números e percentuais de mensagens contendo propagandas eleitorais e não propagandas.

3.2. Propriedades de Conteúdo Analisados

A pergunta que direciona este estudo é: “*Quais são as características textuais das propagandas eleitorais antecipadas no Twitter?*”. Para respondê-la, foram analisados atributos do conteúdo (textual) associado às mensagens rotuladas como propagandas coletadas⁵. Estes atributos estão associados a cinco dimensões: (i) sentimento; (ii) palavras mais informativas, isto é, que melhor caracterizam as propagandas eleitorais compartilhadas; (iii) estruturas mais comuns; (iv) atributos psicolinguísticos; e (v) entidades nomeadas. As diversas análises realizadas se complementam da seguinte forma. A análise de sentimentos permite verificar o tom da corrida eleitoral, revelando se o discurso usado tende a ser mais positivo, negativo ou neutro. Já a análise das palavras mais informativas mostram o conteúdo característico das propagandas antecipadas, medindo o poder semântico de uma

³As palavras-chave e a *blacklist* estão disponíveis em www.dcc.ufmg.br/~samuel.guimaraes/Brasnam2022

⁴Os dados anonimizados poderão ser disponibilizados mediante solicitação por e-mail para o primeiro autor do artigo.

⁵Os dados caracterizados estão descritos na coluna 2 (#Propaganda Eleitoral) da Tabela 1.

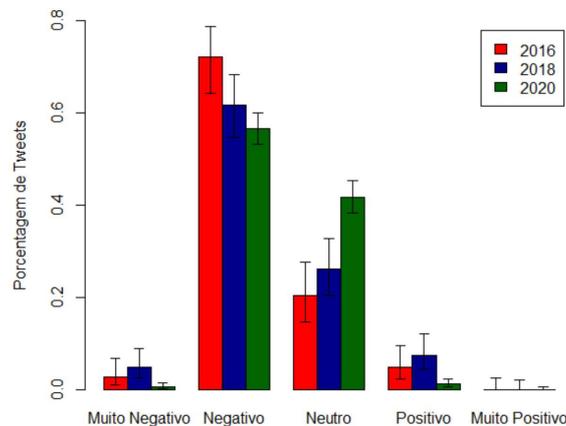


Figura 1. Histograma do sentimento das mensagens com propaganda eleitoral.

palavra em relação ao texto sendo analisado, considerando as palavras individualmente, sem avaliar estruturas textuais completas. Já a análise de estruturas mais comuns evidencia padrões de pedido de voto, bem como de propagandas negativas. Atributos psicolinguísticos das palavras são usados para agrupá-las em categorias, permitindo encontrar padrões mais complexos. E, por fim, a análise de entidades permite obter a frequência de referências a organizações, locais, pessoas e perfis no Twitter, revelando como o contexto político de cada eleição influencia no padrão de discurso destas propagandas.

4. Resultados e Análises

Nesta seção, apresentamos e discutimos o conjunto de análises realizadas que visam entender caracterizar as propagandas eleitorais antecipadas no Twitter. Para cada análise, descrevemos a metodologia utilizada, seguida das principais conclusões alcançadas.

4.1. Análise de Sentimentos

Primeiramente, analisamos o sentimento expresso nas propagandas eleitorais presentes nos conjuntos de dados estudados. A ferramenta SentiStrength (Thelwall et al. 2010), amplamente usada para esta tarefa na literatura (Caetano et al. 2017), foi usada nesta análise. O SentiStrength fornece uma pontuação inteira variando de -4 a +4, de fortemente negativo a fortemente positivo, com 0 indicando um sentimento neutro. Para facilitar a interpretação, optou-se pela discretização dos valores nas seguintes classes: muito negativo (-4 e -3), negativo (-1 e -2), neutro (0), positivo (1 e 2) e muito positivo (3 e 4). A Figura 1 mostra o histograma com o percentual de mensagens em cada classe (com os intervalos de confiança (Newcombe 1998)), para cada um dos três períodos analisados.

Como vemos na figura, os sentimentos negativos e neutros são muito mais frequentes nas propagandas compartilhadas nos três períodos analisados, com uma tendência de diminuição dos sentimentos negativos e aumento do sentimento neutro nas eleições mais recentes (2018 e 2020). Apesar do crescimento dos sentimentos neutros, propagandas negativas são a maioria, com mais de 50% de todas as mensagens com esse sentimento nos três períodos. Exemplos de propagandas com sentimento negativo, neutro e positivo são, respectivamente, “*NOME DO CANDIDATO falso da p*rra, NÃO VOTE NELE LINK*”, “*Vamos trocar todo mundo. Vote em quem NUNCA foi eleito. LINK*” e

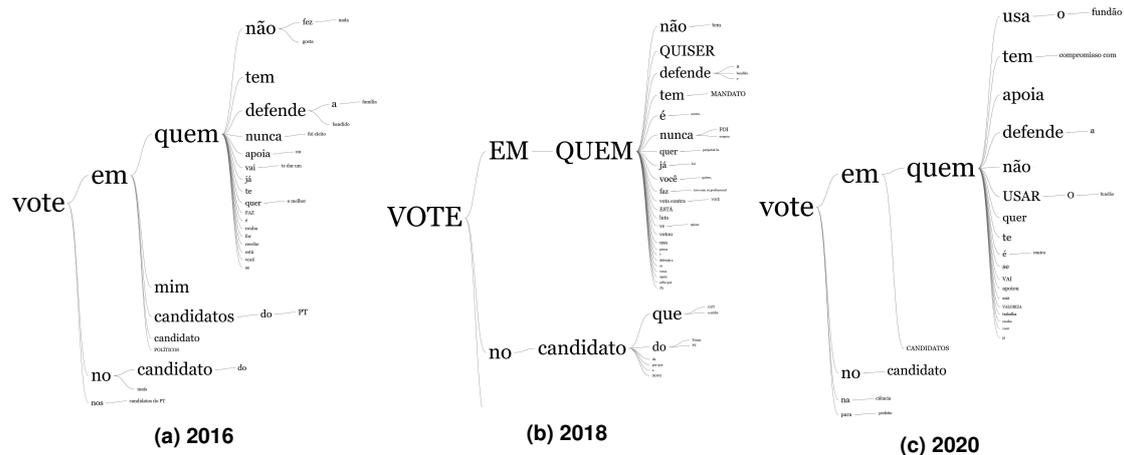


Figura 3. Frases mais comuns nas propagandas eleitorais para cada eleição.

em 2020 (e não em 2016 e 2018). O motivo de 2018 não incluir o nome do atual presidente pode ser explicado pelo fato de que, nos dados coletados, a maioria das propagandas naquele ano, na verdade, ataca sua oposição, o PT, em vez de apoiá-lo diretamente. Confirmando a tendência entre os anos, a similaridade de Jaccard para as top-10 palavras mais informativas para 2016 e 2018 é igual a 0,54 e, tanto para 2016 e 2020 quanto para 2018 e 2020, a similaridade é de 0,43. Embora a similaridade entre as palavras características de 2020 e dos outros dois períodos seja um pouco mais baixa, estes resultados ainda são muito expressivos e indicam que cerca de metade das top-10 palavras mais informativas são as mesmas nos três períodos. Esta observação sugere uma possível estrutura comum das propagandas com links e hashtags sendo usadas para pedir votos. As demais palavras ranqueadas nas dez primeiras posições mudam ao longo dos anos. No entanto, na sua maioria, estas palavras refletem posicionamentos em relação aos candidatos, através de críticas e elogios (*defender*, *apoiar* e *querer*), que também são afetadas pelos contextos de cada período.

4.3. Estruturas Textuais Mais Comuns

Nesta seção, focamos na identificação das frases (i.e., sequências de palavras) que ocorrem com mais frequência nas propagandas eleitorais de cada período analisado. Para tal, foi utilizada a técnica de árvores de palavras (do inglês, *word tree*) que oferece uma visualização intuitiva da estrutura destas frases (Wattenberg and Viégas 2008). Nesta técnica, primeiramente, seleciona-se uma palavra que será considerada como raiz da árvore e, em um processo recursivo, palavras subsequentes que possuem maior frequência são interligadas em uma estrutura de galhos. Ao final, tem-se uma árvore descrevendo o conjunto de frases mais frequentes a partir da raiz. Para esta análise, visando identificar as estruturas textuais mais relevantes e características das propagandas eleitorais, foi escolhida como raiz a palavra com maior informatividade nos três anos, que é a palavra “vote”. As Figuras 3a, 3b e 3c mostram as árvores encontradas para cada período eleitoral.

Nas três árvores apresentadas, as frases mais frequentes iniciam-se por “vote em quem”, indicando a tendência das propagandas antecipadas em recomendar explicitamente candidatos a serem votados, bem como candidatos a serem evitados. Um exemplo de um pedido contra o voto em alguém, é o “não vote em quem usa o fundão”, representada com a omissão da palavra “não” na árvore de 2020. Além disso, tanto em 2016

Tabela 2. Top 5 Frases mais comuns de cada ano com exemplos e frequências.

Frase comum	Exemplo de Propaganda	2016	2018	2020
vote no candidato	Eles se uniram e ferraram o Rio de Janeiro . 3 em 1 -3 PMDB e 1PT. Em 2016 não vote no candidato apoiado põe eles. LINK	5,44%	10,64%	2,80%
vote em quem te	Você é servidor público? Vote em quem te valoriza profissionalmente, não em quem usa a tua profissão politicamente. LINK	6,12%	4,26%	5,73%
vote em quem defende	Minha família vai votar toda em ARROBA. Vote em quem defende a patria e NAO vai deixar roubarem HASHTAG	7,48%	3,19%	2,93%
vote em quem usa	Não vote em quem usa religião pra se eleger !	0,00%	1,60%	10,43%
vote em quem não	Este Partido é Golpista. Não vote em quem não respeita os resultados das urnas. LINK	5,44%	2,66%	3,69%
vote em quem tem	Não vote em quem tem mandato. Decepção total, eles não nos representa. LINK	4,08%	2,66%	2,93%
vote em quem apoia	Não vote em quem apoia Bolsonaro. Não vote em evangélicos.	2,72%	1,06%	5,09%
vote em quem é	Pense bem, nas próximas eleições vote em quem é pró-emprego e empresa, jamais em alguém pró-Estado, lembre de Dilma!	1,36%	3,72%	1,78%
vote em candidatos	Não reeleja ninguém! Não vote em candidatos do PMDB, não vote em quem foi contra o UBER. LINK	4,76%	0,00%	1,27%
vote em quem usar	Não vote em quem usar o fundão partidária LINK	0,00%	0,00%	3,05%

quanto em 2018, ocorreram citações ao PT e pedidos de renovação política, através dos pedidos “não vote em quem tem mandato” e “não vote no(s) candidato(s) do PT”. Semelhante à análise de palavras mais informativas, as posições dos candidatos também são usadas nas propagandas, o que aparece com frequência na frase “votem em quem defende”, com família e bandidos sendo usadas para propagandas positivas e negativas, respectivamente. A Tabela 2 apresenta exemplos das top-5 sequências mais comuns encontradas no nosso conjunto de dados, para cada um dos anos analisados. O conjunto de 10 frases comuns apresentadas foi obtido tomando a união dos top-5 em cada período.

De forma geral, pode-se observar que existe de fato um certo padrão de pedido de votos e recomendação para evitar certos candidatos, em cada eleição, com o contexto de cada período alterando a motivação por trás das propagandas. Em 2016, partidos foram citados com maior frequência. Em 2018, o assunto da Lava Jato e corrupção se torna um foco maior para aconselhar que ninguém vote em um candidato. Por fim, na eleição de 2020, as propagandas são frequentemente contra candidatos, divididas entre grupos anti-bolsonaro e antipetistas, com o fundo partidário sendo também uma pauta relevante.

4.4. Atributos Psicolinguísticos

O *Linguistic Inquiry and Word Count* (LIWC) (Tausczik and Pennebaker 2010) é um dicionário que permite a análise da frequência de uso de palavras em diferentes categorias psicolinguísticas, tanto semânticas quanto sintáticas, ao ser considerada uma determinada frase. Ele é composto por uma hierarquia de categorias, ou atributos, cada uma caracterizada por um conjunto de palavras. Ao avaliar as propagandas eleitorais em termos da frequência das palavras em diferentes atributos do LIWC é possível ter uma visão detalhada e mais aprofundada do seu conteúdo. Esta análise complementa a caracterização de palavras mais informativas bem como de estruturas textuais mais frequentes, apresentada nas seções anteriores, ao enfatizar aspectos relacionados à narrativa utilizada para construir o conteúdo. Cada mensagem com propaganda eleitoral (dos diferentes anos) foi analisada separadamente, de modo a obter a frequência de palavras características de cada atributo do LIWC. As frequências foram então agregadas para todas as mensagens de cada período, para obtenção de valores médios. Ressalta-se que uma mesma palavra pode pertencer a diferentes atributos, representando diferentes significados ou usos da

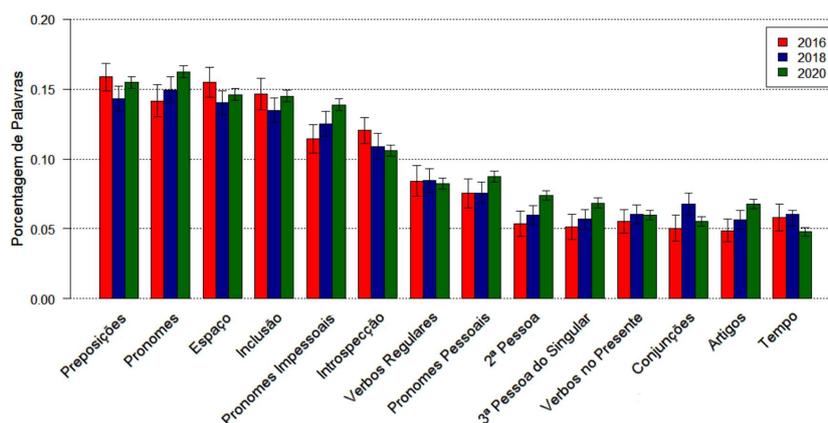


Figura 4. Atributos Psicolinguísticos (LIWC) mais frequentes nas propagandas.

mesma. Desta forma, a soma de todas as frequências pode ultrapassar 100%.

A Figura 4 apresenta a distribuição da frequência média, com os intervalos de confiança, de 14 atributos do LIWC⁶ frequentemente presentes nas propagandas nos três períodos analisados. Os atributos estão ordenados no eixo x do gráfico pela média de ocorrência em todas as mensagens, considerando os três períodos juntamente. Este conjunto de 14 atributos foi obtido tomando a união dos 10 atributos mais frequentes em cada período. Nota-se grande similaridade entre os atributos mais presentes nos três períodos, enfatizando mais uma vez uma certa estrutura sintático-semântica comum às propagandas nos três períodos analisados. Ainda assim, ressaltam-se algumas diferenças importantes, como alguns atributos entre os top-10 mais frequentes em apenas alguns períodos, bem como algumas mudanças no ranking de atributos. Um exemplo é o atributo *verbos no presente* que ocorre apenas no top-10 de 2016. Da mesma forma, o atributo *verbos regulares* e o atributo *pronomes pessoais* tem sua ordem invertida de 2018 para 2020.

Pela figura, nota-se que os atributos mais frequentes são associados a propriedades sintáticas (e.g., preposições e pronomes), e que, em geral, as diferenças entre as frequências observadas entre períodos são bastante sutis (em média $\leq 1\%$) ou sem relevância estatística. A maior diferença estatisticamente significativa (2,45%) é no uso de pronomes impessoais (e.g., *isso*, *aquilo*) entre 2016 e 2020, com uma tendência de aumento mais recentemente. Essas diferenças sutis podem estar relacionadas a mudanças no tipo de discurso usado em cada eleição, resultado similar a análise de árvore de palavras.

A Figura 4 mostra também a presença dos atributos do tipo *introspecção* (do inglês, *insight*), *espaço* e *inclusão*, associados a propriedades semânticas. O atributo *introspecção* é um processo cognitivo caracterizado por palavras relacionadas a decisões como “votar”, “eleger”, “escolher” e “representar”. Estas palavras ocorreram frequentemente em propagandas que recorrem a advertências e verbos no imperativo para incentivar a escolha do voto. Já o atributo *espaço* tem relação com noções de posição, sendo representado por palavras como “cidade”, “envolvido”, “próximo” e “fora”. Por fim, o atributo *inclusão* engloba palavras que descrevem situações ou ações que envolvem um grupo ou que agrupam algo, como, por exemplo, “envolvido”, “ficaremos” e “divulgue”. “Envolvido” é um bom exemplo de palavra em mais de uma categoria, podendo ter relação com

⁶Nomes dos atributos foram traduzidos para o português.

a proximidade entre um grupo de pessoas. Essas palavras de inclusão foram usadas com frequência para descrever como os eleitores deveriam agir, bem como grupos de candidatos atuaram em mandatos anteriores.

4.5. Análise de Entidades

A última análise consiste na caracterização das propagandas em termos da frequência de referências a entidades. Para isso, foi utilizada a técnica de reconhecimento de entidades mencionadas (*named-entity recognition*) (Sang and De Meulder 2003), disponível na biblioteca *Spacy*⁷. Nessa técnica, entidades são categorizadas em termos do domínio ao qual se aplicam. Dado um texto de entrada, é produzida uma lista de entidades identificadas na entrada com suas respectivas categorias. Neste estudo, foram consideradas 3 categorias: LOC para lugares físicos como cidades ou estados, ORG para organizações como partidos e órgãos governamentais e PER para pessoas reais ou personagens, como candidatos ou alguma mascote. Após alguns experimentos iniciais, observou-se que algumas entidades eram categorizadas de forma imprecisa pela ferramenta. Por exemplo, alguns órgãos governamentais foram categorizados como LOC, ao invés de ORG. A categorização foi então manualmente avaliada e eventuais imprecisões corrigidas. Além disso, também foram incluídas na lista de entidades algumas citações a perfis no Twitter. Estas foram associadas a uma nova categoria (AT). Duas análises foram realizadas. Primeiramente, examinamos as porcentagens de tweets que referenciaram alguma entidade e, posteriormente verificamos a distribuição das categorias mencionadas em cada período.

Nos três períodos analisados, em torno da metade das propagandas eleitorais referenciam alguma entidade ou perfil no Twitter, com 2018 tendo um percentual um pouco maior (mais de 60%). Este resultado sugere dois padrões de propagandas: (i) as que fazem referência explícita a candidatos e partidos, e (ii) aquelas sem menções explícitas, mas que sugerem implicitamente um candidato ou partido (para votar ou não) ao descrever suas ideias e/ou opiniões. Para fins de exemplificação, as top-5 entidades mais citadas em cada período são: PT, Brasil, Dilma, Bolsonaro e PMDB (2016); Brasil, PT, Bolsonaro, Temer e Lula (2018); Bolsonaro, Brasil, Fundão, PT e Rodrigo Maia (2020). Nota-se uma grande relação das entidades frequentes com o contexto político de cada eleição, sendo esses fatos possivelmente úteis para explicar parte dos resultados identificados.

Focando nas referências explícitas a entidades, a Figura 5 mostra que organizações (ORG), que incluem partidos políticos e órgãos governamentais, e pessoas (PER) tendem a ser as entidades mais citadas no geral. Locais (LOC) e perfis (AT) vêm em seguida. Mais ainda, conforme mostrado, as organizações eram as entidades mais citadas nas propagandas compartilhadas em 2016, com uma frequência muito maior em comparação com as demais categorias. Exemplos de organizações muito citadas nas propagandas de 2016 são PT e PMDB, dois partidos em maior destaque nas eleições daquele ano. Em 2018, organizações e pessoas passaram a ter proporções mais próximas, com aumento expressivo de citações a perfis do Twitter (AT). Como exemplo, PT, Bolsonaro e Temer foram entidades muito citadas nas propagandas de 2018. Por fim, locais, como Brasil e São Paulo, são as entidades mais citadas nas propagandas em 2020. Mais uma vez, estas diferenças podem ser resultados de aspectos contextuais relacionados às eleições de cada período. Por exemplo, durante as eleições presidenciais de 2018, os nomes dos principais

⁷O modelo *pt_core_news_lg* usado está em <https://github.com/explosion/spaCy>

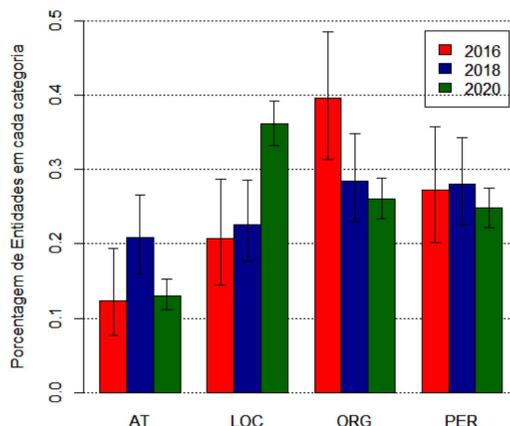


Figura 5. Distribuição das categorias das entidades nas propagandas.

candidatos e partidos políticos ganharam protagonismo. Já nas eleições municipais de 2020, menções a locais específicos ganham uma maior visibilidade.

5. Conclusão

Neste artigo, analisamos os dados de propagandas eleitorais antecipadas compartilhadas no Twitter durante as três últimas eleições brasileiras de 2016, 2018 e 2020. Nossos resultados revelaram que: (i) Os sentimentos mais frequentes são negativos e neutros, com uma tendência de diminuição dos sentimentos negativos e aumento do sentimento neutro ao longo das três eleições; (ii) É comum que propagandas explorem links e hashtags como mecanismo para pedir votos; (iii) Existe uma tendência de recomendar explicitamente candidatos a serem votados, bem como candidatos a serem evitados; (iv) Os atributos mais frequentes nas propagandas são associados a propriedades sintáticas (e.g., preposições e pronomes), e que, em geral, as diferenças entre as frequências observadas entre períodos são bastante sutis (em média $\leq 1\%$) ou sem relevância estatísticas. Essas diferenças sutis podem ser devido a mudanças no tipo de discurso em cada eleição e; (v) Metade das propagandas eleitorais fazem referência à alguma entidade ou perfil no Twitter. Este resultado sugere dois padrões principais de propagandas eleitorais: 1) as que fazem referência explícita a candidatos e partidos, e 2) aquelas sem menções explícitas, mas que sugerem implicitamente um candidato ou partido (para votar ou não). Como trabalhos futuros, esperamos usar a base de dados para criação de ferramentas de identificação de propagandas antecipadas no espaço online, coletar novas propagandas para a eleição de 2022 e possivelmente aplicar análises baseadas em *embeddings*.

Agradecimentos. Este trabalho foi parcialmente financiado pelo Ministério Público de Minas Gerais (projeto Capacidades Analíticas), CNPq, CAPES, FAPEMIG e FAPESP.

Referências

- [Caetano et al. 2017] Caetano, J. A., Lima, H. S., dos Santos, M. F., and Marques-Neto, H. T. (2017). Utilizando análise de sentimentos para definição da homofilia política dos usuários do twitter durante a eleição presidencial americana de 2016. In *Anais do VI BraSNAM*. SBC.
- [Conover et al. 2011] Conover, M., Ratkiewicz, J., Francisco, M., Gonçalves, B., Menczer, F., and Flammini, A. (2011). Political polarization on twitter. In *Proc. of the ICWSM*.

- [Grimaldi 2019] Grimaldi, D. (2019). Can we analyse political discourse using twitter? evidence from spanish 2019 presidential election. *SNAM*, 9(1):1–9.
- [Jaccard 1912] Jaccard, P. (1912). The distribution of the flora in the alpine zone. 1. *New Phytol.*, 11(2):37–50.
- [Kim et al. 2018] Kim, Y. M., Hsu, J., Neiman, D., Kou, C., Bankston, L., Kim, S. Y., Heinrich, R., Baragwanath, R., and Raskutti, G. (2018). The stealth media? groups and targets behind divisive issue campaigns on facebook. *Polit. Commun.*, 35(4):515–541.
- [Kireyev 2009] Kireyev, K. (2009). Semantic-based estimation of term informativeness. In *Proc. of the NAACL HLT*, pages 530–538.
- [Kou et al. 2017] Kou, Y., Kow, Y. M., Gui, X., and Cheng, W. (2017). One social movement, two social media sites: A comparative study of public discourses. *Computer Supported Cooperative Work (CSCW)*, 26(4):807–836.
- [Landis and Koch 1977] Landis, J. R. and Koch, G. G. (1977). The measurement of observer agreement for categorical data. *biometrics*, pages 159–174.
- [Maruyama et al. 2014] Maruyama, M. T., Robertson, S. P., Douglas, S. K., Semaan, B. C., and Faucett, H. A. (2014). Hybrid media consumption: How tweeting during a televised political debate influences the vote decision. In *Proc. of the CSCW*, pages 1422–1432.
- [Newcombe 1998] Newcombe, R. G. (1998). Interval estimation for the difference between independent proportions: comparison of eleven methods. *Stat. Med.*, 17(8):873–890.
- [Oliveira et al. 2020] Oliveira, L. S. D., Vaz-de Melo, P. O. S., Amaral, M. S., and Pinho, J. A. G. (2020). Do politicians talk about politics? assessing online communication patterns of brazilian politicians. *Trans. Soc. Comput.*, 3(4).
- [Ribeiro et al. 2019a] Ribeiro, F. N., Kansaon, D., and Benevenuto, F. (2019a). Leveraging the facebook ads platform for election polling. In *Proc. of the WebMedia*.
- [Ribeiro et al. 2019b] Ribeiro, F. N., Saha, K., Babaei, M., Henrique, L., Messias, J., Benevenuto, F., Goga, O., Gummadi, K. P., and Redmiles, E. M. (2019b). On microtargeting socially divisive ads: A case study of russia-linked ad campaigns on facebook. In *Proc. of the FAT*.
- [Salton and Buckley 1988] Salton, G. and Buckley, C. (1988). Term-weighting approaches in automatic text retrieval. *Inf. Process. Manag.*, 24(5):513–523.
- [Sang and De Meulder 2003] Sang, E. T. K. and De Meulder, F. (2003). Introduction to the conll-2003 shared task: Language-independent named entity recognition. In *Proc. of the CoNLL*, pages 142–145. Morgan Kaufman Publishers.
- [Schubert et al. 2017] Schubert, E., Spitz, A., Weiler, M., Geiß, J., and Gertz, M. (2017). Semantic word clouds with background corpus normalization and t-distributed stochastic neighbor embedding. *arXiv preprint arXiv:1708.03569*.
- [Silva et al. 2021] Silva, M., Guimarães, S., Caetano, J., Araújo, M., Santos, J., Reis, J. C., Silva, A., Benevenuto, F., and Almeida, J. (2021). Propaganda eleitoral antecipada: Uma análise de postagens em mídias sociais. In *Anais do X BraSNAM*. SBC.
- [Silva et al. 2020] Silva, M., Oliveira, L. S. d., Andreou, A., Melo, P. O. V. d., Goga, O., and Benevenuto, F. (2020). Facebook ads monitor: An independent auditing system for political ads on facebook. In *Proc. of the WWW*.
- [Sosnovik and Goga 2021] Sosnovik, V. and Goga, O. (2021). Understanding the complexity of detecting political ads. In *Proc. of the WebConf*, pages 2002–2013.
- [Tausczik and Pennebaker 2010] Tausczik, Y. R. and Pennebaker, J. W. (2010). The psychological meaning of words: Liwc and computerized text analysis methods. *J. Lang. Soc. Psychol.*, 29(1):24–54.
- [Thelwall et al. 2010] Thelwall, M., Buckley, K., Paltoglou, G., Cai, D., and Kappas, A. (2010). Sentiment strength detection in short informal text. *JASIST*, 61(12):2544–2558.
- [Wattenberg and Viégas 2008] Wattenberg, M. and Viégas, F. B. (2008). The word tree, an interactive visual concordance. *IEEE TVCG*, 14(6):1221–1228.