

Análise de Sentimentos em vídeos do YouTube sobre polarização política: uma abordagem híbrida a partir do reconhecimento de entidades

Cláudio Barbosa Silva¹, Daniel Mendes Barbosa¹

¹Instituto de Ciências Exatas e Tecnológicas – Universidade Federal de Viçosa (UFV)
Florestal – MG – Brazil

{claudio.barbosa,danielmendes}@ufv.br

Abstract. *Social networks are increasingly important means of communication, especially in an electoral period. This paper uses this context to perform a study of videos from the YouTube platform in the electoral period. Aspect-level sentiment analysis techniques are applied, using the pre-trained BERTimbau model, by extracting the transcripts and using the main candidates as target entities. It was possible to identify trends in behavior and positioning, in individual videos or channels. Overall, the analysis pointed to a greater presence of candidate Bolsonaro in the YouTube (53.2%), mainly in journalistic channels. The evaluation of the candidates' channels show their approaches and strategies, mostly with mentions of the channel owner. The results also pointed to a trend of neutral manifestations, mainly in news channels, and a large impact of volume of quotes and interactions in Podcasts.*

Resumo. *As redes sociais são meios de comunicação cada vez mais importantes, principalmente em um período eleitoral. O trabalho utiliza esse contexto para realizar um estudo dos vídeos da plataforma YouTube no período eleitoral. São aplicadas técnicas de análise de sentimentos em nível de aspecto, com a utilização do modelo pré-treinado BERTimbau, através da extração das transcrições e utilizando os principais candidatos como entidades-alvo. Foi possível a identificação de tendências de comportamentos e posicionamento, em vídeos individuais ou canais. De maneira geral, a análise apontou para uma maior presença do candidato Bolsonaro no YouTube (53,2%), principalmente em canais jornalísticos. A avaliação dos canais dos candidatos sinalizam suas abordagens e estratégias, majoritariamente com menções ao dono do canal. Os resultados apontaram também para uma tendência de manifestações neutras, principalmente em canais de notícias e um impacto grande de volume de citações e interações em Podcasts.*

Palavras-chave: *Eleições, Análise de Sentimentos, Youtube, BERTimbau, Entidades*

1. Introdução

Verifica-se que o modo como as informações são consumidas no Brasil dá-se de maneira cada vez mais dinâmica e intensa¹, com uma grande quantidade de opções para

¹<https://www.em.com.br/app/noticia/tecnologia/2021/09/28/internatecnologia,1309670/brasil-e-o-terceiro-pais-do-mundo-que-mais-usa-rede-sociais-diz-pesquisa.shtml>

acesso à informação. Uma dessas opções são as redes sociais, um ambiente que muitas vezes é utilizado de maneira livre e com pouca ou quase nenhuma moderação. O estudo feito em [Boxell et al. 2017] relaciona o crescimento da polarização com tais redes, indicando que ultrapassaram a barreira de simples entretenimento para um ecossistema complexo e que nem mesmo a neutralidade das mídias de informação conseguem reduzir [Kobellarz et al. 2021].

A possibilidade de utilização ampla de redes sociais nas campanhas dos candidatos é justificada pelo fato de o Brasil ser o segundo país em que as pessoas mais passam tempo em internet². O tempo médio de utilização das redes pelo brasileiro é de 154 dias por ano, sendo que aproximadamente 56 dias são gastos apenas em redes sociais. Uma das ramificações do uso pelas grandes massas é o consumo e debate de ideias políticas, que abarrotou muitas dessas plataformas com discussões fortemente polarizadas e atraiu várias figuras políticas para essas redes.

Tratando-se de Brasil existem várias mídias e aplicativos que estão presentes no dia-a-dia do brasileiro, sendo as principais: Twitter, Instagram, YouTube, Facebook e até mesmo ferramentas de comunicação como WhatsApp e Telegram, entre várias outras. Uma pesquisa realizada pela Reuters Institute [Newman et al. 2022] mostrou que, no Brasil, a rede social mais utilizada pelos seus usuários para buscar informações é o YouTube com 43% dos usuários, seguida de WhatsApp (41%) e Facebook (40%) . Em um país como o Brasil em que o número de dispositivos com acesso a internet supera a população, ultrapassando os 12,2 bilhões em 2021³, esses dados apontam para a relevância dessa rede social, o que implica também em uma importância política considerável.

As redes sociais oferecem a políticos uma caracterização de conexão permanente, ideia explorada por [Larsson 2016] em uma pesquisa realizada utilizando dados do Facebook. A interação constante é a alternativa mais explorada, principalmente em períodos eleitorais em que o tempo de exposição nas chamadas mídias convencionais é limitado, definido pelo Tribunal Superior Eleitoral (TSE)⁴. Em 2022 ocorreram as eleições presidenciais no Brasil, eleições estas marcadas por uma polarização entre dois candidatos que representam espectros políticos distintos: Luís Inácio Lula da Silva do Partido dos Trabalhadores (PT) e Jair Messias Bolsonaro, candidato pelo Partido Liberal (PL).

Baseando-se nesse contexto, este trabalho propõe uma Análise de Sentimentos considerando os níveis de aspecto e de sentença com o objetivo de identificar a orientação de vídeos de YouTube com base em suas transcrições. A análise em nível de aspecto será feita para classificar os sentimentos sobre cada entidade em um ou mais contextos dentro de um mesmo vídeo, utilizando modelos de computação para identificar e classificar comportamentos, tendências e estratégias de campanha e de publicações no YouTube. Os vídeos publicados nos canais dos candidatos são avaliados individualmente, além de uma análise de vídeos resultantes de buscas por palavras-chave definidas para o contexto.

Seguindo uma linha de estudos de análises de sentimentos em redes sociais e os fatos já mencionados anteriormente, para a obtenção dos resultados este trabalho busca

²<https://www.sortlist.com/blog/your-digital-year/>

³<https://www.poder360.com.br/tecnologia/dispositivos-iot-devem-chegar-a-27-bilhoes-ate-2025>

⁴<https://www.cnnbrasil.com.br/politica/tse-confirma-tempo-de-propaganda-de-candidatos-a-presidencia-na-tv-e-no-radio/>

aplicar um uso híbrido de categorização de Processamento de Linguagem Natural (PLN), identificando as entidades presentes nas transcrições do YouTube tendo como base um modelo capaz de reconhecer entidades, gerado a partir de um modelo pré-treinado em Português chamado BERTimbau. Busca ainda aplicar a classificação sobre as sentenças específicas em que cada entidade é mencionada, caracterizando assim uma Análise de Sentimentos em nível de Sentença e Aspecto utilizando as definições apresentadas em [Liu et al. 2012].

Este artigo está organizado da seguinte forma: na seção 2 são mencionados alguns dos trabalhos relacionados a esse; a seção 3 abrange os principais conceitos teóricos utilizados; a seção seguinte 4, apresenta os materiais e métodos utilizadas no desenvolvimento deste conteúdo; os resultados e suas implicações serão descritos na seção 5; na seção 6, finalizando o trabalho, constarão as considerações finais e possibilidades para trabalhos futuros.

2. Trabalhos Relacionados

São descritos nessa seção alguns dos trabalhos que contemplam o assunto abordado em nosso estudo, sendo que o assunto principal é a literatura que aborda a análise de sentimentos em nível de aspecto no ambiente das redes sociais. O conteúdo dos vídeos ainda é uma fonte de dados pouco estudada, devido à dificuldades encontradas no tratamento das transcrições, visto que muitas das vezes são geradas automaticamente e alguns vídeos nem possuem essa informação.

Ainda assim, existem trabalhos aplicados no contexto brasileiro como o estudo realizado por [Silva and Barbosa 2019], que investigou o comportamento dos usuários em canais do YouTube. Também em um contexto político/eleitoral, os comentários foram analisados buscando determinar a polaridade e traçar uma posterior correlação entre os resultados obtidos e a aprovação dos candidatos.

Uma revisão sistemática realizada em [Kubin and von Sikorski 2021] aponta um interesse crescente em examinar a polarização em redes sociais e a necessidade de obter métricas e ferramentas capazes de medir esse fenômeno. O estudo apresentou também duas características sobre a produção na área. Uma delas é que o Twitter é a rede social mais utilizada nas análises e a outra, que reforça a importância de nosso estudo, que a maioria das produções avaliam apenas o cenário americano. Em [Shah et al. 2021], imagens extraídas do Twitter e Facebook também são analisadas e classificadas. Esse tipo de classificação também possui relevância, visto que a capilaridade de um conteúdo divulgado em massa (virais), com linguagem acessível, também compõem o cenário de disputa política.

Alguns outros estudos apontam para a importância de análises em ambientes políticos, como em [Rodríguez-Ibáñez et al. 2021] que analisou *tweets* nas eleições espanholas de 2019, [Robles et al. 2022] que analisou o impacto de bots e polarização política no debate sobre a COVID-19 e [Yarchi et al. 2021] que apresentou um cruzamento de dados sobre polarização em mídias sociais. Em [Buder et al. 2021], é apontada a relação direta entre polarização e a presença de discursos e ambientes negativos, corroborando com a ideia de que quanto maior a interação com discursos negativos, maior e mais radical será essa polarização.

Alguns estudos guiaram os autores do presente trabalho para o uso do BERTim-

baud, um modelo pré-treinado que utiliza a arquitetura Transformers, sendo a principal ferramenta utilizada nesse estudo. Em [Silva and Freitas 2022] ele foi utilizado para identificar a presença de discurso de ódio em postagens no Twitter, obtendo resultados superiores a outros métodos como NB, SVM, LSTM, entre outros. O estudo apresentado em [Souza et al. 2019] originou o BERTimbau e apresentou como resultado uma performance melhor que o *Multilingual BERT* ocupando o *status* de estado da arte em reconhecimentos de entidades, na língua portuguesa.

3. Fundamentação Teórica

Nessa seção tratam-se os conceitos e fundamentos responsáveis pela condução do presente trabalho. Partindo de alguns dos conceitos mais fundamentais de análise de sentimentos e sua aplicação em redes sociais, este trabalho traz uma abordagem híbrida do modelo de análise de sentimentos, realizando análises em nível de aspecto e sentença e combinando o resultado obtido por elas. A organização de como e em que momento cada atividade e conteúdo foi realizado estão descritos na seção 4.

3.1. Análise de Sentimentos em nível de aspecto

Observando o cenário político brasileiro, uma análise e acompanhamento das redes sociais de maneira quantitativa e qualitativa, torna-se um instrumento interessante também para esse contexto de disputas políticas. Para realizar tal tarefa foi escolhida a análise de sentimento em nível de aspecto, buscando identificar as entidades em seu contexto e, para atingir esse objetivo, elaborou-se uma biblioteca que permita a coleta e tratamento de dados e uma análise de qual o sentimento expresso em diferentes momentos sobre entidades pré-determinadas.

Como dito, este trabalho utiliza o contexto polarizado para entidades nomeadas da classe "Pessoa", correspondentes aos principais candidatos. Nesse grau de especificidade uma sentença pode conter diversos sentimentos relacionados com diferentes entidades e/ou a uma única entidade [Liu et al. 2012]. Para proporcionar uma análise direcionada para o contexto desejado, foi elaborada uma abordagem híbrida que utiliza conceitos de análise léxica baseadas em posicionamento presentes no trabalho feito por [Gu et al. 2018].

O estudo apresentado por [Salas-Zárate et al. 2017] indica que, utilizando-se trigramas em uma abordagem com o método *N-gram around* obteve-se uma maior precisão na classificação de sentimentos. A utilização dessa abordagem foi detalhada na seção 4, que apresenta ainda uma ferramenta que utiliza a arquitetura *Transformer*. Essa arquitetura foi proposta por [Vaswani et al. 2017] e é tida como o estado da arte se tratando de modelos bidirecionais e processamento de linguagem natural (PLN).

3.2. Reconhecimento de Entidades Nomeadas

O reconhecimento de Entidades Nomeadas ou *Named Entity Recognition* (NER) é uma área do Processamento de Linguagem Natural (PLN) que é responsável pela identificação e classificação de entidades nomeadas em um determinado texto ou sentença [Sharnagat 2014]. A principal utilização do NER é a identificação de elementos baseados em entidades pré-definidas, daí o termo nomeadas.

Podem existir diversos grupos de entidades, o que implica que um grande volume de dados deve ser utilizado para treinamento de um modelo efetivo, tendo em vista o

incremento da dificuldade de reconhecimento com tal aumento. Essa relação com as entidades inclui o pré-processamento dos dados, responsável por categorizar de acordo com as entidades definidas. A pesquisa de [Marrero et al. 2013] aponta também para as diversas divergências e falácias existentes no que diz respeito a entidades nomeadas, tratando o problema como não resolvido e a necessidade de profundidade conceitual e em ferramentas. O presente trabalho utilizou as tecnologias consideradas estado da arte e as definições mais difundidas na academia.

De acordo com [Mohit 2014], a definição das Entidades Nomeadas (EN) passa por palavras que são categorizadas em domínios e contextos específicos, normalmente carregando informações que dizem respeito a algo, em específico e que servirão como alvo para o sistema de PLN utilizado. Ou seja, uma única entidade pode englobar diversas palavras. As mais comuns são pessoas, organizações, valores, locais, entre outras e ela é identificada de acordo com sua consistência e menções no texto analisado. Contudo, como possuímos entidades chave já definidas (*features*) para a classe "Pessoa" e, por meio de pré-processamento de dados, as entidades alvo ficam restritas àquelas que fazem referência aos candidatos (nome e apelidos).

4. Materiais e métodos

Nessa seção estão descritos os materiais utilizados e os métodos que foram seguidos para a execução deste trabalho. Os procedimentos descritos a seguir foram aplicados como pesquisa experimental com a análise de sentimentos e transcrições de vídeos do YouTube dos principais expoentes políticos da eleição de 2022 como objeto de estudo.

Para a busca das transcrições avaliadas foram selecionadas palavras-chave capazes de captar o tema escolhido e possibilitar uma análise do comportamento e de como as entidades em destaque pela polarização estão sendo mencionadas. A divisão aplicada nessa seção será: (i) base de dados, (ii) pré-processamento dos dados, (iii) análises utilizando o BERTimbau, (iv) treinamento e refino dos modelos e, por último, (v) as estratégias de normalização e validação dos resultados.

4.1. Base de Dados

A base de dados considera um contexto de polarização em uma disputa presidencial com direcionamento para o conteúdo em que os principais candidatos estão presentes. A partir dessa seção são tratados como Lula (para o candidato Luís Inácio da Silva) e Bolsonaro (para o candidato Jair Messias Bolsonaro). Os dados foram coletados entre agosto e novembro de 2022, ou seja, compreendem a janela eleitoral anterior ao primeiro turno e após a conclusão do segundo turno.

O ponto inicial de obtenção de dados é a definição de palavras-chave utilizadas como parâmetros, buscando encontrar vídeos que estivessem no contexto das eleições e que a menção a algum ou ambos candidatos ocorre. Tais palavras escolhidas foram: "eleições presidenciais 2022", "lula", "bolsonaro" e "segundo turno 2022". Observou-se um distanciamento entre a palavra-chave utilizada para uma análise de contexto e a entidade contida nos vídeos, visto que a máquina de busca utilizada pelo YouTube não contempla o que foi dito no vídeo e sim o que consta nos metadados, como títulos e descrições. Buscando um direcionamento para um assunto que incita polarização [Buder et al. 2021] e manifestação explícita de sentimentos dos interlocutores (jornalistas, comentaristas, autores dos vídeos), também foram coletadas transcrições utilizando

Tabela 1. Vídeos obtidos na busca

Termo	Resultados da busca	Válidos
Lula	656	479
Bolsonaro	631	513
Eleições presidenciais 2022	689	563
Segundo turno 2022	676	609
Corrupção	640	572
Total	3292	2736

Tabela 2. Vídeos obtidos para os canais de Lula e Bolsonaro.

Canal	Nº vídeos
Lula	1576
Bolsonaro	2823

a palavra-chave "corrupção". A tabela 1 indica a quantidade de vídeos coletados e os vídeos em que foi possível realizar a extração das transcrições. A diferença dá-se pela não existência de transcrições para todos os vídeos, já que é uma propriedade adicional.

Para a coleta e tratamento de dados foi importante a criação de uma biblioteca que agrupasse: coleta dos identificadores(ids) dos vídeos, coleta e tratamento das transcrições, processamento de atividades e coleta de detalhes dos vídeos. A linguagem escolhida para esse desenvolvimento foi Python⁵, visto que é a que oferece as melhores ferramentas para as tarefas necessárias (análises e tratamento de dados). O YouTube fornece acesso a uma API (YouTube Data API⁶) para extração de dados por desenvolvedores, assim ela foi integrada à biblioteca criada.

Cada pesquisa retornou um grupo de identificadores únicos para os vídeos que são necessários para a coleta das transcrições. Tal informação não é fornecida pela API própria do YouTube, o que levou à utilização de uma API de terceiros: YouTube Transcript API⁷. Para a utilização do BERTimbau, no reconhecimento das entidades, foi realizado um breve pré-processamento, responsável por simplificar expressões e termos que remetiam aos candidatos e possivelmente poderiam não ser reconhecidas como entidades. Os nomes dos candidatos e alguns apelidos foram substituídos pelos rótulos "Lula" ou "Bolsonaro".

A coleta de dados sofreu impacto devido a qualidade ou inexistência das transcrições. Alguns dos vídeos coletados não possuíam transcrição e foram descartados, o próprio suporte Google⁸ indica diversos motivos, como: sotaque, dialetos, barulhos. Esses fatores podem impactar na legenda gerada automaticamente por meio de aprendizado de máquina, sendo necessário sempre revisar as legendas.

Assim a base de dados para a etapa de reconhecimento de entidades foi estabelecida, com um total de 2736 vídeos válidos. Todo o processo foi elaborado seguindo a

⁵<https://insightlab.ufc.br/por-que-o-python-e-a-linguagem-mais-adotada-na-area-de-data-science>

⁶YouTube Data API - <https://developers.google.com/youtube/v3>

⁷YouTube Transcript API - <https://pypi.org/project/youtube-transcript-api/>

⁸<https://support.google.com/youtube/answer/6373554>

arquitetura apresentada na Figura 1. O mesmo processo de coleta e tratamento foi aplicado, de maneira individual e separada da coleta geral, para os canais proprietários dos candidatos alvo da análise desse trabalho (Tabela 2).

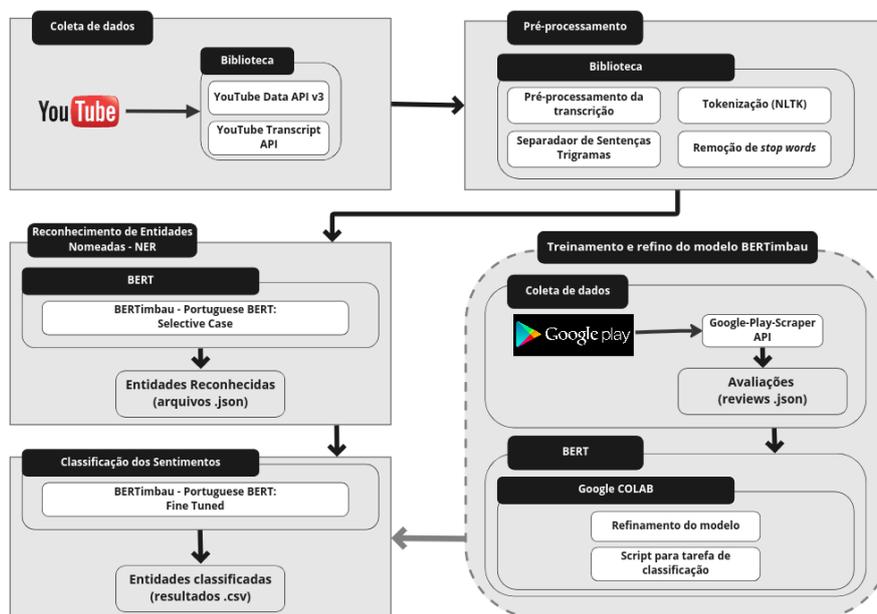


Figura 1. Fluxo completo do processo.

4.2. Pré-processamento dos dados

Efetuamos uma análise que aproxima-se da análise léxica, levando em consideração o sentimento expresso nas sentenças. Por isso foram removidas as chamadas *stop words* e, buscando uma análise em sentenças com elementos subjetivos, grande parte dos nomes próprios foram desconsiderados (utilizando expressões regulares).

As transcrições foram transformadas em *tokens* utilizando a biblioteca de ferramentas de linguagem natural NLTK ⁹, assim cada *token* representa uma palavra, facilitando a manipulação dos elementos para a abordagem utilizada.

Utilizando os resultados propostos em [Salas-Zárate et al. 2017] e considerando-se que a análise de sentimentos em nível de aspecto ainda é um campo aberto para descobertas [Poria et al. 2020], elaborou-se uma abordagem alternativa para a produção das sentenças analisadas. O item I da Figura 2 apresenta como as entidades foram centralizadas e as 3 palavras anteriores e posteriores (*3-around*) foram destacadas. Nos casos em que a entidade alvo estava localizada nas extremidades de um texto seja no final (item II da Figura 2) ou início (item III), eram selecionadas a maior quantidade de palavras possíveis que atendessem aos trigramas.

Cada entidade identificada em uma transcrição, no reconhecimento de entidades nomeadas, passou por esse tratamento antes de sua utilização na classificação de sentimento. Essa abordagem, aliada à remoção de *stop words* e outras possíveis entidades, diminui a chance de erros em que a sentença classificada não corresponda à entidade alvo.

⁹<https://www.nltk.org/>

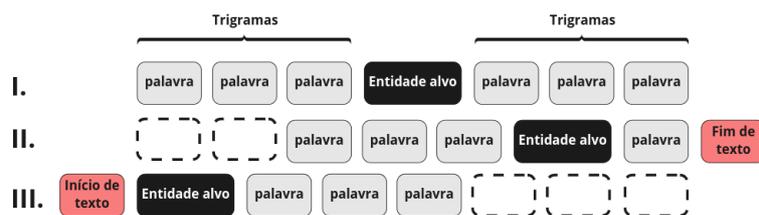


Figura 2. Representação do modelo de trigramas (3-around).

4.3. BERTimbau - Portuguese BERT

Com o propósito de analisar e identificar quais vídeos possuem conteúdos e qual a classificação dos sentimentos contidos nesses conteúdos, de acordo com qual candidato, foi necessário a utilização de uma ferramenta que possibilitasse tais tarefas. Buscou-se a utilização de conteúdos que agregassem em trabalhos na língua portuguesa, evitando a tradução para o inglês para assim utilizar alguma ferramenta para análise [Araújo et al. 2020][Zhang et al. 2021]. A escolhida para a execução desse trabalho foi o BERTimbau.

O BERT (Bidirectional Encoder Representations from Transformers) foi apresentado ao mundo em 2018, criado pela Google ¹⁰. Ele representa o estado da arte em PLN utilizando aprendizado de máquina e redes neurais.

Obter dados em Português para realizar esse tipo de tarefa é um grande desafio, devido à quantidade limitada de materiais. Felizmente, o modelo treinado pela NeuralMind [Souza et al. 2019] apresenta uma alternativa interessante para a língua portuguesa. Recebeu o nome BERTimbau [5] e atualmente é um dos recordistas em downloads no site especializado em modelos Hugging Face ¹¹.

O pré-treinamento é uma das etapas que permite o ajuste do modelo para a redução de resultados divergentes, o que implica na necessidade de um ajuste posterior para que o BERT seja aplicado para a tarefa objetivo (seção 4.4). Para o reconhecimento de entidades ele apresenta duas opções de classes de entidades: total e *selective*. Como o objetivo é identificar pessoas, utilizou-se o BERTimbau *selective*, que é capaz de identificar as classes pessoa, organização, local, tempo e valor. Ela também foi a que obteve melhor resultado no *benchmark* realizado pelos autores [Souza et al. 2019]. Os autores do BERTimbau utilizaram para treino a BrWaC (Brazilian Web as Corpus), um grande corpus resultante do trabalho apresentado por [Wagner Filho et al. 2018]. E, para o treino específico dos modelos de reconhecimentos de entidades, utilizou-se o corpus presente na coleção MiniHAREM¹².

O fluxo das etapas seguidas no processo é demonstrado na Figura 1. A análise final contou com duas camadas do BERTimbau: a primeira identificou as entidades (NER) e a segunda foi utilizada para classificar as sentenças, gerando assim os resultados observados na seção 5.

¹⁰<https://ai.googleblog.com/2018/11/open-sourcing-bert-state-of-art-pre.html>

¹¹<https://neuralmind.ai/2020/11/29/bertimbau-da-neuralmind-e-recordista-em-downloads/>

¹²https://www.linguateca.pt/primeiroHAREM/harem_miniharem.html

Tabela 3. Relatório de classificação

Classificação	Precisão	Revocação	F1-Score
positivo	0.71	0.69	0.70
negativo	0.71	0.65	0.68
neutro	0.66	0.67	0.66
parcialmente positivo	0.51	0.55	0.53
parcialmente negativo	0.49	0.48	0.49

4.4. Treinamento e refino dos modelos

O BERTimbau foi utilizado em dois momentos: aplicação de NER e, posteriormente, classificação de sentimentos das entidades reconhecidas. O modelo de NER foi utilizado como fornecido pelos criadores, sem ajustes ou novos treinamentos.

Para o processo de classificação foram utilizados os hiper parâmetros recomendados pelos criadores do BERTimbau para *fine-tuning* (refinamento) sendo eles: *Batch size*:16; *Learning rate* (Adam): $2e-5$; *Number of epochs*: 10. O modelo apresentou uma acurácia de 62% sendo que o relatório de classificação apontou os dados descritos na Tabela 3.

Com a primeira etapa de coleta e processamento dos dados concluída, iniciou-se a segunda etapa, responsável pela classificação dos sentimentos. Para o refinamento do modelo de BERTimbau foi realizada uma raspagem de dados, buscando por *reviews* de aplicativos da Google Play, utilizando uma API chamada Google-Play-Scraper¹³, em que foram coletados dados dos nove maiores aplicativos da categoria "Food and Drinks".

A escolha dessa base de dados advém da presença de uma característica das revisões, que permite a classificação pelos clientes entre um valor de 1 a 5 estrelas com um conteúdo em Português, o que foi importante para a classificação de sentimentos realizada neste trabalho, a qual considera 5 sentimentos possíveis: negativo, parcialmente negativo, neutro, parcialmente positivo e positivo. Geralmente as bases de dados rotuladas existentes apresentam apenas 3 sentimentos possíveis: negativo, neutro e positivo. Foram extraídas a mesma quantidade de cada avaliação e, buscando uma maior representatividade, foram filtradas por relevância.

4.5. Validação e normalização dos resultados

Os resultados obtidos foram divididos em dois grupos: um com o resultado das buscas obtidas e outro com a análise dos vídeos dos canais próprios dos candidatos. As análises foram realizadas com base na incidência das entidades e suas classificações. Contudo, devido a uma grande ocorrência de sentenças neutras observou-se a necessidade de uma maneira de normalizar estes resultados, possibilitando uma equivalência entre os candidatos e canais.

Buscando uma normalização desses resultados que permita uma análise comparativa, indicando tendências e comportamentos, a estratégia adotada foi a criação de uma fórmula para normalização. Ela leva em consideração fatores intrínsecos do modelo como

¹³<https://github.com/JoMingyu/google-play-scraper>

Tabela 4. Valores para cálculo do Score

Sentimento	Precisão (P)	Peso (W)
positivo	0.71	2.0
negativo	0.71	-2.0
parcialmente positivo	0.51	1.0
parcialmente negativo	0.49	-1.0

a precisão (vide Tabela 3), a ocorrência de cada sentimento para determinada entidade e também a atribuição de pesos diferentes para os sentimentos.

Essa fórmula foi utilizada para a obtenção de uma pontuação utilizada para comparar o comportamento entre os canais e dentro de um mesmo vídeo, conforme a eq. 1. São somadas todas as ocorrências de um sentimento n (E_n) e multiplicadas pelo peso atribuído (W_n) aquele sentimento e pela precisão (P_n) apresentada pelo modelo de classificação (tabela 3), com o somatório dividido pela quantidade total de entidades avaliadas (Et). Ao resultado final denominou-se *Score*.

$$\text{Score} = \left(\sum_{n=1}^n (E_n * W_n * P_n) \right) / Et \quad (1)$$

Grande parte das transcrições obtidas tem origem em canais jornalísticos, que na maior parte do tempo buscam transmitir a informação de maneira objetiva, fato esse que implica em muitas classificações neutras. Como este trabalho busca estudar o contexto de polarização, sentenças subjetivas tem um valor maior e para elas foram atribuídos pesos maiores. Posicionamentos sutis também geram pouco impacto nas redes, levando aos sentimentos parciais uma atribuição de pesos que diminuem seu impacto no *Score* final. Os pesos atribuídos estão representados na tabela 4.

5. Resultados e Discussão

Nas seções até aqui foram descritas as etapas realizadas nesse trabalho. Foram coletadas transcrições que contém o que foi dito nos vídeos de YouTube pesquisados. O objetivo foi extrair informações sobre o posicionamento dos vídeos, no que diz respeito aos candidatos, identificando assim comportamentos e práticas adotadas pelos produtores do conteúdo. Tratando-se de dois candidatos à presidência do país, os resultados e discussões apresentados nessa seção buscam trazer mais informações sobre o cenário de disputa e qual comportamento cada candidato aplicou em seu próprio canal.

O conteúdo foi avaliado e submetido a uma análise de sentimento, responsável por classificar o posicionamento manifestado sobre cada candidato, em diferentes momentos dos vídeos. Ao longo dessa seção serão descritos os comportamentos e resultados obtidos, primeiro utilizando a busca geral, feita com o uso de palavra-chave e, em seguida, a avaliação dos canais de Lula e Bolsonaro.

5.1. Pesquisa por palavra-chave

Em um panorama geral, as entidades foram o objeto de estudo do trabalho, e seu volume (incidência) também é indicativo de comportamento e estratégia. Foram identificados e

Tabela 5. Entidades identificadas por canais.

Canal	Nº vídeos	Nº de entidades	Score-Bolsonaro	Score-Lula
CNN Brasil	123	4238	0.38	0.45
UOL	120	5780	0.29	0.33
Jovem Pan News	73	1044	0.27	0.18
Jornalismo TV Cultura	31	464	0.24	0.14
SBT News	29	1210	0.49	0.48

classificados 36140 sentimentos. Desse total, as mais frequentes foram neutras e parcialmente positivas, sendo que a distribuição total foi de 60,7% das classificações neutras, 28,7% parcialmente positivas, 4,3% parcialmente negativas, 3,2% negativas e 3,1% positivas.

5.1.1. Entidades e principais canais

Nesse panorama é possível observar que o número de menções positivas foi o menor entre todos os sentimentos. A grande presença de análises neutras é justificável pela grande quantidade de sentenças objetivas, ou seja, os vídeos na maioria das vezes foram baseados em fatos, notícias. Vídeos em que opiniões são expressas (sentenças subjetivas) tendem a fugir desse cenário.

No que diz respeito à distribuição das entidades entre os candidatos, foi possível observar uma presença levemente maior do candidato Bolsonaro (53,2%). Este fato é justificável por ser o presidente em exercício durante o período, e um indicativo de sua presença maior nas redes sociais.

Quanto aos canais com maiores expoentes de entidades, na tabela 5 estão descritos os canais com maior número de vídeos publicados e a quantidade de entidades analisadas, além da pontuação (*Score*) obtido por cada candidato. Os canais que apresentaram maior número de menções aos candidatos foram UOL e CNN Brasil, representando 15,9% e 11,7% respectivamente. Na Figura 4 fica indicado visualmente a quantidade superior de entidades reconhecidas dos candidatos em cada canal, com uma presença maior do candidato Bolsonaro em todos eles. É interessante observar que os principais canais analisados são jornalísticos, que originaram grande parte das sentenças classificadas como neutras. O número total de canais que tiveram pelo menos um vídeo que se enquadraram efetivamente nos parâmetros de busca foi 333.

Um fato interessante observado é que em vídeos de opinião (canais particulares, influenciadores, entre outros) o título do vídeo raramente apresentava o nome dos candidatos no título. Uma exceção a essa regra são os *Podcasts*, que utilizam o nome para as chamadas e possuem um tópico próprio a seguir.

5.1.2. *Score*

Calculou-se os *Scores* dos principais canais, conforme apresentado na seção 4, e estão dispostos na tabela 5. Um *Score* próximo a 0 indica um balanceamento entre as opiniões

Tabela 6. Média e desvio padrão de curtidas, visualizações e comentários.

	Visualizações	Comentários	Curtidas
Média	350 337.9	1 723.7	18 538.1
Desvio padrão	1 075 626.0	6 664.5	100 352.6

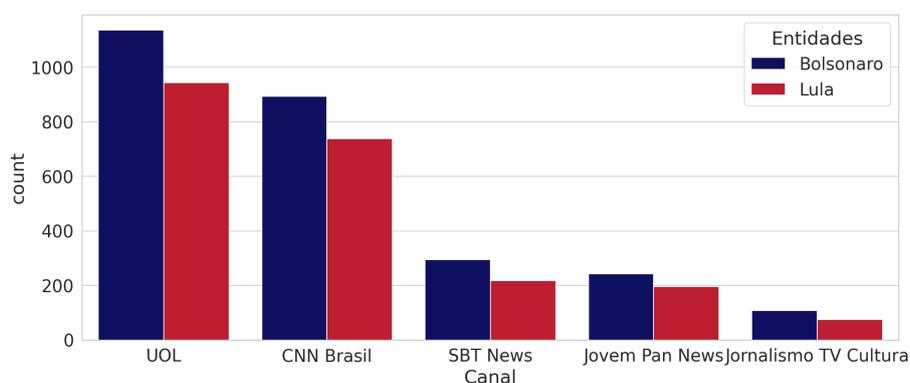


Figura 3. Distribuição das entidades nos principais canais.

expressas no conteúdo das transcrições, esse valor seria um indicativo de neutralidade com relação aos candidatos. Existe uma relação direta entre o crescimento do valor e a quantidade de entidades identificadas, por isso a análise feita aqui se restringe ao comportamento entre os candidatos em um mesmo canal.

Observa-se que para ambos os candidatos os canais apresentaram *Scores* positivos, o que indica que, no geral, o comportamento dos canais indicou manifestações positivas. Um distanciamento entre valores indica uma manifestação mais positiva de um sobre o outro, como pode ser observado nos canais Jovem Pan News e Jornalismo TV Cultura à favor de Bolsonaro e nos canais CNN Brasil e UOL a favor de Lula. O canal SBT News registrou pontuação similar para ambos, o que pode ser considerado também como neutralidade.

5.2. Canais dos candidatos

Considerando que uma fonte de dados confiável para analisar o comportamento dos candidatos no YouTube seriam seus canais proprietários, realizamos a coleta de todos as transcrições dos mesmos. Ao todo foram coletados 5304 vídeos, dos quais 3548 são do canal de Bolsonaro e 1756 do canal de Lula. Contudo, nem todos os vídeos passaram pela classificação de sentimentos.

Dos vídeos coletados dos canais de Lula e Bolsonaro, 90% e 80% deles, respectivamente, tiveram alguma entidade reconhecida correspondente aos candidatos. Este fato é explicado pela ausência de transcrições de alguns vídeos e a incapacidade ou ausência de uma entidade correspondente, no processo de NER.

Assim, foram classificados 9928 sentimentos ficando dispostos em: 54,2% de neutros, 23,8% de parcialmente positivos, 9,7% de parcialmente negativos, 6,2% de negativos e 6,1% classificados como positivos. A maioria das ocorrências foram sobre o candidato Lula, representando 59,6% da base de dados dos canais dos candidatos.

Tabela 7. Score calculado dos canais de Lula e Bolsonaro.

Canal	Nº de entidades	Score-Bolsonaro	Score-Lula
Lula	6677	-0.06	0.45
Bolsonaro	3251	0.31	0.10

Buscando ir além do entendimento geral e identificar a estratégia de redes sociais de cada um, uma análise individual de cada canal foi feita, sendo calculado o *Score*, conforme apresentado na Tabela 7.

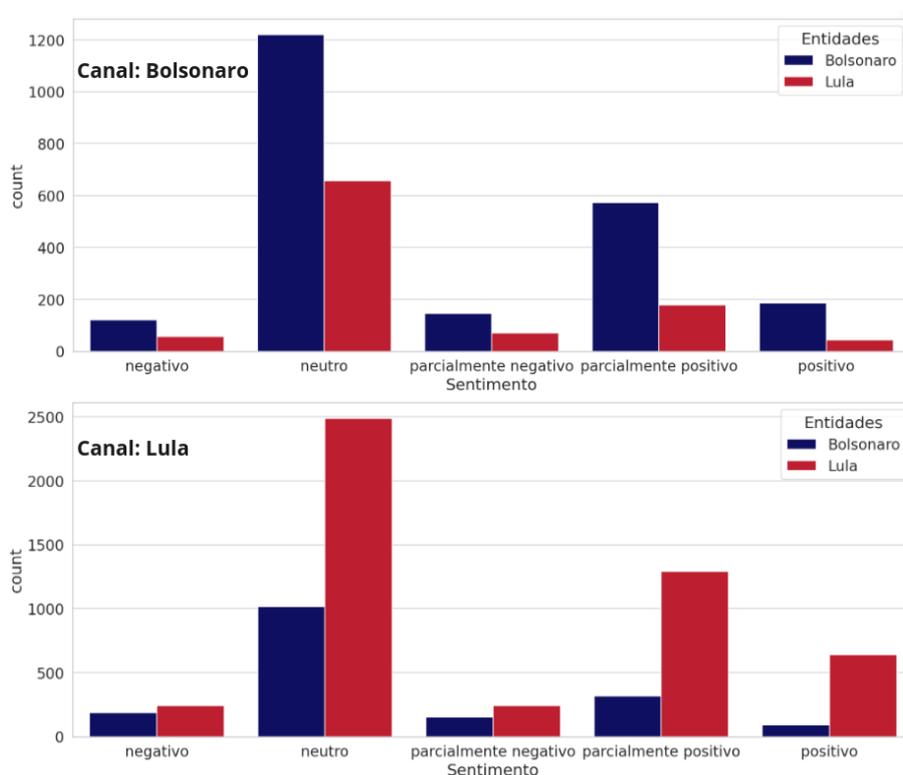


Figura 4. Distribuição dos sentimentos nos canais: Lula x Bolsonaro.

5.2.1. Score

Nas transcrições avaliadas no canal de Lula, as menções a ele mesmo foram de 73,5% contra 26,4% relacionadas a Bolsonaro. Um valor relativamente alto e que indica que o canal busca enaltecer o candidato. Ao observarmos o cálculo do *Score*, o valor calculado para o candidato Bolsonaro é negativo, o que indica a predominância de sentimentos negativos nos vídeos. Para Lula, o resultado de 0.45 reforça a presença da estratégia de enaltecer seus próprios atos e também apontar falhas de seu adversário.

Já para as transcrições avaliadas no canal de Bolsonaro, as menções a ele mesmo foram de 69,0%, ligeiramente menores do que no canal de seu adversário, mas que ainda assim indica a mesma estratégia que o canal de Lula no que diz respeito às menções. Observando o valor do *Score*, o resultado do cálculo para ambos é positivo. Isso pode ser

Tabela 8. Estatísticas dos principais Podcasts.

Candidato	Canal	Visualizações	Comentários	Curtidas
Bolsonaro	Inteligência Ltda.	17 047 970	102 022	2 620 102
Bolsonaro	Flow Podcast	16 179 257	182 007	1 672 106
Lula	Flow Podcast	9 549 930	104 976	955 500

um indicativo de que o modelo não foi eficiente para classificar as sentenças que contém ironia, visto que é uma das características marcantes nos posicionamentos do candidato Bolsonaro. Esse tipo de dificuldade é um dos grandes problemas no que diz respeito à Análise de Sentimentos e abordagens que utilizam Linguagem Natural.

5.3. Podcasts

Podcasts estão presentes no cotidiano dos brasileiros, sendo que quatro em cada 10 usuários da internet já ouviram¹⁴. Sua presença nesta seção vem da necessidade de investigar um comportamento discrepante de alguns dos vídeos coletados.

Os vídeos analisados tiveram outros metadados de controle coletados, além do conteúdo principal que foram suas transcrições. Comentários, curtidas e visualizações também foram obtidos e observou-se um comportamento diferente quando tratava-se de *Podcasts*. Os vídeos de Podcasts analisados nessa seção fazem parte da base de dados da coleta realizada utilizando busca por palavra-chave. As participações de Bolsonaro e Lula nesses programas apresentaram valores exorbitantes. A Tabela 8 exhibe os valores coletados durante o período de desenvolvimento desse trabalho.

Os valores são intrigantes pois ao calcular média e desvio padrão (tabela 6) de toda a base de dados foram encontrados valores muito inferiores e que apontam para um comportamento incomum. Não é possível concluir que a utilização de *bots* foi responsável pelos números absurdos. É válido ressaltar que o contexto político observado é polarizado e marcou uma das disputas presidenciais mais acirradas do Brasil. Pode-se concluir que esse tipo de vídeo é de extrema importância pois, representa um elevado alcance e engajamento para os candidatos, o que levou a análise da participação de ambos no podcast Flow.

A participação de Lula no Flow apresentou 32 menções a ele mesmo e 15 ao candidato Bolsonaro. Considerou-se que as menções próprias foram majoritariamente realizadas pelo próprio entrevistador, indicando ali o sentimento por detrás de suas perguntas ou comentários. Desconsiderando-se manifestações neutras, as menções negativas a Bolsonaro foram de aproximadamente 85%.

Quanto a participação de Bolsonaro, ambos foram citados 28 vezes cada, sendo que os resultados obtidos foram de sentimentos parciais e, portanto, sendo inconclusivos para uma análise. É importante ressaltar que as legendas possuem uma geração automática mas, sua qualidade está também vinculada ao tipo de discurso utilizado pelo interlocutor. Esse fato implica que uma dinâmica de comunicação informal ou excesso de manifestações irônicas impactam no resultado da análise de nossa abordagem.

¹⁴<https://www.t1noticias.com.br/geral/podcasts-explodem-em-popularidade-com-temas-que-voao-do-poker-ao-cinema/104387/>

Foi analisado também o vídeo que apresentou os maiores números, que foi a entrevista de Bolsonaro ao podcast Inteligência Ltda. Nele foram identificadas 72 entidades referentes a Lula e 24 ao próprio entrevistado. Durante o podcast em que Bolsonaro participou, 75% das menções foram ao seu adversário (Lula), sendo que aproximadamente 72% dessas menções foram negativas. Esses valores são compatíveis a estratégia de ataque adotada pelos candidatos durante a campanha eleitoral.

6. Considerações Finais e Trabalhos Futuros

Nesse trabalho, apresentamos uma abordagem híbrida de análises de sentimentos, combinando uma aproximação léxica com análises em nível de aspecto utilizando o BERTimbau, seguindo um processo de coleta e análises de transcrições do YouTube, avaliando as manifestações e incidências dos dois principais candidatos à presidência do Brasil em 2022.

Acreditamos que nossa abordagem de análise obteve resultados satisfatórios, principalmente considerando as dificuldades e a escassez de trabalhos como este utilizando a língua portuguesa. Os resultados indicaram que é possível observar comportamentos em canais do YouTube, identificando as entidades presentes e classificá-las de acordo com o sentimento apresentado.

Apresentamos a fórmula do *Score*, como uma métrica para comparar o comportamento para cada entidade em um mesmo cenário de análise (vídeo ou canal), que permitiu, por exemplo, que observássemos um padrão de comportamento em canais de comunicação e os canais proprietários dos candidatos. Na análise do canal de Lula, o Score obtido para ele foi 0.45, reforçando a presença de uma estratégia de enaltecer seus próprios atos e também apontar falhas de seu adversário, que obteve um resultado negativo (-0.06). Para o canal de Bolsonaro, os resultados indicaram também o enaltecimento do proprietário, mas outras análises não foram conclusivas.

Como trabalhos futuros, identificamos a importância de melhorar o pré-processamento de dados, evitando assim o desperdício de dados coletados e a expansão das entidades avaliadas. Também é possível a implementação de métodos que reduzam as classificações equivocadas, identificando ironia, por exemplo. Outra alternativa é a criação de uma base própria, com dados rotulados sobre política. Por último, pretende-se a integração do processo e a possibilidade de fornecer ontologias como parâmetros de entrada, permitindo assim uma busca e classificação abrangente sobre determinada entidade.

Referências

- Araújo, M., Pereira, A., and Benevenuto, F. (2020). A comparative study of machine translation for multilingual sentence-level sentiment analysis. *Information Sciences*, 512:1078–1102.
- Boxell, L., Gentzkow, M., and Shapiro, J. M. (2017). Greater internet use is not associated with faster growth in political polarization among us demographic groups. *Proceedings of the National Academy of Sciences*, 114(40):10612–10617.
- Buder, J., Rabl, L., Feiks, M., Badermann, M., and Zurstiege, G. (2021). Does negatively toned language use on social media lead to attitude polarization? *Computers in Human Behavior*, 116:106663.

- Gu, S., Zhang, L., Hou, Y., and Song, Y. (2018). A position-aware bidirectional attention network for aspect-level sentiment analysis. In *Proceedings of the 27th international conference on computational linguistics*, pages 774–784.
- Kobellarz, J. K. et al. (2021). Polarização virtual: estudo da dinâmica de cenários politicamente polarizados em sites de redes sociais. Master's thesis, Universidade Tecnológica Federal do Paraná.
- Kubin, E. and von Sikorski, C. (2021). The role of (social) media in political polarization: a systematic review. *Annals of the International Communication Association*, 45(3):188–206.
- Larsson, A. O. (2016). Online, all the time?: a quantitative assessment of the permanent campaign on facebook. *New Media Society*, 18(2):274–292.
- Liu, K., Xu, L., and Zhao, J. (2012). Opinion target extraction using word-based translation model. In *Proceedings of the 2012 Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning*, page 1346–1356.
- Marrero, M., Urbano, J., Sánchez-Cuadrado, S., Morato, J., and Gómez-Berbís, J. M. (2013). Named entity recognition: fallacies, challenges and opportunities. *Computer Standards & Interfaces*, 35(5):482–489.
- Mohit, B. (2014). Named entity recognition. In *Natural language processing of semitic languages*, pages 221–245. Springer.
- Newman, N., Fletcher, R., Robertson, C. T., Eddy, K., and Nielsen, R. K. (2022). Reuters institute digital news report 2022. In *Reuters Institute Digital News Report 2022*, pages 116–117.
- Poria, S., Hazarika, D., Majumder, N., and Mihalcea, R. (2020). Beneath the tip of the iceberg: Current challenges and new directions in sentiment analysis research. *IEEE Transactions on Affective Computing*.
- Robles, J.-M., Guevara, J.-A., Casas-Mas, B., and Gómez, D. (2022). When negativity is the fuel. bots and political polarization in the covid-19 debate. *Comunicar*, 30(71):63–75.
- Rodríguez-Ibáñez, M., Gimeno-Blanes, F.-J., Cuenca-Jiménez, P. M., Soguero-Ruiz, C., and Rojo-Álvarez, J. L. (2021). Sentiment analysis of political tweets from the 2019 spanish elections. *IEEE Access*, 9:101847–101862.
- Salas-Zárate, M. d. P., Medina-Moreira, J., Lagos-Ortiz, K., Luna-Aveiga, H., Rodriguez-García, M. A., and Valencia-García, R. (2017). Sentiment analysis on tweets about diabetes: an aspect-level approach. *Computational and mathematical methods in medicine*, 2017.
- Shah, A., Singh, A., Abhishek, K., and Ramesh, S. S. (2021). Sentimental analysis for political polarization using vader sentiment lexicon.
- Sharnagat, R. (2014). Named entity recognition: A literature survey. *Center For Indian Language Technology*, pages 1–27.

- Silva, C. A. and Barbosa, D. M. (2019). Analyzing the acceptance of the 2018 brazilian presidential election's main candidates based on youtube comments. In *Proceedings of the 25th Brazillian Symposium on Multimedia and the Web*, pages 377–384.
- Silva, F. and Freitas, L. (2022). Brazilian portuguese hate speech classification using bertimbau. In *The International FLAIRS Conference Proceedings*, volume 35.
- Souza, F., Nogueira, R., and Lotufo, R. (2019). Portuguese named entity recognition using bert-crf. *arXiv preprint arXiv:1909.10649*.
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, Ł., and Polosukhin, I. (2017). Attention is all you need. In *Advances in neural information processing systems*, pages 5998–6008.
- Wagner Filho, J. A., Wilkens, R., Idiart, M., and Villavicencio, A. (2018). The brwac corpus: a new open resource for brazilian portuguese. In *Proceedings of the eleventh international conference on language resources and evaluation (LREC 2018)*.
- Yarchi, M., Baden, C., and Kligler-Vilenchik, N. (2021). Political polarization on the digital sphere: A cross-platform, over-time analysis of interactional, positional, and affective polarization on social media. *Political Communication*, 38(1-2):98–139.
- Zhang, W., He, R., Peng, H., Bing, L., and Lam, W. (2021). Cross-lingual aspect-based sentiment analysis with aspect term code-switching. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 9220–9230.