

# Análise Aprimorada das Percepções dos Usuários por Meio de Abordagens de Processamento de Linguagem Natural

Ana Cláudia Machado<sup>1</sup>, Gabriel Prenassi<sup>1</sup>, Elisa Tuler<sup>1,2</sup>, Leonardo Rocha<sup>1</sup>

<sup>1</sup>Universidade Federal de São João del-Rei, <sup>2</sup>Ministério da Gestão e Inovação em Serviços Público  
(anaclaudiamachado211,prenassigabriel)@aluno.ufsj.edu.br;etuler@ufsj.edu.br;lrocha@ufsj.edu.br

## ABSTRACT

In this article, we present a framework for analyzing users' textual comments, going beyond the evaluation carried out using the System Usability Scale, a traditional methodology for system assessment. The framework performs textual analysis through Topic Modeling, Text Summarization, and Sentiment Analysis. Topic Modeling identifies semantic topics in the comments, while Text Summarization generates summaries of each topic, enabling explainability. Sentiment Analysis, in turn, categorizes the sentiment of each topic, identifying its polarity. We propose an intuitive visual interface and apply our framework in a real-world scenario, where it enabled more accurate feedback.

## KEYWORDS

Modelagem de Tópicos, Análise de Sentimento, LLMs, SUS

## 1 INTRODUÇÃO

A Experiência do Usuário refere-se à percepção e às respostas durante a interação com um sistema, produto ou serviço. Sua participação ativa no desenvolvimento é essencial para garantir aderência às especificações e ao contexto de uso, além de identificar problemas. Ainda assim, interpretar o *feedback* dos usuários é desafiador. Instrumentos como o *System Usability Scale* (SUS) [2] avaliam usabilidade com pontuação de 0 a 100. Há relatos na literatura de inconsistências entre avaliações numéricas atribuídas pelos usuários e seus comentários [7], que trazem informações valiosas não expressas numericamente. A análise textual complementa a avaliação, mas a diversidade e o volume dos comentários dificultam sua interpretação manual, exigindo ferramentas para extrair dados estruturados e úteis.

Este trabalho propõe um *framework* de quatro módulos que analisa comentários de usuários com o uso coordenado de estratégias de Processamento de Linguagem Natural (PLN). O **módulo 1** é responsável por extrair os dados. O **módulo 2** aplica Modelagem de Tópicos (i.e. *BERTopic*[5]) para agrupar semanticamente os comentários. O **módulo 3** é uma adaptação da Análise de Sentimento tradicional, adotando uma polaridade contextual (sugestão, *feedback* positivo, crítica e não pertinente), utilizando um Grande Modelo de Linguagem (*Large Language Model - LLM* [1]), em nosso caso o LLaMA 3.1, por ser aberto e amplamente utilizado. O **módulo 4** realiza Sumarização Automática dos comentários associados a cada tópico com o mesmo LLM. O **módulo 5** apresenta uma interface visual com os tópicos e seus impactos na avaliação geral.

O *framework* foi aplicado à funcionalidade “Simulador de Aposentadoria” da plataforma **SouGov**<sup>1</sup>, desenvolvida pelo Ministério da Gestão e Inovação em Serviços Públicos (MGI). Ele identificou pontos críticos e gerou *insights* relevantes para a equipe, incluindo a sumarização e priorização de aspectos positivos, negativos e sugestões. Os resultados mostraram divergências entre as notas do SUS e os comentários, apontando estes últimos como mais representativos [7].

As implementações foram realizadas pela aluna Ana Machado, sob a orientação do professor Leonardo Rocha, com a colaboração do aluno Gabriel Prenassi. Esse trabalho está inserido em um projeto de colaboração técnica entre a UFSJ e o MGI, sob a coordenação da professora Elisa Tuler. O *framework* encontra-se em operação no MGI e resultou na publicação de um artigo na INTERACT (A1).

## 2 FUNDAMENTAÇÃO TEÓRICA

**System Usability Scale (SUS):** Usabilidade é a medida em que usuários específicos utilizam um produto para alcançar objetivos com eficácia, eficiência e satisfação em um contexto. Uma metodologia amplamente usada para avaliá-la é o SUS [2], composto por dez itens em escala Likert de 1 a 5, variando de “discordo totalmente” a “concordo totalmente”. As afirmações ímpares são positivas e as pares, negativas. Para calcular a pontuação do SUS, soma-se a pontuação das afirmações positivas (subtraindo 5) e das negativas (subtraindo de 25). O resultado é multiplicado por 2.5, gerando um valor entre 0 e 100. A média das pontuações individuais representa a nota final, sendo 68 a mediana na abordagem relativa [2]. Embora forneça uma visão quantitativa da usabilidade, não permite capturar percepções específicas. Comentários em campo aberto complementam essa análise, mas sua interpretação manual é complexa e demorada.

**Modelagem de Tópicos:** Agrupa semanticamente uma coleção de documentos  $D$  em  $k$  tópicos, representados pelas  $x$  palavras mais relevantes [11]. Seu desempenho depende da forma de representar os documentos. A abordagem tradicional usa *bag-of-words* com TF-IDF, que pondera termos pela frequência no documento e raridade no corpus. Para incorporar semântica, [11] utiliza a similaridade de cosseno entre *embeddings* para formar *clusters*. Outra abordagem gera *embeddings* via *Transformers*, como BERT. O *BERTopic* [5] agrupa esses *embeddings* em *clusters* e aplica TF-IDF para extrair palavras-chave, oferecendo flexibilidade nas etapas. Mais recentemente, surgiram técnicas baseadas em LLMs, como o *TopicGPT* [10].

**Sumarização Automática de Textos:** Gera versões resumidas de um texto por três abordagens: extrativa (seleção de trechos), abstrativa (geração de novo texto) e híbrida (combinação das anteriores). Recentemente, *Transformers* de 1ª geração têm sido usados pela eficácia em PLN. Com a chegada dos LLMs, houve revolução devido ao maior poder computacional e melhor compreensão semântica e contextual. A forma mais simples de usar LLMs é via *prompt*.

In: V Concurso de Trabalhos de Iniciação Científica (CTIC 2025). Anais Estendidos do XXXI Simpósio Brasileiro de Sistemas Multimídia e Web (CTIC'2025). Rio de Janeiro/RJ, Brasil. Porto Alegre: Brazilian Computer Society, 2025.  
© 2025 SBC – Sociedade Brasileira de Computação.  
ISSN 2596-1683

<sup>1</sup><https://www.gov.br/servidor/pt-br/assuntos/sou-gov>

**Análise de Sentimentos:** Identifica emoções em textos por duas abordagens principais: léxica, que usa vocabulários com PLN, e aprendizado de máquina, que aprende padrões para prever sentimentos [7]. Métodos léxicos recentes, como Lex2Sent [6], combinam léxicos binários e reamostragem para identificar palavras-chave e agregar múltiplos *embeddings*, superando abordagens puramente léxicas. Contudo, LLMs capturam melhor as estruturas semântica e sintática. A escassez de léxicos em línguas como o português ainda é um desafio. Abordagens de aprendizado de máquina exigem grandes bases rotuladas, nem sempre disponíveis. Técnicas baseadas em LLMs, como *in-context learning*, que usam poucos exemplos rotulados, têm mostrado resultados promissores [1].

### 3 TRABALHOS RELACIONADOS

Não há na literatura trabalhos que combinem estratégias de PLN para análise de comentários textuais na avaliação de sistemas. Porém, há estudos que destacam benefícios da análise textual. Nosso trabalho integra a discussão em IHC [3] sobre o impacto dos LLMs na análise de dados e os desafios éticos envolvidos, apresentando uma ferramenta que complementa metodologias existentes, como o SUS. Estudos em [4, 8, 12] indicam direções para nosso trabalho: o primeiro usa PLN para sugerir código automaticamente; o segundo investiga percepções sobre ferramentas baseadas em LLMs; o terceiro analisa entrevistas qualitativas com LLMs.

### 4 FRAMEWORK PROPOSTO

Nosso *framework* compreende cinco módulos, descritos a seguir.

#### 4.1 Extração de dados

O primeiro módulo realiza a extração de dados, configurado para receber um arquivo CSV. Cada linha representa uma avaliação, com a primeira coluna sendo uma identificação anonimizada do usuário e as 10 colunas seguintes correspondem às afirmativas SUS, já na escala Likert (de 1 a 5). A última coluna contém o comentário textual.

#### 4.2 Modelagem de Tópicos

Utilizamos o BERTopic [5] devido à sua representação semanticamente rica e contextualizada, baseada em *embeddings*. Adotamos técnicas de pré-processamento adequadas à linguagem e ao domínio de aplicação dos comentários [9]. Na representação vetorial dos documentos, utilizamos o modelo multilíngue *Sentence-BERT* (SBERT). A redução de dimensionalidade dos *embeddings* foi feita utilizando as configurações padrões do próprio modelo SBERT. Propomos três configurações para o número de tópicos — 5, 10 e 15 — de forma a explorar diferentes níveis de granularidade. Devido ao alto custo computacional, estabelecemos esses valores padrões para otimizar o processo, garantindo eficiência com resultados pré-computados. Para cada tópico, o modelo define as palavras mais representativas por meio do *c-TF-IDF*, uma variação do *TF-IDF*, que foca nos fatores que distinguem os documentos presentes em diferentes *clusters*.

#### 4.3 Sumarização Automática de Texto

O objetivo desse módulo é fornecer sumarizações semântica e contextual de cada tópico. Optamos pelo uso de LLM (i.e. Llama-3.1-8B com quantização de 4 bits). Propomos duas abordagens. A primeira envolve uma sumarização detalhada, onde extraímos um resumo dos comentários em um tópico. Devido a limitações de memória, definimos um limite para o número de comentários no *Prompt 1*. Se o total estiver dentro desse limite, geramos o resumo do tópico. Caso contrário, combinamos essa saída com os comentários restantes e criamos um resumo final pelo *Prompt 2*. A segunda abordagem gera

uma versão mais concisa do resumo anterior, utilizando o *Prompt 3*. Por fim, o *Prompt 4* utiliza o resumo conciso para criar um título para cada tópico. Esses resumos oferecem uma compreensão mais clara dos tópicos, e, junto com a análise de sentimentos, ajudam a identificar rapidamente os pontos positivos e as áreas que necessitam de melhorias no sistema avaliado.

#### Prompt 1 - Sumarização detalhada

You will receive a set of user comments, and your task is to summarize them by extracting the positive points, negative points, and suggestions mentioned.

[Instructions]

Step 1: Determine the main points highlighted by users in a general manner.

- Do not highlight points that have been little commented on, return only those with the most relevance in the set.

- Generalize the points highlighted in the comments as much as possible without losing the context.

Step 2: Return positive, negative and suggestion points according to the response pattern specified below.

- Example of response:

""Positive Points"":

1. Positive Point 1

2. Positive Point 2

""Negative Points"":

1. Negative Point 1

2. Negative Point 2

""Suggestions""

1. Suggestion 1

2. Suggestion 2

- Do not generate anything other than what is requested in the response pattern.

- Summarize the content without directly quoting user comments.

- Do not create new points; only summarize the existing ones.

- Do not use the first person in the generated text.

#### Prompt 2 - Sumarização detalhada (parcial com os comentários restantes)

You will receive a partial summary of user comments. This summary has already integrated some feedback, but additional comments will be provided for you to enrich it further. Your task is to update and rewrite the summary, ensuring that all relevant points from the new comments are integrated while maintaining the original structure.

[Instructions]

Step 1: Analyze the partial summary and the new comments.

- Identify main points in a general manner, considering both existing summary and new feedback.

- Do not include points that were mentioned only once or are not significant in the overall set.

- Maintain the context while ensuring a neutral, structured, and concise synthesis.

Step 2: Rewrite the entire summary following the response pattern below.

- The final response must fully integrate both the partial summary and the new comments.

- Structure the response into three categories: Positive Points, Negative Points, and Suggestions.

- Example of response:

""Positive Points"":

1. Positive Point 1

2. Positive Point 2

""Negative Points"":

1. Negative Point 1

2. Negative Point 2

""Suggestions""

1. Suggestion 1

2. Suggestion 2

- Do not generate anything beyond the requested structure.

- Summarize the content without directly quoting user comments.

- Do not create new points; only summarize the existing ones.

- Do not use the first person in the generated text.

- Ensure that the updated summary remains well-structured, cohesive, and neutral.

#### Prompt 3 - Sumarização concisa

You will receive a summary of user comments, categorized into positive points, negative points, and suggestions.

[Instructions]

Step 1: Summarize the provided summary by extracting the key positive points, negative points, and suggestions mentioned.

- Do not use the first person in the summary.

- Return the summary in a single paragraph of up to 40 words, ensuring it remains clear and neutral.

- The summary should be in {language}.

#### Prompt 4 - Sumarização do tópico em uma única frase

You will receive a brief summary of user comments, which address positive points, negative points, and suggestions. Your task is to generate a concise, informative description that captures the core feedback.

[Instructions]

Step 1: Generate a concise description, up to 6 words, that highlights the key feedback.

- Do not add any new details or explanations—focus only on the relevant points from the summary provided.

- Do not add vague terms like "needs improvement" or "requires adjustments." Instead, highlight specific aspects mentioned in the comments.

- The concise description should be in {language}.

- Avoid redundancy and do not mention the application explicitly.

#### 4.4 Análise de Sentimentos

Permite uma análise das opiniões dos usuários sobre cada tópico. A ausência de dados rotulados é um fator decisivo na escolha da abordagem, sendo LLM uma alternativa viável, alinhando-se com a proposta multilíngue do *framework*. Foi proposto novamente o uso do Llama 3.1, rotulando os comentários como *feedback* positivo, sugestões, críticas ou não pertinentes para fornecer *insights* mais contextualizados. As críticas e *feedback* positivo estão, respectivamente, relacionados às categorias negativa e positiva. Devido à falta de dados previamente classificados, aplicamos *in-context* com 30 comentários rotulados por quatro especialistas, escolhendo três por categoria e assumindo que textos mais longos fornecem mais contexto ao LLM. Utilizamos o *Prompt 5* para rotular os comentários, sendo que a informação entre colchetes é adicionada apenas quando o modelo gera uma resposta fora do padrão, a fim de reclassificar o comentário. A análise é feita para cada comentário, e a polaridade geral de cada tópico é calculada com base na distribuição dessas análises.

##### Prompt 5 - Análise de Sentimentos

Classify the following texts, which are comments in {language} about {context}, into one of the following categories: criticism, suggestion, positive feedback, or not pertinent. Classify as 'not pertinent' only texts that are neither suggestions, positive feedback, nor criticisms, considering that they don't fit those categories but aren't necessarily irrelevant. These comments were provided by users who were encouraged to give suggestions, critiques, or positive feedback. The response must consist solely of the name of one of these categories, with no additional text or information.

Input: {criticism example}  
criticism

Input: {suggestion example}  
suggestion

Input: {positive feedback example}  
positive feedback

Input: {not pertinent example}  
not pertinent

[Attention! Classify only into the categories you were instructed to.]

Input: {text to be labeled}

#### 4.5 Interface Visual

Nosso quinto módulo é composto de metáforas visuais que permitem uma análise combinada de todos os resultados. Na Figura 1, apresentamos uma visão geral da interface. No Bloco 1, selecionamos a análise a ser realizada e o número de tópicos (5, 10 ou 15). No Bloco 2, visualizamos o número de participantes, bem como o número total de comentários. No Bloco 3, definimos a perspectiva de análise. Considerando a “Análise Geral dos Tópicos” como padrão, no Bloco 4 temos o título, as palavras mais relevantes e o número de comentários agrupados em um tópico; no Bloco 5 o resumo conciso; e no Bloco 6 a distribuição das categorias (polaridade) de comentários. Essas análises estão presentes em todos os tópicos para que possam ser classificadas de acordo com o número de críticas ou *feedback* positivo (Bloco 7). Na “Análise das Afirmações (Melhor/Pior)” no Bloco 3, apresentamos a visualização das afirmações do SUS com as melhores e piores pontuações e a distribuição de respostas associada a cada afirmação. Na opção “Análise dos Tópicos (Melhor/Pior)”, apresentamos os tópicos mais positivos e o mais negativos.

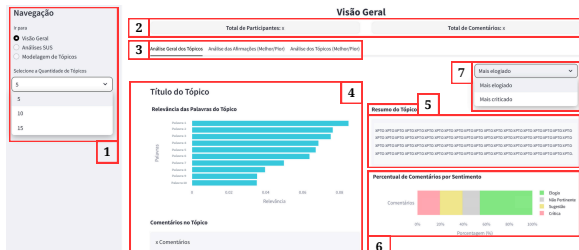


Figura 1: Visão geral da interface visual do *framework*.

A interface permite ainda uma avaliação detalhada das notas obtidas pelo SUS (Figura 2). No Bloco 8, escolhemos entre uma análise global (SUS Global) ou por afirmação (Afirmações de Concordância/Discordância). Na global, apresentamos a média das notas SUS, o desvio padrão e o número total de usuários (Bloco 9). No Bloco 10, as notas dos usuários são distribuídas de Pior Imaginável a Melhor Imaginável, descrevendo a experiência no SUS com adjetivos em vez de números. Já na outra análise, apresentamos o valor médio de cada afirmação do SUS, desvio padrão e a distribuição das respostas.

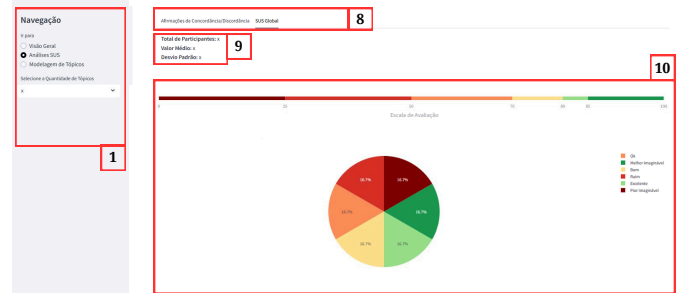


Figura 2: Visão geral da interface para a avaliação do SUS.

Ao selecionar a aba “Modelagem de Tópico” no Bloco 1 (Figura 3), temos as funcionalidades: “Análise do Tópico”, “Resumo do Tópico”, “Análise dos Comentários”, “Afirmações de Concordância/Discordância” e “Métrica SUS do Tópico”, destacadas no Bloco 12. A análise é realizada para um tópico específico selecionado no Bloco 13. Na aba “Análise do Tópico”, temos as mesmas informações presentes nos Blocos 4, 5 e 6 da Figura 1, mas as informações são apresentadas apenas para o tópico selecionado. Na aba “Resumo do Tópico”, visualizamos o resumo detalhado do tópico. Na aba “Análise dos Comentários”, é possível filtrar os comentários de três formas (Bloco 14). O primeiro filtro permite selecionar os comentários com base nas palavras mais relevantes do tópico ou considerando todas elas, enquanto o segundo filtro permite a seleção de uma ou todas as polaridades. Por fim, é possível definir o intervalo de pontuações SUS dos usuários que comentaram.

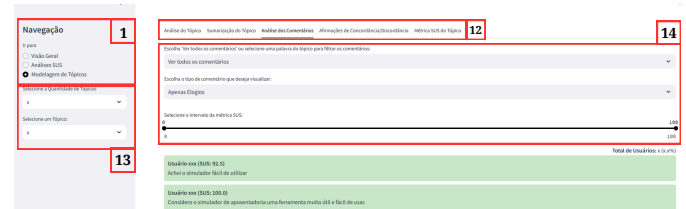


Figura 3: Análise de comentários para um tópico específico.

Assim, podemos definir um grande número de análises que podem ser realizadas em comentários de usuários de acordo com as configurações de filtro e diferentes perspectivas. Além disso, essas análises também podem ser contrastadas com aquelas do SUS. Na seção a seguir, apresentaremos alguns resultados obtidos da aplicação do *framework* proposto a dados de uma avaliação real, demonstrando a praticidade de nossa pesquisa.

## 5 AVALIAÇÃO EXPERIMENTAL

### 5.1 Estudo de Caso

Para validar o *framework*, utilizamos dados de uma avaliação da percepção dos usuários de uma nova funcionalidade do **SouGov**<sup>2</sup>, um “Simulador de Aposentadoria”, desenvolvido pelo MGI. Ainda em fase piloto, foi disponibilizado em agosto de 2024 para alguns órgãos federais que visam auxiliar os servidores públicos federais no acompanhamento de seus processos de aposentadoria. Foi enviado um questionário SUS para 87.555 funcionários de 57 órgãos, ficando disponível por 10 dias em outubro de 2024. Ao todo, foram recebidas 1.343 avaliações que também incluíam comentários textuais. Nosso objetivo é demonstrar que nossa proposta pode complementar os resultados quantitativos do SUS com a análise dos comentários.

### 5.2 Resultados Práticos

A equipe do MGI responsável pela avaliação da nova funcionalidade do SouGov destacou positivamente a avaliação numérica do SUS, mencionando que o *framework* permitiu uma análise aprofundada das respostas de cada afirmativa, visualizando sua pontuação, média e desvio padrão, conforme ilustrado na Figura 4. Percebe-se que a afirmação ilustrada na figura teve uma classificação média de 2.07, com “discordo totalmente” tendo mais de 49% das respostas. Porém, mais da metade dos participantes apresentaram respostas variadas, demonstrando que não há necessariamente um consenso, análise que não é tradicionalmente realizada no SUS.



Figura 4: Distribuição das respostas para afirmativas SUS.

Outros pontos relevantes em nossa ferramenta foram as funcionalidades relacionadas à análise automática de comentários. Considerando o filtro relacionado à pontuação SUS (Figura 5), ao selecionar o intervalo [68, 100], observa-se um padrão já relatado na literatura [7], que indica que os valores numéricos atribuídos pelos usuários, neste caso, por meio do SUS, são divergentes de seus comentários. Valores acima de 68 são considerados acima da mediana. Porém, ao analisar a classificação dos comentários, esta faixa ainda apresenta comentários categorizados como críticas, sendo 4,1% superior à porcentagem de elogios. Mesmo quando os usuários atribuem uma pontuação positiva a uma funcionalidade, eles ainda podem indicar problemas a serem melhorados em seus comentários. Numa avaliação tradicional do SUS, não seria possível identificar essas questões. Outro recurso importante para a equipe de avaliação foram as sumarizações, especialmente a detalhada, que teve como foco distinguir pontos positivos, negativos e sugestões. A combinação dessas informações facilitou diretamente a tomada de decisões, permitindo à equipe listar e priorizar as melhorias no simulador.

<sup>2</sup><https://www.gov.br/servidor/pt-br/assuntos/sou-gov>

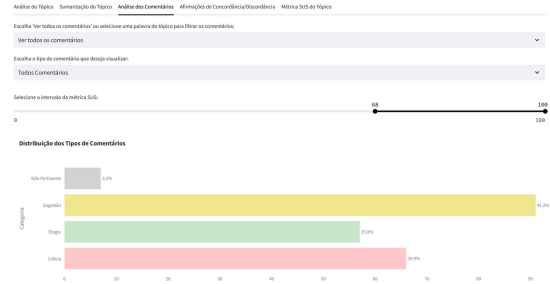


Figura 5: Comentários por categoria e intervalo SUS.

## 6 CONCLUSÕES E TRABALHOS FUTUROS

Neste artigo, apresentamos um novo *framework* que integra um amplo conjunto de técnicas de PLN (Modelagem de Tópicos, Sumarização e Análise de Sentimentos) para aprimorar os *insights* fornecidos pelo SUS, complementando-o com análises qualitativas de comentários textuais. O *framework* inclui uma interface visual que permite explorar e conduzir uma análise abrangente dos resultados das avaliações dos usuários. Aplicamos nosso *framework* em um cenário real para avaliar um sistema do governo federal brasileiro. Com mais de 1.200 avaliações, foi possível demonstrar a robustez do nosso *framework* na extração de *insights* significativos de grandes volumes de texto, permitindo uma compreensão mais aprofundada e precisa dos *feedbacks*. Pretendemos ampliar nosso *framework* incorporando a análise automática de entrevistas e gerando sugestões de mudanças a partir das análises automáticas realizadas.

## AGRADECIMENTOS

Este trabalho foi financiado por CNPq, CAPES, Fapemig e AWS.

## REFERÊNCIAS

- [1] Claudio Moisés Valiente De Andrade, Washington Cunha, Guilherme Fonseca, Ana Clara Pagano, Luana De Castro Santos, Adriana Silvina Pagano, Leonardo Rocha, and Marcos André Gonçalves. 2024. Explaining the Hardest Errors of Contextual Embedding Based Classifiers. In *Proceedings of the 28th ConLL*.
- [2] John Brooke. 1996. A quick and dirty usability scale. *Usability Evaluation in Industry* (1996), 189–194.
- [3] James Eschrich and Sarah Sterman. 2024. A Framework For Discussing LLMs as Tools for Qualitative Analysis. *ArXiv abs/2407.11198* (2024).
- [4] Jie Gao, Yuchen Guo, Gionnieve Lim, Tianqin Zhang, Zheng Zhang, Toby Jia-Jun Li, and Simon Tangi Perrault. 2024. CollabCoder: A Lower-barrier, Rigorous Workflow for Inductive Collaborative Qualitative Analysis with Large Language Models. In *Proceedings of the 2024 CHI*. <https://doi.org/10.1145/3613904.3642002>
- [5] Maarten Grootendorst. 2022. BERTopic: Neural topic modeling with a class-based TF-IDF procedure. *arXiv preprint arXiv:2203.05794* (2022).
- [6] Kai-Robin Lange, Jonas Rieger, and Carsten Jentsch. 2022. Lex2Sent: A bagging approach to unsupervised sentiment analysis. *arXiv:2209.13023* (2022).
- [7] Washington Luiz, Felipe Viegas, Rafael Alencar, Fernando Mourão, Thiago Salles, Dárlinton Carvalho, Marcos Andre Gonçalves, and Leonardo Rocha. 2018. A feature-oriented sentiment rating for mobile app reviews. In *Proc. of the WWW*.
- [8] Ali Nouri, Beatriz Cabrero-Daniel, Fredrik Torner, Hakan Sivencrona, and Christian Berger. 2024. Welcome Your New AI Teammate: On Safety Analysis by Leashing Large Language Models. In *Proceedings of the IEEE/ACM CAIN '24*.
- [9] Antônio Pereira, Pablo Cecilio, Felipe Viegas, Washington Cunha, Elisa Tuler, and Leonardo Rocha. 2022. Evaluating topic modeling pre-processing pipelines for portuguese texts. In *Proceedings of the WebMedia*. 191–201.
- [10] Chau Pham, Alexander Hoyle, Simeng Sun, Philip Resnik, and Mohit Iyyer. 2024. TopicGPT: A Prompt-based Topic Modeling Framework. In *Proceedings of the 2024 Conference of the North American Chapter of the ACL*. 2956–2984.
- [11] Felipe Viegas, Sérgio Canuto, Christian Gomes, Washington Luiz, Thiersson Rosa, Sabir Ribas, Leonardo Rocha, and Marcos André Gonçalves. 2019. CluWords: exploiting semantic word clustering representation for enhanced topic modeling. In *Proceedings of the 12a ACM WSDM*. 753–761.
- [12] Tianhao Zhang, Fu Peiguo, Jie Liu, Yihe Zhang, and Xingmei Chen. 2024. NLDesign: A UI Design Tool for Natural Language Interfaces. In *Proc. of the ACM-TURC*.