



UNIVERSIDADE DO ESTADO DA BAHIA (UNEB)
DEPARTAMENTO DE CIÊNCIAS EXATAS E DA TERRA, CAMPUS I
CURSO DE BACHARELADO EM SISTEMAS DE INFORMAÇÃO

JEOSTON ARAUJO DA CRUZ JÚNIOR

**UM PIPELINE COM BERTOPIC PARA ANÁLISE DE PUBLICAÇÕES
CIENTÍFICAS: UM ESTUDO DE CASO NO OBSERVATÓRIO DE DADOS
PÚBLICOS DE CIÊNCIA E TECNOLOGIA DA BAHIA**

SALVADOR
2024

JEOSTON ARAUJO DA CRUZ JÚNIOR

UM PIPELINE COM BERTOPIC PARA ANÁLISE DE PUBLICAÇÕES
CIENTÍFICAS: UM ESTUDO DE CASO NO OBSERVATÓRIO DE DADOS
PÚBLICOS DE CIÊNCIA E TECNOLOGIA DA BAHIA

Monografia apresentada ao Curso de Bacharelado em Sistemas de Informação do Departamento de Ciências Exatas e da Terra (DCET) - Campus I, da Universidade do Estado da Bahia (UNEB), como requisito à obtenção do grau de bacharel em Sistemas de Informação.
Área de concentração: Ciências da Computação.

Orientador: Prof. Dr. Eduardo Manuel de Freitas Jorge

SALVADOR
2024

TERMO DE ANUÊNCIA DO ORIENTADOR

Declaro para os devidos fins que li e revisei este trabalho e atesto sua qualidade como resultado final desta monografia. Confirmando que o referencial teórico apresentado é completo e suficiente para fundamentar os objetivos propostos e que a metodologia científica utilizada e os resultados finais são consistentes e com qualidade suficiente para submissão à banca examinadora final do Trabalho de Conclusão de Curso II do curso de Bacharelado em Sistemas de Informação.

Prof. Dr. Eduardo Manuel de Freitas Jorge
Orientador

JEOSTON ARAUJO DA CRUZ JÚNIOR

UM PIPELINE COM BERTOPIC PARA ANÁLISE DE PUBLICAÇÕES
CIENTÍFICAS: UM ESTUDO DE CASO NO OBSERVATÓRIO DE DADOS
PÚBLICOS DE CIÊNCIA E TECNOLOGIA DA BAHIA

Monografia apresentada ao Curso de Bacharelado em Sistemas de Informação do Departamento de Ciências Exatas e da Terra (DCET) - Campus I, da Universidade do Estado da Bahia (UNEB), como requisito à obtenção do grau de bacharel em Sistemas de Informação.
Área de concentração: Ciências da Computação.

Aprovada em: .

BANCA EXAMINADORA

Prof. Dr. Eduardo Manuel de Freitas Jorge
Orientador

Prof. Dr. Nome Completo da Pessoa
Examinador interno (DCET-I/UNEB)

Prof. Dra. Nome Completo da Pessoa
Examinador interno (DCET-I/UNEB)

AGRADECIMENTOS

Dedico este trabalho à minha amada Vovó, Cenice, cuja alma resiliente enfrentou os desafios da vida com força e graça incomparáveis. Sua memória permanece viva em mim, provando que aqueles que amamos nunca partem verdadeiramente enquanto os mantemos em nossos corações (In memoriam).

À minha mãe, Elisandra, que me deu a vida e, mesmo com suas asas cortadas pela vida, nunca hesitou em me ajudar a voar. Sua coragem e sacrifício me ensinaram o verdadeiro significado do amor incondicional e da força silenciosa.

Ao meu amor, Beatriz, que nos dias de tempestade foi o meu sol, iluminando meus caminhos e aquecendo minha alma com esperança e carinho.

A presença dessas três mulheres eternas foi o que me deu forças para seguir em frente, mesmo nos momentos mais difíceis. Aprendi que o amor é um ato de vontade. Foi através desse amor que vocês me sustentaram, me inspiraram e me ensinaram a acreditar em mim mesmo. Este trabalho é tão meu quanto de vocês, pois cada conquista minha carrega o peso do incentivo e do sacrifício de cada uma de vocês. Obrigado por nunca desistirem de mim, mesmo quando eu quase desisti.

*“Não importa o quanto a vida possa parecer difícil, há sempre algo que
você pode fazer para ter sucesso.”
(Stephen Hawking)*

RESUMO

A análise de grandes volumes de publicações científicas apresenta desafios complexos, principalmente na organização e categorização de padrões temáticos. Em resposta a esse cenário, este estudo propõe o desenvolvimento de um pipeline que combina o BERTopic e o GPT-4 para a análise de publicações científicas na plataforma SIMCC. O BERTopic é empregado para a modelagem de tópicos através do uso de embeddings contextuais, da redução de dimensionalidade com UMAP e do agrupamento com HDBSCAN. Paralelamente, o GPT-4 é utilizado para enriquecer semanticamente os clusters de tópicos identificados, gerando rótulos descritivos e precisos que complementam a modelagem. A base de dados do projeto provém do SIMCC, uma plataforma da Secretaria Estadual de Ciência, Tecnologia e Inovação da Bahia que centraliza e organiza dados de produção acadêmica de profissionais vinculados a instituições de ensino e pesquisa do estado, integrando informações de diversas fontes como Currículos Lattes, Sucupira e OpenAlex. O sistema oferece funcionalidades para o gerenciamento do conhecimento acadêmico. A integração desse pipeline à base de dados do SIMCC visa facilitar a análise e a visualização das publicações por meio de um modelo de mapeamento visual, semelhante ao WizMap, que organiza os tópicos em clusters. Essa abordagem busca aprimorar a categorização temática, contribuindo para uma compreensão mais estruturada e detalhada do acervo científico disponível na plataforma.

Palavras-chave: Processamento de Linguagem Natural; Inteligência Artificial; Modelagem de Tópicos; BERTopic; GPT-4; SIMCC; Análise de Publicações Científicas.

ABSTRACT

The analysis of large volumes of scientific publications presents complex challenges, mainly in the organization and categorization of thematic patterns. In response to this scenario, this study proposes the development of a pipeline that combines BERTopic and GPT-4 for the analysis of scientific publications on the SIMCC platform. BERTopic is used for topic modeling through the use of contextual embeddings, dimensionality reduction with UMAP, and clustering with HDBSCAN. Simultaneously, GPT-4 is utilized to semantically enrich the identified topic clusters, generating descriptive and precise labels that complement the modeling. The project's database comes from SIMCC, a platform from the State Secretariat for Science, Technology, and Innovation of Bahia, which centralizes and organizes academic production data from professionals affiliated with teaching and research institutions in the state. The system integrates information from various sources such as Lattes Curricula, Sucupira, and OpenAlex. The integration of this pipeline into the SIMCC database aims to facilitate the analysis and visualization of publications through a visual mapping model, similar to WizMap, which organizes topics into clusters. This approach seeks to improve thematic categorization, contributing to a more structured and detailed understanding of the available scientific collection.

Key-words: Natural Language Processing; Artificial Intelligence; Topic Modeling; BERTopic; GPT-4; SIMCC; Scientific Publications Analysis.

LISTA DE ILUSTRAÇÕES

Figura 1 – Arquiteturas CBOW e <i>Skip-gram</i>	17
Figura 2 – Arquitetura do modelo Transformador.	18
Figura 3 – Diferenças nas arquiteturas de pré-treinamento. BERT é bidirecional, GPT é unidirecional (da esquerda para a direita) e ELMo é uma concatenação rasa de LSTMs.	20
Figura 4 – Arquitetura de inferência do SBERT para computar similaridade. . . .	21
Figura 5 – Diagrama esquemático do <i>pipeline</i> BERTopic.	24
Figura 6 – Diagrama ilustrativo do UMAP, demonstrando a relação entre os hiperparâmetros <code>n_neighbors</code> e <code>min_dist</code> e a representação visual dos dados.	25
Figura 7 – Figura ilustrativa de um <i>dataset</i> sintético com quatro <i>clusters</i> e ruído de fundo.	26
Figura 8 – Interface da ferramenta WizMap e seus componentes principais. . . .	29
Figura 9 – Estrutura de dados <i>Quadtree</i> usada pelo WizMap para agregação multi-resolução. (A) Particionamento recursivo do espaço 2d . (B) Representação em árvore.	30

LISTA DE TABELAS

Tabela 1 – Quadro Resumo: Comparativo de Trabalhos Correlatos	33
---	----

SUMÁRIO

1	INTRODUÇÃO	12
2	REFERENCIAL TEÓRICO	14
2.1	Ciência da Informação e Análise de Publicações Científicas . .	14
2.2	Processamento de Linguagem Natural (PLN)	15
2.3	A Evolução das Representações Vetoriais em PLN	16
2.3.1	<i>Embeddings Estáticos: Limitações do Bag-of-Words</i>	16
2.3.2	<i>A Revolução dos Transformadores e o Mecanismo de Atenção</i>	18
2.3.3	<i>Embeddings Contextuais: BERT e SBERT</i>	19
2.4	Abordagens Tradicionais de Modelagem de Tópicos	21
2.5	BERTopic: Uma Abordagem Moderna	23
2.6	Visualização de Dados para Análise Científica	27
3	TRABALHOS CORRELATOS	31
3.1	Síntese Comparativa dos Trabalhos Correlatos	32
	REFERÊNCIAS	34

1 INTRODUÇÃO

O cenário da pesquisa científica global tem testemunhado um crescimento exponencial da produção científica nas últimas décadas, resultando em um vasto volume de dados que desafia os métodos tradicionais de organização e análise desse acervo. Para navegar nessa imensidão de informações, pesquisadores confiam em plataformas de busca, como *Web of Science*¹, *Scopus*² e *IEEE Xplore*³, utilizando principalmente palavras-chave. Contudo, essa abordagem de recuperação de informações é limitada pela ambiguidade e pela diversidade do léxico científico, o que frequentemente resulta em buscas que não retornam a completude esperada e na dificuldade de identificar tendências emergentes na literatura (GALLI *et al.*, 2024).

Segundo Datchanamoorthy, S e B (2023), a complexidade inerente a esses acervos e a necessidade de uma análise mais profunda têm impulsionado a aplicação de técnicas avançadas de Processamento de Linguagem Natural (PLN). Esse avanço, que une ciência da informação, Inteligência Artificial (IA) e linguística computacional, posiciona essas áreas como fundamentais na construção de soluções para a gestão do conhecimento acadêmico. Estudos como os de Mohammadi e Karami (2020) e Xie *et al.* (2020), que analisaram tendências de pesquisa em *Big Data* por meio de mineração de texto, destacam a relevância da integração de técnicas de modelagem de tópicos com modelos baseados em transformadores.

Essa arquitetura de modelos, os Transformadores, introduzida por Vaswani *et al.* (2017), revolucionou o campo da PLN com seu mecanismo de autoatenção (*self-attention*). Modelos subsequentes, como o *Bidirectional Encoder Representations from Transformers* (BERT)⁴, proposto por Devlin *et al.* (2019), passaram a capturar relações contextuais em textos com alta eficiência. A partir dessa base, surgiram os *embeddings*, representações numéricas que codificam o significado semântico de palavras e frases, superando as limitações de modelos tradicionais de *bag-of-words* e de modelagem de tópicos como o *Latent Dirichlet Allocation* (LDA) (GALLI *et al.*, 2024).

Nesse contexto, a técnica de *Bidirectional Encoder Representations Transformers for Topic Modeling* (BERTopic)⁵, proposto por Grootendorst (2022), surge como uma abordagem moderna. Seu diferencial reside na utilização dos *embeddings* contextuais de modelos como o BERT para a modelagem de tópicos. Esta técnica permite identificar tópicos de forma dinâmica e mais coesa, superando as deficiências de modelos tradicionais

¹ Disponível em: <https://access.clarivate.com/login?app=wos>.

² Disponível em: <https://www.scopus.com/home.uri>.

³ Disponível em: <https://ieeexplore.ieee.org/>.

⁴ Disponível em: https://huggingface.co/docs/transformers/model_doc/bert.

⁵ Disponível em: <https://github.com/MaartenGr/BERTopic>.

ao capturar nuances semânticas e lidar com a complexidade de textos interdisciplinares.

Este projeto de pesquisa foca no desenvolvimento de um *pipeline* computacional para o mapeamento interativo de publicações científicas. A proposta central é construir um artefato que combina a modelagem de tópicos do BERTopic Grootendorst (2022) com a ferramenta de visualização *Knowledge Map Visualization Tool* (WizMap)⁶ (WANG; HOHMAN; CHAU, 2023). O estudo de caso é aplicado ao acervo do Observatório de dados públicos de ciência e tecnologia da Bahia, que coleta informações de fontes como Currículos Lattes, Plataforma Sucupira e *OpenAlex*, e tem um papel fundamental na gestão do conhecimento científico regional.

A metodologia *Design Science Research* (DSR) é adotada como a estrutura principal deste estudo, orientando a criação deste artefato. O objetivo é transformar a base de dados textual do Observatório em um mapa de conhecimento navegável. Nessa solução, o BERTopic é empregado para extrair os padrões temáticos e o WizMap é utilizado para a exploração visual e interativa desses tópicos. Essa integração permite a identificação de temas emergentes e a compreensão da estrutura do conhecimento científico da plataforma, indo além das análises estáticas tradicionais. Espera-se que este mecanismo otimize a experiência dos usuários e contribua para a gestão estratégica da pesquisa na plataforma.

Para articular o desenvolvimento deste estudo, a monografia segue uma progressão lógica. O Capítulo 2 estabelece o Referencial Teórico, fundamentando os conceitos de Ciência da Informação, a arquitetura dos Transformadores e as técnicas de modelagem de tópicos, com foco nos componentes do BERTopic. A seguir, o Capítulo 3 analisa os Trabalhos Correlatos, contextualizando esta pesquisa frente ao estado da arte. O Capítulo 4 detalha a Metodologia (DSR), que fornece o rigor científico para a construção do artefato proposto. O Capítulo 5, cerne deste trabalho, apresenta o Projeto de Desenvolvimento, descrevendo a arquitetura completa do *pipeline*: desde a ingestão dos dados do Observatório e a modelagem com BERTopic, até a integração final com a ferramenta de visualização interativa WizMap. Por fim, o Capítulo 6 discute os Resultados Esperados e os métodos de validação aplicados a essa solução de mapeamento de conhecimento.

⁶ Disponível em: <https://github.com/poloclub/wizmap>.

2 REFERENCIAL TEÓRICO

O referencial teórico deste estudo abordará diversos aspectos cruciais relacionados à Ciência da Informação, Análise de Publicações Científicas, Processamento de Linguagem Natural (PLN), Modelagem de Tópicos, e Modelos de Linguagem de Grande Escala (LLMs).

2.1 Ciência da Informação e Análise de Publicações Científicas

A explosão da produção científica global nas últimas décadas, impulsionada pela maior acessibilidade à tecnologia e pela colaboração interdisciplinar, delineia um cenário desafiador para a área da Ciência da Informação. Como destacam Kim, Kogler e Maliphol (2024), o volume crescente de publicações dificulta a atualização contínua de pesquisadores e a identificação de áreas emergentes do conhecimento. Os autores reforçam essa problemática no resumo de seu trabalho:

A produção científica global está se expandindo exponencialmente, o que, por sua vez, exige uma melhor compreensão da ciência da ciência e, especialmente, de como as fronteiras dos campos científicos se expandem através de processos de emergência. Kim, Kogler e Maliphol (2024, Traduzido, p. 1)

Nesse contexto, estratégias tradicionais de busca baseadas em palavras-chave mostram-se limitadas, uma vez que desconsideram a complexidade semântica do léxico científico. Esse fator resulta não apenas na omissão de trabalhos relevantes, mas também na dificuldade de mapear de forma consistente o progresso em determinados campos.

Um aspecto que amplia essa complexidade é a diversidade linguística no ambiente científico. Segundo Xie *et al.* (2020), embora o inglês desempenhe papel predominante, uma parcela significativa da produção ocorre em outros idiomas. Metodologias convencionais de análise revelam-se insuficientes para o tratamento multilíngue, o que pode restringir a circulação global do conhecimento e reduzir a visibilidade de estudos relevantes.

A maioria dos estudos até hoje sobre análise de tópicos tem sido baseada em publicações em língua inglesa e tem dependido fortemente da análise de evolução de tópicos baseada em citações. [...] metodologias baseadas em citações não são adequadas para analisar relações de tópicos de pesquisa multilíngues. Xie *et al.* (2020, Traduzido, p. 1)

Diante desse cenário, técnicas contemporâneas de *Topic Modeling*, em especial aquelas fundamentadas em *embeddings*, têm sido investigadas como alternativas promissoras. De acordo com Galli *et al.* (2024), a utilização de representações densas derivadas de

modelos como o BERT potencializa a análise de grandes volumes textuais, permitindo capturar aspectos semânticos que vão além da simples coincidência lexical.

Um componente essencial para alcançar a compreensão semântica são os *embeddings* — representações numéricas que codificam o significado de palavras ou mesmo de frases — que são fundamentais no PLN para capturar relacionamentos complexos entre palavras e frases usando arquiteturas especiais conhecidas como transformadores. Galli *et al.* (2024, Traduzido, p. 2)

Essa capacidade favorece a identificação de padrões temáticos em documentos que não compartilham necessariamente o mesmo vocabulário. Métodos modernos, como o BERTopic, oferecem uma estrutura metodológica para a extração de tópicos a partir dessas representações vetoriais densas. A literatura aponta que a aplicação dessas ferramentas é particularmente relevante em textos científicos heterogêneos e multilíngues, como os encontrados em grandes repositórios de publicações científicas, dada a sua robustez em capturar nuances semânticas independentemente do idioma (XIE *et al.*, 2020).

2.2 Processamento de Linguagem Natural (PLN)

O Processamento de Linguagem Natural (PLN) é um campo multidisciplinar, situado na interseção da Inteligência Artificial, da Linguística Computacional e da Ciência da Informação. O seu objetivo central é desenvolver métodos computacionais capazes de processar, analisar, compreender e gerar a linguagem humana — seja em formato de texto ou voz — de maneira útil e análoga à humana (JURAFSKY; MARTIN, 2009)¹.

Historicamente, o PLN dependia de abordagens estatísticas e regras linguísticas manuais para modelar a linguagem (MANNING; SCHUTZE, 1999)². As técnicas de PLN são projetadas para extrair significado e estrutura de dados textuais, que são inerentemente não estruturados. Isso envolve uma série de tarefas complexas, desde a análise sintática (a estrutura gramatical) até a análise semântica (o significado por trás das palavras). Tarefas comuns incluem a classificação de textos, a tradução automática, a sumarização de documentos e, de relevância particular para este referencial, a Modelagem de Tópicos (*Topic Modeling*).

A evolução recente do campo foi impulsionada pelo *Deep Learning* (Aprendizado Profundo), que permitiu a criação de representações vetoriais de alta qualidade. Como

¹ Refere-se à obra *Speech and Language Processing*, de Daniel Jurafsky e James H. Martin. É amplamente considerado o livro-texto acadêmico padrão e a referência canônica para o ensino e estudo do Processamento de Linguagem Natural em todo o mundo.

² Refere-se à obra *Foundations of Statistical Natural Language Processing* (Manning e Schütze, 1999), considerada o trabalho seminal que consolidou as abordagens estatísticas como o padrão do PLN antes da ascensão das redes neurais profundas.

destaca Galli *et al.* (2024), o PLN moderno depende fundamentalmente da capacidade de capturar o significado contextual.

Um componente essencial para alcançar a compreensão semântica são os *embeddings* — representações numéricas que codificam o significado de palavras ou mesmo de frases — que são fundamentais no PLN para capturar relacionamentos complexos entre palavras e frases usando arquiteturas especiais conhecidas como transformadores. Galli *et al.* (2024, Traduzido, p. 2)

Dessa forma, o PLN evoluiu de uma análise baseada em contagem de palavras para uma abordagem focada na compreensão semântica, o que viabilizou os avanços em modelagem de tópicos discutidos nas seções seguintes.

2.3 A Evolução das Representações Vetoriais em PLN

O avanço no campo do PLN tem sido marcado pela busca por representações vetoriais que capturem não apenas informações sintáticas, mas também aspectos semânticos e contextuais dos textos.

2.3.1 *Embeddings Estáticos: Limitações do Bag-of-Words*

As primeiras abordagens de sucesso, como o *Word2Vec* proposto por Mikolov *et al.* (2013)³ e o *GloVe* proposto por Pennington, Socher e Manning (2014)⁴, consolidaram a noção de *embeddings*. Estes são vetores em espaços de alta dimensionalidade capazes de representar o significado aproximado de uma palavra.

O avanço conceitual desses modelos foi o de permitir que o significado semântico fosse quantificado. Em vez de tratar palavras como identificadores discretos (como em uma abordagem *bag-of-words*), os *embeddings* posicionam termos com significados similares próximos uns dos outros no espaço vetorial. Isso permite que relações semânticas sejam capturadas matematicamente, como no exemplo clássico “Rei - Homem + Mulher \approx Rainha” (MIKOLOV *et al.*, 2013). Xie *et al.* (2020) na literatura de PLN refere-se a este espaço vetorial como um “espaço semântico”.

Galli *et al.* (2024) destacam que esta capacidade de representação numérica é o alicerce da compreensão semântica no PLN moderno:

³ O *Word2Vec* (2013) foi seminal por introduzir duas arquiteturas eficientes, *Skip-gram* e *CBOW*, que aprendem vetores de palavras prevendo o contexto em que elas aparecem, baseando-se na hipótese distribucional.

⁴ O *GloVe* (2014), ou “Global Vectors”, diferencia-se por combinar as estatísticas globais de coocorrência de palavras (como o LSA) com a modelagem baseada em janelas de contexto (como o *Word2Vec*), capturando relações lineares entre palavras.

Um componente essencial para alcançar a compreensão semântica são os *embeddings* — representações numéricas que codificam o significado de palavras ou mesmo de frases — que são fundamentais no PLN para capturar relacionamentos complexos entre palavras e frases usando arquiteturas especiais [...]. Galli *et al.* (2024, Traduzido, p. 2)

O aprendizado desses vetores ocorre através do treinamento de uma rede neural rasa em uma tarefa de previsão de contexto, conforme ilustrado na Figura 1. O artigo seminal de Mikolov *et al.* (2013) propôs duas arquiteturas principais:

1. **Continuous Bag-of-Words (CBOW):** A arquitetura prevê a palavra atual (saída) com base em uma janela de palavras do contexto (entrada).
2. **Skip-gram:** A arquitetura inverte a lógica e usa a palavra atual (entrada) para prever as palavras do contexto (saída).

É importante notar que os *embeddings* não são o produto final, mas sim um subproduto do treinamento: os vetores aprendidos na camada oculta da rede (*PROJECTION* na figura) tornam-se a representação semântica da palavra, como indica Mikolov *et al.* (2013, p. 4).

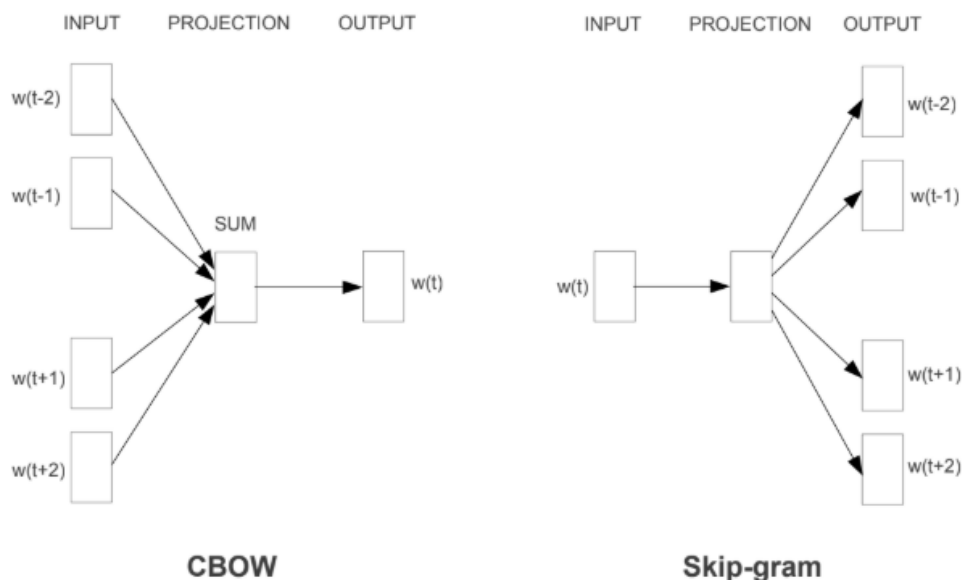


Figura 1 – Arquiteturas CBOW e *Skip-gram*.

Fonte: Mikolov *et al.* (2013, p. 5)

Embora inovadores, esses modelos apresentavam a limitação de atribuir um único vetor fixo a cada termo, independentemente do contexto de ocorrência. Por exemplo, a palavra “banco” teria a mesma representação vetorial em “banco financeiro” e “banco da praça”. Tal restrição, usualmente referida como o problema da ambiguidade do significado

da palavra (*ambiguity of word meaning*), compromete a precisão em tarefas que exigem desambiguação semântica.

2.3.2 A Revolução dos Transformadores e o Mecanismo de Atenção

A verdadeira virada de paradigma ocorreu com a introdução do modelo de Transformadores (*Transformers*), proposto por Vaswani *et al.* (2017) no artigo seminal *Attention Is All You Need*⁵. Essa arquitetura rompeu com o paradigma das RNN e convolucionais, fundamentando-se inteiramente no mecanismo de autoatenção (*self-attention*).

Através dele, o modelo atribui pesos diferenciados a *tokens* em uma sequência, permitindo processar de forma simultânea e bidirecional a totalidade do contexto textual. Essa propriedade conferiu aos modelos baseados em Transformadores a capacidade de gerar representações contextuais, um avanço significativo em relação às técnicas anteriores.

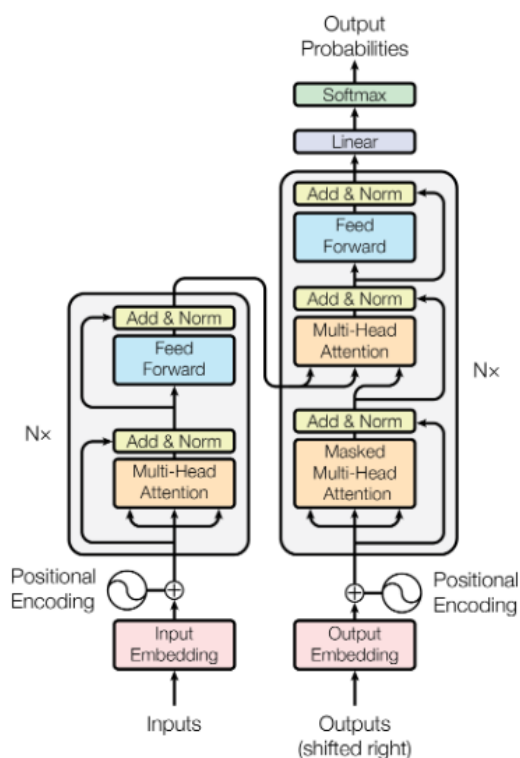


Figura 2 – Arquitetura do modelo Transformador.

Fonte: Vaswani *et al.* (2017, p. 3)

A arquitetura do Transformador, conforme apresentado na Figura 2, segue uma

⁵ Este artigo é considerado um dos trabalhos mais influentes da PLN moderna. Sua principal contribuição foi propor uma arquitetura de rede neural que dispensa totalmente as camadas recorrentes (*Recurrent Neural Network* (RNN)) e convolucionais, baseando-se unicamente em mecanismos de atenção para modelar dependências globais entre a entrada e a saída (VASWANI *et al.*, 2017, p. 1).

estrutura de codificador-decodificador (*encoder-decoder*). O lado esquerdo do diagrama representa o Codificador, enquanto o lado direito representa o Decodificador.

O **Codificador** (*Encoder*) é composto por uma pilha de N camadas idênticas (no artigo original, $N=6$). Cada camada, por sua vez, é composta por duas subcamadas principais: um mecanismo de autoatenção *multi-head* (*multi-head self-attention*) e uma rede neural *feed-forward* (rede neural de alimentação direta) simples e totalmente conectada. Conexões residuais seguidas de normalização de camada (*Add & Norm*) são aplicadas ao redor de cada subcamada.

O **Decodificador** (*Decoder*), de forma similar, é uma pilha de N camadas. Além das duas subcamadas presentes no codificador, o decodificador insere uma terceira subcamada, que realiza a atenção *multi-head* sobre a saída da pilha do codificador. Crucialmente, a subcamada de autoatenção do decodificador é “mascarada” (*Masked Multi-Head Attention*). Esse mascaramento é o que garante que a previsão para uma posição i só possa depender das saídas conhecidas em posições anteriores a i , preservando a propriedade autorregressiva do modelo.

Embora a arquitetura completa do Transformador tenha sido projetada para tarefas de transdução de sequência (como a tradução automática), foi a sua pilha de **Codificadores** (*Encoder*) que se mostrou revolucionária para tarefas de *compreensão* de linguagem. A capacidade do Codificador de processar texto de forma bidirecional e gerar representações numéricas ricas em contexto estabeleceu a base para uma nova classe de modelos focados exclusivamente na representação semântica, como será detalhado a seguir.

2.3.3 *Embeddings Contextuais: BERT e SBERT*

Sobre essa base arquitetônica, foram desenvolvidos os modelos pré-treinados, entre os quais se destaca o BERT, introduzido por Devlin *et al.* (2019). O BERT utiliza a arquitetura do Codificador (*Encoder*) do Transformador para gerar representações de linguagem.

A inovação fundamental do BERT foi o pré-treinamento bidirecional, que diferentemente de abordagens anteriores, como o *Generative Pre-trained Transformer* (GPT) de Radford e Narasimhan (2018), que utilizava um treinamento unidirecional (da esquerda para a direita), o BERT foi projetado para “pré-treinar representações profundamente bidirecionais, condicionando conjuntamente o contexto esquerdo e direito em todas as camadas” como apontam Devlin *et al.* (2019, p. 1, Traduzido).

Para alcançar essa bidirecionalidade sem que o modelo “visse a resposta”, Devlin *et al.* (2019) introduziu o objetivo do *Masked Language Model* (MLM)⁶. A Figura 3 ilustra a

⁶ O MLM é inspirado na tarefa *Cloze* (TAYLOR, 1953), onde o modelo deve prever palavras que foram omitidas (mascaradas) de uma sentença, usando o contexto de ambas as direções (esquerda e direita)

diferença fundamental entre as arquiteturas de pré-treinamento, mostrando como o BERT é capaz de processar informações de toda a sequência em todas as camadas.

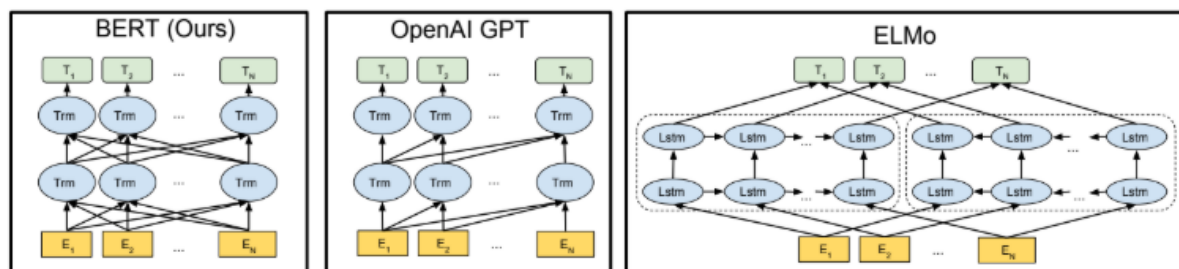


Figura 3 – Diferenças nas arquiteturas de pré-treinamento. BERT é bidirecional, GPT é unidirecional (da esquerda para a direita) e *Embeddings from Language Models* (ELMo) é uma concatenação rasa de LSTMs.

Fonte: Devlin *et al.* (2019, p. 13)

Apesar de sua eficácia em tarefas de classificação, a arquitetura do BERT puro mostrou-se inadequada para tarefas de busca de similaridade semântica ou *clustering*. Como Reimers e Gurevych (2019) explicam, o BERT “requer que ambas as sentenças sejam alimentadas na rede, o que causa um overhead computacional massivo”. Uma busca de similaridade em 10.000 sentenças, por exemplo, exigiria cerca de 50 milhões de inferências (aproximadamente 65 horas), tornando-o inviável para grandes *corpora*.

Além disso, estudos empíricos demonstraram que usar os *embeddings* “crus” do BERT (seja pela média das saídas ou pelo vetor do *token* ‘[CLS]’) produz resultados insatisfatórios, muitas vezes piores do que os *embeddings* estáticos como o *GloVe*.

Para resolver essa limitação, Reimers e Gurevych (2019) propuseram o *Sentence-BERT* (SBERT). Ele modifica o BERT pré-treinado, adicionando uma operação de *pooling* (sendo a média, *MEAN-strategy*, a mais comum) à saída do BERT para criar um *embedding* de sentença de tamanho fixo.

Crucialmente, o SBERT utiliza redes siamesas⁷ para fazer o *fine-tuning* desses *embeddings* de sentença. A Figura 4 ilustra a arquitetura de inferência do SBERT, onde duas sentenças (A e B) são processadas por redes BERT idênticas (com pesos compartilhados), gerando vetores de sentença \mathbf{u} e \mathbf{v} . Esses vetores podem, então, ser comparados eficientemente usando uma medida de similaridade, como a similaridade de cosseno (*cosine-similarity*).

para fazer a previsão (DEVLIN *et al.*, 2019, p. 1).

⁷ Redes siamesas são uma arquitetura onde duas ou mais redes neurais idênticas (com pesos compartilhados) processam entradas diferentes de forma independente. Elas são otimizadas para aprender uma função de similaridade, aproximando os vetores de saída para entradas similares e afastando-os para entradas diferentes.

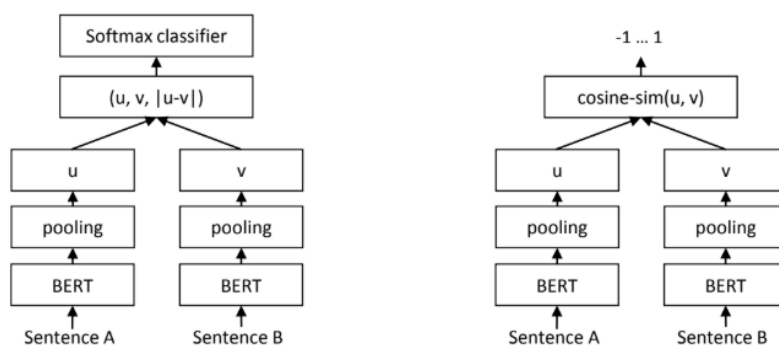


Figura 4 – Arquitetura de inferência do SBERT para computar similaridade.

Fonte: Adaptado de Reimers e Gurevych (2019, p. 3)

Reimers e Gurevych (2019) demonstraram que essa abordagem reduz o custo computacional de encontrar o par mais similar em 10.000 sentenças de 65 horas (com BERT) para cerca de 5 segundos. Essa otimização para similaridade de sentenças permite o uso eficiente desses vetores em cenários multilíngues. A utilização de modelos pré-treinados em múltiplos idiomas (como o *paraphrase-multilingual-MiniLM-L12-v2*⁸) torna-se particularmente relevante, visto que tais modelos produzem *embeddings* semanticamente consistentes mesmo em diferentes idiomas.

2.4 Abordagens Tradicionais de Modelagem de Tópicos

Com o crescimento exponencial de dados textuais e a consequente necessidade de organizar informação em larga escala, a modelagem de tópicos consolidou-se como uma técnica fundamental na área de PLN. Em termos gerais, trata-se de um conjunto de métodos estatísticos cujo objetivo é identificar estruturas semânticas latentes⁹ — denominadas *tópicos* — em coleções de documentos. Assim, essas técnicas permitem inferir distribuições temáticas que não são explicitamente observáveis, mas que emergem a partir de regularidades no uso do vocabulário.

Entre as abordagens iniciais destacam-se três marcos históricos: a *Latent Semantic Analysis* (LSA), a *Probabilistic Latent Semantic Analysis* (PLSA) e a LDA. Esses métodos não apenas moldaram a compreensão inicial sobre a representação semântica de textos, como também estabeleceram fundamentos conceituais e metodológicos que orientaram o desenvolvimento de modelos mais avançados.

A LSA, proposta por Deerwester *et al.* (1990), parte da decomposição de matrizes

⁸ Disponível em: <https://huggingface.co/sentence-transformers/paraphrase-multilingual-MiniLM-L12-v2>.

⁹ O termo “latente” significa que os tópicos não são diretamente observáveis, mas sim inferidos estatisticamente a partir dos padrões de coocorrência de palavras no *corpus*.

termo-documento por meio da técnica de *Singular Value Decomposition* (SVD)¹⁰. Nesse enquadramento, documentos e termos são projetados em um espaço vetorial de dimensionalidade reduzida, o que permite atenuar ruídos lexicais e capturar relações de similaridade latentes. Apesar de sua relevância histórica, a linearidade da LSA e sua insensibilidade a variações contextuais limitam seu desempenho em cenários onde relações semânticas complexas são determinantes (GEORGE; SUMATHY, 2023; XIE *et al.*, 2020).

Com o intuito de superar parte dessas limitações, Hofmann (1999), Hofmann (2001) introduziram a PLSA, que reformulou a representação semântica a partir de um modelo probabilístico. Nessa abordagem, cada ocorrência de palavra em um documento é modelada como proveniente de um tópico latente, de forma que a probabilidade conjunta de palavra w e documento d é expressa como:

$$P(w, d) = \sum_{z \in Z} P(z|d) P(w|z),$$

onde z representa o conjunto de tópicos latentes. Embora tenha representado um avanço em relação à LSA, a PLSA apresenta limitações notáveis, em especial no que se refere à escalabilidade: o número de parâmetros cresce linearmente com a quantidade de documentos, o que compromete sua generalização e a torna suscetível a *overfitting* (DATCHANAMOORTHY; S; B, 2023).

A evolução natural desse paradigma ocorreu com a formulação da LDA, proposta por Blei, Ng e Jordan (2003). Ao contrário da PLSA, a LDA incorpora uma camada Bayesiana por meio da utilização de distribuições de *Dirichlet* como *priors*¹¹. Essa estrutura permite regularizar o modelo e definir uma distribuição de tópicos não apenas a nível de documento, mas também a nível de *corpus*, resultando em maior robustez e interpretabilidade. A LDA parte da premissa de que cada documento é representado como uma mistura de tópicos, e cada tópico, por sua vez, é caracterizado por uma distribuição de palavras. Essa formulação tornou o modelo amplamente aplicável em diferentes domínios, como saúde pública (MIFRAH; BENLAHMAR, 2020) e eficiência energética (POLYZOS; WANG, 2022).

Apesar de sua influência, tanto a LSA quanto a PLSA e a LDA compartilham limitações estruturais. Todas operam no paradigma de *bag-of-words*¹² (saco de palavras), que ignora a ordem e o contexto local das palavras. Segundo George e Sumathy (2023),

¹⁰ A SVD é uma técnica de álgebra linear para a decomposição de matrizes que permite encontrar a melhor aproximação de uma matriz por outra de posto inferior, sendo fundamental para a redução de dimensionalidade em espaços vetoriais de termos.

¹¹ A LDA é um modelo generativo Bayesiano. O uso das distribuições de *Dirichlet* (uma distribuição de probabilidade sobre outras distribuições) permite ao modelo tratar as misturas de tópicos nos documentos e as misturas de palavras nos tópicos como variáveis aleatórias, conferindo maior robustez e melhor generalização.

¹² O *Bag-of-Words* (Saco de Palavras) é um modelo de representação de texto que ignora a ordem e a estrutura gramatical das palavras, tratando um documento apenas como um conjunto (ou multiconjunto) de suas palavras e suas frequências.

Xie *et al.* (2020), isso frequentemente conduz a representações semânticas superficiais em textos técnicos ou multilíngues. Datchanamoorthy, S e B (2023) também reitera que a sensibilidade da LDA à definição do número de tópicos (K) representa um desafio adicional: valores reduzidos podem fundir tópicos distintos em um único, enquanto valores elevados podem fragmentar temas coesos em subtemas artificiais.

A sensibilidade do LDA ao parâmetro do número de temas (K) é uma de suas desvantagens. Encontrar o valor ideal para (K) pode ser desafiador. O modelo pode simplificar excessivamente e combinar diferentes temas em um só se (K) for configurado muito baixo. No entanto, se (K) for configurado muito alto, o modelo pode se tornar muito complexo e produzir temas errôneos (DATCHANAMOORTHY; S; B, 2023, Traduzido).

Essas restrições evidenciam que, embora fundamentais, tais técnicas falham em capturar o significado contextual profundo e a ordem das palavras. Essa limitação estrutural torna-os insuficientes para tarefas que exigem uma compreensão semântica robusta, especialmente em bases textuais heterogêneas ou multilíngues onde a ambiguidade lexical é alta, destacando a necessidade de abordagens que superem o paradigma *bag-of-words*.

2.5 BERTopic: Uma Abordagem Moderna

As limitações das abordagens tradicionais de modelagem de tópicos, especialmente sua dependência do paradigma *bag-of-words* e a falha em capturar o contexto semântico, motivaram o desenvolvimento de novos métodos. Pesquisas recentes demonstraram a viabilidade de tratar a modelagem de tópicos como uma tarefa de *clustering* (agrupamento) de *embeddings*, notavelmente nos trabalhos que introduziram o *Top2Vec* (ANGELOV, 2020) e em estudos comparativos como o de Sia, Dalmia e Mielke (2020).

Nesse contexto, Grootendorst (2022) propôs o BERTopic, um modelo que estende essa abordagem de *clustering* ao introduzir uma variação do *Term Frequency-Inverse Document Frequency* (TF-IDF) baseada em classes para extrair representações de tópicos coerentes. O BERTopic não é um modelo monolítico, mas sim um *pipeline* modular que consiste em três etapas principais: 1) geração de *embeddings* de documentos, 2) *clustering* desses *embeddings* e 3) representação dos tópicos com *Class-based Term Frequency-Inverse Document Frequency* (c-TF-IDF) (GROOTENDORST, 2022, p. 1-2).

A Figura 5 ilustra o fluxo geral dessa arquitetura.

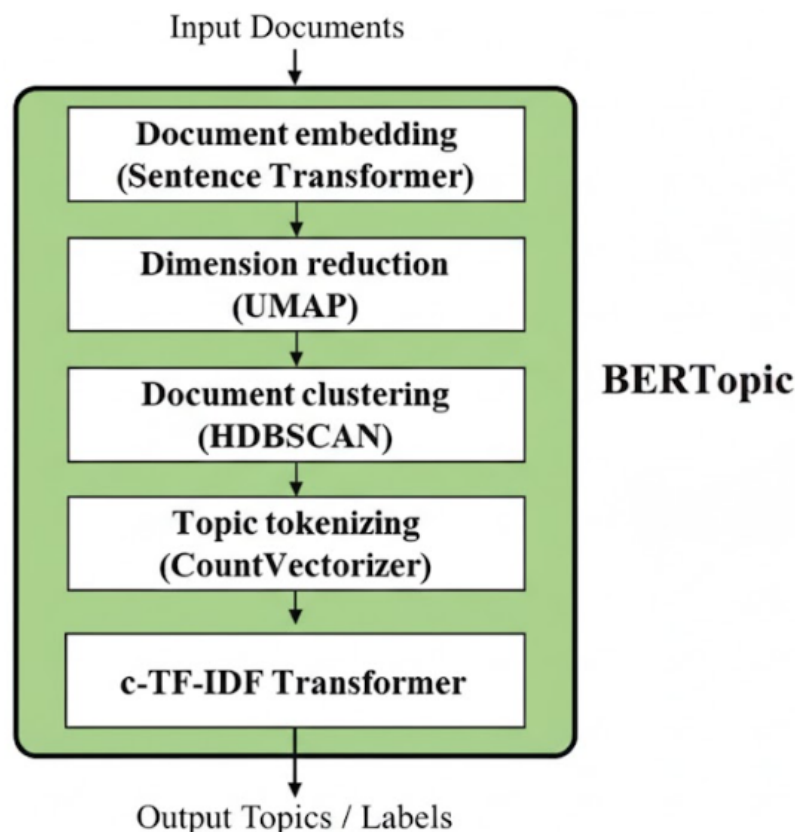


Figura 5 – Diagrama esquemático do *pipeline* BERTopic.

Fonte: Jung *et al.* (2024, p. 7, Traduzido)

Na primeira etapa, os documentos são convertidos em representações vetoriais (embeddings). O BERTopic utiliza nativamente a biblioteca SBERT (*Sentence-BERT*) proposta por Reimers e Gurevych (2019), garantindo que documentos semanticamente similares sejam posicionados próximos no espaço vetorial (GROOTENDORST, 2022, p. 2).

A segunda etapa é o *clustering* desses *embeddings* de alta dimensionalidade. Para que os algoritmos de *clustering* funcionem de forma eficiente, é necessário primeiro combater a “maldição da dimensionalidade” (*curse of dimensionality*), um fenômeno onde as distâncias entre os pontos se tornam pouco significativas em espaços com muitas dimensões (GROOTENDORST, 2022, p. 2). Para isso, o BERTopic emprega o *Uniform Manifold Approximation and Projection* (UMAP) (MCINNES; HEALY; MELVILLE, 2018).

O UMAP é uma técnica de redução de dimensionalidade que se destaca por preservar tanto a estrutura local quanto a estrutura global dos dados em um espaço de dimensão reduzida¹³ (GROOTENDORST, 2022, p. 2-3). A Figura 6 demonstra o impacto de seus dois principais hiperparâmetros.

¹³ O UMAP é fundamentado em geometria Riemanniana e topologia algébrica. Ele constrói uma representação topológica dos dados em alta dimensão e busca uma representação em baixa dimensão que tenha uma estrutura topológica o mais equivalente possível (MCINNES; HEALY; MELVILLE, 2018, p. 3-4).

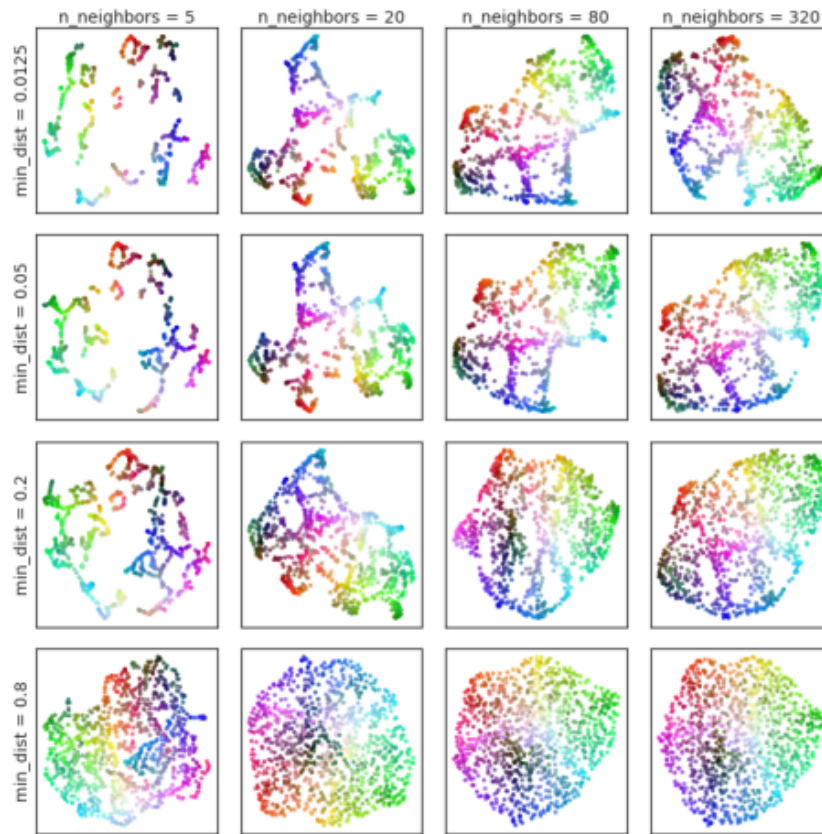


Figura 6 – Diagrama ilustrativo do UMAP, demonstrando a relação entre os hiperparâmetros `n_neighbors` e `min_dist` e a representação visual dos dados.

Fonte: McInnes, Healy e Melville (2018, p. 24)

Conforme ilustrado na Figura 6, o parâmetro `n_neighbors` (número de vizinhos) controla o equilíbrio entre a preservação da estrutura global (valores altos) e local (valores baixos). O parâmetro `min_dist` (distância mínima) ajusta a densidade dos agrupamentos, determinando a proximidade entre os pontos no espaço de baixa dimensionalidade.

Com os vetores em dimensão reduzida pelo UMAP, a etapa seguinte é o agrupamento através do *Hierarchical Density-Based Spatial Clustering of Applications with Noise* (HDBSCAN) (MCINNES; HEALY; MELVILLE, 2018). Diferentemente de métodos clássicos como o *K-Means*¹⁴, o HDBSCAN é um algoritmo baseado em densidade que não assume que os *clusters* sejam esféricos. Sua principal vantagem é a capacidade de identificar agrupamentos de densidades e formas variadas, além de sua robustez em tratar pontos que não se ajustam a nenhum padrão (ruído) como *outliers* (GROOTENDORST, 2022, p. 2-3).

A Figura 7 ilustra essa capacidade, demonstrando o tipo de desafio que o algoritmo

¹⁴ O *K-Means* é um dos algoritmos de *clustering* mais populares. Ele particiona n observações em k agrupamentos, onde cada observação pertence ao *cluster* cujo centro (média) é o mais próximo. Sua simplicidade é uma vantagem, mas ele assume *clusters* de forma esférica e sensibilidade à inicialização dos centroides (MACQUEEN, 1967).

HDBSCAN é capaz de superar, como a identificação de agrupamentos de densidades e formas variadas, além de tratar outliers de forma eficiente. Essa característica é especialmente relevante em contextos de produção científica, onde coexistem tanto publicações centrais com alta densidade de tópicos quanto trabalhos periféricos ou com temas emergentes.

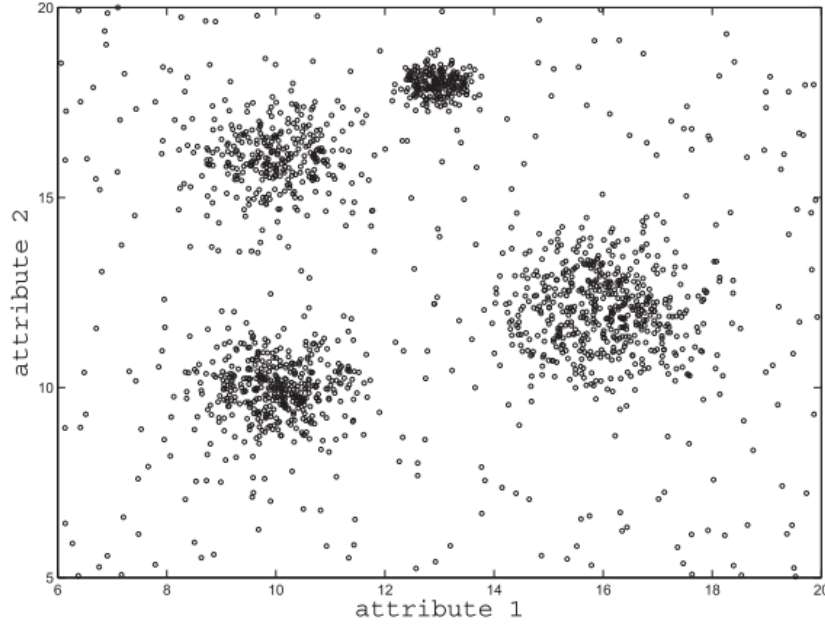


Figura 7 – Figura ilustrativa de um *dataset* sintético com quatro *clusters* e ruído de fundo.

Fonte: Campello, Moulavi e Sander (2013, p. 16)

A etapa final do *pipeline* é a geração da representação dos tópicos. Abordagens anteriores, como o *Top2Vec* (ANGELOV, 2020), baseiam-se em encontrar o centroide (o ponto médio) do *cluster* e identificar as palavras mais próximas a ele. Grootendorst (2022) argumenta que essa abordagem é falha, pois “um *cluster* nem sempre se situa dentro de uma esfera ao redor de um centroide” (GROOTENDORST, 2022, p. 1, Traduzido).

Para resolver isso, o BERTopic introduz o c-TF-IDF (*Class-based Term Frequency-Inverse Document Frequency*). A abordagem primeiro trata todos os documentos dentro de um *cluster* (tópico) como um único documento concatenado. Em seguida, modifica a fórmula padrão do TF-IDF para operar a nível de classe, e não de documento.

O TF-IDF clássico é definido por Joachims (1997) como:

$$W_{t,d} = tf_{t,d} \cdot \log \left(\frac{N}{df_t} \right) \quad (2.1)$$

onde $W_{t,d}$ é a pontuação da palavra t no documento d , $tf_{t,d}$ é a frequência da palavra t no documento d , N é o número total de documentos e df_t é o número de documentos que contêm a palavra t .

O c-TF-IDF adapta essa lógica, onde a frequência do termo (tf) é calculada para a palavra t dentro da classe c inteira (o *cluster* de documentos concatenados). A frequência

inversa do documento (*idf*) é substituída pela “frequência inversa da classe”, que mede a importância da palavra t em relação a todas as outras classes. A fórmula é então ajustada para:

$$W_{t,c} = tf_{t,c} \cdot \log \left(1 + \frac{A}{tf_t} \right) \quad (2.2)$$

onde $tf_{t,c}$ é a frequência da palavra t na classe c , A é o número médio de palavras por classe (total de palavras dividido pelo número de classes), e tf_t é a frequência total da palavra t em todas as classes (GROOTENDORST, 2022, p. 3). O resultado é uma lista de palavras que destaca os termos que são mais distintivos e representativos de um tópico específico.

Uma fraqueza teórica do c-TF-IDF é que, por ser baseado em frequência, os termos no topo da lista podem ser semanticamente redundantes (ex: “modelo”, “modelagem”). O Grootendorst (2022, p. 8) sugere que isso pode ser resolvido “aplicando *Maximal Marginal Relevance* (MMR) às n palavras principais de um tópico”.

O MMR, introduzido por Carbonell e Goldstein (1998), é uma técnica projetada especificamente para otimizar o equilíbrio entre relevância e diversidade na recuperação de informações. O algoritmo funciona de forma iterativa: ele primeiro seleciona o termo de maior relevância (maior pontuação c-TF-IDF); em seguida, para cada termo candidato subsequente, ele aplica uma penalidade com base na similaridade desse candidato com os termos já selecionados. O resultado é um conjunto de palavras-chave que não apenas representa o tema central, mas também cobre diferentes facetas semânticas desse tema, aumentando significativamente a interpretabilidade humana.

Contudo, mesmo com representações de tópicos robustas e rótulos semanticamente diversos, analisar a estrutura latente e as inter-relações de centenas de tópicos em um *corpus* massivo permanece um desafio. A geração de um modelo de tópicos é apenas a primeira etapa; a descoberta de conhecimento emerge da capacidade de explorar esses resultados de forma intuitiva. Isso destaca a necessidade de técnicas que superem listas estáticas e permitam uma análise exploratória, um desafio que é central no campo da Visualização Científica e de Dados.

2.6 Visualização de Dados para Análise Científica

A geração de modelos de tópicos e *embeddings*, conforme discutido nas seções anteriores, produz representações vetoriais de alta dimensionalidade que capturam a semântica do domínio. No entanto, a interpretação e o uso prático desses *embeddings* representam um desafio significativo, dada a sua “opacidade, alta dimensionalidade e o grande tamanho dos conjuntos de dados modernos” (WANG; HOHMAN; CHAU, 2023, p. 1, Traduzido).

Para tornar esses vetores complexos inteligíveis, pesquisadores frequentemente aplicam técnicas de redução de dimensionalidade, como o UMAP (MCINNES; HEALY; MELVILLE, 2018) ou o *t-distributed Stochastic Neighbor Embedding* (t-SNE) (MAATEN; HINTON, 2008), para projetar os *embeddings* em um espaço bidimensional (**2d**) ou tridimensional (**3d**). Embora essa projeção permita a visualização dos dados em um gráfico de dispersão (*scatter plot*), a análise em larga escala permanece um desafio: em conjuntos de dados com milhões de pontos, “é exaustivo ou mesmo implausível inspecionar os dados ponto a ponto para entender a estrutura global” (WANG; HOHMAN; CHAU, 2023, p. 2, Traduzido).

Abordagens alternativas, como gráficos de contorno (*contour plots*), podem resumir a distribuição global, mas “restringem a exploração das estruturas locais de um *embedding*” (WANG; HOHMAN; CHAU, 2023, p. 2, Traduzido). Para preencher a lacuna entre a visão global (contornos) e a exploração local (pontos), ferramentas de visualização interativa tornam-se essenciais.

Neste contexto, surge o WizMap¹⁵, “uma ferramenta de visualização interativa escalável que capacita pesquisadores e especialistas de domínio a explorar e interpretar *embeddings* com milhões de pontos” (WANG; HOHMAN; CHAU, 2023, p. 2, Traduzido). A ferramenta emprega um “design de interação familiar semelhante a um mapa” (*map-like interaction design*), permitindo que o usuário navegue pelo espaço semântico com ações de *pan* e *zoom*.

A interface do WizMap, ilustrada na Figura 8, é dividida em três componentes principais: (A) A Visão de Mapa (*Map View*), que integra as camadas de visualização; (B) O Painel de Busca (*Search Panel*), que permite a filtragem por texto; e (C) O Painel de Controle (*Control Panel*), para customização da visualização (WANG; HOHMAN; CHAU, 2023, p. 1).

¹⁵ O repositório de código aberto do WizMap está disponível em: <https://github.com/poloclub/wizmap>. Uma demonstração interativa da ferramenta, analisando artigos científicos (a que o *dataset* deste estudo de caso se assemelha), pode ser acessada em: <https://fossil-explorer.com/wizmap/?dataURL=https://fossil-explorer.com/wizmap/data/dsm-papers/data.ndjson&gridURL=https://fossil-explorer.com/wizmap/data/dsm-papers/grid.json>

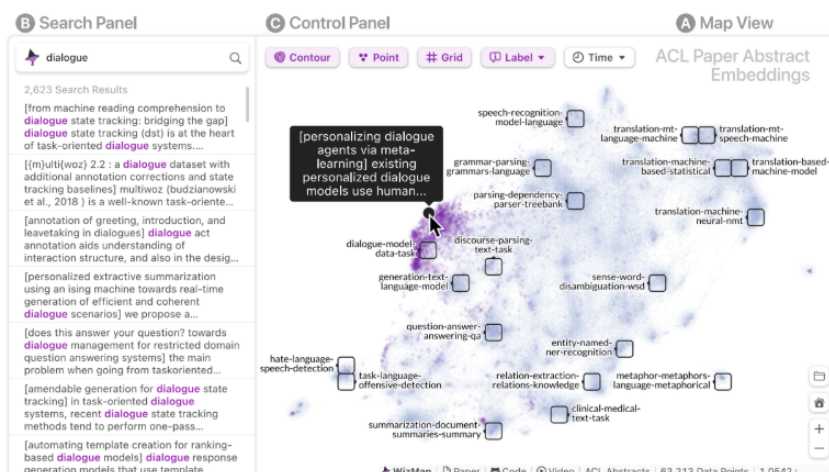


Figura 8 – Interface da ferramenta WizMap e seus componentes principais.

Fonte: Wang, Hohman e Chau (2023, p. 1)

A principal inovação do WizMap é a sua capacidade de escalar para milhões de pontos diretamente no navegador do usuário, sem a necessidade de servidores dedicados. Isso é alcançado através do uso de tecnologias web modernas, como *Web Graphics Library* (WebGL)¹⁶ para renderização gráfica, *Web Workers*¹⁷ para paralelização, e a *Streams API*¹⁸ para o carregamento de dados (WANG; HOHMAN; CHAU, 2023, p. 2, 4).

A “Visão de Mapa” (*Map View*), sua interface primária, integra três camadas de visualização (WANG; HOHMAN; CHAU, 2023, p. 4):

1. **Contorno de Distribuição:** Utiliza *Kernel Density Estimation* (KDE) para fornecer uma visão geral da estrutura global e das áreas de alta densidade.
2. **Gráfico de Dispersão (*Scatter Plot*):** Permite a investigação de *embeddings* individuais em nível local.
3. **Rótulos Multi-Resolução:** Permite uma interpretação contextual em diferentes níveis de granularidade.

Para implementar os Rótulos Multi-Resolução, o WizMap utiliza uma estrutura de dados *quadtree*, conforme detalhado na Figura 9. O *quadtree* particiona recursivamente o espaço **2d** (A) em quadrantes, que são representados como nós em uma árvore (B). A ferramenta então agrega as informações de baixo para cima, permitindo que os rótulos se “ajustem em resolução à medida que os usuários aumentam o *zoom*” (WANG; HOHMAN; CHAU, 2023, p. 2-3).

¹⁶ Disponível em: https://developer.mozilla.org/pt-BR/docs/Web/API/WebGL_API

¹⁷ Disponível em: https://developer.mozilla.org/pt-BR/docs/Web/API/Web_Workers_API.

¹⁸ Disponível em: https://developer.mozilla.org/pt-BR/docs/Web/API/Streams_API.

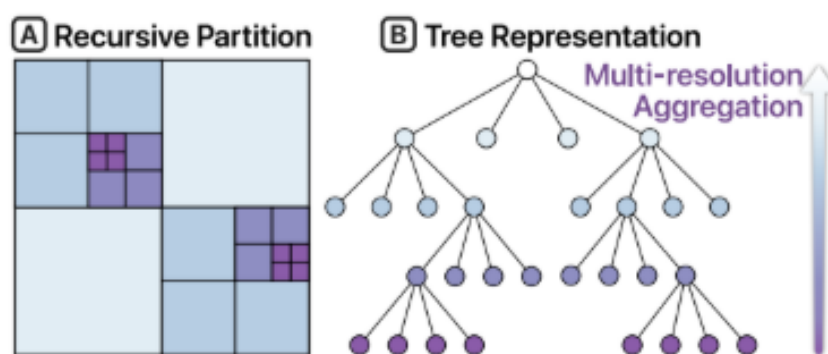


Figura 9 – Estrutura de dados *Quadtree* usada pelo WizMap para agregação multi-resolução. (A) Particionamento recursivo do espaço **2d**. (B) Representação em árvore.

Fonte: Wang, Hohman e Chau (2023, p. 3)

Ferramentas como o WizMap, que combinam redução de dimensionalidade com interfaces interativas de multi-resolução e busca semântica, são fundamentais para traduzir a saída matemática de modelos como o BERTopic em mapas de conhecimento navegáveis, facilitando a descoberta de padrões latentes em grandes base de dados textuais.

3 TRABALHOS CORRELATOS

Este capítulo apresenta uma revisão de literatura conduzida para fundamentar o estudo de metodologias de modelagem de tópicos baseadas em *embeddings*, focando em sua aplicação e avaliação no domínio da análise de tendências científicas.

Estudos comparativos recentes têm se dedicado a avaliar a eficácia de modelos de tópicos modernos frente a abordagens tradicionais. A pesquisa de Jung *et al.* (2024), por exemplo, apresenta uma análise comparativa entre métodos como LDA, *Non-negative Matrix Factorization* (NMF) e o BERTopic, aplicando-os a dados acadêmicos e de mídia. Os autores concluíram que o BERTopic, que combina *embeddings* de texto com técnicas de redução de dimensionalidade e clusterização, “demonstrou predominância em diversidade e coesão de tópicos” (JUNG *et al.*, 2024, p. 27). Essa capacidade de capturar contextos semânticos complexos, superando a abordagem de *bag-of-words* do LDA, é uma capacidade relevante para a análise de produção científica interdisciplinar.

De forma similar, Kim, Kogler e Maliphol (2024) propõem uma combinação de análise de redes e BERTopic para identificar a emergência de campos científicos interdisciplinares. O estudo valida o BERTopic como uma técnica de *embedded topic modeling* (modelagem de tópicos embarcada) que, ao contrário das abordagens baseadas em frequência, “permite considerar o conhecimento contextual de grandes conjuntos de dados de texto” (KIM; KOGLER; MALIPHOL, 2024, p. 3, Traduzido). A arquitetura empregada utiliza os componentes de *embeddings* BERT, UMAP, HDBSCAN e c-TF-IDF.

A aplicação do BERTopic para a análise de publicações científicas, especificamente para a triagem de revisões sistemáticas, também foi explorada por Galli *et al.* (2024). O estudo aplicou o *pipeline* (SBERT/UMAP/HDBSCAN/c-TF-IDF) em *datasets* da literatura médica e concluiu que a ferramenta foi eficaz na segmentação e filtragem de artigos irrelevantes, reduzindo a carga de trabalho manual (GALLI *et al.*, 2024, p. 1, 18). Notavelmente, o trabalho também identificou a representação de tópicos padrão do c-TF-IDF como “muitas vezes obscura” (GALLI *et al.*, 2024, p. 6, Traduzido).

Enquanto os estudos anteriores comparam o BERTopic com modelos clássicos, Gerasimenko *et al.* (2023) o utilizam como *baseline* para uma nova técnica de detecção de tendências científicas em tempo real. A pesquisa oferece uma análise detalhada do desempenho do BERTopic, concluindo que o modelo apresenta alta performance na distinção de documentos por tópicos, extraindo 90 dos 91 tópicos de tendência. No entanto, o estudo também identifica que o BERTopic “tem muita dificuldade na extração de palavras-chave” (GERASIMENKO *et al.*, 2023, p. 10, Traduzido), indicando um desempenho diferenciado entre as tarefas de *clustering* de documentos e de representação

de tópicos.

A arquitetura do BERTopic é também modular, permitindo a exploração de diferentes configurações para otimizar os resultados, um ponto investigado por Wijanto, Widiastuti e Yong (2024). Em seu trabalho, os autores exploraram o ajuste de hiperparâmetros em modelos baseados em BERT, testando combinações variadas de modelos de *embedding* (como *RoBERTa* e SBERT), técnicas de redução de dimensionalidade (UMAP e *Principal Component Analysis* (PCA)) e algoritmos de clusterização (*K-Means* e HDBSCAN). O estudo reforça a importância da seleção criteriosa de cada componente do *pipeline* para garantir a geração de tópicos coerentes e interpretáveis, sendo uma configuração validada o uso de SBERT, UMAP e HDBSCAN (GROOTENDORST, 2022) para documentos heterogêneos.

A literatura também aponta para a validação de *pipelines* coesos de modelagem de tópicos para análise bibliométrica e visualização. Meng *et al.* (2024), por exemplo, propõem uma metodologia que utiliza *BERTopic* para mapear a evolução da pesquisa científica em um grande volume de publicações. O trabalho de Meng *et al.* (2024) culmina no desenvolvimento de uma plataforma *web* de análise bibliométrica para visualização de redes e tópicos, validando a aplicação de *pipelines* de modelagem semântica como base para ferramentas de exploração interativa.

Em suma, a análise dos trabalhos correlatos indica que o BERTopic é uma ferramenta validada pela literatura recente para a análise de publicações científicas. A literatura confirma sua predominância sobre métodos tradicionais em métricas de coerência (JUNG *et al.*, 2024) e sua capacidade de usar contexto semântico como apontam Kim, Kogler e Maliphol (2024), Galli *et al.* (2024). Também aponta para a importância de sua modularidade (WIJANTO; WIDIASTUTI; YONG, 2024) e para um desempenho diferenciado entre a clusterização de documentos (onde é forte) e a extração de palavras-chave (onde é mais fraco) (GERASIMENKO *et al.*, 2023; GALLI *et al.*, 2024). Por fim, a literatura valida o uso de *pipelines* de modelagem como base para o desenvolvimento de plataformas de visualização interativa (MENG *et al.*, 2024).

3.1 Síntese Comparativa dos Trabalhos Correlatos

A fim de consolidar a análise da literatura e posicionar de forma clara o estado da arte, o quadro a seguir (Quadro 1) apresenta uma síntese comparativa dos trabalhos correlatos discutidos. A comparação é estruturada com base em critérios essenciais, como o objetivo principal de cada estudo, o *pipeline* metodológico empregado e as tecnologias de *embedding*. Essa estrutura permite visualizar as sinergias e as particularidades de cada abordagem.

Tabela 1 – Quadro Resumo: Comparativo de Trabalhos Correlatos

Referência	Objetivo Principal	Pipeline/Método Utilizado	Modelo de Embedding	Relação com o Estado da Arte
Jung <i>et al.</i> (2024)	Comparar o desempenho de modelos de tópicos (LDA, NMF, BERTopic) em textos acadêmicos e de notícias sobre LLMs (JUNG <i>et al.</i> , 2024).	Análise comparativa de métricas de coerência e diversidade dos tópicos gerados (JUNG <i>et al.</i> , 2024).	SBERT (implícito no BERTopic) (JUNG <i>et al.</i> , 2024).	Estabelece o BERTopic como uma ferramenta superior aos métodos tradicionais (LDA, NMF) para a análise de textos acadêmicos (JUNG <i>et al.</i> , 2024).
Kim, Kogler e Maliphol (2024)	Identificar a emergência de ciência interdisciplinar em metadados de publicações científicas (KIM; KOGLER; MALIPHOL, 2024).	Combinação de análise de redes de coocorrência (Etapa 1) e modelagem de tópicos com BERTopic (Etapa 2) (KIM; KOGLER; MALIPHOL, 2024).	BERT (usado para <i>embedding vectorization</i>) (KIM; KOGLER; MALIPHOL, 2024).	Valida o BERTopic como ferramenta superior às abordagens baseadas em frequência, por usar “conhecimento contextual”, para analisar publicações científicas (KIM; KOGLER; MALIPHOL, 2024).
Galli <i>et al.</i> (2024)	Explorar como o BERTopic pode ser aplicado para acelerar a triagem de literatura em revisões sistemáticas de publicações científicas (GALLI <i>et al.</i> , 2024).	Pipeline BERTopic padrão (SBERT → UMAP → HDBSCAN → c-TF-IDF) para identificar e filtrar <i>clusters</i> irrelevantes (GALLI <i>et al.</i> , 2024, p. 4).	‘all-mpnet-base-v2’ (SBERT) (GALLI <i>et al.</i> , 2024, p. 4).	Valida o <i>pipeline</i> SBERT/UMAP/HDBSCAN para analisar publicações científicas e corrobora que os rótulos de c-TF-IDF são “muitas vezes obscuros” (GALLI <i>et al.</i> , 2024, p. 6).
Gerasimenko <i>et al.</i> (2023)	Extrair tópicos de tendências científicas (“trend topics”) em tempo real a partir de publicações (GERASIMENKO <i>et al.</i> , 2023).	Propõe um modelo ARTM incremental e o compara com <i>baselines</i> , incluindo PLSA, LDA e BERTopic (GERASIMENKO <i>et al.</i> , 2023).	Sentence-Transformers (para o <i>baseline</i> BERTopic) (GERASIMENKO <i>et al.</i> , 2023).	Fornecer uma análise comparativa do BERTopic, destacando sua alta performance em clusterização de documentos e sua fraqueza na extração de palavras-chave (GERASIMENKO <i>et al.</i> , 2023).
Meng <i>et al.</i> (2024)	Mapear a evolução de um campo de pesquisa científica utilizando uma abordagem integrada de modelagem de tópicos e uma plataforma web (MENG <i>et al.</i> , 2024).	Pipeline integrado: 1. Geração de embeddings (via API de LLM); 2. Clusterização e modelagem com BERTopic; 3. Plataforma de visualização (MENG <i>et al.</i> , 2024, p. 3-4).	GPT-3.5 (text-embedding-ada-002) (MENG <i>et al.</i> , 2024, p. 4).	Valida a aplicação de um <i>pipeline</i> de modelagem de tópicos como base para uma plataforma <i>web</i> de visualização, um objetivo relevante para a análise de grandes <i>corpora</i> (MENG <i>et al.</i> , 2024, p. 18).

REFERÊNCIAS

- ANGELOV, D. **Top2Vec: Distributed Representations of Topics**. 2020. Disponível em: <https://arxiv.org/abs/2008.09470>. Citado nas páginas 23 e 26.
- BLEI, D. M.; NG, A. Y.; JORDAN, M. I. Latent dirichlet allocation. **Journal of Machine Learning Research**, MIT Press, v. 3, n. Jan, p. 993–1022, 2003. ISSN 1532-4435. Disponível em: <http://jmlr.org/papers/v3/blei03a.html>. Citado na página 22.
- CAMPELLO, R. J. G. B.; MOULAVI, D.; SANDER, J. Density-based clustering based on hierarchical density estimates. In: PEI, J. *et al.* (Ed.). **Advances in Knowledge Discovery and Data Mining**. Berlin, Heidelberg: Springer Berlin Heidelberg, 2013. p. 160–172. Citado na página 26.
- CARBONELL, J. G.; GOLDSTEIN, J. The use of MMR, diversity-based reranking for reordering documents and producing summaries. In: **Proceedings of the 21st Annual International ACM SIGIR Conference on Research and Development in Information Retrieval**. Melbourne, Australia: ACM, 1998. p. 335–336. Disponível em: <https://doi.org/10.1145/290941.291025>. Citado na página 27.
- DATCHANAMOORTHY, K.; S, A. M. G.; B, P. Text mining: Clustering using bert and probabilistic topic modeling. **Social Informatics Journal**, 2023. Disponível em: <https://api.semanticscholar.org/CorpusID:267122800>. Citado nas páginas 12, 22 e 23.
- DEERWESTER, S. *et al.* Indexing by latent semantic analysis. **Journal of the American Society for Information Science**, v. 41, n. 6, p. 391–407, 1990. Disponível em: <https://asistdl.onlinelibrary.wiley.com/doi/abs/10.1002/%28SICI%291097-4571%28199009%2941%3A6%3C391%3A%3AAID-ASI1%3E3.0.CO%3B2-9>. Citado na página 21.
- DEVLIN, J. *et al.* **BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding**. 2019. Disponível em: <https://arxiv.org/abs/1810.04805>. Citado nas páginas 12, 19 e 20.
- GALLI, C. *et al.* Topic modeling for faster literature screening using transformer-based embeddings. **Metrics**, v. 1, n. 1, 2024. ISSN 3042-5042. Disponível em: <https://www.mdpi.com/3042-5042/1/1/2>. Citado nas páginas 12, 14, 15, 16, 17, 31, 32 e 33.
- GEORGE, L.; SUMATHY, P. An integrated clustering and bert framework for improved topic modeling. **International Journal of Information Technology**, v. 15, n. 4, p. 2187–2195, 2023. Disponível em: <https://doi.org/10.1007/s41870-023-01268-w>. Citado nas páginas 22 e 23.
- GERASIMENKO, N. *et al.* Incremental topic modeling for scientific trend topics extraction. In: **Computational Linguistics and Intellectual Technologies: Proceedings of the International Conference "Dialogue 2023"**. [S.l.: s.n.], 2023. p. 88–103. Citado nas páginas 31, 32 e 33.

GROOTENDORST, M. **BERTopic: Neural topic modeling with a class-based TF-IDF procedure**. 2022. Disponível em: <https://arxiv.org/abs/2203.05794>. Citado nas páginas 12, 13, 23, 24, 25, 26, 27 e 32.

HOFMANN, T. Probabilistic latent semantic indexing. In: **Proceedings of the 22nd Annual International ACM SIGIR Conference on Research and Development in Information Retrieval**. Berkeley, CA, USA: ACM Press, 1999. p. 50–57. ISBN 1-58113-096-1. Citado na página 22.

HOFMANN, T. Hofmann, t.: Unsupervised learning by probabilistic latent semantic analysis. *machine learning* 42(1-2), 177-196. **Machine Learning**, v. 42, p. 177–196, 01 2001. Citado na página 22.

JOACHIMS, T. A probabilistic analysis of the rocchio algorithm with tfidf for text categorization. In: . [S.l.: s.n.], 1997. p. 143–151. Citado na página 26.

JUNG, H. S. *et al.* Expansive data, extensive model: Investigating discussion topics around llm through unsupervised machine learning in academic papers and news. **PLOS ONE**, Public Library of Science, v. 19, n. 5, p. 1–18, 05 2024. Disponível em: <https://doi.org/10.1371/journal.pone.0304680>. Citado nas páginas 24, 31, 32 e 33.

JURAFSKY, D.; MARTIN, J. **Speech and Language Processing: An Introduction to Natural Language Processing, Computational Linguistics, and Speech Recognition**. Pearson Prentice Hall, 2009. (Prentice Hall series in artificial intelligence). ISBN 9780131873216. Disponível em: <https://books.google.com.br/books?id=fZmj5UNK8AQC>. Citado na página 15.

KIM, K.; KOGLER, D. F.; MALIPHOL, S. Identifying interdisciplinary emergence in the science of science: combination of network analysis and BERTopic. **Palgrave Communications**, v. 11, n. 1, p. 1–15, December 2024. Disponível em: https://ideas.repec.org/a/pal/palcom/v11y2024i1d10.1057_s41599-024-03044-y.html. Citado nas páginas 14, 31, 32 e 33.

MAATEN, L. van der; HINTON, G. Visualizing data using t-sne. **Journal of Machine Learning Research**, v. 9, n. 86, p. 2579–2605, 2008. Disponível em: <http://jmlr.org/papers/v9/vandemaaten08a.html>. Citado na página 28.

MACQUEEN, J. Some methods for classification and analysis of multivariate observations. In: **Proceedings of the 5th Berkeley Symposium on Mathematical Statistics and Probability**. Berkeley: University of California Press, 1967. (Statistics, v. 1), p. 281–297. Disponível em: <http://projecteuclid.org/euclid.bsmsp/1200512992>. Citado na página 25.

MANNING, C.; SCHUTZE, H. **Foundations of Statistical Natural Language Processing**. MIT Press, 1999. (Foundations of Statistical Natural Language Processing). ISBN 9780262133609. Disponível em: <https://books.google.com.br/books?id=YiFDxbEX3SUC>. Citado na página 15.

MCINNES, L.; HEALY, J.; MELVILLE, J. **UMAP: Uniform Manifold Approximation and Projection for Dimension Reduction**. 2018. Disponível em: <https://arxiv.org/abs/1802.03426>. Citado nas páginas 24, 25 e 28.

MENG, F. *et al.* Demand-side energy management reimaged: A comprehensive literature analysis leveraging large language models. **Energy**, v. 291, p. 130303, 2024. ISSN 0360-5442. Disponível em: <https://www.sciencedirect.com/science/article/pii/S0360544224000744>. Citado nas páginas 32 e 33.

MIFRAH, S.; BENLAHMAR, E. H. Topic modeling coherence: A comparative study between lda and nmf models using covid'19 corpus. **International Journal of Advanced Trends in Computer Science and Engineering**, 08 2020. Citado na página 22.

MIKOLOV, T. *et al.* **Efficient Estimation of Word Representations in Vector Space**. 2013. Disponível em: <https://arxiv.org/abs/1301.3781>. Citado nas páginas 16 e 17.

MOHAMMADI, E.; KARAMI, A. Exploring research trends in big data across disciplines: A text mining analysis. **Journal of Information Science**, v. 48, 06 2020. Citado na página 12.

PENNINGTON, J.; SOCHER, R.; MANNING, C. GloVe: Global vectors for word representation. In: MOSCHITTI, A.; PANG, B.; DAELEMANS, W. (Ed.). **Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)**. Doha, Qatar: Association for Computational Linguistics, 2014. p. 1532–1543. Disponível em: <https://aclanthology.org/D14-1162>. Citado na página 16.

POLYZOS, E.; WANG, F. Twitter and market efficiency in energy markets: Evidence using lda clustered topic extraction. **Energy Economics**, v. 114, p. 106264, 2022. ISSN 0140-9883. Disponível em: <https://www.sciencedirect.com/science/article/pii/S0140988322004017>. Citado na página 22.

RADFORD, A.; NARASIMHAN, K. Improving language understanding by generative pre-training. In: . [s.n.], 2018. Disponível em: <https://api.semanticscholar.org/CorpusID:49313245>. Citado na página 19.

REIMERS, N.; GUREVYCH, I. **Sentence-BERT: Sentence Embeddings using Siamese BERT-Networks**. 2019. Disponível em: <https://arxiv.org/abs/1908.10084>. Citado nas páginas 20, 21 e 24.

SIA, S.; DALMIA, A.; MIELKE, S. J. **Tired of Topic Models? Clusters of Pretrained Word Embeddings Make for Fast and Good Topics too!** 2020. Disponível em: <https://arxiv.org/abs/2004.14914>. Citado na página 23.

TAYLOR, W. L. Cloze procedure: A new tool for measuring readability. **Journalism Quarterly**, v. 30, p. 415–433, 1953. Citado na página 19.

VASWANI, A. *et al.* **Attention Is All You Need**. 2017. Disponível em: <https://arxiv.org/abs/1706.03762>. Citado nas páginas 12 e 18.

WANG, Z. J.; HOHMAN, F.; CHAU, D. H. **WizMap: Scalable Interactive Visualization for Exploring Large Machine Learning Embeddings**. 2023. Disponível em: <https://arxiv.org/abs/2306.09328>. Citado nas páginas 13, 27, 28, 29 e 30.

WIJANTO, M. C.; WIDIASTUTI, I.; YONG, H.-S. Topic modeling for scientific articles: Exploring optimal hyperparameter tuning in bert. **International Journal on Advanced Science, Engineering and Information Technology**, v. 14, n. 3, p. 912–919, Jun. 2024.

Disponível em: <https://ijaseit.insightsociety.org/index.php/ijaseit/article/view/19347>. Citado na página 32.

XIE, Q. *et al.* Monolingual and multilingual topic analysis using lda and bert embeddings. **Journal of Informetrics**, v. 14, n. 3, p. 101055, 2020. ISSN 1751-1577. Disponível em: <https://www.sciencedirect.com/science/article/pii/S1751157719305127>. Citado nas páginas 12, 14, 15, 16, 22 e 23.