

UNIVERSIDADE DE SÃO PAULO
INSTITUTO DE MATEMÁTICA E ESTATÍSTICA
BACHARELADO EM CIÊNCIA DA COMPUTAÇÃO

**Análise de Textos Históricos utilizando
LLMs e Modelagem de Tópicos**

João Pedro Lukasavicus Silva

MONOGRAFIA FINAL

MAC 499 — TRABALHO DE
FORMATURA SUPERVISIONADO

Supervisor: Mateus Espadoto

São Paulo

*O conteúdo deste trabalho é publicado sob a licença CC BY 4.0
(Creative Commons Attribution 4.0 International License)*

*Esta seção é opcional e fica numa página separada;
ela pode ser usada para uma dedicatória ou epígrafe.*

[illegible]

Resumo

João Pedro Lukasavicus Silva. **Análise de Textos Históricos utilizando LLMs e Modelagem de Tópicos**. Monografia (Bacharelado). Instituto de Matemática e Estatística, Universidade de São Paulo, São Paulo.

[illegible]

Palavras-chave: Redes neurais. LLM. Grandes modelos de linguagem. Modelagem de tópicos. Textos históricos. Isidoro de Sevilha. Etimologias.

Abstract

João Pedro Lukasavicus Silva. **Analysis of Historical Texts using LLMs and Topic Modeling: a subtitle**. Capstone Project Report (Bachelor). Institute of Mathematics and Statistics, University of São Paulo, São Paulo.

[illegible]

Keywords: Neural Networks. LLM. Large Language Models. Topic Modeling. Historical Texts. Isidore of Seville. Etymologies.

Lista de abreviaturas

LLM	Large Language Model
MLP	Multi-layer Perceptron
PDF	Portable Document Format
HTML	Hyper-text Markup Language
NLP	Natural Language Processing
URL	Localizador Uniforme de Recursos (<i>Uniform Resource Locator</i>)
IME	Instituto de Matemática e Estatística
USP	Universidade de São Paulo

Lista de figuras

2.1	Representação do espaço latente gerado pelo LSA.	6
2.2	Arquiteturas CBOW e Skip-gram.	6
2.3	RNNSearch.	7
2.4	Arquitetura de um Transformer.	8
2.5	Mecanismos de atenção nos Transformers.	9

Lista de tabelas

3.1	Parâmetros usados para o UMAP	17
3.2	Parâmetros usados para o HDBSCAN	18

Lista de programas

Sumário

1	Introdução	1
1.1	Objetivos	1
1.2	Contribuições	2
1.3	Organização	2
2	Background	3
2.1	Redes Neurais e Grandes Modelos de Linguagem	3
2.2	Evolução das redes neurais em PLN	5
2.3	Transformers e desenvolvimentos subsequentes	8
2.3.1	Arquitetura do modelo Transformer	8
2.3.2	Atenção	9
2.3.3	BERT, SBERT e derivados	10
2.4	Modelagem de Tópicos	10
2.4.1	BERTopic	10
2.4.2	UMAP	11
2.4.3	HDBSCAN	11
2.4.4	Representação dos tópicos	11
2.5	Obras e autores estudados	11
3	Experimentos	15
3.1	Conjunto de dados	15
3.2	Preparação dos dados	15
3.2.1	Limpeza	15
3.2.2	Pré-processamento	16
3.2.3	<i>Stop words</i> , <i>stemming</i> e lematização	16
3.2.4	Geração de <i>embeddings</i>	16
3.2.5	Conjuntos de dados	17
3.3	Persistência	17

3.4	Experimentos	17
3.4.1	Redução de dimensionalidade	17
3.4.2	Clustering	18
4	Resultados	21
5	Conclusão	23
	Referências	25

Capítulo 1

Introdução

Tradicionalmente, pesquisadores das áreas de humanidades como história, filosofia, entre outras, dependem da análise de grandes quantidades de fontes para a realização do seu trabalho. Esta análise tipicamente é realizada de forma manual, e com grande custo, tanto em termos de tempo quanto de recursos humanos, custo este que pode se tornar um fator limitador da produção científica desses pesquisadores.

A área de processamento de linguagem natural (NLP) possui métodos bem estabelecidos para análise semântica de textos, como por exemplo, modelagem de tópicos (DEERWESTER *et al.*, 1990; BLEI *et al.*, 2003; JELODAR *et al.*, 2019), que visa agrupar textos que tratam de assuntos semelhantes. Com o surgimento dos grandes modelos de linguagem (LLM) baseados em transformers (VASWANI *et al.*, 2017; DEVLIN *et al.*, 2019), é possível observar um grande salto em termos de qualidade das ferramentas e métodos para análise de texto, o que pode ser atribuído à maior capacidade de mapeamento de conceitos em um espaço latente de atributos, que, por sua vez, possibilita um melhor agrupamento de textos em termos de semântica. Em todo caso, o uso destas técnicas requer certa familiaridade com o seu funcionamento e com os textos analisados, de modo que sejam feitos os ajustes necessários para a obtenção de bons resultados.

1.1 Objetivos

Neste trabalho pretendemos, por meio do uso de técnicas computacionais de processamento de texto, estudar os tópicos existentes na obra conhecida como *Etimologias*, de Isidoro de Sevilha (c.560-636). Esta obra, que é considerada como a primeira grande enciclopédia da Idade Média, é formada por 20 livros que tratam das origens das palavras agrupadas em diversos grandes temas. O autor, que viveu em uma época de grandes mudanças culturais, buscava registrar o conhecimento de escritores latinos da Antiguidade Clássica, como Varrão, Catão e Plínio o Velho, entre outros. A obra, que foi escrita em latim medieval, foi copiada exaustivamente ao longo de cerca de 700 anos, sendo utilizada como uma espécie de livro-texto básico nas instituições de ensino da época. Por conveniência, para este trabalho será utilizada uma tradução recente para a língua inglesa (BARNEY *et al.*, 2006).

1.2 Contribuições

A principal contribuição pretendida com este trabalho é conseguir mapear temas transversais discutidos na obra estudada, para além dos temas gerais propostos pelo autor. Adicionalmente, propomos um experimento de análise de similaridade semântica entre a obra estudada e possíveis fontes da antiguidade, para identificar potenciais citações sem atribuição.

1.3 Organização

O texto está organizado da seguinte forma: no Capítulo 2 tratamos do background histórico dos autores estudados, bem como das ferramentas utilizadas; no Capítulo 3 apresentamos os dados analisados, o *pipeline* de processamento criado para os dados, e a configuração dos experimentos realizados; no Capítulo 4 apresentamos e discutimos os resultados obtidos para cada experimento realizado. O Capítulo 5 conclui o texto.

Capítulo 2

Background

Neste trabalho, utilizaremos ferramentas de modelagem de tópicos e grandes modelos de linguagem com o objetivo de demonstrar sua utilidade no estudo de textos históricos. Nas seções a seguir detalhamos as ferramentas utilizadas e as obras a serem estudadas.

2.1 Redes Neurais e Grandes Modelos de Linguagem

Desde o surgimento dos primeiros computadores digitais no período logo após a Segunda Guerra Mundial (1939 - 1945), pesquisadores trabalharam para desenvolver ferramentas que pudessem executar tarefas que até então eram realizadas apenas por humanos. Cientistas como Warren McCulloch (1898 - 1969) e Walter Pitts (1923 - 1969), um médico e um lógico, estão entre os primeiros a discutirem modelos do cérebro humano, e como estes poderiam ser implementados computacionalmente. Em 1943 publicaram um artigo (McCULLOCH e PITTS, 1943) considerado seminal na área, aonde propuseram um modelo matemático do cérebro representado como uma rede de elementos simples interconectados. Um elemento particular recebe sinais dos elementos conectados à sua entrada, e produz uma soma ponderada como sinal de saída, que, dependendo de um limiar pré-determinado, é enviado para elementos conectados à sua saída como um sinal ativo (valor 1) ou inativo (valor 0). Os autores demonstram no artigo que é possível calcular funções lógicas utilizando diferentes configurações de conexões entre elementos. Posteriormente, estes elementos passaram a ser conhecidos como neurônios artificiais, que serviram de base para o desenvolvimento das redes neurais artificiais, anos mais tarde.

Alan Turing (1912 - 1954), matemático e teórico da computação que ajudou a estabelecer as bases da ciência da computação como disciplina, foi um dos primeiros a se dedicar ao que veio a ser chamado posteriormente de inteligência artificial, de forma mais abrangente. Turing apresenta (TURING, 1950) um conjunto de questões e definições que seriam essenciais para o desenvolvimento da área. Em primeiro lugar, ele pergunta se “máquinas podem pensar”, o que traz o primeiro problema importante, o de falta de definições claras sobre o que é “pensar” e, naquele momento, sobre o que seria a “máquina” em questão. Como definições do que é pensar, e de forma mais abrangente, do que é inteligência, representam problemas filosóficos importantes e sem resposta objetiva, Turing propôs

uma abordagem prática para isso: no lugar de “máquinas podem pensar”, ele passa a perguntar se “máquinas podem agir de forma indistinguível de um humano”. Em outras palavras, ele ignora o processo interno pelo qual uma ação é produzida, e se preocupa apenas com o efeito da ação por aqueles que a percebem. Para demonstrar esta ideia, Turing propôs um experimento chamado de Jogo da Imitação. Na versão mais simples do jogo, há três participantes: um humano e um computador que devem responder a perguntas de um juiz humano que só pode fazer perguntas e receber respostas escritas em papel, sem contato direto com os outros dois. Se o juiz, após uma série de perguntas, não for capaz de distinguir o humano da máquina, a máquina venceu o jogo, ou seja, conseguiu se passar por um humano. Apesar de ser uma definição prática, simples, e que serve aos propósitos do autor, esta ideia de que basta ser percebido como humano para ser considerado inteligente ou pensante é problemática para pensadores de outras áreas, que tratam de aspectos metafísicos da mente e do pensamento.

Em 1980, o filósofo John Searle (1932 -) propôs o argumento do Quarto Chinês (SEARLE, 1980). Neste experimento, o autor imagina uma pessoa que não fala chinês isolada em um quarto com um livro contendo instruções sobre como interpretar símbolos chineses. Se alguém colocar um texto em chinês por debaixo da porta, por exemplo, a pessoa poderia seguir as instruções do livro para produzir símbolos que representem uma resposta coerente para falantes de chinês. No entanto, esta pessoa estaria apenas seguindo regras sintáticas sem nenhuma compreensão semântica do texto que está produzindo. O autor busca com isso se contrapor a ideias funcionalistas, como a de Turing, de que a mente é apenas um sistema de processamento de informações. Ou seja, para Searle, não basta ser capaz de se comportar como se entendesse uma conversa para uma entidade ser considerada como pensante ou inteligente.

Apesar das críticas, a busca pela inteligência artificial seguiu com maior ou menor intensidade ao longo da segunda metade do século XX. Por exemplo, partindo da ideia de neurônios artificiais de McCulloch e Pitts em 1943, Frank Rosenblatt (1928 - 1971), psicólogo, desenvolveu o Perceptron (ROSENBLATT, 1958) em 1957, que é considerada a primeira rede neural artificial. A ideia do Perceptron foi desenvolvida ao longo dos anos, passando por diversos altos e baixos, e com a evolução do hardware existente, cada vez mais capaz, culminou na criação do que passou a ser chamado de Deep Learning (KRIZHEVSKY *et al.*, 2012), por volta de 2012 com o trabalho de Alex Krizhevsky (1986 -), que consiste no uso de grandes redes neurais com milhões e até bilhões de elementos aplicado a problemas de processamento de texto e imagem. Mais recentemente, por volta de 2017, a partir do trabalho de Ashish Vaswani (1986 -) e outros foram criadas novas arquiteturas de conexão entre os neurônios artificiais, batizadas de Transformers (VASWANI *et al.*, 2017), com o objetivo de permitir que mais dados pudessem ser armazenados e processados em conjunto, desta forma possibilitando fornecer mais informações de contexto para um problema computacional de processamento de texto ou imagem. A ideia dos transformers é que possibilitou o surgimento de modelos mais recentes como ChatGPT e Gemini, que são exemplares do que é considerado inteligência artificial hoje em dia.

Modelos baseados em transformers possuem grande capacidade de mapeamento semântico de dados e de geração de texto verossímil, que podem ser consideradas como suas características mais salientes. Estas capacidades são obtidas com base na observação e processamento de grande quantidade de dados de exemplo, processo que é chamado

de “treinamento” no jargão da área, realizados com o objetivo de se identificar padrões existentes nos dados para poder determinar, por exemplo, quais palavras que aparecem comumente próximas a outras em certos contextos. Concluída esta etapa de treinamento, o modelo é capaz de produzir texto em resposta a uma questão: primeiro é feito o mapeamento semântico da questão para encontrar termos e sentenças com significado parecido na memória do modelo, e segundo, com base nos termos e sentenças encontrados, o modelo faz a síntese de um texto que tenha verossimilhança. Note o uso do termo verossimilhança no lugar de corretude: como os modelos são criados de acordo com a definição de Turing, basta que forneçam respostas indistinguíveis das de um humano, e cabe a quem coloca a questão avaliar a corretude das respostas fornecidas.

2.2 Evolução das redes neurais em PLN

Hoje, modelos de redes neurais artificiais são prevalentes em diversas áreas de aplicação de machine learning, como classificação e geração de imagens, Processamento de Linguagem Natural (PLN), reconhecimento de fala, etc. A seguir, apresentaremos uma breve (e incompleta) história da evolução desses modelos em algumas sub-áreas de PLN, desde desenvolvimentos anteriores à adoção de redes neurais, até o advento dos Transformers.

HARRIS (1954) e FIRTH (1957) argumentam que o sentido de uma palavra pode ser deduzido, em parte, a partir dos contextos onde ela é comumente utilizada - “conhecerás uma palavra pela companhia que ela mantém”. Essa ideia, que podemos chamar de Hipótese Distribucional, direta ou indiretamente, guiou inúmeros avanços em diferentes campos na área de PLN. DEERWESTER *et al.* (1990), com seu modelo de Análise Semântica Latente (LSA), foi um dos pioneiros em operacionalizar essa ideia, obtendo resultados promissores na área de indexação e recuperação de informação. Neste modelo, aplica-se uma técnica de álgebra linear, denominada Decomposição em Valores Singulares (SVD), a uma matriz de frequência (ou alguma outra métrica, como TF-IDF) termo-documento, para se obter uma representação mais densa e de menor dimensionalidade desta matriz. O resultado é uma representação vetorial que consegue capturar relações semânticas entre termos e documentos, mapeando itens com significados similares a lugares comuns em um espaço latente, configurando uma forma primitiva do conceito que hoje conhecemos como *word embeddings*.

BENGIO *et al.* (2003) propôs o uso de redes neurais artificiais para modelar distribuições de probabilidade conjunta de sequências de palavras em um determinado idioma. Na arquitetura proposta, uma rede neural é treinada para executar duas tarefas simultaneamente: Associar cada palavra de um vocabulário a um vetor de \mathbb{R}^m (*embedding*), e determinar a distribuição de probabilidade condicional de sequências de palavras do vocabulário (isto é, dada uma sequência de palavras, determinar a distribuição de probabilidade da próxima palavra), expressas em termos de suas representações vetoriais. Embora poderosa, essa abordagem ainda tinha uma grande complexidade computacional, inviabilizando o seu uso em grandes conjuntos de dados. Um grande avanço na área, que possibilitou a criação de modelos de embedding com um número bem maior de dimensões, e o uso de volumes muito maiores de dados no treinamento desses modelos, foi a chegada dos modelos *Continuous Skip-gram* e *Continuous Bag-of-Words* (CBOW) (MIKOLOV, CHEN *et al.*, 2013). Estas arquiteturas dispensaram o uso de camadas ocultas, e tinham como objetivo

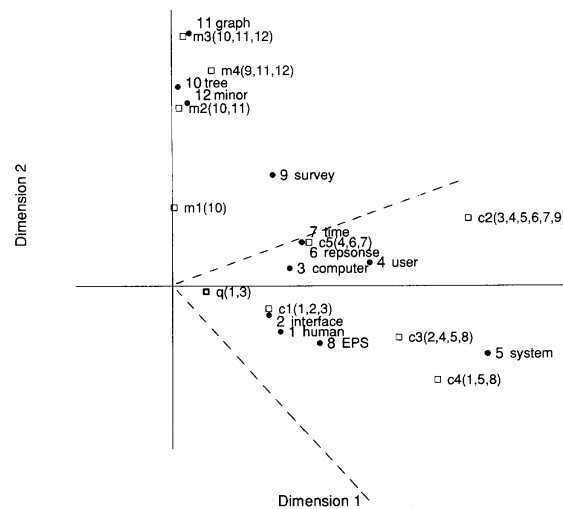


Figura 2.1: Representação do espaço latente gerado pelo LSA.

simplesmente prever uma palavra dadas as palavras ao seu redor (arquitetura CBOW), ou, dada uma palavra, prever as palavras ao seu redor (arquitetura Skip-gram), usando uma janela deslizante de contexto, cujo tamanho é dado como parâmetro da arquitetura. Embora, nessa arquitetura, a ordem das palavras na janela de contexto deixasse de ter importância (no caso do CBOW, usa-se a média dos vetores dessas palavras, por exemplo), os embeddings gerados conseguiam capturar bem melhor o sentido global das palavras, principalmente devido à baixa complexidade computacional do modelo, em relação ao proposto por Bengio, o que possibilitou o uso de um volume muito grande de dados em seu treinamento, e a criação de embeddings de dimensionalidades muito maiores (até 1000 dimensões, comparados com as 50-100 do modelo proposto por Bengio). Porém, vale dizer que este modelo, diferentemente do anterior, não é um modelo de linguagem probabilístico, no sentido de que o seu objetivo não é de modelar distribuições de probabilidades de sequências de palavras, mas de aprender representações semânticas eficientes. Posteriormente, usando técnicas como subamostragem e amostragem negativa, Mikolov mostrou como acelerar o processo de aprendizado e inferência dos modelos, e criar melhores representações (MIKOLOV, SUTSKEVER *et al.*, 2013).

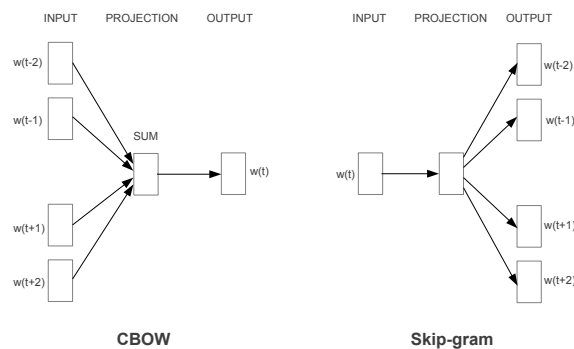


Figura 2.2: Arquiteturas CBOW e Skip-gram.

Na área de tradução automática de texto, um trabalho pioneiro propôs o uso de uma arquitetura do tipo *encoder-decoder* para traduzir sequências de palavras do Inglês para o Francês (SUTSKEVER *et al.*, 2014). O era composto por duas redes do tipo *Long Short-Term Memory*, ou LSTM, que é um tipo de rede neural recorrente (RNN) dotada de mecanismos de “controle de memória”, para lidar com problemas de dependência de longo alcance em sequências. O modelo funcionava em duas etapas: primeiro, o *encoder*, consistindo de uma RNN mapeava sequências de palavras de tamanho variado para um vetor de dimensão fixa, denominado *estado oculto* (também chamado de *vetor de contexto*) iterativamente. Depois, o *decoder* usa o último estado oculto do *encoder* para gerar a sequência de saída. Assim, a predição de cada palavra na saída da rede é influenciada tanto por todas as palavras na entrada, quanto pelas previsões anteriores.

Uma evolução significativa dessa última abordagem foi introduzida no trabalho de BAHDANAU *et al.* (2016), onde foi introduzida uma versão precursora dos mecanismos de atenção modernos. No modelo anterior, sequências de palavras eram mapeadas para um vetor de dimensão fixa, o que fazia com que a performance do modelo deteriorasse rapidamente para sequências longas (CHO *et al.*, 2014). O novo modelo proposto buscava livrar o *encoder* de ter que comprimir todo o contexto em um único vetor de tamanho fixo, gerando uma sequência de estados ocultos que o *decoder* poderia acessar. Além disso, o modelo dispunha de um mecanismo para “procurar”, entre os estados ocultos, informações relevantes para a predição de uma determinada palavra dentro da sequência, e *selecionar* quais palavras deveriam influenciar a palavra atual. O resultado foi uma melhora significativa na transdução de sequências longas.

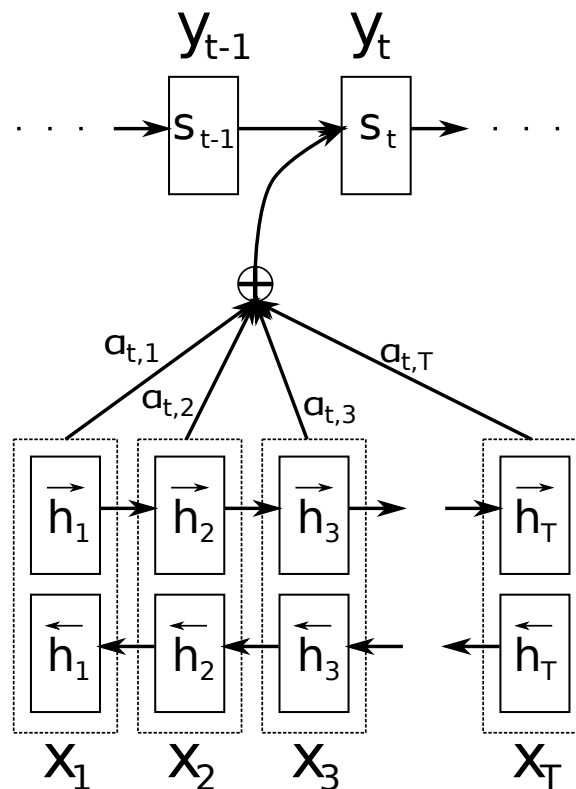


Figura 2.3: RNNSearch.

2.3 Transformers e desenvolvimentos subsequentes

Todos estes desenvolvimentos culminaram na criação da arquitetura batizada de Transformer (VASWANI *et al.*, 2017). Todos os modelos mencionados anteriormente dependem do processamento sequencial de tokens da entrada, devido ao emprego de redes neurais recorrentes. Na arquitetura Transformer, os tokens de entrada são processados em paralelo, o que permitiu uma enorme escalabilidade e aceleração de seus processos de treinamento e inferência. A seguir mostramos uma descrição superficial da sua arquitetura.

2.3.1 Arquitetura do modelo Transformer

Assim como os modelos de transdução já citados (BENGIO *et al.*, 2003; SUTSKEVER *et al.*, 2014; BAHNDANAU *et al.*, 2016), um Transformer conta com um *encoder* e um *decoder*, porém, com a diferença de não haver RNNs envolvidas.

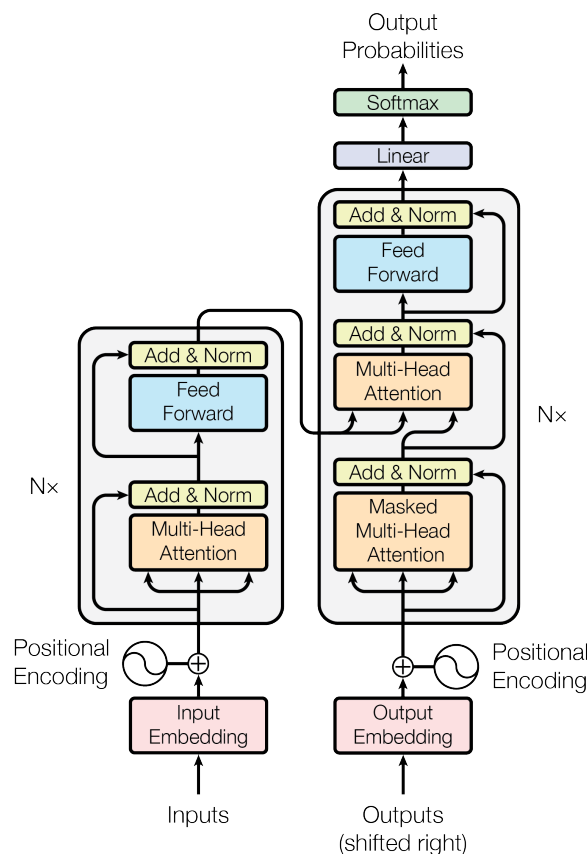


Figura 2.4: Arquitetura de um Transformer.

Encoder

O encoder começa com uma camada de *embedding*, juntamente com uma camada de codificação posicional, que serve para imbuir os embeddings de informação sobre suas posições, já que não estamos mais utilizando RNNs. Em cima disto, múltiplas camadas idênticas são empilhadas. Cada uma destas camadas é composta por duas sub-camadas. A

primeira implementa um mecanismo de atenção paralela (*multi-head attention*), e a segunda é uma simples rede do tipo *multi-layer perceptron* (também chamada de rede *feed-forward*).

Decoder

O decoder também é composto por várias camadas idênticas empilhadas, onde cada camada é composta por três subcamadas, duas iguais às subcamadas do encoder, e uma extra, *masked self-attention*, que consiste em uma camada de atenção modificada para evitar que uma palavra gerada pelo modelo “preste atenção” em uma palavra subsequente (isto é, seja modificada por ela).

2.3.2 Atenção

Podemos entender os mecanismos de atenção presentes na arquitetura Transformer como uma maneira de palavras “informarem” umas às outras quais palavras elas podem influenciar, e de que forma. Mais especificamente, sobre cada embedding (que representa um token, na entrada ou na saída), um bloco de atenção calcula três valores, Q , K e V (*query*, *key* e *value*) usando transformações lineares simples. Para cada embedding, o seu valor de Q é comparado com o valor de K de si mesmo e todos os outros (ou dos embeddings anteriores, no caso de *masked self-attention*), usando uma função de similaridade (no caso, o produto escalar), seguida de uma normalização (*softmax*). O resultado então é usado para calcular a influência dos embeddings uns sobre os outros (efetivamente, soma-se o valor da atenção para um embedding ao valor do embedding, em conexões residuais):

$$\text{Atencao}(Q, K, V) = \text{softmax} \left(\frac{QK^T}{\sqrt{d_k}} \right) V$$

onde d_k é a dimensionalidade de Q e K .

Após uma camada de atenção há uma rede *feed-forward*, com funções de ativação do tipo ReLU (*Rectified Linear Unit*).

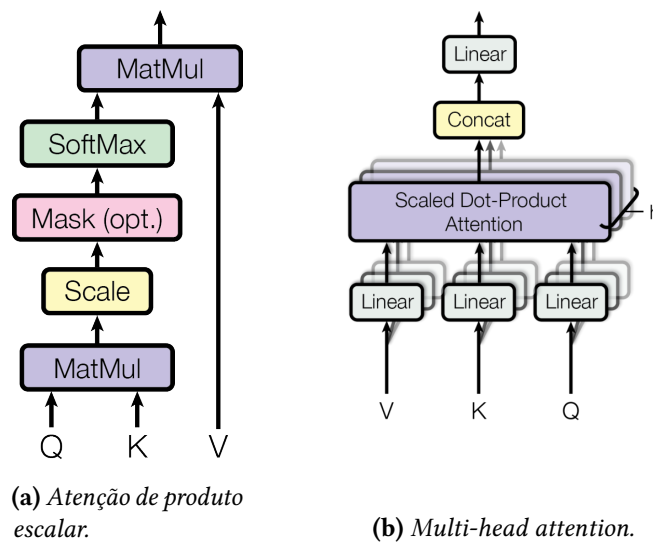


Figura 2.5: Mecanismos de atenção nos Transformers.

Multi-head attention

Ao invés de calcular uma única função de atenção, os valores Q , K e V são projetados linearmente h vezes em espaços de menor dimensionalidade, e vários valores de atenção são calculados em paralelo, e depois esses valores são concatenados e projetados novamente no espaço original, resultando nos valores finais de atenção. Isso permite ao modelo considerar influências entre palavras em diferentes espaços latentes, em diferentes posições.

2.3.3 BERT, SBERT e derivados

A chegada dos Transformers representou um novo paradigma para o uso de redes neurais em PLN, e em outras áreas de aplicação. Pela sua capacidade de modelagem de dependências de longo alcance, e pela sua escalabilidade, modelos baseados em Transformers logo dominaram o cenário de machine learning, impulsionados principalmente por modelos generativos como ChatGPT (RADFORD e NARASIMHAN, 2018), Gemini (TEAM *et al.*, 2025), entre outros. Estes modelos apresentam algumas variações arquitetônicas (como uso apenas de decoders em modelos generativos) mas ainda mantêm a estrutura básica de um Transformer, com seus mecanismos de atenção. Um exemplo de modelo que conta apenas com encoders é o BERT (DEVLIN *et al.*, 2019). Este modelo é pré-treinado em diferentes tarefas, como predição de tokens ocultos, e predição de próxima sentença, e depois passa por um ajuste fino para tarefas específicas, como detecção de paráfrases, inferência textual, classificação de texto, etc. O seu pré-treinamento faz com que o modelo obtenha um “entendimento” básico da linguagem, que melhora consideravelmente o seu desempenho em tarefas específicas posteriores.

Apesar de alcançar bons resultados, o modelo BERT se torna computacionalmente inviável para algumas tarefas, como busca por similaridade semântica em um grande corpus, ou clustering de documentos. Para essas tarefas, é desejável um modelo capaz de gerar embeddings de frases, o que o BERT não foi projetado para fazer. Para mitigar essa limitação surgiu o SBERT (REIMERS e GUREVYCH, 2019), que é uma modificação dos modelos BERT pré-treinados, usando redes neurais siamesas e estruturas de rede do tipo triplet para gerar boas representações semânticas de frases. Essas representações então podem ser comparadas usando distância de cosseno, por exemplo. Os embeddings usados neste trabalho são gerados por descendentes desse modelo.

2.4 Modelagem de Tópicos

TODO: escrever +- 0.5 pagina sobre modelagem de tópicos. É bom ter algumas figuras.

2.4.1 BERTopic

O BERTopic é uma técnica de modelagem de tópicos baseada em modelos de sentence embedding. Mais especificamente, ele usa agrupamentos de representações vetoriais semânticas de trechos de documentos para representar tópicos subjacentes no texto.

Ele é constituído, basicamente, de 4 etapas:

1. Extração de text embeddings

2. Redução de dimensionalidade
3. Clustering
4. Extração de representações dos tópicos

Primeiramente, extrai-se representações vetoriais das sentenças do corpus, utilizando modelos de *sentence embedding*. Esses modelos são capazes de capturar o sentido de frases em linguagem natural, projetando-as em um espaço latente. Essas representações, ou *embeddings*, possuem diversas propriedades interessantes e úteis. Uma delas é que textos com significado similar ficam próximos no espaço latente. Podemos então encontrar agrupamentos de frases similares semanticamente nesse espaço latente, e naturalmente, interpretar esses agrupamentos como assuntos presentes no texto.

A etapa de redução de dimensionalidade é importante para o desempenho da etapa posterior do método, de clusterização, por conta de um fenômeno conhecido como a Maldição da Dimensionalidade: em espaços de alta dimensionalidade, os conceitos de proximidade, distância, ou vizinhos mais próximos perdem sua significância (RADOVANOVIĆ *et al.*, 2010). Este fenômeno então acaba por prejudicar a performance de algoritmos e técnicas que utilizem esses conceitos, como algoritmos de clusterização populares (RADOVANOVIĆ *et al.*, 2010; AGGARWAL *et al.*, 2001).

Além disso, nessa etapa também computamos reduções desses embeddings a 2 dimensões, para serem utilizadas posteriormente em visualizações, também utilizando o UMAP. Fora o número de dimensões do espaço resultante, todos os outros parâmetros permanecem os mesmos.

2.4.2 UMAP

2.4.3 HDBSCAN

2.4.4 Representação dos tópicos

2.5 Obras e autores estudados

Neste trabalho utilizaremos edições em inglês das seguintes obras: as Etimologias (*Etymologiae*), de Isidoro de Sevilha, Sobre Agricultura (*De Agri Cultura*), de Catão, e Sobre Ciência Agrícola (*De Re Rustica*), de Varrão. A escolha destas obras não é por acaso: estudiosos da obra de Isidoro consideram que ele se baseou em diversas obras da antiguidade para escrever as Etimologias, particularmente as obras de Catão e Varrão sobre agricultura. Sendo assim, tentaremos encontrar similaridades de ideias entre os textos a partir do uso de modelos computacionais. A seguir apresentaremos detalhes sobre os autores e seu contexto.

Isidoro de Sevilha (c. 560 - 636) foi um clérigo, teólogo e pensador da alta idade média, que é considerado um dos intelectuais mais importantes do seu tempo, e cuja influência foi sentida por muitos séculos após a sua morte. Nasceu em Cartagena, que à época era parte do Reino Visigótico, estado que ocupou as regiões da Península Ibérica e atual sul da França no período seguinte ao fim do Império Romano do ocidente, e fazia parte de uma família que percentia à elite hispano-romana. Seus três irmãos ocuparam funções

importantes na igreja, com destaque para seu irmão mais velho, Leandro de Sevilha (c. 534 - c. 600) foi Bispo de Sevilha, cargo que Isidoro ocuparia após a morte de Leandro. Todos os quatro irmãos são venerados como santos pela Igreja Católica.

Como Bispo de Sevilha, Isidoro exerceu grande influência no seu tempo, presidindo sínodos e concílios importantes, como os de Sevilha e Toledo, protegendo os monastérios, e ainda se envolveu na conversão dos reis Visigodos do Arianismo, uma doutrina cristã não-trinitária, para o cristianismo Calcedoniano, que veio a se tornar a doutrina dominante na Igreja Católica.

Como intelectual, produziu diversas obras, dentre as quais se destacam as Etimologias, que são um conjunto de livros que formam uma enciclopédia etimológica, que resume e organiza o conhecimento de diversos autores da antiguidade clássica. A obra segue a tradição de enciclopedistas clássicos, como Plínio, o Velho (c. 23 - 79), com o uso de ordenação alfabética de tópicos e de uma abordagem literária para o conhecimento, baseada no pensamento analógico.

As Etimologias tratam de temas diversos como gramática, retórica, matemática, direito, a Igreja, heresias, guerra, agricultura, entre outros. A sua influência foi tão grande nos séculos seguintes que algumas das obras clássicas utilizadas como base deixaram de ser lidas e copiadas e acabaram se perdendo no tempo. Era considerado o texto base para a educação sobre o período clássico durante a maior parte da idade média.

O estudo crítico das Etimologias revela suas possíveis fontes clássicas, que na maioria das vezes não são citadas por Isidoro. De acordo com Stephen A. Barney, tradutor para o inglês (BARNEY *et al.*, 2006) das Etimologias, é possível identificar que o material dos livros I e II, que tratam de gramática, retórica e dialética (as disciplinas do *trivium*) provavelmente foram extraídos dos Institutos, de Cassiodoro (c. 485 - c. 585), o livro III, sobre matemática (contendo as disciplinas do *quadrivium*), provavelmente foi inspirado em Boécio (c. 480 - 524), e o livro XVII, sobre agricultura, deriva de Catão, o Velho (234 - 149 a.C.) e Varrão (116 - 27 a.C.), para citar alguns exemplos.

Marco Pórcio Catão, o Velho (234 - 149 a.C.), foi um soldado, senador e historiador romano. Nascido em uma família de plebeus, que era a classe baixa de cidadãos livres em Roma, descendia de gerações de soldados com reputação de bravura, como seu pai e seu bisavô. Ainda na infância, com a morte de seu pai, passou a cuidar das atividades da fazenda família. Como jovem soldado, especula-se que aos 20 anos tenha participado de campanhas das Guerras Púnicas no papel de tribuno militar, uma patente de oficial do exército romano. Ao retornar do campo de batalha para a sua fazenda, e com o apoio do seu vizinho e amigo Lúcio Valério Flaco, iniciou carreira política como *questor*, cargo que possuía diversas atribuições, entre elas a de cobrança de impostos e de supervisão financeira. Daí se seguiram diversos cargos políticos importantes, como *pretor* (magistrado), *consul* (o cargo mais alto da República Romana) e *censor* (magistrado de nível superior). Em paralelo a suas atividades políticas, escreveu diversas obras, cuja maioria infelizmente foi perdida. Escreveu uma história de Roma em sete livros chamada de Origens, uma obra sobre assuntos militares, e a obra sobre agricultura que utilizaremos neste trabalho, a única preservada na íntegra. Além disso, foi famoso orador e teve cerca de 150 discursos registrados.

Marco Terêncio Varrão (116 - 27 a.C.) foi um intelectual e polímata romano, descrito

por Petrarca como a “terceira grande luz” de Roma, depois de Virgílio e Cícero. Nascido em família pertencente à classe equestre de Roma, abaixo somente da classe senatorial, ocupou cargos políticos ao longo da vida, como *questor*, *pretor* e tribuno do povo. Estudou com o filólogo romano Lúcio Élio Estilo e com o filósofo platonista Antíoco de Ascalão. Foi também um líder militar sem grande prestígio durante a Guerra Civil Cesariana. Foi um escritor prolífico, que produziu cerca de 74 obras sobre temas diversos, entre as quais se destacam a Cronologia Varroniana, que lista as datas de eventos importantes de Roma, e os nove livros das Disciplinas, organizados de acordo com os temas das artes liberais da época, e que serviram de exemplo para enciclopedistas que vieram posteriormente, como Plínio, o Velho, e o próprio Isidoro de Sevilha.

Capítulo 3

Experimentos

Neste capítulo detalhamos os dados utilizados neste trabalho juntamente com o *pipeline* de processamento empregado para a sua coleta e preparação, e a configuração experimental utilizada no estudo.

3.1 Conjunto de dados

Como conjunto de dados a ser explorado neste trabalho, temos as seguintes obras:

- Etimologias (*Etymologiae*) (BARNEY *et al.*, 2006), por Isidoro de Sevilha (c. 560 - 636);
- Sobre Agricultura (*De Agri Cultura*) (CATO e VARRO, 1934), por Marco Pórcio Catão, o Velho (234 - 149 a.C.);
- Sobre Ciência Agrícola (*De Re Rustica*) (CATO e VARRO, 1934), por Marco Terêncio Varrão (116 - 27 a.C.).

No caso das Etimologias, foi utilizada a tradução (BARNEY *et al.*, 2006) do latim para o inglês em formato *pdf*. Para os textos de Catão e Varrão foram utilizadas edições da *Loeb Classical Library* (CATO e VARRO, 1934), atualmente em domínio público e disponíveis online (THAYER, 2025) em formato *html*.

3.2 Preparação dos dados

(TODO: Adicionar figura do esquema geral do processo)

3.2.1 Limpeza

No caso das Etimologias, o texto foi extraído do arquivo *pdf* do livro e convertido para formato de texto puro, removendo-se cabeçalhos, rodapés, números de página, e outros artefatos dessa conversão que não são necessários para a análise. Devido a peculiaridades do formato *pdf* e dificuldades encontradas nas bibliotecas utilizadas para a extração do texto, parte do processo de limpeza teve de ser realizado manualmente. No caso das outras

obras, disponíveis em formato *html*, o processamento foi realizado de forma totalmente automática.

3.2.2 Pré-processamento

Primeiramente, pelo modo como cada livro está formatado, fizemos um processamento inicial para separar cada livro em capítulos, e juntar palavras com hífen. Em seguida, cada capítulo foi subdividido em sentenças usando o módulo *SentenceRecognizer* da biblioteca Spacy (HONNIBAL *et al.*, 2020), que possui diversos recursos para processamento de linguagem natural (NLP).

Uma particularidade da edição escolhida das Etimologias é que cada capítulo é subdividido em seções relativamente curtas, a maioria tendo somente uma oração. Com isso, podemos fazer a modelagem dos tópicos usando dois níveis diferentes de granularidade, a depender de como definimos um “documento”, unidade básica de nossa análise:

- **Seções**, já presentes no texto;
- **Sentenças**, delimitadas automaticamente;

Em nossa análise, utilizamos as duas abordagens e comparamos os resultados.

TODO: Exemplos das seções.

3.2.3 Stop words, stemming e lematização

Tarefas rotineiras de pré-processamento de linguagem natural incluem: remoção de *stop words*, que são palavras muito comuns e irrelevantes, *stemming*, e lematização, que são formas de padronizar palavras, reduzindo-as a sua forma mais básica. Apesar de outras técnicas de modelagem de tópicos adotarem essas tarefas como parte do seu processo de preparação dos dados, não utilizamos essas técnicas para este trabalho, pois os modelos de *embedding* que usamos utilizam informações contextuais de cada palavra em uma sentença, e remover palavras ou modificá-las poderia prejudicar a performance de tais modelos (CHAERUL HAVIANA *et al.*, 2023).

3.2.4 Geração de *embeddings*

Depois da etapa de pré-processamento dos textos, geramos *embeddings* para cada documento, utilizando modelos do tipo *sentence embedders* pré-treinados. Os modelos usados para a geração desses *embeddings* foram: *sentence-transformers/LaBSE*, *jinaai/jina-embeddings-v3*, *intfloat/multilingual-e5-large-instruct*, *nomic-ai/nomic-embed-text-v2-moe*. Estes modelos foram desenhados e treinados para que sentenças semanticamente similares em línguas diferentes, ou traduções, estejam próximas umas das outras em um espaço latente. Podemos usar estes modelos para comparar textos em diferentes idiomas e analisar suas conexões.

3.2.5 Conjuntos de dados

Após essas etapas iniciais, definimos nossos conjuntos de dados. Cada conjunto é definido por um modelo de embedding e um nível de granularidade, já explicado. TODO: exemplos, visualizações dos espaços de embeddings.

3.3 Persistência

Os embeddings gerados, suas reduções, e diversos metadados foram armazenados em um banco de dados do *ChromaDB*, juntamente com cada documento.

TODO: versão final do modelo de dados no ChromaDB.

Dessa forma, podemos fazer diversos tipos de busca, como busca por similaridade, busca textual, filtrar resultados baseados em metadados de cada documento, como autor, idioma, etc.

3.4 Experimentos

Como já foi explicado, a técnica de modelagem de tópicos usada (BERTopic) é composta de diferentes etapas, com várias possibilidades de escolha de algoritmos para cada etapa, e hiperparâmetros para estes algoritmos. Neste trabalho, foram usados o UMAP para redução de dimensionalidade, e o HDBSCAN para o clustering dos documentos e detecção dos tópicos.

3.4.1 Redução de dimensionalidade

A seguir, fizemos uma redução de dimensionalidade de cada um dos conjuntos de embeddings gerados. Para isso, utilizamos o UMAP, com os seguintes parâmetros:

Tabela 3.1: *Parâmetros usados para o UMAP*

Nome do parâmetro	Descrição	Valor usado	Nota
n_neighbors	Controla como o UMAP equilibra estrutura local versus global nos dados. Valores maiores levam a uma visão mais global da estrutura dos dados, e valores menores a uma visão mais local.	10	O valor padrão para o BERTopic é 15.
n_components	Dimensionalidade do <i>dataset</i> resultante.	5	Valor padrão para o BERTopic.

Continua na próxima página

Tabela 3.1 – continuação

Nome do parâmetro	Descrição	Valor usado	Nota
min_dist	Controla o quão dispersos os pontos estarão na projeção de baixa dimensionalidade. Valores baixos podem ser interessantes para tarefas de clustering.	0	Valor padrão para o BERTopic.
low_memory	Restringe o uso de memória em detrimento da velocidade da computação. Útil quando há pouca memória disponível.	false	
metric	Métrica usada para calcular distâncias entre pontos.	"cosine"	Valor padrão para o BERTopic.
random_state	Usado para garantir determinismo.	42	

TODO: normalmente o BERTopic usa `n_neighbors = 15`, então por que eu mudei esse valor? A implementação usada foi a da biblioteca `umap-learn`.¹

3.4.2 Clustering

Com os embeddings gerados e suas reduções já computadas, inicializamos o BERTopic de modo a pular essas etapas iniciais e começar pela etapa de clustering. O algoritmo utilizado foi o HDBSCAN, disponível na biblioteca Scikit-learn,² e utilizamos os seguintes parâmetros:

Tabela 3.2: Parâmetros usados para o HDBSCAN

Nome do parâmetro	Descrição	Valor usado	Nota
min_cluster_size	Controla o tamanho mínimo e a quantidade de <i>clusters</i> gerados.	10	Valor padrão para o BERTopic.
min_samples	Controla a quantidade de pontos classificados como ruído.	1	Escolhemos um valor baixo para tentar diminuir a quantidade pontos classificados como outliers.

Para melhorar a representação dos tópicos gerados, fizemos alguns ajustes nas etapas finais do método:

¹ <https://umap-learn.readthedocs.io/en/latest/>

² <https://scikit-learn.org/stable/modules/generated/sklearn.cluster.HDBSCAN.html>

- Remoção de stop-words, no CountVectorizer
- Remoção de palavras comuns e uso de BM-25 weighting measure.

Capítulo 4

Resultados

Capítulo 5

Conclusão

Referências

- [AGGARWAL *et al.* 2001] Charu C. AGGARWAL, Alexander HINNEBURG e Daniel A. KEIM. “On the surprising behavior of distance metrics in high dimensional spaces”. In: *Proceedings of the 8th International Conference on Database Theory. ICDT '01*. Berlin, Heidelberg: Springer-Verlag, 2001, pp. 420–434. ISBN: 3540414568 (citado na pg. 11).
- [BAHDANAU *et al.* 2016] Dzmitry BAHKANAU, Kyunghyun CHO e Yoshua BENGIO. *Neural Machine Translation by Jointly Learning to Align and Translate*. 2016. arXiv: 1409.0473 [cs.CL]. URL: <https://arxiv.org/abs/1409.0473> (citado nas pgs. 7, 8).
- [BARNEY *et al.* 2006] Stephen A BARNEY, Wendy J LEWIS, Jennifer A BEACH e Oliver BERGHOF. *The etymologies of Isidore of Seville*. Cambridge University Press, 2006 (citado nas pgs. 1, 12, 15).
- [BENGIO *et al.* 2003] Yoshua BENGIO, Réjean DUCHARME, Pascal VINCENT e Christian JANVIN. “A neural probabilistic language model”. In: *Journal of machine learning research*. 2003. URL: <https://api.semanticscholar.org/CorpusID:221275765> (citado nas pgs. 5, 8).
- [BLEI *et al.* 2003] David M BLEI, Andrew Y NG e Michael I JORDAN. “Latent dirichlet allocation”. *Journal of machine Learning research* 3.Jan (2003), pp. 993–1022 (citado na pg. 1).
- [CATO e VARRO 1934] Marcus Porcius CATO e Marcus Terentius VARRO. *On Agriculture; On Agriculture*. Ed. por William Davis HOOPER e Harrison Boyd ASH. Vol. 283. Loeb Classical Library. Latin text with facing English translation; revised edition. Cambridge, MA: Harvard University Press, 1934. ISBN: 978-0674993136. DOI: 10.4159/DLCL.cato-agriculture.1934 (citado na pg. 15).
- [CHAERUL HAVIANA *et al.* 2023] Sam Farisa CHAERUL HAVIANA, Sri MULYONO e BADI'AH. “The effects of stopwords, stemming, and lemmatization on pre-trained language models for text classification: a technical study”. In: *2023 10th International Conference on Electrical Engineering, Computer Science and Informatics (EECSI)*. 2023, pp. 521–527. DOI: 10.1109/EECSI59885.2023.10295797 (citado na pg. 16).

- [CHO *et al.* 2014] Kyunghyun CHO, Bart van MERRIENBOER, Dzmitry BAHDAU e Yoshua BENGIO. *On the Properties of Neural Machine Translation: Encoder-Decoder Approaches*. 2014. arXiv: 1409.1259 [cs.CL]. URL: <https://arxiv.org/abs/1409.1259> (citado na pg. 7).
- [DEERWESTER *et al.* 1990] Scott DEERWESTER, Susan T DUMAIS, George W FURNAS, Thomas K LANDAUER e Richard HARSHMAN. “Indexing by latent semantic analysis”. *Journal of the American society for information science* 41.6 (1990), pp. 391–407 (citado nas pgs. 1, 5).
- [DEVLIN *et al.* 2019] Jacob DEVLIN, Ming-Wei CHANG, Kenton LEE e Kristina TOUTANOVA. “Bert: pre-training of deep bidirectional transformers for language understanding”. In: *Proceedings of the 2019 conference of the North American chapter of the association for computational linguistics*. 2019, pp. 4171–4186 (citado nas pgs. 1, 10).
- [FIRTH 1957] J. R. FIRTH. “A synopsis of linguistic theory 1930-55.” 1952-59 (1957), pp. 1–32 (citado na pg. 5).
- [HARRIS 1954] Zellig S. HARRIS. “Distributional structure”. *WORD* 10.2-3 (1954), pp. 146–162. DOI: 10.1080/00437956.1954.11659520. eprint: <https://doi.org/10.1080/00437956.1954.11659520>. URL: <https://doi.org/10.1080/00437956.1954.11659520> (citado na pg. 5).
- [HONNIBAL *et al.* 2020] Matthew HONNIBAL, Ines MONTANI, Sofie VAN LANDEGHEM e Adam BOYD. “spaCy: industrial-strength natural language processing in python”. *arXiv preprint arXiv:2002.06201* (2020). URL: <https://arxiv.org/abs/2002.06201> (citado na pg. 16).
- [JELODAR *et al.* 2019] Hamed JELODAR *et al.* “Latent dirichlet allocation (lda) and topic modeling: models, applications, a survey”. *Multimedia tools and applications* 78 (2019), pp. 15169–15211 (citado na pg. 1).
- [KRIZHEVSKY *et al.* 2012] Alex KRIZHEVSKY, Ilya SUTSKEVER e Geoffrey E. HINTON. “Imagenet classification with deep convolutional neural networks”. In: *Proceedings of the 26th International Conference on Neural Information Processing Systems - Volume 1*. NIPS’12. Lake Tahoe, Nevada: Curran Associates Inc., 2012, pp. 1097–1105 (citado na pg. 4).
- [McCULLOCH e PITTS 1943] Warren S. McCULLOCH e Walter PITTS. “A logical calculus of the ideas immanent in nervous activity”. *Bulletin of Mathematical Biophysics* 5 (1943), pp. 115–133 (citado na pg. 3).
- [MIKOLOV, CHEN *et al.* 2013] Tomas MIKOLOV, Kai CHEN, Greg CORRADO e Jeffrey DEAN. *Efficient Estimation of Word Representations in Vector Space*. 2013. arXiv: 1301.3781 [cs.CL]. URL: <https://arxiv.org/abs/1301.3781> (citado na pg. 5).

- [MIKOLOV, SUTSKEVER *et al.* 2013] Tomas MIKOLOV, Ilya SUTSKEVER, Kai CHEN, Greg CORRADO e Jeffrey DEAN. *Distributed Representations of Words and Phrases and their Compositionality*. 2013. arXiv: 1310.4546 [cs.CL]. URL: <https://arxiv.org/abs/1310.4546> (citado na pg. 6).
- [RADFORD e NARASIMHAN 2018] Alec RADFORD e Karthik NARASIMHAN. “Improving language understanding by generative pre-training”. In: 2018. URL: <https://api.semanticscholar.org/CorpusID:49313245> (citado na pg. 10).
- [RADOVANOVIĆ *et al.* 2010] Miloš RADOVANOVIĆ, Alexandros NANOPOULOS e Mirjana IVANOVIĆ. “Hubs in space: popular nearest neighbors in high-dimensional data”. *J. Mach. Learn. Res.* 11 (dez. de 2010), pp. 2487–2531. ISSN: 1532-4435 (citado na pg. 11).
- [REIMERS e GUREVYCH 2019] Nils REIMERS e Iryna GUREVYCH. *Sentence-BERT: Sentence Embeddings using Siamese BERT-Networks*. 2019. arXiv: 1908.10084 [cs.CL]. URL: <https://arxiv.org/abs/1908.10084> (citado na pg. 10).
- [ROSENBLATT 1958] F. ROSENBLATT. *The Perceptron: A Theory of Statistical Separability in Cognitive Systems (Project Para)*. Cornell Aeronautical Laboratory, Inc. Cornell Aeronautical Laboratory, 1958. URL: <https://books.google.com.br/books?id=7Q4-AQAIAAJ> (citado na pg. 4).
- [SEARLE 1980] John R. SEARLE. “Minds, brains, and programs”. *Behavioral and Brain Sciences* 3.3 (1980), pp. 417–457 (citado na pg. 4).
- [SUTSKEVER *et al.* 2014] Ilya SUTSKEVER, Oriol VINYALS e Quoc V. LE. *Sequence to Sequence Learning with Neural Networks*. 2014. arXiv: 1409.3215 [cs.CL]. URL: <https://arxiv.org/abs/1409.3215> (citado nas pgs. 7, 8).
- [TEAM *et al.* 2025] Gemini TEAM *et al.* *Gemini: A Family of Highly Capable Multimodal Models*. 2025. arXiv: 2312.11805 [cs.CL]. URL: <https://arxiv.org/abs/2312.11805> (citado na pg. 10).
- [THAYER 2025] William P. THAYER. *LacusCurtius: A Gateway to Ancient Rome*. penelope.uchicago.edu/Thayer/E/Roman/home.html. Bill Thayer’s Web Site, hosted by the University of Chicago. Last updated: 20 October 2025 (according to the main site). 2025. (Acesso em 08/11/2025) (citado na pg. 15).
- [TURING 1950] Alan Mathison TURING. “Computing machinery and intelligence”. *Mind* 49 (1950), pp. 433–460 (citado na pg. 3).
- [VASWANI *et al.* 2017] Ashish VASWANI *et al.* “Attention is all you need”. *Advances in neural information processing systems* 30 (2017) (citado nas pgs. 1, 4, 8).