

Aula 2 – O Pipeline Clássico de PLN

Bem-vindos à segunda etapa da nossa jornada pelo fascinante mundo do Processamento de Linguagem Natural (PLN)! Se você já se perguntou como os computadores conseguem "entender" e interagir com a linguagem humana, esta aula é o seu ponto de partida para desvendar os segredos por trás dessa capacidade. Imagine poder transformar textos complexos em dados estruturados que máquinas podem processar, ou até mesmo preparar-se para os desafios de um concurso público que exige conhecimento em inteligência artificial.

Nesta aula, vamos mergulhar nos fundamentos que pavimentaram o caminho para as inovações que vemos hoje. Compreender o "pipeline clássico" de PLN não é apenas revisitar o passado; é construir uma base sólida para entender como as abordagens mais recentes, como os Modelos de Linguagem de Grande Escala (LLMs), funcionam e quais problemas eles vieram resolver. É como aprender a mecânica básica de um carro antes de pilotar um veículo autônomo.

Ao final desta aula, você será capaz de identificar e explicar as principais etapas do pipeline clássico de PLN, como tokenização, stemming e lematização. Além disso, compreenderá os modelos de representação de palavras como Bag-of-Words e TF-IDF, e a importância da análise sintática. Prepare-se para conectar esses conceitos com as limitações que levaram ao surgimento das abordagens neurais, preparando o terreno para as próximas aulas.

A Jornada do Texto: Desvendando o Pipeline Clássico de PLN

Imagine que você tem um texto, talvez um artigo de jornal ou um e-mail, e precisa que um computador o "leia" e extraia informações importantes. Para nós, humanos, isso é relativamente fácil, mas para uma máquina, um texto é apenas uma sequência de caracteres sem sentido inerente. É aqui que entra o pipeline clássico de Processamento de Linguagem Natural: uma série de etapas sequenciais que transformam o texto bruto em um formato que os algoritmos podem entender e processar.



Analogia: Pense nesse pipeline como uma linha de montagem em uma fábrica, onde cada estação tem uma função específica para transformar a matéria-prima (o texto) em um produto final utilizável (dados processados).

Cada etapa constrói sobre a anterior, adicionando mais estrutura e significado ao texto. Sem essa organização, seria impossível para os computadores realizarem tarefas como tradução automática, análise de sentimentos ou até mesmo responder a perguntas.

Nosso objetivo é entender cada uma dessas "estações" e como elas contribuem para a compreensão da linguagem. Desde a quebra do texto em suas menores unidades até a identificação de suas funções gramaticais, cada passo é crucial. Vamos começar pela primeira e fundamental etapa: a tokenização.

Tokenização: A Primeira Quebra do Texto

Antes que um computador possa fazer qualquer coisa com um texto, ele precisa saber onde uma "palavra" termina e outra começa. Para nós, isso é intuitivo: espaços e pontuações geralmente delimitam as palavras. No entanto, para uma máquina, "Olá, mundo!" é apenas uma sequência de 12 caracteres. A **tokenização** é o processo de dividir uma sequência de texto em unidades menores, chamadas **tokens**.

O que são tokens?

Tokens podem ser palavras, números, símbolos de pontuação ou até mesmo subpalavras, dependendo da complexidade do tokenizador.

Exemplo prático

Na frase "Eu gosto de PLN.", um tokenizador simples pode gerar os tokens ["Eu", "gosto", "de", "PLN", "."].

Parece simples, mas a tokenização é mais complexa do que parece, especialmente em idiomas com escrita contínua (como o chinês) ou com contrações e hifenizações.

A escolha do tokenizador pode impactar significativamente as etapas subsequentes do pipeline. Um tokenizador que separa "não-verbal" em ["não", "-", "verbal"] terá um resultado diferente de um que o mantém como um único token.

Essa primeira decisão é crucial, pois define as "peças" com as quais todo o processamento futuro será construído, influenciando diretamente a qualidade da análise.

Stemming e Lematização: Reduzindo a Complexidade Morfológica

A linguagem humana é rica em variações. Uma mesma palavra pode aparecer em diferentes formas gramaticais: "correr", "correndo", "correu", "corremos". Para um computador, essas são quatro palavras distintas, mas para nós, todas se referem à mesma ação fundamental. Lidar com essa variabilidade é um desafio para o PLN, pois aumenta a complexidade e a esparsidade dos dados.

Para resolver isso, utilizamos técnicas de normalização como **stemming** e **lematização**. Ambas têm o objetivo de reduzir as palavras às suas formas base, mas o fazem de maneiras diferentes. Pense nisso como tentar encontrar a "essência" de uma palavra, ignorando suas flexões gramaticais. Isso é crucial para que o sistema reconheça que "gatos" e "gato" são a mesma entidade.

Stemming

O stemming é um processo mais rudimentar, que "corta" os sufixos das palavras para chegar a uma raiz comum, que nem sempre é uma palavra válida.

Exemplo: "correndo" e "correu" podem ser reduzidos a "corr".

Lematização

Já a lematização é mais sofisticada, utilizando um vocabulário e análise morfológica para retornar a forma base ou "lema" da palavra, que é sempre uma palavra válida.


Exemplo: "correndo" e "correu" seriam lematizados para "correr".

| Conceito | Âmbito/Aplicação | Base/Origem | Exemplo (Português) |
|--------------------|------------------------------------------------|----------------------------------------------|-----------------------|
| Stemming | Redução rápida de palavras para uma raiz comum | Algoritmos heurísticos de remoção de sufixos | "Correndo" → "corr" |
| Lematização | Redução para a forma base (lema) da palavra | Dicionários e análise morfológica | "Correndo" → "correr" |

A lematização é mais precisa, mas também mais computacionalmente intensiva.

Modelos de Representação de Palavras: Bag-of-Words (BoW)

Uma vez que temos nossos tokens normalizados, o próximo desafio é como transformar essas palavras em algo que um algoritmo de aprendizado de máquina possa entender: números. Os computadores não processam texto diretamente; eles precisam de representações numéricas. Um dos modelos mais simples e fundamentais para isso é o **Bag-of-Words (BoW)**, ou "Saco de Palavras".

 **Analogia do Supermercado:** Imagine que você está em um supermercado e quer saber o que há em sua cesta de compras. Você não se importa com a ordem em que os itens foram colocados, apenas com quais itens estão lá e quantos de cada um. O modelo BoW funciona de forma semelhante!

O modelo BoW representa um documento como uma coleção (um "saco") de suas palavras, ignorando completamente a ordem em que elas aparecem, mas registrando a frequência de cada palavra.

01

Criar vocabulário

Primeiro criamos um vocabulário de todas as palavras únicas em nosso conjunto de documentos.

02

Contar ocorrências

Para cada documento, contamos a ocorrência de cada palavra desse vocabulário.

03

Gerar vetor

O resultado é um vetor numérico onde cada posição corresponde a uma palavra do vocabulário, e o valor é a contagem dessa palavra no documento.


Exemplo: Se o vocabulário for ["gato", "cachorro", "brinca", "com"], e o documento for "O gato brinca com o cachorro", o vetor BoW poderia ser [1, 1, 1, 1].

Embora simples, o BoW é a base para muitas análises textuais, mas sua simplicidade também é sua maior limitação, pois perde informações valiosas sobre a estrutura e o contexto da frase.

TF-IDF: Pesando a Importância das Palavras

O modelo Bag-of-Words nos dá uma contagem de palavras, mas nem todas as palavras são igualmente informativas. Palavras muito comuns, como "o", "a", "de", "e", aparecem em quase todos os documentos e, embora sejam frequentes, não carregam muito significado sobre o *tópico* do documento. Por outro lado, palavras raras e específicas podem ser cruciais para identificar o assunto principal. Como podemos dar mais peso às palavras que realmente importam?

É aí que entra o **TF-IDF**, uma sigla para *Term Frequency-Inverse Document Frequency* (Frequência do Termo-Frequência Inversa do Documento). Essa técnica é uma forma inteligente de atribuir um peso numérico a cada palavra em um documento, refletindo sua importância não apenas dentro daquele documento, mas também em relação a todo o conjunto de documentos (o corpus).

 **Analogia do Detetive:** Pense em um detetive que busca pistas: ele não se importa com objetos comuns que estão em todo lugar, mas sim com aqueles itens raros que podem apontar para algo específico.

Term Frequency (TF)

Quão frequentemente uma palavra aparece em um documento específico. Quanto maior o TF, mais relevante a palavra parece ser para aquele documento.

Inverse Document Frequency (IDF)

Quão rara uma palavra é em todo o corpus. Palavras que aparecem em poucos documentos têm um IDF alto, indicando que são mais discriminativas.

Ao multiplicar TF por IDF, o TF-IDF penaliza palavras muito comuns e valoriza palavras que são frequentes em um documento, mas raras no corpus geral. Isso nos dá uma representação numérica mais sofisticada e útil para tarefas como busca de informações, sumarização e classificação de textos.

Análise de Sintaxe: Part-of-Speech (POS) Tagging

Até agora, tratamos as palavras como unidades isoladas ou como parte de uma "saco". No entanto, a linguagem humana é muito mais do que apenas uma coleção de palavras; ela tem estrutura. Para um computador começar a entender essa estrutura, ele precisa saber a função gramatical de cada palavra. É aqui que entra o **Part-of-Speech (POS) Tagging**, ou "Marcação de Classes Gramaticais".

O POS tagging é o processo de atribuir uma categoria gramatical (como substantivo, verbo, adjetivo, advérbio, preposição, etc.) a cada palavra em uma frase. Pense nisso como rotular cada palavra com sua "função" na oração.



Exemplo de POS Tagging

"O gato preto dorme rapidamente"


- "O" → Artigo
- "gato" → Substantivo
- "preto" → Adjetivo
- "dorme" → Verbo
- "rapidamente" → Advérbio

Essa etapa é fundamental porque a função gramatical de uma palavra pode mudar seu significado ou a forma como ela se relaciona com outras palavras. Por exemplo, "banco" pode ser um substantivo (instituição financeira ou assento) ou um verbo (de "bancar"). O POS tagging ajuda a desambiguar essas situações e é um pré-requisito para análises sintáticas mais complexas, como o parsing, que veremos a seguir.

É uma ponte essencial entre a simples identificação de palavras e a compreensão da estrutura da frase.

Análise de Sintaxe: Parsing (Análise Sintática)

Com o POS tagging, já sabemos a função de cada palavra. Mas uma frase é mais do que uma sequência de palavras com rótulos; ela tem uma estrutura hierárquica que define o relacionamento entre elas. A **análise sintática**, ou **parsing**, é o processo de analisar a estrutura gramatical de uma frase, identificando as relações entre as palavras e agrupando-as em constituintes sintáticos (como sintagmas nominais, sintagmas verbais) ou relações de dependência.

 **Analogia do Quebra-Cabeça:** Imagine que você está montando um quebra-cabeça complexo, onde cada peça é uma palavra e as conexões entre elas formam a imagem completa da frase.

O parsing constrói uma representação da estrutura da frase, geralmente na forma de uma árvore sintática. Existem dois tipos principais:

1. Parsing de Constituintes

Identifica grupos de palavras que formam unidades gramaticais (sintagmas).

Exemplo: "[O gato] [dorme [rapidamente]]"

2. Parsing de Dependência



Identifica as relações de dependência entre as palavras, mostrando qual palavra "modifica" ou "depende" de qual.

Exemplo: "dorme" é o núcleo, "gato" é o sujeito de "dorme", "rapidamente" é um advérbio que modifica "dorme".

Essa etapa é vital para tarefas que exigem uma compreensão profunda da estrutura da frase, como tradução automática (onde a estrutura de uma frase em um idioma precisa ser mapeada para outro), extração de informações (identificar quem fez o quê a quem) e até mesmo para sistemas de perguntas e respostas. O parsing nos permite ir além do significado individual das palavras e começar a entender o significado da frase como um todo.

Limitações dos Modelos Clássicos: Onde Eles Falham?

Os modelos clássicos de PLN que exploramos – tokenização, stemming, lematização, BoW, TF-IDF, POS tagging e parsing – foram e ainda são ferramentas valiosas. Eles nos permitiram dar os primeiros passos na compreensão computacional da linguagem. No entanto, a linguagem humana é incrivelmente complexa e cheia de nuances que esses modelos, por sua natureza, não conseguem capturar plenamente.

  **Analogia da Foto:** Pense em uma foto em preto e branco. Ela nos dá uma ideia da cena, mas falta a riqueza de detalhes e a profundidade que as cores proporcionam.

Da mesma forma, os modelos clássicos têm limitações significativas:

✗ Ambiguidade Semântica

Palavras como "banco" (instituição financeira vs. assento) ou "manga" (fruta vs. parte da roupa) são um desafio. Os modelos clássicos têm dificuldade em desambiguar o significado com base no contexto.

✗ Contexto e Sinônimos

BoW e TF-IDF ignoram a ordem das palavras, perdendo informações contextuais cruciais. Além disso, eles tratam sinônimos como "carro" e "automóvel" como palavras completamente diferentes, mesmo que tenham significados semelhantes.

✗ Ironia e Sarcasmo

A linguagem humana é cheia de figuras de linguagem que dependem de um entendimento contextual e cultural profundo, algo que os modelos baseados em contagem e regras não conseguem processar.

✗ Relações Complexas

Embora o parsing ajude, ele ainda pode ter dificuldades com frases muito longas ou estruturas gramaticais incomuns, e não captura o significado subjacente ou as intenções do autor.

Essas limitações mostram que, para realmente "entender" a linguagem como os humanos, precisaríamos de abordagens que pudessem capturar o significado contextual e as relações semânticas de forma mais rica.

A Necessidade de Abordagens Neurais: Um Salto Adiante

As limitações dos modelos clássicos nos levaram a uma encruzilhada: como podemos fazer com que os computadores não apenas processem palavras, mas também compreendam seu significado e contexto de forma mais humana? A resposta veio com o advento das **abordagens neurais**, inspiradas no funcionamento do cérebro humano e impulsionadas pelo avanço do aprendizado de máquina e da capacidade computacional.

Imagine que, em vez de apenas ver a foto em preto e branco, agora podemos ver uma imagem colorida, em 3D, e até mesmo interativa.

As redes neurais, especialmente as redes neurais recorrentes (RNNs) e, mais tarde, os Transformers, revolucionaram o PLN ao permitir que os modelos aprendessem representações densas e contextuais das palavras, conhecidas como **word embeddings**.



Word Embeddings

Esses embeddings são vetores numéricos que capturam o significado semântico e as relações entre as palavras. Palavras com significados semelhantes, como "rei" e "rainha", ou "homem" e "mulher", teriam vetores próximos no espaço multidimensional.



Compreensão Contextual

Isso permite que os modelos entendam que "carro" e "automóvel" são sinônimos, ou que "rei" se relaciona com "rainha" da mesma forma que "homem" se relaciona com "mulher".



Nova Era do PLN

Essa capacidade de capturar nuances e contexto abriu as portas para uma nova era no PLN, superando muitos dos gargalos dos modelos clássicos e preparando o terreno para os poderosos Modelos de Linguagem de Grande Escala (LLMs) que conhecemos hoje.

O Impacto dos LLMs: Uma Nova Era para o PLN

Avançando rapidamente para os dias atuais, somos testemunhas de uma revolução no PLN impulsionada pelos **Modelos de Linguagem de Grande Escala (LLMs)**. Nomes como GPT (Generative Pre-trained Transformer), Llama e Claude tornaram-se sinônimos de inteligência artificial, capazes de gerar texto coerente, traduzir idiomas, responder a perguntas complexas e até mesmo escrever código. Mas o que os torna tão poderosos e como eles se conectam com o que vimos até agora?

Os LLMs representam o ápice das abordagens neurais no PLN. Eles são treinados em quantidades massivas de dados textuais da internet, aprendendo padrões linguísticos, fatos e até mesmo raciocínio. Sua capacidade de gerar texto que é indistinguível do humano, ou de realizar tarefas de compreensão de linguagem com uma precisão sem precedentes, transformou a forma como interagimos com a tecnologia e abriu novas fronteiras para a pesquisa e aplicação de IA.



Treinamento Massivo

Treinados em quantidades massivas de dados textuais da internet



Aprendizado de Padrões

Aprendem padrões linguísticos, fatos e até mesmo raciocínio




Capacidades Avançadas

Geram texto coerente, traduzem, respondem perguntas e escrevem código

A arquitetura subjacente a muitos desses modelos é o **Transformer**, que introduziu um mecanismo de atenção que permite aos modelos ponderar a importância de diferentes partes da entrada ao processar uma sequência. Essa inovação foi um divisor de águas, permitindo que os LLMs capturassem dependências de longo alcance no texto de forma muito mais eficiente do que as arquiteturas anteriores, como as Redes Neurais Recorrentes (RNNs). O impacto é vasto, desde assistentes virtuais mais inteligentes até ferramentas de criação de conteúdo e análise de dados em escala.

Arquitetura Transformer: O Coração dos LLMs

Para entender a verdadeira magia por trás dos LLMs, precisamos olhar para a [Arquitetura Transformer](#). Publicada em 2017 por pesquisadores do Google em um artigo seminal "Attention Is All You Need", o Transformer revolucionou o PLN ao introduzir o mecanismo de [atenção \(self-attention\)](#). Antes dele, as redes neurais recorrentes (RNNs) e suas variantes (LSTMs, GRUs) eram o estado da arte para sequências, mas tinham limitações significativas, como a dificuldade de processar informações em paralelo e de capturar dependências de longo alcance de forma eficiente.

 **Analogia da Atenção Humana:** Quando lemos uma frase complexa, nosso cérebro não processa cada palavra isoladamente. Em vez disso, ele foca em certas palavras ou partes da frase que são mais relevantes para entender o significado de outras palavras.

O mecanismo de self-attention faz algo semelhante: ele permite que o modelo pondere a importância de cada palavra na sequência de entrada em relação a todas as outras palavras, calculando um "peso" de atenção. Isso significa que, ao processar uma palavra, o modelo pode "olhar" para qualquer outra palavra na frase e decidir o quão relevante ela é para o contexto atual.

Captura de Dependências

Permite capturar dependências entre palavras que estão muito distantes na frase, algo que era um desafio para as RNNs.

Paralelização

A arquitetura Transformer é altamente paralelizável, o que significa que pode ser treinada de forma muito mais eficiente em grandes conjuntos de dados.

Escalabilidade

Acelera o desenvolvimento de modelos cada vez maiores e mais capazes, permitindo o surgimento dos LLMs que hoje dominam o cenário do PLN.

É essa inovação que permitiu o surgimento dos LLMs que hoje dominam o cenário do PLN.

Vieses e Aplicações Éticas dos LLMs

Com o poder sem precedentes dos LLMs, surge uma responsabilidade igualmente grande. Esses modelos são treinados em vastas quantidades de dados da internet, que, infelizmente, contêm e refletem os vieses e preconceitos presentes na sociedade humana. Isso significa que os LLMs podem, inadvertidamente, perpetuar ou até amplificar esses **vieses** em suas saídas, seja em termos de gênero, raça, cultura, religião ou outras características demográficas.



Exemplos de Vieses

Um LLM pode associar certas profissões a um gênero específico ("o médico e a enfermeira") ou gerar respostas que reforçam estereótipos prejudiciais.



Detecção e Mitigação

Detectar e mitigar esses vieses é um campo de pesquisa ativo e crucial. A questão ética vai além dos vieses; envolve também a transparência sobre como esses modelos funcionam.



Uso Responsável



A explicabilidade de suas decisões e o uso responsável em aplicações críticas são fundamentais para garantir que essas ferramentas sejam usadas de forma ética.

As **aplicações éticas** dos LLMs exigem que desenvolvedores e usuários estejam cientes de suas limitações e potenciais impactos negativos. Isso inclui a implementação de filtros de conteúdo, a auditoria regular dos modelos para vieses, o desenvolvimento de diretrizes para seu uso e a educação sobre suas capacidades e falhas.

A discussão sobre o impacto social e ético dos LLMs é tão importante quanto o seu desenvolvimento tecnológico, garantindo que essas ferramentas poderosas sejam usadas para o bem da humanidade.

O Pipeline Clássico no Mundo dos LLMs: Ainda Relevante?

Diante da ascensão dos LLMs e sua capacidade de processar linguagem de forma tão sofisticada, é natural questionar: o pipeline clássico de PLN ainda tem relevância? A resposta é um retumbante **sim!** Embora os LLMs tenham absorvido muitas das funções que antes eram etapas separadas no pipeline clássico, os conceitos e técnicas fundamentais continuam sendo a base para muitas aplicações e para a compreensão de como os LLMs funcionam.

  **Analogia das Ferramentas:** Pense no pipeline clássico como as ferramentas básicas de um artesão. Mesmo que ele agora tenha máquinas de alta tecnologia, o conhecimento sobre como usar um martelo, uma serra ou um cinzel ainda é fundamental.

Integração com LLMs

- A tokenização ainda é uma etapa crucial, mesmo que os LLMs usem tokenizadores mais avançados (como subpalavras)
- O stemming e a lematização, embora não sejam explicitamente aplicados como etapas separadas dentro de um LLM pré-treinado, os princípios de normalização de palavras são inerentes à forma como os modelos aprendem representações de palavras

Aplicações Específicas

- Para muitas tarefas específicas, como a extração de entidades nomeadas ou a classificação de texto em domínios muito específicos com poucos dados, as abordagens clássicas ainda podem ser mais eficientes
- Servem como um pré-processamento valioso para modelos menores

O conhecimento do pipeline clássico também é essencial para o **fine-tuning** de LLMs, onde entendemos como preparar os dados de entrada e interpretar as saídas. Em suma, o pipeline clássico não foi substituído; ele foi integrado e, em muitos casos, aprimorado pelas abordagens neurais, formando uma sinergia poderosa.

Consolidação: A Base para o Futuro do PLN

Chegamos ao fim da nossa exploração pelo pipeline clássico de Processamento de Linguagem Natural. Vimos como cada etapa – desde a tokenização e normalização de palavras com stemming e lematização, passando pela representação numérica com BoW e TF-IDF, até a análise sintática com POS tagging e parsing – contribui para transformar o texto bruto em dados compreensíveis para as máquinas. Reconhecemos as limitações desses modelos, que nos levaram à necessidade de abordagens neurais e, finalmente, à revolução dos Modelos de Linguagem de Grande Escala (LLMs) e à arquitetura Transformer.

Tokenização
Quebra do texto em unidades menores

Abordagens Neurais
LLMs e Transformers para compreensão profunda



Normalização

Stemming e lematização para reduzir variações

Representação

BoW e TF-IDF para vetorização

Análise Sintática

POS tagging e parsing para estrutura

Em prática: A compreensão desses fundamentos permite que você não apenas utilize ferramentas de PLN de forma mais eficaz, mas também entenda os desafios e as oportunidades que surgem com as novas tecnologias. Seja para otimizar a busca de informações, desenvolver assistentes virtuais ou analisar grandes volumes de texto, o conhecimento do pipeline clássico é a base para qualquer profissional da área.

Autoavaliação

1 Qual das seguintes etapas do pipeline clássico de PLN tem como principal objetivo reduzir as palavras às suas formas base, utilizando um vocabulário e análise morfológica para garantir que a forma resultante seja uma palavra válida?

- a) Tokenização
- b) Stemming
- c) Lematização
- d) Part-of-Speech Tagging

2 O modelo Bag-of-Words (BoW) representa um documento como uma coleção de suas palavras, ignorando qual característica fundamental da linguagem?

- a) A frequência das palavras
- b) A ordem das palavras
- c) A presença de pontuação
- d) O tamanho do vocabulário

3 Qual das seguintes técnicas é utilizada para atribuir um peso numérico a cada palavra em um documento, valorizando palavras que são frequentes em um documento específico, mas raras em todo o conjunto de documentos?

- a) Bag-of-Words (BoW)
- b) Part-of-Speech (POS) Tagging
- c) Term Frequency-Inverse Document Frequency (TF-IDF)
- d) Parsing

4 A arquitetura Transformer, fundamental para os Modelos de Linguagem de Grande Escala (LLMs), revolucionou o PLN principalmente pela introdução de qual mecanismo?

- a) Redes Neurais Recorrentes (RNNs)
- b) Convoluções (CNNs)
- c) Mecanismo de atenção (self-attention)
- d) Algoritmos de stemming

5 Explique brevemente por que, mesmo com o avanço dos Modelos de Linguagem de Grande Escala (LLMs), o conhecimento sobre o pipeline clássico de PLN ainda é considerado relevante e útil.

Gabarito

Questão 1

c)
Lematização

Questão 2

**b) A ordem
das palavras**

Questão 3

**c) Term
Frequency-
Inverse
Document
Frequency
(TF-IDF)**

Questão 4

**c) Mecanismo
de atenção
(self-
attention)**

Conexão com a Próxima Aula



Próxima Aula: Na próxima aula, "Aula 3 – Representação Vetorial de Palavras: Word Embeddings", aprofundaremos como as abordagens neurais superam as limitações dos modelos clássicos, explorando a criação e o uso de representações densas de palavras que capturam o significado semântico e as relações contextuais.

Recursos Adicionais



Artigo Fundamental

"Attention Is All You Need"
(Vaswani et al., 2017): Para entender a base do Transformer.



Documentação Técnica

Biblioteca NLTK (Natural Language Toolkit): Para exemplos práticos de tokenização, stemming e lematização em Python.



Publicações de Pesquisa

OpenAI, Meta AI, Google AI: Para acompanhar as últimas tendências e desenvolvimentos em LLMs.

NOTA IMPORTANTE: As informações técnicas desta aula estão atualizadas até 2025. Consulte sempre fontes oficiais e publicações científicas para verificar alterações e novos desenvolvimentos na área.