

Nota: Este material complementar, disponível em https://www.rettore.com.br/public_data/lectures/ representa uma cópia resumida de conteúdos bibliográficos disponíveis gratuitamente na Internet.

Processamento de Linguagem Natural e Modelos de Linguagem

[Introdução](#)

[Conceitos Fundamentais](#)

[Técnicas de PLN](#)

[Modelos de Linguagem](#)

[LaMDA \(Language Model for Dialogue Applications\)](#)

[BERT \(Bidirectional Encoder Representations from Transformers\)](#)

[BUM \(Bidirectional Unsupervised Model\)](#)

[GPT \(Generative Pre-trained Transformer\)](#)

[Comparação Geral](#)

[Conclusão](#)

[Aplicações práticas](#)

[Conclusão](#)

[Referências](#)

Introdução

O **Processamento de Linguagem Natural (PLN)** é uma área da Inteligência Artificial que visa a interação entre computadores e a linguagem humana. Ele permite que máquinas compreendam, interpretem e respondam ao texto ou à fala de maneira semelhante a como as pessoas fazem. Para que isso ocorra, o PLN se apoia em vários conceitos e técnicas que otimizam o processamento de grandes volumes de dados textuais.

Conceitos Fundamentais

- **Corpus:** É o conjunto de dados textuais usados para treinamento e análise dos modelos de linguagem. Serve como base para construir o entendimento da máquina sobre a linguagem.
- **Tokenization (Tokenização):** É o processo de dividir o texto em unidades menores, como palavras ou subpalavras, chamadas *tokens*. Isso facilita a análise subsequente do texto.
- **Normalization (Normalização):** Consiste em padronizar o texto, como transformar tudo em minúsculas, remover pontuação ou acentuação, para evitar variações desnecessárias.
- **n-grams:** São sequências de n palavras ou caracteres que aparecem juntas em um texto, usadas para identificar padrões e relações frequentes.
- **Lexicons:** Conjuntos de palavras e suas informações associadas, como significado ou sentimento, usados para enriquecer a análise de texto.

- **Pré-processamento e Limpeza de Dados:** Inclui etapas como remoção de ruído (ex.: pontuação ou números) e tratamento de erros gramaticais para garantir que o texto seja adequado para análise.

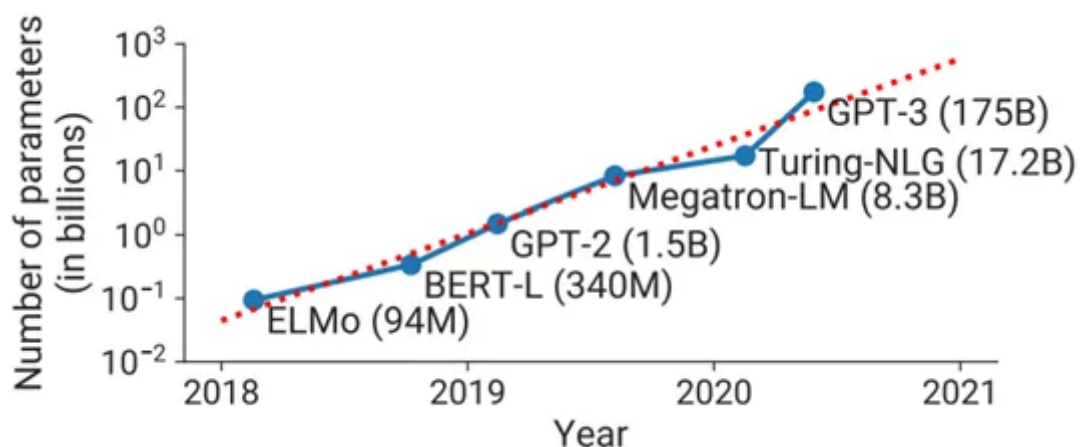
Técnicas de PLN

- **Stop Words:** São palavras muito frequentes e geralmente sem valor informativo (como "de", "o", "e") que podem ser removidas para simplificar a análise.
- **Stemming e Lemmatization:** Processos que reduzem as palavras às suas formas básicas. *Stemming* corta as palavras para suas raízes (ex.: "correr" vira "corr"), enquanto *lemmatization* reduz as palavras às suas formas canônicas (ex.: "corri", "correndo" viram "correr").
- **Bag of Words:** Um modelo simples que trata o texto como uma coleção de palavras, sem considerar a ordem, e conta a frequência de cada palavra.
- **TF-IDF (Term Frequency - Inverse Document Frequency):** Técnica que calcula a importância de uma palavra em um texto levando em conta sua frequência no documento e na coleção total de documentos.
- **Word Embeddings (Word2Vec, GloVe):** Representações vetoriais de palavras que capturam relações semânticas entre elas, permitindo que palavras com significados semelhantes fiquem próximas no espaço vetorial.
- **Named Entity Recognition (NER):** Técnica que identifica entidades nomeadas em um texto, como pessoas, lugares ou organizações.
- **Parts-of-speech (POS) Tagging:** Processo que atribui uma categoria gramatical (como substantivo, verbo, adjetivo) a cada palavra do texto.

Modelos de Linguagem

- **Transformers (BERT, GPT):** Modelos avançados que revolucionaram o PLN. *BERT* (Bidirectional Encoder Representations from Transformers) entende o contexto das palavras em ambas as direções (antes e depois). *GPT* (Generative Pre-trained Transformer) é um modelo gerador capaz de produzir textos coerentes com base em contextos anteriores.

Fonte: <http://arxiv.org/abs/2104.04473>



LaMDA (Language Model for Dialogue Applications)

LaMDA, criado pelo Google, é um modelo especializado em **conversas interativas**, focado em entender o contexto das interações humanas e gerar respostas mais naturais e relevantes em diálogos. Ele é treinado para capturar nuances do diálogo e manter conversas de longa duração de forma coerente.

- **Finalidade:** Diferente de outros modelos de linguagem que são projetados para realizar tarefas gerais de PLN, o LaMDA foi projetado especificamente para manter **diálogos fluentes** e naturais, entendendo contextos e trocas complexas, o que o torna ideal para assistentes virtuais e chatbots.
- **Diferencial:** Enquanto modelos como BERT ou GPT focam no entendimento do texto ou geração de linguagem geral, o LaMDA foca em compreender a **intenção do usuário** ao longo de uma conversa, garantindo respostas mais humanas e precisas.
- **Aplicações:** Utilizado em chatbots avançados e assistentes virtuais que precisam se engajar em **conversas naturais** e fluídas, respondendo de maneira inteligente e coerente ao longo de múltiplos turnos de diálogo.

BERT (Bidirectional Encoder Representations from Transformers)

BERT é um modelo da Google introduzido em 2018, revolucionando o PLN ao ser o primeiro grande modelo de linguagem bidirecional.

- **Como funciona:** O BERT é treinado para entender o **contexto de uma palavra** olhando para os termos que vêm antes e depois dela. Ele usa uma arquitetura de **transformer** que processa o texto bidirecionalmente (diferente de modelos unidimensionais que leem da esquerda para a direita ou vice-versa).
- **Objetivo:** O foco do BERT é **compreensão de linguagem**. Ele é excelente em tarefas como resposta a perguntas, classificação de textos, e preenchimento de lacunas em sentenças.
- **Tarefas:** Ele é amplamente utilizado em mecanismos de busca, assistentes virtuais, e outras aplicações onde a **compreensão precisa do significado** das palavras no contexto é crucial. Por exemplo, o Google Search usa BERT para melhorar a compreensão das consultas dos usuários.

BUM (Bidirectional Unsupervised Model)

O modelo **BUM** (Bidirectional Unsupervised Model) é menos amplamente discutido em comparação com LaMDA, BERT, ou GPT, mas segue uma estrutura bidirecional similar ao BERT, com o foco em aprendizado não supervisionado para a compreensão de texto.

- **Como funciona:** Assim como o BERT, o BUM utiliza uma abordagem bidirecional, o que significa que ele processa o texto observando tanto o contexto anterior quanto o posterior de uma palavra.
- **Objetivo:** A ideia do BUM é capturar a essência do texto sem necessidade de grandes volumes de dados rotulados, permitindo que ele seja usado de forma mais eficiente em aplicações que envolvem **aprendizado não supervisionado**.
- **Aplicações:** Pode ser utilizado para pré-treinar modelos que realizam tarefas como classificação de texto, sumarização e até mesmo tradução, especialmente em situações onde há uma **escassez de dados anotados**.

GPT (Generative Pre-trained Transformer)

GPT, desenvolvido pela OpenAI, é uma série de modelos de linguagem que evoluíram ao longo dos anos, com versões como GPT, GPT-2, GPT-3 e o mais recente GPT-4, trazendo melhorias significativas em suas capacidades.

- **Como funciona:** Ao contrário do BERT, que é bidirecional, o **GPT é um modelo unidirecional** (no início de sua formação), o que significa que ele processa o texto palavra por palavra, prevendo a próxima palavra com base nas anteriores. Isso o torna excelente para **geração de texto**.
- **Finalidade:** O GPT é treinado para gerar texto de alta qualidade e coerência, produzindo desde frases simples até artigos longos ou histórias. Ele é altamente eficaz em **tarefas de geração**, como redação de textos, criação de conteúdo e diálogo.
- **Transformers:** O GPT usa a arquitetura de **transformers**, que lhe dá a capacidade de processar grandes quantidades de dados de forma eficiente e gerar respostas que parecem muito naturais e contextualizadas.
- **Aplicações:** GPT é usado em **assistentes de escrita**, chatbots avançados, ferramentas de automação de conteúdo, tradução, resumo de textos e muitas outras áreas. A OpenAI também introduziu o **ChatGPT**, uma aplicação prática do GPT, que interage com os usuários em conversas.

Comparação Geral

- **LaMDA** é especializado em conversas e diálogos, focando em manter uma comunicação fluida e relevante em longas interações.
- **BERT** é otimizado para entender o contexto e o significado das palavras, sendo amplamente utilizado em mecanismos de busca e compreensão de linguagem.
- **BUM** se concentra no aprendizado não supervisionado e na captura de contexto bidirecional para otimizar tarefas de PLN.
- **GPT** é um modelo poderoso de geração de texto que é capaz de criar conteúdo fluido e complexo com base em uma arquitetura unidirecional.

Conclusão

Cada um desses modelos tem uma abordagem única para o **processamento de linguagem natural**, e suas diferenças os tornam adequados para aplicações distintas. BERT é excelente para tarefas de compreensão de linguagem, enquanto GPT se destaca na geração de texto. LaMDA e BUM têm seus próprios nichos, com LaMDA focando no diálogo e BUM em aprendizado não supervisionado.

Aplicações práticas

1. **Análise de Sentimentos:** Esta aplicação usa técnicas de PLN para identificar as emoções ou opiniões em textos, como comentários em redes sociais, avaliações de produtos ou feedbacks de clientes. É útil para monitorar a percepção do público sobre marcas ou produtos.
2. **Tradução Automática:** Modelos de linguagem são utilizados para traduzir textos de um idioma para outro de maneira automática. Ferramentas como Google Translate fazem uso de modelos avançados, como transformers, para melhorar a precisão das traduções.
3. **Chatbots:** Chatbots inteligentes, como assistentes virtuais, utilizam PLN para entender e responder às perguntas dos usuários de forma natural e eficiente. Com o uso de modelos como o GPT, eles conseguem gerar respostas mais humanas e adequadas ao contexto.

Conclusão

O PLN está no centro de muitas inovações tecnológicas que transformam a maneira como interagimos com as máquinas. Com técnicas avançadas como tokenização, embeddings e transformers, modelos de linguagem são capazes de realizar desde tarefas simples, como a análise de sentimentos, até desafios mais complexos, como a tradução automática e o desenvolvimento de chatbots.

Referências

1. Jurafsky, D., & Martin, J. H. (2020). Speech and Language Processing (3rd ed.).
2. Caseli, H.M.; Nunes, M.G.V. (org.) Processamento de Linguagem Natural: Conceitos, Técnicas e Aplicações em Português. 2 ed. BPLN, 2024. Disponível em: <https://brasileiraspln.com/livro-pln/2a-edicao>.
3. Conteúdo gratuito disponível na Internet:
https://www.youtube.com/watch?v=LCJ_PfmeO5Q;
<https://www.youtube.com/watch?v=OS5iRU0ZhWY>

Isenção de Responsabilidade:

Os autores deste documento não reivindicam a autoria do conteúdo original compilado das fontes mencionadas. Este documento foi elaborado para fins educativos e de referência, e todos os créditos foram devidamente atribuídos aos respectivos autores e fontes originais.

Qualquer utilização comercial ou distribuição do conteúdo aqui compilado deve ser feita com a devida autorização dos detentores dos direitos autorais originais. Os compiladores deste documento não assumem qualquer responsabilidade por eventuais violações de direitos autorais ou por quaisquer danos decorrentes do uso indevido das informações contidas neste documento.

Ao utilizar este documento, o usuário concorda em respeitar os direitos autorais dos autores originais e isenta os compiladores de qualquer responsabilidade relacionada ao conteúdo aqui apresentado.