# Básico



### Fundamentos de modelos de linguagem

Os modelos de linguagem são uma classe de modelos de aprendizado de máquina que são treinados para entender e gerar texto em linguagem natural. Eles são fundamentais para uma variedade de aplicações, desde assistentes de voz e chatbots até tradução automática e análise de sentimentos. Aqui estão alguns dos conceitos fundamentais associados aos modelos de linguagem.

O que são modelos de linguagem?: Os modelos de linguagem são sistemas que são treinados para prever a próxima palavra em uma sequência, dadas as palavras anteriores. Eles aprendem a estrutura da linguagem durante o treinamento, o que os permite gerar texto coerente e compreensível.

**Tokenização:** A tokenização é o processo de dividir o texto em "tokens", que são tipicamente palavras ou caracteres individuais. A tokenização permite que o modelo processe o texto em partes menores e mais gerenciáveis.

**Vetorização:** A vetorização é o processo de transformar texto ou tokens em vetores numéricos que podem ser alimentados em um modelo de aprendizado de máquina. Isso geralmente envolve atribuir a cada palavra um vetor numérico exclusivo em um espaço de alta dimensão.

**Embeddings de palavras:** Os embeddings de palavras são uma forma específica de vetorização onde palavras semelhantes são mapeadas para pontos próximos em um espaço vetorial. Isso permite ao modelo entender a semântica das palavras, incluindo sinônimos, antônimos e outras relações semânticas.

**Arquiteturas de modelagem de linguagem:** Existem várias arquiteturas diferentes usadas para construir modelos de linguagem, incluindo redes neurais recorrentes (RNNs), long short-term memory networks (LSTMs), e mais recentemente, Transformers. Cada um tem seus próprios pontos fortes e fracos, e a escolha da arquitetura depende do problema específico que se está tentando resolver.

Treinamento e fine-tuning: Os modelos de linguagem são geralmente treinados em um grande corpus de texto, como a Wikipédia ou a web em geral. Eles aprendem a estrutura da linguagem prevendo a próxima palavra em uma sequência. Após o treinamento inicial, os modelos de linguagem podem ser "fine-tuned" em um conjunto de dados mais específico para se especializar em uma tarefa particular.

Modelos de aprendizado profundo: A maioria dos modelos de linguagem modernos usa alguma forma de aprendizado profundo, que envolve alimentar os dados através de múltiplas camadas de redes neurais para aprender representações complexas. Modelos como o BERT e o GPT-3 são exemplos de modelos de aprendizado profundo que alcançaram desempenho de ponta em uma série de tarefas de processamento de linguagem natural.

Limitações e desafios: Embora os modelos de linguagem tenham feito grandes progressos, ainda há muitos desafios a serem enfrentados. Estes incluem lidar com ambiguidade e nuances na linguagem humana, compreendendo o contexto e o sentido de uma conversa, e garantindo que os modelos são éticos e justos em suas respostas.

Os fundamentos dos modelos de linguagem estão em constante evolução à medida que os pesquisadores continuam a explorar novas técnicas e abordagens. No entanto, o objetivo final permanece o mesmo: desenvolver modelos que possam entender e gerar linguagem humana com a maior precisão e naturalidade possível.

## Arquitetura Transformer e sua importância para o GPT

A arquitetura Transformer é um componente crucial no desenvolvimento de modelos de linguagem modernos, como o Generative Pretrained Transformer (GPT) desenvolvido pela OpenAI. Ela representa um avanço significativo no processamento de linguagem natural, introduzindo uma nova maneira de lidar com a sequencialidade e a dependência de contexto na linguagem humana.

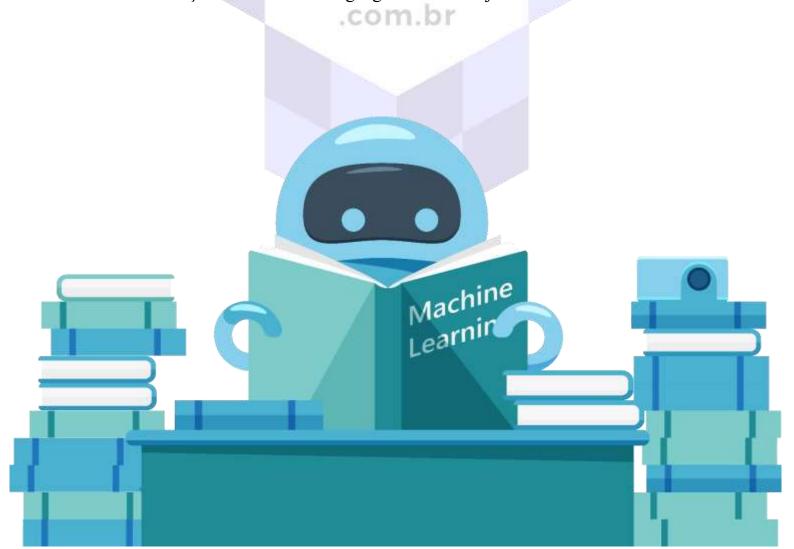
O que é a arquitetura Transformer? A arquitetura Transformer foi introduzida em 2017 por Vaswani et al. no artigo "Attention is All You Need". O Transformer descarta completamente a recorrência e, em vez disso, depende inteiramente de mecanismos de atenção auto-regressiva para extrair a estrutura sequencial do texto. Isso torna a arquitetura muito mais paralelizável e, portanto, mais eficiente em termos de computação, em comparação com as arquiteturas de rede neural recorrentes e de memória de curto e longo prazo (LSTM) usadas anteriormente para tarefas de processamento de linguagem.

Atenção e Auto-Atenção A atenção é um conceito que permite ao modelo se concentrar em diferentes partes do texto de entrada com base em sua relevância para a tarefa em questão. No contexto da arquitetura Transformer, a auto-atenção, também conhecida como atenção de múltiplas cabeças, permite que o modelo pondere a importância de cada palavra em relação a todas as outras palavras na sequência, resultando em uma representação mais rica e contextualizada do texto.

**Importância para o GPT** Os modelos GPT, do GPT-1 ao GPT-4, baseiamse na arquitetura Transformer. Eles utilizam uma variação conhecida como "Transformer de decodificador único", que emprega somente a parte de decodificação do Transformer original. A arquitetura Transformer permite que o GPT processe grandes quantidades de texto de maneira eficiente e capture dependências de longo alcance entre palavras. Além disso, o mecanismo de auto-atenção permite que o GPT modele o contexto em que uma palavra aparece, o que é essencial para tarefas como prever a próxima palavra em uma frase ou gerar texto coerente e relevante.

Em particular, a capacidade do GPT de gerar texto realista e coerente é amplamente atribuída à eficácia da arquitetura Transformer em modelar dependências contextuais. Essa característica permitiu aplicações de geração de texto em uma variedade de domínios, desde a composição de poemas e histórias até a geração de código de programação.

A arquitetura Transformer desempenha um papel crucial nos modelos GPT, permitindo-lhes processar eficientemente grandes quantidades de texto e capturar a complexidade e a dependência do contexto da linguagem humana. Ela representou um grande avanço na área de processamento de linguagem natural e continua a ser um componente chave em muitos dos mais avançados modelos de linguagem em uso hoje.



### Entendimento dos aspectos técnicos ficação de posição, atenção de múltiplas cabe

### Codificação de posição, atenção de múltiplas cabeças, camadas de transformadores

Os modelos baseados na arquitetura Transformer, como o GPT, incorporam várias técnicas avançadas de aprendizado profundo para lidar com a complexidade da linguagem natural. Aqui, vamos nos concentrar em três aspectos-chave: codificação de posição, atenção de múltiplas cabeças e camadas de transformadores.

Codificação de Posição A codificação de posição é usada para fornecer informações de sequência aos Transformers, que são inerentemente não sequenciais. A codificação de posição permite que o modelo saiba não apenas quais palavras estão presentes em uma sequência, mas também onde essas palavras ocorrem. Isso é importante porque a ordem das palavras pode mudar o significado de uma frase.

A codificação de posição é implementada adicionando vetores de posição aos embeddings de palavras antes que eles sejam passados para as camadas de auto-atenção. Esses vetores de posição são criados de tal forma que eles permitem que o modelo distingua a posição de cada palavra na sequência e também capture a distância relativa entre as palavras.

Atenção de Múltiplas Cabeças A atenção é um mecanismo que permite que os modelos ponderem a importância relativa de diferentes palavras em uma sequência. A atenção de múltiplas cabeças é uma extensão deste conceito que permite ao modelo concentrar-se em diferentes partes da sequência para diferentes aspectos da tarefa.

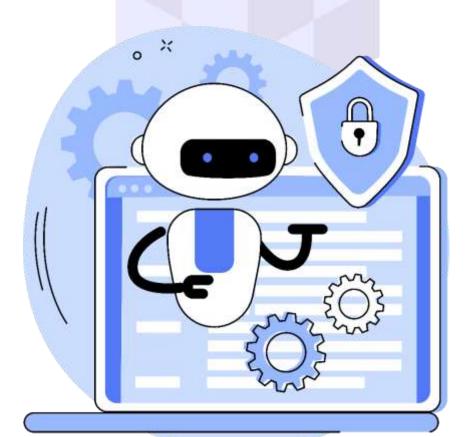
Por exemplo, ao prever a próxima palavra em uma frase, o modelo pode precisar prestar atenção em uma palavra anterior que é diretamente relevante, bem como em outras palavras que fornecem contexto geral. A atenção de múltiplas cabeças permite ao modelo fazer isso de forma eficaz, prestando atenção em várias palavras ao mesmo tempo.

Camadas de Transformadores Um Transformer é composto por várias camadas, cada uma das quais contém sua própria atenção de múltiplas cabeças e redes neurais feed-forward. Cada camada recebe a saída da camada anterior, permitindo que o modelo aprenda representações cada vez mais complexas.

As camadas inferiores tendem a aprender padrões mais simples, como a sintaxe, enquanto as camadas superiores aprendem padrões mais complexos, como a semântica e o contexto. Isso permite que os Transformers capturem a complexidade da linguagem natural.

Cada camada de um Transformer pode ser vista como um módulo de aprendizado que contribui para a capacidade geral do modelo de entender e gerar texto. O número de camadas em um Transformer pode variar, mas os modelos como o GPT-3 têm dezenas de camadas, permitindo-lhes aprender representações muito ricas e complexas do texto.

A arquitetura Transformer, com suas técnicas de codificação de posição, atenção de múltiplas cabeças e múltiplas camadas de transformadores, oferece uma maneira poderosa e flexível de modelar a linguagem humana. Essas técnicas permitem que o modelo capture tanto a estrutura sequencial quanto o contexto da linguagem, tornando-o capaz de gerar texto que é surpreendentemente coerente e humano em seu estilo e conteúdo.



## Conceito de tokens e como o GPT-4 manipula sequências de tokens

Tokens são uma unidade fundamental no processamento de linguagem natural (NLP). Eles podem ser pensados como as peças que compõem uma linguagem. Dependendo do contexto, um token pode ser uma palavra, uma parte de uma palavra (por exemplo, subpalavras ou caracteres), ou mesmo uma sequência mais longa de palavras.

#### Conceito de Tokens

No contexto da NLP, a tokenização é o processo de dividir o texto em pedaços ou "tokens". Cada token pode ser considerado como uma unidade independente de significado. Por exemplo, em inglês, uma maneira comum de tokenizar uma frase é dividi-la em palavras. Assim, a frase "ChatGPT é incrível!" seria dividida em quatro tokens: ["ChatGPT", "é", "incrível", "!"].

No entanto, a tokenização nem sempre é feita ao nível da palavra. Às vezes, pode ser mais útil tokenizar em subpalavras ou caracteres. Por exemplo, a palavra "incrível" poderia ser dividida em ["incr", "í", "vel"].

#### **Tokens e GPT-4**

O GPT-4, como os modelos de linguagem anteriores da série GPT, utiliza tokens como a unidade básica de processamento. O GPT-4 é treinado para prever o próximo token em uma sequência, dado os tokens anteriores. Assim, para gerar uma resposta, o modelo começa com uma sequência de tokens de entrada, produz um token de saída, adiciona esse token à sequência de entrada e repete o processo.

Para tokenizar o texto, o GPT-4 usa uma técnica chamada Byte-Pair Encoding (BPE), que tokeniza o texto em subpalavras. O BPE permite que o modelo lide com palavras que não viu antes, dividindo-as em subpalavras que conhece. Também ajuda a reduzir o tamanho do vocabulário do modelo, o que pode melhorar a eficiência computacional.

Além disso, para preservar a informação da sequência no processamento de tokens, o GPT-4 usa codificações de posição, que são adicionadas aos embeddings de token antes que eles sejam passados para as camadas de auto-atenção. Isso permite que o modelo saiba onde cada token ocorre na sequência.

### Importância dos Tokens

A manipulação de tokens é fundamental para o funcionamento do GPT-4. Permite que o modelo entenda a estrutura da linguagem e faça previsões coerentes. Também torna possível para o modelo lidar com uma variedade de tarefas de linguagem, desde a geração de texto até a tradução automática e a resposta a perguntas, simplesmente alterando a sequência de tokens de entrada.

Os tokens são uma unidade fundamental no processamento de linguagem natural, e a manipulação de sequências de tokens é central para o funcionamento do GPT-4. Através de técnicas como a tokenização BPE e as codificações de posição, o GPT-4 é capaz de entender e gerar texto de uma maneira que é surpreendentemente coerente e semelhante à linguagem humana.

