

# Processamento de Linguagem Natural por Redes Neurais

Prof. Eduardo Bezerra  
CEFET/RJ

22 de outubro de 2025

## Conteúdo

1	Linguagem como sequência discreta de tokens	2
2	Representação One-hot Encoding	3
3	Representações Distribuídas: Word2Vec e GloVe	5
4	Redes Neurais Recorrentes (RNN) e LSTMs	7
5	Mecanismo de Atenção e o Surgimento dos Transformers	10

# 1 Linguagem como sequência discreta de tokens

## A natureza discreta da linguagem

A linguagem natural é composta por símbolos discretos (palavras, caracteres ou subpalavras) organizados em sequência. Cada sentença pode ser vista como uma lista ordenada de unidades elementares denominadas **tokens**. O processo de transformar texto em tokens é chamado **tokenização**.

Por exemplo, a frase:

“As redes neurais aprendem representações de linguagem.”

pode ser tokenizada de diferentes formas, dependendo da estratégia adotada:

- **Tokenização por palavra:** [As, redes, neurais, aprendem, representações, de, linguagem]
- **Tokenização por caractere:** [A, s, \_, r, e, d, e, s, \_, n, e, ...]
- **Tokenização subpalavra (BPE):** [As, redes, ne, ur, ais, apren, dem, repre, sent, ações]

O vocabulário resultante de um corpus é o conjunto de todos os tokens distintos observados. Para modelos neurais, cada token é associado a um índice inteiro no intervalo  $[0, |V| - 1]$ , onde  $|V|$  é o tamanho do vocabulário.

## Implicações para o aprendizado

Essa discretização impõe um desafio fundamental: os tokens não possuem estrutura numérica inerente. Modelos de aprendizado de máquina operam sobre números reais, não sobre símbolos. É necessário, portanto, definir um mapeamento entre tokens e vetores numéricos — o que dá origem ao problema da representação, explorado na próxima seção.

## Representação de sequência

Formalmente, uma sentença pode ser representada como uma sequência ordenada:

$$S = [t_1, t_2, \dots, t_n],$$

onde cada  $t_i$  é um token pertencente ao vocabulário  $V$ . Essa estrutura sequencial é o ponto de partida para qualquer modelo neural de linguagem.

Modelos simples tratam os tokens como independentes. Já os modelos de linguagem modernos capturam dependências entre elementos da sequência, estimando a probabilidade conjunta:

$$P(S) = P(t_1, t_2, \dots, t_n) = \prod_{i=1}^n P(t_i \mid t_1, \dots, t_{i-1}).$$

Essa decomposição autorregressiva é central para o aprendizado de representações linguísticas em redes neurais.

## Síntese

A tokenização traduz o texto para uma forma manipulável por modelos computacionais, preservando sua natureza sequencial. O próximo passo é definir como cada token será representado numericamente, o que nos leva à questão das representações (codificações) vetoriais.

# 2 Representação One-hot Encoding

## Do símbolo ao número

Após a tokenização, cada palavra ou subpalavra é representada por um identificador inteiro. Contudo, redes neurais exigem entradas vetoriais em  $\mathbb{R}^n$ . É necessário, portanto, converter cada token em um vetor numérico.

A forma mais direta de fazer isso é por meio da **codificação one-hot**. Suponha um vocabulário com  $V$  palavras distintas:

$$V = \{\text{gato}, \text{cachorro}, \text{pássaro}, \text{peixe}\}.$$

Cada palavra é representada por um vetor de dimensão  $|V|$ , contendo um único 1 na posição correspondente à palavra, e 0 nas demais. Por exemplo:

$$\text{gato} = [1, 0, 0, 0], \quad \text{cachorro} = [0, 1, 0, 0].$$

## Limitações da codificação one-hot

Embora simples, essa representação apresenta problemas fundamentais:

1. **Ortogonalidade absoluta:** Todos os vetores são ortogonais, implicando que as palavras são igualmente distantes entre si. Não há noção de *similaridade semântica*. Isso significa que essa codificação não representa o fato de que *gato* e *cachorro* são mais próximos (no sentido semântico) do que *gato* e *avião*.
2. **Alta dimensionalidade:** O vetor tem dimensão igual ao tamanho do vocabulário, que pode conter dezenas ou centenas de milhares de palavras. Isso gera vetores extremamente esparsos, de custo computacional elevado.
3. **Ausência de generalização:** Como cada palavra é tratada de forma independente, o modelo não consegue generalizar o conhecimento aprendido sobre uma palavra para outra de significado próximo.

## Consequências práticas

Em tarefas de PLN, a codificação one-hot obriga o modelo a aprender do zero todas as relações entre palavras, sem qualquer estrutura prévia. Esse tipo de representação é suficiente apenas para modelos lineares simples ou experimentos didáticos, mas se torna inviável em aplicações reais de larga escala.

## Motivação para representações distribuídas

A limitação do one-hot levou ao surgimento das **representações distribuídas**, que são vetores densos e de baixa dimensão que capturam propriedades semânticas a partir dos dados. Esses vetores, chamados de *word embeddings*, são aprendidos de modo que palavras usadas em contextos semelhantes tenham representações próximas no espaço vetorial.

O próximo tópico introduz precisamente essa ideia.

## 3 Representações Distribuídas: Word2Vec e GloVe

### Ideia central

A principal limitação das representações one-hot é a ausência de semelhança semântica. A solução proposta pelas **representações distribuídas** é aprender vetores densos de baixa dimensão que capturem relações de significado entre palavras a partir dos próprios dados.

Cada palavra  $w$  é associada a um vetor  $v_w \in \mathbb{R}^d$ , onde  $d \ll V$ . Esses vetores são aprendidos de modo que palavras que ocorrem em contextos semelhantes fiquem próximas no espaço vetorial.

### Hipótese distribucional

A base teórica dessa abordagem é a **hipótese distribucional**:

*Palavras com significados semelhantes tendem a ocorrer em contextos semelhantes.*

Logo, em vez de tentar definir o significado de uma palavra explicitamente, deixamos o modelo inferi-lo estatisticamente a partir de seus contextos de ocorrência.

### Modelos clássicos de embeddings

**Word2Vec (Mikolov et al., 2013).** Aprende representações por meio de uma tarefa de predição local. Duas variantes principais:

- **CBOW (Continuous Bag-of-Words):** prevê uma palavra a partir do contexto vizinho;
- **Skip-gram:** prevê o contexto dado uma palavra central.

O treinamento é feito em grandes corpora textuais, ajustando os vetores para maximizar a probabilidade de coocorrência entre palavras próximas.

**GloVe (Global Vectors, Pennington et al., 2014).** Em vez de aprender localmente, o GloVe utiliza estatísticas globais de coocorrência. Define uma função de custo que aproxima o produto escalar entre vetores de palavras da razão entre frequências de coocorrência observadas:

$$v_i^\top v_j \approx \log(P_{ij}),$$

onde  $P_{ij}$  é a probabilidade de que a palavra  $j$  ocorra no contexto da palavra  $i$ .

## Propriedades notáveis

Os embeddings resultantes exibem **relações semânticas lineares**. Um exemplo clássico é:

$$v_{\text{rei}} - v_{\text{homem}} + v_{\text{mulher}} \approx v_{\text{rainha}}.$$

Essas relações emergem naturalmente porque os vetores capturam padrões de coocorrência e analogia entre palavras.

## Visualização

É possível projetar os vetores para duas dimensões (por PCA ou t-SNE) e observar agrupamentos semânticos: animais próximos de outros animais, verbos próximos entre si, etc. Essas estruturas revelam como o modelo aprende aspectos latentes da semântica sem instrução explícita.

## Limitação dos embeddings estáticos

Tanto Word2Vec quanto GloVe produzem um vetor fixo por palavra, independentemente do contexto. Assim, a palavra “banco” terá a mesma representação em:

“o banco fechou às quatro” e “o pescador sentou-se no banco”.

Modelos mais recentes, como o ELMo e o BERT, resolvem esse problema ao gerar **representações contextuais**, em que o significado depende da frase completa.

## Síntese

Os embeddings distribuídos substituíram as representações simbólicas por vetores aprendidos. Esse avanço permitiu que redes neurais processassem texto de maneira contínua, capturando semelhanças semânticas e estruturais — uma base essencial para os modelos modernos de linguagem.

## 4 Redes Neurais Recorrentes (RNN) e LSTMs

### Motivação

Modelos baseados em embeddings tratam cada palavra como independente, ignorando a ordem e as dependências entre tokens. No entanto, o significado de uma frase depende fortemente da sequência em que as palavras ocorrem. Para capturar essa estrutura sequencial, surgiram as **redes neurais recorrentes** (RNNs).

### Ideia central da recorrência

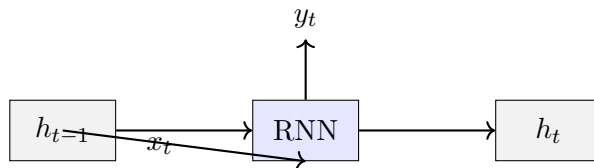
A RNN introduz a noção de **estado oculto** (*hidden state*), que é atualizado a cada novo token da sequência. Esse estado atua como uma memória que acumula informações do passado.

Para uma sequência de entrada  $x_1, x_2, \dots, x_T$ , a RNN define:

$$h_t = f(W_{xh}x_t + W_{hh}h_{t-1} + b_h),$$

$$y_t = W_{hy}h_t + b_y,$$

onde  $h_t$  é o estado oculto no instante  $t$ ,  $f$  é uma função de ativação não linear (geralmente tanh ou ReLU), e  $y_t$  é a saída correspondente.



*Diagrama simplificado de uma célula RNN.*

Essa estrutura permite que o modelo mantenha dependências temporais e processe sequências de comprimento variável.

#### Observação importante: célula $\times$ neurônio

Nas redes neurais recorrentes, o termo **célula** não se refere a um neurônio individual, mas sim a uma *unidade computacional composta*.

Um **neurônio** é a operação básica de uma rede neural: ele recebe um vetor de entradas, calcula uma combinação linear com pesos e aplica uma função de ativação. Formalmente:

$$y = \phi(w^\top x + b).$$

Uma **célula recorrente**, por outro lado, é um *bloco funcional* que contém vários neurônios e mantém um estado interno que evolui ao longo do tempo. Ela define como o estado  $h_t$  deve ser atualizado a partir do estado anterior  $h_{t-1}$  e da entrada atual  $x_t$ :

$$h_t = \phi(W_{xh}x_t + W_{hh}h_{t-1} + b_h).$$

No caso de uma LSTM, por exemplo, a célula inclui múltiplas portas (entrada, saída, esquecimento), cada uma composta por diversos neurônios com pesos próprios. Assim, uma célula é uma **estrutura que organiza um conjunto de neurônios** para realizar uma operação recorrente específica.

Em resumo:

- **Neurônio:** unidade elementar que aplica uma transformação não linear.
- **Célula:** conjunto de neurônios que, juntos, definem uma etapa de processamento temporal.



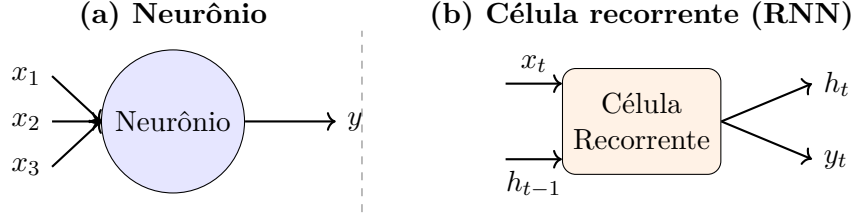


Figura 1: Diferença conceitual entre um neurônio e uma célula recorrente. O neurônio realiza uma transformação pontual, enquanto a célula recorrente mantém estado e processa sequências temporais.

## Treinamento e retropropagação no tempo

O treinamento das RNNs é feito via **Backpropagation Through Time** (BPTT), que estende o algoritmo de retropropagação para sequências. Durante o cálculo dos gradientes, o erro é propagado de volta através dos passos temporais.

Entretanto, como o gradiente é multiplicado repetidamente pelas derivadas das funções de ativação, ele tende a:

- **explodir**, levando a valores numéricos muito grandes; ou
- **desvanecer** (*vanish*), tornando-se próximo de zero.

Esses fenômenos dificultam o aprendizado de dependências longas.

## LSTMs e GRUs

Para contornar esse problema, foram propostas variantes recorrentes com mecanismos de controle de fluxo de informação.

**LSTM (Long Short-Term Memory).** Introduz *portas* (*gates*) que regulam o que deve ser armazenado, esquecido ou propagado:

$$\begin{aligned}
 f_t &= \sigma(W_f x_t + U_f h_{t-1} + b_f), & i_t &= \sigma(W_i x_t + U_i h_{t-1} + b_i), \\
 \tilde{c}_t &= \tanh(W_c x_t + U_c h_{t-1} + b_c), & c_t &= f_t \odot c_{t-1} + i_t \odot \tilde{c}_t, \\
 h_t &= o_t \odot \tanh(c_t),
 \end{aligned}$$

onde  $f_t, i_t, o_t$  são as portas de esquecimento, entrada e saída, respectivamente.

**GRU (Gated Recurrent Unit).** É uma versão mais simples, que combina algumas dessas portas em menos parâmetros, mantendo desempenho semelhante.

## Limitações das RNNs

Apesar dos avanços, as redes recorrentes apresentam limitações importantes:

- Dificuldade de paralelização: o processamento é sequencial.
- Dificuldade em capturar dependências de longo prazo.
- Alto custo de treinamento em grandes corpora.

Essas limitações motivaram o surgimento de arquiteturas que substituem a recorrência pelo mecanismo de atenção, base dos modelos Transformer.

## Síntese

As RNNs introduziram a noção de memória em redes neurais e foram o principal modelo de PLN por mais de uma década. Contudo, o custo computacional e o problema do desvanecimento de gradientes abriram caminho para uma nova abordagem, o **mecanismo de atenção**, tema do próximo tópico.

# 5 Mecanismo de Atenção e o Surgimento dos Transformers

## Limitação das arquiteturas recorrentes

As redes recorrentes (RNNs, LSTMs, GRUs) processam sequências de forma sequencial, i.e., um passo por vez. Isso impede o paralelismo e dificulta o aprendizado de dependências longas: quanto mais distante uma palavra está de outra, mais fraco se torna o gradiente que conecta suas representações.

Por exemplo, em frases como:

“O livro que o professor recomendou foi interessante.”

a relação entre *livro* e *interessante* envolve múltiplas etapas intermediárias, o que torna a captura dessa dependência difícil para RNNs.

## Ideia central da atenção

O mecanismo de **atenção** resolve esse problema permitindo que o modelo se concentre dinamicamente nas partes relevantes da sequência, independentemente da distância. A atenção introduz uma operação que calcula o quanto cada token deve “olhar” para os outros tokens da entrada.

Para cada posição  $i$ , o modelo combina as representações dos demais tokens, ponderadas por coeficientes de relevância:

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

onde:

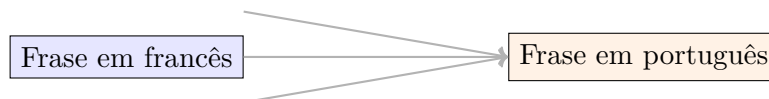
- $Q$  — *queries* (consultas);
- $K$  — *keys* (chaves);
- $V$  — *values* (valores);
- $d_k$  — dimensão das chaves.

Essa operação permite que cada palavra gere uma representação contextualizada, levando em conta as demais.

## Atenção como mecanismo de alinhamento

Atenção foi inicialmente proposta em tarefas de tradução automática, como uma forma de *alinhamento*: enquanto gera cada palavra da tradução, o modelo decide quais partes da frase original merecem mais foco.

### Atenção na tradução



*O modelo aprende a atribuir pesos de atenção que ligam partes correspondentes das frases.*

## Autoatenção (Self-Attention)

O passo seguinte foi aplicar a atenção dentro da própria sequência: cada palavra “olha” para todas as outras. Essa técnica é chamada de **autoatenção** e substitui completamente a recorrência.

A autoatenção permite paralelizar o processamento e capturar dependências de qualquer distância em uma única etapa de cálculo.

## Surgimento dos Transformers

Em 2017, Vaswani et al. apresentaram o artigo “*Attention is All You Need*”, que propôs a arquitetura **Transformer**. Essa arquitetura baseia-se inteiramente na autoatenção, sem qualquer estrutura recorrente ou convolucional.

As principais inovações foram:

- substituição da recorrência por múltiplas camadas de autoatenção;
- uso de **atenção multi-cabeça**, permitindo que o modelo aprenda múltiplos padrões de dependência simultaneamente;
- adição de **codificação posicional** para preservar a ordem dos tokens.

O resultado foi um salto de desempenho em tradução automática e, posteriormente, em praticamente todas as tarefas de PLN.

## Síntese

O mecanismo de atenção transformou a maneira como as redes neurais tratam sequências. Ao permitir que cada token interaja diretamente com todos os outros, a atenção eliminou as limitações das RNNs e tornou possível o treinamento em larga escala de modelos como BERT, GPT e T5.

*Na atenção, o modelo decide o que lembrar e o que ignorar, um princípio que aproxima as redes neurais da flexibilidade do pensamento humano.*

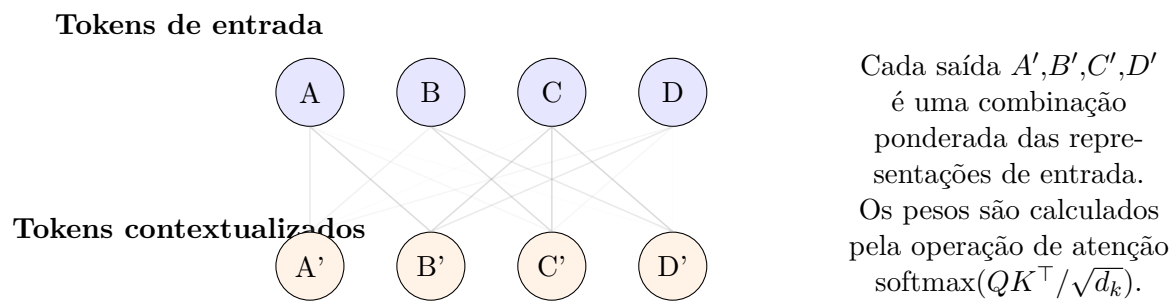


Figura 2: Mecanismo de autoatenção. Cada token combina informações de todos os outros tokens da sequência, com pesos aprendidos que indicam o grau de relevância contextual.