



Processamento de Linguagem Natural

Modelos Neurais em PLN

Marlo Souza¹

¹Universidade Federal da Bahia - Brasil

31 de julho de 2024



Sumário

Redes Neurais
Artificiais

Aplicando RNA para PLN

Arquitetura
Codificador-
Decodificador
(Encoder-
Decoder)

Mecanismo de
Atenção

Arquitetura
Transformer

Redes Neurais Artificiais
Aplicando RNA para PLN

Arquitetura Codificador-Decodificador (Encoder-Decoder)

Mecanismo de Atenção

Arquitetura Transformer



3

Redes Neurais Artificiais

Aplicando RNA para PLN

Arquitetura
Codificador-
Decodificador
(Encoder-
Decoder)

Mecanismo de
Atenção

Arquitetura
Transformer

Redes Neurais Artificiais

26



O neurônio

Redes Neurais Artificiais

Aplicando RNA para PLN

Arquitetura Codificador-Decodificador (Encoder-Decoder)

Mecanismo de Atenção

Arquitetura Transformer

4

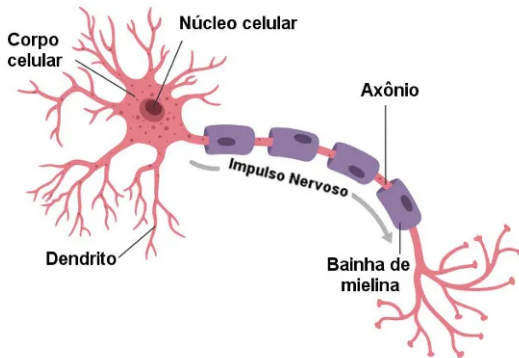


Figura: Representação simplificada de um neurônio



5

Redes Neurais
Artificiais

Aplicando RNA para PLN

Arquitetura
Codificador-
Decodificador
(Encoder-
Decoder)

Mecanismo de
Atenção

Arquitetura
Transformer

26

Um modelo matemático para o neurônio

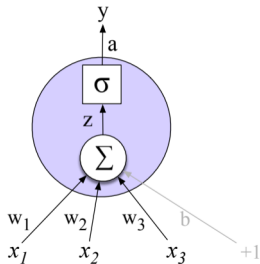


Figura: Um modelo matemático para o funcionamento do neurônio



5

Redes Neurais
Artificiais

Aplicando RNA para PLN

Arquitetura
Codificador-
Decodificador
(Encoder-
Decoder)

Mecanismo de
Atenção

Arquitetura
Transformer

26

Um modelo matemático para o neurônio

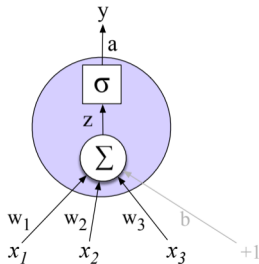


Figura: Um modelo matemático para o funcionamento do neurônio

$$y = \sigma(w_1 \cdot x_1 + w_2 \cdot x_2 + w_3 \cdot x_3 + b)$$



Um modelo matemático para o neurônio

Redes Neurais
Artificiais

Aplicando RNA para PLN

Arquitetura
Codificador-
Decodificador
(Encoder-
Decoder)

Mecanismo de
Atenção

Arquitetura
Transformer

6

Podemos então compreender um neurônio como uma função $f : \mathbb{R}^n \mapsto \mathbb{R}$ tal que $y = f(x_1, \dots, x_n) = f_1(\sum_i^n w_i x_i + b)$, com $f_1 : \mathbb{R} \rightarrow \mathbb{R}$.

26



Um modelo matemático para o neurônio

Redes Neurais Artificiais

Aplicando RNA para PLN

Arquitetura
Codificador-
Decodificador
(Encoder-
Decoder)

Mecanismo de
Atenção

Arquitetura
Transformer

6

Podemos então compreender um neurônio como uma função $f : \mathbb{R}^n \mapsto \mathbb{R}$ tal que $y = f(x_1, \dots, x_n) = f_1(\sum_i^n w_i x_i + b)$, com $f_1 : \mathbb{R} \rightarrow \mathbb{R}$. f_1 é chamada então de função de ativação do neurônio e descreve a condição de disparada (o processamento realizado) por aquele neurônio.



Um modelo matemático para o neurônio

Redes Neurais Artificiais

Aplicando RNA para PLN

Arquitetura
Codificador-
Decodificador
(Encoder-
Decoder)

Mecanismo de
Atenção

Arquitetura
Transformer

6

Podemos então compreender um neurônio como uma função $f : \mathbb{R}^n \mapsto \mathbb{R}$ tal que $y = f(x_1, \dots, x_n) = f_1(\sum_i^n w_i x_i + b)$, com $f_1 : \mathbb{R} \rightarrow \mathbb{R}$.

f_1 é chamada então de função de ativação do neurônio e descreve a condição de disparada (o processamento realizado) por aquele neurônio.

Podemos também descrever o termo $z = \sum_i^n w_i x_i + b$ de forma vetorial como $z = w \cdot x + b$, com $w = \langle w_1, \dots, w_n \rangle$, $x = \langle x_1, \dots, x_n \rangle^T$.



Funções de ativação

Redes Neurais Artificiais

Aplicando RNA para PLN

Arquitetura
Codificador-
Decodificador
(Encoder-
Decoder)

Mecanismo de
Atenção

Arquitetura
Transformer

7

As funções de ativação descrevem o processo de funcionamento de um neurônio, numa rede neural, ajudam a definir o viés de aprendizado das mesmas.

26



Funções de ativação

Redes Neurais Artificiais

Aplicando RNA para PLN

Arquitetura
Codificador-
Decodificador
(Encoder-
Decoder)

Mecanismo de
Atenção

Arquitetura
Transformer

7

As funções de ativação descrevem o processo de funcionamento de um neurônio, numa rede neural, ajudam a definir o viés de aprendizado das mesmas. As propriedades dessas funções vão definir aspectos computacionais importantes, como a expressividade da rede e sua complexidade computacional.

26



Funções de ativação

Redes Neurais Artificiais

Aplicando RNA para PLN

Arquitetura
Codificador-
Decodificador
(Encoder-
Decoder)

Mecanismo de
Atenção

Arquitetura
Transformer

7

As funções de ativação descrevem o processo de funcionamento de um neurônio, numa rede neural, ajudam a definir o viés de aprendizado das mesmas. As propriedades dessas funções vão definir aspectos computacionais importantes, como a expressividade da rede e sua complexidade computacional. É necessário que a função de ativação escolhida não seja linear! (por quê?)

26



Funções de ativação

Redes Neurais
Artificiais

7

Aplicando RNA para PLN

Arquitetura
Codificador-
Decodificador
(Encoder-
Decoder)

Mecanismo de
Atenção

Arquitetura
Transformer

As funções de ativação descrevem o processo de funcionamento de um neurônio, numa rede neural, ajudam a definir o viés de aprendizado das mesmas. As propriedades dessas funções vão definir aspectos computacionais importantes, como a expressividade da rede e sua complexidade computacional.

É necessário que a função de ativação escolhida não seja linear! (por quê?)

Algumas funções comuns:

- ▶ Função sigmóide (σ)
- ▶ Função tangente hiperbólico (\tanh)
- ▶ Função Linear Unitária Retificada ($ReLU$)

26



Funções de ativação: sigmóide

Redes Neurais
Artificiais

Aplicando RNA para PLN

Arquitetura
Codificador-
Decodificador
(Encoder-
Decoder)

Mecanismo de
Atenção

Arquitetura
Transformer

$$y = \sigma(z) = \frac{1}{1 + e^{-z}}$$



Funções de ativação: sigmóide

$$y = \sigma(z) = \frac{1}{1 + e^{-z}}$$

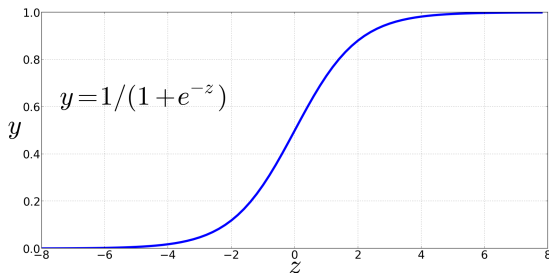


Figura: Gráfico da função sigmóide



Funções de ativação: tangente hiperbólico

Redes Neurais
Artificiais

Aplicando RNA para PLN

Arquitetura
Codificador-
Decodificador
(Encoder-
Decoder)

Mecanismo de
Atenção

Arquitetura
Transformer

$$y = \tanh(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$



Funções de ativação: tangente hiperbólico

$$y = \tanh(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

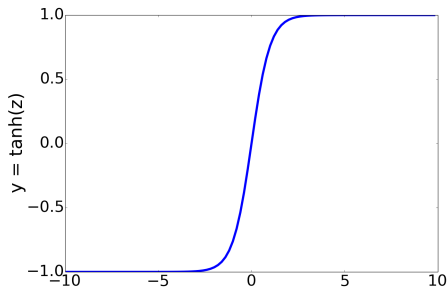


Figura: Gráfico da função tangente hiperbólico



Funções de ativação: ReLU

Redes Neurais
Artificiais

10

Aplicando RNA para PLN

Arquitetura
Codificador-
Decodificador
(Encoder-
Decoder)

Mecanismo de
Atenção

Arquitetura
Transformer

26

$$y = \text{ReLU}(z) = \max(z, 0)$$



Funções de ativação: ReLU

$$y = \text{ReLU}(z) = \max(z, 0)$$

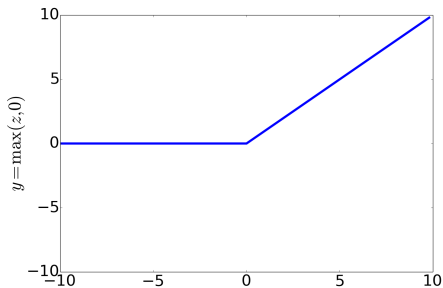


Figura: Gráfico da função ReLU



Redes neurais artificiais

Redes Neurais Artificiais

Aplicando RNA para PLN

Arquitetura
Codificador-
Decodificador
(Encoder-
Decoder)

Mecanismo de
Atenção

Arquitetura
Transformer

11

Uma rede neural artificial é um grafo direcionado e rotulado em cada cada nó representa um neurônio (descrito por uma função de ativação) e cada aresta é rotulada com um peso de associação.

26



Redes neurais artificiais

Redes Neurais Artificiais

Aplicando RNA para PLN

Arquitetura
Codificador-
Decodificador
(Encoder-
Decoder)

Mecanismo de
Atenção

Arquitetura
Transformer

11

Uma rede neural artificial é um grafo direcionado e rotulado em cada cada nó representa um neurônio (descrito por uma função de ativação) e cada aresta é rotulada com um peso de associação.

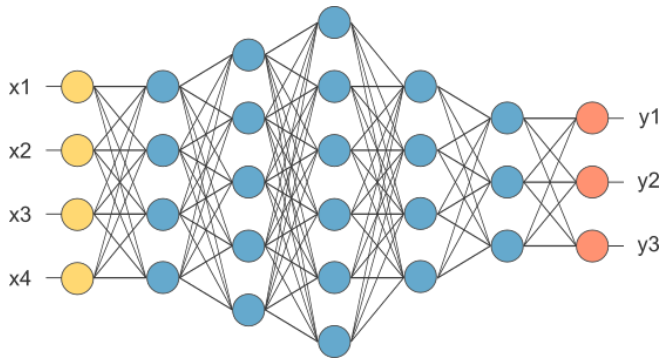


Figura: Representação de uma RNA

26



Redes Feedforward

Redes Neurais Artificiais

Aplicando RNA para PLN

Arquitetura
Codificador-
Decodificador
(Encoder-
Decoder)

Mecanismo de
Atenção

Arquitetura
Transformer

12

Redes *feedforward*, redes densas, perceptron multicamadas, ou redes lineares, são RNA acíclicas com múltiplas camadas completamente conectadas.

26



Redes Feedforward

Redes Neurais
Artificiais

Aplicando RNA para PLN

Arquitetura
Codificador-
Decodificador
(Encoder-
Decoder)

Mecanismo de
Atenção

Arquitetura
Transformer

12

Redes *feedforward*, redes densas, perceptron multicamadas, ou redes lineares, são RNA acíclicas com múltiplas camadas completamente conectadas.

São o tipo mais básico de rede neural e, comumente, uma camada densa faz parte da maioria das arquiteturas empregadas na prática.

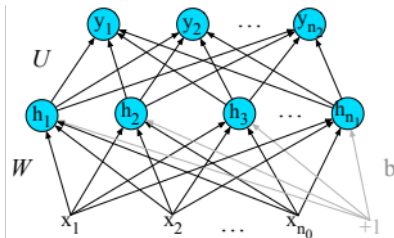


Figura: Rede Feedforward

26



Regressão Multinomial

Redes Neurais Artificiais

Aplicando RNA para PLN

Arquitetura
Codificador-
Decodificador
(Encoder-
Decoder)

Mecanismo de
Atenção

Arquitetura
Transformer

13

Podemos implementar uma regressão multinomial usando redes neurais empregando funções softmax

26



Regressão Multinomial

Redes Neurais
Artificiais

Aplicando RNA para PLN

Arquitetura
Codificador-
Decodificador
(Encoder-
Decoder)

Mecanismo de
Atenção

Arquitetura
Transformer

13

Podemos implementar uma regressão multinomial usando redes neurais empregando funções softmax

$$y = \text{softmax}(Wx + b)$$

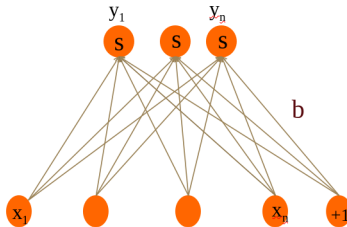


Figura: Camada Softmax

26



Regressão Multinomial

Redes Neurais
Artificiais

Aplicando RNA para PLN

Arquitetura
Codificador-
Decodificador
(Encoder-
Decoder)

Mecanismo de
Atenção

Arquitetura
Transformer

13

Podemos implementar uma regressão multinomial usando redes neurais empregando funções softmax

$$y = \text{softmax}(Wx + b)$$

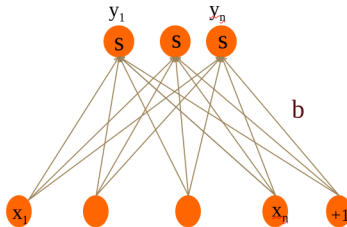


Figura: Camada Softmax

$$\text{softmax}(z) = \left[\frac{\exp(z_1)}{\sum_{i=1}^k \exp(z_i)}, \frac{\exp(z_2)}{\sum_{i=1}^k \exp(z_i)}, \dots, \frac{\exp(z_k)}{\sum_{i=1}^k \exp(z_i)} \right]$$

26



Redes Neurais
Artificiais

Aplicando RNA para PLN

Arquitetura
Codificador-
Decodificador
(Encoder-
Decoder)

Mecanismo de
Atenção

Arquitetura
Transformer

14

26

Aplicando RNA para PLN



Classificação de texto: Análise de sentimentos

Redes Neurais
Artificiais

Aplicando RNA para PLN

Arquitetura
Codificador-
Decodificador
(Encoder-
Decoder)

15

Mecanismo de
Atenção

Arquitetura
Transformer

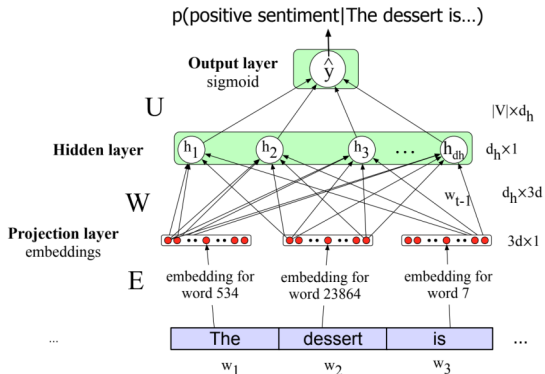


Figura: Uma rede neural para análise de sentimentos

26



Classificação de texto: Análise de sentimentos

Redes Neurais
Artificiais

Aplicando RNA para PLN

Arquitetura
Codificador-
Decodificador
(Encoder-
Decoder)

Mecanismo de
Atenção

Arquitetura
Transformer

16

```
class AS(nn.Module):
    def __init__(self):
        self.embeddings = nn.Embeddings(size_vocab,size_repr)
        self.hidden = nn.Linear(size_repr*number_words, number_neurons)
        self.hidden_act = nn.ReLU()
        self.out = nn.Linear(number_neurons,1)
        self.out_act = nn.Sigmoid()
    def forward(self,data)
        #data.shape == (1,n_words)
        embeddings = self.embeddings(data) #embeddings.shape ==(1,n_words,s
        embeddings = torch.cat(embeddings) #embeddings.shape ==(1,n_words*s
        z = self.hidden (embeddings) #z.shape == (1,number_neurons)
        z = self.hidden_act(z) #z.shape == (1,number_neurons)
        y = self.out # y.shape ==(1,1)
        y = self.out_act(y) #y.shape == (1,1)
        return y
```

26



Um modelo de linguagem neural

Redes Neurais
Artificiais

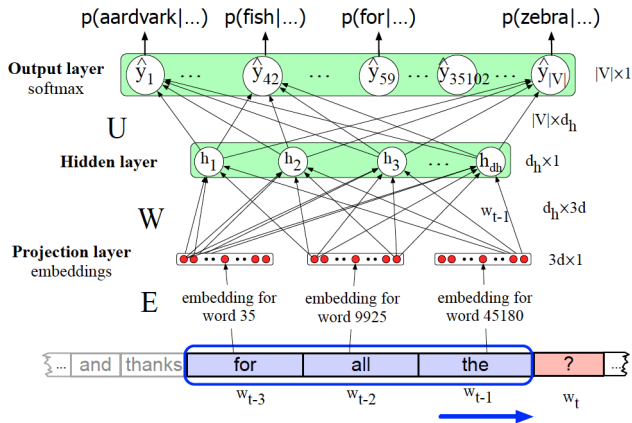
Aplicando RNA para PLN

Arquitetura
Codificador-
Decodificador
(Encoder-
Decoder)

Mecanismo de
Atenção

Arquitetura
Transformer

17



26



Vantagens sobre os modelos de n-gramas

Redes Neurais
Artificiais

Aplicando RNA para PLN

Arquitetura
Codificador-
Decodificador
(Encoder-
Decoder)

Mecanismo de
Atenção

Arquitetura
Transformer

18

Note que como o modelo aprende a partir das representações vetoriais de palavras (embeddings), não das sequências de palavras em si, o modelo pode explorar as similaridades entre palavras codificadas nelas para realizar suas predições.

26



Vantagens sobre os modelos de n-gramas

Redes Neurais
Artificiais

Aplicando RNA para PLN

Arquitetura
Codificador-
Decodificador
(Encoder-
Decoder)

Mecanismo de
Atenção

Arquitetura
Transformer

18

Note que como o modelo aprende a partir das representações vetoriais de palavras (embeddings), não das sequências de palavras em si, o modelo pode explorar as similaridades entre palavras codificadas nelas para realizar suas predições.

Assim, mesmo que uma sequência não tenha sido vista explicitamente nos dados de treino, o modelo pode avaliar sua probabilidade com base nessas similaridades.

26



Vantagens sobre os modelos de n-gramas

Redes Neurais
Artificiais

Aplicando RNA para PLN

Arquitetura
Codificador-
Decodificador
(Encoder-
Decoder)

Mecanismo de
Atenção

Arquitetura
Transformer

18

Note que como o modelo aprende a partir das representações vetoriais de palavras (embeddings), não das sequências de palavras em si, o modelo pode explorar as similaridades entre palavras codificadas nelas para realizar suas predições.

Assim, mesmo que uma sequência não tenha sido vista explicitamente nos dados de treino, o modelo pode avaliar sua probabilidade com base nessas similaridades.

- ▶ Tenho que lembrar de alimentar o gato hoje a noite (visto nos dados)
- ▶ Tenho que lembrar de alimentar o cachorro (não visto)

26



Redes Neurais
Artificiais

Aplicando RNA para PLN

Arquitetura
Codificador-
Decodificador
(Encoder-
Decoder)

Mecanismo de
Atenção

Arquitetura
Transformer

19

26

Redes Recorrentes e Sequências



Redes Neurais Recorrentes

Redes Neurais
Artificiais

Aplicando RNA para PLN

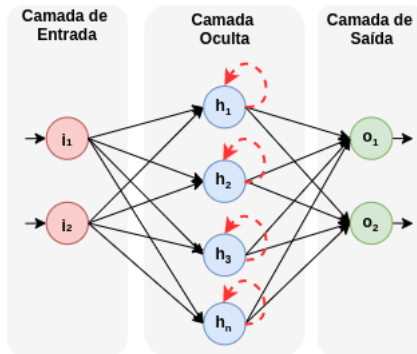
Arquitetura
Codificador-
Decodificador
(Encoder-
Decoder)

Mecanismo de
Atenção

Arquitetura
Transformer

20

Redes Recorrentes são aqueles em que acontecem ciclos na rede, i.e. em que a saída de um neurônio (ou camada) retroalimenta ele(a) mesmo(a).



Rede Neural Recorrente

26



Redes Neurais Recorrentes

Redes Neurais
Artificiais

Aplicando RNA para PLN

Arquitetura
Codificador-
Decodificador
(Encoder-
Decoder)

Mecanismo de
Atenção

Arquitetura
Transformer

os enlaces recorrentes numa RNN funcionam como um mecanismo de memória, que mantém uma noção de estado na rede.

21

26



Redes Neurais Recorrentes

Redes Neurais
Artificiais

Aplicando RNA para PLN

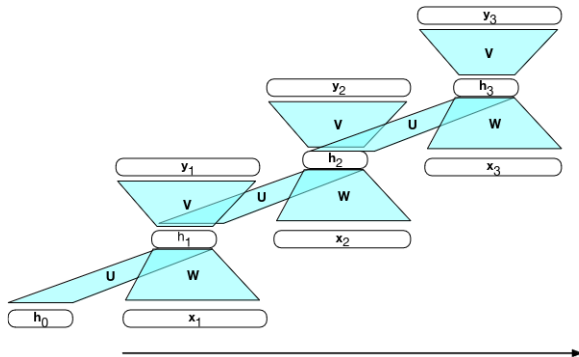
Arquitetura
Codificador-
Decodificador
(Encoder-
Decoder)

Mecanismo de
Atenção

Arquitetura
Transformer

21

os enlaces recorrentes numa RNN funcionam como um mecanismo de memória, que mantém uma noção de estado na rede. A execução de uma RNN é implementada através da expansão iterada da rede a partir da saída da execução anterior.



26



Redes Neurais
Artificiais

Aplicando RNA para PLN

Arquitetura
Codificador-
Decodificador
(Encoder-
Decoder)

Mecanismo de
Atenção

Arquitetura
Transformer

22

26

treinamento de uma RNA



Treinando uma rede neural

Redes Neurais
Artificiais

Aplicando RNA para PLN

Arquitetura
Codificador-
Decodificador
(Encoder-
Decoder)

Mecanismo de
Atenção

Arquitetura
Transformer

O treinamento de uma rede neural se dá num processo iterativo sobre os dados de treinamento em duas etapas:

- ▶ Passo de programação direta (forward propagation)
- ▶ Passo de propagação retrograda (backward propagation)

23

26



Treinando uma rede neural

Redes Neurais
Artificiais

Aplicando RNA para PLN

Arquitetura
Codificador-
Decodificador
(Encoder-
Decoder)

Mecanismo de
Atenção

Arquitetura
Transformer

O treinamento de uma rede neural se dá num processo iterativo sobre os dados de treinamento em duas etapas:

- ▶ Passo de programação direta (forward propagation)
- ▶ Passo de propagação retrograda (backward propagation)

23 É necessário definir alguns ingredientes da metodologia de treino:

- ▶ Quantidade de épocas de treino, i.e. quantas vezes a rede processará o conjunto de treino para aprender
- ▶ A função de perda (*loss*), que representa a computação do erro cometido pela rede em seu processamento



Treinando uma rede neural: grafo de computação

Redes Neurais
Artificiais

Aplicando RNA para PLN

Arquitetura
Codificador-
Decodificador
(Encoder-
Decoder)

Mecanismo de
Atenção

Arquitetura
Transformer

Um grafo de computação determina uma ordenação topológica das operações necessárias a serem realizadas para se obter um determinado resultado.

24

26



Treinando uma rede neural: grafo de computação

Redes Neurais
Artificiais

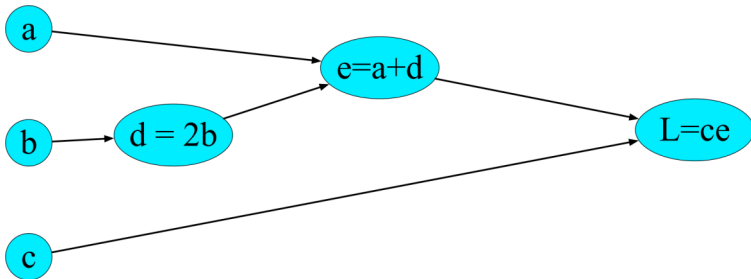
Aplicando RNA para PLN

Arquitetura
Codificador-
Decodificador
(Encoder-
Decoder)

Mecanismo de
Atenção

Arquitetura
Transformer

Um grafo de computação determina uma ordenação topológica das operações necessárias a serem realizadas para se obter um determinado resultado.



24

26



Treinando uma rede neural: propagação direta

Redes Neurais
Artificiais

Aplicando RNA para PLN

Arquitetura
Codificador-
Decodificador
(Encoder-
Decoder)

Mecanismo de
Atenção

Arquitetura
Transformer

Na propagação direta, as ações do grafo de computação são executadas em ordem, i.e. uma operação é executada somente quando todas as operações das quais ela depende já o foram.

25

26



Treinando uma rede neural: propagação direta

Redes Neurais
Artificiais

Aplicando RNA para PLN

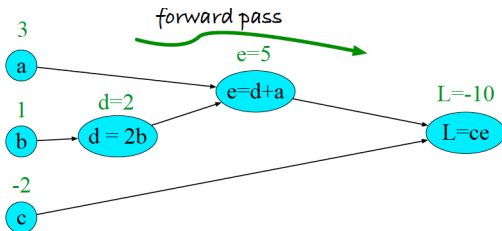
Arquitetura
Codificador-
Decodificador
(Encoder-
Decoder)

Mecanismo de
Atenção

Arquitetura
Transformer

25

Na propagação direta, as ações do grafo de computação são executadas em ordem, i.e. uma operação é executada somente quando todas as operações das quais ela depende já o foram.



26



Treinando uma rede neural: cálculo do erro

Redes Neurais
Artificiais

Aplicando RNA para PLN

Arquitetura
Codificador-
Decodificador
(Encoder-
Decoder)

Mecanismo de
Atenção

Arquitetura
Transformer

Ao final da execução, a função de perda $L(y, \hat{y})$ é aplicada sobre o resultado encontrado e o resultado conhecido (rótulo).

26

26



Treinando uma rede neural: cálculo do erro

Redes Neurais
Artificiais

Aplicando RNA para PLN

Arquitetura
Codificador-
Decodificador
(Encoder-
Decoder)

Mecanismo de
Atenção

Arquitetura
Transformer

Ao final da execução, a função de perda $L(y, \hat{y})$ é aplicada sobre o resultado encontrado e o resultado conhecido (rótulo).

Para ajustar os pesos da rede neural, precisamos entender a contribuição de cada componente da rede ao erro obtido. Fazemos isso computando a derivada parcial de cada peso para a perda, i.e.

$$\frac{\partial L(y, \hat{y})}{\partial w_i}$$

26

26



Treinando uma rede neural: cálculo do erro

Redes Neurais
Artificiais

Aplicando RNA para PLN

Arquitetura
Codificador-
Decodificador
(Encoder-
Decoder)

Mecanismo de
Atenção

Arquitetura
Transformer

Ao final da execução, a função de perda $L(y, \hat{y})$ é aplicada sobre o resultado encontrado e o resultado conhecido (rótulo).

Para ajustar os pesos da rede neural, precisamos entender a contribuição de cada componente da rede ao erro obtido. Fazemos isso computando a derivada parcial de cada peso para a perda, i.e.

$$\frac{\partial L(y, \hat{y})}{\partial w_i}$$

Tal valor será computado de forma numérica e propagado para trás usando a regra da cadeia $\frac{df}{dx} = \frac{df}{du} \cdot \frac{du}{dx}$

26

26



Treinando uma rede neural: propagação retrograda

Redes Neurais
Artificiais

Aplicando RNA para PLN

Arquitetura
Codificador-
Decodificador
(Encoder-
Decoder)

Mecanismo de
Atenção

Arquitetura
Transformer

Na propagação retrograda, calculamos as derivadas parciais do erro em relação aos pesos e atualizamos os pesos da rede, propagando para trás as mudanças.

27

26



Treinando uma rede neural: propagação retrograda

Redes Neurais
Artificiais

Aplicando RNA para PLN

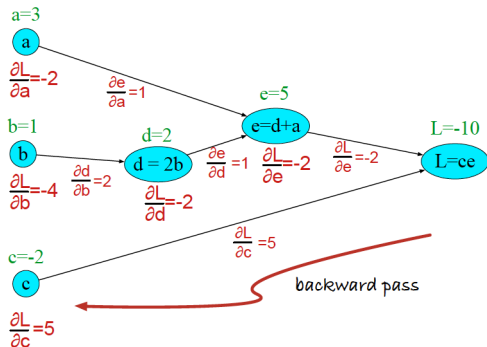
Arquitetura
Codificador-
Decodificador
(Encoder-
Decoder)

Mecanismo de
Atenção

Arquitetura
Transformer

27

Na propagação retrograda, calculamos as derivadas parciais do erro em relação aos pesos e atualizamos os pesos da rede, propagando para trás as mudanças.



26



Redes Neurais
Artificiais

Aplicando RNA para PLN

Arquitetura
Codificador-
Decodificador
(Encoder-
Decoder)

Mecanismo de
Atenção

Arquitetura
Transformer

28

26

Aplicação: Geração de texto

Vamos para o Colab!