

(6) Redes neurais para dados sequenciais

Redes Neurais e Arquiteturas Profundas

Moacir A. Ponti
CeMEAI/ICMC, Universidade de São Paulo
MBA em Ciência de Dados

www.icmc.usp.br/~moacir — moacir@icmc.usp.br

São Carlos-SP/Brasil – 2020

Agenda

Dados sequenciais: recorrência

- Camada recorrente básica (RNN)

- Long Short Term Memory (LSTM)

- Gated Recurrent Unit (GRU)

Word2Vec: representações para texto

Sequence-to-Sequence e Mecanismo de Atenção

- Transformer Network

Agenda

Dados sequenciais: recorrência

- Camada recorrente básica (RNN)

- Long Short Term Memory (LSTM)

- Gated Recurrent Unit (GRU)

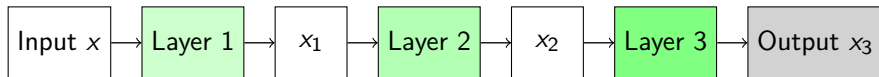
Word2Vec: representações para texto

Sequence-to-Sequence e Mecanismo de Atenção

- Transformer Network

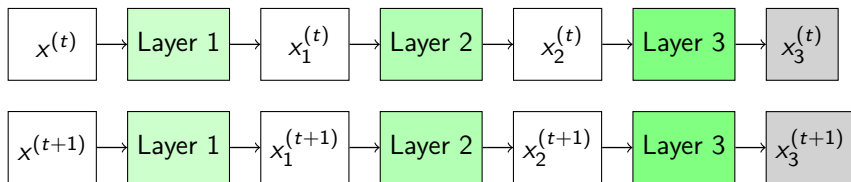
Para dados não sequenciais

- ▶ Camadas densas e convolucionais consideram apenas o exemplo atual para computar a saída
- ▶ Em cada iteração, cada entrada vai passando pelas camadas até atingir a saída



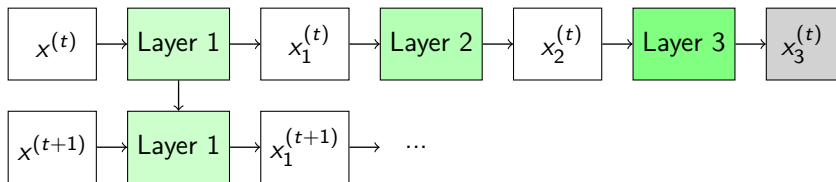
Para dados não sequenciais

- ▶ Na iteração $t + 1$ usamos os dados de $t + 1$ para adaptar os parâmetros



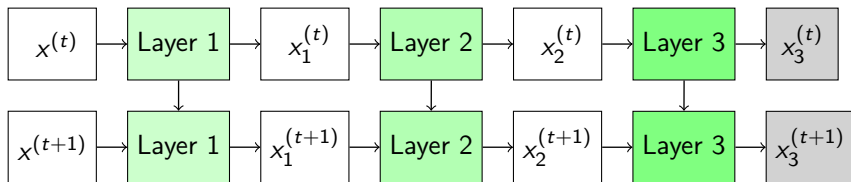
Para dados sequenciais

- Se a iteração $t + 1$ depende da anterior t , usamos a saída de cada camada para alimentar a camada na entrada da iteração $t + 1$



Para dados sequenciais

- Dessa forma, a saída (após a primeira), dependerá não apenas da entrada atual, mas das saídas computadas anteriormente para cada unidade

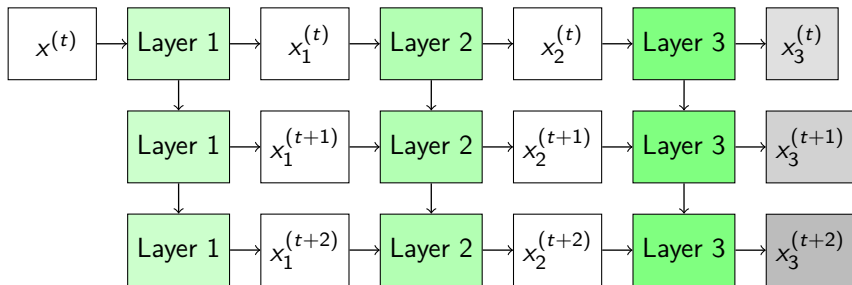


Configurações de sequências

- ▶ Uma entrada, saída sequencial
- ▶ Entrada sequencial, uma saída
- ▶ Entrada sequencial, saída sequencial

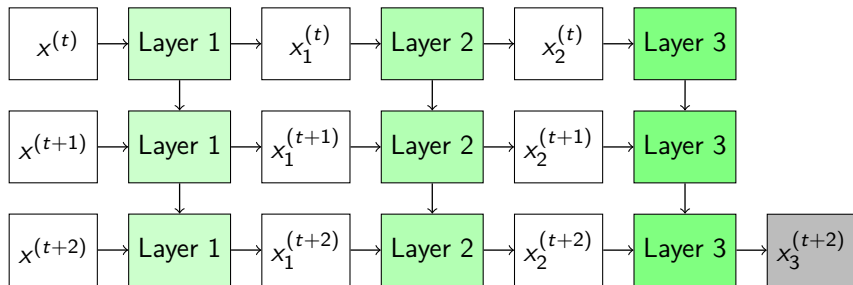
Configurações de seqüências

- **Uma entrada, saída seqüencial:** e.g. um áudio ou imagem é dado como entrada e a rede produz uma seqüência de palavras que os descrevem



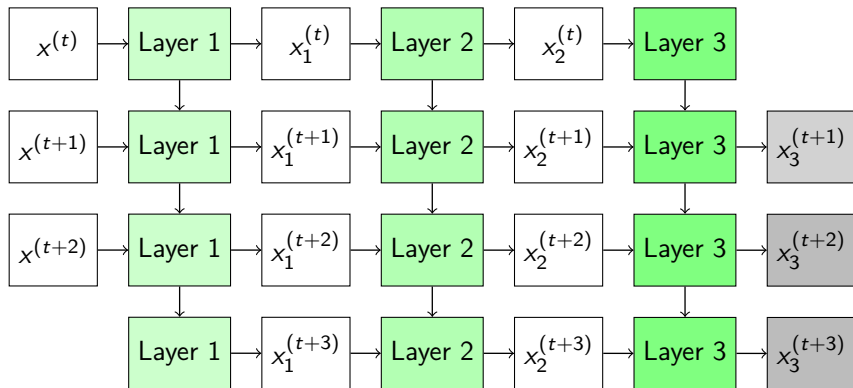
Configurações de seqüências

- **Entrada sequencial, uma saída** e.g. um texto é dado como entrada e a saída é sua análise de sentimentos: conteúdo positivo ou negativo.



Configurações de seqüências

- **Entrada sequencial, saída sequencial** e.g. tradução automática de sentenças entre diferentes idiomas (com atraso k)



Agenda

Dados sequenciais: recorrência

Camada recorrente básica (RNN)

Long Short Term Memory (LSTM)

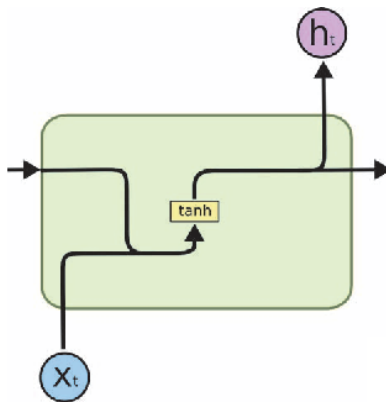
Gated Recurrent Unit (GRU)

Word2Vec: representações para texto

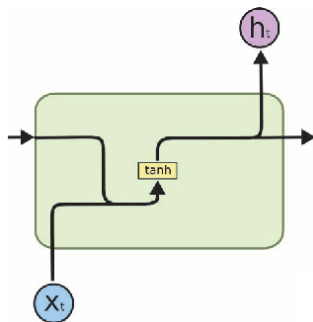
Sequence-to-Sequence e Mecanismo de Atenção

Transformer Network

Aprende um tipo de "memória"



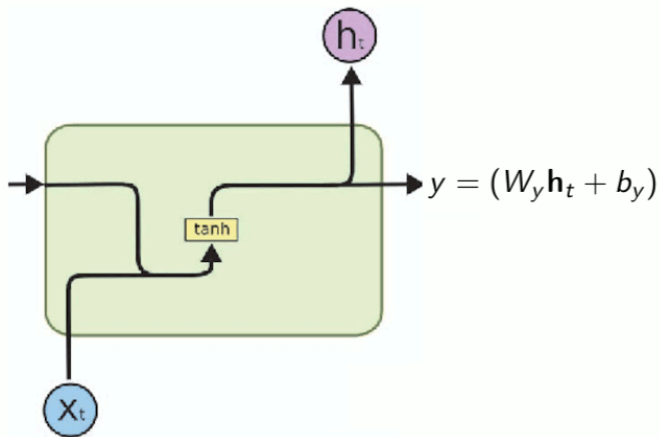
Componentes são combinações lineares



$$h_t = \tanh(W_h h_{t-1} + W_x x_t + b_h)$$
$$y = (W_y h_t + b_y)$$

Saída recorrente (sumário) e saída da rede

$$\mathbf{h}_t = \tanh(W_h \mathbf{h}_{t-1} + W_x \mathbf{x}_t + b_h)$$



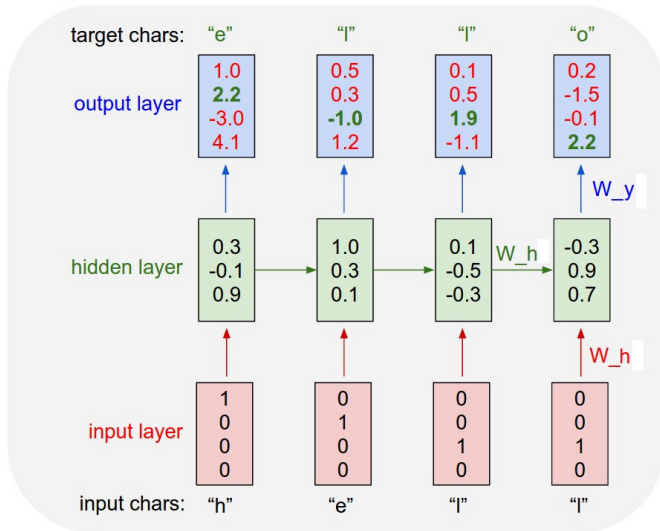
Exemplo: prever próximo caractere

Definimos uma codificação one-hot para os caracteres:

- ▶ $h = [1, 0, 0, 0]$
- ▶ $e = [0, 1, 0, 0]$
- ▶ $l = [0, 0, 1, 0]$
- ▶ $o = [0, 0, 0, 1]$

<http://karpathy.github.io/2015/05/21/rnn-effectiveness/>

Exemplo: prever próximo caractere



Agenda

Dados sequenciais: recorrência

- Camada recorrente básica (RNN)

- Long Short Term Memory (LSTM)

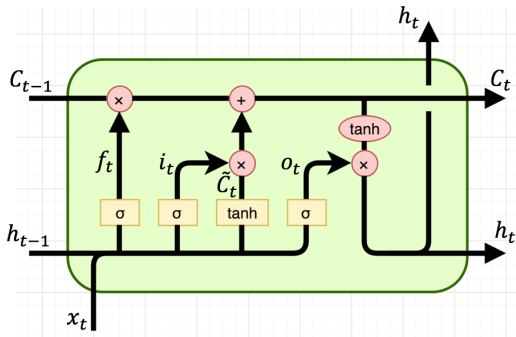
- Gated Recurrent Unit (GRU)

Word2Vec: representações para texto

Sequence-to-Sequence e Mecanismo de Atenção

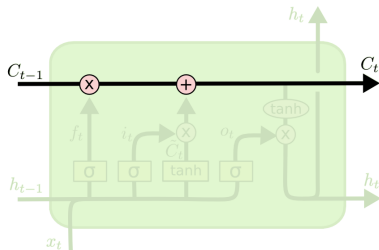
- Transformer Network

Long Short Term Memory Unit (LSTM)



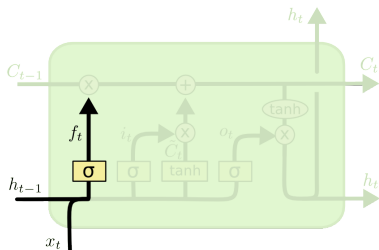
Adaptados de <http://colah.github.io/posts/2015-08-Understanding-LSTMs/>

LSTM network: Cell state



- ▶ Responsável pela memória longa da unidade, adiciona contribuições para além da iteração anterior.
- ▶ Esse estado pode ser modificado por 2 portões/**gates**

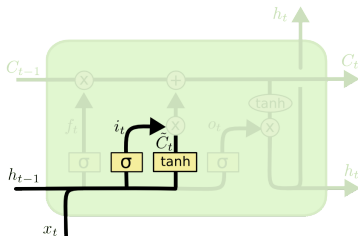
LSTM network: forget gate



$$f_t = \sigma(W_f \cdot [h_{t-1}, x_t] + b_f)$$

- ▶ decide o que cancelar de C com base no sumário anterior e a entrada atual
- ▶ saída entre 0 (esquecer) e 1 (manter totalmente) para cada dimensão de C

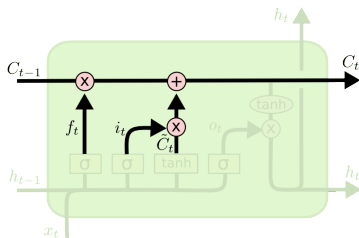
LSTM network: input gate



$$i_t = \sigma(W_i \cdot [h_{t-1}, x_t] + b_i)$$
$$\tilde{C}_t = \tanh(W_C \cdot [h_{t-1}, x_t] + b_C)$$

- ▶ primeiro, combina o sumário anterior h_{t-1} e a entrada x_t
- ▶ então, aprende um filtro \tilde{C}_t que indica quais partes devem ser mantidas na "memória longa", sendo depois somado a C_{t-1}

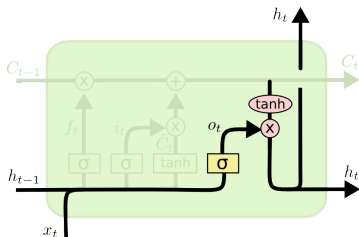
LSTM network: update Cell state



$$C_t = f_t * C_{t-1} + i_t * \tilde{C}_t$$

- ▶ agora temos uma combinação entre os estados atual e anterior
- ▶ acima o * significa uma multiplicação ponto-a-ponto

LSTM network: output gate



$$o_t = \sigma (W_o [h_{t-1}, x_t] + b_o)$$

$$h_t = o_t * \tanh (C_t)$$

- ▶ decide qual será o sumário, transformado a partir do anterior
- ▶ e ponderado de acordo com o estado de célula atual, C_t

Agenda

Dados sequenciais: recorrência

- Camada recorrente básica (RNN)

- Long Short Term Memory (LSTM)

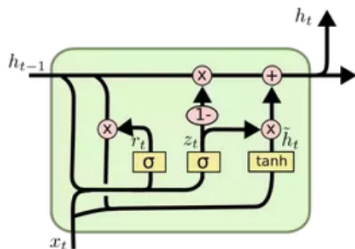
- Gated Recurrent Unit (GRU)

Word2Vec: representações para texto

Sequence-to-Sequence e Mecanismo de Atenção

- Transformer Network

Gated Recurrent Unit (GRU)



$$z_t = \sigma(W_z \cdot [h_{t-1}, x_t])$$

$$r_t = \sigma(W_r \cdot [h_{t-1}, x_t])$$

$$\tilde{h}_t = \tanh(W \cdot [r_t * h_{t-1}, x_t])$$

$$h_t = (1 - z_t) * h_{t-1} + z_t * \tilde{h}_t$$

- ▶ Não possui Cell State
- ▶ Reset gate r : filtra qual parte de h_{t-1} será utilizada para compor o novo sumário candidato em conjunto com x_t
- ▶ Update gate z : pondera partes do sumário anterior de forma complementar ao novo estado candidato
- ▶ \tilde{h}_t é o sumário "candidato"

GRU x LSTM

- ▶ Não há um consenso de "melhor" método
- ▶ GRU em muitos casos tem resultados similares à LSTM, com menos parâmetros
- ▶ Há uma versão recente, JANET, que simplificou ainda mais o modelo, removendo o "Reset gate"
- ▶ Temporal Convolutional Networks, que utilizam convoluções 1D para aprender posicionamento local, também se mostraram eficientes em alguns cenários

Agenda

Dados sequenciais: recorrência

Camada recorrente básica (RNN)

Long Short Term Memory (LSTM)

Gated Recurrent Unit (GRU)

Word2Vec: representações para texto

Sequence-to-Sequence e Mecanismo de Atenção

Transformer Network

- ▶ Representação (embedding) para palavras
- ▶ A função de custo para aprender essa representação:

$$p(w_{t+j}|w_t)$$

t é a palavra, $t + j$ são todas as palavras no contexto de t

- ▶ Otimiza em função de palavras que devem estar próximas se estiverem no mesmo contexto

- ▶ Representação (embedding) para palavras
- ▶ A função de custo para aprender essa representação:

$$p(w_{t+j}|w_t)$$

t é a palavra, $t + j$ são todas as palavras no contexto de t

- ▶ Otimiza em função de palavras que devem estar próximas se estiverem no mesmo contexto

Skip-grams (SG)

Predição de palavras em uma certa "janela" de proximidade m de uma palavra t

- Formulação "softmax":

$$\frac{\exp(u_o^T v_c)}{\sum_w^V \exp(u_w^T v_c)}$$

V é o total de palavras no vocabulário

u_o é a representação de uma palavra de "saída" (Que queremos prever)

v_c é a representação de uma palavra de entrada (central)

Word2Vec: skipgram

Dada uma representação one-hot de uma palavra $w_t \in R^V$, calculamos sua representação vetorial $v_c \in R^d$ (central)

$$W \cdot w_t = v_c$$

$$\begin{bmatrix} \dots & \dots & 0.1 & \dots & \dots \\ \dots & \dots & -0.3 & \dots & \dots \\ \dots & \dots & 1.4 & \dots & \dots \\ \dots & \dots & 0.2 & \dots & \dots \\ \dots & \dots & 0.5 & \dots & \dots \end{bmatrix} \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 1 \\ 0 \end{bmatrix} = \begin{bmatrix} 0.1 \\ -0.3 \\ 1.4 \\ 0.2 \\ 0.5 \end{bmatrix}$$

Word2Vec: skipgram

v_c é filtrada por representações u_o das palavras de saída (no contexto, que queremos prever) em diferentes posições $t - i$

$$u_o^T \cdot v_c$$

Para todas as palavras do vocabulário isso é codificado em uma matriz:

$$U_o \cdot v_c$$

$$\begin{bmatrix} 0.0 & 2.0 & 0.1 & 2.0 & 0.1 \\ 0.0 & 1.0 & 2.0 & -0.5 & 1.0 \\ \dots & \dots & \dots & \dots & \dots \\ \dots & \dots & \dots & \dots & \dots \\ \dots & \dots & \dots & \dots & \dots \\ \dots & \dots & \dots & \dots & \dots \end{bmatrix} \begin{bmatrix} 0.1 \\ -0.3 \\ 1.4 \\ 0.2 \\ 0.5 \end{bmatrix} = \text{softmax} \left(\begin{bmatrix} 0.0 \\ 2.9 \\ 0.1 \\ 1.4 \\ -0.5 \\ 0.0 \\ 0.0 \end{bmatrix} \right) = \begin{bmatrix} 0.04 \\ 0.67 \\ 0.04 \\ 0.15 \\ 0.02 \\ 0.04 \\ 0.04 \end{bmatrix}$$

Word2Vec: skipgram

- ▶ W aprende representações (nas colunas) para cada palavra quando são "centrais"
- ▶ U_o aprende representações (nas linhas) para cada palavra quando são "contexto"

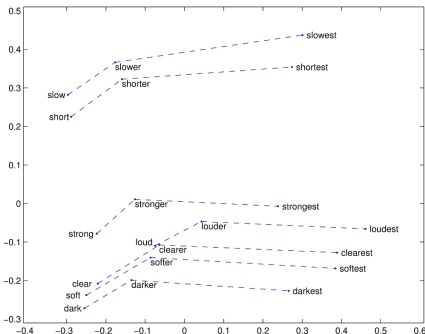
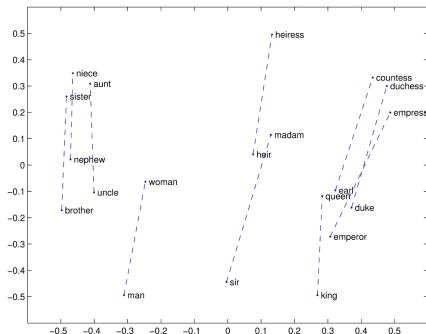
$$W = \begin{bmatrix} \dots & \dots & 0.1 & \dots & \dots \\ \dots & \dots & -0.3 & \dots & \dots \\ \dots & \dots & 1.4 & \dots & \dots \\ \dots & \dots & 0.2 & \dots & \dots \\ \dots & \dots & 0.5 & \dots & \dots \end{bmatrix} \quad U_o = \begin{bmatrix} \dots & \dots & \dots & \dots & \dots \\ 0.0 & 1.0 & 2.0 & -0.5 & 1.0 \\ \dots & \dots & \dots & \dots & \dots \\ \dots & \dots & \dots & \dots & \dots \\ \dots & \dots & \dots & \dots & \dots \\ \dots & \dots & \dots & \dots & \dots \\ \dots & \dots & \dots & \dots & \dots \end{bmatrix}$$

- ▶ Palavras que aparecem num mesmo contexto terão representações similares

Deixa o menino jogar
Deixa o moleque jogar
Deixa o piá jogar
Deixa seu filho jogar

Word2Vec: GloVe (Global Vectors for Word Representation)

<https://nlp.stanford.edu/projects/glove/>



Word Embeddings em Português - NILC/ICMC

<http://www.nilc.icmc.usp.br/embeddings>

Agenda

Dados sequenciais: recorrência

Camada recorrente básica (RNN)

Long Short Term Memory (LSTM)

Gated Recurrent Unit (GRU)

Word2Vec: representações para texto

Sequence-to-Sequence e Mecanismo de Atenção

Transformer Network

RNNs e Sequence-to-Sequence (seq2seq)

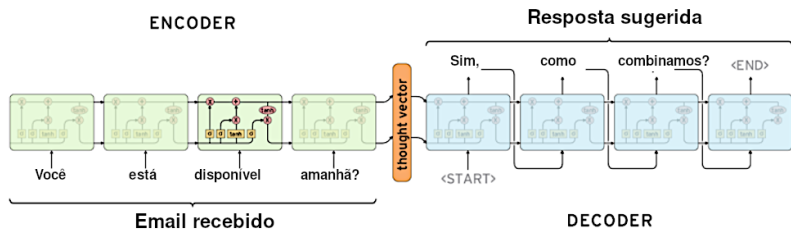


Figura adaptada de: Sachin Abeywardana

- (3 vídeos de Jay Alammar)

Mecanismo de atenção: intuição e motivação

- ▶ Encontrar qual parte de uma sequência é mais importante para prever uma certa saída
- ▶ Em unidades recorrentes, cada entrada perturba a memória prejudicando conhecimento de dados anteriores

Mecanismo de Atenção: imagens

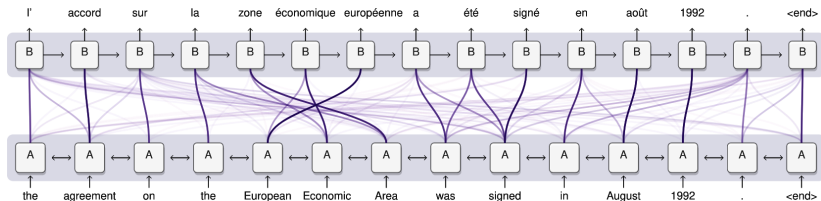


A man wearing a hat and
a hat on a skateboard.

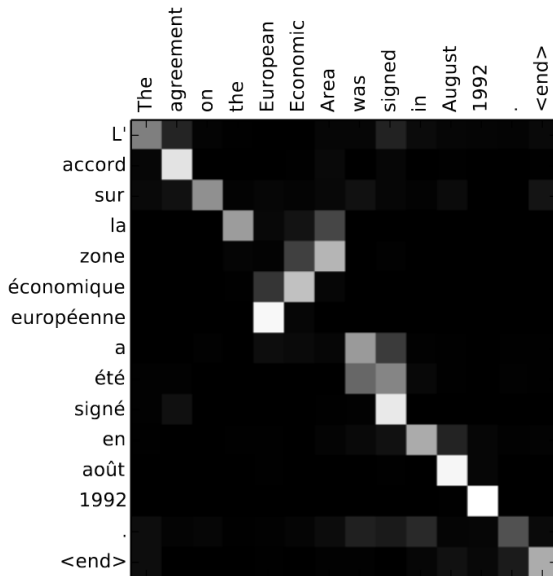


A man is talking on his cell phone
while another man watches.

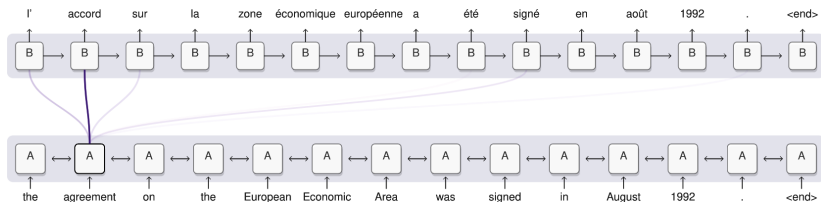
Mecanismo de Atenção: texto



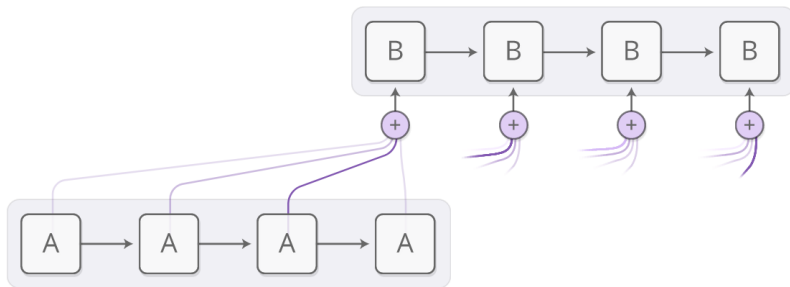
Mecanismo de Atenção: texto



Mecanismo de Atenção: texto



RNNs seq2seq e atenção global



Adaptado de Olah & Carter, "Attention and Augmented Recurrent Neural Networks", Distill, 2016.

<http://doi.org/10.23915/distill.00001>

► (vídeos de Jay Alammar)

Mecanismo de atenção: implementação básica

- Computar o alinhamento/similaridade entre o sumário atual do decoder, s_i , com sumários anteriores do encoder, h_j

Usa softmax para obter pesos na forma de probabilidades

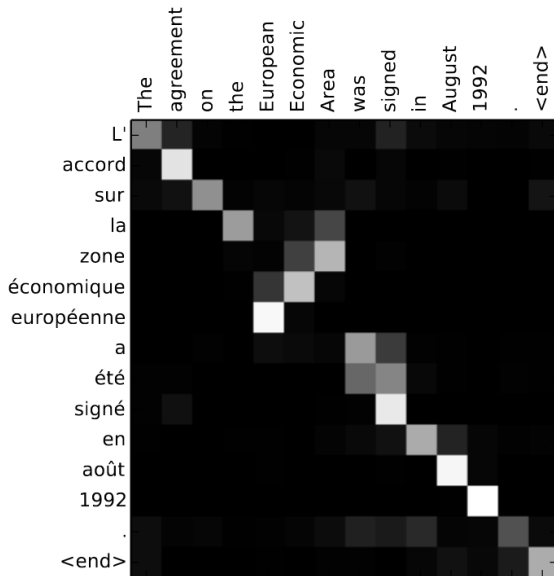
$$a_{i,j} = \frac{\exp(\text{alinhamento}(s_i, h_k))}{\sum_k \exp(\text{alinhamento}(s_i, h_k))},$$

"alinhamento" é um tipo de similaridade, e.g. produto interno:

$$\text{alinhamento}(s_i, h_k) = s_i^T h_j$$

- Atenção produz um vetor de "contexto" $c_i = \sum_j a_{i,j} h_j$ a ser usado para produzir a saída atual.

Mecanismo de Atenção: exemplo de vetores de alinhamento



Agenda

Dados sequenciais: recorrência

Camada recorrente básica (RNN)

Long Short Term Memory (LSTM)

Gated Recurrent Unit (GRU)

Word2Vec: representações para texto

Sequence-to-Sequence e Mecanismo de Atenção Transformer Network

RNNs vs Transformer Networks

RNNs

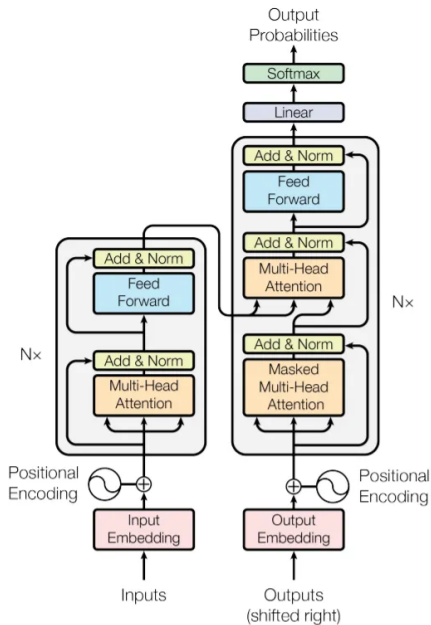
- ▶ podem não funcionar com dependências longas
- ▶ recorrência dificulta computação paralela
- ▶ podem sofrer com explosão ou desaparecimento de gradiente

Transformer Networks

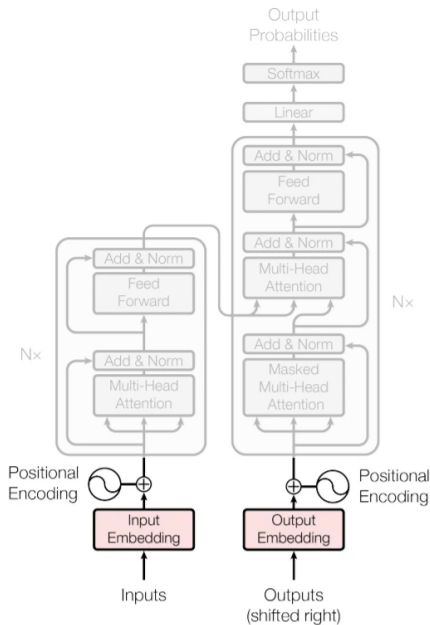
- ▶ não possui recorrência, apenas atenção
- ▶ facilita capturar dependências longas
- ▶ facilita processamento paralelo

Artigo: "Attention is all you need" Vaswani, NeurIPS 2017.

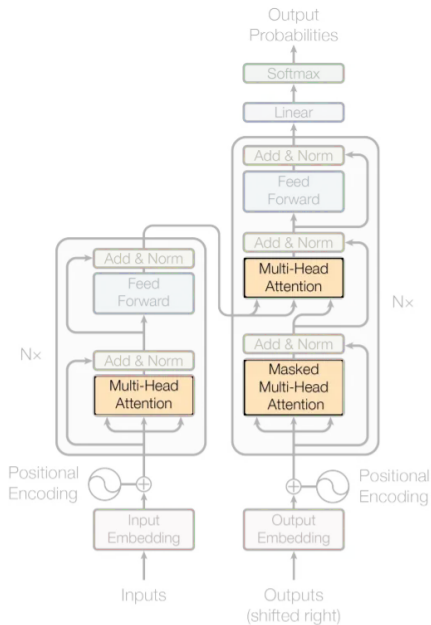
Arquitetura Transformer



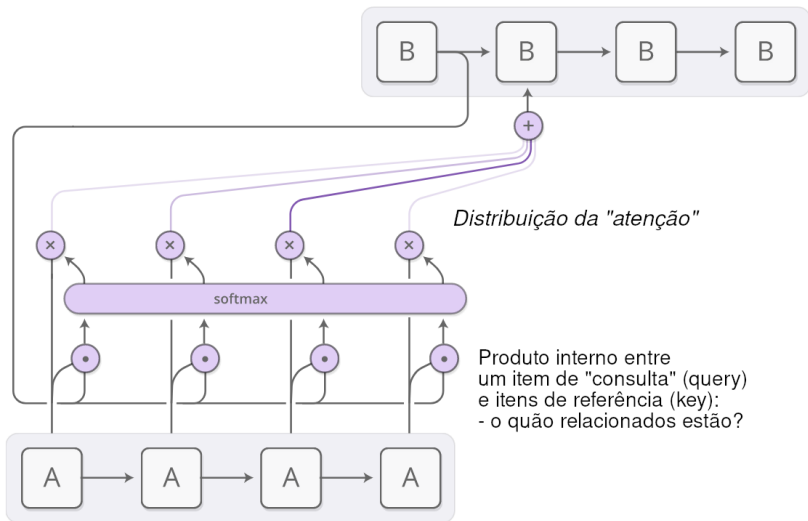
Positional Encoding: adiciona informação de posicionamento



Arquitetura Transformer: multi-atenção



Distribuição da atenção com relação ao estado atual



Mecanismo de atenção em Transformer Networks

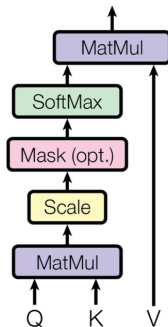
- ▶ Recuperar um valor v_i para uma consulta/query q baseada numa chave/key k_i

$$Attention(q, k, v) = \sum_i similarity(q, k_i) \times v_i$$

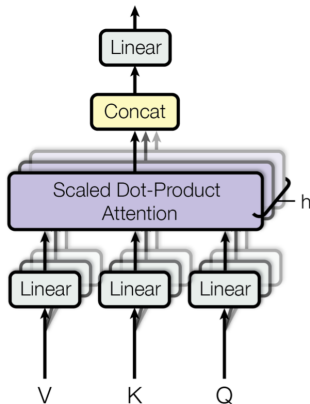
- ▶ A **similaridade** entre uma consulta e todas as chaves, ponderadas pelos valores
- ▶ Somar ao longo de todas as chaves/valores, produz uma **distribuição** de pesos relacionando consulta e todos os valores

Multi-head attention

Scaled Dot-Product Attention



Multi-Head Attention



Representações para cada palavra

1. Word embeddings: w
2. Word embeddings + Positional encoding: $w + p = e$
3. Query: $Q_h = W_Q \cdot e$
4. Key: $K_h = W_K \cdot e$
5. Value: $V_h = W_V \cdot e$

h indexa as H "cabeças de atenção"

Attention e Masked Attention

Attention

$$\text{attention}(Q, K, V) = \text{softmax} \left(\frac{Q^T K}{\sqrt{d_k}} \right) V$$

d_k dimensões da chave

Masked Attention

- ▶ o decoder deve anular atenção com relação a palavras de entrada futuras, senão não há aprendizado

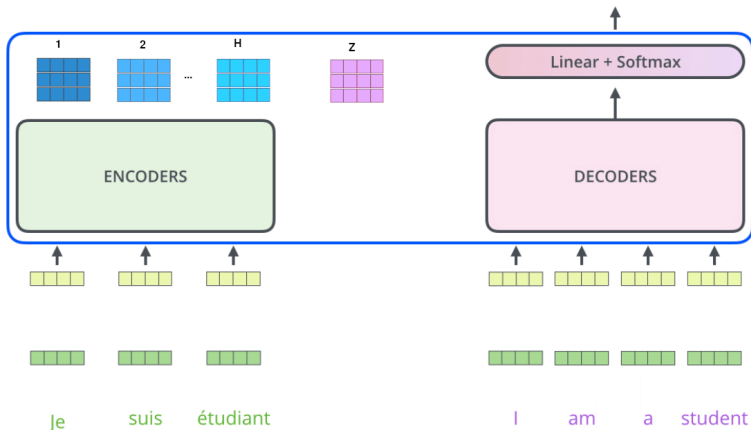
$$\text{maskedattention}(Q, K, V) = \text{softmax} \left(\frac{Q^T K + M}{\sqrt{d_k}} \right) V$$

M é uma matriz com $-\infty$ nas posições de palavras futuras

Representações para cada palavra

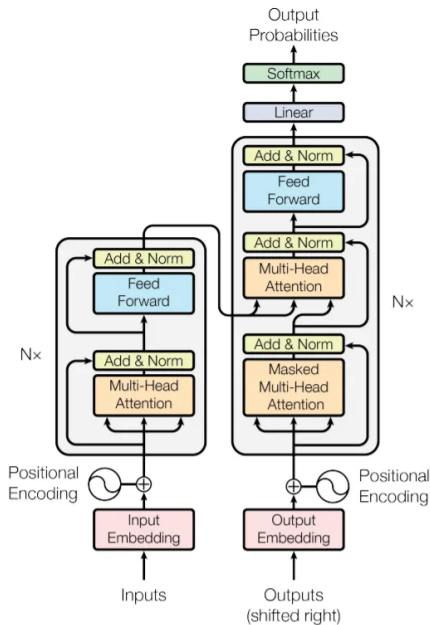
1. Word embeddings: w
2. com positional encoding: $w + p = e$
3. Query: $Q_h = W_Q \cdot e_h$
4. Key: $K_h = W_K \cdot e_h$
5. Value: $V_h = W_V \cdot e_h$
6. Attention heads: $Z_h = \text{softmax} \left(\frac{Q_h^T K_h + M}{\sqrt{d_k}} \right) V_h$
7. Attention output: $Z = W_O \text{concat}((Z_1, Z_2, \dots, Z_H))$

Recorrência via "Teacher forcing" no treinamento



Adaptado de Jay Alammar

Arquitetura Transformer



Arquitetura Transformer: alguns resultados

One day, I decided to bake a loaf of bread on my own oven.

<https://transformer.huggingface.co/>

Arquitetura Transformer: alguns resultados

One day, I decided to bake a loaf of bread on my own oven. It was so simple and simple and easy.

<https://transformer.huggingface.co/>

Arquitetura Transformer: alguns resultados

One day, I decided to bake a loaf of bread on my own oven. It was so simple and simple and easy. This bread is my first time baking it in my oven.

<https://transformer.huggingface.co/>

Arquitetura Transformer: alguns resultados

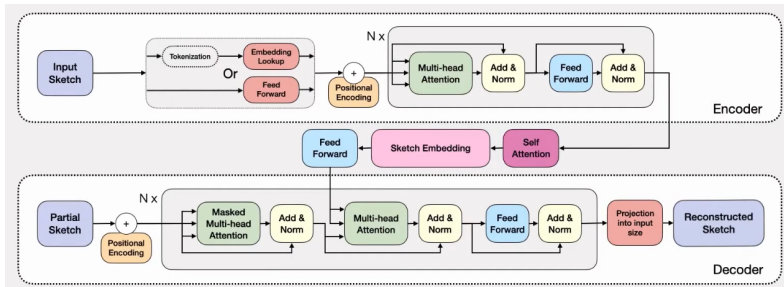
One day, I decided to bake a loaf of bread on my own oven. It was so simple and simple and easy. This bread is my first time baking it in my oven. I don't know what you will have to do to make this loaf.

<https://transformer.huggingface.co/>

Após Transformer

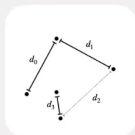
- ▶ GPT e GPT2
- ▶ BERT
- ▶ Lambda Networks
- ▶ ...

Outras aplicações: para dados visuais



Outras aplicações: para dados visuais

Continuous



$(dx_0, dy_0, 0)$

$(dx_1, dy_1, 0)$

$(dx_2, dy_2, 1)$

$(dx_1, dy_1, 0)$



Original Sample

Grid Tokens



SOS

03

01

SEP

07

04

EOS



n=1000

References

- ▶ Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, Illia Polosukhin. Attention is all you need, NeurIPS 2017
- ▶ Dzmitry Bahdanau, KyungHyun Cho and Yoshua Bengio. Neural Machine Translation By Jointly Learning To Align And Translate. ICLR 2015.
- ▶ A. Karpathy. Understanding LSTM Networks.
<http://karpathy.github.io/2015/05/21/rnn-effectiveness/>
- ▶ C. Olah. Understanding LSTM Networks
<http://colah.github.io/posts/2015-08-Understanding-LSTMs/>