

Deep Learning Model for NLP Task - July 5'th, 2019

Attention Model

申恒恒 - <https://github.com/rh01>

MyToc

- Attention
 - Attention 背景知识
 - Attention 模型
 - Attention 模型在自然语言处理和计算机视觉上的应用
- Transformer
 - Transformer 背景知识
 - Transformer 架构介绍
 - Transformer 模型在自然语言处理和计算机视觉上的应用
- Bert
 - Bert 背景知识
 - Bert 模型介绍
 - Bert 模型在自然语言处理上的应用
- XLNet
 - XLNet 背景知识
 - XLNet 模型介绍
 - XLNet 模型在自然语言处理上的应用
- Transformer XLNet
 - Transformer XLNet 背景知识
 - Transformer XLNet 模型介绍
 - Transformer XLNet 模型在自然语言处理上的应用

Attention 机制

Attention 背景知识

机器翻译解决的是输入是一串源语言表达的一句话，输出是目标语言相对应的话的问题，如将法语中的一段话翻译成合适的英语。

$x^{<1>} \quad x^{<2>} \quad x^{<3>} \quad x^{<4>} \quad x^{<5>}$
Jane visite l'Afrique en septembre
→ Jane is visiting Africa in September.
 $y^{<1>} \quad y^{<2>} \quad y^{<3>} \quad y^{<4>} \quad y^{<5>} \quad y^{<6>}$

图 1：机器翻译例子（法语 → 英语）

之前的机器翻译(简称NMT)模型中,通常的配置是Encoder-Decoder结构,即Encoder读取输入的句子并将其转换为定长的一个向量,然后Decoder再将这个向量翻译成对应的目标语言的文字。通常encoder及decoder均采用RNN结构,如LSTM或GRU等如下图所示,我们利用encoder RNN将输入语句信息总结到最后一个hidden vector中,并将其作为decoder初始的hidden vector,利用decoder解码成对应的其他语言中的文字。

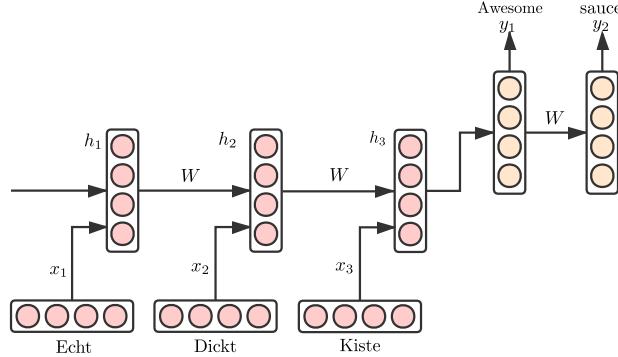


图 2: NMT 的 Encoder-Decoder 架构

但是这个结构有些问题,尤其是RNN机制实际中存在**长程梯度消失**的问题,对于较长的句子,我们**很难寄希望于将输入的序列转化为定长的向量而保存所有的有效信息**,所以随着所需翻译句子的长度的增加,这种结构的效果会显著下降。

以英语 - 法语翻译为例,给定一对英语输入序列“*They*”、“*are*”、“*watching*”、“.”和法语输出序列“*Ils*”、“*regardent*”、“.”。解码器可以在输出序列的时间步1使用更集中编码了“*They*”、“*are*”信息的背景变量来生成“*Ils*”,在时间步2使用更集中编码了“*watching*”信息的背景变量来生成“*regardent*”,在时间步3使用更集中编码了“.”信息的背景变量来生成“.”。这看上去就像是在解码器的每一时间步对输入序列中不同时间步编码的信息分配不同的注意力。这也是注意力机制的由来。它最早由 Bahdanau 等提出。

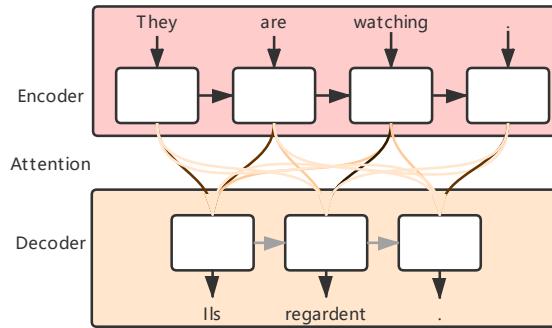


图 3: Attention + Encoder-Decoder 架构

上面是直观上的理解,下面来详细解释一下数学上有哪些具体的计算过程,和详细的算法流程。

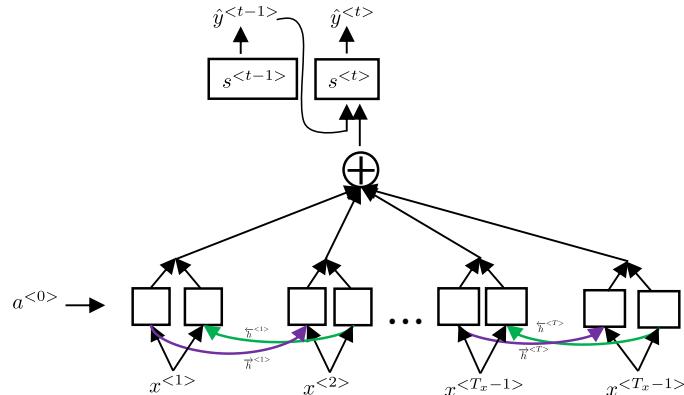


图 4: Attention 算子

- 首先我们利用 RNN 结构得到 encoder 中的 hidden state (h_1, h_2, \dots, h_T) ,
- 假设当前 decoder 的 hidden state 是 s_{t-1} , 我们可以计算每一个输入位置 j 与当前输出位置的关联性, $e_{tj} = a(s_{t-1}, h_j)$, 写成相应的向量形式即为 $\vec{e}_t = (a(s_{t-1}, h_1), \dots, a(s_{t-1}, h_T))$, 其中 a 是一种相关性的算符, 例如常见的有点乘形式 $\vec{e}_t = s_{t-1}^T \vec{h}$, 加权点乘 $\vec{e}_t = s_{t-1}^T W \vec{h}$, 加和 $\vec{e}_t = \vec{v}^T \tanh(W_1 \vec{h} + W_2 s_{t-1})$ 等等。
- 对于 \vec{e}_t 进行 softmax 操作将其 normalize 得到 attention 的分布, $\vec{\alpha}_t = \text{softmax}(\vec{e}_t)$, 展开形式为 $\alpha_{tj} = \frac{\exp(e_{tj})}{\sum_{k=1}^T \exp(e_{tk})}$
- 利用 $\vec{\alpha}_t$ 我们可以进行加权求和得到相应的 context vector $\alpha_{tj} = \frac{\exp(e_{tj})}{\sum_{k=1}^T \exp(e_{tk})}$
- 由此, 我们可以计算 decoder 的下一个 hidden state \vec{a}_t 以及该位置的输出 $\vec{c}_t = \sum_{j=1}^T \alpha_{tj} h_j$ 。

这里关键的操作是计算 encoder 与 decoder state 之间的关联性的权重, 得到 Attention 分布, 从而对于当前输出位置得到比较重要的输入位置的权重, 在预测输出时相应的会占较大的比重。

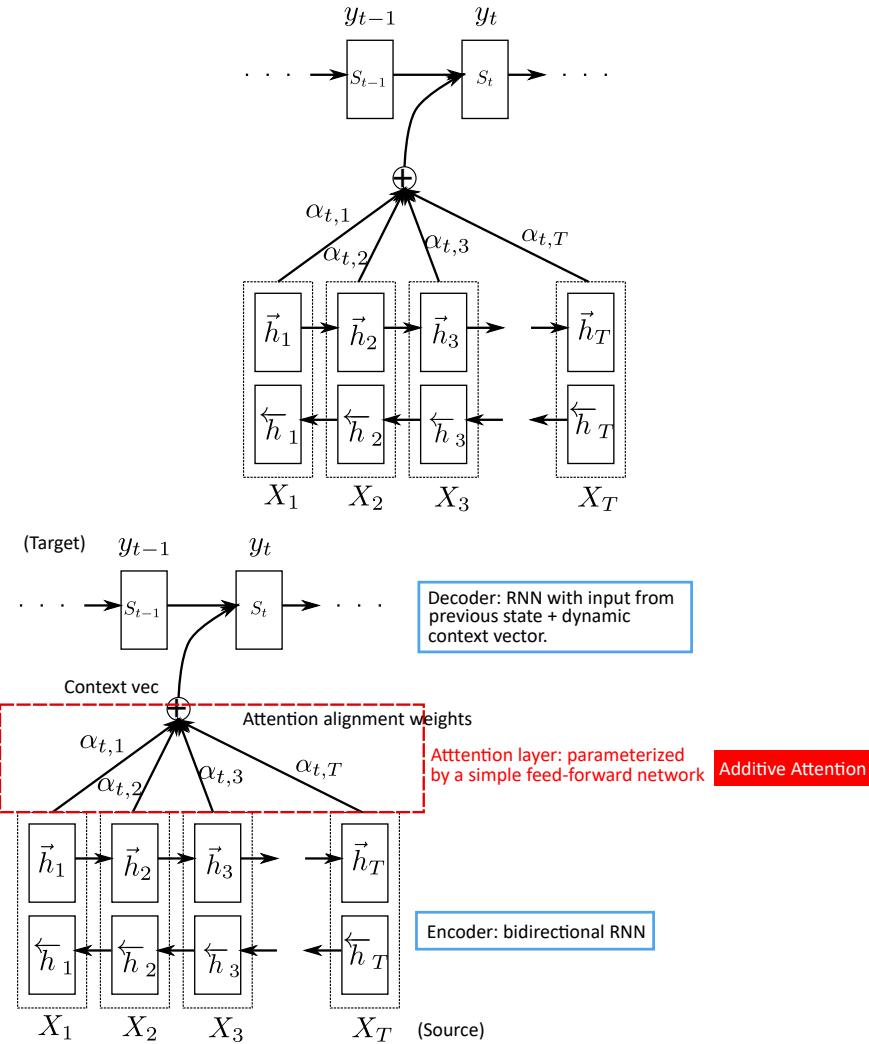


图 5: The encoder-decoder model with additive attention mechanism in Bahdanau et al., 2015.

Attention Model Definition

我们有一个长度为 n 的源序列 \mathbf{x} , 并尝试输出长度为 m 的目标序列 \mathbf{y} :

$$\mathbf{x} = [x_1, x_2, \dots, x_n] \quad (1)$$

$$\mathbf{y} = [y_1, y_2, \dots, y_m] \quad (2)$$

(粗体变量表示它们是向量; 本文中其他所有内容都相同。)

编码器是双向 RNN (或您选择的其他循环网络设置), 具有前向隐藏状态 $\vec{\mathbf{h}}_i$ 和后向隐藏状态 $\overleftarrow{\mathbf{h}}_i$ 。两个方向隐藏状态的简单“连结”表示一个编码器状态。这样做的动机是因为在一个单词的标注中包括前面和后面的单词。

$$\mathbf{h}_i = [\vec{\mathbf{h}}_i^\top; \overleftarrow{\mathbf{h}}_i^\top]^\top, i = 1, \dots, n$$

解码器网络具有隐藏状态 $\mathbf{s}_t = f(\mathbf{s}_{t-1}, y_{t-1}, \mathbf{c}_t)$ 用于在时刻 t 处来输出单词, $t = 1, \dots, m$ 其中上下文向量 \mathbf{c}_t 是输入序列的隐藏状态的总和, 由联配值 (alignment scores) 加权得到:

$$\mathbf{c}_t = \sum_{i=1}^n \alpha_{t,i} \mathbf{h}_i \quad ; \text{ Context vector for output } y_t \quad (3)$$

$$\alpha_{t,i} = \text{align}(y_t, x_i) \quad ; \text{ How well two words } y_t \text{ and } x_i \text{ are aligned.} \quad (4)$$

$$= \frac{\exp(\text{score}(\mathbf{s}_{t-1}, \mathbf{h}_i))}{\sum_{i'=1}^n \exp(\text{score}(\mathbf{s}_{t-1}, \mathbf{h}_{i'}))} \quad ; \text{ Softmax of some predefined alignment score..} \quad (5)$$

配分模型基于它们匹配的程度, 将得分 $\alpha_{t,i}$ 分配给位置 i 处的输入对和位置 t 处的输出 (y_t, x_i) 。 $\{\alpha_{t,i}\}$ 的集合是定义每个输出应考虑每个源隐藏状态的多少的权重。在 Bahdanau 的论文中, 联配值 $\boldsymbol{\alpha}$ 由具有单个隐藏层的前馈神经网络来进行参数化学习, 并且该网络与模型的其他部分共同训练。因此, 得分函数采用以下形式, 因为 \tanh 用作非线性激活函数:

$$\text{score}(\mathbf{s}_t, \mathbf{h}_i) = \mathbf{v}_a^\top \tanh(\mathbf{W}_a[\mathbf{s}_t; \mathbf{h}_i]) \quad (6)$$

其中 \mathbf{v}_a 和 \mathbf{W}_a 都是在联配模型中学习的权重矩阵。

联配值矩阵是一个很好的副产品, 可以明确显示源词和目标词之间的相关性。如下图所示:

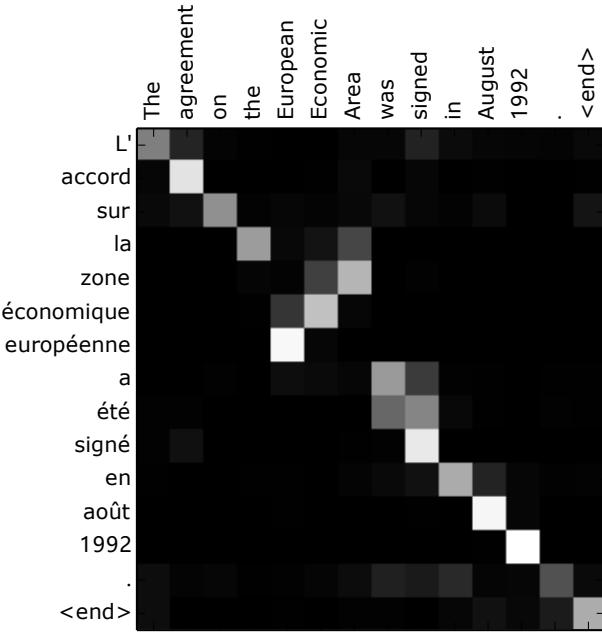


图 6: Alignment matrix of “L’ accord sur l’ Espace économique européen a été signé en août 1992” (French) and its English translation “The agreement on the European Economic Area was signed in August 1992” . (Image source: Fig 3 in Bahdanau et al., 2015.)

Attention 模型设计

在时间步 t ，设解码器的背景变量为 $c_{t'}$ ，输出 $y_{t'}$ 的特征向量为 $\mathbf{y}_{t'}$ 。和输入的特征向量一样，这里每个输出的特征向量也是模型参数。解码器在时间步 t 的隐藏状态

$$\mathbf{s}_{t'} = g(\mathbf{y}_{t'-1}, \mathbf{c}_{t'}, \mathbf{s}_{t'-1}).$$

令编码器在时间步 t 的隐藏状态为 \mathbf{h}_t ，且总时间步数为 T 。解码器在时间步 t 的背景变量为

$$\mathbf{c}_{t'} = \sum_{t=1}^T \alpha_{t't} \mathbf{h}_t,$$

其中 $\alpha_{t't}$ 是权值。也就是说，给定解码器的当前时间步 t ，我们需要对编码器中不同时间步 t 的隐藏状态求加权平均。这里的权值也称注意力权重。它的计算公式是

$$\alpha_{t't} = \frac{\exp(e_{t't})}{\sum_{k=1}^T \exp(e_{t'k})},$$

其中 $e_{t't} \in \mathbb{R}$ 的计算为

$$e_{t't} = a(\mathbf{s}_{t'-1}, \mathbf{h}_t).$$

上式中的函数 a 有多种设计方法。Bahdanau 等使用了多层感知机：

$$e_{t't} = \mathbf{v}^\top \tanh(\mathbf{W}_s \mathbf{s}_{t'-1} + \mathbf{W}_h \mathbf{h}_t),$$

其中 v 、 \mathbf{W}_s 、 \mathbf{W}_h 以及编码器与解码器中的各个权重和偏差都是模型参数。

Bahdanau 等在编码器和解码器中分别使用了门控循环单元。在解码器中，我们需要对门控循环单元的设计稍作修改。解码器在 t' 时间步的隐藏状态为

$$\mathbf{s}_{t'} = \mathbf{z}_{t'} \odot \mathbf{s}_{t'-1} + (1 - \mathbf{z}_{t'}) \odot \tilde{\mathbf{s}}_{t'},$$

其中的重置门、更新门和候选隐含状态分别为

$$\begin{aligned}\mathbf{r}_{t'} &= \sigma(\mathbf{W}_{yr}\mathbf{y}_{t'-1} + \mathbf{W}_{sr}\mathbf{s}_{t'-1} + \mathbf{W}_{cr}\mathbf{c}_{t'} + \mathbf{b}_r), \\ \mathbf{z}_{t'} &= \sigma(\mathbf{W}_{yz}\mathbf{y}_{t'-1} + \mathbf{W}_{sz}\mathbf{s}_{t'-1} + \mathbf{W}_{cz}\mathbf{c}_{t'} + \mathbf{b}_z), \\ \tilde{\mathbf{s}}_{t'} &= \tanh(\mathbf{W}_{ys}\mathbf{y}_{t'-1} + \mathbf{W}_{ss}(\mathbf{s}_{t'-1} \odot \mathbf{r}_{t'}) + \mathbf{W}_{cs}\mathbf{c}_{t'} + \mathbf{b}_s).\end{aligned}$$

A Family of Attention Mechanisms

Self-Attention

下图为 attention 以及 Transformer 研究进展的大概趋势。

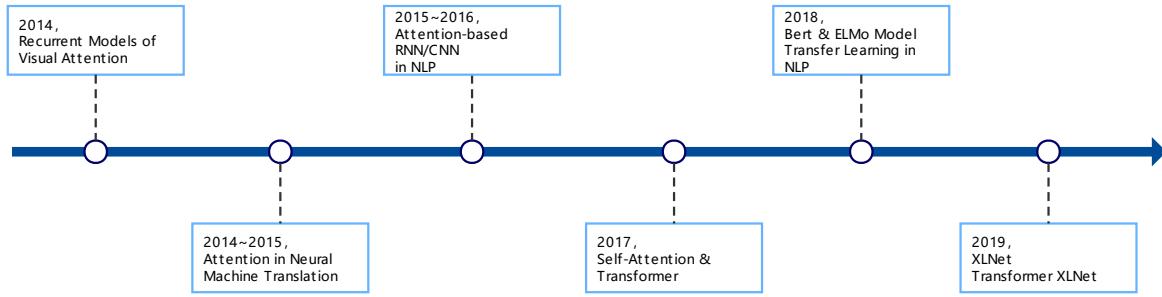


图 7：近年来 NLP 领域深度学习模型和技术发展里程表

Attention 机制的本质来自于人类视觉注意力机制。人们视觉在感知东西的时候一般不会是一个场景从头看到尾每次全部都看，而往往是根据需求观察注意特定的一部分。而且当人们发现一个场景经常在某部分出现自己想观察的东西时，人们会进行学习在将来再出现类似场景时把注意力放到该部分上。

自注意力 (Self-attention)，也叫做内部注意力，是关联单个序列 (Source) 内部元素的不同位置的注意力机制，以便计算相同序列的表示。它已被证明在机器阅读，抽象摘要或计算机图像描述生成中非常有用。

Vaswani et al., 2017 在论文 **Attention is all you need** 中提出不使用 RNN 或 CNN 等复杂的模型，仅依赖于 Attention 模型可以使训练并行化且拥有全局信息，该模型在机器翻译领域不同规模的公开数据集上都表现良好。并且，Google 最新的机器翻译模型内部大量采用了 Self-Attention 机制。

普通的 Attention 机制都发生在输出 Target 句子中某个单词和输入 Source 句子每个单词之间的相似度。而 Self-Attention 顾名思义，指的不是 Target 和 Source 之间的 Attention 机制，而是 Source 内部元素之间或者 Target 内部元素之间发生的 Attention 机制。

Cheng et al., 2016 的 LSTMN 论文使用自我注意力进行机器阅读。在下面的例子中，自注意力机制使我们能够学习当前单词和句子前一部分之间的相关性。

The FBI is chasing a criminal on the run .
 The FBI is chasing a criminal on the run .
 The FBI is chasing a criminal on the run .
 The FBI is chasing a criminal on the run .
 The FBI is chasing a criminal on the run .
 The FBI is chasing a criminal on the run .
 The FBI is chasing a criminal on the run .
 The FBI is chasing a criminal on the run .
 The FBI is chasing a criminal on the run .

图 8: The current word is in red and the size of the blue shade indicates the activation level. (Image source: Cheng et al., 2016)

在Xu et al., 2015 的 Show, Attend and Tell 论文中, 自注意力机制应用在图像上来生成合适的自然语言描述。该图像首先由卷积神经网络进行编码, 并且使用具有自注意力机制的循环神经网络作为解码器, 喂给其编码网络产生的卷积特征图以逐一生成描述性词。注意力权重的可视化清楚地表明了模型会注意哪些图像区域以便输出对应的单词。



图 9: “A woman is throwing a frisbee in a park.” (Image source: Fig. 6(b) in Xu et al., 2015)

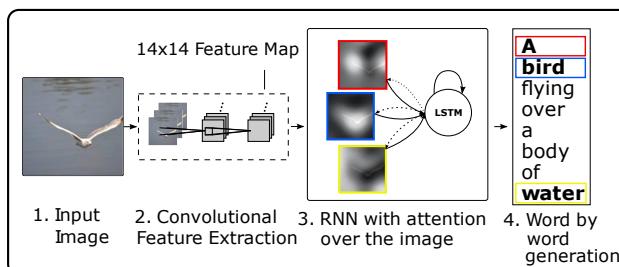


图 10: Architecture of Image-to-Text: Caption Generation with Attention (Image source: Fig. 1 in Xu et

Attention 模型在自然语言处理和计算机视觉上的应用

- NMT with Recurrent Nets and Attention Mechanism Bahdanau, ICLR' 15
- Image-to-Text: Caption Generation with Attention Xu, Kelvin, et al. Arxiv' 15
- Teaching Machines to Read and Comprehend Hermann et al. NIPS' 15
- Neural Attention Model for Sentence Summarization Rush et al. EMNLP' 15

参考文献

- Volodymyr Mnih, Nicolas Heess, Alex Graves, Koray Kavukcuoglu. *Recurrent Models of Visual Attention*. 2014. In Advances in Neural Information Processing Systems.
- Bahdanau D, Cho K, Bengio Y. *Neural machine translation by jointly learning to align and translate*[J]. arXiv preprint arXiv:1409.0473, 2014.
- Sutskever, Ilya, Oriol Vinyals, and Quoc V. Le. *Sequence to sequence learning with neural networks*. Advances in neural information processing systems. 2014.
- Hermann K M, Kočiský T, Grefenstette E, et al. *Teaching machines to read and comprehend* [C]Advances in Neural Information Processing Systems. 2015: 1693-1701.
- Xu, Kelvin, et al. *Show, attend and tell: Neural image caption generation with visual attention*. arXiv preprint arXiv:1502.03044 2.3 (2015): 5.
- Zaremba, Wojciech, and Ilya Sutskever. *Learning to execute*. arXiv preprint arXiv:1410.4615 (2014).
- Ashish Vaswani, et al. *Attention is all you need*. NIPS 2017.
- Jianpeng Cheng, Li Dong, and Mirella Lapata. *Long short-term memory-networks for machine reading*. EMNLP 2016.
- <http://www.wildml.com/2016/01/attention-and-memory-in-deep-learning-and-nlp/>
- <http://yanran.li/peppypapers/2015/10/07/survey-attention-model-1.html#fn:1>
- <https://www.cnblogs.com/robert-dlut/p/8638283.html>
- https://mp.weixin.qq.com/s/3911D_FkTWrtKwBo30vENg
- <https://lilianweng.github.io/lil-log/2018/06/24/attention-attention.html>