

# Table of Contents

- 1 命名实体识别任务
- ▼ 2 从Encoder/Decoder架构到Attention机制
  - 2.1 Encoder/Decoder架构
  - 2.2 Attention机制
  - 2.3 自注意力机制
- ▼ 3 从Transformer到BERT模型
  - 3.1 Transformer架构
  - 3.2 BERT模型
- 4 BERT-BiLSTM-CRF模型

## 命名实体识别任务的 BERT-BiLSTM-CRF模型

### 1 命名实体识别任务

输入向量序列：

$$\mathbf{X} = (\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_t, \dots, \mathbf{x}_n)$$

其中 $\mathbf{x}_t \in \mathbb{R}^m$ 为t时刻输入向量。

输出序列：

$$\mathbf{y} = (y_1, y_2, \dots, y_t, \dots, y_n)$$

其中 $y_t \in \{1, \dots, k\}$ 为t时刻输出。

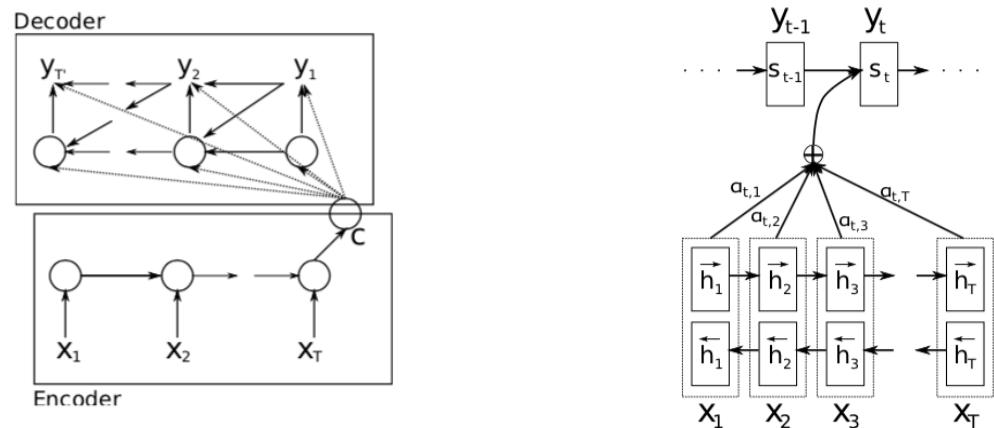
数据样例：

迈/O 向/O 充/O 满/O 希/O 望/O 的/O 新/O 世/O 纪/O —/O —/O —/O 九/O 九/O 八/O 年/O 新/O 年/O 讲/O 话/O ( /O 附/O 图/O 片/O 1 /O 张/O ) /O  
中/B\_nt 共/M\_nt 中/M\_nt 央/E\_nt 总/O 书/O 记/O  
国/O 家/O 主/O 席/O 江/B\_nr 泽/M\_nr 民/E\_nr  
( /O —/O 九/O 九/O 七/O 年/O 十/O 二/O 月/O 三/O 十/O —/O 日/O ) /O  
1 /O 2 /O 月/O 3 /O 1 /O 日/O  
中/B\_nt 共/M\_nt 中/M\_nt 央/E\_nt 总/O 书/O 记/O  
国/O 家/O 主/O 席/O 江/B\_nr 泽/M\_nr 发/O 表/O 1 /O 9 /O 9 /O 8 /O 年/O 新/O 年/O 讲/O 话/O 《 /O 迈/O 向/O 充/O 满/O 希/O 望/O 的/O 新/O 世/O 纪/O 》 /O  
( /O 新/B\_nt 华/M\_nt 社/E\_nt 记/O 者/O 兰/B\_nr 红/M\_nr 光/E\_nr 摄/O ) /O

人名：nr；地名：ns；组织名：nt。

### 2 从Encoder/Decoder架构到Attention机制

#### 2.1 Encoder/Decoder架构

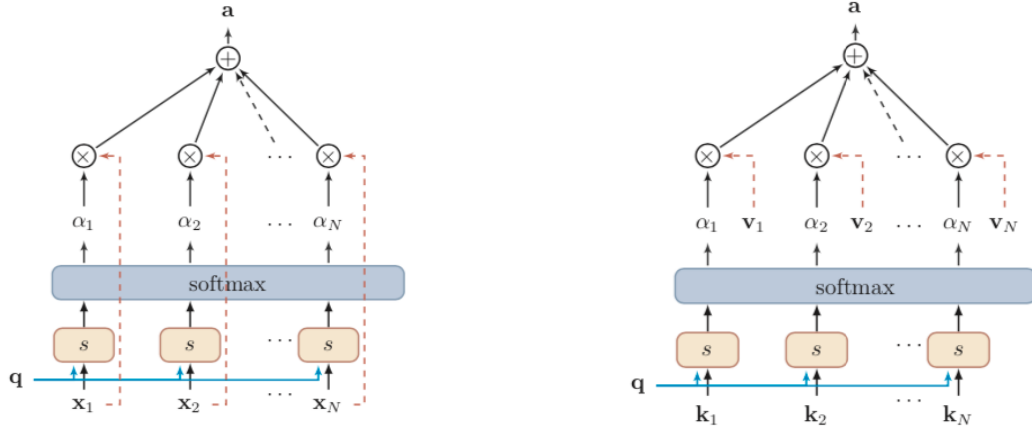


RNN Encoder-Decoder神经网络架构使用循环神经网络学习将变长源序列 $X$ 编码成定长向量表示 $\mathbf{c}$ ，并将学习的定长向量表示

c解码成变长目标序列 $\mathbf{y}$ 。模型的编码器和解码器被联合训练，以最大化给定源序列的目标序列的条件概率。

seq2seq with Attention神经网络架构中，编码器采用双向循环神经网络学习将输入序列 $\mathbf{x}$ 编码成每个时刻的上下文向量（注意力分布） $c_i$ ，解码器学习将上下文向量 $c_i$ 解码为输出序列 $\mathbf{y}$ 。

## 2.2 Attention机制



输入序列  $X = [\mathbf{x}_1, \dots, \mathbf{x}_N]$

输出序列  $H = [\mathbf{h}_1, \dots, \mathbf{h}_N]$

注意力机制的计算：

1. 在输入信息上计算注意力分布；
2. 根据注意力分布计算输入信息的加权平均。

给定一个和任务相关的查询向量 $\mathbf{q}$ ，用注意力变量 $z \in [1, N]$ 表示被选择信息的索引位置，即 $z = i$ 表示选择了第 $i$ 个输入信息。其中，查询向量 $\mathbf{q}$ 可以是动态生成的，也可以是可学习的参数。

软性注意力的注意力分布

在给定输入信息 $X$ 和查询变量 $\mathbf{q}$ 下，选择第 $i$ 个输入信息的概率

$$\begin{aligned}\alpha_i &= p(z = i | X, \mathbf{q}) \\ &= \text{softmax}(s(\mathbf{x}_i, \mathbf{q})) \\ &= \frac{\exp(s(\mathbf{x}_i, \mathbf{q}))}{\sum_{j=1}^N \exp(s(\mathbf{x}_j, \mathbf{q}))}\end{aligned}$$

其中， $\alpha_i$ 称为注意力分布， $s(\mathbf{x}_i, \mathbf{q})$ 称为注意力打分函数。

注意力打分函数

- 加性模型  $s(\mathbf{x}_i, \mathbf{q}) = \mathbf{v}^\top \tanh(W\mathbf{x}_i + U\mathbf{q})$
- 点积模型  $s(\mathbf{x}_i, \mathbf{q}) = \mathbf{x}_i^\top \mathbf{q}$
- 缩放点积模型  $s(\mathbf{x}_i, \mathbf{q}) = \frac{\mathbf{x}_i^\top \mathbf{q}}{\sqrt{d}}$
- 双线性模型  $s(\mathbf{x}_i, \mathbf{q}) = \mathbf{x}_i^\top W \mathbf{q}$

其中， $W, U, \mathbf{v}$ 为可学习的网络参数， $d$ 为输入信息的维度。

加性模型和点积模型的复杂度近似，但点积模型可利用矩阵乘积，计算效率更高。当输入信息的维度 $d$ 比较高，点积模型值方差较大，导致softmax函数的梯度较小，缩放点积模型可以解决。双线性模型是泛化的点积模型。若假设 $W = U^\top V$ ，则 $s(\mathbf{x}_i, \mathbf{q}) = \mathbf{x}_i^\top U^\top V \mathbf{q} = (U\mathbf{x}_i)^\top (V\mathbf{q})$ ，即分别对 $\mathbf{x}_i$ 和 $\mathbf{q}$ 进行线性变换后进行点积。相比点积模型，双线性模型在计算相似度是引入了非对称性。

注意力函数

$$\begin{aligned}\text{att}(X, \mathbf{q}) &= \sum_{i=1}^N \alpha_i \mathbf{x}_i \\ &= \mathbb{E}_{z \sim p(z | X, \mathbf{q})} [\mathbf{x}]\end{aligned}$$

## 2.3 自注意力机制

输入序列  $X = [\mathbf{x}_1, \dots, \mathbf{x}_N] \in \mathbb{R}^{d_1 \times N}$

输出序列  $H = [\mathbf{h}_1, \dots, \mathbf{h}_N] \in \mathbb{R}^{d_2 \times N}$

通过线性变换得到向量序列：

$$Q = W_Q X \in \mathbb{R}^{d_3 \times N}$$

$$K = W_K X \in \mathbb{R}^{d_3 \times N}$$

$$V = W_V X \in \mathbb{R}^{d_2 \times N}$$

其中， $Q$ 为查询向量序列， $K$ 为键向量序列， $V$ 为值向量序列， $W_Q, W_K, W_V$ 分别为可学习参数矩阵。

预测输出向量

$$\begin{aligned}\hat{\mathbf{h}}_i &= att((K, V), \mathbf{q}_i) \\ &= \sum_{j=1}^N \alpha_{ij} \mathbf{v}_j \\ &= \sum_{j=1}^N softmax(s(\mathbf{k}_j, \mathbf{q}_i)) \mathbf{v}_j\end{aligned}$$

其中， $i, j \in [1, N]$ 为输出和输入向量序列的位置，连接权重 $\alpha_{ij}$ 由注意力机制动态生成。

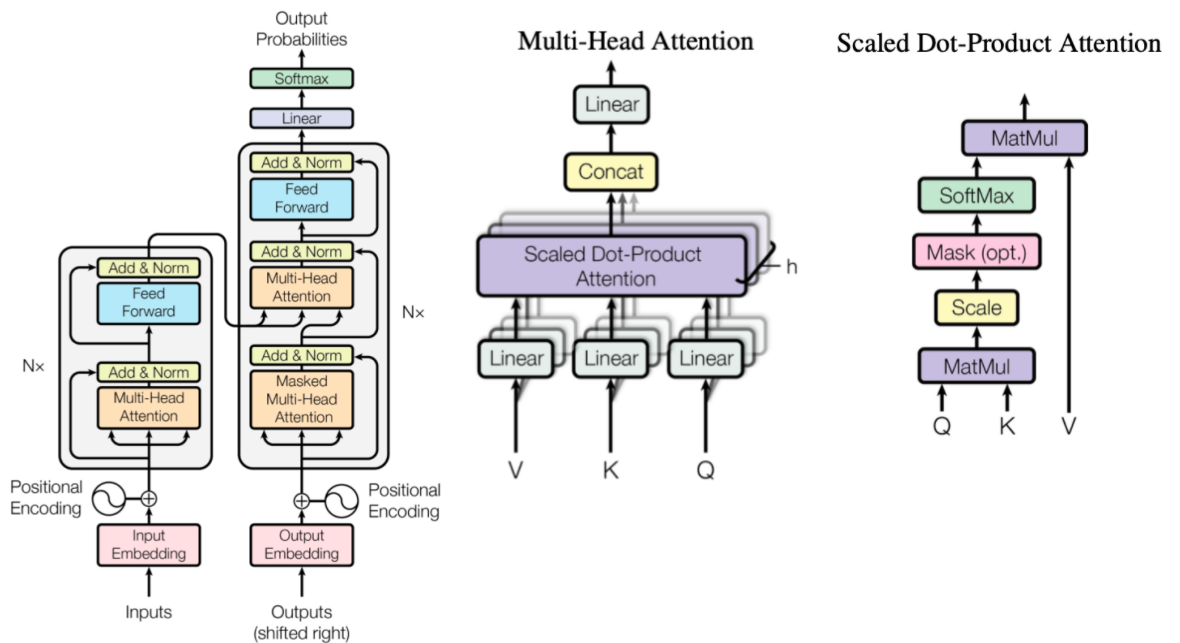
若使用缩放点积模型作为打分函数，则输出向量序列

$$\begin{aligned}H_{d_2 \times N} &= softmax\left(\frac{K^\top Q}{\sqrt{d_3}}\right) V_{d_2 \times N} \\ &= softmax\left(\frac{K^\top Q}{\sqrt{d_3}}\right) W_V X\end{aligned}$$

其中，softmax为按列归一化的函数。

## 3 从Transformer到BERT模型

### 3.1 Transformer架构



输入序列  $inputs = (i_1, i_2, \dots, i_p, \dots, i_N)$ ，其中  $i_p \in \mathbb{N}^*$  为输入符号表中的序号。

目标序列  $targets = (t_1, t_2, \dots, t_q, \dots, t_M)$ ，其中  $t_q \in \mathbb{N}^*$  为目标符号表中的序号。

$$outputs\_probabilities = Transformer(inputs, targets)$$

其中,  $outputs\_probabilities = (o_1, o_2, \dots, o_q, \dots, o_M)$  为预测序列,  $o_q \in \mathbb{N}^*$  为目标符号表中的序号。

输入序列词嵌入  $Embedding(inputs) \in \mathbb{R}^{N \times d_{model}}$ , 其中,  $N$  为输入序列长度,  $d_{model}$  为词嵌入维度。

编码器结构:

$$e_0 = Embedding(inputs) + Pos\_Enc(inputs\_position)$$

$$e_l = EncoderLayer(e_{l-1}), l \in [1, n]$$

其中,  $e_0 \in \mathbb{R}^{N \times d_{model}}$  为编码器输入,  $EncoderLayer(\cdot)$  为编码器层,  $n$  为层数,  $e_l \in \mathbb{R}^{N \times d_{model}}$  为第  $l$  层编码器层输出。

编码器层  $EncoderLayer$ :

$$e_{mid} = LayerNorm(e_{in} + MultiHeadAttention(e_{in}))$$

$$e_{out} = LayerNorm(e_{mid} + FFN(e_{mid}))$$

其中,  $e_{in} \in \mathbb{R}^{N \times d_{model}}$  为编码器层输入,  $e_{out} \in \mathbb{R}^{N \times d_{model}}$  为编码器层输出,  $MultiHeadAttention(\cdot)$  为多头注意力机制,  $FFN(\cdot)$  为前馈神经网络,  $LayerNorm(\cdot)$  为层归一化。

输入向量序列  $e_{in} = (e_{in1}, e_{in2}, \dots, e_{inN}) \in \mathbb{R}^{N \times d_{model}}$ , 分别得到查询向量序列  $Q = e_{in}$ , 键向量序列  $K = e_{in}$ , 值向量序列  $V = e_{in}$ 。

多头注意力机制

$$MultiHeadAttention(e_{in}) = MultiHead(Q, K, V) = Concat(head_1, \dots, head_h) W^O$$

其中, 多头输出  $head_i = Attention(QW_i^Q, KW_i^K, VW_i^V)$ , 可学习的参数矩阵

$$W_i^Q \in \mathbb{R}^{d_{model} \times d_k}, W_i^K \in \mathbb{R}^{d_{model} \times d_k}, W_i^V \in \mathbb{R}^{d_{model} \times d_v}, W^O \in \mathbb{R}^{hd_v \times d_{model}}$$

使用缩放点积作为打分函数的自注意力机制

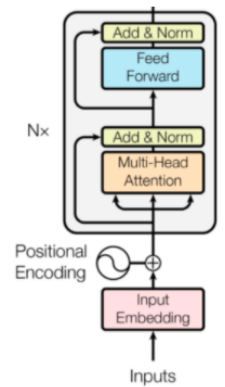
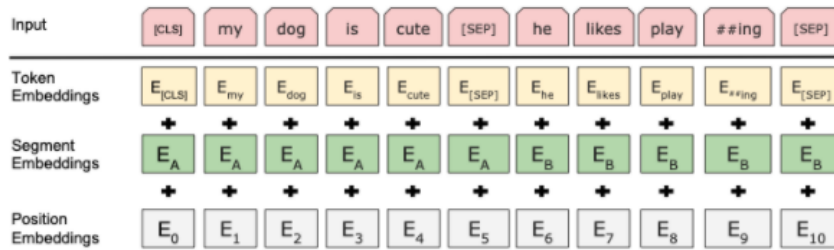
$$Attention(QW_i^Q, KW_i^K, VW_i^V) = softmax\left(\frac{QW_i^Q(KW_i^K)^T}{\sqrt{d_k}}\right) VW_i^V$$

$$FFN(e_{mid}) = ReLU(e_{mid}W_1 + b_1)W_2 + b_2$$

$$= \max(0, e_{mid}W_1 + b_1)W_2 + b_2$$

其中, 参数矩阵  $W_1 \in \mathbb{R}^{d_{model} \times d_{ff}}, W_2 \in \mathbb{R}^{d_{ff} \times d_{model}}$ , 偏置  $b_1 \in \mathbb{R}^{d_{ff}}, b_2 \in \mathbb{R}^{d_{model}}$ 。

### 3.2 BERT模型

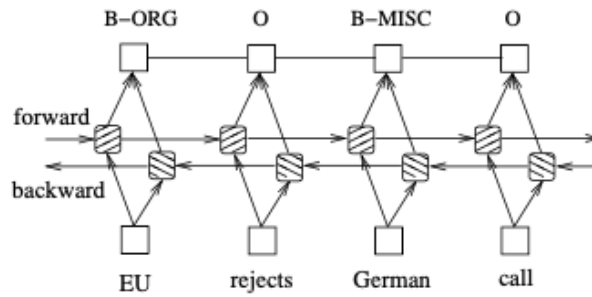


输入序列

$$input = ([CLS], s_1, s_2, \dots, s_m, [SEP], p_1, p_2, \dots, p_n, [SEP])$$

其中,  $s_i, p_j \in \mathbb{N}$  为输入符号表中的序号; 子序列  $(s_1, \dots, s_m)$  为句子对中前序句子; 子序列  $(p_1, \dots, p_n)$  为句子对中后续句子; 输入序列首标记  $[CLS]$  用作分类任务表示; 特殊标记  $[SEP]$  用作区分句子对各子句。

## 4 BERT-BiLSTM-CRF模型



BiLSTM输出分值矩阵：

$$\mathbf{P} = [P_{i,j}]_{n \times k}$$

其中 $P_{i,j}$ 为第 $i$ 个单词对应第 $j$ 个标签的分数。

转移分值矩阵：

$$\mathbf{A} = [A_{i,j}]_{(k+2) \times (k+2)}$$

其中 $A_{i,j} = p(y_i | y_{i-1})$ 。

模型分值：

$$s(\mathbf{X}, \mathbf{y}) = \sum_{i=1}^{n+1} A_{i,j} + \sum_{i=1}^n P_{i,j}$$

所有可能的标签序列上的softmax产生序列 $\mathbf{y}$ 的概率

$$p(\mathbf{y}|\mathbf{X}) = \frac{e^{s(\mathbf{X}, \mathbf{y})}}{\sum_{\tilde{\mathbf{y}} \in \mathbf{Y}} e^{s(\mathbf{X}, \tilde{\mathbf{y}})}}$$

预测输出：

$$\begin{aligned} \mathbf{y}^* &= \arg \max_{\tilde{\mathbf{y}} \in \mathbf{Y}} p(\mathbf{y}|\mathbf{X}) \\ &= \arg \max_{\tilde{\mathbf{y}} \in \mathbf{Y}} \log p(\mathbf{y}|\mathbf{X}) \\ &= \arg \max_{\tilde{\mathbf{y}} \in \mathbf{Y}} s(\mathbf{X}, \mathbf{y}) \\ &= - \arg \min_{\tilde{\mathbf{y}} \in \mathbf{Y}} s(\mathbf{X}, \mathbf{y}) \end{aligned}$$

BERT模型生成的Embedding词嵌入表示作为双向LISTM模型的输入。

In [ ]:

1