

# 结合 Bi-LSTM 和注意力模型的问答系统研究

邵 曦 陈 明

(南京邮电大学通信与信息工程学院 江苏 南京 210003)

**摘 要** 针对传统的问答系统普遍存在回答准确率不高、语义识别能力差等问题,提出一种结合双向长短时记忆网络(Bi-LSTM)和注意力(Attention)模型的问答系统。利用生成的句向量,学习句子中的语义特征以及问答之间的匹配关系,获取上下文信息;融合注意力模型,能够找到对话的主题信息,从而为用户做出精准的回答。实验结果表明,该系统的回答准确率高于其他模型,可达到 80.76%。

**关键词** 深度学习 Bi-LSTM 注意力模型 句向量 问答系统

中图分类号 TP389.1 文献标志码 A DOI: 10.3969/j.issn.1000-386x.2020.10.009

## A QUESTION AND ANSWER SYSTEM BASED ON BI-LSTM AND ATTENTION MODEL

Shao Xi Chen Ming

(College of Telecommunications and Information Engineering, Nanjing University of Posts and Telecommunications,  
Nanjing 210003, Jiangsu, China)

**Abstract** Traditional question and answer systems generally have problems such as low accuracy of answering and poor semantic recognition. In order to solve the above shortcomings, this paper proposes a question and answer system based on Bi-LSTM and attention model. By using the generated sentence vector, the system learned the semantic features of sentences and the matching relationship between questions and answers to get the context information. It integrated the attention model to find the topic information of the dialogue, so as to make accurate answers for the users. The experimental results show that the answer accuracy of our question and answer system is higher than other models, which can reach 80.76%.

**Keywords** Deep learning Bi-LSTM Attention model Sentence vector Question and answer system

## 0 引 言

随着人工智能技术的快速发展,其实际应用场景也变得越来越广泛,在人们生活中的各个方面都有所体现。信息爆炸的今天,人们对于搜索引擎简单地返回一个相关网页感到不满,而问答系统能够与用户一对一进行交互,精确理解用户意图,从而能够高效快速地完成用户的需求<sup>[1-2]</sup>。最初对于问答系统的研究受限于语料数据的限制,并没有取得很好的效果。但随着互联网的发展,微博、Twitter 等聊天工具的兴起,为问答系统模型训练提供了大量的文本数据。1966 年,

Weizenbaum<sup>[3]</sup>开发了最早的智能问答机器人 ELIZA;2004 年,Knill 等<sup>[4]</sup>研发了问答机器人 Sofia。这些早期的智能问答大都是基于检索技术和机器学习算法来实现的。常用的机器学习算法包括潜在狄克雷分配模型(Latent Dirichlet Allocation, LDA)、支持向量机(Support Vector Machine, SVM)等,并利用贝叶斯法和 K 近邻等方法进行分类,以此构建问答之间的对应关系<sup>[5]</sup>。此类算法对于数据提出了很高的要求,必须有充足的数据才能保证匹配的准确性,且对于不同场景算法的泛化性较低,使用性能较差。

随着深度学习技术的快速发展,其在图像处理、语音交互等领域取得了优异的成绩<sup>[6]</sup>。近年来,深度学

收稿日期:2019-06-16。国家自然科学基金面上项目(61872199, 61872424)。邵曦,副教授,主研领域:多媒体信息处理系统。陈明,硕士生。

习在自然语言处理领域也大放异彩<sup>[7]</sup>。2012 年, Mikolov 等<sup>[8]</sup>发现了一种基于循环神经网络(Recurrent Neural Network, RNN)包含上下文信息的语言模型, 将 RNN 应用到上下文信息的获取中去。Schuster 等<sup>[9]</sup>提出了 Bi-RNN 模型, 可以利用句子的未来信息进行预测。因此, 深度学习模型对于问答系统的研究具有十分重要的作用。鉴于深度学习技术在自然语言处理方向上具有不错的效果, 本文提出了一种基于 Bi-LSTM 和注意力模型的问答系统。通过生成句向量以获取句子上下文之间的语义信息和匹配关系, 结合注意力模型, 找到主题信息, 从而生成最佳的回答。

目前, 对于问答系统的研究主要分为两个方向: 基于统计特征的机器学习方法和基于深度学习的方法。随着深度学习技术的发展和各种深度神经网络模型的提出, 基于深度学习的方法成为了当前的研究热点。Kim<sup>[10]</sup>提出一种基于卷积神经网络(Convolution Neural Network, CNN)的分类模型, 利用训练好的词向量模型进行文本分类。Shi 等<sup>[11]</sup>提出了基于长短时记忆网络(Long Short-Term Memory, LSTM)的映射分类模型。为了提高回答的准确率以及效率, Sutskever 等<sup>[12]</sup>提出了序列到序列(Sequence to Sequence, Seq2Seq)框架, 通过输入端编码形成中间语义, 再解码出相应的回答。Feng 等<sup>[13]</sup>提出了基于共享卷积神经网络用于进行训练问答模型, 通过该模型进行语义相似度计算, 并且在英文数据集上取得了优异的成绩。注意力模型的提出使得自然语言处理领域又有了一个新的研究方向, 该模型通过模拟人脑的机制, 对语句中的信息进行加权处理, 从而对语句主题信息进行重点关注。Yin 等<sup>[14]</sup>提出了一种基于注意力模型的多层卷积神经网络模型, 实现了对文本语义的建模, 并且在问答匹配和语义识别上都取得了很好的效果。目前, 关于注意力模型大致可以分为两类: Soft Attention 和 Hard Attention。Soft Attention 是在求注意力概率分布的时候, 对输入句子中每个单词都给出权重, 其为概率分布; 而 Hard Attention 在进行权重分配时, 只会对句中某个或某几个单词进行分配, 把目标句子单词和这个单词进行对齐, 句中其他单词硬性认为对齐概率为 0。本文提出的注意力模型是基于 Soft Attention 机制实现的。

## 1 系统设计

图 1 为本文提出的结合 Bi-LSTM 和注意力模型的问答系统框图。在本系统中, 输入问句会首先经过句向量(doc2vec)层生成相应的句向量, 然后将生成的向量作为 Bi-LSTM 的网络输入, 最后将网络模型的输出

通过 Attention 机制形成最后的输出结果。

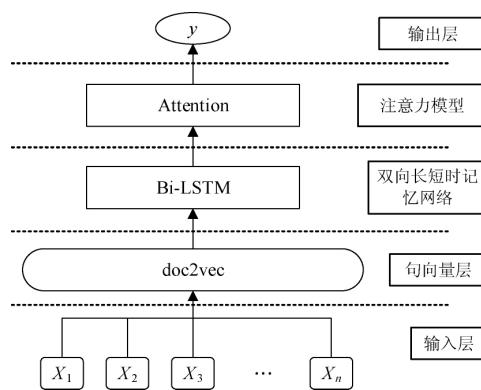


图 1 基于 Bi-LSTM 和注意力模型问答系统

### 1.1 句向量层

在词向量(word2vec)技术占据主流时, Le 等<sup>[15]</sup>在 word2vec 的基础上进行拓展, 提出了句向量(doc2vec)技术。在 word2vec 技术中, 主要分为 Continuous Bag Of Words(CBOW)模型和 Skip-gram 模型。本文主要讨论 CBOW 模型, 其模型结构如图 2 所示。该模型分为输出层和输入层, 相比传统的语言模型, 由于其去掉了隐藏层, 所以运算速度得到大幅提升。CBOW 模型使用一段文本的中间词作为目标词, 即利用该中间词的上下词来预测该词, 通过利用上下文各词的词向量的平均值来替代之前模型各个拼接的词向量, 可以提高模型预测的准确性。

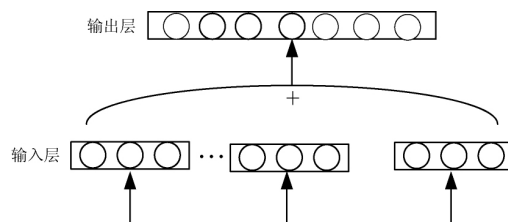


图 2 CBOW 模型结构图

与 word2vec 相对应, doc2vec 也存在两种模型: DM (Distributed Memory) 模型和 DBOW (Distributed Bag of Words) 模型, 本文采用的是 DM 模型。DM 模型框架如图 3 所示。

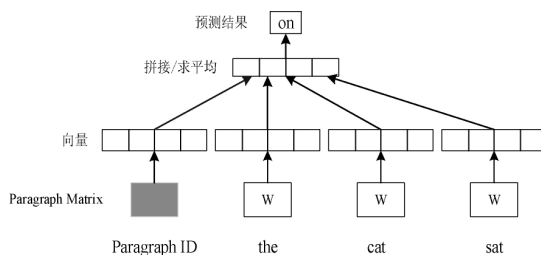


图 3 DM 模型

DM 模型增加了一个 Paragraph ID, 在训练过程中, Paragraph ID 会先映射成一个向量, 且与词向量的维数相同, 但二者属于不同的向量空间。在之后的计算中, paragraph vector 和 word vector 累加起来, 输入

softmax 层,从而输出预测结果。以“the cat sat”为例,通过 DM 模型可以预测出下一个词为“on”。在一个句子的训练过程中,paragraph ID 始终保持不变,共享同一个 paragraph vector,相当于每次在预测单词概率时,都利用了整个句子的语义。在进行预测时,给每一个句子分配一个 paragraph ID,词向量和输出层 softmax 参数保持训练得到的不变,通过随机梯度下降法训练预测语句。本文利用句向量代替词向量进行预测,从而可以充分利用语句的语序信息,准确理解语句意图,并将其输入下一层的网络。

## 1.2 双向长短时记忆网络

Bi-LSTM 是由 LSTM 演化而来,LSTM 的提出是为了解决神经网络在面对长序列时产生的梯度消失问题<sup>[16]</sup>。LSTM 网络模型由各个记忆单元组成,通过输入门、遗忘门和输出门来控制记忆单元的存储内容,通过门的控制可以在新的状态下不断叠加输入序列,从而对前面的信息具有记忆功能。LSTM 的网络结构如图 4 所示。

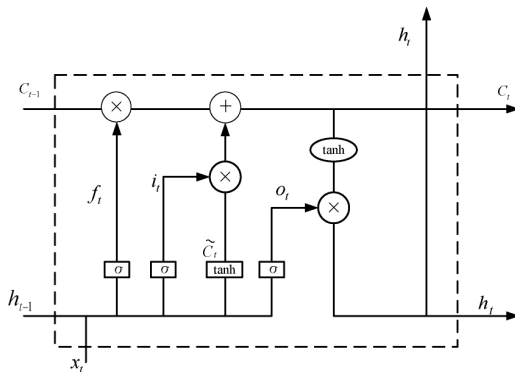


图4 LSTM网络结构

$x = \{x_1, x_2, \dots, x_t\}$  表示输入序列  $h = \{h_1, h_2, \dots, h_t\}$  表示记忆单元的输出  $C_t$  表示记忆单元的记忆内容。网络具体计算方式如下:

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

$$i_t = \sigma(w_i \cdot [h_{t-1}, x_t] + b_i) \quad (2)$$

$$\tilde{C}_t = \tanh[w_c \cdot [h_{t-1}, x_t] + b_c] \quad (3)$$

$$C_t = f_t * C_{t-1} + i_t * \tilde{C}_t \quad (4)$$

$$o_t = \sigma(w_o \cdot [h_{t-1}, x_t] + b_o) \quad (5)$$

$$h_t = o_t * \tanh(C_t) \quad (6)$$

式中:  $w$  表示权重矩阵;  $b$  表示偏置向量;  $\sigma$  表示 sigmoid 函数。上一单元的输出和当前输入信息经过遗忘门后,网络决定需要抛弃哪种信息,然后通过输入门进行信息更新,将  $C_{t-1}$  更新为  $C_t$ ,最后通过输出门确定需要输出的值  $o_t$ 。

由于单向 LSTM 网络在进行训练时只考虑到句子的时序信息而忽略了上下文之间的关系,因此在进行

多句对话时,往往不能取得很好的效果。Bi-LSTM 网络是由前向 LSTM 和后向 LSTM 组成,可以充分利用序列的上下文信息。本文提出用 Bi-LSTM 对问句进行训练,融合前向 LSTM 和后向 LSTM 的结果进行输出,从而可以提高模型回答的准确率。

## 1.3 注意力模型

注意力模型是由 Bahdanau 等<sup>[17]</sup>提出,该模型借鉴了人脑的思维模式,即人的注意力一定是集中在目光看到的事物上,随着目光的转移,注意力也在转移。在自然语言处理领域中,注意力模型往往是附着在 Encoder-Decoder(即 Seq2Seq) 框架下使用。该框架是处理由一个句子生成另一个句子的通用处理模型,对于输入问句  $X$  经过 Encoder 模块编码形成中间内容向量  $C$ ,Decoder 模块根据中间内容向量  $C$  和之前已经生成的历史信息解码出该时刻的单词。在该框架中,解码生成目标单词时,采用的都是同一个内容向量  $C$ ,因此无法获取关键信息,例如以下对话:

Q: 南京有什么好玩的地方?

A: 南京是江苏省会,著名的六朝古都,是一座文化名城,有新街口、总统府和夫子庙等旅游景点。

当根据这个问题去生成答案时,“新街口”“总统府”“夫子庙”等回答与问题语义更加贴切,所以在模型中应当突出这些关键词语的作用。采用将注意力模型应用到 Encoder-Decoder 框架中,可以有效地为不同词语分配不同的权重,达到获取对话主题信息的目的。融合注意力模型的 Encoder-Decoder 框架如图 5 所示。

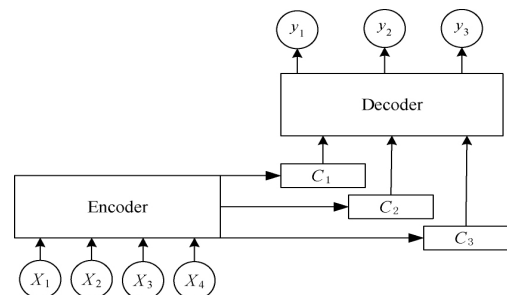


图5 融合注意力模型的 Encoder-Decoder 框架

图 5 中的 Encoder-Decoder 框架融合了注意力模型,在解码时条件概率可以写为:

$$p(y_1, y_2, \dots, y_{i-1}, x) = g(y_{i-1}, s_i, c_i) \quad (7)$$

式中:  $s_i$  可以是一个非线性的多层神经网络,表示解码器在  $i$  时刻的隐藏状态,其计算公式如下:

$$s_i = f(s_{i-1}, y_{i-1}, c_i) \quad (8)$$

式中:  $f(\cdot)$  表示某种非线性函数。可以看出,目标输出与相对应的内容向量  $c_i$  有关,相比于传统的 Encoder-Decoder 框架只有一个内容向量  $C$ ,图 5 模型具有更好的获取主题信息的能力。 $c_i$  由编码时的隐藏向量序列加权得到:

$$c_i = \sum_{j=1}^L \alpha_{ij} h_j \quad (9)$$

$$\alpha_{ij} = \frac{\exp(e_{ij})}{\sum_{k=1}^L \exp(e_{ik})} \quad (10)$$

式中:  $e_{ij} = a(s_{i-1}, h_j)$ ;  $L$  表示输入源语句单词的个数;  $\alpha_{ij}$  表示在输入第  $i$  个单词时源输入语句第  $j$  个单词的注意力分配系数,  $\alpha_{ij}$  的值越高, 表明在生成第  $i$  个输出时受第  $j$  个输入的影响越大;  $h_j$  表示源输入语句第  $j$  个单词的语义编码。将第  $i-1$  时刻的隐藏状态  $s_{i-1}$  和  $h_j$  通过前馈神经网络  $a$  计算得到一个数值, 然后使用 softmax 得到  $i$  时刻输出在  $L$  个输入隐藏状态中的分配系数。

本文研究的问答系统就是采用融合了注意力模型的 Encoder-Decoder 框架, 对问句进行编码并解码生成出相应的回答, 具有回答准确、多样等特点。

## 2 实验

### 2.1 实验环境

本文的实验都是在 TensorFlow 框架下进行的。TensorFlow 是一个经典的深度学习框架, 具有一个很强大的库以支持大规模的数值计算, 并在后端使用 C++ 加快其计算速度, 拥有丰富的高级机器学习应用程序接口 (API) 可以使其更容易地配置、训练模型。

### 2.2 实验数据收集及处理

本文实验数据来源于百度客服推广的对话录音, 利用科大讯飞的语音转文字工具, 将对话录音转换成文本形式。将对话文本整理成两个 txt 文档, 分别是问句文档 question.txt 和答句文档 answer.txt。整个问答文档包括大约 80 000 条对话数据, 其中 60 000 条作为训练数据集, 20 000 条作为测试数据集。采用 jieba 分词器进行文档的中文分词, 然后进行去标点处理等操作, 为句向量的生成做准备。

### 2.3 结果与分析

利用 TensorFlow 框架对该神经网络模型迭代训练, 直至网络收敛并保存模型, 模型参数如表 1 所示。

表 1 模型的相关参数

参数名	值
初始学习率	0.001
迭代次数	80 000
批尺寸	256
句向量维度	256

在对话问答系统领域, 主要采用计算系统生成的回答与参考回答之间的余弦相似度, 以此来对问答系统的性能进行评价。

本文的问答系统采用上述方法进行性能评价。通过使用句向量模型分别计算生成回答与参考回答的句向量, 然后计算二者之间的余弦相似度确定生成回答的准确性。为了更好地验证本文所提出的方法, 另外设置了两组对照实验, 分别是: 未融合注意力模型的 Seq2Seq 问答模型和基于检索技术的问答模型。不同模型的回答准确率如表 2 所示。

表 2 各个模型实验结果对比

模型	方法	准确率/%
检索模型	词向量	71.34
Seq2Seq 模型	LSTM	69.56
本文模型	Bi-LSTM、Attention	80.76

可以看出, 本文提出的基于 Bi-LSTM 和注意力模型的问答系统在回答的准确率上较传统的两种模型都有较大的提升, 其准确率可达到 80.76%。对于检索模型, 其只能依靠训练过的语料库来进行回答, 不能对训练语料库中未出现的问句做出有效的回答, 但对于训练语料库中出现过的问题, 可以给出准确的回答; 而 Seq2Seq 模型虽然能够解决检索模型的一些缺陷, 但是在解决长对话问题和理解用户语义方面仍然存在不足, 往往不能做出准确的回答。本文提出的模型充分利用了 Bi-LSTM 对上下文信息的获取能力以及注意力模型对主题信息的获取能力, 从而提高了问答系统的准确率。图 6 为本文提出的模型的问答效果。

Q: 你好  
A: 您好, 我们是做百度推广的, 您有需要吗  
Q: 什么是百度推广  
A: 就是帮您做业务推广, 网站建设一系列相关工作  
Q: 那你们怎么收费呢  
A: 我们每年向您收取八千四百元包括两千元的服务费  
Q: 太高了  
A: 我们这边已经给出最优惠的价格了  
Q: 再考虑考虑  
A: 好的, 如果您后面还有这方面的需要的话, 我们再联系  
Q: 再见  
A: 再见, 祝您生活愉快

图 6 对话效果演示

可以看出, 本文设计的问答系统在面对提问时, 能够理解问题的语义并结合上下文对话信息, 做出较为

准确的回答。例如:在面向“太高了”这个简短的问句时,系统能够联系到上文的价格信息,最后生成相应的回答。

### 3 结 语

本文提出一种结合 Bi-LSTM 和注意力模型的问答系统,采用了 Encoder-Decoder 框架实现系统的问答,并进行了模型训练和测试。设置对照实验,通过回答的准确率对系统的性能进行评估。实验表明,本文系统在问答结果上能够达到较高的准确率,可以较为流畅地回答问题,解决了传统问答系统回答僵硬、理解能力不足等缺陷。在今后的研究工作中,将着重研究句法分析、数据预处理等,以降低中文问答语料库不足对模型训练带来的影响。同时,尝试利用不同的神经网络模型构建系统,观察是否能够更好地提高系统回答准确率,以及如何将注意力模型应用到整个问句上也是未来研究的方向。

### 参 考 文 献

- [1] 冯升. 聊天机器人问答系统现状与发展[J]. 机器人技术与应用, 2016(4): 34-36.
- [2] 王元卓, 贾岩涛, 刘大伟, 等. 基于开放网络知识的信息检索与数据挖掘[J]. 计算机研究与发展, 2015, 52(2): 456-474.
- [3] Weizenbaum J. ELIZA-a computer program for the study of natural language communication between man and machine[J]. Communications of the ACM, 1966, 9(1): 36-45.
- [4] Knill O, Carlsson J, Chi A, et al. An artificial intelligence experiment in college math education[EB]. [2019-03-24]. <http://www.math.harvard.edu/~knill/preprints/sofia.pdf>.
- [5] 李航. 统计学习方法[M]. 北京: 清华大学出版社, 2018.
- [6] 尹宝才, 王文通, 王立春. 深度学习研究综述[J]. 北京工业大学学报, 2015, 41(1): 48-59.
- [7] Hu B T, Lu Z D, Li H, et al. Convolutional neural network architectures for matching natural language sentences[C]//NIPS Proceedings of the 27th International Conference on Neural Information Processing Systems. ACM, 2014: 2042-2050.
- [8] Mikolov T, Zweig G. Context dependent recurrent neural network language Model[C]//2012 IEEE Spoken Language Technology Workshop, 2012: 234-239.
- [9] Schuster M, Paliwal K K. Bidirectional recurrent neural networks[J]. IEEE Transactions on Signal Processing, 1997, 45(11): 2673-2681.
- [10] Kim Y. Convolutional neural networks for sentence classification[EB]. [2019-03-24]. arXiv: 1408.5882, 2014.
- [11] Shi Y Y, Yao K S, Tian L, et al. Deep LSTM based feature mapping for query classification[C]//Proceedings of the 15th Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, 2016: 1501-1511.
- [12] Sutskever I, Vinyals O, Quoc L. Sequence to Sequence learning with neural networks[C]//NIPS Proceedings of the 27th International Conference on Neural Information Processing Systems. ACM, 2014: 3104-3112.
- [13] Feng M W, Xiang B, Glass M R, et al. Applying deep learning to answer selection: a study and an open task[C]//Proceedings of the IEEE Workshop on Automatic Speech Recognition and Understanding. IEEE, 2015: 813-820.
- [14] Yin W P, Schytze H, Xiang B, et al. ABCNN: attention-based convolutional neural network for modeling sentence pairs[EB]. [2019-03-24]. arXiv: 1512.05193, 2018.
- [15] Le Q V, Mikolov T. Distributed representations of sentences and documents[C]//Proceedings of the 31st International Conference on Machine Learning, 2014: 1188-1196.
- [16] 荣光辉, 黄震华. 基于深度学习的问答匹配方法[J]. 计算机应用, 2017, 37(10): 2861-2865.
- [17] Bahdanau D, Cho K, Bengio Y. Neural machine translation by jointly learning to align and translate[EB]. [2019-03-24]. arXiv: 1409.0473, 2016.

(上接第33页)

- [7] 李志超. 两轮自平衡机器人 LQR-模糊控制算法研究[D]. 哈尔滨: 哈尔滨理工大学, 2014.
- [8] 宋松. 两轮自平衡机器人自适应模糊神经网络控制研究[D]. 哈尔滨: 哈尔滨理工大学, 2014.
- [9] 楚焱芳, 张瑞华. 模糊控制理论综述[J]. 科技信息, 2009(20): 161-162.
- [10] 肖建, 赵涛. T-S 模糊控制综述与展望[J]. 西南交通大学学报, 2016, 51(3): 462-474.
- [11] Zadeh L A. The concept of a linguistic variable and its application to approximate reasoning—I[J]. Information sciences, 1975, 8(3): 199-249.
- [12] 阮晓钢, 蔡建美, 李欣源, 等. 两轮自平衡机器人的研究与设计[M]. 北京: 科学出版社, 2012.
- [13] 高为炳. 变结构控制理论基础[M]. 北京: 中国科学技术出版社, 1990.
- [14] 姚娟. 滑模变结构中抖振消除方法研究[J]. 机电工程技术, 2016, 45(9): 90-92.
- [15] 自平衡小车建模及仿真手册 GBOT2001 V2014B[M]. 深圳: 固高科技, 2014.
- [16] 宋雅伟, 孙文, 寇恒静. 步态运动学及动力学研究方法[J]. 中国组织工程研究与临床康复, 2010, 14(2): 321-324.
- [17] 王艳敏. 柔性机械手非奇异终端滑模控制方法的研究[D]. 哈尔滨: 哈尔滨工业大学, 2009.