

基于深度学习的问题分类方法研究

李 超 柴玉梅 南晓斐 高明磊
(郑州大学信息工程学院 郑州 450001)

摘要 问题分类是问答系统中的重要组成部分。但现阶段的问题分类需要人工制定提取特征的策略和不断优化特征规则。深度学习方法在问题分类上具有可行性,通过自我学习特征的方式表示和理解问题,避免人工特征的制定,从而减少人工代价。针对问题分类,改进了长短期记忆神经网络(LSTM)和卷积神经网络(CNN)模型,并结合两者的优势组合成为一种新的学习框架(LSTM-MFCNN),加强对词序语义和深度特征的学习。实验结果表明,该方法在不需要制定繁琐的特征规则的前提下,仍然有较好的表现,准确率达到了 93.08%。

关键词 问题分类,深度学习,卷积神经网络,LSTM,机器学习

中图分类号 TP391 **文献标识码** A **DOI** 10.11896/j.issn.1002-137X.2016.12.020

Research on Problem Classification Method Based on Deep Learning

LI Chao CHAI Yu-mei NAN Xiao-fei GAO Ming-lei
(School of Information Engineering, Zhengzhou University, Zhengzhou 450001, China)

Abstract Question classification is an important part of question answering system. But question classification requires the strategy of extracting features and the continuous optimization of characteristic rules at the present stage. The method of deep learning is feasible in the question classification by the way of self learning question characteristics to represent and understand the problem so as to avoid formulating artificial features and reduce labor costs. For question classification, the long-short term memory(LSTM) model and the convolution neural network (CNN) model were improved, combining the advantages of these two models into a new learning framework (LSTM-MFCNN) to strengthen the semantic study of word order and study of depth characteristics. Experimental results show that the proposed method still has good performance under the condition of no need to formulate the characteristic rules, and the accuracy of this method is 93.08%.

Keywords Question classification, Deep learning, CNN, LSTM, Machine learning

1 引言

随着互联网的发展,人们可以通过搜索引擎从互联网中检索到越来越多的信息,传统搜索引擎只是返回大量的相关网页,而返回的相关信息越多,找到所需要的关键信息就会越困难。问答系统是允许用户以自然语言形式进行提问,并能够从大量的互联网信息中为用户返回更加简洁准确的信息,其查询方式能够更好地满足用户快速、准确地获取信息的需求。目前,问答系统(QA)是自然语言处理领域的一个研究热点^[1,2]。

问答系统主要分为 3 部分:问题分析、信息检索和答案抽取。问题分析的核心是问题分类,其性能会直接影响后期答案抽取的准确性^[3],主要体现在以下两个方面:

(1) 减少了候选答案的选择空间。例如,问题“第一个飞跃大西洋的人是谁”,是一个问人名的问题,问题答案的类型应与问题的类型一致,答案抽取范围也应限制在人名的范围内。

(2) 针对不同的问题类型,问答系统后续流程可以制定不同的策略。例如,问题“红茶和绿茶有什么不同”,问题的答案是描述红茶和绿茶的不同点。答案抽取应选取针对描述类别的策略。

以上两个方面表明,不同的问题类型,答案抽取的范围以及抽取方法也是不同的,问题分类的好坏会影响整个问答系统的性能。

早期的问题分类主要采用基于规则匹配的方法^[4],需要人工制定大量规则,建立规则库。近年来,问题分类方法主要是针对问题制定词法、句法、语义等特征提取策略^[5],再借助机器学习^[6-11](例如 K 近邻、SVM、贝叶斯等)的方法对问题进行分类。其分类结果的准确率在于特征提取的好坏,提取的特征越丰富,分类的精度越高。但与文本分类相比,问题分类中的问题是简短的句子,含有较少的特征信息,分类的难度也会更大。基于规则和特征提取的方法有以下 3 点不足:

(1) 人工制定的特征提取策略具有一定的主观性,不能全面理解问题。

(2) 为取得较好的分类效果,需要不断地调整、优化特征提取策略,以便在句法、语义方面更好地表示问题,灵活性不高。

(3) 问题的句法复杂度高或问题的类别粒度较小时,制定特征规则的难度增大,分类效果不好。

利用深度学习的方法可以借助大量语料,让模型主动学习到问题中潜在的句法和语义特征,更好地理解问题,有效弥

到稿日期:2016-02-02 返修日期:2016-05-10

李 超(1989-),男,硕士生,主要研究方向为机器学习、自然语言处理等,E-mail: struggle_more@163.com;柴玉梅(1964-),女,教授,主要研究方向为机器学习、自然语言处理等,E-mail: ieymchai@zzu.edu.cn(通信作者);南晓斐(1983-),女,博士,副教授,主要研究方向为机器学习、数据挖掘等,E-mail: iexfnan@zzu.edu.cn;高明磊(1963-),女,实验师,主要研究方向为数据分析、工控系统,E-mail: iemlgao@zzu.edu.cn。

补了人工提取特征在表示问题方面的不足,具有更好的灵活性、鲁棒性。

2 研究背景

2.1 传统问题分类方法

近年来问题分类主要是基于机器学习的方法,关键在于问题特征提取的好坏。张宇等人^[6]提取词频及词性的特征,通过改进贝叶斯模型对问题进行分类。在词频词性的基本特征之上,文勋等人^[7]利用句法结构来提取主干词和疑问词及其附属成分作为分类的特征。孙景广等人^[8]采用知网(How Net)作为语义资源提取分类特征。Silva 等人^[9]和 Ioni 等人^[10]使用特征组合的方式,采用线性核(linear SVM)对问题分类。此外,Li 等人^[11]提出将词性、词袋和句法依存树结合,并计算其核函数值的方法来探索问题的结构。上述方法主要根据问题分类数据集、人工提取问题的某种特征或者某些特征的组合来表示问题,具有一定的主观性,且语言表达的多样性^[12]使得手工制定较为准确的特征提取方法的代价变得更昂贵。

2.2 深度学习

深度学习是机器学习研究的新热点,其思想是模拟人脑的机制对问题进行分析学习。近年来,深度学习已在图像处理和语音识别领域成功应用,在自然语言处理领域,也有越来越多的学者利用深度学习的方法解决问题。索赫尔等人^[13]利用改进的递归自动编码模型处理情感分析问题中的语义合成的问题。崔等人^[14]将深度神经网络学习到的主题表达用于统计机器翻译消歧。Nal 等人^[15]为更加准确地理解句子,提出动态卷积神经网络(DCNN)对句子的语义建模。Li 等人^[16]将多柱卷积神经网络(multi-column)的方法用于处理基于知识库的问答系统问题。Zhang 等人^[17]在处理句子中的词与词之间相关分类的问题中,利用卷积深度信念网(DNN)学习词汇和句子层面的特征。对于问题分类,可以借助深度学习的思想主动学习隐含在句中的句法和语义特征^[18],深度分析问题结构。

2.3 问题分类体系

分类体系是问题分类的前提和基础。目前,中文问题分类没有统一的分类体系结构,哈尔滨工业大学信息检索和社会计算中心在国外研究的基础上,根据汉语独有的特点所提出的分类体系被大多数学者认可。其分类体系根据问题类别粒度的不同分为大类和小类,现已包括描述类(DES)、人物类(HUM)、地点类(LOC)、数字类(NUM)、时间类(TIME)、实体类(OBJ) 6 个大类,每个大类又划分为若干小类,共计 84 个小类,如表 1 所列。

表 1 哈工大中文问题分类体系

大类	小类
描述(DES)	简写、表达、意思、方式、原因、定义、判断、其它
人物(HUM)	特定人物、机构团体、人物描述、其他
地点(LOC)	宇宙、城市、大陆、国家、省、河流、湖泊、山脉、大洋、岛屿、建筑、地址、其它
数字(NUM)	温度、面积、体积、重量、速度、频率、距离、钱数、数量、顺序、倍数、百分比、号码、时间长度、范围、其它
时间(TIME)	年、月、日、季节、时代、星期、节气、节假日、时间、时间范围、其它
实体(OBJ)	物质、动物、植物、微生物、身体、材料、机具、衣物、食物药品、货币、票据、语言、事件、疾病、艺术品、服务、文字作品、学术学科、计划规划、法律法规、职位头衔、职业行业、符号、奖励、刑法、类别、权利义务、颜色、宗教、运动娱乐、术语、其它

3 基于深度学习的问题分类方法

本文用多种深度学习的方法对问题分类进行探索,试图发掘一种更加适合问题分类的深度学习框架。

地点类问题:“有一条江是黑龙江的最大支流,它的名字叫什么”。

人物类问题:“请问这个女孩的名字叫什么”。

简单的特征提取根据疑问词(叫什么)和其属性词(名字),会将问题判别为物人类。问题中的隐含的句法语义特征不可忽视。深度学习的方法可通过自我学习的方式学习到句子中内在的句法和语义特征,更具有准确性、客观性。卷积神经网络(CNN)利用卷积核可以学习到问题的内在特征;递归神经网络,特别是长短期记忆人工神经网络(LSTM)具有时序的特点,更加符合语言在听说读写中所表现出的自左到右的词序顺序。

本文采用卷积神经网络和长短期记忆人工神经网络围绕着问题分类来探索方法,针对问题分类的特点改进模型,并根据两种深度学习模型各自的优势将卷积神经网络和长短期记忆人工神经网络结合,在问题分类上表现出较好的成绩。

3.1 基于多粒度卷积核卷积神经网络(MFCNN)的问题分类

传统的卷积神经网络在同一隐藏层只存在单一粒度的卷积核,问题分类中的问题所含信息较少,因此提出用多粒度的卷积核挖掘更多隐含在问题中的特征。如图 1 所示,每一种粒度的卷积核都可学习到问题的一种特征表示,重组所有特征共同表示问题。模型是以词向量作为系统的输入,用多粒度卷积核来学习问题的不同表示,将学习到的问题特征作为 Softmax 分类器的输入。

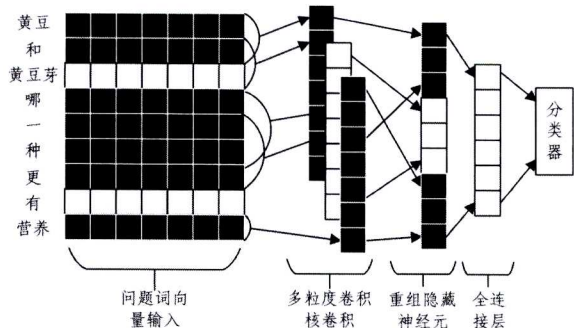


图 1 多粒度卷积核卷积神经网络(MFCNN)

MFCNN 模型描述和分类算法如下。

输入:问题 $Q=q_1, q_2, \dots, q_n$, 其中 q_i 表示问题的词向量。

数据预处理:截取问题为相同的长度,句子词较少的,用 0 填充。

创建具有 n 个隐藏单元, m 个不同粒度的卷积核的 MFCNN 模型:输入层经过 m 个粒度的卷积核映射到 m 个隐藏层,隐藏层经过池化层后通过含有 k 个隐藏单元的全连接层连接,最后将学习到问题的 k 个特征作为 Softmax 分类器的输入。

初始化所有模型参数为小的随机值。

在遇到终止条件前:

对于训练样例 training_examples 中的每个 (\vec{Q}, \vec{y}) :

1. 把输入沿网络前向传播:

隐藏神经元 h_i 并进行池化: $h_i = t(\text{pooling}(\sigma(wQ))) + b$

重组 m 个隐藏层的隐藏单元 $h; h \leftarrow h_1, h_2, \dots, h_m$

全连接层隐藏单元 $H_k; H_k = \sigma(wb + b)$

Softmax 层问题分类的概率计算: $P(y=i | w, b) = \frac{e^{w_i H + b_i}}{\sum_j e^{w_j H + b_j}}$

代价函数: $J(w, b) \leftarrow -\sum \log(P(y=y^i | Q^i, w, b) + \lambda |w|)$

2. 反向传播更新模型权值 w, b :

$$w \leftarrow w + \Delta w, b \leftarrow b + \Delta b$$

其中, $\Delta w = \frac{\partial J(w, b)}{\partial w}, \Delta b = \frac{\partial J(w, b)}{\partial b}$.

MFCNN 的算法描述中, 参数 w, b 为模型各层的权值和偏置, λ 是代价函数的正则参数, σ 是激活函数 sigmoid, t 是激活函数 tanh, pooling 为池化层的 max-pooling 操作。

3.2 基于长短期记忆神经网络(LSTM)的问题分类

长短期记忆神经网络(LSTM)可以对时间进行显式建模(见图2), 更加适合句子中的词序序列的语义表示, 符合自然语言的表示规律。LSTM 单元关键组成是记忆单元(见图3), 记忆单元是由3个控制门控制其所记忆的信息, 分别是输入门 i_t , 遗忘门 f_t 和输出门 o_t 。Graves 等人^[19]深入剖析了 LSTM 单元的结构及其性能表现。

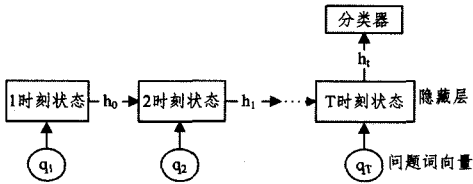


图2 长短期记忆神经网络(LSTM)

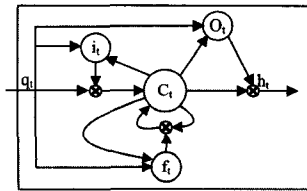


图3 LSTM单元组成结构

记忆单元可以记忆某时间段的信息, 对于问题分类可以记忆句子中前一时间词语的信息。例如上面提到的问题“有一条江是黑龙江的最大支流, 它的名字叫什么”, 其中词语“名字”所对应的隐藏单元的输入不止是来自“名字”本身的信息, 还有“名字”之前重要词语的信息, 比如词语“江”、“支流”等 LSTM 单元经过学习所认为重要的词语信息。长短期记忆神经网络学习问题中的词序相关, 进而更加贴切地表示问题的句法和语义特征。自然语言所提的问题具有表达多样性的特点, 本文梯度更新使用 AdaDelta 方法, 以便训练快速收敛; 并使用了 Dropout^[20] 思想, 以防止模型过拟合。

LSTM 模型描述和分类算法如下。

输入: 问题 $Q = q_1, q_2, \dots, q_T$, 其中 q_t 表示问题的词向量。

数据预处理: 便于并行矩阵计算, 输入的长度选定最长句子的长度, 句子词较少的, 用 0 填充。

创建具有 n 个 LSTM 单元的 LSTM 模型: t 时刻隐藏层中的 LSTM 单元的输入由词向量 q_t 和前一个 LSTM 单元的输出 h_{t-1} 组成, T 时刻隐藏层的输出作为 Softmax 分类器的输入。

初始化所有模型参数为小的随机值。

在遇到终止条件前:

对于训练样例 training_examples 中的每个 (\vec{Q}, \vec{y}) :

1. 把输入沿网络前向传播:

计算 t 时刻 LSTM 单元的 3 个控制门:

$$\text{输入门: } i_t = \sigma(w_i q_t + u_i h_{t-1} + b_i)$$

$$\text{遗忘门: } f_t = \sigma(w_f q_t + u_f h_{t-1} + b_f)$$

$$\text{输出门: } o_t = \sigma(w_o q_t + u_o h_{t-1} + b_o)$$

计算 t 时刻记忆单元 c_t : $c_t = i_t * \tanh(w_c q_t + u_c h_{t-1} + b_c) + f_t c_{t-1}$

t 时刻 LSTM 单元的输出 h_t : $h_t = o_t \tanh(c_t)$

Softmax 层问题分类的概率计算: $P(y=i | w, b) = \frac{e^{w_i h_T + b_i}}{\sum_j e^{w_j h_T + b_j}}$

代价函数: $J(w, b) \leftarrow -\sum \log(P(y=y^i | Q^i, w, b))$

2. 反向传播更新模型权值 w, b :

$$w \leftarrow w + \Delta w, b \leftarrow b + \Delta b$$

其中, $\Delta w = \frac{\partial J(w, b)}{\partial w}, \Delta b = \frac{\partial J(w, b)}{\partial b}$.

LSTM 模型的算法描述中, 参数 w, b 为模型权值和偏置, σ 是激活函数 sigmoid。

3.3 基于长短期记忆卷积神经网络(LSTM-MFCNN)的问题分类

在上述的讨论中, 两种改进模型对问题学习方法各有优势, 卷积神经网络通过利用多粒度卷积核对问题的词向量进行卷积操作, 更好地挖掘问题的特征, LSTM 可以更好地表示问题中词序序列的语义。因此借助两种学习框架对问题学习的优势, 将卷积神经网络和长短期记忆网络混合组成一种新的学习框架, 即先利用 LSTM 表示自然语言中的时序规律, 初步学习问题中的句法和语义特征, 再将其学习到的特征作为卷积神经网络的输入, 进一步挖掘问题中的深度特征, 最后用 Softmax 分类器对问题进行分类(见图4)。根据两种学习框架和数据集的特点, 对 3.1 节和 3.2 节描述的学习框架进行改进。首先为了提供给上层卷积神经网络更多的问题信息, 将下层 LSTM 模型所有时刻隐藏层的输出作为上层卷积神经网络的输入。由于所收集的数据集数量有限, 为避免训练数据的过拟合, 去掉了全连接层, 减少了上层卷积神经网络模型的复杂度。

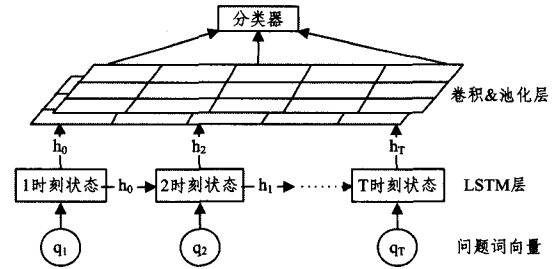


图4 长短期记忆卷积神经网络(LSTM-MFCNN)

如图4所示, 问题的词向量 (q_1, q_2, \dots, q_T) 作为整个系统的输入, 通过 LSTM 层学习问题的浅层句法语义特征, 隐藏层的输出 $(h_0 \dots h_T)$ 作为多粒度卷积神经网络的输入, 进一步学习问题的深度特征, 其输出的深度特征进入 Softmax 分类器, 对问题进行分类。

4 实验

4.1 词的向量表示

分布式表示最早由 Hinton^[21] 于 1986 年提出。其基本思想是通过训练将每个词映射成 K 维实数向量, 通过词之间的距离来判断它们之间的语义相似度。word2vec 使用的就是这种分布式表示的词向量表示方式, Google 在 2013 年发布开源的将词表征为实数值向量的高效工具 word2vec。本文用 word2vec 来训练词向量, 训练语料来自中文维基百科, 去除多余标签, 语料大小为 850 兆, 所训练词向量的维数为 150 维。

4.2 数据集

数据集由 3 部分组成, 哈尔滨工业大学问题分类数据集 6312 个、NLPCC 2015 QA 评测的问题集 888 个、复旦大学问

题分类数据集 2400 个,共 9600 个问题。其中训练集 7200 个,测试集 2400 个,数据集大类分布情况如表 2 所列。

表 2 训练语料和测试语料的问题分布

数据集	描述类	人物类	地点类	数字类	时间类	实体类	总计
训练集	687	1544	1234	1411	1281	1043	7200
测试集	252	476	451	460	368	393	2400

4.3 本文 3 种方法与基本模型的比较

为验证深度学习在问题分类上的有效性,本文采用 3 种评价标准:准确率(P)、召回率(R)和 F1 值(F1)。表 3 列出了

本文方法以及普通卷积神经网络在数据集以 6 个大类和 84 个小类分类的方式的准确率对比。表 4 给出了 6 大类问题 3 种评价的具体实验结果。表 3 中 B_P 代表大类准确率,S_P 代表小类准确率。

表 3 本文方法以大类和小类分类的方式的准确率对比(%)

方法	B_P	S_P
CNN	84.96	61.45
MFCNN	88.5	64.95
LSTM	91.05	80.7
LSTM-MFCNN	93.08	78.95

表 4 6 大类问题的 3 种评价实验结果(%)

方法	描述			人物			地点		
	P	R	F1	P	R	F1	P	R	F1
CNN	66.27	72.61	69.29	89.92	85.77	87.79	90.24	88.67	89.45
MFCNN	72.61	75.93	74.24	93.07	87.2	90.04	93.79	92.97	93.38
LSTM	74.6	82.1	78.17	92	93.88	92.93	95.81	94.71	95.26
LSTM-MFCNN	78.57	88	83.02	94.54	94.34	94.44	95.12	95.97	95.55
方法	数字			时间			实体		
	P	R	F1	P	R	F1	P	R	F1
CNN	89.35	91.13	90.23	91.03	84.81	87.81	74.04	79.51	76.68
MFCNN	91.74	94.83	93.26	92.66	87.89	90.21	79.39	85.95	82.54
LSTM	93.7	96.75	95.2	97.13	92.22	94.61	89.31	84.38	86.77
LSTM-MFCNN	96.09	97.36	96.72	97.55	94.23	95.86	90.59	85.58	88.01

表 3 表明,在问题分类上,多粒度卷积神经网络(MFCNN)更优于普通的卷积神经网络(CNN);与卷积神经网络相比,长短期记忆人工神经网络(LSTM)学习到的词序语义特征可以更好地对问题进行表示;长短期记忆卷积神经网络(LSTM-MFCNN)在以大类分类的方式上有较好的成绩,准确率为 93.08%,与 CNN 模型相比准确率提高了 8 个百分点以上,说明 LSTM-MFCNN 模型在学习到词序特征的同时,可对特征再次深度挖掘,更有利于对问题的理解。LSTM-MFCNN 模型在以小类分类的方式上比 LSTM 模型稍差,原因主要有两点:小类的类别较多,数据集相对较少;一些类别有较少的语料,类型的实例数目分布不均匀(例如 NUM_MULTIPLE 类有 3 条语料、TIME_HOLIDAY 有 4 条语料、HUM_PERSON 有 146 条语料),且 LSTM-MFCNN 模型复杂度相对较高。

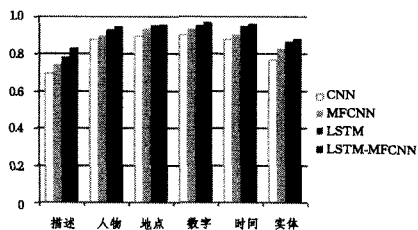


图 5 本文方法以大类分类的方式的 F1 值对比

图 5 展示了本文 3 种方法及 CNN 在 6 大类上的 F1 值分布,LSTM-MFCNN 方法在 6 大类都表现出较好的优势,较适合各种类型问题。表 4 显示,人物、地点、数字、时间 4 个大类 3 种评分方式的得分都很高,其中 LSTM-MFCNN 方法的 F1 值都在 94.44%以上,描述和实体类的 F1 值在 90%以下。原因主要有以下两点。

(1)描述类和实体类语料较少,其他类语料相对较多(见表 2)。更多的语料可以更好地训练模型,模型越复杂,需要的语料应越多。

(2)人物、地点、数字、时间 4 大类的特征相对较明显,例如时间类的问题中会含有“年”、“月”、“日”、“时候”、“何时”等较为明显特征的实体。

4.4 本文方法与现有同类工作的比较

中文问题分类数据集并没有统一的问题集,大多数学者是对哈尔滨工业大学问题集进行扩展,例如孙景广、田卫东使用了中国科学院自动化研究所模式识别国家重点实验室的问题集资源。表 5 列出了近年来在中文问题分类上的一些代表性工作。

表 5 中文分类问题工作比较(%)

方法	数据集	B_P	S_P
张宇等	TF-IDF、词频、词性+Bayes	4280	72.40 (65 类)
余正涛等 ^[22]	词性、语块、词义+SVM	1500	88.7 (6 类)
文勳等	句法结构、主干词、疑问词、附属+Bayes	6565	86.62 (7 类) 71.92 (60 类)
孙景广等	知网(How Net)、句法结构、疑问词+ME	5613	92.18 (7 类) 83.86 (60 类)
田卫东等 ^[23]	自学习方法、疑问词+中心词-类别+Bayes	4280	84 (6 类)
李茹等 ^[24]	疑问词、框架元素及中心词+ME	2011	91.38 (7 类) 83.20 (73 类)
本文方法	深度学习方法(LSTM-MFCNN)	9600	93.08 (6 类) 78.95 (84 类)

从表 5 可以看出:

(1)中文问题分类大都采用人工制定提取特征的策略,并采用机器学习的方法对问题进行分类。

(2)小类分类精度不高,主要原因是中文问题分类的语料相对较少,资源缺乏,制定精度较高的特征提取规则的难度增大。

(3)问题类别的划分逐年趋向于更多较小粒度的类别,较小粒度的类别有利于增加后续答案抽取的限制,但太小粒度的类别会增加问题分类的难度。

结束语 本文分析了问题分类的现状以及深度学习在自然语言处理领域的成功应用,提出用深度学习的方法来改善问题分类中存在的不足。本文尝试用多种深度学习的方法对问题分类进行应用与研究,试图发掘一种适合问题分类的深度学习框架。借助实验,提出混合长短时记忆网络和卷积神经网络的学习框架(LSTM-MFCNN),结合两者的优点,可以更好地理解问题。LSTM-MFCNN的学习方法在不需要制定繁琐的特征规则的前提下,准确率可以达到93.08%。但从现有的中文问题分类的工作来看,分类精度特别是小类的分类精度还需要进一步提高。对于语料不足的问题,除了收集和标记更多语料之外,还可以采用传统提取特征方法和深度学习结合的方式来进一步提高在现有语料基础之上的分类精度,这也将是下一步研究的重点。

参考文献

- [1] Hong L, Davison B D. A classification-based approach to question answering in discussion boards[C]// Proceedings of the 32nd International ACM SIGIR Conference on Research and Development in Information Retrieval. ACM, 2009; 171-178
- [2] Wang Y Z, Jia Y T, Liu D W, et al. Open Web Knowledge Aided Information Search and Data Mining[J]. Computer Research and Development, 2015, 52(2): 456-474 (in Chinese)
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] Mudgal R, Madaan R, Sharma A K, et al. A Novel architecture for question classification based indexing scheme for efficient question answering[J]. International Journal of Computer Engineering, 2013, 2(2): 27-43
- [4] Li X, Roth D. Learning question classifiers: the role of semantic information[J]. Natural Language Engineering, 2006, 12(3): 229-249
- [5] Mishra M, Mishra V K, Sharma H R. Question Classification using Semantic, Syntactic and Lexical features[J]. International Journal of Web & Semantic Technology, 2013, 4(3): 39-47
- [6] Zhang Y, Liu T, Wen X. Modified Bayesian Model Based Question Classification[J]. Journal of Chinese Information Processing, 2005, 19(2): 100-105 (in Chinese)
张宇, 刘挺, 文勳. 基于改进贝叶斯模型的问题分类[J]. 中文信息学报, 2005, 19(2): 100-105
- [7] Wen X, Zhang Y, Liu T, et al. Syntactic Structure Parsing Based Chinese Question Classification[J]. Journal of Chinese Information Processing, 2006, 20(2): 33-39 (in Chinese)
文勳, 张宇, 刘挺, 等. 基于句法结构分析的中文问题分类[J]. 中文信息学报, 2006, 20(2): 33-39
- [8] Sun J G, Cai D F, LV D X, et al. HowNet based Chinese question automatic classification[J]. Journal of Chinese Information Processing, 2007, 21(1): 90-95 (in Chinese)
孙景广, 蔡东风, 吕德新, 等. 基于知网的中文问题自动分类[J]. 中文信息学报, 2007, 21(1): 90-95
- [9] Silva J, Coheur L, Mendes A C, et al. From symbolic to sub-symbolic information in question classification[J]. Artificial Intelligence Review, 2011, 35(2): 137-154
- [10] Loni B, Van Tulder G, Wiggers P, et al. Question classification by weighted combination of lexical, syntactic and semantic features[C]// Text, Speech and Dialogue. Springer Berlin Heidelberg, 2011; 243-250
- [11] Liu L, Yu Z, Guo J, et al. Chinese question classification based on question property kernel[J]. International Journal of Machine Learning and Cybernetics, 2014, 5(5): 713-720
- [12] Yadav R, Mishra M, Bhilai S. Question Classification Using Naive Bayes Machine Learning Approach[J]. International Journal of Engineering and Innovative Technology (IJEIT), 2013, 2(8): 291-294
- [13] Socher R, Pennington J, Huang E H, et al. Semi-supervised recursive auto encoders for predicting sentiment distributions[C]// Proceedings of the Conference on Empirical Methods in Natural Language Processing. Association for Computational Linguistics, 2011; 151-161
- [14] Cui L, Zhang D, Liu S, et al. Learning topic representation for smt with neural networks[C]// Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics. 2014; 133-143
- [15] Blunsom P, Grefenstette E, Nal K, et al. A convolutional neural network for modelling sentences[C]// Proceeding for the 52nd Annual Meeting of the Association for Computational Linguistics. 2014; 655-665
- [16] Dong L, Wei F, Zhou M, et al. Question answering over freebase with multi-column convolutional neural networks[C]// Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing. 2015; 260-269
- [17] Zhang D, Wang D. Relation Classification via Recurrent Neural Network[J]. arXiv preprint arXiv: 1508. 01006, 2015
- [18] Kim Y. Convolutional neural networks for sentence classification [C] // Empirical Methods in Natural Language Processing. 2014; 1746-1751
- [19] Graves A. Generating sequences with recurrent neural networks [J]. arXiv preprint arXiv: 1308. 0850, 2013
- [20] Hinton G E, Srivastava N, Krizhevsky A, et al. Improving neural networks by preventing coadaptation of feature detectors[J]. arXiv preprint arXiv: 1207. 0580, 2012
- [21] Hinton G E. Learning distributed representations of concepts [C]// Proceedings of the eighth annual conference of the cognitive science society. 1986
- [22] Yu Zheng-tao, Fan Xiao-zhong, Guo Jian-yi. Chinese Question Classification Based on Support Vector Machine[J]. Journal of South China University of Technology, 2005, 33(9): 25-29 (in Chinese)
余正涛, 樊孝忠, 郭剑毅. 基于支持向量机的汉语问句分类[J]. 华南理工大学学报, 2005, 33(9): 25-29
- [23] Tian W D, Gao Y Y, Zu Y L. Question classification based on self-learning rules and modified Bayes[J]. Application Research of Computers, 2010, 27(8): 2869-2871 (in Chinese)
田卫东, 高艳影, 祖永亮. 基于自学习规则和改进贝叶斯结合的问题分类[J]. 计算机应用研究, 2010, 27(8): 2869-2871
- [24] Li R, Song X X, Wang W J. Chinese question classification based on Chinese FrameNet[J]. Computer Engineering and Applications 2009, 45(31): 111-114 (in Chinese)
李茹, 宋小香, 王文晶. 基于汉语框架网的中文问题分类[J]. 计算机工程与应用, 2009, 45(31): 111-114