

DOI: 10.13718/j.cnki.xdzk.2023.06.020

# 基于特征增强的多方位农业问句语义匹配

王奥<sup>1,2,3</sup>, 吴华瑞<sup>2,3,4</sup>, 朱华吉<sup>2,3,4</sup>

- 广西大学 计算机与电子信息学院, 南宁 530004; 2. 北京市农林科学院 信息技术研究中心, 北京 100097;
3. 国家农业信息化工程技术研究中心, 北京 100097; 4. 农业农村部 数字乡村技术重点实验室, 北京 100097

**摘要:** 农业问句文本数据具有专业名词多、特征稀疏、语句规范性差等特征, 难以深入挖掘句间交互关系。为改善农业相似问句的匹配性能, 提出一种基于特征增强的多方位农业问句语义匹配模型。模型通过共享参数的双向长短期记忆网络提取上下文向量, 分别引入自注意力机制、多维注意力机制增强农业问句文本语义推断特征和文本距离特征, 通过多特征增强聚焦语义特征信息, 将增强特征嵌入到多方位匹配函数中, 从向量值、方向和元素等角度进行句间相似度对比, 以捕获句子多样性特征。从农业问答社区导出农业问答文本数据, 人工标注相似问句构建试验数据集。试验结果表明: 基于特征增强的多方位农业问句语义匹配模型可以增强文本特征之间的交互, 获取更多的关系特征信息, 在构建的农业问句数据集上正确率及 F1 值达 95.3% 和 97.3%, 与其他 5 种问句语义匹配模型相比, 效果提升明显。

**关键词:** 农业问句语义匹配; 特征增强; 自然语言处理;

双向长短期记忆网络; 自注意力机制

中图分类号: TP391.1

文献标志码: A

开放科学(资源服务)标识码(OSID):



文章编号: 1673-9868(2023)06-0201-10

## Multi-Level Semantic Matching of Agricultural Questions Based on Feature Enhancement

WANG Ao<sup>1,2,3</sup>, WU Huarui<sup>2,3,4</sup>, ZHU Huaji<sup>2,3,4</sup>

- School of Computer, Electronics and Information, Guangxi University, Nanning 530004, China;
- Information Technology Research Center, Beijing Academy of Agriculture and Forestry Sciences, Beijing 100097, China;
- National Engineering Research Center for Information Technology in Agriculture, Beijing 100097, China;
- Key Laboratory of Digital Village Technology, Ministry of Agriculture and Rural Affairs, Beijing 100097, China

**Abstract:** To improve the performance of similarity calculation in agricultural Q & A community, accord-

收稿日期: 2023-02-07

基金项目: 科技创新 2030——“新一代人工智能”重大项目(2021ZD0113605); 国家重点研发计划项目(2019YFD1101105, 2020YFD1100602)。

作者简介: 王奥, 硕士研究生, 主要从事农业智能问答研究。

通信作者: 朱华吉, 博士, 研究员。

ing to the characteristics of agricultural question which are many professional nouns, sparse and poor sentence standardization, a semantic matching model of agricultural question sentences based on features enhancement was proposed. The model extracts context vectors through a bidirectional long-short term memory network that shares parameters. The self-attention mechanism and multi-dimensional attention mechanism are used to enhance the semantic inference features and distance features of agricultural question text data, respectively. Through multi-feature enhancement, the semantic feature information is focused, the enhanced features are embedded in the multi-directional matching function, and the similarity is compared from the perspectives of vector value, direction and element to capture the diversity characteristics of sentences. Agricultural Q & A text data is exported from the agricultural Q & A community, and similar questions are manually labelled to construct experimental datasets. The experimental results showed that the agricultural question semantic matching model based on enhanced multi-feature can enhance the interaction between text features, get more relationship feature information. The accuracy and F1 values of the proposed model were 95.3% and 97.3%. Compared with the other five semantic matching models, the experimental results showed obvious advantages.

**Key words:** agricultural question semantic matching; feature enhancement; natural language processing; bi-long-short term memory network; self-attention

农业复杂交互式问答平台为农户提供专家在线指导、在线学习、农业技术交流多种功能<sup>[1-2]</sup>,在协助用户解决农业生产生活和日常信息需求中发挥着重要作用。平台农户和专家实时在线互动,问答文本海量增长,但经常出现不同表达方式表达相同语义的情况,相似问题解答消耗大量人力、物力,因此构建能够快速准确给出答案的问答系统就显得十分必要。相似度匹配是语音、人脸识别<sup>[3]</sup>、问答等系统的基础任务,其相似度计算的精度直接影响问答系统回复的准确率,利用问句相似度匹配<sup>[4]</sup>开展高精度的农业智能问答模型研究,是农业智能化的重要发展方向。

以往的语义匹配研究集中在短语、语法和词汇匹配,如文献[5]提出一种语法驱动文本匹配方法,通过融合具有鲁棒性的非词汇语法和由对数驱动的词汇语法的线性模型进行文本匹配。随着深度学习的蓬勃发展<sup>[6-8]</sup>,语义匹配从基础的文本嵌入到相似度计算,再到复杂的神经网络,有效解决了人工设计特征提取量少、泛化性差的问题。卜维琼等<sup>[9]</sup>针对农业领域特征,提出一种多重信息融合的相似度算法,首次将深度学习与农业问句匹配结合。孪生神经网络在文本匹配领域表现出良好的性能<sup>[10]</sup>。刘志超等<sup>[11]</sup>采用孪生神经网络架构,结合双向长短期神经网络和卷积神经网络进行水稻问句语义匹配。这种网络结构减少训练模型参数,提高了训练效率。金宁等<sup>[12]</sup>采用孪生神经网络结构,运用双向长短期记忆网络、卷积神经网络和密集连接网络从深度语义、词语共现、最大匹配度 3 个层面实现农业短文本匹配,但是直接进行句子表示的相似度匹配,忽略了句间交互,导致交互特征信息的损失,无法有效学习句子关系特征。

注意力机制<sup>[13]</sup>可有效解决上述问题,利用注意力机制对特征信息进行聚合或增强匹配信息,挖掘丰富的句子关联信息<sup>[14-16]</sup>。融入注意力机制的交互模型通过赋予词不同的权重,能快速获得有效信息,有效提升文本匹配模型性能,文献[17]针对农业文本特征,利用基于协同注意力机制的紧密连接 BiGRU(双向门控循环单元)实现农业问句相似度匹配。在注意力机制基础上从字、词、句的角度研究文本相似度计算<sup>[18-21]</sup>,细粒度对比句子差异能够提高相似度计算的效率和准确率。但农业文本数据存在词汇总量较少、专有名词多,具有冗余性、稀疏性、规范性差等特点,导致传统语义匹配方法提取句子间关联特征信息不够充分,忽略了句间推理关系。如何实现农业相似问句语义智能检索仍是农业问答需要解

决的一个重要问题。

针对农业文本句子关联特征信息难以深入挖掘, 句子多样性捕获不足等问题, 构建双向长短期循环神经网络提取特征, 融合自注意力机制、多维注意力机制增强的文本语义推断特征和距离特征, 通过多特征增强聚焦语义特征, 将增强特征嵌入多方位匹配层, 多角度对比句子特征信息, 捕获句子的多样性, 以期实现农业问句精准、自动的语义匹配。

## 1 特征增强语义匹配模型

如图 1 所示, 农业文本专业名词多、规范性差和高度依赖上下文等特点导致句子交互信息提取不足, 句间关系推理不够深入; 本文构建适用于农业问句文本的特征增强文本匹配模型, 由特征提取层、特征增强层、多方位匹配层构成。特征增强层利用自注意力机制和多维注意力机制提取不同粒度的局部特征, 获取具有丰富语义的交互向量特征。将两种增强特征信息嵌入多方位匹配函数中, 由 3 种匹配函数实现文本特征的多角度对比。作为问答的基础任务, 相似度匹配精度的提升能有效提高问答系统的答案返回效率和准确率。利用 BILSTM(双向长短期记忆网络)提取农业文本输入上下文向量, 获得农业问句文本前后关联语义, 设置 2 层 BILSTM 网络, 每层 LSTM 的隐藏神经单元为 128 个。

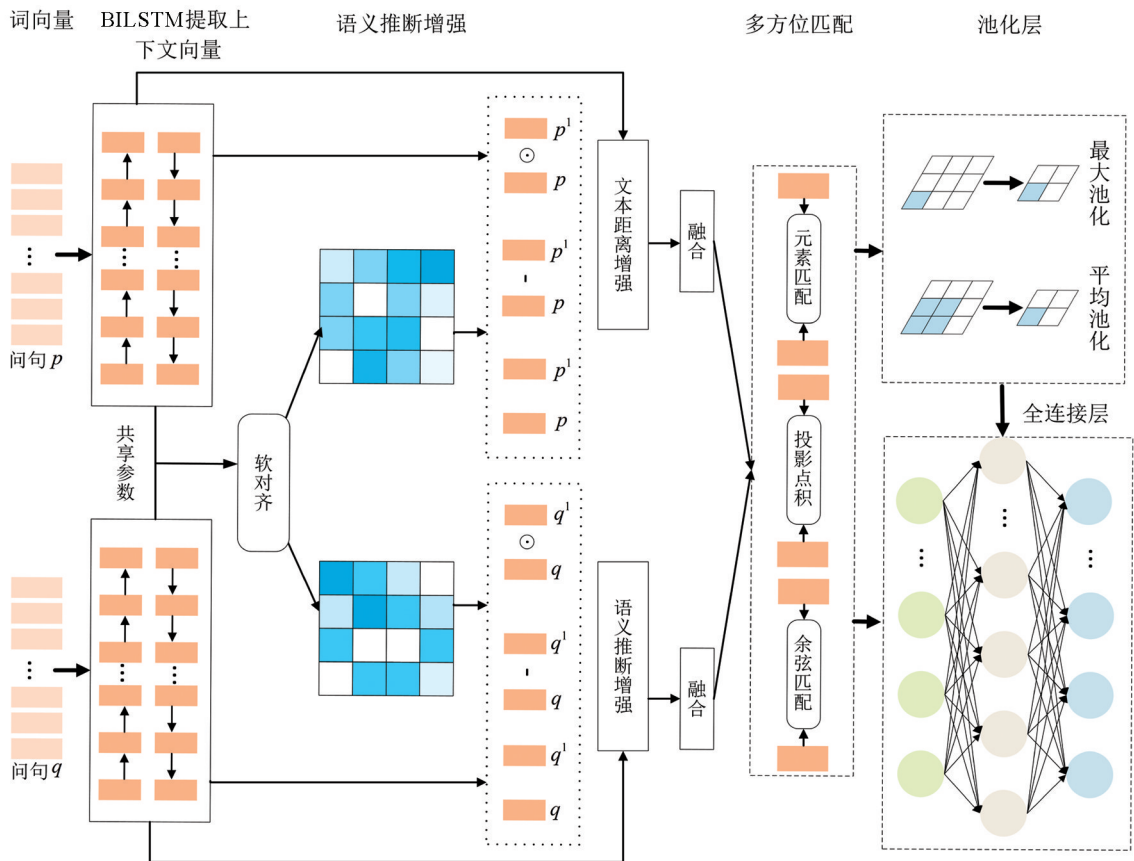


图 1 特征增强语义匹配模型架构图

### 1.1 特征增强层

传统文本匹配模型获取文本特征后直接进行相似度对比, 缺少关联特征信息或挖掘不够深入。农户提问问句存在文本数据专业名词多、规范性差等特点, 文本关系特征信息难以挖掘。利用自注意力机制、多维注意力机制分别增强语义推理特征和文本距离特征, 准确聚焦语义特征, 合理建模上下文信息, 提高问答匹配精度。

### 1.1.1 语义推断特征增强

通过自注意力机制计算注意力权重, 获得文本向量间的对齐关系. 公式(1)为权重计算公式, 作为隐藏状态的相似性矩阵.

$$k_{ij} = \frac{\overline{p_i^T q_j}}{\sqrt{d_H}} \quad (1)$$

式中 $\overline{p_i}$ 表示句子 $p$ 的第 $i$ 个时间步的隐藏状态, $\overline{q_j}$ 同理; $d_H$ 为缩放因子.

通过上式的注意权重获取向量间的局部相关性, 隐藏状态 $\overline{p_i}$ 在 $\overline{q_j}$ 中的相关语义在 $k_{ij}$ 中表示. 两个句子之间的相互关联和组合表达由下式得到:

$$\widetilde{p}_i = \sum_{j=1}^{l_q} \frac{e^{k_{ij}} \overline{q_j}}{\sum_{h=1}^{l_p} e^{k_{ih}}} \overline{q_j}, \forall i \in [1, \dots, l_p] \quad (2)$$

$$\widetilde{q}_j = \sum_{i=1}^{l_p} \frac{e^{k_{ij}} \overline{p_i}}{\sum_{h=1}^{l_q} e^{k_{jh}}} \overline{p_i}, \forall j \in [1, \dots, l_q] \quad (3)$$

式中, $\widetilde{q}_j$ 是 $\overline{p_j}$ 的权重和, 即 $\overline{q_j}$ 中与 $\overline{p_i}$ 相关的内容被抽出, 用 $\widetilde{p}_i$ 表示, $\widetilde{q}_j$ 同理.

农业问句文本中存在“蕴含”“矛盾”“中性”等多种关系, “蕴含”指句子 $p$ 能推断出句子 $q$ , 例如“玉米栽培技术要点是什么”推断出“我想学习玉米栽培技术”. “矛盾”指句子 $p$ 能推断出句子 $q$ 的否定, 例如“红薯的主要病害有哪些”与“红薯的形态特征是什么”互为否定. “中性”为其他所有情况. 对于农业文本数据的稀疏性, 传统自然语言匹配方法很难推理得到文本的多种关系和相关语义. 通过 $\overline{p_i}$ 和 $\widetilde{q}_j$ ,  $\overline{p_j}$ 和 $\widetilde{q}_i$ 之间的向量运算来锐化局部推理特征信息, 捕获局部推理过程中比较明显或突出的特征信息, 或者获取矛盾关系的推理信息. 最后将差值向量、原始隐藏状态向量和句子间的关联表示拼接起来, 得到增强语义推断特征.

$$r_p = [\overline{p_i}; \widetilde{p}_i; |\overline{p_i} - \widetilde{p}_i|; \overline{p_i} - \widetilde{p}_i] \quad (4)$$

$$r_q = [\overline{q_j}; \widetilde{q}_j; |\overline{q_j} - \widetilde{q}_j|; \overline{q_j} - \widetilde{q}_j] \quad (5)$$

### 1.1.2 文本距离特征增强

农业特定领域中, 农业专业术语比常用词承载了更多信息, 可作为问句中关键词. 循环神经网络提取特征时, 直接提取句子的每个词向量, 忽略了关键词在句子语义表示中的重要作用. 引入多维自注意力机制捕获每个词的上下文表示, 强调关键词的重要性, 增强句子中原始语义特征的提取. 多维自注意力将传统自注意力中的权重向量替换为权重矩阵, 使向量特征获得独自的权重. 注意力权重公式如下:

$$s(g_i, g_j) = \tanh(g_i \mathbf{W}_1 + g_j \mathbf{W}_2 + b) \quad (6)$$

式中, $g_i, g_j$ 为句中的隐藏状态, $\mathbf{W}_1$ 和 $\mathbf{W}_2$ 为可学习的权重矩阵, $b$ 是大小为隐藏节点个数的偏置.

农业问句文本中词间的距离能代表其相关性, 引入距离感知掩码使相近词获得更多关注, 距离更远的词关注更少. 计算相似度时词间距离越远所加负数越小, 经过 softmax 函数后, 距离越远的词权重越小, 词间的依赖也随之削弱. 在公式(6)中加上掩码 $\mathbf{M}$ ,  $\mathbf{M}$ 维度为 $1 \times 1$ , 矩阵中的值在 $\{0, -\infty\}$ 之间, 由此构建适用农业问句文本函数(7):

$$s(g_i, g_j) = \tanh((g_i \mathbf{W}_1 + g_j \mathbf{W}_2 + b)) + \mathbf{M}_{ij} \quad (7)$$

忽略 $\mathbf{M}_{ij}$ 为负无穷的情况,  $i = j$ 时,  $\mathbf{M}_{ij}$ 为负无穷, 即词向量与自身相比. 其余情况取决 $i$ 和 $j$ 的关系, 公式(8)为掩码矩阵的取值范围.  $f(i, j)$ 为 $i$ 和 $j$ 的距离函数,  $k$ 是超参数, 为正标量.  $|i - j| < k$ 时,  $\mathbf{M}_{ij}$ 为0, 表示当词间距离小于 $k$ 时,  $s(g_i, g_j)$ 为原始相似度; 当 $|i - j| \geq k$ 且 $i \neq j$ 时,  $\mathbf{M}_{ij}$ 为距离函数 $f(i, j)$ , 词间距离呈负相关. 通过文本距离特征增强获得句子表示 $\overline{p'}$ 和 $\overline{q'}$ . 文本距离增强如图2所示. 为防止问句过长导致过分削弱词之间的注意分数, 当句子长度小于 $t$ 时距离函数采用线性函数, 如

式(9). 句子长度大于  $t$  时, 采用对数函数, 如式(10).  $t$  为正标量的超参数.

$$M_{ij} = \begin{cases} -\infty & \text{不相关} \\ 0 & |i - j| < k \\ f(i, j) & i \neq j \text{ and } |i - j| \geq k \end{cases} \quad (8)$$

$$f(i, j) = -|i - j| \quad (9)$$

$$f(i, j) = -\exp |i - j| \quad (10)$$

将上述增强特征信息输入特征融合层, 融合增强局部推理特征表示和增强距离感知特征表示, 不仅增强语义特征, 而且保留了句子间的交互特征, 获得具有丰富语义特征信息的对齐特征向量.

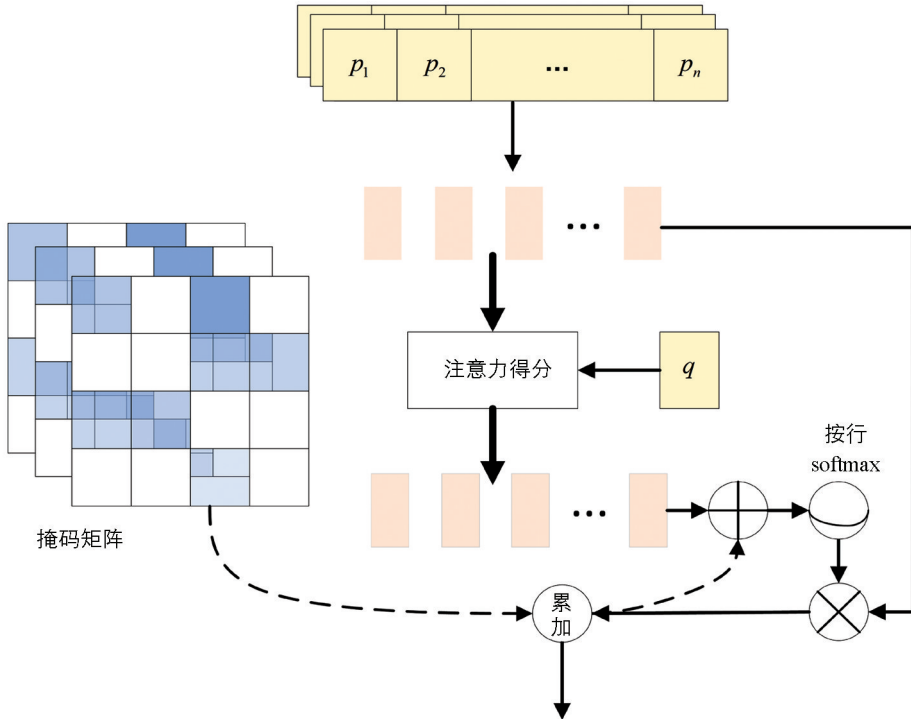


图 2 文本距离特征增强

### 1.2 多方位匹配

为了解决农业文本数据的稀疏性和文本词汇总量少导致的句子间关系信息获取不充分的问题, 使用 3 种匹配函数从不同角度获取更丰富的聚合信息和更准确的句子关系. 余弦相似匹配根据词频对比匹配相似程度, 对单字分组匹配时准确率高.

$$A = v_i^p - R \quad (11)$$

$$B = v_i^q - R \quad (12)$$

$$m_1 = \frac{\sum_{i=1}^l (A \times B)}{\sqrt{\sum_i A^2} \sqrt{\sum_i B^2}} \quad (13)$$

式中, 特征在所有维度上均减去均值  $R$ , 减少余弦匹配因仅进行向量对比的影响.

余弦相似度是对向量空间的度量, 忽略了排序和重叠词影响. 利用投影点积相似度匹配将向量进行投影, 通过点积乘法, 同时进行大小和角度的对比, 考虑整个句子值对相似度的影响. 其中  $W_p$  和  $W_q$  为可学习参数,  $\sigma$  为 sigmoid 函数.

$$m_2 = \sigma((v^p W_p)(v^q W_q)) \quad (14)$$

元素匹配则从元素角度比较向量异同. 词在句子中的重要程度不同, 其向量值也不同. 计算向量差

异和能更好地学习句间关系. 同时使用 3 种匹配函数, 从不同角度捕捉句子间的特征关联得到最终的匹配向量.

$$m_3 = \text{concat}(v^p, v^q, v^p * v^q, v^p - v^q, v^p + v^q) \quad (15)$$

3 种输出向量通过平均池化和最大池化聚合全局语义, 导入语义特征, 对最终匹配向量进行聚合. 输入到多层感知机(MLP)分类器, 使用 tanh 激活函数和 softmax 函数输出, 模型采用端对端的训练方式, 损失函数为交叉熵损失函数.

## 2 试验与分析

### 2.1 试验数据

通过农业领域最大的知识问答社区“中国农技推广信息平台”后台导出涉及 5 个种类的 20 000 个问答对来构建农业问句匹配数据集. 采用 jieba 分词工具加载停用词表, 剔除文本中的停用词、特殊符号等冗余信息. 人工筛选出信息不完整和无效问答的问句, 标注相似问句, 相同语义的问句占比为 54%, 不同语义的问句标注占比为 46%, 问答对包含病虫草害、土壤肥料、栽培管理、动物疫病、养殖管理等 5 类. 训练集和测试集比例为 8:2, 利用 Adam 优化器迭代更新神经网络权重, 采用准确率、精确率、召回率和 F1 值作为评价指标. 表 1 为训练集样本示例.

表 1 训练集样本示例

问句 1	问句 2	标签	类别
小青菜霜霉病病菌产生卵孢子的适宜温湿度是多少?	小青菜霜霉病病菌什么时候产生卵孢子?	0	病虫草害
大豆高产栽培技术是什么?	请问咋进行大豆高产栽培?	1	栽培管理
如何防治豆角炭疽病?	豆角炭疽病的发病条件是?	0	病虫草害
如何搞好莴笋病虫害防治?	怎么才能种植好莴笋?	0	栽培管理
玉米生长期需要施什么肥?	玉米生长期什么样的肥料好	1	土壤肥料
羊低镁血病的症状是什么?	羊低镁血病预防措施是什么?	0	动物疫病
水稻的田间管理技术要点有哪些?	如何进行水稻的田间管理?	1	栽培管理
肉牛养殖管理要点是啥?	肉牛夏季养殖如何防暑	0	养殖管理

模型训练迭代次数设置为 70, batchsize 为 110, BiLSTM 模型输出特征维度为 128 维, 全连接层隐藏单元设置为 128, 学习率设置为 0.001, 问句中有效词语使用 300 维的词向量表示, 句子最大长度为 20, 孪生网络共享参数. 为防止过拟合, 模型使用 dropout 函数, 随机使神经元失活, dropout 设置为 0.2.

### 2.2 试验结果与分析

掩码矩阵由超参数  $k$  来决定距离函数的取值, 词距离大于等于  $k$  时使用距离掩码限制注意力权重, 如表 2 所示, 为验证  $k$  值对模型性能的影响, 将  $k$  值分别设置为 0, 1, 2, 3, 4, 5, 6. 分别在农业文本数据和 lcqmc 数据集上进行试验对比,  $k$  为 0 时, 使用掩码矩阵限制注意力权重, 此时距离当前词较远的词注意力权重较低,  $k$  为 1 时相邻词会得到更多的注意力权重. 试验表明,  $k$  为 3 时性能较好,  $k$  继续增大, 对模型性能影响变差, 因此关注距离当前词 2 或 3 个词时获取信息更多. 在 lacqmc 数据集上准确率均在 90% 以上, 仍不及在农业数据上的表现, 说明模型具有一定的泛用性, 但更适合处理农业文本.

表 2  $K$  值对模型性能的影响

$k$ 值	0	1	2	3	4	5	6	%
农业文本	94.1	94.3	94.5	95.3	94.3	94.0	93.8	
lcqmc	91.1	91.4	91.6	91.9	91.7	91.5	90.7	

通过一组试验验证本文模型各个模块的有效性, 删除语义推断增强和文本距离增强得到模型 2 和模型 3. 由表 3 可知, 正确率和 F1 值下降了 1.7, 1.1 个百分点和 1.5, 0.9 个百分点, 表明单独的特征增强无法充分挖掘农业问句文本的交互信息. 同时删除两种特征增强策略得到模型 5, 可以看出两种策略融合更能提高模型的效果. 模型 6 为共享参数的 BILSTM 模型, 删除多角度匹配后正确度和 F1 值下降了 0.6, 0.5 个百分点, 因为单一角度的匹配无法获取足够的句子多样性. 表 4 为试验部分预测结果展示, 语义相同的问句标签记为 1, 反之标签记为 0, 预测与标签值相同时则为预测成功.

表 3 消融试验

序号	模型	正确率	精确率	召回率	F1 值	%
1	本文模型	95.3	97.1	97.5	97.3	
2	删除距离增强	93.8	95.7	97.1	96.4	
3	删除推断增强	93.6	96.9	95.1	96.2	
4	删除多方位匹配	94.6	97.0	96.7	96.8	
5	删除距离增强和推断增强	90.1	89.6	89.1	89.3	
6	Siamese-BILSTM	88.1	87.3	88.2	87.7	

表 4 部分预测结果

问句 1	问句 2	标签	预测
大棚茄子 6 月管理技术要点有哪些?	6 月份温室茄子的管理要点是什么?	1	1
大豆带状种植要点?	大豆高产栽培要点?	0	0
番茄晚疫病有什么症状?	如何防治番茄晚疫病?	0	0
肉牛养殖管理要点是啥?	肉牛夏季养殖如何防暑	1	1

图 3 为来自农业文本数据集的一个实例的注意力权重热力图, 问句 1 为“土豆早疫病有哪些症状表现?”, 问句 2 为“土豆早疫病发病原因是什么?”. 图 3a 是两个句子自注意力的对齐情况, 其中“土豆-土豆”, “早疫病-早疫病”, “症状-发病”有很强的对齐关系, 这些为句子关键词, 可明确表示句子语义, 通过捕捉两句话对齐关系可一定程度上判断词间关系. 图 3b 是同一问句注意力权重的可视化结果, 可以看出距离更近的词间注意力权重更大, 融合语义推断特征和文本距离特征进一步捕获句子的语义对齐信息, 获取丰富的交互信息, 提升语义匹配任务的性能.

本文与相似度匹配常用 5 个深度学习模型进行对比, ESIM<sup>[22]</sup> 使用 BILSTM 提取文本特征, 计算两个句子向量特征的相似度矩阵, 对向量特征加权, 再由一层 BILSTM 整合向量特征, 获得新的文本向量表示进行相似度匹配; DIIN<sup>[23]</sup> 是一种交互推理网络, 使用密集连接的卷积神经网络在交互空间中分层提取语义特征实现句子对的理解; ABCNN<sup>[24]</sup> 在 CNN 的基础上引入注意力机制, 在卷积计算和池化计算之前进行注意力权重计算, 判断文本相似情况; BIMPM<sup>[25]</sup> 在 BILSTM 提取文本特征后, 根据两句话不同的时间进行多角度对比; TextCNN<sup>[26]</sup> 通过不同大小的内核获取句子信息, 使用 CNN 完成句子的匹配和分类.

表 5 展示了 6 种模型针对农业问句数据集的试验结果, 本文模型在正确率、精确率、召回率、F1 值均超过了 95%, 较对比模型均有明显提升, 对比模型中, ESIM 模型 4 项指标均超过 91%. 本文模型 F1 值较其他模型提高接近 5 个百分点, 说明该模型能较为全面地捕捉文本间的交互信息, 相似度计算总体性能较好. 以卷积神经网络框架为基础的模型评价指标均低于以循环神经网络框架为基础的模型, 这是由网络结构所决定的, 卷积神经网络结构更擅长局部特征信息的提取, 并非文本序列化方向的特征提取, 且会丢失一些距离较远的文本特征向量. 5 种对比模型中 ESIM 模型召回率为 93.8%, 但仍与本文模型有些差距.

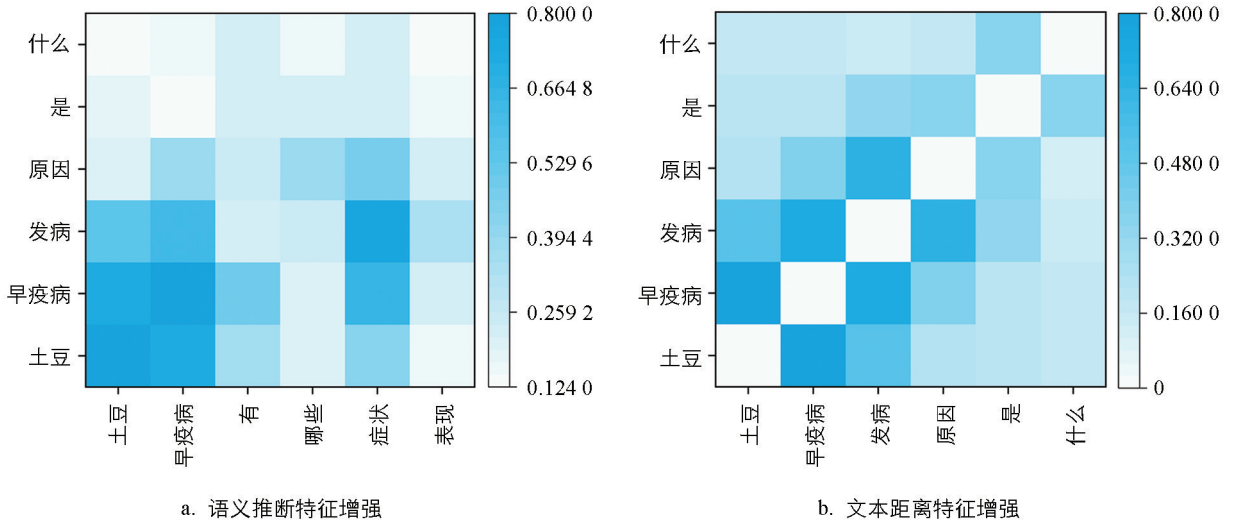


图 3 注意权重可视化图

表 5 不同模型对比结果

试验模型	正确率	精确率	召回率	F1 值
ESIM	91.7	91.1	93.8	92.5
DIIN	89.3	88.7	89.4	89.1
TextCNN	82.2	76.1	88.6	81.9
BIMPM	88.7	89.1	88.6	88.8
ABCNN	87.1	87.4	86.7	87.0
本文模型	95.3	97.1	97.4	97.3

如图 4, 与 ESIM, DIIN, ABCNN, BIMPM, TextCNN 5 种文本匹配模型相比, 本文模型在病虫害害、家畜疫病、栽培管理、养殖管理、土壤肥料 5 个类别的问句数据集上均有最高的匹配准确率, 整体匹配效

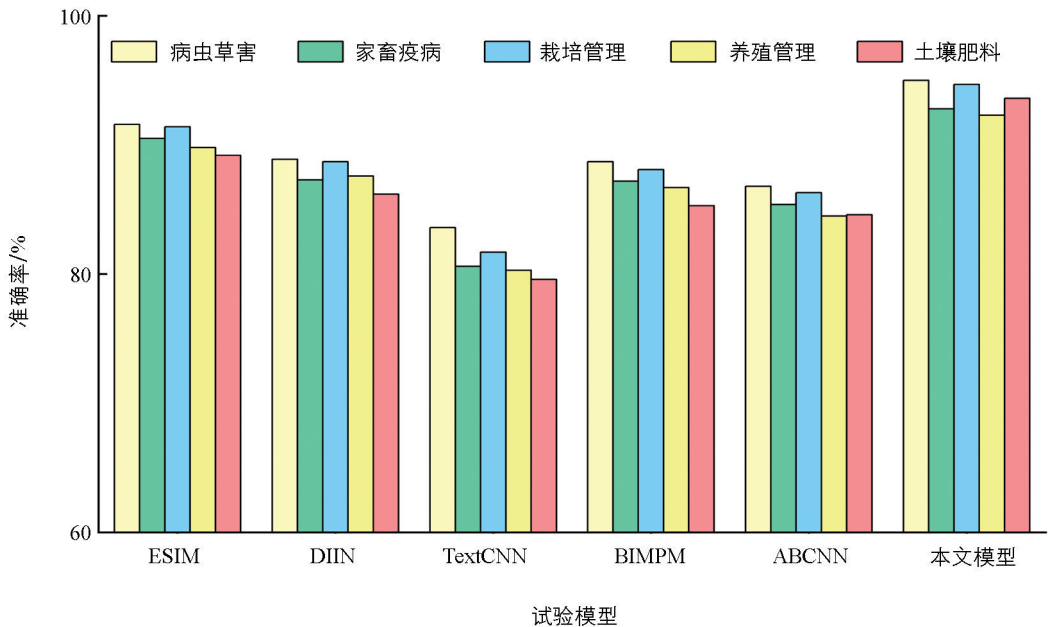


图 4 不同模型在农业问句数据集不同类别的准确率

果优于对比匹配模型。在病虫害害和栽培管理两个试验数据量充足的类别上准确率为95.0%和94.7%，因为数据集越充分，对深度学习模型迭代训练的效果提升越高。在养殖管理和土壤肥料两个数据量较少的类别中也高于其他模型的精确率，说明本文模型鲁棒性较强，在数据量不充足时也能有效提取文本特征进行相似度匹配。

### 3 结 语

为提高农户和农技工作者对农业问题检索的效率，减轻农业专家回复相似问题的压力及人工回复的延时性，构建了包含5个类别的农业问句语料库，提出一种基于多特征增强的农业问句语义匹配模型，在特征增强层增强语义推断特征和文本距离特征，深层次挖掘出农业文本交互特征信息，进一步获取丰富的文本间关联特征信息，由多方位匹配获取更丰富的聚合信息和句子关系。试验证明，在构建的农业问句数据集上较其他模型对语义匹配的计算性能有进一步提升，实现农业问句快速自动检测，有效提高农业智能问答中海量问句匹配效率和问答结果的准确率，进一步发挥智能问答在农技推广领域中的作用。由于农业具有地域性，在未来的工作中可考虑开展对方言问句和非规范的口语化问句语义匹配的相关研究。

#### 参考文献:

- [1] FENGSHI, JING. Knowledge-Enhanced Attentive Learning for Answer Selection in Community Question Answering Systems [J]. Knowledge-Based Systems, 2022, 250: 109117.
- [2] 马满福, 刘元喆, 李勇, 等. 基于 LCN 的医疗知识问答模型 [J]. 西南大学学报(自然科学版), 2020, 42(10): 25-36.
- [3] 施志刚. 基于改进协同表示的二级分类人脸识别方法 [J]. 西南大学学报(自然科学版), 2017, 39(1): 172-178.
- [4] LIU Y, TANG A H, SUN Z B, et al. An Integrated Retrieval Framework for Similar Questions: Word-Semantic Embedded Label Clustering-LDA with Question Life Cycle [J]. Information Sciences, 2020, 537: 227-245.
- [5] 王寒茹, 张仰森. 文本相似度计算研究进展综述 [J]. 北京信息科技大学学报(自然科学版), 2019, 34(1): 68-74.
- [6] WANG M, SMITH N A, TERUKO M. What is the Jeopardy Model? A Quasi-Synchronous Grammar for QA [C] // Proceedings of the 2007 Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning. Prague, Czech Republic: Association for Computational Linguistics Press, 2007: 22-32.
- [7] KALCHBRENNER N, GREFFENSTETTE E, BLUNSON P. A Convolutional Neural Network for Modelling Sentences [EB/OL]. 2014: arXiv: 1404. 2188. <https://arxiv.org/abs/1404.2188>.
- [8] GREFF K, SRIVASTAVA R K, KOUTNIK J, et al. LSTM: a Search Space Odyssey [J]. IEEE Transactions on Neural Networks and Learning Systems, 2017, 28(10): 2222-2232.
- [9] 卜伟琼, 方遼, 陈益能. 农业知识问答系统句子相似度算法研究 [J]. 农业网络信息, 2012(10): 17-20.
- [10] MUELLER J, THYAGARAJAN A. Siamese Recurrent Architecture for Learning Sentence Similarity [C] // Proceedings of the Thirtieth AAAI Conference on Artificial Intelligence. Phoenix Arizona, USA: AAAI Press, 2016: 2786-2792.
- [11] 刘志超, 王晓敏, 吴华瑞, 等. 基于 BiLSTM-CNN 的水稻问句相似度匹配方法研究 [J]. 中国农机化学报, 2022, 43(12): 125-132.
- [12] 金宁, 赵春江, 吴华瑞, 等. 基于多语义特征的农业短文本匹配技术 [J]. 农业机械学报, 2022, 53(5): 325-331.
- [13] VASWANI A, SHAZZER N, PARMAR N, et al. Attention is All You Need [C] // Proceedings of the 31st International Conference on Neural Information Processing Systems. New York, USA: Curran Associates Inc Press, 2017: 6000-6010.

- [14] TAN C Q, WEI F R, WANG W H, et al. Multiway Attention Networks for Modelling Sentence Pairs [C] //Proceedings of the 27th International Joint Conference on Artificial Intelligence. Stockholm Sweden: AAAI Press, 2018: 4411-4417.
- [15] KIM S, KANG I, KWAK N. Semantic Sentence Matching with Densely-Connected Recurrent and Co-Attentive Information [J]. Proceedings of the AAAI Conference on Artificial Intelligence, 2019, 33(1): 6586-6593.
- [16] LU W P, ZHANG X, LU H M, et al. Deep Hierarchical Encoding Model for Sentence Semantic Matching [J]. Journal of Visual Communication and Image Representation, 2020, 71: 102794.
- [17] 王郝日钦, 王晓敏, 缪祎晟, 等. 基于 BERT-Attention-DenseBiGRU 的农业问答社区问句相似度匹配 [J]. 农业机械学报, 2022, 53(1): 244-252.
- [18] 于碧辉, 王加存. 孪生网络中文语义匹配方法的研究 [J]. 小型微型计算机系统, 2021, 42(2): 231-234.
- [19] 冯月春, 陈惠娟. 改进 Bi-LSTM 的文本相似度计算方法 [J]. 计算机工程与设计, 2022, 43(5): 1397-1403.
- [20] 石彩霞, 李书琴, 刘斌. 多重检验加权融合的短文本相似度计算方法 [J]. 计算机工程, 2021, 47(2): 95-102.
- [21] 刘继明, 于敏敏, 袁野. 基于句向量的文本相似度计算方法 [J]. 科学技术与工程, 2020, 20(17): 6950-6955.
- [22] CHEN Q, ZHU X, LING Z, et al. Enhanced LSTM for Natural Language Inference [EB/OL]. (2017-04-26) [2023-02-27]. 2016: arXiv: 1609. 06038. <https://arxiv.org/abs/1609.06038>.
- [23] GONG Y, LUO H, ZHANG J. Natural Language Inference over Interaction Space [EB/OL]. (2017-09-13) [2023-02-07]. 2017: arXiv: 1709. 04348. <https://arxiv.org/abs/1709.04348>.
- [24] YIN W P, SCHÜTZE H, XIANG B, et al. ABCNN: Attention-Based Convolutional Neural Network for Modeling Sentence Pairs [J]. Transactions of the Association for Computational Linguistics, 2016, 4: 259-272.
- [25] WANG Z G, HAMZA W, FLORIAN R. Bilateral Multi-Perspective Matching for Natural Language Sentences [C] // Proceedings of the Twenty-Sixth International Joint Conference on Artificial Intelligence. August 19-26, 2017. Melbourne, Australia. California: International Joint Conferences on Artificial Intelligence Organization, 2017: 4144-4150.
- [26] ZHANG Y, WALLACE B. A Sensitivity Analysis of (and Practitioners' Guide to) Convolutional Neural Networks for Sentence Classification [C] // Proceedings of the Eighth International Joint Conference on Natural Language Processing. Taiwan, China: Asian Federation of Natural Language Processing, 2016: 253-263.

责任编辑 王新娟