



(12) 发明专利申请

(10) 申请公布号 CN 114416947 A

(43) 申请公布日 2022.04.29

(21) 申请号 202210047822.X

G06Q 30/02 (2012.01)

(22) 申请日 2022.01.17

G06N 3/04 (2006.01)

G06N 3/08 (2006.01)

(71) 申请人 中国科学技术大学

地址 230026 安徽省合肥市包河区金寨路
96号

(72) 发明人 陈恩红 刘淇 陈彦敏 王皓
黄振亚

(74) 专利代理机构 北京凯特来知识产权代理有
限公司 11260

代理人 郑立明 韩珂

(51) Int. Cl.

G06F 16/332 (2019.01)

G06F 40/289 (2020.01)

G06F 40/30 (2020.01)

G06F 16/35 (2019.01)

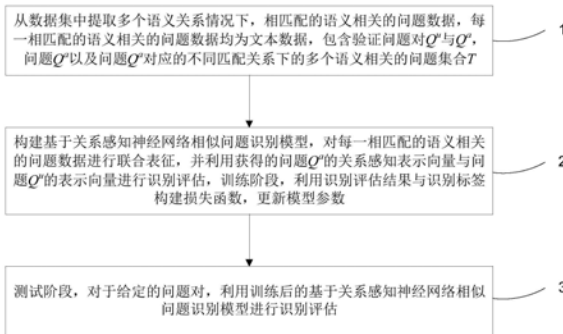
权利要求书3页 说明书8页 附图2页

(54) 发明名称

关系感知相似问题识别评估方法、系统、设备
及存储介质

(57) 摘要

本发明公开了一种关系感知相似问题识别
评估方法、系统、设备及存储介质,相关方案使用
基于关系感知神经网络相似问题识别模型来进行
问题对的相似问题识别,相比于传统模型,使用
了多个语义匹配的语义关系相关信息。对于预测
的结果,在多个评价指标上有一定的提高。



1. 一种关系感知的相似问题识别评估方法,其特征在於,包括:

从数据集中提取多个语义关系情况下,相匹配的语义相关的问题数据,每一相匹配的语义相关的问题数据均为文本数据,包含验证问题对 Q^u 与 Q^a ,以及问题 Q^a 对应的不同匹配关系下的多个语义相关的问题集合 T ;

构建基于关系感知神经网络相似问题识别模型,对每一相匹配的语义相关的问题数据进行联合表征,并利用获得的问题 Q^a 的关系感知表示向量与问题 Q^u 的表示向量进行识别评估,训练阶段,利用识别评估结果与识别标签构建损失函数,更新模型参数;

测试阶段,对于给定的问题对,利用训练后的基于关系感知神经网络相似问题识别模型进行识别评估。

2. 根据权利要求1所述的一种关系感知的相似问题识别评估方法,其特征在於,所述从数据集中提取多个语义关系情况下,相匹配的语义相关的问题数据包括:

从收集到的数据中提取具有匹配相关的问题,每一问题都与其他多个问题构成相似问题匹配的关系,所述相似问题匹配的关系包含两种类别,第一种分类标签为1,表示两个问题具有相同的语义,属于语义匹配关系;第二种分类标签为0,表示两个问题具有相同的关键词,但语义不同,属于语义相似关系;

将提取出的问题划分为两部分:第一部分作为匹配问题知识库;另一部分作为验证模型的训练数据;提取训练数据中的问题 Q^u 与匹配问题知识库中的单个问题 Q^a 构成一个验证问题对,对于问题 Q^a ,如果匹配问题知识库中的其他任意问题与 Q^a 的分类标签为1,则将放入集合 T^p ,如果分类标签为0,则放入集合 T^q ;集合 $T = \{T^p, T^q\}$ 即为问题 Q^a 对应的不同匹配关系下的多个语义相关的问题集合;

之后,对每一相匹配的语义相关的问题数据进行分词,获得验证问题对 Q^u 与 Q^a 的文本内容分词结果以及问题集合 T 中各问题的文本内容分词结果。

3. 根据权利要求1所述的一种关系感知的相似问题识别评估方法,其特征在於,所述基于关系感知神经网络相似问题识别模型包括:问题表示层、关系感知表示层和问题识别评估层;其中:

所述问题表示层,用于提取每一验证问题对 Q^u 与 Q^a 的表示向量,以及问题 Q^a 对应的不同匹配关系下的多个语义相关的问题集合中每一问题的表示向量;

所述关系感知表示层,用于结合问题 Q^a 与问题集合 T 中每一问题的语义关系,获得问题 Q^a 的关系感知表示向量;

所述问题识别评估层,用于利用问题 Q^u 与问题 Q^a 的表示向量,以及问题 Q^a 的关系感知表示向量,对验证问题对 Q^u 与 Q^a 的语义关系进行识别评估。

4. 根据权利要求3所述的一种关系感知的相似问题识别评估方法,其特征在於,所述问题表示层,对每个问题进行独立编码,生成问题的表示向量:

$$h = \text{Sentence-BERT}(Q)$$

其中, h 为问题 Q 的表示向量;Sentence-BERT表示Sentence BERT模型;

对于验证问题对 Q^u 与 Q^a 通过问题表示层获得相应的表示向量 h_u 与 h_a ;对于问题集合 $T = \{T^p, T^q\}$ 通过问题表示层获得语义匹配关系向量集 H^p 和语义相似关系的向量集 H^q ;其中,集合 T^p 与 T^q 中的问题与 Q^a 的分类标签分别为1与0,分类标签为1表示两个问题属于语义匹配关系,分类标签为0表示两个问题属于语义相似关系。

5. 根据权利要求3所述的一种关系感知的相似问题识别评估方法,其特征在於,所述结合问题 Q^a 与问题集合 T 中每一问题的语义关系,提取出问题 Q^a 的关系感知表示向量的步骤包括:

使用 R_N s关系表示网络表征问题 Q^a 与问题集合 T 中每一问题的关系表示:

$$H_r = \sum g_\theta(h_k^a)$$

其中, g_θ 是多层感知机网络, $h_k^a = h_a \odot h_k$ 表示采用element-wise product方法进行运算得到的 h_a 与 h_k 之间的交互特征, h_a 为问题 Q^a 的表示向量, h_k 表示问题集合 T 中单个问题的表示向量; $h_k \in H^p$ or H^q , H^p 、 H^q 分别为集合 T^p 、 T^q 中问题的表示向量集合,集合 T^p 中每一问题与问题 Q^a 的分类标签为1,表示两个问题属于语义匹配关系;集合 T^q 中每一问题与问题 Q^a 的分类标签为0,表示两个问题属于语义相似关系;

结合问题集合 T 中每一问题的标签,计算如下多标签下的匹配表示函数:

$$\tilde{H}_r = \sum_{r \in R} \sum_{o \in |N^r|} g_\theta(h_a, h_o^r)$$

其中, $h_o^r \in N^r$, N^r 为在第 r 个关系下多个语义匹配问题的语义向量集合; $|N^r|$ 表示语义向量集合 N^r 中的向量数目, R 表示标签的类别集合;

基于表示向量集合 H^p 与 H^q 将上述多标签下的匹配表示函数展开,得到:

$$\tilde{H}_r = \sum_{h_i^p \in N^p} g_\theta(W_{pi}(h_a, h_i^p) + b_{pi}) + \sum_{h_j^q \in N^q} g_\theta(W_{qj}(h_a, h_j^q) + b_{qj})$$

其中, W_{pi} 和 W_{qj} 分别是在两个关系下的学习权重, b_{pi} 和 b_{qj} 是偏差参数; $h_i^p \in H^p$, $h_j^q \in H^q$,各自为相应问题的表示向量;

采用element-wise product方法进行运算,分别得到的 h_a 与 h_i^p 之间、 h_a 与 h_j^q 之间的交互特征:

$$h_i^{ap} = h_a \odot h_i^p$$

$$h_j^{aq} = h_a \odot h_j^q$$

得到:

$$v_r = \sum_{h_i^{ap} \in N^p} g_\theta(W_{pi} h_i^{ap} + b_{pi}) + \sum_{h_j^{aq} \in N^q} g_\theta(W_{qj} h_j^{aq} + b_{qj})$$

其中, v_r 为问题 Q^a 的关系感知表示向量。

6. 根据权利要求3所述的一种关系感知的相似问题识别评估方法,其特征在於,所述利用问题 Q^u 与问题 Q^a 的表示向量,以及问题 Q^a 的关系感知表示向量,对验证问题对 Q^u 与 Q^a 的语义关系进行识别评估包括:

将问题 Q^u 与问题 Q^a 的表示向量,以及问题 Q^a 的关系感知表示向量连接,表示为:

$$z_{au} = h_u \oplus h_a \oplus v_r$$

其中, h_u 、 h_a 、 v_r 依次表示问题 Q^u 的表示向量、问题 Q^a 的表示向量、问题 Q^a 的关系感知表示向量;

通过RELU激活函数和sigmoid函数的运算获得验证问题对 Q^u 与 Q^a 的语义关系 $R(Q^u, Q^a)$,

表示为:

$$o_{au} = \text{ReLU}(W_1 z_{au} + b_1),$$

$$R(Q^u, Q^a) = \sigma(W_2 o_{au} + b_2)$$

其中, o_{au} 为RELU激活函数的运算结果, $\sigma(\cdot)$ 为sigmoid函数, W_1, W_2, b_1, b_2 为网络参数。

7. 根据权利要求1~6任一项所述的一种关系感知的相似问题识别评估方法, 其特征在于, 所述损失函数表示为:

$$L = -\frac{1}{N} \sum_{i=1}^N y_i \log P(y_i | R(Q_i^a, Q_i^u)) + \lambda_\theta \|\theta\|^2$$

其中, y_1 是对于第1个验证问题对的真实识别标签, N 是训练时验证问题对的数量; (Q_i^a, Q_i^u) 表示第1个验证问题对, $R(Q_i^a, Q_i^u)$ 表示1个验证问题对的识别评估结果, λ_θ 为正则化超参数, θ 为待更新的模型参数。

8. 一种关系感知的相似问题识别评估系统, 其特征在于, 用于实现权利要求1~7任一项所述的方法, 该系统包括:

数据提取单元, 用于从数据集中提取多个语义关系情况下, 相匹配的语义相关的问题数据, 每一相匹配的语义相关的问题数据均为文本数据, 包含验证问题对 Q^u 与 Q^a , 问题 Q^a 以及问题 Q^a 对应的不同匹配关系下的多个语义相关的问题集合 T ;

模型构建与训练单元, 用于构建基于关系感知神经网络相似问题识别模型, 对每一相匹配的语义相关的问题数据进行联合表征, 并利用获得的问题 Q^a 的关系感知表示向量与问题 Q^u 的表示向量进行识别评估, 训练阶段, 利用识别评估结果与识别标签构建损失函数, 更新模型参数;

识别评估测试单元, 用于对于给定的问题对, 利用训练后的基于关系感知神经网络相似问题识别模型进行识别评估。

9. 一种处理设备, 其特征在于, 包括: 一个或多个处理器; 存储器, 用于存储一个或多个程序;

其中, 当所述一个或多个程序被所述一个或多个处理器执行时, 使得所述一个或多个处理器实现如权利要求1~7任一项所述的方法。

10. 一种可读存储介质, 存储有计算机程序, 其特征在于, 当计算机程序被处理器执行时实现如权利要求1~7任一项所述的方法。

关系感知相似问题识别评估方法、系统、设备及存储介质

技术领域

[0001] 本发明涉及自然语言处理领域,尤其涉及一种关系感知相似问题识别评估方法、系统、设备及存储介质。

背景技术

[0002] 相似问题识别是智能客服问答研究领域的一个核心问题。当用户提出一个新问题,智能客服需要对用户的新问题进行理解,找到和用户问题可能匹配的相似问题,通过匹配模型得到最匹配的问题,反馈给用户相应的答案。因此用户相似问题识别任务也可以建模为相似文本匹配任务。相似问题识别任务应用在很多领域,例如在社区问答查询,信息检索和智能客户服务系统等,都将相似问题识别作为该应用领域的核心问题来研究。因此如何解决相似问题识别已成为一个非常重要的基础问题。

[0003] 围绕这个研究课题,研究者们提出了多种解决方法,大部分相关的研究主要集中在两个问题之间的词法、句法或问题结构,通过建模两个问题的语义关系,来判断两个问题的相似程度。

[0004] 然而由于问题长度较短和自然语言表达的灵活性和宽泛性,相似问题匹配面临着多样性的挑战。为了解决问题的多样性,部分的研究是引入外部知识例如知识图谱,问题答案等方法,用来解决多样性不足的情况。但是这些外部知识包含范围领域广泛或者针对性不强,并不能完全符合问题多样性表达的补充,因此,识别准确度还有待提升。

发明内容

[0005] 本发明的目的是提供一种关系感知相似问题识别评估方法、系统、设备及存储介质,可以充分利用多个问题对之间的语义关系信息来解决问题之间的相似问题识别,并具有较高的预测精度。

[0006] 本发明的目的是通过以下技术方案实现的:

[0007] 一种关系感知的相似问题识别评估方法,包括:

[0008] 从数据集中提取多个语义关系情况下,相匹配的语义相关的问题数据,每一相匹配的语义相关的问题数据均为文本数据,包含验证问题对 Q^u 与 Q^a ,以及问题 Q^a 对应的不同匹配关系下的多个语义相关的问题集合 T ;

[0009] 构建基于关系感知神经网络相似问题识别模型,对每一相匹配的语义相关的问题数据进行联合表征,并利用获得的问题 Q^a 的关系感知表示向量与问题 Q^u 的表示向量进行识别评估,训练阶段,利用识别评估结果与识别标签构建损失函数,更新模型参数;

[0010] 测试阶段,对于给定的问题对,利用训练后的基于关系感知神经网络相似问题识别模型进行识别评估。

[0011] 一种关系感知的相似问题识别评估系统,用于实现前述的方法,该系统包括:

[0012] 数据提取单元,用于从数据集中提取多个语义关系情况下,相匹配的语义相关的问题数据,每一相匹配的语义相关的问题数据均为文本数据,包含验证问题对 Q^u 与 Q^a ,问题

Q^a 以及问题 Q^a 对应的不同匹配关系下的多个语义相关的问题集合T;

[0013] 模型构建与训练单元,用于构建基于关系感知神经网络相似问题识别模型,对每一相匹配的语义相关的问题数据进行联合表征,并利用获得的问题 Q^a 的关系感知表示向量与问题 Q^u 的表示向量进行识别评估,训练阶段,利用识别评估结果与识别标签构建损失函数,更新模型参数;

[0014] 识别评估测试单元,用于对于给定的问题对,利用训练后的基于关系感知神经网络相似问题识别模型进行识别评估。

[0015] 一种处理设备,包括:一个或多个处理器;存储器,用于存储一个或多个程序;

[0016] 其中,当所述一个或多个程序被所述一个或多个处理器执行时,使得所述一个或多个处理器实现前述的方法。

[0017] 一种可读存储介质,存储有计算机程序,其特征在于,当计算机程序被处理器执行时实现前述的方法。

[0018] 由上述本发明提供的技术方案可以看出,使用基于关系感知神经网络相似问题识别模型来进行问题对的相似问题识别,相比于传统模型,使用了多个语义匹配的语义关系相关信息。对于预测的结果,在多个评价指标上有一定的提高。

附图说明

[0019] 为了更清楚地说明本发明实施例的技术方案,下面将对实施例描述中所需要使用的附图作简单地介绍,显而易见地,下面描述中的附图仅仅是本发明的一些实施例,对于本领域的普通技术人员来讲,在不付出创造性劳动的前提下,还可以根据这些附图获得其他附图。

[0020] 图1为本发明实施例提供了一种关系感知相似问题识别评估方法的流程图;

[0021] 图2为本发明实施例提供了一种关系感知相似问题识别评估系统的示意图;

[0022] 图3为本发明实施例提供了一种处理设备的示意图。

具体实施方式

[0023] 下面结合本发明实施例中的附图,对本发明实施例中的技术方案进行清楚、完整地描述,显然,所描述的实施例仅仅是本发明一部分实施例,而不是全部的实施例。基于本发明的实施例,本领域普通技术人员在没有做出创造性劳动前提下所获得的所有其他实施例,都属于本发明的保护范围。

[0024] 首先对本文中可能使用的术语进行如下说明。

[0025] 术语“包括”、“包含”、“含有”、“具有”或其它类似语义的描述,应被解释为非排它性的包括。例如:包括某技术特征要素(如原料、组分、成分、载体、剂型、材料、尺寸、零件、部件、机构、装置、步骤、工序、方法、反应条件、加工条件、参数、算法、信号、数据、产品或制品等),应被解释为不仅包括明确列出的某技术特征要素,还可以包括未明确列出的本领域公知的其它技术特征要素。

[0026] 其次,对已有的相似结果记录进行分析说明。

[0027] 经分析发现,现有数据集中存在如下情况,即,同一个问题和多个问题都有语义关系。例如表1所示,在某银行的客服问题数据集中存在这样的多匹配问题的语义关系。其中

Q_1 分别和4个问题(即表1中的 $Q_2 \sim Q_5$)存在语义关系。分类标签为1代表两个问题具有相同的语义,分类标签为0代表两个问题相似,但是语义不相同。具体而言, Q_2 和 Q_3 与 Q_1 的标签为1,分析问题内容可知, Q_2 和 Q_3 虽然表达的句式不一样,但是所表达的基本语义都是相同的,都是关于“电话审核时间”的信息,因此 Q_2 和 Q_3 是 Q_1 的多个语义匹配问题。 Q_4 和 Q_5 与 Q_1 的标签为0,分析 Q_4 和 Q_5 可知, Q_4 和 Q_5 这两个问题和 Q_1 有相同的关键词,但是表达的语义却不同,因此 Q_4 和 Q_5 是 Q_1 的多个语义相似但不匹配的问题。

问题	多匹配问题	分类标签
[0028] Q_1 : 电话审核时间	Q_2 : 晚上你们也会电话审核吗	1
	Q_3 : 现在借款大概什么时候能得到审核	1
Q_1 : 电话审核时间	Q_4 : 为什么说我没有通过审核呢	0
	Q_5 : 未通过审核标准是什么意思	0

[0029] 表1多个语义匹配问题示例

[0030] 因此,通过观察数据集可以发现,在基于标签分类的关系下,存在对 Q_1 的多个语义匹配的语义相关问题。以上两个不同的关系下的多个语义匹配情况,从更细粒度上反映了多个语义匹配关系对这个基本问题的语义影响。相比较于在之前的研究中匹配标签只是用来区分问题,实际上问题的匹配标签反映出问题对象之间具有的隐式语义关系。这些匹配关系具有很强的语义关联性,可以对问题的语义理解起到帮助作用。为了展现多个匹配语义关系在相似问题识别上所带来的积极作用,因此利用匹配语义关系来体现问题对之间的语义信息,是本发明方法研究的重点。

[0031] 下面对本发明所提供的一种关系感知的相似问题识别评估方法进行详细描述。本发明实施例中未作详细描述的内容属于本领域专业技术人员公知的现有技术。本发明实施例中未注明具体条件者,按照本领域常规条件或制造商建议的条件进行。本发明实施例中所用仪器未注明生产厂商者,均为可以通过市售购买获得的常规产品。

[0032] 如图1所示,一种关系感知的相似问题识别评估方法,主要包括如下步骤:

[0033] 步骤1、从数据集中提取多个语义关系情况下,相匹配的语义相关的问题数据,每一相匹配的语义相关的问题数据均为文本数据,包含验证问题对 Q^u 与 Q^a ,以及问题 Q^a 对应的不同匹配关系下的多个语义相关的问题集合T。

[0034] 本步骤的优选实施方式如下:

[0035] 1) 寻找存在多个语义匹配的问题。

[0036] 本发明实施例中,从收集到的数据中提取具有匹配相关的问题,每一问题都与其他多个问题构成相似问题匹配的关系,所述相似问题匹配的关系包含两种类别,第一种类别标签为1,表示两个问题具有相同的语义,属于语义匹配关系;第二种类别标签为0,表示两个问题具有相同的关键词,但语义不同,属于语义相似关系。

[0037] 2) 划分数据集。

[0038] 本发明实施例中,为了进一步划分出多个语义匹配问题集合和验证模型数据集,将已有的数据划分为两部分:第一部分作为匹配问题知识库;另一部分作为验证模型的训练数据;其中,匹配问题知识库提供 Q^a 与集合T,验证模型的训练数据,提供问题 Q^u ;在训练

时,验证数据从匹配问题知识库中找到对应问题的多个语义匹配问题集合,来实现对应问题的语义统一表示。

[0039] 具体来说:训练数据中的问题 Q^u 与匹配问题知识库中的单个问题 Q^a 进行匹配时,对于问题 Q^a ,如果匹配问题知识库 $M=(S_1, S_2, \dots)$ 中的其他任意问题与问题 Q^a 的分类标签为1,则将相应问题放入集合 $T^p = \{t_1^p, t_2^p, \dots, t_n^p\}$, T^p 表示问题 Q^a 的语义匹配关系数据集, n 为集合 T^p 中问题数目;如果分类标签为0,则放入集合 $T^q = \{t_1^q, t_2^q, \dots, t_m^q\}$, T^q 表示问题 Q^a 的语义相似关系数据集, m 为集合 T^q 中问题数目;集合 $T = \{T^p, T^q\}$ 即为问题 Q^a 对应的不同匹配关系下的多个语义相关的问题集合。

[0040] 3) 文本内容处理。

[0041] 本发明实施例中,对上述过程得到的每一个匹配的语义相关的问题数据进行分词,获得验证问题对 Q^u 与 Q^a 的文本内容分词结果,以及问题集合 T 中各问题的文本内容分词结果,以上两类分词结果都用作模型训练的输入,用于之后每个问题对的相似问题识别的评估。

[0042] 步骤2、构建基于关系感知神经网络相似问题识别模型,对每一相匹配的语义相关的问题数据进行联合表征,并利用获得的问题 Q^a 的关系感知表示向量与问题 Q^u 的表示向量进行识别评估,训练阶段,利用识别评估结果与识别标签构建损失函数,更新模型参数。

[0043] 本步骤的优选实施方式如下:

[0044] 1) 构建基于关系感知神经网络相似问题识别模型。

[0045] 本发明实施例中,基于关系感知神经网络相似问题识别模型主要包括:问题表示层、关系感知表示层和问题识别评估层;其中:a) 所述问题表示层,用于提取每一验证问题对 Q^u 与 Q^a 的表示向量,以及问题 Q^a 对应的不同匹配关系下的多个语义相关的问题集合中每一问题的表示向量;示例性的,可以处理阶段可以使用Sentence-BERT模型。b) 所述关系感知表示层,用于结合问题 Q^a 与问题集合 T 中每一问题的语义关系,获得问题 Q^a 的关系感知表示向量,具体来说,此阶段是将多个语义匹配关系信息整合到问题 Q^a 的表示向量中,得到关系感知的问题统一表示向量。c) 所述问题识别评估层,用于利用问题 Q^u 与问题 Q^a 的表示向量,以及问题 Q^a 的关系感知表示向量,对验证问题对 Q^u 与 Q^a 的语义关系进行识别评估,预测验证问题对的相似问题识别标签。

[0046] 下面针对各层的原理及相关过程进行介绍。

[0047] a) 问题表示层。

[0048] 所有的问题都需要输入至问题表示层,包括前文提到的问题 Q^u ,问题匹配知识库中问题 Q^a 和关于问题 Q^a 的问题集合 $T = \{T^p, T^q\}$ 。在问题表示层,通过已经预训练初始化单词向量将问题表示成 d 维度的稠密向量。因此,对于任意一个问题为了捕捉问题单词之间的上下文关系,采用Sentence BERT模型来对每个问题进行独立编码,并将它们映射到一个密集的向量,生成问题的语义表示。

[0049] $h = \text{Sentence-BERT}(Q)$

[0050] 其中,Sentence-BERT表示Sentence BERT模型, h 为问题 Q 的表示向量, $h \in \mathbb{R}^{l_q \times d_h}$, l_q 是问题 Q 的长度, d_h 是输出的维度。

[0051] 对于验证问题对 Q^u 与 Q^a 通过问题表示层获得相应的表示向量 h_u 与 h_a ;对于问题集

合 $T = \{T^p, T^q\}$ 通过问题表示层获得语义匹配关系向量集 $H^p = \{h_1^p, h_2^p, \dots, h_i^p, \dots, h_n^p\}$ 和语义相似关系的向量集 $H^q = \{h_1^q, h_2^q, \dots, h_j^q, \dots, h_m^q\}$ 。

[0052] b) 关系感知表示层。

[0053] 在关系感知表示层中, 多个语义关系中包括语义匹配关系和语义相似关系, 每个语义关系都是关于问题 Q^a 的多个语义匹配问题。对于多个语义匹配问题, 首先使用RNs(关系网络)来捕获两个问题的关系向量, 然后, 扩展RNs, 探讨多关系标签对语义匹配和语义相似的整合。

[0054] 为了明确问题及其关联句中的关系语义, 采用RNs关系网络结构来表征两个对象之间的关系。以问题为关系对象, 利用RNs来表示问题 Q^a 和该问题相应的多关系向量 h_k 。依照RNs模型, 模型关系表征的函数:

$$[0055] \quad \text{RNs}(H_r) = f_\phi(\sum g_\theta(h_a, h_k))$$

[0056] 其中, g_θ 和 f_ϕ 是多层感知机网络(Multi-Layer Perceptron, MLP), h_a 为问题 Q^a 的表示向量, h_k 表示问题集合 T 中单个问题的表示向量; $h_k \in H^p$ or H^q , H^p 、 H^q 分别为集合 T^p 、 T^q 中问题的表示向量集合。

[0057] 为了得到 h_a 和 h_k 之间的交互特征, 采用element-wise product(元素的点乘积)方法, 利用 h_a 和 h_k 对应位置元素进行乘积进行运算:

$$[0058] \quad h_k^a = h_a \odot h_k$$

[0059] 因此, 问题的关系表示的函数可以改写为:

$$[0060] \quad H_r = \sum g_\theta(h_k^a)$$

[0061] 其中, H_r 通过RNs网络得到关于问题 Q^a 的多语义匹配向量, 它利用元素的点乘积运算获取两句之间对应词的相似语义特征, 利用RNs网络结构获取多个关系表示的核心交互特征。

[0062] 为了体现不同标签之间的匹配关系的隐含语义, 本部分将基本的问题关系的函数扩展到多关系语义匹配表示。结合问题集合 T 中每一问题的标签, 计算如下多标签下的匹配表示函数:

$$[0063] \quad \tilde{H}_r = \sum_{r \in R} \sum_{o \in |N^r|} g_\theta(h_a, h_o^r)$$

[0064] 其中, $h_o^r \in N^r$, $o \in |N^r|$, N^r 为在第 r 个关系下多个语义匹配问题的语义向量集合, $|N^r|$ 表示语义向量集合 N^r 中的向量数目, R 表示标签的类别集合。直观上, 上式通过对多句的标签特征向量进行归一化编码。

[0065] 基于表示向量集合 H^p 与 H^q 将上述多标签下的匹配表示函数展开, 得到:

$$[0066] \quad \tilde{H}_r = \sum_{h_i^p \in N^p} g_\theta(W_{pi}(h_a, h_i^p) + b_{pi}) + \sum_{h_j^q \in N^q} g_\theta(W_{qj}(h_a, h_j^q) + b_{qj})$$

[0067] 其中, W_{pi} 和 W_{qj} 分别是在两个关系下的学习权重, b_{pi} 和 b_{qj} 是偏差参数; $h_i^p \in H^p$, $h_j^q \in H^q$, 各自为相应问题的表示向量。

[0068] 采用element-wise product方法进行运算, 分别得到的 h_a 与 h_i^p 之间、 h_a 与 h_j^q 之间的交互特征:

$$[0069] \quad h_i^{ap} = h_a \odot h_i^p$$

$$[0070] \quad h_j^{aq} = h_a \odot h_j^q$$

[0071] 得到:

$$[0072] \quad v_r = \sum_{h_i^{ap} \in N^p} g_\theta(W_{pi}h_i^{ap} + b_{pi}) + \sum_{h_j^{aq} \in N^q} g_\theta(W_{qj}h_j^{aq} + b_{qj})$$

[0073] 其中, v_r 为问题 Q^a 的关系感知表示向量。通过上述原理可知, 问题 Q^a 通过 RNs 网络获取了多个关系的语义匹配的表征。

[0074] c) 问题识别评估层。

[0075] 问题识别评估层的目标是评估每个问题对的语言匹配情况。首先, 将问题 Q^u 与问题 Q^a 的表示向量, 以及问题 Q^a 的关系感知表示向量连接, 表示为:

$$[0076] \quad z_{au} = h_u \oplus h_a \oplus v_r$$

[0077] 其中, h_u 、 h_a 、 v_r 依次表示问题 Q^u 的表示向量、问题 Q^a 的表示向量、问题 Q^a 的关系感知表示向量。

[0078] 然后, 通过 RELU 激活函数和 sigmoid 函数的运算获得验证问题对 Q^u 与 Q^a 的语义关系 $R(Q^u, Q^a)$, 表示为:

$$[0079] \quad o_{au} = \text{ReLU}(W_1 z_{au} + b_1),$$

$$[0080] \quad R(Q^u, Q^a) = \sigma(W_2 o_{au} + b_2)$$

[0081] 其中, o_{au} 为 RELU 激活函数的运算结果, $\sigma(\cdot)$ 为 sigmoid 函数, W_1, W_2, b_1, b_2 为网络参数。语义关系 $R(Q^u, Q^a)$ 即为识别评估结果, 为验证问题对的相似标签。

[0082] 2) 模型训练。

[0083] 本发明实施例中, 对上一部分构建的基于关系感知的神经网络模型中的所有上面介绍中所涉及的全部的 W 与 b 参数矩阵或向量 (即模型参数) 进行训练, 利用交叉熵损失函数作为最终优化目标, 损失函数表示为:

$$[0084] \quad L = -\frac{1}{N} \sum_{i=1}^N y_i \log P(y_i | R(Q_i^a, Q_i^u)) + \lambda_\theta \|\theta\|^2$$

[0085] 其中, θ 为待更新的模型参数, y_1 是对于第 1 个验证问题对的真实识别标签, N 是训练实例 (即验证问题对) 的数量; (Q_i^a, Q_i^u) 表示第 1 个验证问题对, $R(Q_i^a, Q_i^u)$ 表示 1 个验证问题对的识别评估结果, 考虑到模型的复杂性, 加入 l_2 -范数作为训练参数, λ_θ 为正则化超参数。

[0086] 通过优化上述损失函数, 能够学习到最优状态。示例性的, 在整个训练过程中, 采用 Adam 作为优化器, 学习率为 0.0005; 在训练过程中, 首先将获得的数据按照 20%、80% 的原则划分成训练数据和问题匹配知识库, 然后在训练数据中, 按照 60%、20%、20% 的比例划分训练集, 验证集和测试集, 训练集和验证集用于优化模型的参数, 测试集用来验证模型。

[0087] 步骤 3、测试阶段, 对于给定的问题对, 利用训练后的基于关系感知神经网络相似问题识别模型进行识别评估。

[0088] 通过前述步骤完成模型训练后, 对于给定的新的问题对 Q^u 与 Q^a , 同样的, 对于问题 Q^a 通过匹配问题知识库查找不同匹配关系下的多个语义相关的问题集合 T , 并进行文本

内容处理后,输入至训练后的基于关系感知神经网络相似问题识别模型中,然后按照上面的方式进行处理。得到新的问题对的相似标签,实现问题对的相似问题识别评估。

[0089] 本发明实施例上述方案,留言基于关系感知神经网络相似问题识别模型来进行问题对的相似问题识别,相比于传统模型,使用了多个语义匹配的语义关系相关信息。对于预测的结果,在多个评价指标上有一定的提高。

[0090] 本发明另一实施例还提供一种关系感知的相似问题识别评估系统,其主要用于实现前述实施例提供的方法,如图2所示,该系统主要包括:

[0091] 数据提取单元,用于从数据集中提取多个语义关系情况下,相匹配的语义相关的问题数据,每一相匹配的语义相关的问题数据均为文本数据,包含验证问题对 Q^u 与 Q^a ,以及问题 Q^a 对应的不同匹配关系下的多个语义相关的问题集合T;

[0092] 模型构建与训练单元,用于构建基于关系感知神经网络相似问题识别模型,对每一相匹配的语义相关的问题数据进行联合表征,并利用获得的问题 Q^a 的关系感知表示向量与问题 Q^u 的表示向量进行识别评估,训练阶段,利用识别评估结果与识别标签构建损失函数,更新模型参数;

[0093] 识别评估测试单元,用于对于给定的问题对,利用训练后的基于关系感知神经网络相似问题识别模型进行识别评估。

[0094] 所属领域的技术人员可以清楚地了解到,为描述的方便和简洁,仅以上述各功能模块的划分进行举例说明,实际应用中,可以根据需要而将上述功能分配由不同的功能模块完成,即将系统的内部结构划分成不同的功能模块,以完成以上描述的全部或者部分功能。

[0095] 此外,该系统各单元所涉及的具体技术细节在之前的方法实施例中已经做了详细的介绍,故不再赘述。

[0096] 本发明另一实施例还提供一种处理设备,如图3所示,其主要包括:一个或多个处理器;存储器,用于存储一个或多个程序;其中,当所述一个或多个程序被所述一个或多个处理器执行时,使得所述一个或多个处理器实现前述实施例提供的方法。

[0097] 进一步的,所述处理设备还包括至少一个输入设备与至少一个输出设备;在所述处理设备中,处理器、存储器、输入设备、输出设备之间通过总线连接。

[0098] 本发明实施例中,所述存储器、输入设备与输出设备的具体类型不做限定;例如:

[0099] 输入设备可以为触摸屏、图像采集设备、物理按键或者鼠标等;

[0100] 输出设备可以为显示终端;

[0101] 存储器可以为随机存取存储器(Random Access Memory, RAM),也可为非不稳定的存储器(non-volatile memory),例如磁盘存储器。

[0102] 本发明另一实施例还提供一种可读存储介质,存储有计算机程序,当计算机程序被处理器执行时实现前述实施例提供的方法。

[0103] 本发明实施例中可读存储介质作为计算机可读存储介质,可以设置于前述处理设备中,例如,作为处理设备中的存储器。此外,所述可读存储介质也可以是U盘、移动硬盘、只读存储器(Read-Only Memory, ROM)、磁碟或者光盘等各种可以存储程序代码的介质。

[0104] 以上所述,仅为本发明较佳的具体实施方式,但本发明的保护范围并不局限于此,任何熟悉本技术领域的技术人员在本发明披露的技术范围内,可轻易想到的变化或替换,

都应涵盖在本发明的保护范围之内。因此,本发明的保护范围应该以权利要求书的保护范围为准。

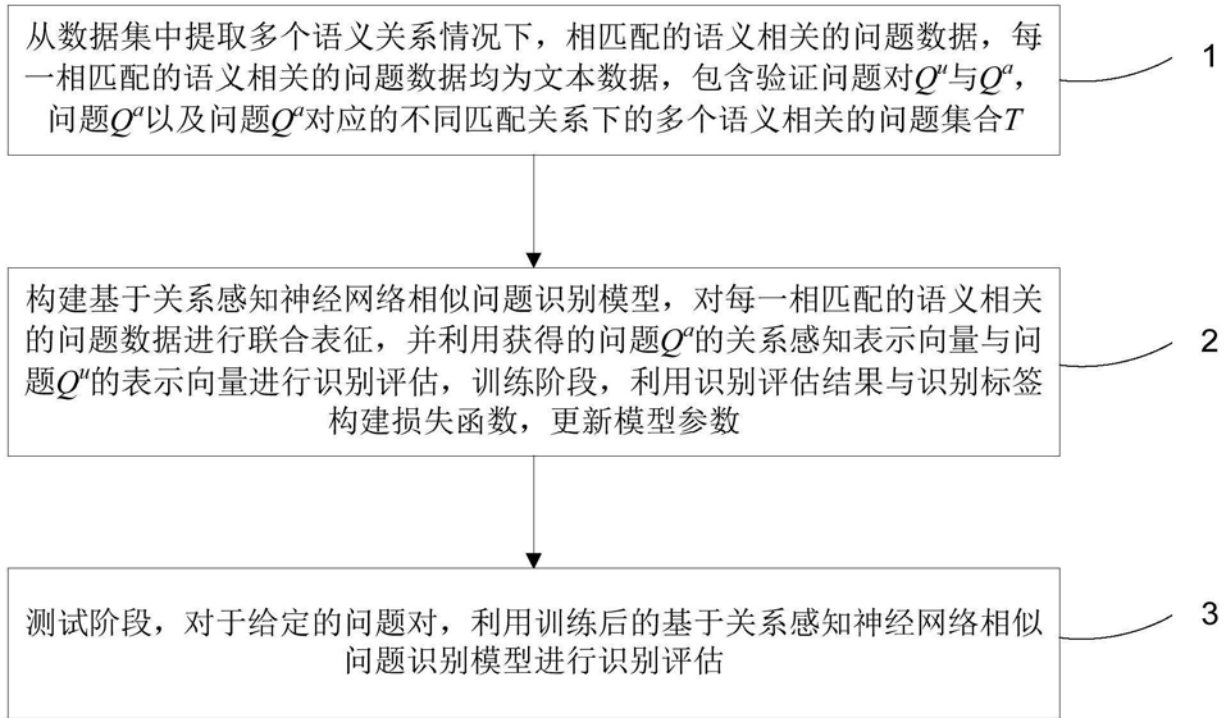


图1



图2



图3