



(12) 发明专利

(10) 授权公告号 CN 111339780 B

(45) 授权公告日 2020.11.06

(21) 申请号 202010408398.8

G06N 3/04 (2006.01)

(22) 申请日 2020.05.14

审查员 白露霜

(65) 同一申请的已公布的文献号
申请公布号 CN 111339780 A

(43) 申请公布日 2020.06.26

(73) 专利权人 北京金山数字娱乐科技有限公司
地址 100085 北京市海淀区小营西路33号
金山软件大厦2层西区

(72) 发明人 李长亮 白静 唐剑波

(74) 专利代理机构 北京智信禾专利代理有限公司
11637

代理人 王治东

(51) Int. Cl.

G06F 40/295 (2020.01)

G06N 3/08 (2006.01)

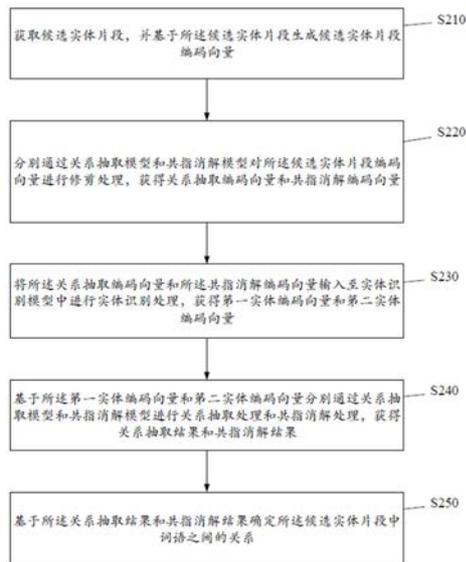
权利要求书2页 说明书16页 附图4页

(54) 发明名称

一种基于多任务模型的词语处理方法及装置

(57) 摘要

本申请提供一种基于多任务模型的词语处理方法及装置。其中所述方法包括：获取候选实体片段，并基于候选实体片段生成候选实体片段编码向量；分别通过关系抽取模型和共指消解模型对所述候选实体片段编码向量进行修剪处理，获得关系抽取编码向量和共指消解编码向量；将所述关系抽取编码向量和所述共指消解编码向量输入至实体识别模型中进行实体识别处理，获得第一实体编码向量和第二实体编码向量；再分别通过关系抽取模型和共指消解模型进行关系抽取处理和共指消解处理，获得关系抽取结果和共指消解结果，并确定候选实体片段中词语之间的关系。本申请所述的基于多任务模型的词语处理方法可以有效提高词语关系确定的准确率。



1. 一种基于多任务模型的词语处理方法,其特征在于,包括:

获取候选实体片段,并基于所述候选实体片段生成候选实体片段编码向量;

分别通过关系抽取模型和共指消解模型对所述候选实体片段编码向量进行修剪处理,所述关系抽取模型和所述共指消解模型分别与实体识别模型共用前馈神经网络,通过所述前馈神经网络对所述候选实体片段编码向量进行打分,获得候选实体片段编码向量的分数,并将分数大于或等于预设阈值的候选实体片段编码向量分别作为关系抽取编码向量和共指消解编码向量;

将所述关系抽取编码向量和共指消解编码向量分别输入至所述实体识别模型中,通过前馈神经网络进行基础打分与分类打分,分别获得关系抽取编码向量的分数和共指消解编码向量的分数,基于所述关系抽取编码向量的分数生成第一实体编码向量,以及基于所述共指消解编码向量的分数生成第二实体编码向量;

基于所述第一实体编码向量和第二实体编码向量分别通过所述关系抽取模型和所述共指消解模型进行关系抽取处理和共指消解处理,获得关系抽取结果和共指消解结果;

基于所述关系抽取结果和共指消解结果确定所述候选实体片段中词语之间的关系。

2. 根据权利要求1所述的基于多任务模型的词语处理方法,其特征在于,基于所述关系抽取编码向量的分数生成第一实体编码向量,包括:

基于所述关系抽取编码向量的分数对所述关系抽取编码向量进行分类预测处理,获得所述关系抽取编码向量的分类标签;

基于所述关系抽取编码向量以及关系抽取编码向量的分类标签,生成第一实体编码向量。

3. 根据权利要求2所述的基于多任务模型的词语处理方法,其特征在于,基于所述关系抽取编码向量以及关系抽取编码向量的分类标签,生成第一实体编码向量,包括:

将所述关系抽取编码向量的分类标签进行编码处理,生成第一标签向量;

基于所述关系抽取编码向量和所述第一标签向量,生成第一实体编码向量。

4. 根据权利要求1所述的基于多任务模型的词语处理方法,其特征在于,基于所述共指消解编码向量的分数生成第二实体编码向量,包括:

基于所述共指消解编码向量的分数对所述共指消解编码向量进行分类预测处理,获得所述共指消解编码向量的分类标签;

基于所述共指消解编码向量以及共指消解编码向量的分类标签,生成第二实体编码向量。

5. 根据权利要求4所述的基于多任务模型的词语处理方法,其特征在于,基于所述共指消解编码向量和所述共指消解编码向量的分类标签,生成第二实体编码向量,包括:

将所述共指消解编码向量的分类标签进行编码处理,生成第二标签向量;

基于所述共指消解编码向量和所述第二标签向量,生成第二实体编码向量。

6. 一种基于多任务模型的词语处理装置,其特征在于,包括:

片段获取模块,被配置为获取候选实体片段,并基于所述候选实体片段生成候选实体片段编码向量;

片段修剪模块,被配置为分别通过关系抽取模型和共指消解模型对所述候选实体片段编码向量进行修剪处理,所述关系抽取模型和所述共指消解模型分别与实体识别模型共用

前馈神经网络,通过所述前馈神经网络对所述候选实体片段编码向量进行打分,获得候选实体片段编码向量的分数,并将分数大于或等于预设阈值的候选实体片段编码向量分别作为关系抽取编码向量和共指消解编码向量;

实体识别模块,被配置为将所述关系抽取编码向量和共指消解编码向量分别输入至所述实体识别模型中,通过前馈神经网络进行基础打分与分类打分,分别获得关系抽取编码向量的分数和共指消解编码向量的分数,基于所述关系抽取编码向量的分数生成第一实体编码向量,以及基于所述共指消解编码向量的分数生成第二实体编码向量;

关系处理模块,被配置为基于所述第一实体编码向量和第二实体编码向量分别通过所述关系抽取模型和所述共指消解模型进行关系抽取处理和共指消解处理,获得关系抽取结果和共指消解结果;

关系确定模块,被配置为基于所述关系抽取结果和共指消解结果确定所述候选实体片段中词语之间的关系。

7.一种计算设备,包括存储器、处理器及存储在存储器上并可在处理器上运行的计算机指令,其特征在于,所述处理器执行所述指令时实现权利要求1-5任意一项所述方法的步骤。

8.一种计算机可读存储介质,其存储有计算机指令,其特征在于,该指令被处理器执行时实现权利要求1-5任意一项所述方法的步骤。

一种基于多任务模型的词语处理方法及装置

技术领域

[0001] 本申请涉及计算机技术领域,特别涉及一种基于多任务模型的词语处理方法及装置、计算设备及计算机可读存储介质。

背景技术

[0002] 实体识别是指在非结构化的文本中识别并抽出具有特定意义或指代性强的实体,比如人名、地名、组织结构名、日期时间、专有名词等。

[0003] 关系是两个或多个实体之间的某种联系,关系抽取是从文本中检测和识别出实体与实体之间具有的某种语义关系,比如句子“北京是中国的首都、政治中心和文化中心”,其中表述的关系可以为(中国,首都,北京)、(中国,政治中心,北京)或(中国,文化中心,北京)。

[0004] 共指消解是特殊的关系抽取,共指消解的其中一个实体通常是另外一个实体在当前语境下的不同说法,两个实体之间的关系可以表示为(实体1,共指,实体2)。

[0005] 目前,对于语句的实体识别任务、关系抽取任务、共指消解任务均是分别单独进行的,进而导致实体识别、关系抽取、共指消解的效果均不理想。

发明内容

[0006] 有鉴于此,本申请实施例提供了一种基于多任务模型的词语处理方法及装置、计算设备及计算机可读存储介质,以解决现有技术中存在的技术缺陷。

[0007] 本申请实施例公开了一种基于多任务模型的词语处理方法及装置、计算设备及计算机可读存储介质,包括:

[0008] 获取候选实体片段,并基于所述候选实体片段生成候选实体片段编码向量;

[0009] 分别通过关系抽取模型和共指消解模型对所述候选实体片段编码向量进行修剪处理,获得关系抽取编码向量和共指消解编码向量;

[0010] 将所述关系抽取编码向量和所述共指消解编码向量输入至实体识别模型中进行实体识别处理,获得第一实体编码向量和第二实体编码向量;

[0011] 基于所述第一实体编码向量和第二实体编码向量分别通过所述关系抽取模型和所述共指消解模型进行关系抽取处理和共指消解处理,获得关系抽取结果和共指消解结果;

[0012] 基于所述关系抽取结果和共指消解结果确定所述候选实体片段中词语之间的关系。

[0013] 进一步地,所述关系抽取模型与所述实体识别模型共用前馈神经网络;

[0014] 通过关系抽取模型对所述候选实体片段编码向量进行修剪处理,包括:

[0015] 通过所述前馈神经网络对所述候选实体片段编码向量进行打分,获得候选实体片段编码向量的分数,并将分数大于或等于预设阈值的候选实体片段编码向量作为关系抽取编码向量。

- [0016] 进一步地,所述共指消解模型与所述实体识别模型共用前馈神经网络;
- [0017] 通过共指消解模型对所述候选实体片段编码向量进行修剪处理,包括:
- [0018] 通过所述前馈神经网络对所述候选实体片段编码向量进行打分,获得候选实体片段编码向量的分数,并将分数大于或等于预设阈值的候选实体片段编码向量作为共指消解编码向量。
- [0019] 进一步地,将所述关系抽取编码向量输入至实体识别模型中进行实体识别处理获得第一实体编码向量,包括:
- [0020] 将所述关系抽取编码向量输入至实体识别模型中,通过所述前馈神经网络进行基础打分与分类打分,获得关系抽取编码向量的分数;
- [0021] 基于所述关系抽取编码向量的分数生成第一实体编码向量。
- [0022] 进一步地,基于所述关系抽取编码向量的分数生成第一实体编码向量,包括:
- [0023] 基于所述关系抽取编码向量的分数对所述关系抽取编码向量进行分类预测处理,获得所述关系抽取编码向量的分类标签;
- [0024] 基于所述关系抽取编码向量以及关系抽取编码向量的分类标签,生成第一实体编码向量。
- [0025] 进一步地,基于所述关系抽取编码向量以及关系抽取编码向量的分类标签,生成第一实体编码向量,包括:
- [0026] 将所述关系抽取编码向量的分类标签进行编码处理,生成第一标签向量;
- [0027] 基于所述关系抽取编码向量和所述第一标签向量,生成第一实体编码向量。
- [0028] 进一步地,将所述共指消解编码向量输入至实体识别模型中进行实体识别处理获得第二实体编码向量,包括:
- [0029] 将所述共指消解编码向量输入至实体识别模型中,通过所述前馈神经网络进行基础打分与分类打分,获得共指消解编码向量的分数;
- [0030] 基于所述共指消解编码向量的分数生成第二实体编码向量。
- [0031] 进一步地,基于所述共指消解编码向量的分数生成第二实体编码向量,包括:
- [0032] 基于所述共指消解编码向量的分数对所述共指消解编码向量进行分类预测处理,获得所述共指消解编码向量的分类标签;
- [0033] 基于所述共指消解编码向量以及共指消解编码向量的分类标签,生成第二实体编码向量。
- [0034] 进一步地,基于所述共指消解编码向量和所述共指消解编码向量的分类标签,生成第二实体编码向量,包括:
- [0035] 将所述共指消解编码向量的分类标签进行编码处理,生成第二标签向量;
- [0036] 基于所述共指消解编码向量和所述第二标签向量,生成第二实体编码向量。
- [0037] 本申请还提供一种基于多任务模型的词语处理装置,包括:
- [0038] 片段获取模块,被配置为获取候选实体片段,并基于所述候选实体片段生成候选实体片段编码向量;
- [0039] 片段修剪模块,被配置为分别通过关系抽取模型和共指消解模型对所述候选实体片段编码向量进行修剪处理,获得关系抽取编码向量和共指消解编码向量;
- [0040] 实体识别模块,被配置为将所述关系抽取编码向量和所述共指消解编码向量输入

至实体识别模型中进行实体识别处理,获得第一实体编码向量和第二实体编码向量;

[0041] 关系处理模块,被配置为基于所述第一实体编码向量和第二实体编码向量分别通过所述关系抽取模型和所述共指消解模型进行关系抽取处理和共指消解处理,获得关系抽取结果和共指消解结果;

[0042] 关系确定模块,被配置为基于所述关系抽取结果和共指消解结果确定所述候选实体片段中词语之间的关系。

[0043] 本申请还公开了一种计算机可读存储介质,其存储有计算机指令,该指令被处理器执行时实现所述基于多任务模型的词语处理方法的步骤。

[0044] 本申请提供的基于多任务模型的词语处理方法及装置,通过关系抽取模型和共指消解模型分别对候选实体片段编码向量进行修剪处理,获得关系抽取编码向量和共指消解编码向量,实现了基于不同的任务需求对候选实体片段的初步筛选,再通过对关系抽取编码向量和共指消解编码向量进行实体识别处理,获得第一实体编码向量和第二实体编码向量,可以进一步地、更深层次的加强对于关系抽取编码向量、共指消解编码向量对应的候选实体片段的理解,为关系抽取任务和共指消解任务的执行提供基础,最终通过对第一实体编码向量、第二实体编码向量分别进行关系抽取处理和共指消解处理,并基于获得的关系抽取结果和共指消解结果确定候选实体片段中词语的关系,实现了关系抽取模型、共指消解模型、实体识别模型的有机结合,也即实现了关系抽取任务、共指消解任务、实体识别任务三个任务中信息的共享,可以有效提高关系抽取、共指消解、实体识别的正确率、召回率,有效提高词语关系确定的准确率。

附图说明

[0045] 图1是本申请一实施例所述的多任务模型的结构示意图;

[0046] 图2是本申请一实施例所述的基于多任务模型的词语处理方法的步骤流程示意图;

[0047] 图3是本申请一实施例所述的基于多任务模型的词语处理方法的步骤流程示意图;

[0048] 图4是本申请一实施例所述的基于多任务模型的词语处理装置的结构示意图;

[0049] 图5是本申请一实施例所述的计算设备的结构示意图。

具体实施方式

[0050] 在下面的描述中阐述了很多具体细节以便于充分理解本申请。但是本申请能够以很多不同于在此描述的其它方式来实施,本领域技术人员可以在不违背本申请内涵的情况下做类似推广,因此本申请不受下面公开的具体实施的限制。

[0051] 在本说明书一个或多个实施例中使用的术语是仅仅出于描述特定实施例的目的,而非旨在限制本说明书一个或多个实施例。在本说明书一个或多个实施例和所附权利要求书中所使用的单数形式的“一种”、“所述”和“该”也旨在包括多数形式,除非上下文清楚地表示其他含义。还应当理解,本说明书一个或多个实施例中使用的术语“和/或”是指并包含一个或多个相关联的列出项目的任何或所有可能组合。

[0052] 应当理解,尽管在本说明书一个或多个实施例中可能采用术语第一、第二等来描

述各种信息,但这些信息不应限于这些术语。这些术语仅用来将同一类型的信息彼此区分开。例如,在不脱离本说明书一个或多个实施例范围的情况下,第一也可以被称为第二,类似地,第二也可以被称为第一。取决于语境,如在此所使用的词语“如果”可以被解释成为“在……时”或“当……时”或“响应于确定”。

[0053] 首先,对本发明一个或多个实施例涉及的名词术语进行解释。

[0054] 实体识别:是指在非结构化的文本中识别并抽出具有特定意义或指代性强的实体,比如人名、地名、组织结构名、日期时间、专有名词等。

[0055] 实体识别模型:用于执行实体识别任务的模型。

[0056] 关系抽取:从文本中检测和识别出实体与实体之间具有的某种语义关系,比如句子“北京是中国的首都、政治中心和文化中心”,其中表述的关系可以为(中国,首都,北京)、(中国,政治中心,北京)或(中国,文化中心,北京)。

[0057] 关系抽取模型:用于执行关系抽取任务的模型。

[0058] 共指消解:特殊的关系抽取,共指消解的其中一个实体通常是另外一个实体在当前语境下的不同说法,两个实体之间的关系可以表示为(实体1,共指,实体2)。

[0059] 共指消解模型:用于执行共指消解任务的模型。

[0060] 候选实体片段(span):由语句中的一个词或多个词组成的片段。

[0061] 候选实体片段编码向量(span embedding):候选实体片段经过编码器的编码处理生成的向量。

[0062] 修剪:根据预设的规则进行筛选。

[0063] 关系抽取编码向量:基于关系抽取处理的结果以及相应的候选实体片段编码向量的分数对候选实体片段编码向量进行修剪后剩余的编码向量。

[0064] 共指消解编码向量:基于共指消解处理的结果相应的候选实体片段编码向量的分数对候选实体片段编码向量进行修剪后剩余的编码向量。

[0065] 第一实体编码向量:由关系抽取编码向量与第一标签向量组成的编码向量。

[0066] 第二实体编码向量:由共指消解编码向量与第二标签向量组成的编码向量。

[0067] 第一标签向量:对关系抽取编码向量的分类标签进行编码处理得到的编码向量。

[0068] 第二标签向量:对共指消解编码向量的分类标签进行编码处理得到的编码向量。

[0069] 预设阈值:判断候选实体片段编码向量是否可以作为关系抽取编码向量或共指消解编码向量的分数临界值。

[0070] 分类标签:用于标识编码向量类型的标识。

[0071] 前馈神经网络(FeedForward Neural Network,FFNN):是一种最简单的神经网络,各神经元分层排列,每个神经元只与前一层的神经元相连,接收前一层的输出,并输出给下一层.各层间没有反馈,是目前应用最广泛、发展最迅速的人工神经网络之一。本申请的实体识别模型、关系抽取模型与共指消解模型共用一个用于打分的前馈神经网络。

[0072] 卷积神经网络(Convolutional Neural Networks, CNN):是一类包含卷积计算且具有深度结构的前馈神经网络,是深度学习(deep learning)的代表算法之一,在本申请中,通过cnn网络编码得到字符级别的特征向量。

[0073] 正确率:是指识别出的正确实体数与识别出的实体数的比值,取值在0-1之间,数值越大,正确率越高。

[0074] 召回率:是指识别出的正确实体数与样本的实体数的比值,取值在0-1之间,数值越大,找回率越高。

[0075] 加权调和平均值:又称F1值, $F1值 = (2 * 正确率 * 召回率) / (正确率 + 召回率)$ 。

[0076] 在本申请中,提供了一种基于多任务模型的词语处理方法及装置、计算设备及计算机可读存储介质,在下面的实施例中逐一进行详细说明。

[0077] 如图1所示,本实施例提供了一种多任务模型,所述多任务模型用于本申请所述的基于多任务模型的词语处理方法,包括编码器、关系抽取模型、实体识别模型与共指消解模型,其中,关系抽取模型、实体识别模型与共指消解模型共用一个前馈神经网络。

[0078] 关系抽取模型,是用于检测、识别、抽取实体与实体之间语义关系的模型,在本实施例中,关系抽取模型首先对候选实体片段编码向量进行打分,根据打分结果对候选实体编码向量进行修剪,获得关系抽取编码向量,在上述关系抽取编码向量经过实体识别模型的处理并获得第一识别编码向量后,对第一识别编码向量进行打分,基于上述打分结果对第一识别编码向量进行分类预测处理,得到关系抽取结果。

[0079] 实体识别模型,是用于识别非结构化输入文本中的实体的模型,在本实施例中,关系抽取编码向量、共指消解编码向量输入至实体识别模型中进行实体识别处理,分别对关系抽取编码向量以及共指消解编码向量进行打分,并基于打分结果分别对上述关系抽取编码向量以及共指消解编码向量进行分类预测处理,获得每一个编码向量对应的分类标签,基于关系抽取编码向量及其分类标签生成第一实体编码向量,基于共指消解编码向量及其分类标签生成第二实体编码向量。

[0080] 共指消解模型,是用于检测、识别、抽取存在共指关系的实体的模型,在本实施例中,共指消解模型首先对候选实体片段编码向量进行打分,根据打分结果对候选实体编码向量进行修剪,获得共指消解编码向量,在上述共指消解编码向量经过实体识别模型的处理获得第二识别编码向量后,对第二识别编码向量进行打分,基于上述打分结果对第二识别编码向量进行分类预测处理,得到共指消解结果。

[0081] 本实施例提供的多任务模型,通过将关系抽取模型、实体识别模型、共指消解模型进行有机结合,三者共用一个用于打分的前馈神经网络,可以实现关系抽取模型、实体识别模型、共指消解模型彼此之间的信息共享,提高上述关系抽取模型、实体识别模型、共指消解模型的正确率和召回率。

[0082] 如图2所示,本实施例提供了一种基于多任务模型的词语处理方法,包括步骤S210至步骤S250。

[0083] S210、获取候选实体片段,并基于所述候选实体片段生成候选实体片段编码向量。

[0084] 其中,候选实体片段是由目标语句或目标段落、目标文章中的一个或多个词语组合成的词语集合,每一个词语均表示一个实体。具体地,可以通过对目标语句或目标段落、目标文章等进行分词处理,再在上述分词处理的结果中抽取一个或多个目标词语组合成词语集合,即得到候选实体片段。

[0085] 例如,假设对目标语句进行分词处理后,得到分词处理结果A1-A10在内的10个词语,在上述分词处理结果中进行抽取得到A1-A6组成的词语集合,并将上述词语集合作为候选实体片段。

[0086] 在实际应用中,可以将候选实体片段输入至编码器中进行编码处理,生成候选实体片段编码向量。

[0087] 在本实施例中,编码器包括双向LSTM、预训练的BERT模型、CNN网络及其任意组合。

[0088] 优选地,首先通过预训练的BERT模型对包含若干候选实体片段的语句进行编码处理后得到该语句词级别的特征向量,通过CNN网络进行编码处理后得到该语句字符级别的特征向量,将上述词级别的特征向量以及字符级别的特征向量进行拼接,得到拼接向量,再将上述拼接向量通过双向LSTM网络进行编码处理后得到带有上下文特征的特征向量,最后基于抽取得到的候选实体片段利用注意力机制计算得到每个候选实体片段编码向量,候选实体片段的编码向量可以通过下式表示:

$$[0089] \quad \mathbf{g}_i = [\mathbf{x}_{\text{START}(i)}^*, \mathbf{x}_{\text{END}(i)}^*, \hat{\mathbf{x}}_i, \phi(i)] \quad (1)$$

[0090] 其中, \mathbf{g}_i 表示候选实体片段的编码向量, $\mathbf{x}_{\text{START}(i)}^*$ 、 $\mathbf{x}_{\text{END}(i)}^*$ 表示候选实体片段起止位置的向量, $\phi(i)$ 表示额外的特征, $\hat{\mathbf{x}}_i$ 表示基于注意力机制对每一个候选实体片段中的词进行计算得到的结果, $\hat{\mathbf{x}}_i$ 的具体计算过程如下:

$$[0091] \quad \mathbf{x}_t^* = [\mathbf{h}_{t,1}, \mathbf{h}_{t,-1}] \quad (2)$$

$$[0092] \quad \alpha_t = \mathbf{w}_\alpha \cdot \text{FFNN}_\alpha(\mathbf{x}_t^*) \quad (3)$$

$$[0093] \quad a_{i,t} = \frac{\exp(\alpha_t)}{\sum_{k=\text{START}(i)}^{\text{END}(i)} \exp(\alpha_k)} \quad (4)$$

$$[0094] \quad \hat{\mathbf{x}}_i = \sum_{t=\text{START}(i)}^{\text{END}(i)} a_{i,t} \cdot \mathbf{x}_t \quad (5)$$

[0095] 具体地, t 表示候选实体片段, i 表示候选实体片段中的词,公式(2)表示候选实体片段中每一个词对应的编码向量 \mathbf{x}_t^* 由经过双向lstm的正向传播输出的向量($\mathbf{h}_{t,1}$)以及反向传播输出的向量($\mathbf{h}_{t,-1}$)组成,公式(3)表示候选实体片段 t 的参数 α 通过其参数 w 与前馈神经网络对该候选实体片段打出的分数点乘得到,公式(4)表示候选实体片段中每一个词的权重 $a_{i,t}$ 基于其所在的候选实体片段的参数 α 以及该词在候选实体片段的总参数得到,公式(4)表示候选实体片段中每一个词对应的编码向量 $\hat{\mathbf{x}}_i$ 该词在该候选实体片段中的权重参数 $a_{i,t}$ 与该候选实体片段编码向量 \mathbf{x}_t 得到。

[0096] 本实施例通过获取候选实体片段,并对候选实体片段进行编码处理,以为后续其他任务的执行做好准备,提高后续任务执行的效率。

[0097] S220、分别通过关系抽取模型和共指消解模型对所述候选实体片段编码向量进行修剪处理,获得关系抽取编码向量和共指消解编码向量。

[0098] 其中,关系抽取是指通过关系抽取模型检测和识别候选实体片段中的词语即实体之间的语义关系,语义关系的类型包括但不限于原因、特征、上位、场所、方式、材料、方法、部分、所有者、意图、同义、时间、受事、施事、使用者关系。

[0099] 共指消解是指通过共指消解模型检测和识别候选实体片段中的词语即实体之间的共指关系,比如词语“蓉城”、“天府之国”均是指代“成都”,所以词语“蓉城”、“天府之国”之间存在共指关系。

[0100] 具体地,所述关系抽取模型与所述实体识别模型共用一个用于打分的前馈神经网络。

[0101] 在实际应用中,可以通过所述前馈神经网络对所述候选实体片段编码向量进行打分,获得候选实体片段编码向量的分数,并将分数大于或等于预设阈值的候选实体片段编码向量作为关系抽取编码向量。

[0102] 具体地,所述共指消解模型与所述实体识别模型共用一个用于打分的前馈神经网络。

[0103] 在实际应用中,可以通过所述前馈神经网络对所述候选实体片段编码向量进行打分,获得候选实体片段编码向量的分数,并将分数大于或等于预设阈值的候选实体片段编码向量作为共指消解编码向量。

[0104] 其中,每一个候选实体片段编码向量的分数均由基础分数(Mention score)和分类分数(classifier score)组成,并且上述基础分数和分类分数均通过前馈神经网络打分得到。前馈神经网络是利用深度学习的原理对候选实体片段编码向量进行打分的,具体而言,通过利用前馈神经网络对候选实体片段编码向量再次进行计算或编码,并映射出相应的分数,即获得候选实体片段编码向量的分数。需要说明的是,前馈神经网络对于分数的映射可以通过后续任务的执行、损失值的计算、梯度的回传等不断进行调整。候选实体片段编码向量的分数可以为十分制分数、百分制分数、千分制分数等,本申请对此不做限制。

[0105] 例如,假设候选实体片段编码向量分数的预设阈值为60分,存在a1-a6共6个候选实体编码向量。

[0106] 通过前馈神经网络对候选实体片段编码向量进行打分,得到候选实体片段编码向量a1-a6分数分别为85分、72分、40分、33分、68分、45分,那么候选实体片段编码向量a1、a2、a5的分数大于预设阈值,将候选实体片段编码向量a1、a2、a5作为关系抽取编码向量。

[0107] 通过前馈神经网络对候选实体片段编码向量进行打分,得到候选实体片段编码向量a1-a6分数分别为74分、49分、60分、74分、68分、30分,那么候选实体片段编码向量a1、a3、a4、a5的分数大于或等于预设阈值,将候选实体片段编码向量a1、a3、a4、a5作为共指消解编码向量。

[0108] 本实施例将候选实体片段编码向量分别对候选实体片段编码向量进行打分,以获得复合任务需求的候选实体片段编码向量,可以针对不同的任务需求对候选实体片段编码向量进行相应的初步筛选,为后续的步骤做铺垫,以提高关系抽取任务、实体识别任务、共指消解任务的准确率。

[0109] S230、将所述关系抽取编码向量和所述共指消解编码向量输入至实体识别模型中进行实体识别处理,获得第一实体编码向量和第二实体编码向量。

[0110] 具体地,所述步骤S230包括步骤S231至步骤S234。需要说明的是,步骤S231与步骤

S232为并列执行的步骤,步骤S233与步骤S234为并列执行的步骤。

[0111] S231、将所述关系抽取编码向量输入至实体识别模型中,通过所述前馈神经网络进行基础打分与分类打分,获得关系抽取编码向量的分数。

[0112] 具体地,实体识别模型与关系抽取模型、共指消解模型共用一个用于打分的前馈神经网络,其中基础打分是基于上一步关系抽取处理的结果通过前馈神经网络对关系抽取编码向量的再次打分,分类打分是通过对关系抽取编码向量进行实体识别处理,再基于实体识别处理的结果通过前馈神经网络对关系抽取编码向量进行的打分,最终基于关系抽取编码向量基础打分的分数与分类打分的分数计算得到关系抽取编码向量的分数,上述计算方式可以为求和、求平均值、求加权平均值等等,可视具体情况而定,本申请对此不做限制。

[0113] 例如,假设关系抽取编码向量a1、a2、a5的基础打分分数分别为90分、70分、70分,分类打分的分数分别为69分、73分、81分,在采用求和的方式计算关系抽取编码向量分数的情况下,关系抽取编码向量a1、a2、a5的分数分别为159分、143分、151分。

[0114] 本实施例中分别对关系抽取编码向量进行基础打分和分类打分以得到该关系抽取编码向量最终的分数,可以有效提高打分的准确率,提高关系抽取任务的准确率。

[0115] 在步骤S231执行完毕后,执行步骤S233。

[0116] S232、将所述共指消解编码向量输入至实体识别模型中,通过所述前馈神经网络进行基础打分与分类打分,获得共指消解编码向量的分数。

[0117] 具体地,实体识别模型与关系抽取模型、共指消解模型共用一个用于打分的前馈神经网络,其中基础打分是基于上一步共指消解处理的结果通过前馈神经网络对共指消解编码向量的再次打分,分类打分是通过对共指消解编码向量进行实体识别处理,再基于实体识别处理的结果通过前馈神经网络对共指消解编码向量进行的打分,最终基于共指消解编码向量基础打分的分数与分类打分的分数计算得到共指消解编码向量的分数,上述计算方式可以为求和、求平均值、求加权平均值等等,可视具体情况而定,本申请对此不做限制。

[0118] 例如,假设共指消解编码向量a1、a3、a4、a5的基础打分分数分别为75分、59分、74分、70分,分类打分的分数分别为60分、70分、65分、75分,在采用求和的方式计算共指消解编码向量分数的情况下,共指消解编码向量a1、a3、a4、a5的分数分别为135分、129分、139分、145分。

[0119] 本实施例中分别对共指消解编码向量进行基础打分和分类打分以得到该共指消解编码向量最终的分数,可以有效提高打分的准确率,提高共指消解任务的准确率。

[0120] 在步骤S232执行完毕后,执行步骤S234。

[0121] S233、基于所述关系抽取编码向量的分数生成第一实体编码向量。

[0122] 具体地,所述步骤S233包括步骤S2331至步骤S2332。

[0123] S2331、基于所述关系抽取编码向量的分数对所述关系抽取编码向量进行分类预测处理,获得所述关系抽取编码向量的分类标签。

[0124] 具体地,分类预测处理是指通过softmax函数对关系抽取编码向量基于其分数按照其对应词语的属性进行分类,并获得分类标签。

[0125] 更为具体地,softmax函数的公式如下:

$$[0126] \quad s_i = \frac{e^i}{\sum_j e^j} \quad (6)$$

[0127] 其中, S_i 表示第*i*个关系抽取编码向量对应的softmax值; *i*代表第*i*个关系抽取编码向量; *j*代表关系抽取编码向量的总个数。

[0128] 比如,若第一个关系抽取编码向量的分类标签为“方法”,第二个关系抽取编码向量的分类标签为“任务”,那么二者之间的语义关系为“用于”的关系。

[0129] 例如,基于关系抽取编码向量a1、a2和a5的分数159分、143分和151分按照每个关系抽取编码向量对应词语的属性对上述关系抽取编码向量进行分类预测处理,获得关系抽取编码向量a1、a2和a5的分类标签分别为M1、M2和M5。

[0130] 本实施例通过分类预测处理获得关系抽取编码向量的分类标签,有助于加深模型对候选实体片段的认知,进而提高关系抽取的准确率。

[0131] S2332、基于所述关系抽取编码向量以及关系抽取编码向量的分类标签,生成第一实体编码向量。

[0132] 在实际应用中,可以将所述关系抽取编码向量的分类标签进行编码处理,生成第一标签向量;基于所述关系抽取编码向量和所述第一标签向量,生成第一实体编码向量。

[0133] 进一步地,关系抽取编码向量与第一标签向量相结合,首先生成包含标签信息的关系抽取编码向量,再将具有语义关系的两个词语对应的包含标签信息的关系抽取编码向量相结合,生成第一实体编码向量,如下所示:

[0134] `span_pair_embeddings = torch.cat([span1_embeddings, span2_embeddings, span1_embeddings*span2_embeddings, span1_label_embedding, span2_label_embedding], -1)`。

[0135] 其中, `torch.cat`是用于将两个或多个向量拼接在一起的函数, `span_pair_embeddings`表示第一实体编码向量, `span1_embeddings`表示关系抽取编码向量1, `span2_embeddings`表示关系抽取编码向量2, `span1_label_embedding`表示关系抽取编码向量1的标签向量即第一标签向量1, `span2_label_embedding`表示关系抽取编码向量2的标签向量即第一标签向量2。

[0136] 例如,将关系抽取编码向量a1、a2和a5的分类标签M1、M2和M5输入至编码器中进行编码处理,分别生成第一标签向量m1、m2和m5,将关系抽取编码向量a1、a2、a5分别与第一标签向量m1、m2、m5相结合生成包含标签信息的关系抽取编码向量am1、am2、am5,再将具有语义关系的包含标签信息的关系抽取编码向量am1和am2相结合,生成第一实体编码向量(am1+ am2),将具有语义关系的包含标签信息的关系抽取编码向量am1和am5相结合,生成第一实体编码向量(am1+ am5)。

[0137] 本实施例基于所述关系抽取编码向量以及关系抽取编码向量的分类标签,生成第一实体编码向量,有助于提高关系抽取任务的执行效率与结果准确率。

[0138] S234、基于所述共指消解编码向量的分数生成第二实体编码向量。

[0139] 具体地,所述步骤S234包括步骤S2341至步骤S2342。

[0140] S2341、基于所述共指消解编码向量的分数对所述共指消解编码向量进行分类预测处理,获得所述共指消解编码向量的分类标签。

[0141] 具体地,分类预测处理是指对关系抽取编码基于其分数按照其对应词语的属性进行分类,并获得分类标签。

[0142] 例如,基于共指消解编码向量a1、a3、a4和a5的分数135分、129分、139分和145分按

照每个共指消解编码向量对应词语的属性对上述共指消解编码向量进行分类预测处理,获得共指消解编码向量a1、a3、a4和a5的分类标签分别为N1、N3、N4和N5。

[0143] 本实施例通过分类预测处理获得共指消解编码向量的分类标签,有助于加深模型对候选实体片段的认知,进而提高共指消解的准确率。

[0144] S2342、基于所述共指消解编码向量以及共指消解编码向量的分类标签,生成第二实体编码向量。

[0145] 在实际应用中,可以将所述共指消解编码向量的分类标签进行编码处理,生成第二标签向量;基于所述共指消解编码向量和所述第二标签向量,生成第二实体编码向量。

[0146] 进一步地,共指消解编码向量与第二标签向量相结合,首先生成包含标签信息的共指消解编码向量,再将具有语义关系的两个词语对应的包含标签信息的共指消解编码向量相结合,生成第二实体编码向量,如下所示:

[0147] `span_pair_embeddings = torch.cat([span1_embeddings, span2_embeddings, span1_embeddings*span2_embeddings, span1_label_embedding, span2_label_embedding], -1)`。

[0148] 其中, `torch.cat` 是用于将两个或多个向量拼接在一起的函数, `span_pair_embeddings` 表示第二实体编码向量, `span1_embeddings` 表示共指消解编码向量1, `span2_embeddings` 表示共指消解编码向量2, `span1_label_embedding` 表示共指消解编码向量1的标签向量即第二标签向量1, `span2_label_embedding` 表示共指消解编码向量2的标签向量即第二标签向量2。

[0149] 例如,将共指消解编码向量a1、a3、a4和a5的分类标签N1、N3、N4和N5输入至编码器中进行编码处理,分别生成第二标签向量n1、n3、n4和n5,将共指消解编码向量a1、a3、a4和a5分别与第二标签向量n1、n3、n4和n5相结合生成包含标签信息的共指消解编码向量an1、an3、an4和an5,再将具有语义关系的包含标签信息的共指消解编码向量an1和an3相结合,生成第二实体编码向量(an1+ an3),将具有语义关系的包含标签信息的共指消解编码向量an4和an5相结合,生成第二实体编码向量(an4+ an5)。

[0150] 本实施例基于所述共指消解编码向量以及共指消解编码向量的分类标签,生成第二实体编码向量,有助于提高共指消解任务的执行效率与结果准确率。

[0151] 需要说明的是,分类标签属于候选实体片的特征信息的一种,除此之外,在生成第一实体编码向量、第二实体编码向量时还可以结合其他类型的特征信息,如距离等,可视具体情况而定,本申请对此不做限制。

[0152] 本实施例通过对关系抽取编码向量进行实体识别处理和多层次打分,并结合关系抽取编码向量的分类标签信息生成实体识别编码向量,将实体识别融入至关系抽取和共指消解的任务中,三者相辅相成、内容共享,互相为彼此提供更为丰富的内容信息、特征信息,可以有效提高实体识别任务、关系抽取任务及共指消解任务的准确率。

[0153] S240、基于所述第一实体编码向量和第二实体编码向量分别通过关系抽取模型和共指消解模型进行关系抽取处理和共指消解处理,获得关系抽取结果和共指消解结果。

[0154] 具体地,经过关系抽取模型对第一实体编码向量进行打分、分类预测处理后,生成实体关系三元组形式的关系抽取结果,如(实体1、关系1、实体2),即表示实体1与实体2之间存在关系1的关系;经过共指消解模型第二实体编码向量进行打分、共指消解处理后,生成

实体共指关系三元组形式的共指消解结果,如(实体3、共指、实体4),即表示实体3与实体4之间存在共指关系。

[0155] 例如,将第一实体编码向量(am1+ am2)、(am1+ am5)输入至关系抽取模型中进行打分,得到第一实体编码向量(am1+ am2)的分数为70分,(am1+ am5)的分数为73分,经过分类预测处理,得到第一实体编码向量(am1+ am2)属于“x1关系”类别,(am1+ am5)属于“x3关系”类别,生成最终的关系抽取处理结果(a1,x1关系,a2)、(a1,x3关系,a5),将第二实体编码向量(an1+ an3)、(an4+ an5)输入至共指消解模型中进行打分,得到第二实体编码向量(an1+ an3)的分数为66分、(an4+ an5)的分数为49分再次进行共指消解处理,经过分类预测处理,得到第二实体编码向量(an1+ an3)属于“共指关系”类别,生成最终的共指消解处理结果(a1,共指,a3)。

[0156] 本实施例通过将实体编码向量进行关系抽取、共指消解处理,可以基于候选实体片段的内容信息、特征信息等对候选实体片段编码向量首次进行关系抽取、共指消解处理得到的结果进行修正,提高关系抽取任务以及共指消解任务的准确率。

[0157] S250、基于所述关系抽取结果和共指消解结果确定所述候选实体片段中词语之间的关系。

[0158] 具体地,候选片段词语之间关系可以为时间关系、位置关系、用于关系、共指关系等等,可视具体情况而定,本申请对此不做限制。

[0159] 例如,基于关系抽取处理结果(a1,x1关系,a2)得到词语A1与词语A2之间的关系为x1关系,基于关系抽取处理结果(a1,x3关系,a5)得到词语A1与词语A5之间的关系为x3关系,基于共指消解处理结果(a1,共指,a3)得到词语A1与词语A3之间的关系为共指关系。

[0160] 本实施例提供的基于多任务模型的词语处理方法,通过关系抽取模型和共指消解模型分别对候选实体片段编码向量进行修剪处理,获得关系抽取编码向量和共指消解编码向量,实现了基于不同的任务需求对候选实体片段的初步筛选,再通过对关系抽取编码向量和共指消解编码向量进行实体识别处理,获得第一实体编码向量和第二实体编码向量,可以进一步地、更深层次的加强对于关系抽取编码向量、共指消解编码向量对应的候选实体片段的理解,为关系抽取任务和共指消解任务的执行提供基础,最终通过对第一实体编码向量、第二实体编码向量分别进行关系抽取处理和共指消解处理,并基于获得的关系抽取结果和共指消解结果确定候选实体片段中词语的关系,实现了关系抽取模型、共指消解模型、实体识别模型的有机结合,实现了关系抽取任务、共指消解任务、实体识别任务的有机结合,实现了上述三个任务中信息的共享,有效提高关系抽取、共指消解、实体识别的正确率、召回率和正确率与召回率之间的加权调和平均值,有效提高词语关系确定的准确率。

[0161] 如图3所示,本申请提供一种基于多任务模型的词语处理方法,包括步骤S310至步骤S3100,在本实施例中结合具体的例子进行详细说明。

[0162] S310、获取候选实体片段,并基于所述候选实体片段生成候选实体片段编码向量。

[0163] 在本实施例中,假设目标语句包括“小明在图书馆看书,他六点才回家”,经过分词处理后得到分词处理结果“小明”、“在”、“图书馆”、“看”、“书”、“他”、“六点”、“才”、“回家”,在上述分词处理结果中抽取目标词语组成词语集合[小明、在、图书馆、看、书、他、六点、回家],并将上述词语结合作为候选实体片段。

[0164] 将候选实体片段[小明、在、图书馆、看、书、他、六点、回家]输入至编码器中,生成

候选实体片段编码向量[b1、b2、b3、b4、b5、b6、b7、b8]。

[0165] S320、通过关系抽取模型对候选实体片段编码向量进行修剪处理,获得关系抽取编码向量。

[0166] 在本实施例中,通过关系抽取模型的前馈神经网络对候选实体片段编码向量[b1、b2、b3、b4、b5、b6、b7、b8]进行打分,得到候选实体片段编码向量b1为88分、b2为49分、b3为79分、b4为54分、b5为67分、b6为70分、b7为50分、b8为61分。

[0167] 选择分数大于60分的候选实体编码向量作为关系抽取编码向量,那么在本实施例中关系抽取编码向量包括[b1、b3、b5、b6、b8]。

[0168] S330、通过共指消解模型对候选实体片段编码向量进行修剪处理,获得共指消解编码向量。

[0169] 在本实施例中,通过共指消解模型的前馈神经网络对候选实体片段编码向量[b1、b2、b3、b4、b5、b6、b7、b8]进行打分,得到候选实体片段编码向量b1为88分、b2为40分、b3为44分、b4为50分、b5为52分、b6为83分、b7为50分、b8为51分。

[0170] 选择分数大于60分的候选实体编码向量作为共指消解编码向量,那么在本实施例中共指消解编码向量包括[b1、b6]。

[0171] S340、将所述关系抽取编码向量输入至实体识别模型中,通过所述前馈神经网络进行基础打分与分类打分,获得关系抽取编码向量的分数。

[0172] 在本实施例中,将上述关系抽取编码向量输入至实体识别模型中,通过前馈神经网络进行基础打分,得到关系抽取编码向量b1、b3、b5、b6、b8的基础分数分别为60分、61分、63分、63分、65分。

[0173] 通过前馈神经网络对上述关系抽取编码向量进行分类打分,得到关系抽取编码向量b1、b3、b5、b6、b8的分类分数分别为65分、63分、60分、66分、64分。

[0174] 将每一个关系抽取编码向量的基础分数与分类分数相加,得到关系抽取编码向量b1、b3、b5、b6、b8的分数分别为125分、124分、123分、129分、129分。

[0175] S350、将所述共指消解编码向量输入至实体识别模型中,通过所述前馈神经网络进行基础打分与分类打分,获得共指消解编码向量的分数。

[0176] 在本实施例中,将上述共指消解编码向量输入至实体识别模型中,通过前馈神经网络进行基础打分,得到共指消解编码向量b1、b6的基础分数分别为76分、67分。

[0177] 通过前馈神经网络对上述共指消解编码向量进行分类打分,得到共指消解编码向量b1、b6的分类分数分别为65分、64分。

[0178] 将每一个共指消解编码向量的基础分数与分类分数相加,得到共指消解编码向量b1、b6的分数分别为141分、131分。

[0179] S360、基于所述关系抽取编码向量的分数对所述关系抽取编码向量进行分类预测处理,获得所述关系抽取编码向量的分类标签,基于所述关系抽取编码向量以及关系抽取编码向量的分类标签,生成第一实体编码向量。

[0180] 在本实施例中,基于关系抽取编码向量b1、b3、b5、b6、b8对应的词语“小明”、“图书馆”、“书”、“他”、“回家”及其上一步得到的分数进行分类预测处理,得到关系抽取编码向量b1、b3、b5、b6、b8的分类标签分别为“施事者”、“地点”、“受事者”、“施事者”、“行为”。

[0181] 将上述分类标签输入至编码器中进行编码处理,得到关系抽取编码向量b1、b3、

b5、b6、b8对应的分类标签向量 p_1 、 p_3 、 p_5 、 p_6 、 p_8 ，将每一个关系抽取编码向量与其对应的分类标签向量相结合，得到包含标签信息的关系抽取编码向量 bp_1 、 bp_3 、 bp_5 、 bp_6 、 bp_8 。

[0182] 基于步骤S320中关系抽取处理的结果将包含标签信息的关系抽取编码向量进行组合，得到第一实体编码向量 (bp_1+bp_3) 、 (bp_1+bp_5) 、 (bp_6+bp_8) 。

[0183] S370、基于所述共指消解编码向量的分数对所述共指消解编码向量进行分类预测处理，获得所述共指消解编码向量的分类标签，基于所述共指消解编码向量以及共指消解编码向量的分类标签，生成第二实体编码向量。

[0184] 在本实施例中，基于共指消解编码向量 b_1 、 b_6 对应的词语“小明”、“他”及其上一步得到的分数进行分类预测处理，得到共指消解编码向量 b_1 、 b_6 的分类标签均为“施事者”。

[0185] 将上述分类标签输入至编码器中进行编码处理，得到共指消解编码向量 b_1 、 b_6 对应的分类标签向量 q_1 、 q_6 ，将每一个共指消解编码向量与其对应的分类标签向量相结合，得到包含标签信息的共指消解编码向量 bq_1 、 bq_6 。

[0186] 基于步骤S330中共指消解处理的结果将包含标签信息的共指消解编码向量进行组合，得到第二实体编码向量 (bq_1+bq_6) 。

[0187] S380、基于所述第一实体编码向量通过关系抽取模型进行关系抽取处理，获得关系抽取结果。

[0188] 在本实施例中，基于第一实体编码向量 (bp_1+bp_3) 、 (bp_1+bp_5) 、 (bp_6+bp_8) 再次进行关系抽取处理，得到关系抽取结果 $(b_1, \text{场所}, b_3)$ 、 $(b_1, \text{施事}, b_5)$ 、 $(b_1, \text{时间}, b_8)$ 。

[0189] S390、基于所述第二实体编码向量通过共指消解模型进行共指消解处理，获得共指消解结果。

[0190] 在本实施例中，基于第二实体编码向量 (bq_1+bq_6) 再次进行共指消解处理，得到共指消解结果 $(b_1, \text{共指}, b_6)$ 。

[0191] S3100、基于所述关系抽取结果和共指消解结果确定所述候选实体片段中词语之间的关系。

[0192] 在本实施例中，基于关系抽取处理结果 $(b_1, \text{场所}, b_3)$ 可以确定“小明”与“图书馆”之间存在场馆关系，基于关系抽取处理结果 $(b_1, \text{施事}, b_5)$ 可以确定“小明”与“书”之间存在施事关系，基于关系抽取处理结果 $(b_1, \text{时间}, b_8)$ 可以确定“小明”与“回家”之间存在时间关系，基于共指消解处理结果 $(b_1, \text{共指}, b_6)$ 可以确定“小明”与“他”之间存在共指关系。

[0193] 本实施例提供的基于多任务模型的词语处理方法，实现了关系抽取模型、共指消解模型、实体识别模型的有机结合，实现了关系抽取任务、共指消解任务、实体识别任务的有机结合，实现了上述三个任务中信息的共享，有效提高关系抽取、共指消解、实体识别的正确率、召回率和加权调和平均值，有效提高词语关系确定的准确率。

[0194] 如图4所示，本实施例提供了一种基于多任务模型的词语处理装置，包括：

[0195] 片段获取模块410，被配置为获取候选实体片段，并基于所述候选实体片段生成候选实体片段编码向量；

[0196] 片段修剪模块420，被配置为分别通过关系抽取模型和共指消解模型对所述候选实体片段编码向量进行修剪处理，获得关系抽取编码向量和共指消解编码向量；

[0197] 实体识别模块430，被配置为将所述关系抽取编码向量和所述共指消解编码向量输入至实体识别模型中进行实体识别处理，获得第一实体编码向量和第二实体编码向量；

[0198] 关系处理模块440,被配置为基于所述第一实体编码向量和第二实体编码向量分别通过关系抽取模型和共指消解模型进行关系抽取处理和共指消解处理,获得关系抽取结果和共指消解结果;

[0199] 关系确定模块450,被配置为基于所述关系抽取结果和共指消解结果确定所述候选实体片段中词语之间的关系。

[0200] 可选地,所述关系抽取模型与所述实体识别模型共用一个前馈神经网络;

[0201] 所述片段修剪模块420,进一步被配置为:

[0202] 通过所述前馈神经网络对所述候选实体片段编码向量进行打分,获得候选实体片段编码向量的分数,并将分数大于或等于预设阈值的候选实体片段编码向量作为关系抽取编码向量。

[0203] 可选地,所述共指消解模型与所述实体识别模型共用前馈神经网络;

[0204] 所述片段修剪模块420,进一步被配置为:

[0205] 通过所述前馈神经网络对所述候选实体片段编码向量进行打分,获得候选实体片段编码向量的分数,并将分数大于或等于预设阈值的候选实体片段编码向量作为共指消解编码向量。

[0206] 可选地,所述实体识别模块430,进一步被配置为:

[0207] 将所述关系抽取编码向量输入至实体识别模型中,通过所述前馈神经网络进行基础打分与分类打分,获得关系抽取编码向量的分数;

[0208] 基于所述关系抽取编码向量的分数生成第一实体编码向量。

[0209] 可选地,所述实体识别模块430,更进一步被配置为:

[0210] 基于所述关系抽取编码向量的分数对所述关系抽取编码向量进行分类预测处理,获得所述关系抽取编码向量的分类标签;

[0211] 基于所述关系抽取编码向量以及关系抽取编码向量的分类标签,生成第一实体编码向量。

[0212] 可选地,所述实体识别模块430,更进一步被配置为:

[0213] 将所述关系抽取编码向量的分类标签进行编码处理,生成第一标签向量;

[0214] 基于所述关系抽取编码向量和所述第一标签向量,生成第一实体编码向量。

[0215] 可选地,所述实体识别模块430,更进一步被配置为:

[0216] 所述共指消解编码向量输入至实体识别模型中,通过所述前馈神经网络进行基础打分与分类打分,获得共指消解编码向量的分数;

[0217] 基于所述共指消解编码向量的分数生成第二实体编码向量。

[0218] 可选地,所述实体识别模块430,更进一步被配置为:

[0219] 基于所述共指消解编码向量的分数对所述共指消解编码向量进行分类预测处理,获得所述共指消解编码向量的分类标签;

[0220] 基于所述共指消解编码向量以及共指消解编码向量的分类标签,生成第二实体编码向量。

[0221] 可选地,所述实体识别模块430,更进一步被配置为:

[0222] 将所述共指消解编码向量的分类标签进行编码处理,生成第二标签向量;

[0223] 基于所述共指消解编码向量和所述第二标签向量,生成第二实体编码向量。

[0224] 本申请提供的词语关系确定装置,通过关系抽取模型和共指消解模型分别对候选实体片段编码向量分别进行修剪处理,获得关系抽取编码向量和共指消解编码向量,实现了基于不同的任务需求对候选实体片段的初步筛选,再通过对关系抽取编码向量和共指消解编码向量进行实体识别处理,获得第一实体编码向量和第二实体编码向量,可以进一步地、更深层次的加强对于关系抽取编码向量、共指消解编码向量对应的候选实体片段的理解,为关系抽取任务和共指消解任务的执行提供基础,最终通过对第一实体编码向量、第二实体编码向量分别进行关系抽取处理和共指消解处理,并基于获得的关系抽取结果和共指消解结果确定候选实体片段中词语的关系,实现了关系抽取模型、共指消解模型、实体识别模型的有机结合,也即实现了关系抽取任务、共指消解任务、实体识别任务的有机结合,实现了上述三个任务中信息的共享,有效提高关系抽取、共指消解、实体识别的确率、召回率和加权调和平均值,有效提高词语关系确定的准确率。

[0225] 如图5所示,图5是示出了根据本说明书一实施例的计算设备500的结构框图。该计算设备500的部件包括但不限于存储器510和处理器520。处理器520与存储器510通过总线530相连接,数据库550用于保存数据。

[0226] 计算设备500还包括接入设备540,接入设备540使得计算设备500能够经由一个或多个网络560通信。这些网络的示例包括公用交换电话网(PSTN)、局域网(LAN)、广域网(WAN)、个域网(PAN)或诸如因特网的通信网络的组合。接入设备540可以包括有线或无线的任何类型的网络接口(例如,网络接口卡(NIC))中的一个或多个,诸如IEEE802.11无线局域网(WLAN)无线接口、全球微波互联接入(Wi-MAX)接口、以太网接口、通用串行总线(USB)接口、蜂窝网络接口、蓝牙接口、近场通信(NFC)接口,等等。

[0227] 在本说明书的一个实施例中,计算设备500的上述部件以及图5中未示出的其他部件也可以彼此相连接,例如通过总线。应当理解,图5所示的计算设备结构框图仅仅是出于示例的目的,而不是对本说明书范围的限制。本领域技术人员可以根据需要,增添或替换其他部件。

[0228] 计算设备500可以是任何类型的静止或移动计算设备,包括移动计算机或移动计算设备(例如,平板计算机、个人数字助理、膝上型计算机、笔记本计算机、上网本等)、移动电话(例如,智能手机)、可佩戴的计算设备(例如,智能手表、智能眼镜等)或其他类型的移动设备,或者诸如台式计算机或PC的静止计算设备。计算设备500还可以是移动式或静止式的服务器。

[0229] 本申请一实施例还提供一种计算设备,包括存储器、处理器及存储在存储器上并可在处理器上运行的计算机指令,所述处理器执行所述指令时实现以下步骤:

[0230] 获取候选实体片段,并基于所述候选实体片段生成候选实体片段编码向量;

[0231] 分别通过关系抽取模型和共指消解模型对所述候选实体片段编码向量进行修剪处理,获得关系抽取编码向量和共指消解编码向量;

[0232] 将所述关系抽取编码向量和所述共指消解编码向量输入至实体识别模型中进行实体识别处理,获得第一实体编码向量和第二实体编码向量;

[0233] 基于所述第一实体编码向量和第二实体编码向量分别通过关系抽取模型和共指消解模型进行关系抽取处理和共指消解处理,获得关系抽取结果和共指消解结果;

[0234] 基于所述关系抽取结果和共指消解结果确定所述候选实体片段中词语之间的关

系。

[0235] 本申请一实施例还提供一种计算机可读存储介质,其存储有计算机指令,该指令被处理器执行时实现如前所述基于多任务模型的词语处理方法的步骤。

[0236] 上述为本实施例的一种计算机可读存储介质的示意性方案。需要说明的是,该存储介质的技术方案与上述的基于多任务模型的词语处理的技术方案属于同一构思,存储介质的技术方案未详细描述的细节内容,均可以参见上述基于多任务模型的词语处理的技术方案的描述。

[0237] 所述计算机指令包括计算机程序代码,所述计算机程序代码可以为源代码形式、对象代码形式、可执行文件或某些中间形式等。所述计算机可读介质可以包括:能够携带所述计算机程序代码的任何实体或装置、记录介质、U盘、移动硬盘、磁碟、光盘、计算机存储器、只读存储器(ROM,Read-Only Memory)、随机存取存储器(RAM,Random Access Memory)、电载波信号、电信信号以及软件分发介质等。需要说明的是,所述计算机可读介质包含的内容可以根据司法管辖区内立法和专利实践的要求进行适当的增减,例如在某些司法管辖区,根据立法和专利实践,计算机可读介质不包括电载波信号和电信信号。

[0238] 需要说明的是,对于前述的各方法实施例,为了简便描述,故将其都表述为一系列的动作组合,但是本领域技术人员应该知悉,本申请并不受所描述的动作顺序的限制,因为依据本申请,某些步骤可以采用其它顺序或者同时进行。其次,本领域技术人员也应该知悉,说明书中所描述的实施例均属于优选实施例,所涉及的动作和模块并不一定都是本申请所必须的。

[0239] 在上述实施例中,对各个实施例的描述都各有侧重,某个实施例中沒有详述的部分,可以参见其它实施例的相关描述。

[0240] 以上公开的本申请优选实施例只是用于帮助阐述本申请。可选实施例并没有详尽叙述所有的细节,也不限制该发明仅为所述的具体实施方式。显然,根据本说明书的内容,可作很多的修改和变化。本说明书选取并具体描述这些实施例,是为了更好地解释本申请的原理和实际应用,从而使所属技术领域技术人员能很好地理解和利用本申请。本申请仅受权利要求书及其全部范围和等效物的限制。

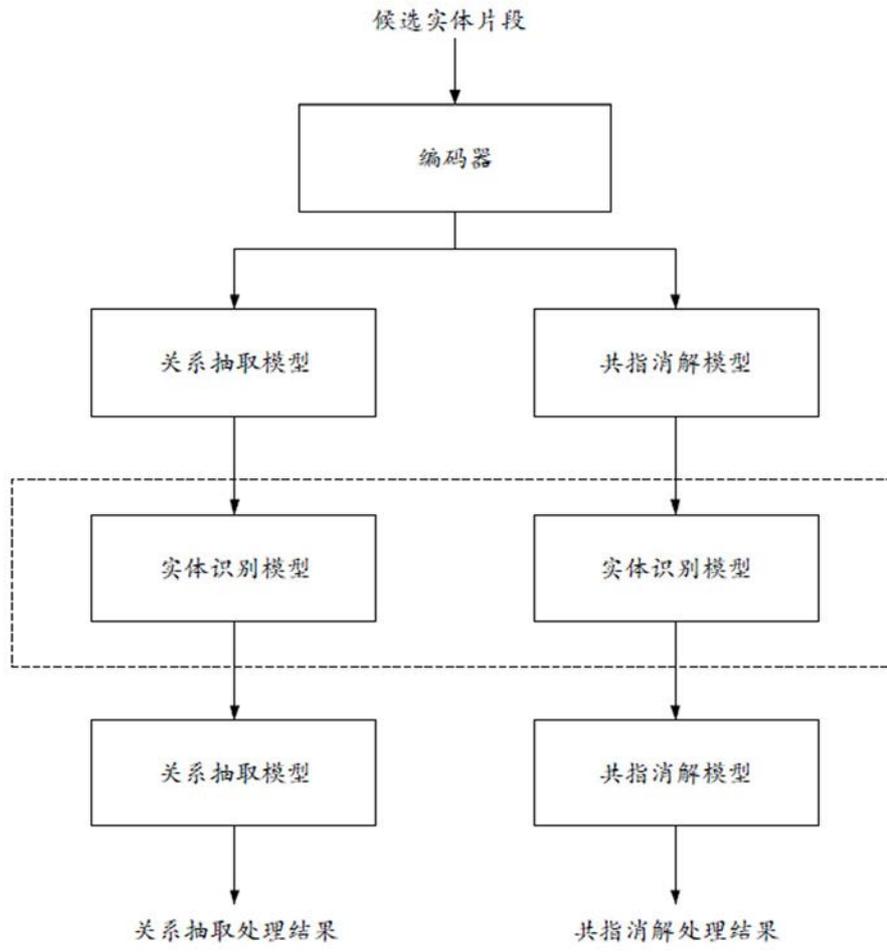


图1

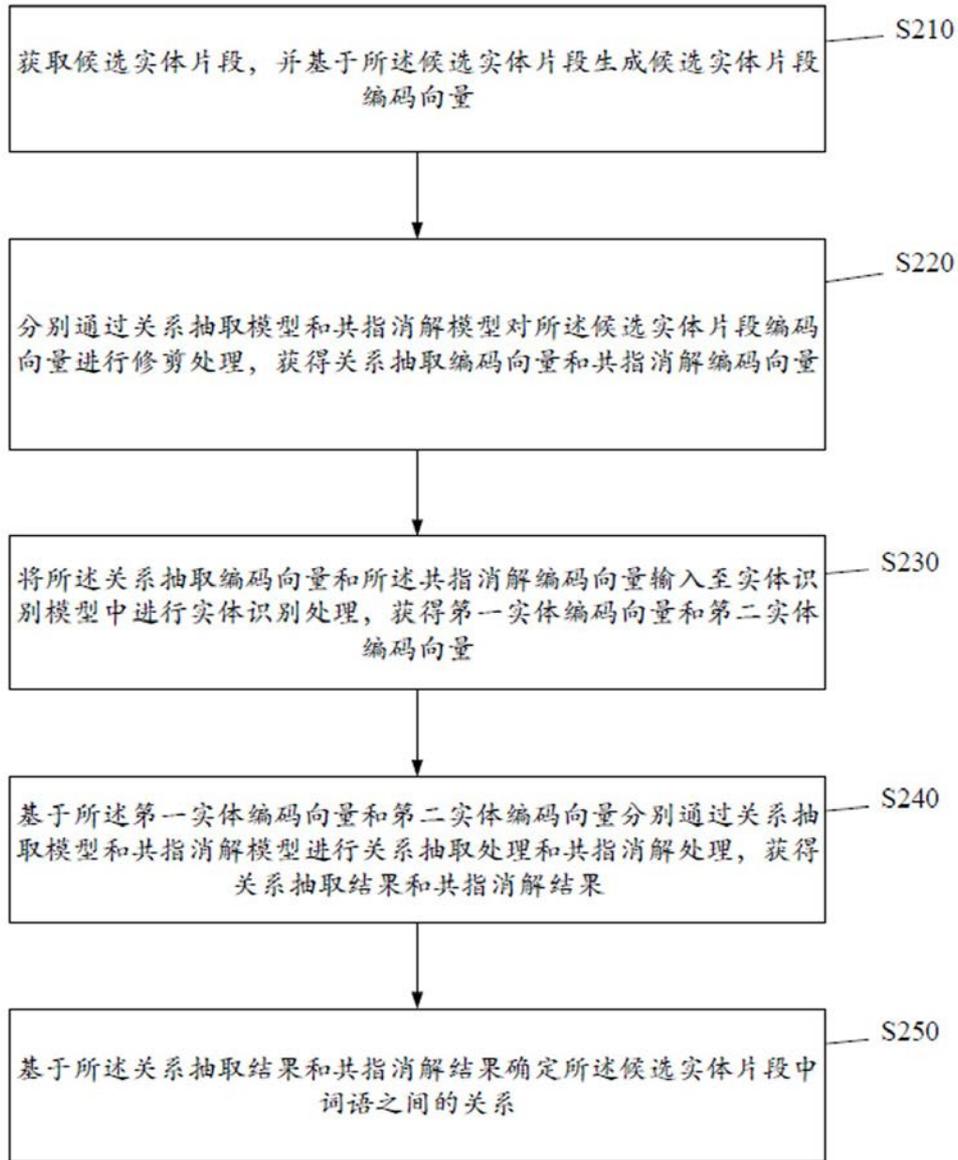


图2

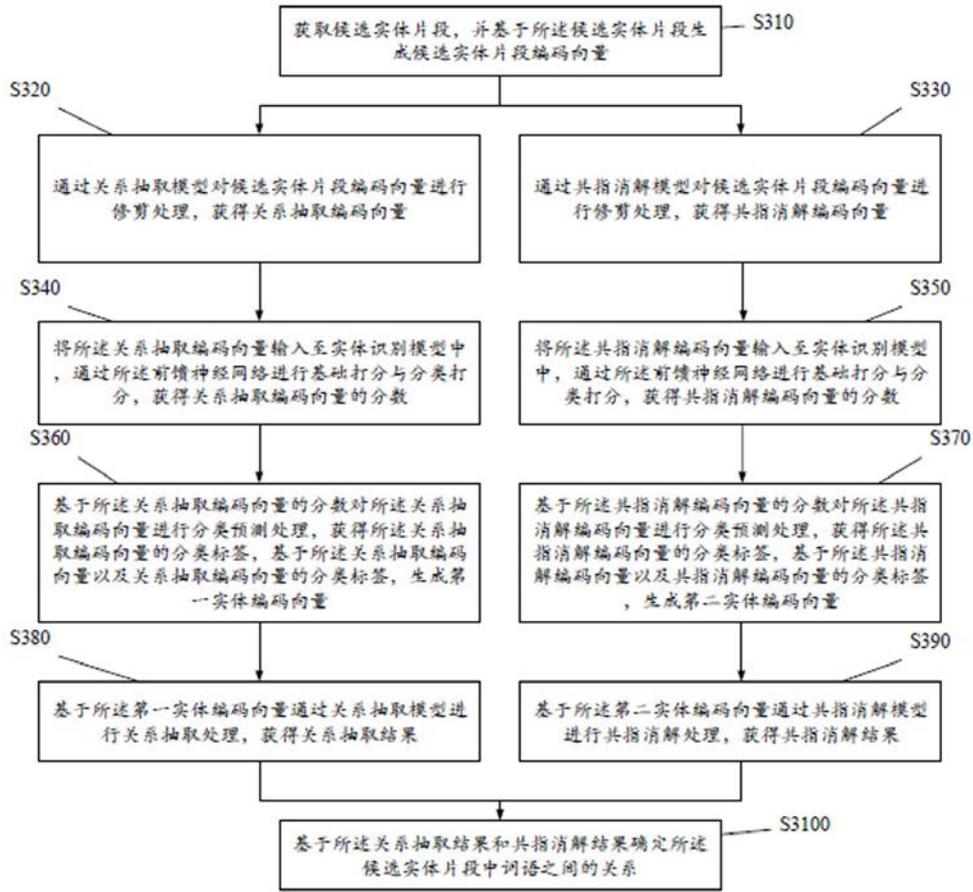


图3

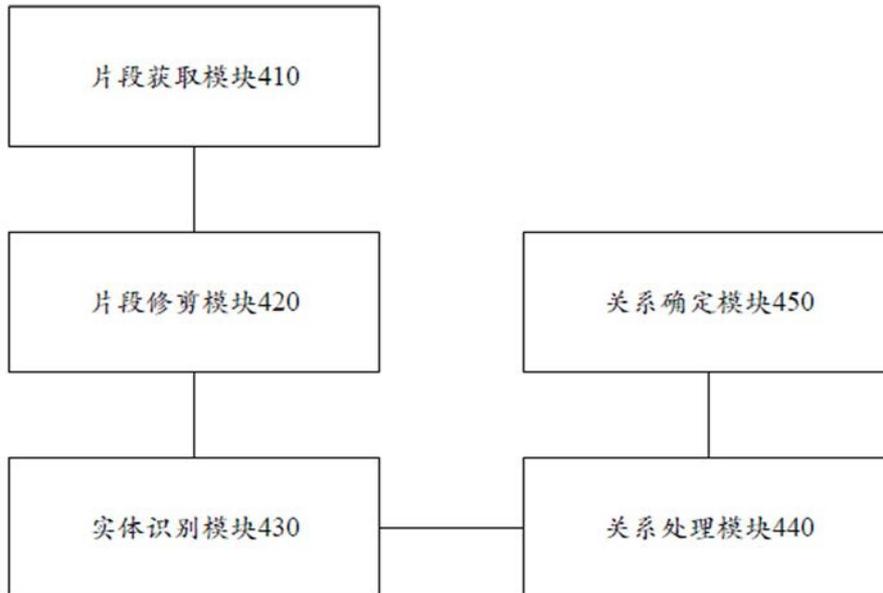


图4

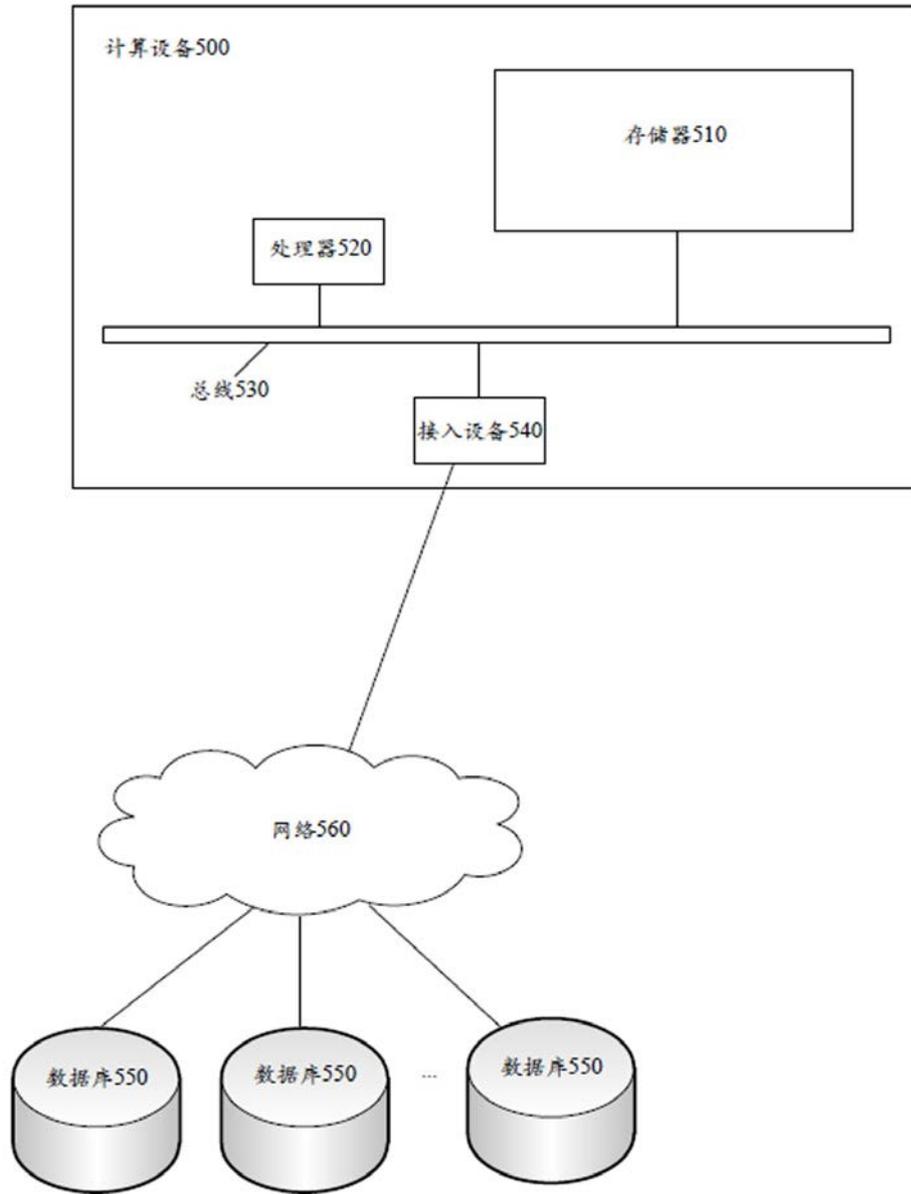


图5