



(12) 发明专利申请

(10) 申请公布号 CN 113312445 A

(43) 申请公布日 2021.08.27

(21) 申请号 202110866035.3

(22) 申请日 2021.07.29

(71) 申请人 阿里云计算有限公司

地址 310012 浙江省杭州市西湖区转塘科技经济区块12号

(72) 发明人 汪诚愚

(74) 专利代理机构 北京太合九思知识产权代理有限公司 11610

代理人 刘戈 曹威

(51) Int.Cl.

G06F 16/31 (2019.01)

G06F 16/332 (2019.01)

G06F 16/33 (2019.01)

G06F 16/35 (2019.01)

G06K 9/62 (2006.01)

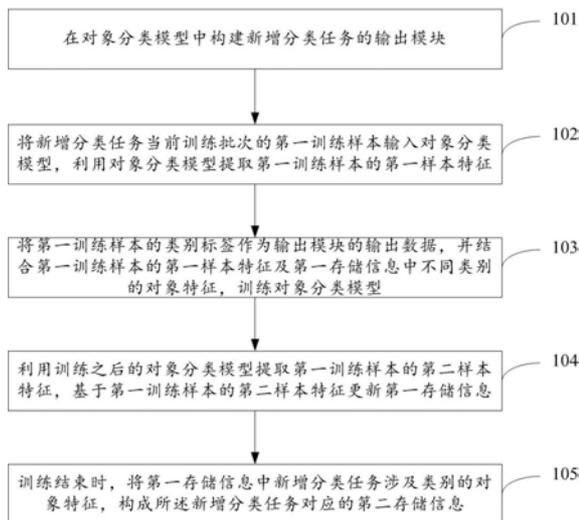
权利要求书3页 说明书21页 附图9页

(54) 发明名称

数据处理方法、模型构建方法、分类方法及计算设备

(57) 摘要

本申请实施例提供一种数据处理方法、模型构建方法、分类方法及计算设备。其中,在对象分类模型中构建新增分类任务的输出模块;将当前训练批次的第一训练样本输入对象分类模型,提取第一训练样本的第一样本特征;将第一训练样本的类别标签作为输出模块的输出数据,结合第一训练样本的第一样本特征及第一存储信息中不同类别的对象特征训练对象分类模型;利用训练之后的对象分类模型提取第一训练样本的第二样本特征,基于第一训练样本的第二样本特征更新第一存储信息;训练结束时,将第一存储信息中新增分类任务涉及类别的对象特征,构成新增分类任务对应的第二存储信息。本申请实施例提供的技术方案降低了模型训练成本并保证了模型分类准确度。



1. 一种数据处理方法,其特征在于,包括:

在对象分类模型中构建新增分类任务的输出模块;

将所述新增分类任务当前训练批次的第一训练样本输入所述对象分类模型,利用所述对象分类模型提取所述第一训练样本的第一样本特征;

将所述第一训练样本的类别标签作为所述输出模块的输出数据,并结合所述第一训练样本的第一样本特征及第一存储信息中不同类别的对象特征,训练所述对象分类模型;

利用训练之后的所述对象分类模型提取所述第一训练样本的第二样本特征,并基于所述第一训练样本的第二样本特征更新所述第一存储信息;

针对所述新增分类任务的训练结束时,将所述第一存储信息中所述新增分类任务涉及类别的对象特征,构成所述新增分类任务对应的第二存储信息;其中,所述第二存储信息中不同类别的对象特征用于参与利用所述对象分类模型对属于所述新增分类任务的待分类对象的分类操作。

2. 根据权利要求1所述的方法,其特征在于,所述将所述第一训练样本的类别标签作为所述输出模块的输出数据,并结合所述第一训练样本的第一样本特征以及第一存储信息中不同类别的对象特征,训练所述对象分类模型包括:

将所述第一训练样本的第一样本特征分别与第一存储信息中不同类别的对象特征进行融合,获得第一融合特征;

将所述第一融合特征作为所述新增分类任务的输出模块的输入数据,以及将所述第一训练样本的类别标签作为所述新增分类任务的输出模块的输出数据,训练所述对象分类模型。

3. 根据权利要求1所述的方法,其特征在于,所述对象分类模型包括至少一个原分类任务分别对应的输出模块;

所述对象分类模型针对所述至少一个原分类任务按照如下方式预先训练获得:

利用所述对象分类模型提取所述至少一个原分类任务对应的第二训练样本的第一样本特征;

计算属于同一个类别的第二训练样本的平均第一样本特征,作为该类别的对象特征,并将计算获得的不同类别的对象特征构成第一存储信息;

将所述第二训练样本对应的类别标签作为其对应输出模块的输出数据,并结合所述第二训练样本的第一样本特征以及所述第一存储信息,训练所述对象分类模型;

针对所述至少一个原分类任务训练结束时,将所述第一存储信息中所述至少一个原分类任务各自涉及类别的对象特征,构成所述至少一个原分类任务各自对应的第二存储信息,或者,针对所述至少一个原分类任务训练结束时,抽取所述第一存储信息中不同类别的对象特征,构成所述至少一个原分类任务对应的第二存储信息。

4. 根据权利要求3所述的方法,其特征在于,所述将所述第二训练样本对应的类别标签作为其对应输出模块的输出数据,并结合所述第二训练样本的第一样本特征以及所述第一存储信息,训练所述对象分类模型包括:

将当前训练批次的第二训练样本的第一样本特征与第一存储信息中不同类别的对象特征进行融合,获得第二融合特征;

将所述第二融合特征输入对应的输出模块,以及将所述第二训练样本的类别标签作为

对应输出模块的输出数据,训练所述对象分类模型;

利用训练之后的所述对象分类模型提取所述第二训练样本的第二样本特征,并基于所述第二训练样本的第二样本特征,更新所述第一存储信息中相应类别的对象特征。

5. 根据权利要求1所述的方法,其特征在于,所述基于所述第一训练样本的第二样本特征更新所述第一存储信息包括:

基于所述第一训练样本的第二样本特征,计算属于同一类别的第一训练样本的平均第二样本特征;

将属于同一类别的平均第二样本特征与第一存储信息中的对象特征进行加权求和,并利用加权求和结果替换所述第一存储信息中对应类别的对象特征。

6. 根据权利要求2所述的方法,其特征在于,所述将所述第一训练样本的第一样本特征与第一存储信息中不同类别的对象特征进行融合,获得第一融合特征包括:

根据所述第一训练样本的第一样本特征,确定第一存储信息中不同类别的对象特征对应的第二权重系数;

按照各自的第二权重系数,将所述第一存储信息中不同类别的对象特征进行加权求和,获得第一加权特征;

将所述第一加权特征与所述第一训练样本的第一样本特征累加,获得第一融合特征。

7. 一种模型构建方法,其特征在于,包括:

依次搭建输入模块、特征提取模块、特征融合模块及至少一个原分类任务对应的输出模块,获得对象分类模型;

对应设置所述对象分类模型的第一存储信息及所述至少一个原分类任务各自的第二存储信息;其中,所述第一存储信息存储不同类别的对象特征,基于所述至少一个原分类任务的第二训练样本提取的样本特征获得;任一原分类任务的第二存储信息存储从所述第一存储信息中抽取的所述原分类任务涉及类别的对象特征;

根据分类任务扩展需求,在所述对象分类模型中搭建新增分类任务的输出模块;

对应设置所述新增分类任务的第二存储信息;其中,所述第一存储信息根据所述新增分类任务的第一训练样本所提取的样本特征进行更新,所述新增分类任务的第二存储信息存储从所述第一存储信息中抽取的所述新增分类任务涉及类别的对象特征。

8. 一种分类方法,其特征在于,包括:

确定待分类对象所属的目标分类任务;

将所述待分类对象输入对象分类模型,利用所述对象分类模型提取所述待分类对象的目标对象特征;

基于所述目标对象特征及所述目标分类任务对应第二存储信息中不同类别的对象特征,利用所述目标分类任务对应的输出模块,识别所述待分类对象的分类结果。

9. 一种数据处理方法,其特征在于,包括:

在对象分类模型中构建新增分类任务的输出模块;

提取所述新增分类任务的第一训练样本的第一样本特征;

基于所述第一训练样本的第一样本特征,更新第一存储信息中相应类别的对象特征;

将所述第一训练样本的类别标签作为所述输出模块的输出数据,并结合所述第一训练样本的第一样本特征及所述第一存储信息中不同类别的对象特征,训练所述对象分类模

型；

针对所述新增分类任务训练结束时，将所述第一存储信息中所述新增分类任务涉及类别的对象特征，构成所述新增分类任务对应的第二存储信息；其中，所述第二存储信息中不同类别的对象特征用于参与利用所述对象分类模型对所述新增分类任务的待分类对象的分类操作。

10. 一种文本分类方法，其特征在于，包括：

确定待分类文本所属的目标分类任务；

将所述待分类文本输入文本分类模型，利用所述文本分类模型提取所述待分类文本的文本特征；

基于所述文本特征及所述目标分类任务对应第二存储信息中不同类别的文本特征，利用所述目标分类任务对应的输出模块，识别所述待分类文本的分类结果。

11. 根据权利要求10所述的方法，其特征在于，所述待分类文本为针对电商产品的评论数据；所述分类结果为所述待分类文本所属的情感类别；所述确定待分类文本所属的目标分类任务包括：根据确定所述评论数据所属的目标产品类目，确定所述目标产品类目对应的目标分类任务；

所述方法还包括：统计所述目标产品类目中属于同一情感类别的评论数据数量；基于不同情感类别的评论数据数量，生成提示信息；

或者，

所述待分类文本为人机对话中的用户输入文本，所述分类结果为所述用户输入文本匹配的标准文本；所述方法还包括：

基于所述标准文本查找对应的应答内容；输出所述应答内容；

或者，

所述待分类文本为人机对话中的用户输入文本，所述分类结果为所述用户输入文本匹配的应答内容；所述方法还包括：

输出所述应答内容。

12. 一种计算设备，其特征在于，包括处理组件以及存储组件；

所述存储组件存储一个或多个计算机指令；所述一个或多个计算机指令用以被所述处理组件调用执行，以实现如权利要求1~6任一项所述的数据处理方法，或者实现如权利要求7所述的模型构建方法、或者实现如权利要求8所述的分类方法或者实现如权利要求9所述的数据处理方法。

13. 一种计算机存储介质，其特征在于，存储计算机程序，所述计算机程序被计算机执行时实现如权利要求1~6任一项所述的数据处理方法，或者实现如权利要求7所述的模型构建方法、或者实现如权利要求8所述的分类方法或者实现如权利要求9所述的数据处理方法。

14. 一种计算机程序产品，其特征在于，包括计算机程序，所述计算机程序被计算机执行时实现如权利要求1~6任一项所述的数据处理方法，或者实现如权利要求7所述的模型构建方法、或者实现如权利要求8所述的分类方法或者实现如权利要求9所述的数据处理方法。

数据处理方法、模型构建方法、分类方法及计算设备

技术领域

[0001] 本申请实施例涉及计算机应用技术领域,尤其涉及一种数据处理方法、模型构建方法、分类方法及计算设备。

背景技术

[0002] 在计算机应用技术领域,常涉及文本分类、图像分类、音频分类等数据对象的分类需求,目前,多是采用机器学习方式进行对象分类,通过学习已有数据,训练得到相应的对象分类模型,利用对象分类模型即可以对新数据进行分类。

[0003] 然而,随着分类需求的不断扩展,针对同一种类的分类需求也会出现很多分类任务,如果针对每一个分类任务都训练一个对象分类模型,训练成本将会非常大。

发明内容

[0004] 本申请实施例提供一种数据处理方法、分类方法及计算设备,用以解决现有技术中模型训练成本大的技术问题。

[0005] 第一方面,本申请实施例中提供了一种数据处理方法,包括:

在对象分类模型中构建新增分类任务的输出模块;

将所述新增分类任务当前训练批次的第一训练样本输入对象分类模型,利用所述对象分类模型提取所述第一训练样本的第一样本特征;

将所述第一训练样本的类别标签作为所述输出模块的输出数据,并结合所述第一训练样本的第一样本特征及第一存储信息中不同类别的对象特征,训练所述对象分类模型;

利用训练之后的所述对象分类模型提取所述第一训练样本的第二样本特征,并基于所述第一训练样本的第二样本特征更新所述第一存储信息;

针对所述新增分类任务的训练结束时,将所述第一存储信息中所述新增分类任务涉及类别的对象特征,构成所述新增分类任务对应的第二存储信息;其中,所述第二存储信息中不同类别的对象特征用于参与利用所述对象分类模型对所述新增分类任务的待分类对象的分类操作。

[0006] 可选地,所述将所述第一训练样本应的类别标签作为所述输出模块的输出数据,并结合所述第一训练样本的第一样本特征以及第一存储信息中不同类别的对象特征,训练所述对象分类模型包括:将所述第一训练样本的第一样本特征分别与第一存储信息中不同类别的对象特征进行融合,获得第一融合特征;将所述至第一融合特征作为所述新增分类任务的输出模块的输入数据,以及将所述第一训练样本的类别标签作为所述新增分类任务的输出模块的输出数据,训练所述对象分类模型。

[0007] 可选地,所述对象分类模型包括至少一个原分类任务分别对应的输出模块;所述对象分类模型针对所述至少一个原分类任务按照如下方式预先训练获得:

利用所述对象分类模型提取所述至少一个原分类任务对应的第二训练样本的第

一样本特征;计算属于同一个类别的第二训练样本的平均第一样本特征,作为该类别的对象特征,并将计算获得的不同类别的对象特征构成第一存储信息;将所述第二训练样本对应的类别标签作为其对应输出模块的输出数据,并结合所述第二训练样本的第一样本特征以及所述第一存储信息,训练所述对象分类模型;针对所述至少一个原分类任务训练结束时,将所述第一存储信息中所述至少一个原分类任务各自涉及类别的对象特征,构成所述至少一个原分类任务各自的第二存储信息,或者,针对所述至少一个原分类任务训练结束时,抽取所述第一存储信息中不同类别的对象特征,构成所述至少一个原分类任务对应的第二存储信息。

[0008] 可选地,所述将所述第二训练样本对应的类别标签作为其对应输出模块的输出数据,并结合所述第二训练样本的第一样本特征以及所述第一存储信息,训练所述对象分类模型可以包括:

将当前训练批次的第二训练样本的第一样本特征与第一存储信息中不同类别的对象特征进行融合,获得第二融合特征;将所述第二融合特征输入对应的输出模块,以及将所述第二训练样本的类别标签作为对应输出模块的输出数据,训练所述对象分类模型;利用训练之后的所述对象分类模型提取所述第二训练样本的第二样本特征,并基于所述第二训练样本的第二样本特征,更新所述第一存储信息中相应类别的对象特征。

[0009] 可选地,所述基于所述第一训练样本的第二样本特征更新所述第一存储信息包括:

基于所述第一训练样本的第二样本特征,计算属于同一类别的第一训练样本的平均第二样本特征;将属于同一类别的平均第二样本特征与第一存储信息中的对象特征进行加权求和,并利用加权求和结果替换所述第一存储信息中对应类别的对象特征。

[0010] 可选地,所述将所述第一训练样本的第一样本特征与第一存储信息中不同类别的对象特征进行融合,获得第一融合特征包括:

根据所述第一训练样本的第一样本特征,确定第一存储信息中不同类别的对象特征对应的第二权重系数;按照各自的第二权重系数,将所述第一存储信息中不同类别的对象特征进行加权求和,获得第一加权特征;将所述第一加权特征与所述第一训练样本的第一样本特征累加,获得第一融合特征。

[0011] 可选地,将所述第二训练样本的第一样本特征分别与第一存储信息中不同类别的对象特征进行融合,获得第二融合特征包括:

根据所述第二训练样本的第一样本特征,确定所述第一存储信息中不同类别的对象特征对应的第三权重系数;按照各自的第三权重系数,将所述第一存储信息中不同类别的对象特征的进行加权求和,获得第二加权特征;将所述第二加权特征与所述第二训练样本的第一样本特征累加,获得第二融合特征。

[0012] 可选地,根据所述第一训练样本的第一样本特征,确定第一存储信息中不同类别的对象特征对应的第二权重系数包括:

根据所述第一训练样本的第一样本特征,计算第一存储信息中不同类别的对象特征分别与所述第一样本特征的内积;计算不同类别的对象特征对应的内积和;根据每个类别的对象特征的内积在所述内积和中的占比,获得每个类别的对象特征所对应的第二权重系数。

[0013] 可选地,计算属于同一个类别的第二训练样本的平均第一样本特征,作为该类别的对象特征包括:

计算每个原分类任务对应的属于同一个类别的第二训练样本的第一样本特征的平均特征;将至少一个原分类任务分别对应的同一个类别的平均特征进行平均计算,获得同一类别的第二训练样本的平科第一样本特征,并作为该类别对应的对象特征。

[0014] 可选地,所述将所述第一训练样本应的类别标签作为所述输出模块的输出数据,并结合所述第一训练样本的第一样本特征以及第一存储信息中不同类别的对象特征,训练所述对象分类模型之前,所述方法还包括:

基于所述第一训练样本的第一样本特征,更新所述第一存储信息。

[0015] 第二方面,本申请实施例提供了一种模型构建方法,包括:

依次搭建输入模块、特征提取模块、特征融合模块及至少一个原分类任务分别对应的输出模块,获得对象分类模型;

对应设置所述对象分类模型的第一存储信息及所述至少一个原分类任务各自的第二存储信息;其中,所述第一存储信息存储不同类别的对象特征,基于所述至少一个原分类任务的第二训练样本提取的样本特征获得;任一原分类任务的第二存储信息存储从所述第一存储信息中抽取的所述原分类任务涉及类别的对象特征;

根据分类任务扩展需求,在所述对象分类模型中搭建新增分类任务的输出模块;

对应设置所述新增分类任务的第二存储信息;其中,所述第一存储信息根据所述新增分类任务的第一训练样本所提取的样本特征进行更新,所述新增分类任务的第二存储信息存储从所述第一存储信息中抽取的所述新增分类任务涉及类别的对象特征。

[0016] 第三方面,本申请实施例提供了一种分类方法,包括:

确定待分类对象所属的目标分类任务;

将所述待分类对象输入对象分类模型,利用所述对象分类模型提取所述待分类对象的目标对象特征;

基于所述目标对象特征及所述目标分类任务对应第二存储信息中不同类别的对象特征,利用所述目标分类任务对应的输出模块,识别所述待分类对象的分类结果。

[0017] 第四方面,本申请实施例中提供了一种数据处理方法,包括:

在对象分类模型中构建新增分类任务的输出模块;

提取所述新增分类任务的第一训练样本的第一样本特征;

基于所述第一训练样本的第一样本特征,更新第一存储信息中相应类别的对象特征;

将所述第一训练样本的类别标签作为所述输出模块的输出数据,并结合所述第一训练样本的第一样本特征及所述第一存储信息中不同类别的对象特征,训练所述对象分类模型;

针对所述新增分类任务训练结束时,将所述第一存储信息中所述新增分类任务涉及类别的对象特征,构成所述新增分类任务对应的第二存储信息;其中,所述第二存储信息中不同类别的对象特征用于参与利用所述对象分类模型对所述新增分类任务的待分类对象的分类操作。

[0018] 第五方面,本申请实施例提供了一种文本分类方法,包括:

确定待分类文本所属的目标分类任务；

将所述待分类文本输入文本分类模型，利用所述文本分类模型提取所述待分类文本的文本特征；

基于所述文本特征及所述目标分类任务对应第二存储信息中不同类别的文本特征，利用所述目标分类任务对应的输出模块，识别所述待分类文本的分类结果。

[0019] 可选地，所述待分类文本为针对电商产品的评论数据；所述分类结果为所述待分类文本所属的情感类别；所述确定待分类文本所属的目标分类任务包括：根据确定所述评论数据所属的目标产品类目，确定所述目标产品类目对应的目标分类任务；

所述方法还包括：统计所述目标产品类目中属于同一情感类别的评论数据数量；基于不同情感类别的评论数据数量，生成提示信息；

或者，所述待分类文本为人机对话中的用户输入文本，所述分类结果为所述用户输入文本匹配的标准文本；所述方法还包括：基于所述标准文本查找对应的应答内容；输出所述应答内容；

或者，所述待分类文本为人机对话中的用户输入文本，所述分类结果为所述用户输入文本匹配的应答内容；所述方法还包括：输出所述应答内容。

[0020] 第六方面，本申请实施例提供了一种计算设备，包括处理组件以及存储组件；

所述存储组件存储一个或多个计算机指令；所述一个或多个计算机指令用以被所述处理组件调用执行，以实现如上述第一方面所述的数据处理方法，或者实现如上述第二方面所述的模型构建方法、或者实现如上述第三方面所述的分类方法或者实现如上述第四方面所述的数据处理方法。

[0021] 第七方面，本申请实施例提供了一种计算机存储介质，存储计算机程序，所述计算机程序被计算机执行时实现如上述第一方面所述的数据处理方法，或者实现如上述第二方面所述的模型构建方法、或者实现如上述第三方面所述的分类方法或者实现如上述第四方面所述的数据处理方法。

[0022] 第八方面，本申请实施例提供了一种计算机程序产品，包括计算机程序，所述计算机程序被计算机执行时实现如上述第一方面所述的数据处理方法，或者实现如上述第二方面所述的模型构建方法、或者实现如上述第三方面所述的分类方法或者实现如上述第四方面所述的数据处理方法。

[0023] 本申请实施例中训练获得的对象分类模型可支持扩展，存在新增分类任务时，只需在对象分类模型中搭建对应输出模块即可，对象分类模型对应设置有第一存储信息以及多个第二存储信息；第一存储信息保存不同类别的对象特征，并可以不断根据新增分类任务的训练样本的样本特征进行更新，第二存储信息为从第一存储信息中抽取出的其所对应的分类任务所涉及类别的对象特征，进行对象分类时，针对待分类对象所属的目标分类任务，利用对象分类模型中对应的输出模块以及所对应的第二存储信息即可以实现对象分类，本申请实施例的对象分类模型不仅支持扩展，可以同时处理多个分类任务，无需单独训练模型，降低了训练成本，并且新增分类任务可以利用历史分类任务的训练知识，保证分类准确度。

[0024] 本申请的这些方面或其他方面在以下实施例的描述中会更加简明易懂。

附图说明

[0025] 为了更清楚地说明本申请实施例或现有技术中的技术方案,下面将对实施例或现有技术描述中所需要使用的附图作一简单地介绍,显而易见地,下面描述中的附图是本申请的一些实施例,对于本领域普通技术人员来讲,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的附图。

[0026] 图1示出了本申请提供了一种数据处理方法一个实施例的流程图;
图2示出了本申请提供了一种数据处理方法又一个实施例的流程图;
图3示出了本申请提供了一种模型构建方法一个实施例的结构图;
图4示出了本申请实施例在一个实际应用中的对象分类模型的模型结构图;
图5示出了本申请提供了一种分类方法一个实施例的流程图;
图6示出了本申请实施例在一个实际应用中对象分类模型对应的分类处理示意图;

图7示出了本申请提供了一种文本分类方法一个实施例的流程图;
图8示出了本申请提供了一种数据处理装置一个实施例的结构示意图;
图9示出了本申请提供了一种数据处理装置又一个实施例的结构示意图;
图10示出了本申请提供了一种计算设备一个实施例的结构示意图;
图11示出了本申请提供了一种模型构建装置一个实施例的结构示意图;
图12示出了本申请提供了一种计算设备又一个实施例的结构示意图;
图13示出了本申请提供了一种分类装置一个实施例的结构示意图;
图14示出了本申请提供了一种计算设备又一个实施例的结构示意图。

具体实施方式

[0027] 为了使本技术领域的人员更好地理解本申请方案,下面将结合本申请实施例中的附图,对本申请实施例中的技术方案进行清楚、完整地描述。

[0028] 在本申请的说明书和权利要求书及上述附图中的描述的一些流程中,包含了按照特定顺序出现的多个操作,但是应该清楚了解,这些操作可以不按照其在本文中出现的顺序来执行或并行执行,操作的序号如101、102等,仅仅是用于区分各个不同的操作,序号本身不代表任何的执行顺序。另外,这些流程可以包括更多或更少的操作,并且这些操作可以按顺序执行或并行执行。需要说明的是,本文中的“第一”、“第二”等描述,是用于区分不同的消息、设备、模块等,不代表先后顺序,也不限定“第一”和“第二”是不同的类型。

[0029] 本申请的技术方案适用于对文本、图像、音频等数据对象的分类应用场景中。在下文一个或多个实施例中可能多以文本分类为例对本申请的技术方案进行介绍。

[0030] 以文本分类为例,实际应用中可能存在多个种类的文本分类需求,比如情感分类、意图识别、问答匹配等,而对于同一种类的文本分类需求也会存在多个分类任务,比如电子商务领域中,针对电商商品的商品评论数据进行好评以及差评等的情感分类,由于不同商品的类目不同,评论的主题和关注点均可能不一样,因此,对于不同类目的情感分类就会生成多个分类任务。传统的实现方式,针对每个分类任务,均会对应训练一个文本分类模型,随着任务数量的不断增加,需要训练的模型数量相应增加,这无疑会导致训练成本非常大。

[0031] 为了降低模型训练成本,同时保证模型训练准确度,发明人经过一系列研究提出

了本申请的技术方案,在本申请实施例中,对象分类模型支持扩展,当存在新增分类任务时,在对象分类模型中构建新增分类任务对应的输出模块,基于新增分类任务的第一训练样本,首先利用对象分类模型提取出第一样本特征,结合第一样本特征以及第一存储中不同类别的对象特征,重新训练对象分类模型,之后再利用训练之后的对象分类模型重新提取第一训练样本的样本特征,得到第二样本特征,并基于第二样本特征更新该第一存储信息,在新增分类任务训练结束时,将第一存储信息中新增分类任务涉及类别的对象特征,构成新增分类任务对应的第二存储信息;其中,第二存储信息中不同类别的对象特征用于参与利用对象分类模型对新增分类任务的待分类对象的分类操作。本申请实施例对象分类模型支持扩展,使得对象分类模型可以支持多个分类任务,无需单独训练模型,降低了训练成本以及训练复杂度,并且新增分类任务可以利用历史分类任务的训练知识,保证分类准确度。

[0032] 下面将结合本申请实施例中的附图,对本申请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申请一部分实施例,而不是全部的实施例。基于本申请中的实施例,本领域技术人员在没有作出创造性劳动前提下所获得的所有其他实施例,都属于本申请保护的范围。

[0033] 图1为本申请实施例提供的一种数据处理方法一个实施例的流程图,本实施例主要从模型训练角度对本申请技术方案进行介绍,该方法可以包括以下几个步骤:

101:在对象分类模型中构建新增分类任务的输出模块。

[0034] 对象分类模型可以用于对文本、图像、或音频等数据对象进行分类,例如可以是对文本进行分类的文本分类模型,对图像进行分类的图像分类模型或者对音频数据进行分类的音频分类模型。

[0035] 对象分类模型可以为机器学习模型,可以采用神经网络模型实现,其主要由输入模块、输出模块以及中间模块构成,也分别被称为输入层、输出层及中间层,分别可以由一个或多个神经网络层构成,本实施例中,对于新增分类任务,在对象分类模型中会构建仅属于新增分类任务的输出模块。

[0036] 该对象分类模型可以是针对原分类任务训练获得的模型,具体训练过程在下文实施例中会详细描述,也可以是指未经训练的模型,针对任一个分类任务,均可以作为新增分类任务按照本实施例的技术方案对模型进行训练。

[0037] 102:将新增分类任务当前训练批次的第一训练样本输入对象分类模型,利用对象分类模型提取第一训练样本的第一样本特征。

[0038] 实际应用中,由于模型的训练通常是分批次(batch)进行,每一训练批次使用新增分类任务对应训练数据集中的部分数据对模型进行一次训练和参数更新。

[0039] 本申请实施例中,为了方便描述,将新增分类任务对应训练数据集包含的训练样本命名为第一训练样本,每一训练批次可以使用一个或多个第一训练样本。该第一训练样本的数据类型根据数据对象的数据类型不同而不同,例如可以是文本、图像或音频等。

[0040] 其中,新增分类任务的训练数据集中包括新增分类任务所涉及类别分别对应的第一训练样本,每个第一训练样本对应设置有类别标签。

[0041] 将每一训练批次中的第一训练样本输入对象分类模型之后,首先利用对象分类模型提取第一样本特征。实际应用中,第一训练样本经由输入模块输入,中间模块通常包括至

少一个中间层,中间模块可以对经由输入模块输入的第一训练样本进行逐级计算处理,获得第一训练样本的深入表达,得到第一样本特征。第一样本特征可以具体采用向量形式表示。

[0042] 103:将第一训练样本的类别标签作为输出模块的输出数据,并结合第一训练样本的第一样本特征及第一存储信息中不同类别的对象特征,训练对象分类模型。

[0043] 第一存储信息中存储了不同类别的对象特征,比如对象分类模型用于对商品评论数据进行分类,涉及类别可以包括好评、中评以及差评,第一存储信息中可以存储好评、中评以及差评分别对应的对象特征。

[0044] 其中,对象分类模型预先针对原分类任务训练获得的情况下,第一存储信息中不同类别的对象特征可以根据基于原分类任务的训练样本得到,在下文会详细介绍;若对象分类模型为未经训练的模型,第一存储信息中的初始数据可以为空。

[0045] 可以将类别标签作为新增分类任务对应输出模块的输出数据,第一样本特征及不同类别的对象特征作为输出模块的输入数据,以此对对象分类模型进行训练。类别标签即为第一训练样本所属类别。

[0046] 可选地,可以是将第一训练样本的第一样本特征分别与第一存储信息中不同类别的对象特征进行融合,获得第一融合特征;再将至第一融合特征作为新增分类任务的输出模块的输入数据,以及将第一训练样本的类别标签作为新增分类任务的输出模块的输出数据,训练对象分类模型。对象分类模型的中间模块可以包括特征提取模块以及特征融合模块等,特征提取模块用于提取第一样本特征,而特征融合模块用于将第一样本特征与第一存储信息中各个类别的对象特征进行融合。

[0047] 其中,对象分类模型的特征提取模块可以是采用通用数据训练得到预训练模型,当然,也可以是利用训练样本及样本特征预先训练获得的模型。在文本分类场景中,该特征提取模块可以为文本编码(text encoder)模型,用以对文本数据编码获得文本特征,可选地,文本编码模型例如可以采用bert(Bidirectional Encoder Representations from Transformer,双向编码Transformer,一种预训练语言模型)实现。

[0048] 104:利用训练之后的对象分类模型提取第一训练样本的第二样本特征,并基于第一训练样本的第二样本特征更新第一存储信息。

[0049] 基于第一训练样本对对象分类模型训练之后,利用该对象分类模型提取重新提取第一训练样本的样本特征,为了描述上的区分,将重新提取的样本特征作为第二样本特征,可以基于该第二样本特征更新该第一存储信息,具体可以是,基于第二样本特征更新第一存储信息中,其对应类别标签指代类别所对应的对象特征。

[0050] 可选地,针对每一训练批次中的第一训练样本均可以按照步骤102~步骤104的操作执行直至所有训练批次训练结束。

[0051] 其中,作为另一种可选方式,执行步骤103之前,还可以利用102提取得到的第一样本特征,更新第一存储信息,具体更新第一存储信息中,第一样本特征所对应类别的对象特征,从而使得输入输出模块之前即对第一存储信息进行一次更新,使得第一存储信息会更加准确,进一步保证模型准确度。

[0052] 此外,可选地,基于第一训练样本的第二样本特征更新第一存储信息可以包括:基于第一训练样本的第二样本特征,计算属于同一类别的第一训练样本的平均第二样本特

征;将属于同一类别的平均第二样本特征与第一存储信息中的对象特征进行加权求和,并利用加权求和结果替换第一存储信息中对应类别的对象特征。

[0053] 也即首先计算当前训练批次中的属于同一类别的所有第一训练样本的第二样本特征的平均第二样本特征,可以得到多个类别对应的平均第二样本特征,属于同一类别的平均第二样本特征和第一存储信息中对象特征进行加权求和,再利用加权求和得到的特征替换掉第一存储信息中同一类别的对象特征。

[0054] 其中,可以基于属于同一类别的平均第二样本特征和第一存储信息分别对应的第一权重系数进行加权求和,其中,平均第二样本特征和第一存储信息分别对应的第一权重系数小于1,可以结合实际情况进行预先设定,并可以进行调整等。

[0055] 可选地,可以具体按照如下加权计算公式获得加权求和结果,意即第一存储信息中更新之后的对象特征:

$$G_j^{(m)} = (1 - \gamma)G_{j-1}^{(m)} + \frac{\gamma}{|\mathcal{D}_n^{(m)}|} \sum_{(x_{n,i}, y_{n,i}) \in \mathcal{D}_n^{(m)}} Q(x_{n,i});$$

其中,在该加权计算公式中, $Q(x_{n,i})$ 表示属于第m个类别的第一训练样本的第二样本特征; $G_{j-1}^{(m)}$ 表示第一存储信息中第m个类别对应原始的对象特征; $G_j^{(m)}$ 即为加权求和结果,意即第一存储信息中第m个类别更新之后的对象特征; $\mathcal{D}_n^{(m)}$ 表示当前训练批次中属于

第m个类别的第一训练样本的总数量; $\sum_{(x_{n,i}, y_{n,i}) \in \mathcal{D}_n^{(m)}} Q(x_{n,i})$ 表示第m个类别的第一训练样

本的第二样本特征之和; $\sum_{(x_{n,i}, y_{n,i}) \in \mathcal{D}_n^{(m)}} Q(x_{n,i}) / \mathcal{D}_n^{(m)}$ 即表示第m个类别对应的平均第二样

本特征; γ 表示第m个类别的第二样本特征的第一权重系数, $1 - \gamma$ 即为第一存储信息中第m个类别的原始对象特征所对应的第一权重系数。可选地,不同类别的平均第二样本特征所对应的第一权重系数可以相同,不同类别的对象特征所对应的第一权重系数也可以相同。

[0056] 105:针对新增分类任务的训练结束时,将第一存储信息中新增分类任务涉及类别的对象特征,构成新增分类任务对应的第二存储信息。

[0057] 其中,第二存储信息中不同类别的对象特征用于参与利用对象分类模型对新增分类任务的待分类对象的分类操作。

[0058] 新增分类任务训练结束意即利用训练数据集中的所有数据对对象分类模型完成训练,此时,可以将第一存储信息中新增分类任务所涉及类别的对象特征,单独保存为新增分类任务对应的第二存储信息。从而再利用对象分类模型对新增分类任务的待分类对象进行分类识别时,利用第二存储信息而非第一存储信息参与分类操作。具体可以是利用第二存储信息中各个类别的对象特征与提取的待分类对象的对象特征所进行融合。

[0059] 其中,第一存储信息和第二存储信息可以跟随对象分类模型而进行存储,以用于参与对象分类模型的训练和使用等。

[0060] 可选地,对象分类模型中还可以设置用于存储的第一记忆网络以及不同分类任务对应的第二记忆网络,第一存储信息存储至第一记忆网络中,不同分类任务的第二存储信息存储至各自对应的第二记忆网络。

[0061] 本实施例中,对象分类模型支持扩展,使得对象分类模型可以支持多个分类任务,

从而无需针对每个分类任务单独训练模型,降低了训练成本,且第一存储信息会参与新增分类任务对应的模型训练,同时也会根据新增分类任务的训练结果进行更新,用于参与下一个分类任务的训练,同时使得每加入一个分类任务,可以利用历史分类任务的训练知识,由于不同分类任务均是针对同一个分类种类,本申请充分利用了历史分类任务的训练知识,从而可以保证分类准确度。

[0062] 其中,对于每一个新增分类任务,均可以按照图1所示实施例的技术方案进行模型训练,若同时存在多个新增分类任务,则可以逐一针对每个新增分类任务执行图1所示实施例的技术方案,多个新增分类任务的训练顺序可以任意指定,当然也可以结合实际应用情况进行确定等。

[0063] 其中,将第一训练样本的第一样本特征与第一存储信息中不同类别的对象特征进行融合,可以由多种实现方式,比如可以将第一样本特征与不同类别的对象特征进行累加,又如首先将不同类别的对象特征进行加权求和之后再于第一样本特征进行累加等,通过将第一样本特征与不同类别的对象特征进行融合,使得训练过程可以融合历史训练知识,以此提高模型准确度。

[0064] 作为又一种可选方式,将第一训练样本的第一样本特征与第一存储信息中不同类别的对象特征进行融合,获得第一融合特征可以包括:

根据第一训练样本的第一样本特征,确定第一存储信息中不同类别的对象特征对应的第二权重系数;按照各自的第二权重系数,将第一存储信息中不同类别的对象特征进行加权求和,获得第一加权特征;将第一加权特征与第一训练样本的第一样本特征累加,获得第一融合特征。

[0065] 也就是说,第一存储信息中不同类别的对象特征各自对应的第二权重系数,可以根据第一样本特征来确定,可以根据与第一样本特征的关联程度来确定。比如,其中一种确定方式可以为:

根据第一训练样本的第一样本特征,计算第一存储信息中不同类别的对象特征分别与第一样本特征的内积;计算不同类别的对象特征对应的内积和;根据每个类别的对象特征的内积在内积和中的占比,获得每个类别的对象特征所对应的第二权重系数。

[0066] 为了便于理解,第一存储信息中不同类别的对象特征各自对应的第二权重系数可以按照如下权重系数计算公式获得:

$$\alpha^{(m)}(x_{n,i}) = \frac{Q(x_{n,i})^T \cdot G_n^{(m)}}{\sum_{y^{(\tilde{m})} \in \mathcal{Y}_n} \alpha^{(\tilde{m})}(x_{n,i})};$$

其中,在该权重系数计算公式中,第一存储信息中包括 y_n 个类别的对象特征; $\alpha^{(m)}(x_{n,i})$ 表示第 m 个类别的对象特征对应的第二权重系数; $Q(x_{n,i})$ 表示任一个类别的第一训练样本的第一样本特征; $G_n^{(m)}$ 表示第 m 个类别的对象特征; $Q(x_{n,i})^T \cdot G_n^{(m)}$ 表示第 m 个类别的对象特征与第一样本特征的内积,其中, m 为正整数; $\sum_{y^{(\tilde{m})} \in \mathcal{Y}_n} \alpha^{(\tilde{m})}(x_{n,i})$ 表示不同类别的对象特征对应的内积和。

[0067] 基于上式计算获得的第二权重系数,第一融合特征即可以按照如下融合计算公式

获得：

$$Att(x_{n,i}) = Q(x_{n,i}) + \sum_{y^{(m)} \in \mathcal{Y}_n} \alpha^{(m)}(x_{n,i}) \cdot G_n^{(m)} ;$$

其中,在该融合计算公式中,第一存储信息中包括 y_n 个类别的对象特征; $Att(x_n, i)$ 表示第一融合特征; $\sum_{y^{(m)} \in \mathcal{Y}_n} \alpha^{(m)}(x_{n,i}) \cdot G_n^{(m)}$ 即表示第一加权特征; $\alpha^{(m)}(x_n, i)$ 表示第 m 个类别的对象特征对应的第二权重系数; $Q(x_n, i)$ 表示任一个类别的第一训练样本的第一样本特征; $G_n^{(m)}$ 表示第 m 个类别的对象特征。

[0068] 此外,由前文描述可知,对象分类模型可以预先针对原始分类任务训练获得,因此,在对象分类模型可以包括至少一个原分类任务分别对应的输出模块。

[0069] 作为一种可选方式,对象分类模型针对至少一个原分类任务可以按照如下方式预先训练获得：

利用对象分类模型提取至少一个原分类任务对应的第二训练样本的第一样本特征；

计算属于同一个类别的第二训练样本的平均第一样本特征作为该类别的对象特征,并将计算获得的不同类别的对象特征构成第一存储信息；

将第二训练样本对应的类别标签作为其对应输出模块的输出数据,并结合第二训练样本的第一样本特征以及第一存储信息,训练对象分类模型；

针对至少一个原分类任务训练结束时,将第一存储信息中至少一个原分类任务各自对应类别的对象特征,构成至少一个原分类任务各自的第二存储信息。

[0070] 其中,每个原分类任务对应的第二存储信息用于参与利用对象分类模型对属于该原分类任务的待分类对象的分类操作。

[0071] 其中,对象分类模型中构建有每一个原分类任务对应的输出模块,针对至少一个原分类任务训练结束时,将第一存储信息中每个原分类任务所涉及类别的对象特征单独保存为每个原分类任务的第二存储信息。

[0072] 作为另一种可选方式,对象分类模型针对至少一个原分类任务按照如下方式预先训练获得：

利用对象分类模型提取至少一个原分类任务对应的第二训练样本的第一样本特征；

计算属于同一个类别的第二训练样本的平均第一样本特征作为该类别的对象特征,并将计算获得的不同类别的对象特征构成第一存储信息；

将第二训练样本对应的类别标签作为其对应输出模块的输出数据,并结合第二训练样本的第一样本特征以及第一存储信息,训练对象分类模型；

针对至少一个原分类任务训练结束时,抽取第一存储信息中不同类别的对象特征,构成至少一个原分类任务对应的第二存储信息。

[0073] 与上一种可选方式不同之处在于,该至少一个原分类任务可以对应一个第二存储信息,将针对至少一个原分类任务训练结束时的第一存储信息,保存作为该至少一个原分类任务对应的第二存储信息,而该至少一个原分类任务对应的第二存储信息用于参与利用

对象分类模型对该至少一个原分类任务的待分类对象的分类操作。

[0074] 针对上述对象分类模型针对原分类任务进行训练的两种可选方式,在某些实施例中,计算属于同一个类别的第二训练样本的平均第一样本特征作为该类别的对象特征可以包括:

计算每个原分类任务对应的属于同一个类别的第二训练样本的第一样本特征的平均特征;

将至少一个原分类任务分别对应的同一个类别的平均特征进行平均计算,获得同一类别的第二训练样本的平均第一样本特征,并作为该类别对应的对象特征。

[0075] 具体的,可以按照如下特征计算公式,计算获得第一存储信息中不同类别所对应的对象特征:

$$G_N^{(m)} = \frac{1}{|\mathcal{T}^{(m)}|} \sum_{\mathcal{T}_n \in \mathcal{T}^{(m)}} \frac{1}{|\mathcal{D}_n^{(m)}|} \sum_{(x_{n,i}, y_{n,i}) \in \mathcal{D}_n^{(m)}} Q(x_{n,i});$$

其中, $G_N^{(m)}$ 表示第 m 个类别对应的第二训练样本的平均第一样本特征,意即第 m 个类别的对象特征; $Q(x_{n,i})$ 表示属于第 \mathcal{T}_n 个原分类任务中的第二训练样本的第一样本特征; $|\mathcal{D}_n^{(m)}|$ 表示属于第 \mathcal{T}_n 个原分类任务中第 m 个类别的第二训练样本的数量;

$\frac{1}{|\mathcal{D}_n^{(m)}|} \sum_{(x_{n,i}, y_{n,i}) \in \mathcal{D}_n^{(m)}} Q(x_{n,i})$ 表示第 \mathcal{T}_n 个原分类任务中属于同一个类别的第二训练样本

的第一样本特征的平均特征;其中, \mathcal{T}_n 为正整数, $\mathcal{T}^{(m)}$ 表示原分类任务的总个数。

[0076] 另外,由于模型训练分批次进行,因此,在某些实施例中,将第二训练样本对应的类别标签作为其对应输出模块的输出数据,并结合第二训练样本的第一样本特征以及第一存储信息,训练对象分类模型可以包括:

将当前训练批次的第二训练样本的第一样本特征与第一存储信息中不同类别的对象特征进行融合,获得第二融合特征;

将第二融合特征输入对应的输出模块,以及将第二训练样本的类别标签作为对应输出模块的输出数据,训练对象分类模型;

利用训练之后的对象分类模型提取第二训练样本的第二样本特征,并基于第二训练样本的第二样本特征,更新第一存储信息中相应类别的对象特征。

[0077] 其中,第一存储信息中初始数据可以为空。基于第二训练样本的第二样本特征,更新第一存储信息中相应类别的对象特征的具体更新方式,与新增分类任务中基于第一训练样本的第一样本特征更新第一存储信息中相应类别的对象特征的方式类似,比如可以是:

基于第二训练样本的第二样本特征,计算属于同一类别的第二训练样本的平均第二样本特征;将属于同一类别的平均第二样本特征与第一存储信息中的对象特征进行加权求和,并利用加权求和结果替换第一存储信息中对应类别的对象特征。

[0078] 另外,在某些实施例,将第二训练样本的第一样本特征分别与第一存储信息中不同类别的对象特征进行融合,获得第二融合特征可以包括:

根据第二训练样本的第一样本特征,确定第一存储信息中不同类别的对象特征对应的第三权重系数;

按照各自的第三权重系数,将第一存储信息中不同类别的对象特征的进行加权求和,获得第二加权特征;

将第二加权特征与第二训练样本的第一样本特征累加,获得第二融合特征。

[0079] 其中,根据第二训练样本的第一样本特征,确定第一存储信息中不同类别的对象特征对应的第三权重系数,与前文相应描述中根据第一训练样本的第一样本特征,确定第一存储信息中不同类别的对象特征对应的第二权重系数的确定方式相似,比如可以是:

根据第二训练样本的第一样本特征,计算第一存储信息中不同类别的对象特征分别与第二训练样本的第一样本特征的内积;计算不同类别的对象特征对应的内积和;根据每个类别的对象特征的内积在内积和中的占比,获得每个类别的对象特征所对应的第三权重系数。

[0080] 图2为本申请实施例提供的一种数据处理方法又一个实施例的流程图,该方法可以包括以下几个步骤:

201:在对象分类模型中构建新增分类任务的输出模块。

[0081] 202:提取新增分类任务的第一训练样本的第一样本特征。

[0082] 203:基于第一训练样本的第一样本特征,更新第一存储信息中相应类别的对象特征。

[0083] 可选地,可以利用特征提取模型提取新增分类任务对应每一个第一训练样本的第一样本特征。本实施例与图1所示实施例不同之处在于,可以由单独的特征提取模型提取第一训练样本的第一样本特征。

[0084] 可选地,可以利用新增分类任务所对应的所有第一训练样本的第一样本特征,更新第一存储信息中相应类别的对象特征。具体的,可以是首先计算属于同一类别的第一训练样本的平均第一样本特征,然后再将属于同一类别的平均第一样本特征与第一存储信息中的对象特征进行加权求和,利用加权求和结果替换第一存储信息中对应类别的对象特征。

[0085] 204:将第一训练样本的类别标签作为输出模块的输出数据,并结合第一训练样本的第一样本特征及第一存储信息中不同类别的对象特征,训练对象分类模型。

[0086] 模型训练过程可以分批次进行,针对每一训练批次的每个第一训练样本均可以按照步骤204的操作执行,以实现对象分类模型的训练。

[0087] 可选地,可以首先将第一训练样本的第一样本特征分别与第一存储信息中不同类别的对象特征进行融合,获得第一融合特征;将第一融合特征作为新增分类任务的输出模块的输入数据,以及将第一训练样本的类别标签作为新增分类任务的输出模块的输出数据,训练对象分类模型。

[0088] 可选地,可以根据第一训练样本的第一样本特征,确定第一存储信息中不同类别的对象特征对应的第二权重系数;按照各自的第二权重系数,将第一存储信息中不同类别的对象特征进行加权求和,获得第一加权特征;将第一加权特征与第一训练样本的第一样本特征累加,获得第一融合特征。

[0089] 第二权重系数的计算方式可以详见前文所述,在此不再赘述。

[0090] 205:针对新增分类任务训练结束时,将第一存储信息中新增分类任务涉及类别的对象特征,构成新增分类任务对应的第二存储信息。

[0091] 其中,第二存储信息中不同类别的对象特征用于参与利用对象分类模型对新增分类任务的待分类对象的分类操作。

[0092] 本实施例与图1所示实施例不同之处在于,针对新增分类任务训练结束之后,可以直接将第一存储信息中新增分类任务涉及类别的对象特征,构成新增分类任务对应的第二存储信息。

[0093] 可选地,对象分类模型包括至少一个原分类任务分别对应的输出模块;对象分类模型针对至少一个原分类任务可以按照如下方式预先训练获得:

利用特征提取模型提取至少一个原分类任务对应的第二训练样本的第一样本特征;

计算属于同一个类别的第二训练样本的平均第一样本特征,作为该类别的对象特征,并将计算获得的不同类别的对象特征构成第一存储信息;

将第二训练样本对应的类别标签作为其对应输出模块的输出数据,并结合第二训练样本的第一样本特征以及第一存储信息,训练对象分类模型;

针对至少一个原分类任务训练结束时,将第一存储信息中至少一个原分类任务各自对应类别的对象特征,构成至少一个原分类任务各自的第二存储信息;或者抽取第一存储信息中不同类别的对象特征,构成至少一个原分类任务对应的第二存储信息。

[0094] 在某些实施例中,计算属于同一个类别的第二训练样本的平均第一样本特征作为该类别的对象特征可以包括:

计算每个原分类任务对应的属于同一个类别的第二训练样本的第一样本特征的平均特征;

将至少一个原分类任务分别对应的同一个类别的平均特征进行平均计算,获得同一类别的第二训练样本的平均第一样本特征,并作为该类别对应的对象特征。

[0095] 另外,由于模型训练分批次进行,因此,在某些实施例中,将第二训练样本分别对应的类别标签作为各自对应输出模块的输出数据,并结合第二训练样本的第一样本特征以及第一存储信息,训练对象分类模型可以包括:

将当前训练批次的第二训练样本输入对象分类模型,利用对象分类模型提取第二训练样本的第一样本特征;

将第二训练样本的第一样本特征分别与第一存储信息中不同类别的对象特征进行融合,获得第二融合特征;

将第二融合特征输入对应的输出模块,以及将第二训练样本的类别标签作为对应输出模块的输出数据,训练对象分类模型;

利用训练之后的对象分类模型提取第二训练样本的第二样本特征,并基于第二训练样本的第二样本特征,更新第一存储信息中相应类别的对象特征。

[0096] 其中,第一存储信息中初始数据可以为空。基于第二训练样本的第二样本特征,更新第一存储信息中相应类别的对象特征的具体更新方式,与新增分类任务中基于第一训练样本的第一样本特征更新第一存储信息中相应类别的对象特征的方式类似,比如可以是:

基于第二训练样本的第二样本特征,计算属于同一类别的第二训练样本的平均第二样本特征;将属于同一类别的平均第二样本特征与第一存储信息中的对象特征进行加权求和,并利用加权求和结果替换第一存储信息中对应类别的对象特征。

[0097] 另外,在某些实施例,将第二训练样本的第一样本特征分别与第一存储信息中不同类别的对象特征进行融合,获得第二融合特征可以包括:

根据第二训练样本的第一样本特征,确定第一存储信息中不同类别的对象特征对应的第三权重系数;

按照各自的第三权重系数,将第一存储信息中不同类别的对象特征的进行加权求和,获得第二加权特征;

将第二加权特征与第二训练样本的第一样本特征累加,获得第二融合特征。

[0098] 其中,根据第二训练样本的第一样本特征,确定第一存储信息中不同类别的对象特征对应的第三权重系数,与前文相应描述中根据第一训练样本的第一样本特征,确定第一存储信息中不同类别的对象特征对应的第二权重系数的确定方式相似,比如可以是:根据第二训练样本的第一样本特征,计算第一存储信息中不同类别的对象特征分别与第二训练样本的第一样本特征的内积;计算不同类别的对象特征对应的内积和;根据每个类别的对象特征的内积在内积和中的占比,获得每个类别的对象特征所对应的第三权重系数。

[0099] 本申请实施例中可以利用单独的特征提取模型提取样本特征,该特征提取模型例如可以采用bert模型实现,对于对象分类模型的训练可以按照图2所示实施例的技术方案执行,当然,特征提取模型也可以集成在对象分类模型中,作为特征提取模块用于进行特征提取,具体即可以按照图1所示实施例的技术方案训练对象分类模型。

[0100] 图3为本申请实施例提供的一种模型构建方法一个实施例的流程图,本实施例主要从模型构建角度对本申请技术方案进行介绍,该方法可以包括以下几个步骤:

301:依次搭建输入模块、特征提取模块、特征融合模块及至少一个原分类任务对应的输出模块,获得对象分类模型。

[0101] 其中,每个原分类任务可以分别对应一个输出模块,当然也可以是多个原分类任务对应一个输出模块,意即可以是为每个原分类任务分别搭建一个输出模块,当然也可以是多个原分类任务对应搭建一个输出模块。

[0102] 为了进一步保证模型准确度,可以是为至少一个原分类任务分别搭建各自对应的输出模块。

[0103] 302:对应设置对象分类模型的第一存储信息及至少一个原分类任务各自的第二存储信息。

[0104] 其中,第一存储信息存储不同类别的对象特征,基于至少一个原分类任务的第二训练样本提取的样本特征获得,具体获得方式可以详见前文相应实施例中,此处不再赘述;

其中,任一原分类任务的第二存储信息存储从第一存储中抽取的原分类任务涉及类别的对象特征。

[0105] 303:根据分类任务扩展需求,在对象分类模型中搭建新增分类任务的输出模块。

[0106] 304:对应设置新增分类任务的第二存储信息。

[0107] 其中,第一存储信息根据新增分类任务的第一训练样本所提取的样本特征进行更新,具体更新方式可以详见前文相应实施例中所述,此处不再赘述。

[0108] 新增分类任务的第二存储信息存储从第一存储信息中抽取的新增分类任务涉及类别的对象特征。

[0109] 对于构建获得的对象分类模型具体训练方式可以详见前文图1或图2所示实施例

中所述,此处不再赘述。

[0110] 可选地,在对象分类模型中还可以搭建第一记忆网络以及不同分类任务对应的第二记忆网络,第一记忆网络用以存储第一存储信息,第二记忆网络用于存储器对应分类任务的第二存储信息。

[0111] 为了便于理解,下面以用于进行文本分类的文本分类模型为例,下面结合图4所示的模型架构图对文本分类模型进行简要介绍,如图5中所示,文本分类模型主要由输入模块401、特征提取模块402、特征融合模块403、输出模块构成,此外还可以包括第一记忆网络404以及针对每个分类任务(Task)的第二记忆网络;第一记忆网络404中保存第一存储信息;第二记忆网络中保存相应分类任务的第二存储信息。其中,特征提取模块可以采用Bert模型实现,当然也可以采用其它语言模型实现,例如ALBert(A lite BERT for self-supervised learning of language representations,一种精简的语言特征自监督学习的Bert模型),RoBERTa(A Robustly Optimized BERT Pretraining Approach,一种鲁棒优化的预训练Bert模型)等。

[0112] 文本分类模型的训练过程可以包括两个阶段:初始学习阶段及终身学习阶段,初始学习阶段针对原分类任务进行训练,终身学习阶段针对新增分类任务进行训练。

[0113] 假设包括两个原分类任务Task1以及Task2,文本分类模型处于初始学习阶段时,文本分类模型中构建有原分类任务Task1对应的输出模块Out1以及原分类任务Task2对应的输出模块Out2,当然,初始学习阶段多个原分类任务也可以对应设置一个输出模块,之后即可以利用两个原分类任务的训练文本对文本分类模型进行训练,具体训练方式详见前文所述,并在训练结束之后,在文本分类模型中对应设置原分类任务Task1对应的原分类任务Task2对应的第二记忆网络LM1以及第二记忆网络LM2;

存在分类任务扩展需求时,由于不同分类任务针对的分类种类相同,因此可以直接对文本分类模型进行再训练,文本分类模型即处于终身学习阶段,逐一针对每个新增分类任务进行再训练,假设新增分类任务为Task3,首先在文本分类模型构建新增分类任务为Task3对应的输出模块Out3,之后即可以基于新增分类任务的训练文本对文本分类模型进行训练,具体训练方式详见前文所述,并在训练结束之后,对应设置新增分类任务为Task3对应的第二记忆网络LM3,文本分类模型中的其它模块架构不变。后续若继续存在新增分类任务,文本分类模型即可以按照新增分类任务为Task3的方式进行扩建以及训练等。

[0114] 第一记忆网络中的第一存储信息会不断根据新增分类任务而进行更新,第二记忆网络中的第二存储信息一旦生成即被冻结,以参与对待分类文本的分类操作。

[0115] 针对前文相应实施例中搭建并训练获得的对象分类模型即可以用于对多个分类任务的分类操作,如图5所示,为本申请实施例提供的一种分类方法一个实施例的流图,该方法可以包括以下几个步骤:

501:确定待分类对象所属的目标分类任务。

[0116] 实际应用中,该待分类对象可以为文本、音频数据或者图像等数据对象。

[0117] 502:将待分类对象输入对象分类模型,利用对象分类模型提取待分类对象的目标对象特征。

[0118] 其中,对象分类模型的具体训练方式可以详见前文相应实施例中所述,此处不再赘述。

[0119] 503:基于目标对象特征及目标分类任务对应第二存储信息中不同类别的对象特征,利用目标分类任务对应的输出模块,识别待分类对象的分类结果。

[0120] 在某些实施例中,基于目标对象特征及目标分类任务对应第二存储信息中不同类别的对象特征,利用目标分类任务对应的输出模块,识别待分类对象的分类结果可以包括:

将目标对象特征与目标分类任务对应第二存储信息中不同类别的对象特征进行融合,获得目标融合特征;

将目标融合特征输入目标分类任务对应的输出模块,以识别获得待分类对象的分类结果。

[0121] 可选地,将目标对象特征与目标分类任务对应第二存储信息中不同类别的对象特征进行融合,获得目标融合特征可以包括:

根据目标对象特征,确定第一存储信息中不同类别的对象特征对应的第二权重系数;按照各自对应的第四权重系数,将第一存储信息中不同类别的对象特征进行加权求和,获得目标加权特征;将目标加权特征与目标对象特征进行累加,获得目标融合特征。

[0122] 其中,根据目标对象特征,确定第一存储信息中不同类别的对象特征对应的第四权重系数可以包括:

根据目标对象特征,计算第一存储信息中不同类别的对象特征分别与第一样本特征的内积;计算不同类别的对象特征对应的内积和;根据每个类别的对象特征的内积在内积和中的占比,获得每个类别的对象特征所对应的第四权重系数。

[0123] 为了便于理解,仍以文本分类模型为例,参见图6所示,以图4所示的模块架构图为例,假设文本分类模型当前支持x个分类任务,x为正整数,下面对分类过程进行简要介绍。首先,待分类文本经由输入模块401输入文本分类模型,之后由特征提取模块402提取目标文本特征,假设待分类文本属于分类任务Task1,则从第二记忆网络LM1获得第二存储信息,经由特征融合模块403将目标文本特征与第二存储信息中各个类别的文本特征进行融合,获得目标融合特征,之后目标融合特征输入分类任务Task1对应的输出模块Out1,经由输出模块Out1获得最终的分类结果。

[0124] 在实际应用中,由前文描述可知,对象分类模型可具体为文本分类模型,用于进行文本分类,其中文本分类的种类通常可以包括情感分类、意图识别、问答匹配等。如图7所述,本申请实施例还提供了一种文本分类方法,可以包括以下几个步骤:

701:确定待分类文本所属的目标分类任务。

[0125] 702:将待分类文本输入文本分类模型,利用文本分类模型提取待分类文本的目标文本特征。

[0126] 703:基于目标文本特征及目标分类任务对应第二存储信息中不同类别的文本特征,利用目标分类任务对应的输出模块,识别待分类文本的分类结果。

[0127] 图8与图5所示实施例不同之处在于,待分类对象具体为待分类文本,其它相同或相似步骤可以详见图5所示实施例中所述,在此不再赘述。文本分类模型的处理过程可以详见图6中所示。

[0128] 而文本分类模型的训练过程可以详见前文相应实施例中的对象分类模型的训练过程,区别仅在于,训练样本具体为文本形式,此处将不再赘述。

[0129] 在一个具体应用中,文本分类模型可以用于进行情感分类;待分类文本可以为针

对电商产品的评论数据;分类结果为待分类文本所属的情感类别,比如可以包括好评、中评以及插差评等;针对不同产品类目需要单独进行情感类别的识别,因此,不同产品类目对应不同分类任务。则确定待分类文本所属的目标分类任务可以包括:

根据确定评论数据所属的目标产品类目,确定目标产品类目对应的目标分类任务;

此外,该方法还可以包括:

统计目标产品类目中属于同一情感类别的评论数据数量;

基于不同情感类别的评论数据数量,生成提示信息。

[0130] 该提示信息可以发送至相关人员,以便于相关人员进行相应的业务处理,比如根据差评数量作出相应的商品调整策略等等。

[0131] 在另一个具体应用中,文本分类模型可以具体用于进行意图识别,待分类文本可以为人机对话中的用户输入文本,分类结果为用户输入文本匹配的标准文本;该标准文本即代表用户输入文本的用户意图,其中,用户输入文本可以通过用户输入语音识别获得,由于受限于用户口语化以及教育程度不同,需要对用户输入文本进行意图识别,得到标准文本,以了解用户真正意图。获得标准文本之后,该方法还可以包括:

基于标准文本查找对应的应答内容;

输出应答内容。

[0132] 可选地,该应答内容可以发送至客户端,由客户端向用户展示该应答内容。

[0133] 在人机对话场景中,可以预先保存了不同标准文本对应的应答内容,因此通过将用户输入文本利用文本分类模型转化为标准文本,即可以准确找到与之匹配的应答内容,并输出给用户。

[0134] 在又一个具体应用中,文本分类模型可以具体用于进行问答匹配,待分类文本为人机对话中的用户输入文本,分类结果为用户输入文本匹配的应答内容。

[0135] 该方法还可以包括:

输出应答内容。

[0136] 可选地,该应答内容可以发送至客户端,由客户端向用户展示该应答内容。

[0137] 需要说明的是,上述仅是举例说明对象分类模型的几种实现场景,本申请并不局限于此。

[0138] 图8为本申请提供的一种数据处理装置一个实施例的结构示意图,该装置可以包括:

第一构建单元801,用于在对象分类模型中构建新增分类任务的输出模块;

第一提取单元802,用于将新增分类任务当前训练批次的第一训练样本输入对象分类模型,利用对象分类模型提取第一训练样本的第一样本特征;

第一训练单元803,用于将第一训练样本的类别标签作为输出模块的输出数据,并结合第一训练样本的第一样本特征及第一存储信息中不同类别的对象特征,训练对象分类模型;

第一更新单元804,用于利用训练之后的对象分类模型提取第一训练样本的第二样本特征,并基于第一训练样本的第二样本特征更新第一存储信息;

第一存储单元805,用于针对新增分类任务的训练结束时,将第一存储信息中新增

分类任务涉及类别的对象特征,构成新增分类任务对应的第二存储信息;其中,第二存储信息中不同类别的对象特征用于参与利用对象分类模型对新增分类任务的待分类对象的分类操作。

[0139] 图8所述的数据处理装置可以执行图1所示实施例所述的数据处理方法,其实现原理和技术效果不再赘述。对于上述实施例中的数据处理装置其中各个模块、单元执行操作的具体方式已经在有关该方法的实施例中进行了详细描述,此处将不做详细阐述说明。

[0140] 图9为本申请实施例提供的一种数据处理装置又一个实施例的结构示意图,该装置可以包括:

第二构建单元901,用于在对象分类模型中构建新增分类任务的输出模块;

第二提取单元902,用于提取新增分类任务的第一训练样本的第一样本特征;

第二更新单元903,用于基于第一训练样本的第一样本特征,更新第一存储信息中相应类别的对象特征;

第二训练单元904,用于将第一训练样本的类别标签作为输出模块的输出数据,并结合第一训练样本的第一样本特征及第一存储信息中不同类别的对象特征,训练对象分类模型;

第二存储单元905,用于针对新增分类任务训练结束时,将第一存储信息中新增分类任务涉及类别的对象特征,构成新增分类任务对应的第二存储信息;其中,第二存储信息中不同类别的对象特征用于参与利用对象分类模型对新增分类任务的待分类对象的分类操作。

[0141] 图9所述的数据处理装置可以执行图2所示实施例所述的数据处理方法,其实现原理和技术效果不再赘述。对于上述实施例中的数据处理装置其中各个模块、单元执行操作的具体方式已经在有关该方法的实施例中进行了详细描述,此处将不做详细阐述说明。

[0142] 在一个可能的设计中,图8或图9所示实施例的数据处理装置可以实现为计算设备,如图10所示,该计算设备可以包括存储组件1001以及处理组件1002;

存储组件1001存储一条或多条计算机指令,其中,一条或多条计算机指令供处理组件1002调用执行,以实现如图1或图2所示的数据处理方法。

[0143] 当然,计算设备必然还可以包括其他部件,例如输入/输出接口、通信组件等。输入/输出接口为处理组件和外围接口模块之间提供接口,上述外围接口模块可以是输出设备、输入设备等。通信组件被配置为便于计算设备和其他设备之间有线或无线方式的通信等。

[0144] 需要说明的是,该计算设备可以为物理设备或者云计算平台提供的弹性计算主机等,此时计算设备即可以是指云服务器,上述处理组件、存储组件等可以从云计算平台租用或购买的基础服务器资源。

[0145] 此外,上述计算设备也可以实现成多个服务器或终端设备组成的分布式集群,也可以实现成单个服务器或单个终端设备。

[0146] 此外,本申请实施例还提供了一种计算机可读存储介质,存储有计算机程序,所述计算机程序被计算机执行时可以实现上述图1所示实施例的数据处理方法。

[0147] 此外,本申请实施例还提供了一种计算机可读存储介质,存储有计算机程序,所述计算机程序被计算机执行时可以实现上述图2所示实施例的数据处理方法。

[0148] 此外,本申请实施例提供了一种计算机程序产品,包括计算机程序,所述计算机程序被计算机执行时可以实现上述图1所示实施例的数据处理方法。

[0149] 此外,本申请实施例提供了一种计算机程序产品,包括计算机程序,所述计算机程序被计算机执行时可以实现上述图2所示实施例的数据处理方法。

[0150] 图11为本申请实施例提供的一种模型构建装置一个实施例的结构示意图,该装置可以包括:

第三构建单元1101,用于依次搭建输入模块、特征提取模块、特征融合模块及对应至少一个原分类任务的输出模块,获得对象分类模型;

第一设置单元1102,对应设置对象分类模型的第一存储信息及至少一个原分类任务各自的第二存储信息;其中,第一存储信息存储不同类别的对象特征,基于至少一个原分类任务的第二训练样本提取的样本特征获得;任一原分类任务的第二存储信息存储从第一存储中抽取的原分类任务涉及类别的对象特征;

第四构建单元1103,用于根据分类任务扩展需求,在对象分类模型中搭建新增分类任务的输出模块;

第二设置单元1104,用于对应设置新增分类任务的第二存储信息;其中,第一存储信息根据新增分类任务的第一训练样本所提取的样本特征进行更新,新增分类任务的第二存储信息存储从第一存储信息中抽取的新增分类任务涉及类别的对象特征。

[0151] 图11所述的模型构建装置可以执行图3所示实施例所述的模型构建方法,其实现原理和技术效果不再赘述。对于上述实施例中的数据处理装置其中各个模块、单元执行操作的具体方式已经在有关该方法的实施例中进行了详细描述,此处将不做详细阐述说明。

[0152] 在一个可能的设计中,图11所示实施例的模型构建装置可以实现为计算设备,如图12所示,该计算设备可以包括存储组件1201以及处理组件1202;

存储组件1201存储一条或多条计算机指令,其中,所述一条或多条计算机指令供处理组件1202调用执行,以实现如图3所示的模型构建方法。

[0153] 当然,计算设备必然还可以包括其他部件,例如输入/输出接口、通信组件等。输入/输出接口为处理组件和外围接口模块之间提供接口,上述外围接口模块可以是输出设备、输入设备等。通信组件被配置为便于计算设备和其他设备之间有线或无线方式的通信等。

[0154] 需要说明的是,该计算设备可以为物理设备或者云计算平台提供的弹性计算主机等,此时计算设备即可以是指云服务器,上述处理组件、存储组件等可以从云计算平台租用或购买的基础服务器资源。

[0155] 此外,上述计算设备也可以实现成多个服务器或终端设备组成的分布式集群,也可以实现成单个服务器或单个终端设备。

[0156] 此外,本申请实施例还提供了一种计算机可读存储介质,存储有计算机程序,所述计算机程序被计算机执行时可以实现上述图3所示实施例的模型构建方法。

[0157] 此外,本申请实施例提供了一种计算机程序产品,包括计算机程序,所述计算机程序被计算机执行时可以实现上述图3所示实施例的模型构建方法。

[0158] 图13为本申请实施例提供的一种分类装置一个实施例的结构示意图,该装置可以包括:

任务确定单元1301,用于确定待分类对象所属的目标分类任务;

特征提取单元1302,用于将待分类对象输入对象分类模型,利用对象分类模型提取待分类对象的目标对象特征;

分类单元1303,用于基于目标对象特征及目标分类任务对应第二存储信息中不同类别的对象特征,利用目标分类任务对应的输出模块,识别待分类对象的分类结果。

[0159] 在一个实际应用中,对象分类模型可以具体为文本分类模型,用于进行文本分类,因此,任务确定单元,可以具体用于确定待分类文本所属的目标分类任务;

特征提取单元,可以具体用于将待分类文本输入文本分类模型,利用文本分类模型提取待分类文本的文本特征;

分类单元,可以具体用于基于文本特征及目标分类任务对应第二存储信息中不同类别的文本特征,利用目标分类任务对应的输出模块,识别待分类文本的分类结果。

[0160] 图13所述的分类装置可以执行图5所示实施例所述的分类方法,其实现原理和技术效果不再赘述。对于上述实施例中的数据处理装置其中各个模块、单元执行操作的具体方式已经在有关该方法的实施例中进行了详细描述,此处将不做详细阐述说明。

[0161] 在一个可能的设计中,图13所示实施例的分类装置可以实现为计算设备,如图14所示,该计算设备可以包括存储组件1401以及处理组件1402;

存储组件1401存储一条或多条计算机指令,其中,所述一条或多条计算机指令供处理组件1402调用执行,以实现如图5所示的分类方法。

[0162] 当然,计算设备必然还可以包括其他部件,例如输入/输出接口、通信组件等。输入/输出接口为处理组件和外围接口模块之间提供接口,上述外围接口模块可以是输出设备、输入设备等。通信组件被配置为便于计算设备和其他设备之间有线或无线方式的通信等。

[0163] 需要说明的是,该计算设备可以为物理设备或者云计算平台提供的弹性计算主机等,此时计算设备即可以是指云服务器,上述处理组件、存储组件等可以从云计算平台租用或购买的基础服务器资源。

[0164] 此外,上述计算设备也可以实现成多个服务器或终端设备组成的分布式集群,也可以实现成单个服务器或单个终端设备。

[0165] 此外,本申请实施例还提供了一种计算机可读存储介质,存储有计算机程序,所述计算机程序被计算机执行时可以实现上述图5所示实施例的分类方法。

[0166] 此外,本申请实施例提供了一种计算机程序产品,包括计算机程序,所述计算机程序被计算机执行时可以实现上述图5所示实施例的分类方法。

[0167] 其中,上述相应实施例中所涉及的处理组件可以包括一个或多个处理器来执行计算机指令,以完成上述的方法中的全部或部分步骤。当然处理组件也可以为一个或多个应用专用集成电路(ASIC)、数字信号处理器(DSP)、数字信号处理设备(DSPD)、可编程逻辑器件(PLD)、现场可编程门阵列(FPGA)、控制器、微控制器、微处理器或其他电子元件实现,用于执行上述方法。

[0168] 存储组件被配置为存储各种类型的数据以支持在终端的操作。存储组件可以由任何类型的易失性或非易失性存储设备或者它们的组合实现,如静态随机存取存储器(SRAM),电可擦除可编程只读存储器(EEPROM),可擦除可编程只读存储器(EPROM),可编程

只读存储器 (PROM), 只读存储器 (ROM), 磁存储器, 快闪存储器, 磁盘或光盘。

[0169] 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, 上述描述的系统, 装置和单元的具体工作过程, 可以参考前述方法实施例中的对应过程, 在此不再赘述。

[0170] 以上所描述的装置实施例仅仅是示意性的, 其中所述作为分离部件说明的单元可以是或者也可以不是物理上分开的, 作为单元显示的部件可以是或者也可以不是物理单元, 即可以位于一个地方, 或者也可以分布到多个网络单元上。可以根据实际的需要选择其中的部分或者全部模块来实现本实施例方案的目的。本领域普通技术人员在不付出创造性的劳动的情况下, 即可以理解并实施。

[0171] 通过以上的实施方式的描述, 本领域的技术人员可以清楚地了解到各实施方式可借助软件加必需的通用硬件平台的方式来实现, 当然也可以通过硬件。基于这样的理解, 上述技术方案本质上或者说对现有技术做出贡献的部分可以以软件产品的形式体现出来, 该计算机软件产品可以存储在计算机可读存储介质中, 如ROM/RAM、磁碟、光盘等, 包括若干指令用以使得一台计算机设备 (可以是个人计算机, 服务器, 或者网络设备等) 执行各个实施例或者实施例的某些部分所述的方法。

[0172] 最后应说明的是: 以上实施例仅用以说明本申请的技术方案, 而非对其限制; 尽管参照前述实施例对本申请进行了详细的说明, 本领域的普通技术人员应当理解: 其依然可以对前述各实施例所记载的技术方案进行修改, 或者对其中部分技术特征进行等同替换; 而这些修改或者替换, 并不使相应技术方案的本质脱离本申请各实施例技术方案的精神和范围。

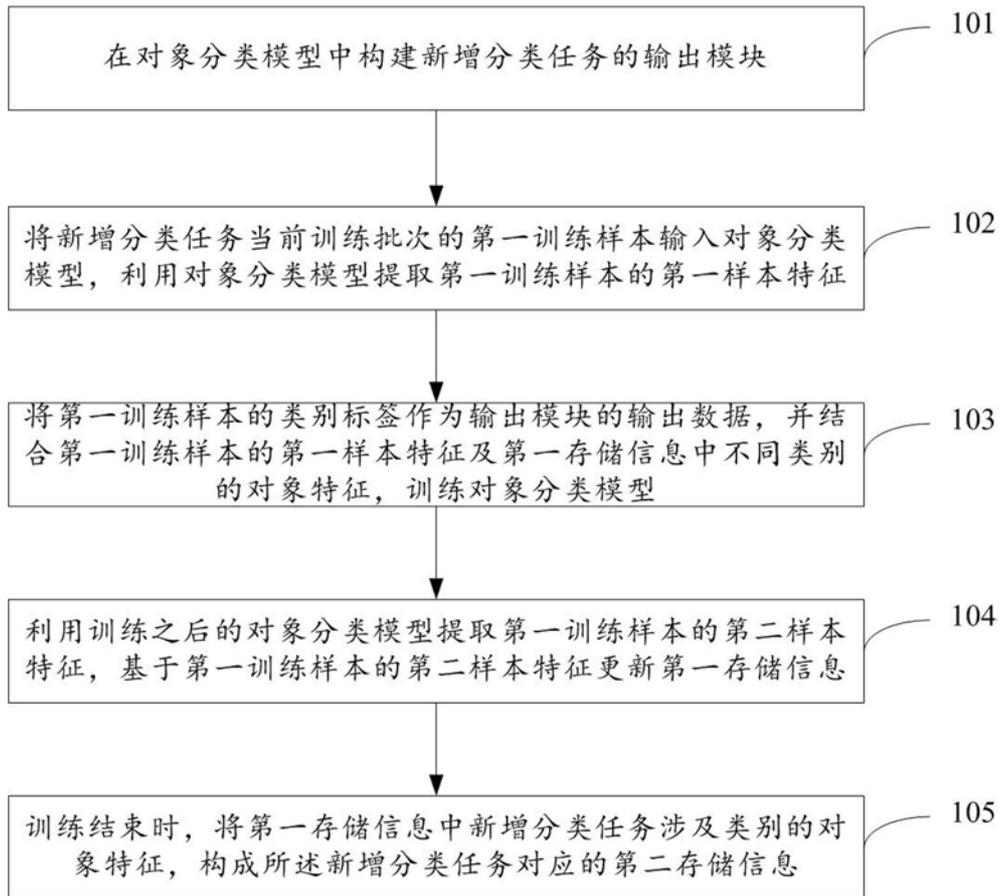


图1

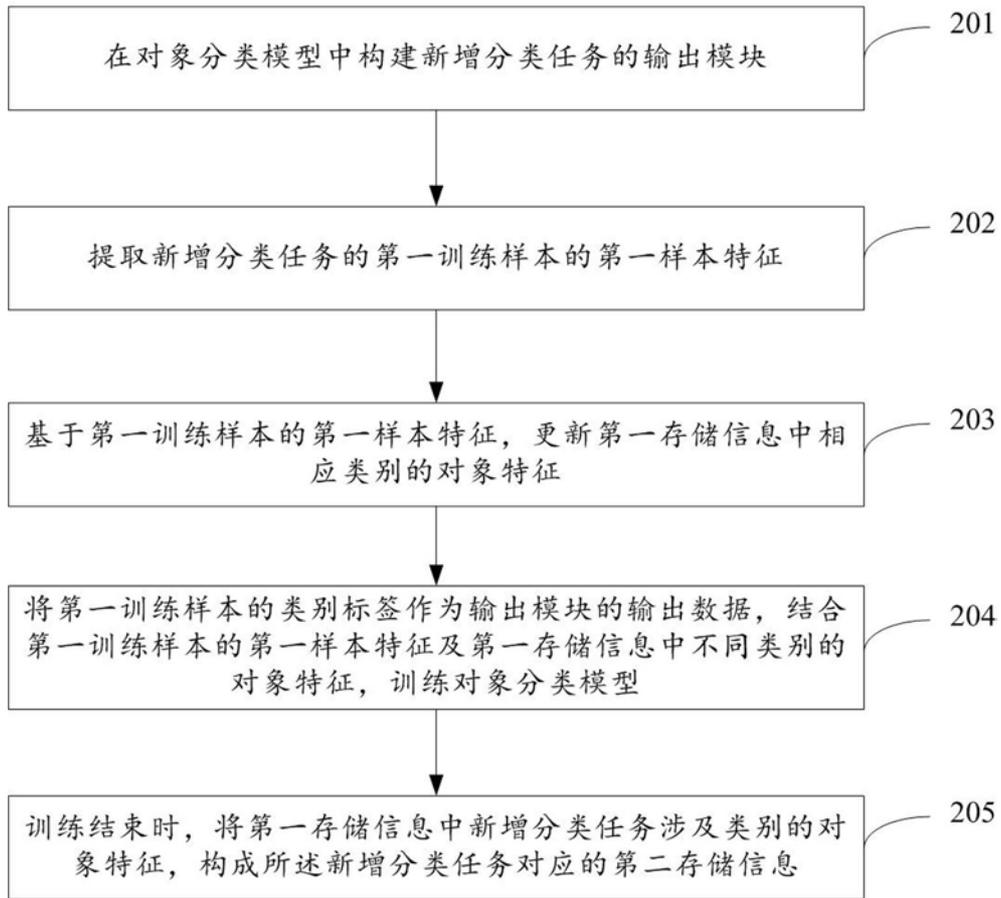


图2

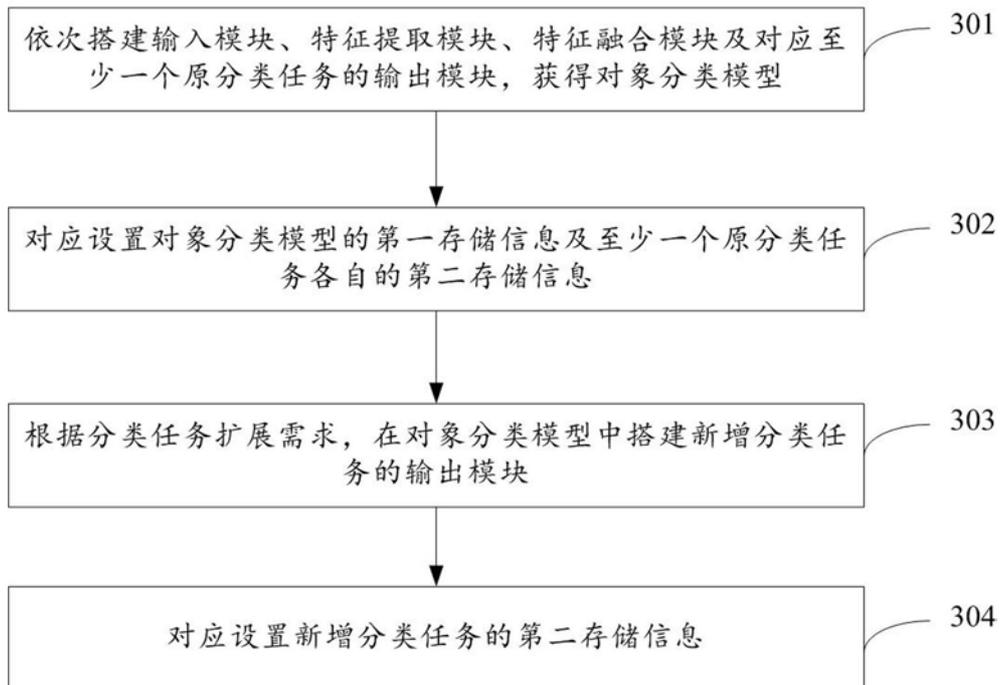


图3

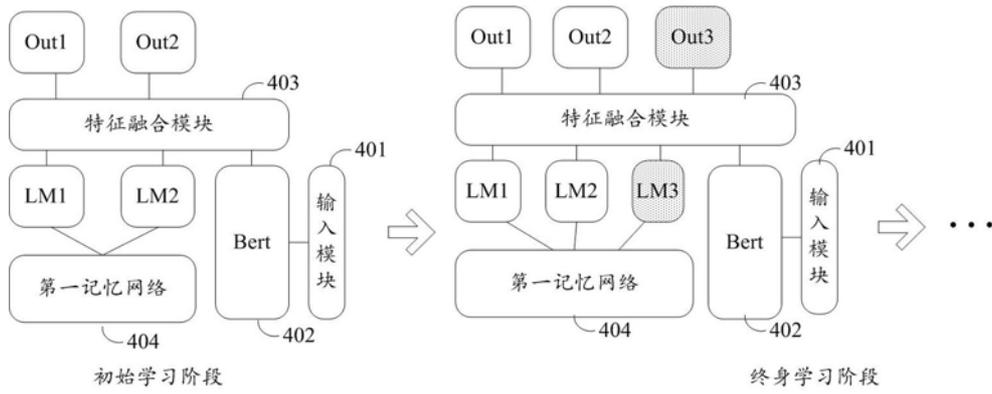


图4

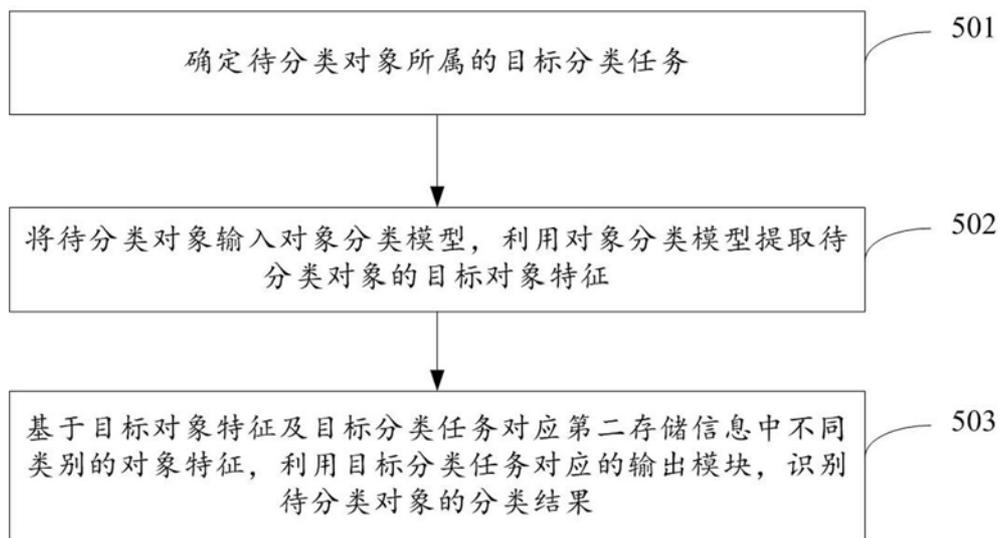


图5

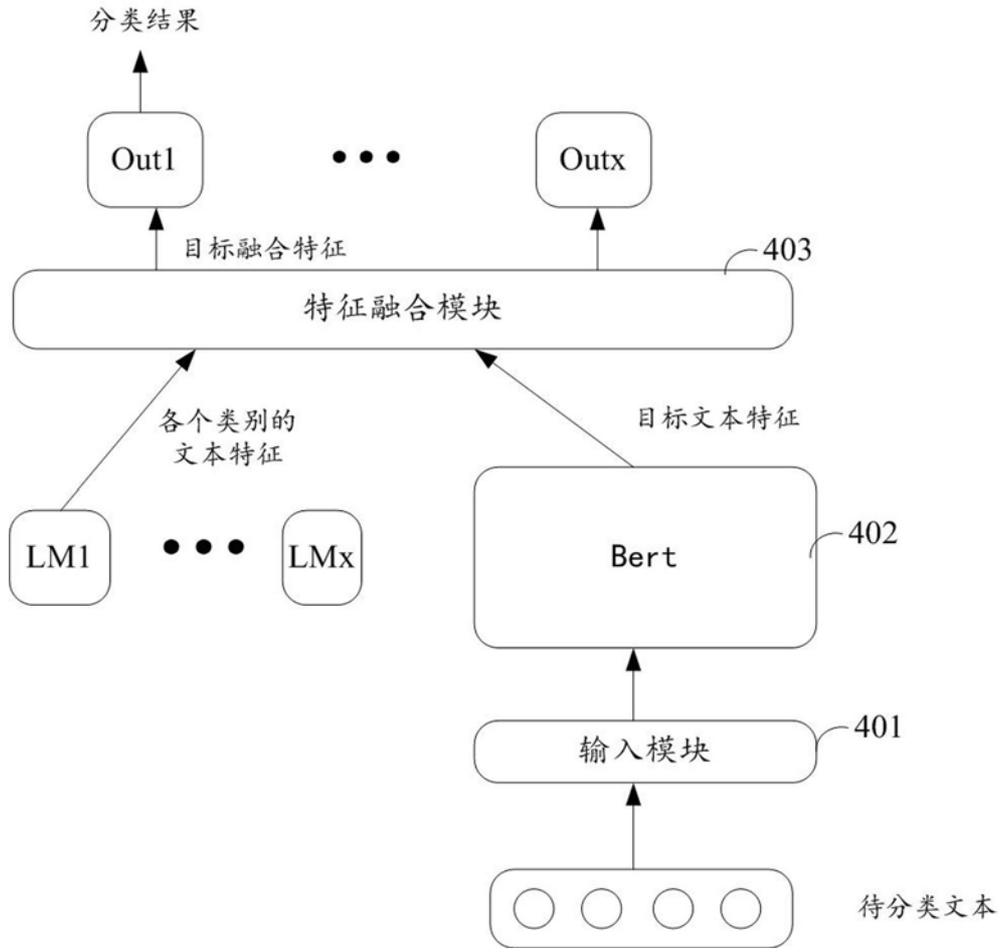


图6

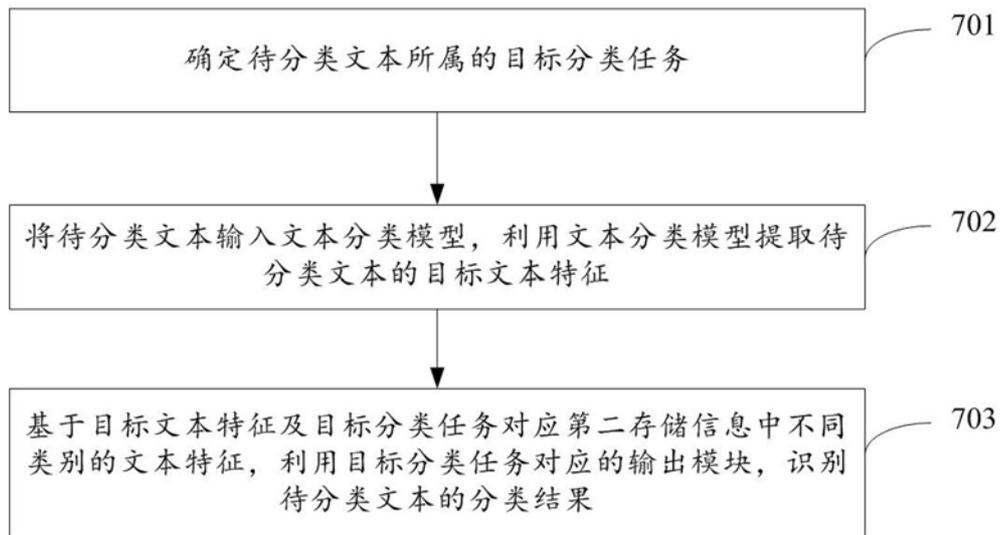


图7

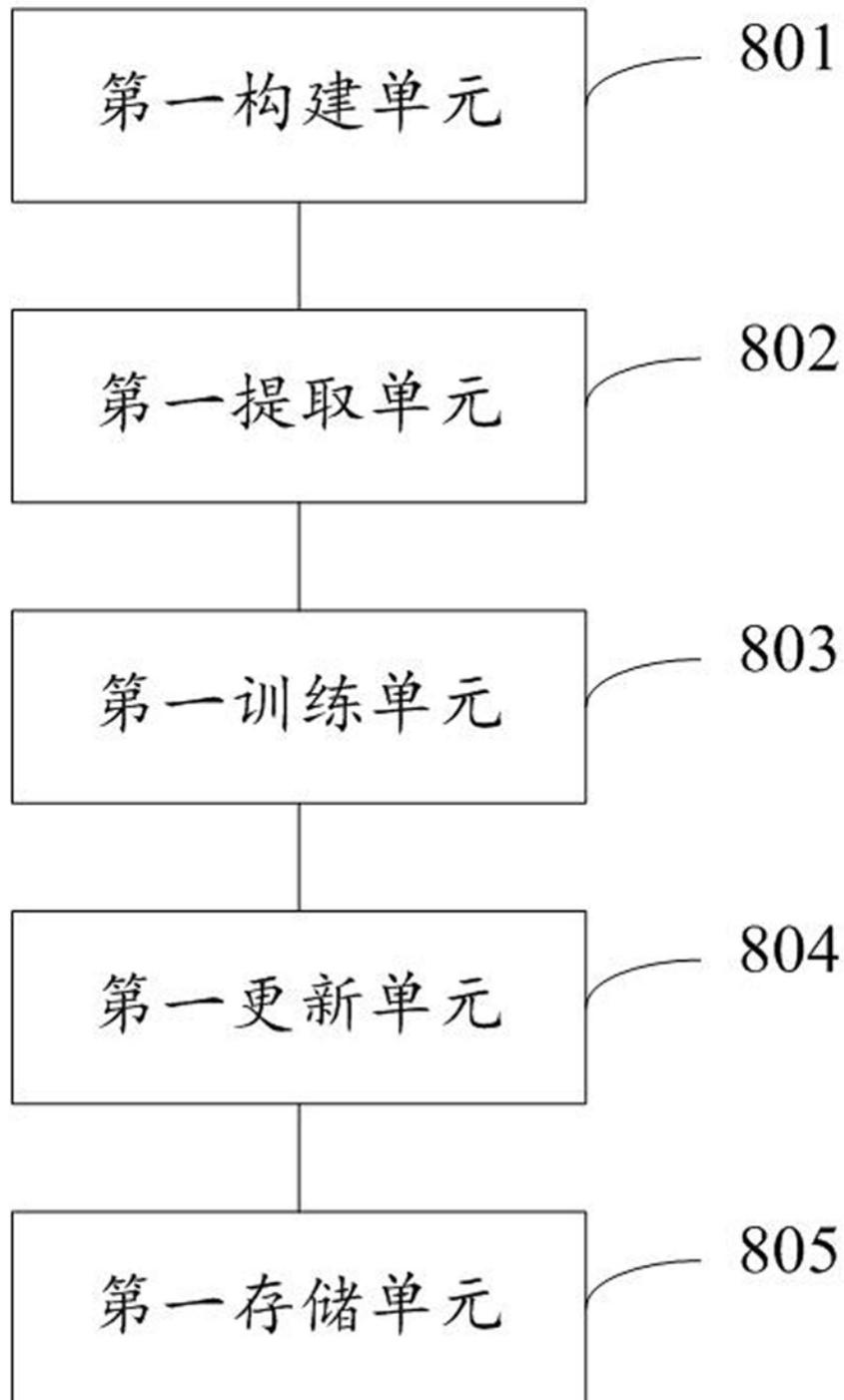


图8

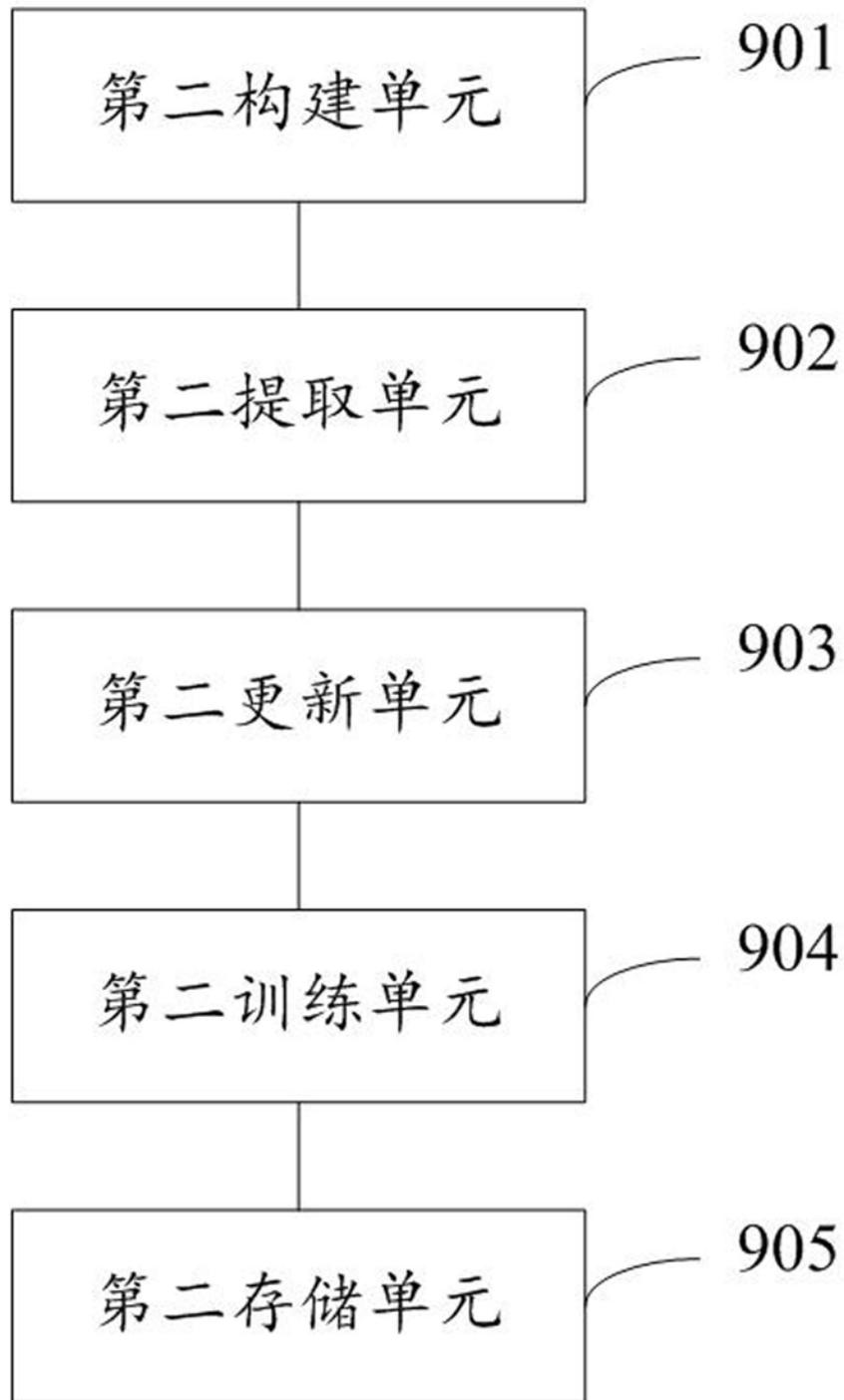


图9

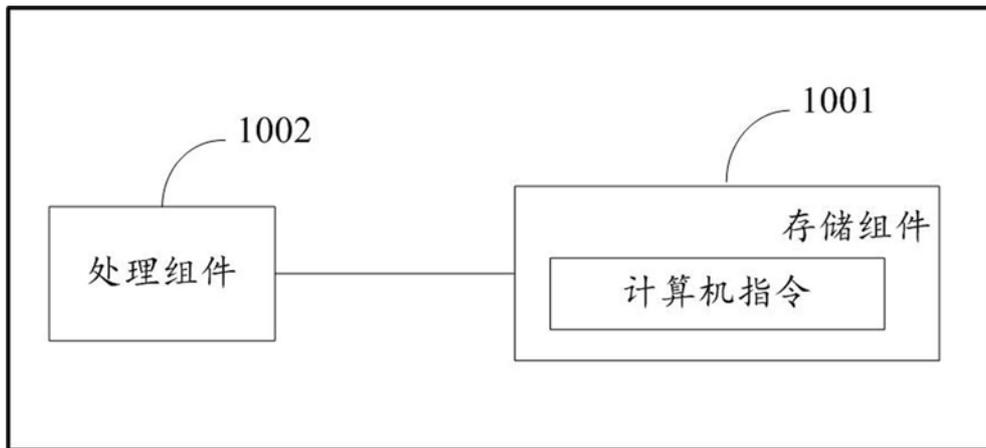


图10

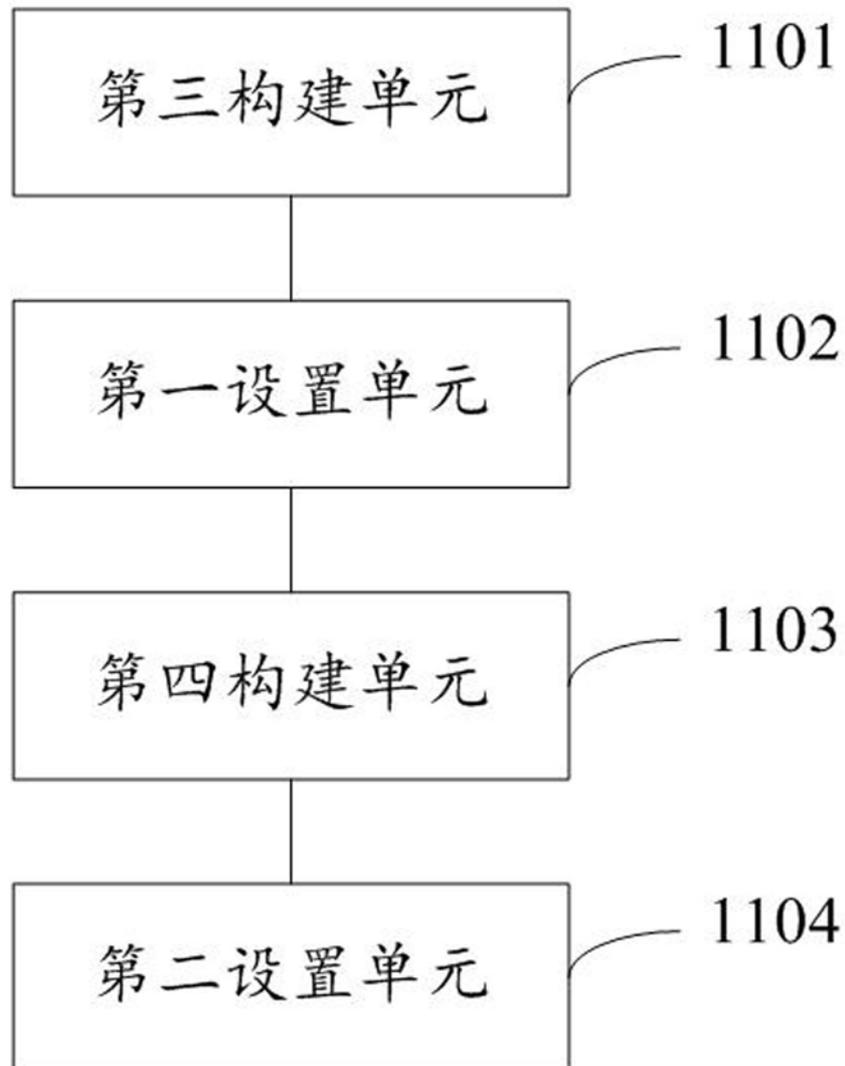


图11

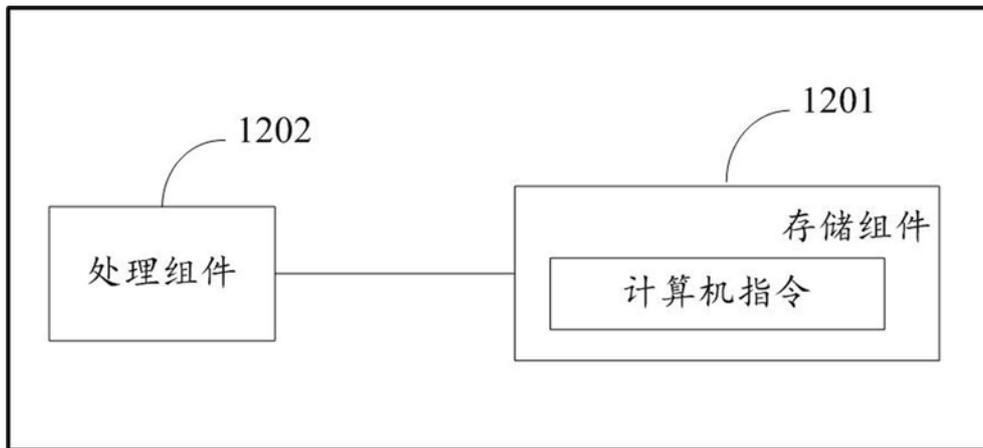


图12

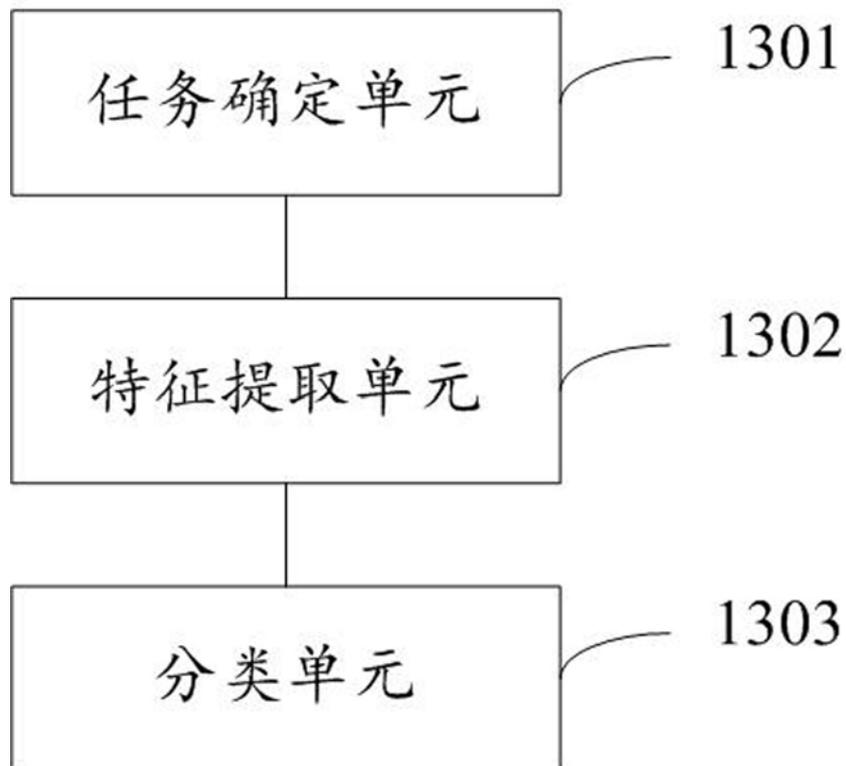


图13

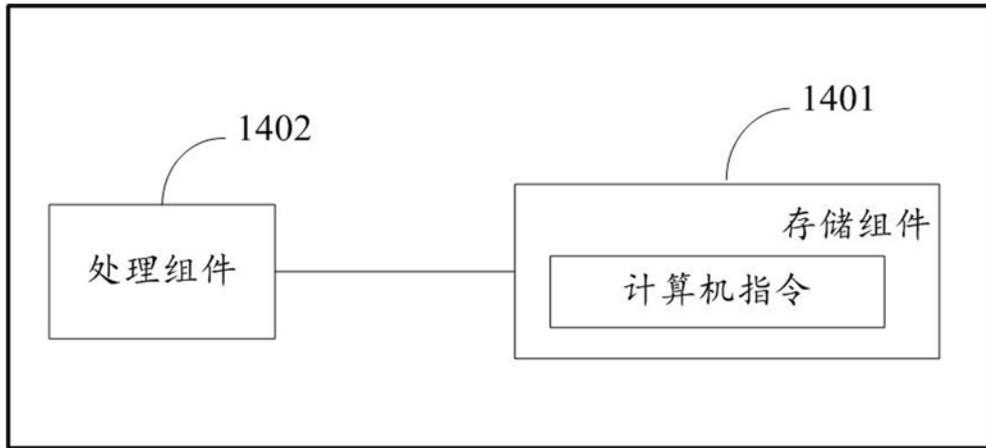


图14