



(12)发明专利申请

(10)申请公布号 CN 108921207 A

(43)申请公布日 2018. 11. 30

(21)申请号 201810637638.4

(22)申请日 2018.06.20

(71)申请人 中诚信征信有限公司

地址 100011 北京市东城区东四南大街礼士胡同54号

(72)发明人 何博睿 李映坤

(74)专利代理机构 北京柏杉松知识产权代理事务所(普通合伙) 11413

代理人 赵元 马敬

(51)Int.Cl.

G06K 9/62(2006.01)

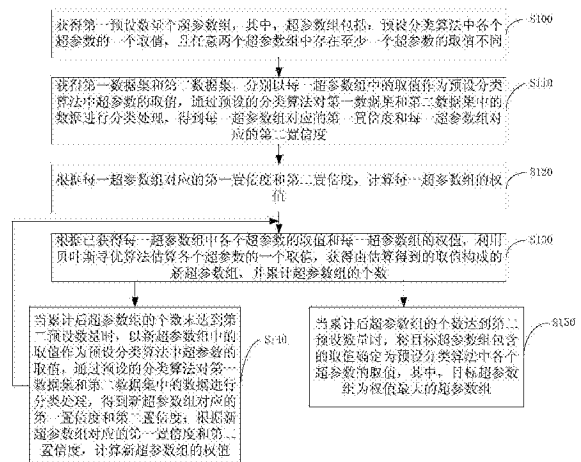
权利要求书3页 说明书13页 附图2页

(54)发明名称

一种超参数确定方法、装置及设备

(57)摘要

本发明实施例提供了一种超参数确定方法、装置及设备,该方法包括:获得第一预设数量个超参数组;获得第一数据集和第二数据集,分别以每一超参数组中的取值作为预设分类算法中超参数的取值,对第一数据集和第二数据集中的数据进行分类处理,得到每一超参数组对应的第一置信度和第二置信度;根据每一超参数组对应的置信度,计算每一超参数组的权值;根据已获得每一超参数组中的取值和每一超参数组的权值,利用贝叶斯寻优算法估算各个超参数的一个取值,获得新超参数组并累计超参数组个数;当累计后超参数组个数达到第二预设数量时将权值最大的超参数组中的取值作为预设分类算法中超参数的取值。应用本发明实施例提供的方法能够提高超参数确定效率。



1. 一种超参数确定方法,其特征在于,所述方法包括:

获得第一预设数量个超参数组,其中,所述超参数组包括:预设分类算法中各个超参数的一个取值,且任意两个超参数组中存在至少一个超参数的取值不同;

获得第一数据集和第二数据集,分别以每一超参数组中的取值作为所述预设分类算法中超参数的取值,通过所述预设的分类算法对所述第一数据集和第二数据集中的数据进行分类处理,得到每一超参数组对应的第一置信度和每一超参数组对应的第二置信度,其中,第一置信度为:第一数据集分类结果的置信度,第二置信度为:第二数据集分类结果的置信度;

根据每一超参数组对应的第一置信度和第二置信度,计算每一超参数组的权值;

根据已获得每一超参数组中各个超参数的取值和每一超参数组的权值,利用贝叶斯寻优算法估算各个超参数的一个取值,获得由估算得到的取值构成的新超参数组,并累计超参数组的个数;

当累计后超参数组的个数未达到第二预设数量时,以新超参数组中的取值作为所述预设分类算法中超参数的取值,通过所述预设的分类算法对所述第一数据集和第二数据集中的数据进行分类处理,得到新超参数组对应的第一置信度和第二置信度;根据新超参数组对应的第一置信度和第二置信度,计算新超参数组的权值;返回所述根据已获得每一超参数组中各个超参数的取值和每一超参数组的权值,利用贝叶斯寻优算法估算各个超参数的一个取值,获得由估算得到的取值构成的新超参数组,并累计超参数组的个数的步骤;

当累计后超参数组的个数达到第二预设数量时,将目标超参数组包含的取值确定为所述预设分类算法中各个超参数的取值,其中,所述目标超参数组为权值最大的超参数组。

2. 如权利要求1所述的方法,其特征在于,所述获得第一预设数量个超参数组的步骤,包括:

获取所述预设分类算法中各个超参数的取值范围;

以任意两个超参数组中存在至少一个超参数的取值不同为选取原则,分别在各个超参数的取值范围内选取一个值,并获得由所选取的各个值构成的一个超参数组;

返回所述以任意两个超参数组中存在至少一个超参数的取值不同为选取原则,分别在各个超参数的取值范围内选取一个值,并获得由所选取的各个值构成的一个超参数组的步骤,直至获得所述第一预设数量个超参数组。

3. 如权利要求1所述的方法,其特征在于,所述获得第一数据集和第二数据集,分别以每一超参数组中的取值作为所述预设分类算法中超参数的取值,通过所述预设的分类算法对所述第一数据集和第二数据集中的数据进行分类处理,得到每一超参数组对应的第一置信度和每一超参数组对应的第二置信度的步骤,包括:

利用分层抽样法将包含待处理数据的待处理数据集划分为 k 个子集,其中, $k > 1$,且任意两个子集之间不包含相同的待处理数据;

从 k 个子集中选取 $k-1$ 个子集合并作为第一数据集、剩余子集作为第二数据集的方式,确定 k 个第一数据集和 k 个第二数据集;

针对每一超参数组,以该超参数组中的取值作为所述预设分类算法中超参数的取值,通过所述预设的分类算法分别对 k 个第一数据集和 k 个第二数据集中的数据进行分类处理,得到该超参数组对应的 k 个第一数据集分类结果的置信度和 k 个第二数据集分类结果的置

信度；

将每一超参数组对应的k个第一数据集分类结果的置信度的平均值确定为该超参数组对应的第一置信度，将每一超参数组对应的k个第二数据集分类结果的置信度的平均值确定为该超参数组对应的第二置信度。

4. 如权利要求1所述的方法，其特征在于，所述将目标超参数组包含的取值确定为所述预设分类算法中各个超参数的取值的步骤，包括：

计算目标超参数组对应的第一置信度和第二置信度的差值；

当所述差值小于预设第一阈值时，将所述目标超参数组包含的取值确定为所述预设分类算法中各个超参数的取值；

当所述差值不小于预设第一阈值时，按照预设第一步长增大所述第一预设数量，按照预设第二步长增大所述第二预设数量，并返回所述获得第一预设数量个超参数组的步骤。

5. 如权利要求1-4任一所述的方法，其特征在于，所述根据每一超参数组对应的第一置信度和第二置信度，计算每一超参数组的权值的步骤，包括：

通过以下表达式计算每一超参数组的权值：

$$y(AUC) = \begin{cases} \frac{AUC_{test}}{abs(AUC_{train} - AUC_{test})} & abs(AUC_{train} - AUC_{test}) > M \\ \frac{AUC_{test}}{M} & abs(AUC_{train} - AUC_{test}) \leq M \end{cases}$$

其中，y(AUC)表示超参数组的权值，abs()表示取绝对值函数， AUC_{test} 表示超参数组对应的第二置信度， AUC_{train} 表示超参数组对应的第一置信度，M为预设数值。

6. 一种超参数确定装置，其特征在于，所述装置包括：

获取模块，用于获得第一预设数量个超参数组，其中，所述超参数组包括：预设分类算法中各个超参数的一个取值，且任意两个超参数组中存在至少一个超参数的取值不同；

处理模块，用于获得第一数据集和第二数据集，分别以每一超参数组中的取值作为所述预设分类算法中超参数的取值，通过所述预设的分类算法对所述第一数据集和第二数据集中的数据进行分类处理，得到每一超参数组对应的第一置信度和每一超参数组对应的第二置信度，其中，第一置信度为：第一数据集分类结果的置信度，第二置信度为：第二数据集分类结果的置信度；

计算模块，用于根据每一超参数组对应的第一置信度和第二置信度，计算每一超参数组的权值；

累计模块，用于根据已获得每一超参数组中各个超参数的取值和每一超参数组的权值，利用贝叶斯寻优算法估算各个超参数的一个取值，获得由估算得到的取值构成的新超参数组，并累计超参数组的个数；

返回模块，用于当累计后超参数组的个数未达到第二预设数量时，以新超参数组中的取值作为所述预设分类算法中超参数的取值，通过所述预设的分类算法对所述第一数据集和第二数据集中的数据进行分类处理，得到新超参数组对应的第一置信度和第二置信度；根据新超参数组对应的第一置信度和第二置信度，计算新超参数组的权值，并触发所述累计模块；

确定模块，用于当累计后超参数组的个数达到第二预设数量时，将目标超参数组包含

的取值确定为所述预设分类算法中各个超参数的取值,其中,所述目标超参数组为权值最大的超参数组。

7. 如权利要求6所述的装置,其特征在于,所述获得模块包括:

取值范围获取子模块,用于获取所述预设分类算法中各个超参数的取值范围;

超参数组获得子模块,用于以任意两个超参数组中存在至少一个超参数的取值不同为选取原则,分别在各个超参数的取值范围内选取一个值,并获得由所选取的各个值构成的一个超参数组;返回所述以任意两个超参数组中存在至少一个超参数的取值不同为选取原则,分别在各个超参数的取值范围内选取一个值,并获得由所选取的各个值构成的一个超参数组的步骤,直至获得所述第一预设数量个超参数组。

8. 如权利要求6所述的装置,其特征在于,所述处理模块具体用于,

利用分层抽样法将包含待处理数据的待处理数据集划分为k个子集,其中, $k > 1$,且任意两个子集之间不包含相同的待处理数据;

从k个子集中选取k-1个子集合并作为第一数据集、剩余子集作为第二数据集的方式,确定k个第一数据集和k个第二数据集;

针对每一超参数组,以该超参数组中的取值作为所述预设分类算法中超参数的取值,通过所述预设的分类算法分别对k个第一数据集和k个第二数据集中的数据进行分类处理,得到该超参数组对应的k个第一数据集分类结果的置信度和k个第二数据集分类结果的置信度;

将每一超参数组对应的k个第一数据集分类结果的置信度的平均值确定为该超参数组对应的第一置信度,将每一超参数组对应的k个第二数据集分类结果的置信度的平均值确定为该超参数组对应的第二置信度。

9. 如权利要求6所述的方法,其特征在于,所述确定模块具体用于,

计算目标超参数组对应的第一置信度和第二置信度的差值;

当所述差值小于预设第一阈值时,将所述目标超参数组包含的取值确定为所述预设分类算法中各个超参数的取值;

当所述差值不小于预设第一阈值时,按照预设第一步长增大所述第一预设数量,按照预设第二步长增大所述第二预设数量,并触发所述获取模块。

10. 一种电子设备,其特征在于,包括处理器、通信接口、存储器和通信总线,其中,处理器,通信接口,存储器通过通信总线完成相互间的通信;

存储器,用于存放计算机程序;

处理器,用于执行存储器上所存放的程序时,实现权利要求1-5任一所述的方法步骤。

一种超参数确定方法、装置及设备

技术领域

[0001] 本发明涉及计算机技术领域,特别是涉及一种超参数确定方法、装置、及设备。

背景技术

[0002] 分类算法在智能识别领域的应用越来越广泛,其中,主流的分类算法包括:决策树、随机森林、梯度提升决策树(GBDT)和xgboost(Extreme Gradient Boosting)等。程序员可以通过分类算法构建分类模型,然后利用分类模型对数据进行快速、准确的分类,分类后的数据作为其他应用的基础。例如,通过分类模型对用户的数据进行分类,提取出用于表征用户信誉的贷款数据、还款数据等,进而根据所提取的数据判断用户是否为存在违约风险的用户。

[0003] 然而,在构建分类模型过程中调整分类算法中包含的超参数时,通常需要程序员凭借经验进行反复调整,才能使得构建的分类模型对数据进行分类时得到理想的分类结果,这一调整方式不仅非常耗费程序员的精力,而且繁琐、效率低下。

发明内容

[0004] 本发明实施例的目的在于提供一种超参数确定方法、装置及设备,以实现降低程序员的工作量,并提高超参数的确定效率。具体技术方案如下:

[0005] 本发明实施的一方面,提供了一种超参数确定方法,所述方法包括:

[0006] 获得第一预设数量个超参数组,其中,所述超参数组包括:预设分类算法中各个超参数的一个取值,且任意两个超参数组中存在至少一个超参数的取值不同;

[0007] 获得第一数据集和第二数据集,分别以每一超参数组中的取值作为所述预设分类算法中超参数的取值,通过所述预设的分类算法对所述第一数据集和第二数据集中的数据进行分类处理,得到每一超参数组对应的第一置信度和每一超参数组对应的第二置信度,其中,第一置信度为:第一数据集分类结果的置信度,第二置信度为:第二数据集分类结果的置信度;

[0008] 根据每一超参数组对应的第一置信度和第二置信度,计算每一超参数组的权值;

[0009] 根据已获得每一超参数组中各个超参数的取值和每一超参数组的权值,利用贝叶斯寻优算法估算各个超参数的一个取值,获得由估算得到的取值构成的新超参数组,并累计超参数组的个数;

[0010] 当累计后超参数组的个数未达到第二预设数量时,以新超参数组中的取值作为所述预设分类算法中超参数的取值,通过所述预设的分类算法对所述第一数据集和第二数据集中的数据进行分类处理,得到新超参数组对应的第一置信度和第二置信度;根据新超参数组对应的第一置信度和第二置信度,计算新超参数组的权值;返回所述根据已获得每一超参数组中各个超参数的取值和每一超参数组的权值,利用贝叶斯寻优算法估算各个超参数的一个取值,获得由估算得到的取值构成的新超参数组,并累计超参数组的个数的步骤;

[0011] 当累计后超参数组的个数达到第二预设数量时,将目标超参数组包含的取值确定

为所述预设分类算法中各个超参数的取值,其中,所述目标超参数组为权值最大的超参数组。

[0012] 可选的,所述获得第一预设数量个超参数组的步骤,包括:

[0013] 获取所述预设分类算法中各个超参数的取值范围;

[0014] 以任意两个超参数组中存在至少一个超参数的取值不同为选取原则,分别在各个超参数的取值范围内选取一个值,并获得由所选取的各个值构成的一个超参数组;

[0015] 返回所述以任意两个超参数组中存在至少一个超参数的取值不同为选取原则,分别在各个超参数的取值范围内选取一个值,并获得由所选取的各个值构成的一个超参数组的步骤,直至获得所述第一预设数量个超参数组。

[0016] 可选的,所述获得第一数据集和第二数据集,分别以每一超参数组中的取值作为所述预设分类算法中超参数的取值,通过所述预设的分类算法对所述第一数据集和第二数据集中的数据进行分类处理,得到每一超参数组对应的第一置信度和每一超参数组对应的第二置信度的步骤,包括:

[0017] 利用分层抽样法将包含待处理数据的待处理数据集划分为k个子集,其中, $k > 1$,且任意两个子集之间不包含相同的待处理数据;

[0018] 以从k个子集中选取k-1个子集合并作为第一数据集、剩余子集作为第二数据集的方式,确定k个第一数据集和k个第二数据集;

[0019] 针对每一超参数组,以该超参数组中的取值作为所述预设分类算法中超参数的取值,通过所述预设的分类算法分别对k个第一数据集和k个第二数据集中的数据进行分类处理,得到该超参数组对应的k个第一数据集分类结果的置信度和k个第二数据集分类结果的置信度;

[0020] 将每一超参数组对应的k个第一数据集分类结果的置信度的平均值确定为该超参数组对应的第一置信度,将每一超参数组对应的k个第二数据集分类结果的置信度的平均值确定为该超参数组对应的第二置信度。

[0021] 可选的,所述将目标超参数组包含的取值确定为所述预设分类算法中各个超参数的取值的步骤,包括:

[0022] 计算目标超参数组对应的第一置信度和第二置信度的差值;

[0023] 当所述差值小于预设第一阈值时,将所述目标超参数组包含的取值确定为所述预设分类算法中各个超参数的取值;

[0024] 当所述差值不小于预设第一阈值时,按照预设第一步长增大所述第一预设数量,按照预设第二步长增大所述第二预设数量,并返回所述获得第一预设数量个超参数组的步骤。

[0025] 可选的,所述根据每一超参数组对应的第一置信度和第二置信度,计算每一超参数组的权值的步骤,包括:

[0026] 通过以下表达式计算每一超参数组的权值:

$$[0027] \quad y(AUC) = \begin{cases} \frac{AUC_{test}}{abs(AUC_{train} - AUC_{test})} & abs(AUC_{train} - AUC_{test}) > M \\ \frac{AUC_{test}}{M} & abs(AUC_{train} - AUC_{test}) \leq M \end{cases}$$

[0028] 其中, y (AUC) 表示超参数组的权值, $\text{abs}()$ 表示取绝对值函数, AUC_{test} 表示超参数组对应的第二置信度, $\text{AUC}_{\text{train}}$ 表示超参数组对应的第一置信度, M 为预设数值。

[0029] 本发明实施的又一方面, 还提供了一种超参数确定装置, 所述装置包括:

[0030] 获取模块, 用于获得第一预设数量个超参数组, 其中, 所述超参数组包括: 预设分类算法中各个超参数的一个取值, 且任意两个超参数组中存在至少一个超参数的取值不同;

[0031] 处理模块, 用于获得第一数据集和第二数据集, 分别以每一超参数组中的取值作为所述预设分类算法中超参数的取值, 通过所述预设的分类算法对所述第一数据集和第二数据集中的数据进行分类处理, 得到每一超参数组对应的第一置信度和每一超参数组对应的第二置信度, 其中, 第一置信度为: 第一数据集分类结果的置信度, 第二置信度为: 第二数据集分类结果的置信度;

[0032] 计算模块, 用于根据每一超参数组对应的第一置信度和第二置信度, 计算每一超参数组的权值;

[0033] 累计模块, 用于根据已获得每一超参数组中各个超参数的取值和每一超参数组的权值, 利用贝叶斯寻优算法估算各个超参数的一个取值, 获得由估算得到的取值构成的新超参数组, 并累计超参数组的个数;

[0034] 返回模块, 用于当累计后超参数组的个数未达到第二预设数量时, 以新超参数组中的取值作为所述预设分类算法中超参数的取值, 通过所述预设的分类算法对所述第一数据集和第二数据集中的数据进行分类处理, 得到新超参数组对应的第一置信度和第二置信度; 根据新超参数组对应的第一置信度和第二置信度, 计算新超参数组的权值, 并触发所述累计模块;

[0035] 确定模块, 用于当累计后超参数组的个数达到第二预设数量时, 将目标超参数组包含的取值确定为所述预设分类算法中各个超参数的取值, 其中, 所述目标超参数组为权值最大的超参数组。

[0036] 可选的, 所述获得模块包括:

[0037] 取值范围获取子模块, 用于获取所述预设分类算法中各个超参数的取值范围;

[0038] 超参数组获得子模块, 用于以任意两个超参数组中存在至少一个超参数的取值不同为选取原则, 分别在各个超参数的取值范围内选取一个值, 并获得由所选取的各个值构成的一个超参数组; 返回所述以任意两个超参数组中存在至少一个超参数的取值不同为选取原则, 分别在各个超参数的取值范围内选取一个值, 并获得由所选取的各个值构成的一个超参数组的步骤, 直至获得所述第一预设数量个超参数组。

[0039] 可选的, 所述处理模块具体用于,

[0040] 利用分层抽样法将包含待处理数据的待处理数据集划分为 k 个子集, 其中, $k > 1$, 且任意两个子集之间不包含相同的待处理数据;

[0041] 以从 k 个子集中选取 $k-1$ 个子集合并作为第一数据集、剩余子集作为第二数据集的方式, 确定 k 个第一数据集和 k 个第二数据集;

[0042] 针对每一超参数组, 以该超参数组中的取值作为所述预设分类算法中超参数的取值, 通过所述预设的分类算法分别对 k 个第一数据集和 k 个第二数据集中的数据进行分类处理, 得到该超参数组对应的 k 个第一数据集分类结果的置信度和 k 个第二数据集分类结果的

置信度；

[0043] 将每一超参数组对应的k个第一数据集分类结果的置信度的平均值确定为该超参数组对应的第一置信度,将每一超参数组对应的k个第二数据集分类结果的置信度的平均值确定为该超参数组对应的第二置信度。

[0044] 可选的,所述确定模块具体用于,

[0045] 计算目标超参数组对应的第一置信度和第二置信度的差值;

[0046] 当所述差值小于预设第一阈值时,将所述目标超参数组包含的取值确定为所述预设分类算法中各个超参数的取值;

[0047] 当所述差值不小于预设第一阈值时,按照预设第一步长增大所述第一预设数量,按照预设第二步长增大所述第二预设数量,并触发所述获取模块。

[0048] 可选的,所述计算模块具体用于,

[0049] 通过以下表达式计算每一超参数组的权值:

$$[0050] \quad y(AUC) = \begin{cases} \frac{AUC_{test}}{abs(AUC_{train} - AUC_{test})} & abs(AUC_{train} - AUC_{test}) > M \\ \frac{AUC_{test}}{M} & abs(AUC_{train} - AUC_{test}) \leq M \end{cases}$$

[0051] 其中,y(AUC)表示超参数组的权值,abs()表示取绝对值函数,AUC_{test}表示超参数组对应的第二置信度,AUC_{train}表示超参数组对应的第一置信度,M为预设数值。

[0052] 本发明实施的又一方面,还提供了一种电子设备,其特征在于,包括处理器、通信接口、存储器和通信总线,其中,处理器,通信接口,存储器通过通信总线完成相互间的通信;

[0053] 存储器,用于存放计算机程序;

[0054] 处理器,用于执行存储器上所存放的程序时,实现上述任一所述的超参数确定方法。

[0055] 在本发明实施的又一方面,还提供了一种计算机可读存储介质,所述计算机可读存储介质中存储有指令,当其在计算机上运行时,使得计算机执行上述任一所述的超参数确定方法。

[0056] 在本发明实施的又一方面,本发明实施例还提供了一种包含指令的计算机程序产品,当其在计算机上运行时,使得计算机执行上述任一所述的超参数确定方法。

[0057] 本发明实施例提供的超参数确定方法、装置及设备,可以根据所获得的每一超参数组中各个超参数的取值和所获得的每一超参数组的权值,通过贝叶斯寻优算法估算各个超参数的一个取值,获得由估算得到的取值构成的新超参数组,并累计超参数组的个数,当累计后超参数组的数量达到第二预设数量时,将权值最大的超参数组包含的取值确定为预设分类算法中各个超参数的取值。应用本发明实施例提供的技术方案无需程序员对超参数进行反复调整,因而,能够降低程序员的工作量,并且通过自动确定超参数的方式能够提高超参数的确定效率。

附图说明

[0058] 为了更清楚地说明本发明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。

[0059] 图1为本发明实施例提供的一种超参数确定方法的流程示意图;

[0060] 图2为本发明实施例提供的一种超参数确定装置的结构示意图;

[0061] 图3为本发明实施例提供的一种电子设备的结构示意图。

具体实施方式

[0062] 下面将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。基于本发明中的实施例,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,都属于本发明保护的范围。

[0063] 参见图1,示出了本发明实施例提供的一种超参数确定方法的流程示意图,该方法包括:

[0064] S100,获得第一预设数量个超参数组,其中,超参数组包括:预设分类算法中各个超参数的一个取值,且任意两个超参数组中存在至少一个超参数的取值不同。

[0065] 第一预设数量可以根据实际需要来进行设定,第一预设数量越大,根据最终确定的超参数组中的各个值作为预设分类算法中超参数的取值时,利用预设分类算法对待处理数据集进行分类处理时得到的处理结果越符合预期。

[0066] 预设分类算法不同相应的超参数也不同。例如,预设分类算法为xgboost时,超参数包括:max_depth、colsample_bytree、min_child_weight、scale_pos_weight等超参数。

[0067] 超参数组可以理解为预设分类算法中各个超参数的一个取值的集合,当两个超参数组中各个超参数的取值完全相同时,以这两组超参数组中的取值作为预设分类算法中超参数的取值时,对待处理数据进行分类处理时得到的结果是相同的,也就是进行了重复处理,因此,为了避免这一情况的发生,可以设定任意两个超参数组中存在至少一个超参数的取值不同。例如,预设分类算法中超参数包括:A、B、C;那么,如果两个超参数组中超参数A和B的取值相同,那么这两个超参数组中超参数C的取值则必须不同。

[0068] 一种实现方式中,可以通过以下步骤来获得第一预设数量个超参数组:

[0069] 步骤一,获取预设分类算法中各个超参数的取值范围;

[0070] 步骤二,以任意两个超参数组中存在至少一个超参数的取值不同为选取原则,分别在各个超参数的取值范围内选取一个值,并获得由所选取的各个值构成的一个超参数组;

[0071] 步骤三,返回步骤二,直至获得第一预设数量个超参数组

[0072] 实际应用中,各个超参数的取值可以根据实际需要来进行设定,例如,预设分类算法为xgboost时,超参数max_depth的取值范围可以为[3,11]。

[0073] S110,获得第一数据集和第二数据集,分别以每一超参数组中的取值作为预设分类算法中超参数的取值,通过预设的分类算法对第一数据集和第二数据集中的数据进行分类处理,得到每一超参数组对应的第一置信度和每一超参数组对应的第二置信度,其中,第一置信度为:第一数据集分类结果的置信度,第二置信度为:第二数据集分类结果的置信

度。

[0074] 置信度可以理解为根据分类算法对数据进行分类处理得到的处理结果,评价处理结果是否可信的指标,置信度越高则表明分类算法对数据处理的处理结果越可信,也就表示该分类算法对数据处理的效果越好。一种实现方式中,可以利用AUC(Area Under Curve)的值来作为置信度。

[0075] 一种实现方式中,可以将待处理数据集划分为第一数据集和第二数据集,将第一数据集作为训练集,将第二数据集作为测试集;然后,利用一超参数组中的取值作为预设分类算法中超参数的取值,通过预设的分类算法对训练集中的数据进行分类处理,得到训练集分类结果的AUC值;利用通过预设分类算法对测试集中的数据进行分类处理,得到测试集分类结果的AUC值,通过训练集分类结果的AUC值和测试集分类结果的AUC值进行比较,来验证该超参数组中的取值作为预设分类算法中超参数的取值对数据处理的效果。

[0076] 一种实现方式中,可以按照预设的比例将待处理数据集划分为第一数据集和第二数据集,例如,第一数据集中包括的数据与第二数据集中包括的数据的比例为2:1。

[0077] 为了能够更全面的检验每一超参数组中的取值作为预设分类算法中超参数的取值时,对待处理数据集进行分类处理效果是否符合预期,本发明实施例一种实现方式中,可以通过以下步骤得到第一数据集的分类结果和第二数据集的分类结果:

[0078] 步骤一,利用分层抽样法将包含待处理数据的待处理数据集划分为k个子集,其中, $k > 1$,且任意两个子集之间不包含相同的待处理数据;

[0079] 步骤二,以从k个子集中选取k-1个子集合并作为第一数据集、剩余子集作为第二数据集的方式,确定k个第一数据集和k个第二数据集;

[0080] 步骤三,针对每一超参数组,以该超参数组中的取值作为所述预设分类算法中超参数的取值,通过预设的分类算法分别对k个第一数据集和k个第二数据集中的数据进行分类处理,得到该超参数组对应的k个第一数据集分类结果的置信度和k个第二数据集分类结果的置信度;

[0081] 步骤四,将每一超参数组对应的k个第一数据集分类结果的置信度的平均值确定为该超参数组对应的第一置信度,将每一超参数组对应的k个第二数据集分类结果的置信度的平均值确定为该超参数组对应的第二置信度。

[0082] 分层抽样法可以理解为按照待处理数据集中不同类型数据的比例,将待处理数据集划分为k个子集,每个子集中不同类型数据的比例与待处理数据集中不同类型数据的比例相同。例如,待处理数据集中包括:贷款类型数据、还款类型数据,两个类型数据的比例为3:1,则将待处理数据集划分为k个子集后,每个子集中贷款类型数据和还款类型数据的比例为3:1。

[0083] 在k个子集中选取k-1个子集存在k中选择方式,因此,能够得到k组第一数据集和k组第二数据集。例如,k为3,3个子集分别为为:1、2、3;第一种选择方式:当第一数据集为:1、2时,相应地第二数据集为:3;第二种选择方式:当第一数据集为:1、3时,相应地第二数据集为:2;第三种选择方式:当第一数据集为2、3时,相应地第二数据集为:1。

[0084] S120,根据每一超参数组对应的第一置信度和第二置信度,计算每一超参数组的权值。

[0085] 一种实现方式中,通过以下表达式计算每一超参数组的权值:

$$[0086] \quad y(AUC) = \begin{cases} \frac{AUC_{test}}{abs(AUC_{train} - AUC_{test})} & abs(AUC_{train} - AUC_{test}) > M \\ \frac{AUC_{test}}{M} & abs(AUC_{train} - AUC_{test}) \leq M \end{cases}$$

[0087] 其中, $y(AUC)$ 表示超参数组的权值, $abs()$ 表示取绝对值函数, AUC_{test} 表示超参数组对应的第二置信度, AUC_{train} 表示超参数组对应的第一置信度, M 为预设数值。

[0088] M 可以根据实际需要来进行设定, 一种实现方式中, M 可以取 $[0.01, 0.03]$ 。一种更优选的实现方式中, M 可以取 0.02 。

[0089] S130, 根据已获得每一超参数组中各个超参数的取值和每一超参数组的权值, 利用贝叶斯寻优算法估算各个超参数的一个取值, 获得由估算得到的取值构成的新超参数组, 并累计超参数组的个数。

[0090] 一种实现方式中, 可以以每一个超参数组包含的取值为自变量, 以每一个超参数组的权值为因变量, 通过贝叶斯寻优算法建立两者之间的对应关系, 然后利用建立的对应关系估算各个超参数的一个取值, 从而获得由估算得到的取值构成的新超参数组, 并累计超参数组的个数。例如, 获得的超参数为 5 个, 则根据这 5 个超参数组包含的超参数和 5 个超参数组的权值, 利用贝叶斯寻优算法估算一个新的超参数组并进行累计, 累计后超参数组的个数则为 6 组。

[0091] S140, 当累计后超参数组的个数未达到第二预设数量时, 以新超参数组中的取值作为预设分类算法中超参数的取值, 通过预设的分类算法对第一数据集和第二数据集中的数据进行分类处理, 得到新超参数组对应的第一置信度和第二置信度; 根据新超参数组对应的第一置信度和第二置信度, 计算新超参数组的权值; 并返回执行步骤 S130。

[0092] 第二预设数量可以根据实际需要来进行设定第二预设数量越大, 根据最终确定的超参数组中的各个值作为预设分类算法中超参数的取值时, 利用预设分类算法对待处理数据集进行分类处理时得到的处理结果越符合预期。

[0093] 每次利用贝叶斯寻优算法进行估算得到一个超参数组时, 根据上一次累计后的超参数组中的取值和超参数组的权值来进行预估, 例如, 本次进行估算时根据上一次累计后的 6 个超参数组和 6 个超参数的权值估算第 7 个超参数组, 在进行下一次预估时根据本次累计后的 7 个超参数组和 7 个超参数组的权值估算第 8 个超参数组。基于此, 经过累计, 超参数组的数量逐渐增多, 也就是通过贝叶斯寻优算法建立超参数组和超参数的权值之间的对应关系时, 所使用的自变量和因变量也越来越多, 使得所建立的对应关系越来越丰富。

[0094] S150, 当累计后超参数组的个数达到第二预设数量时, 将目标超参数组包含的取值确定为预设分类算法中各个超参数的取值, 其中, 目标超参数组为权值最大的超参数组。

[0095] 一种实现方式中, 可以直接将累计后超参数组中权值最大的超参数组包含的取值确定为预设分类算法中各个超参数的取值。

[0096] 而为了最终确定的目标超参数包含的取值作为预设分类算法中各个超参数的取值时, 对待处理数据集进行分类处理得到更好的处理结果。本发明实施例一种实现方式中, 将目标超参数组包含的取值确定为预设分类算法中各个超参数的取值的步骤, 可以包括:

[0097] 步骤一, 计算目标超参数组对应的第一置信度和第二置信度的差值;

[0098] 步骤二, 当差值小于预设第一阈值时, 将目标超参数组包含的取值确定为预设分

类算法中各个超参数的取值；

[0099] 步骤三,当差值不小于预设第一阈值时,按照预设第一步长增大第一预设数量,按照预设第二步长增大第二预设数量,并返回步骤S100。

[0100] 预设第一阈值可以根据实际需要来进行设定,预设第一阈值越小,根据最终确定的超参数组中的各个值作为预设分类算法中超参数的取值时,利用预设分类算法对待处理数据集进行分类处理时得到的处理结果越符合预期。

[0101] 预设第一步长和预设第二步长可以根据实际需要进行设定,两者可以相同,也可以不同。在实际应用中,可以按照预设第一步长将第一预设数量增大5、10等;可以按照第二步预设步长将第二预设数量增大到20、35、50等。一般情况下,按照预设第一步长将第一预设数量增大到10,按照第二步预设步长将第二预设数量增大到50,以最终得到的50组超参数组中权值最大的超参数组中的取值作为分类算法中超参数的取值时,对待处理数据集进行分类处理时即可得到符合预期的处理结果。

[0102] 本发明实施例提供的超参数确定方法,可以根据所获得的每一超参数组中各个超参数的取值和所获得的每一超参数组的权值,通过贝叶斯寻优算法估算各个超参数的一个取值,获得由估算得到的取值构成的新超参数组,并累计超参数组的个数,当累计后超参数组的数量达到第二预设数量时,将权值最大的超参数组包含的取值确定为预设分类算法中各个超参数的取值。应用本发明实施例提供的技术方案无需程序员对超参数进行反复调整,因而,能够降低程序员的工作量,并且通过自动确定超参数的方式能够提高超参数的确定效率。

[0103] 参照图2,示出了本发明实施例提供的一种超参数确定装置的结构示意图,该装置包括:

[0104] 获取模块200,用于获得第一预设数量个超参数组,其中,所述超参数组包括:预设分类算法中各个超参数的一个取值,且任意两个超参数组中存在至少一个超参数的取值不同;

[0105] 处理模块210,用于获得第一数据集和第二数据集,分别以每一超参数组中的取值作为所述预设分类算法中超参数的取值,通过所述预设的分类算法对所述第一数据集和第二数据集中的数据进行分类处理,得到每一超参数组对应的第一置信度和每一超参数组对应的第二置信度,其中,第一置信度为:第一数据集分类结果的置信度,第二置信度为:第二数据集分类结果的置信度;

[0106] 计算模块220,用于根据每一超参数组对应的第一置信度和第二置信度,计算每一超参数组的权值;

[0107] 累计模块230,用于根据已获得每一超参数组中各个超参数的取值和每一超参数组的权值,利用贝叶斯寻优算法估算各个超参数的一个取值,获得由估算得到的取值构成的新超参数组,并累计超参数组的个数;

[0108] 返回模块240,用于当累计后超参数组的个数未达到第二预设数量时,以新超参数组中的取值作为所述预设分类算法中超参数的取值,通过所述预设的分类算法对所述第一数据集和第二数据集中的数据进行分类处理,得到新超参数组对应的第一置信度和第二置信度;根据新超参数组对应的第一置信度和第二置信度,计算新超参数组的权值,并触发所述累计模块;

[0109] 确定模块250,用于当累计后超参数组的个数达到第二预设数量时,将目标超参数组包含的取值确定为所述预设分类算法中各个超参数的取值,其中,所述目标超参数组为权值最大的超参数组。

[0110] 本发明实施例一种实现方式中,获得模块200包括:

[0111] 取值范围获取子模块,用于获取所述预设分类算法中各个超参数的取值范围;

[0112] 超参数组获得子模块,用于以任意两个超参数组中存在至少一个超参数的取值不同为选取原则,分别在各个超参数的取值范围内选取一个值,并获得由所选取的各个值构成的一个超参数组;返回所述以任意两个超参数组中存在至少一个超参数的取值不同为选取原则,分别在各个超参数的取值范围内选取一个值,并获得由所选取的各个值构成的一个超参数组的步骤,直至获得所述第一预设数量个超参数组。

[0113] 本发明实施例一种实现方式中,处理模块220具体用于,

[0114] 利用分层抽样法将包含待处理数据的待处理数据集划分为k个子集,其中, $k > 1$,且任意两个子集之间不包含相同的待处理数据;

[0115] 以从k个子集中选取k-1个子集合并作为第一数据集、剩余子集作为第二数据集的方式,确定k个第一数据集和k个第二数据集;

[0116] 针对每一超参数组,以该超参数组中的取值作为所述预设分类算法中超参数的取值,通过所述预设的分类算法分别对k个第一数据集和k个第二数据集中的数据进行分类处理,得到该超参数组对应的k个第一数据集分类结果的置信度和k个第二数据集分类结果的置信度;

[0117] 将每一超参数组对应的k个第一数据集分类结果的置信度的平均值确定为该超参数组对应的第一置信度,将每一超参数组对应的k个第二数据集分类结果的置信度的平均值确定为该超参数组对应的第二置信度。

[0118] 本发明实施例一种实现方式中,所述确定模块260具体用于,

[0119] 计算目标超参数组对应的第一置信度和第二置信度的差值;

[0120] 当所述差值小于预设第一阈值时,将所述目标超参数组包含的取值确定为所述预设分类算法中各个超参数的取值;

[0121] 当所述差值不小于预设第一阈值时,按照预设第一步长增大所述第一预设数量,按照预设第二步长增大所述第二预设数量,并触发获取模块200。

[0122] 本发明实施例一种实现方式中,计算模块230具体用于,

[0123] 通过以下表达式计算每一超参数组的权值:

$$[0124] \quad y(AUC) = \begin{cases} \frac{AUC_{test}}{abs(AUC_{train} - AUC_{test})} & abs(AUC_{train} - AUC_{test}) > M \\ \frac{AUC_{test}}{M} & abs(AUC_{train} - AUC_{test}) \leq M \end{cases}$$

[0125] 其中, $y(AUC)$ 表示超参数组的权值, $abs()$ 表示取绝对值函数, AUC_{test} 表示超参数组对应的第二置信度, AUC_{train} 表示超参数组对应的第一置信度, M 为预设数值。

[0126] 本发明实施例提供的各个方案中,超参数确定装置可以根据所获得的每一超参数组中各个超参数的取值和所获得的每一超参数组的权值,通过贝叶斯寻优算法估算各个超参数的一个取值,获得由估算得到的取值构成的新超参数组,并累计超参数组的个数,当累

计后超参数组的数量达到第二预设数量时,将权值最大的超参数组包含的取值确定为预设分类算法中各个超参数的取值。应用本发明实施例提供的技术方案无需程序员对超参数进行反复调整,因而,能够降低程序员的工作量,并且通过自动确定超参数的方式能够提高超参数的确定效率。

[0127] 本发明实施例还提供了一种电子设备,如图3所示,包括处理器001、通信接口002、存储器003和通信总线004,其中,处理器001,通信接口002,存储器003通过通信总线004完成相互间的通信,

[0128] 存储器003,用于存放计算机程序;

[0129] 处理器001,用于执行存储器003上所存放的程序时,实现本发明实施例提供的超参数确定方法。

[0130] 具体的,上述方法包括:

[0131] 获得第一预设数量个超参数组,其中,所述超参数组包括:预设分类算法中各个超参数的一个取值,且任意两个超参数组中存在至少一个超参数的取值不同;

[0132] 获得第一数据集和第二数据集,分别以每一超参数组中的取值作为所述预设分类算法中超参数的取值,通过所述预设的分类算法对所述第一数据集和第二数据集中的数据进行分类处理,得到每一超参数组对应的第一置信度和每一超参数组对应的第二置信度,其中,第一置信度为:第一数据集分类结果的置信度,第二置信度为:第二数据集分类结果的置信度;

[0133] 根据每一超参数组对应的第一置信度和第二置信度,计算每一超参数组的权值;

[0134] 根据已获得每一超参数组中各个超参数的取值和每一超参数组的权值,利用贝叶斯寻优算法估算各个超参数的一个取值,获得由估算得到的取值构成的新超参数组,并累计超参数组的个数;

[0135] 当累计后超参数组的个数未达到第二预设数量时,以新超参数组中的取值作为所述预设分类算法中超参数的取值,通过所述预设的分类算法对所述第一数据集和第二数据集中的数据进行分类处理,得到新超参数组对应的第一置信度和第二置信度;根据新超参数组对应的第一置信度和第二置信度,计算新超参数组的权值;返回所述根据已获得每一超参数组中各个超参数的取值和每一超参数组的权值,利用贝叶斯寻优算法估算各个超参数的一个取值,获得由估算得到的取值构成的新超参数组,并累计超参数组的个数的步骤;

[0136] 当累计后超参数组的个数达到第二预设数量时,将目标超参数组包含的取值确定为所述预设分类算法中各个超参数的取值,其中,所述目标超参数组为权值最大的超参数组。

[0137] 需要说明的是,上述处理器011执行存储器013上所存放的程序实现超参数确定方法的其他实施例,与前述方法实施例部分提供的实施例相同,这里不再赘述。

[0138] 本发明实施例提供的各个方案中,电子设备可以根据所获得的每一超参数组中各个超参数的取值和所获得的每一超参数组的权值,通过贝叶斯寻优算法估算各个超参数的一个取值,获得由估算得到的取值构成的新超参数组,并累计超参数组的个数,当累计后超参数组的数量达到第二预设数量时,将权值最大的超参数组包含的取值确定为预设分类算法中各个超参数的取值。应用本发明实施例提供的技术方案无需程序员对超参数进行反复调整,因而,能够降低程序员的工作量,并且通过自动确定超参数的方式能够提高超参数的

确定效率。

[0139] 上述电子设备提到的通信总线可以是外设部件互连标准 (Peripheral Component Interconnect, PCI) 总线或扩展工业标准结构 (Extended Industry Standard Architecture, EISA) 总线等。该通信总线可以分为地址总线、数据总线、控制总线等。为便于表示,图中仅用一条粗线表示,但并不表示仅有一根总线或一种类型的总线。

[0140] 通信接口用于上述电子设备与其他设备之间的通信。

[0141] 存储器可以包括随机存取存储器 (Random Access Memory, RAM),也可以包括非易失性存储器 (Non-Volatile Memory, NVM),例如至少一个磁盘存储器。可选的,存储器还可以是至少一个位于远离前述处理器的存储装置。

[0142] 上述的处理器可以是通用处理器,包括中央处理器 (Central Processing Unit, CPU)、网络处理器 (Network Processor, NP) 等;还可以是数字信号处理器 (Digital Signal Processing, DSP)、专用集成电路 (Application Specific Integrated Circuit, ASIC)、现场可编程门阵列 (Field-Programmable Gate Array, FPGA) 或者其他可编程逻辑器件、分立门或者晶体管逻辑器件、分立硬件组件。

[0143] 在本发明提供的又一实施例中,还提供了一种计算机可读存储介质,该计算机可读存储介质中存储有指令,当其在计算机上运行时,本发明实施例提供的超参数确定方法。

[0144] 具体的,上述超参数确定方法,包括:

[0145] 获得第一预设数量个超参数组,其中,所述超参数组包括:预设分类算法中各个超参数的一个取值,且任意两个超参数组中存在至少一个超参数的取值不同;

[0146] 获得第一数据集和第二数据集,分别以每一超参数组中的取值作为所述预设分类算法中超参数的取值,通过所述预设的分类算法对所述第一数据集和第二数据集中的数据进行分类处理,得到每一超参数组对应的第一置信度和每一超参数组对应的第二置信度,其中,第一置信度为:第一数据集分类结果的置信度,第二置信度为:第二数据集分类结果的置信度;

[0147] 根据每一超参数组对应的第一置信度和第二置信度,计算每一超参数组的权值;

[0148] 根据已获得每一超参数组中各个超参数的取值和每一超参数组的权值,利用贝叶斯寻优算法估算各个超参数的一个取值,获得由估算得到的取值构成的新超参数组,并累计超参数组的个数;

[0149] 当累计后超参数组的个数未达到第二预设数量时,以新超参数组中的取值作为所述预设分类算法中超参数的取值,通过所述预设的分类算法对所述第一数据集和第二数据集中的数据进行分类处理,得到新超参数组对应的第一置信度和第二置信度;根据新超参数组对应的第一置信度和第二置信度,计算新超参数组的权值;返回所述根据已获得每一超参数组中各个超参数的取值和每一超参数组的权值,利用贝叶斯寻优算法估算各个超参数的一个取值,获得由估算得到的取值构成的新超参数组,并累计超参数组的个数的步骤;

[0150] 当累计后超参数组的个数达到第二预设数量时,将目标超参数组包含的取值确定为所述预设分类算法中各个超参数的取值,其中,所述目标超参数组为权值最大的超参数组。

[0151] 需要说明的是,通过上述计算机可读存储介质实现超参数确定方法的其他实施例,与前述方法实施例部分提供的实施例相同,这里不再赘述。

[0152] 本发明实施例提供的各个方案中,通过运行上述计算机可读存储介质中存储的指令,可以根据所获得的每一超参数组中各个超参数的取值和所获得的每一超参数组的权值,通过贝叶斯寻优算法估算各个超参数的一个取值,获得由估算得到的取值构成的新超参数组,并累计超参数组的个数,当累计后超参数组的数量达到第二预设数量时,将权值最大的超参数组包含的取值确定为预设分类算法中各个超参数的取值。应用本发明实施例提供的技术方案无需程序员对超参数进行反复调整,因而,能够降低程序员的工作量,并且通过自动确定超参数的方式能够提高超参数的确定效率。

[0153] 在本发明提供的又一实施例中,还提供了一种包含指令的计算机程序产品,当其在计算机上运行时,实现本发明实施例提供的超参数确定方法。

[0154] 具体的,上述超参数确定方法,包括:

[0155] 获得第一预设数量个超参数组,其中,所述超参数组包括:预设分类算法中各个超参数的一个取值,且任意两个超参数组中存在至少一个超参数的取值不同;

[0156] 获得第一数据集和第二数据集,分别以每一超参数组中的取值作为所述预设分类算法中超参数的取值,通过所述预设的分类算法对所述第一数据集和第二数据集中的数据进行分类处理,得到每一超参数组对应的第一置信度和每一超参数组对应的第二置信度,其中,第一置信度为:第一数据集分类结果的置信度,第二置信度为:第二数据集分类结果的置信度;

[0157] 根据每一超参数组对应的第一置信度和第二置信度,计算每一超参数组的权值;

[0158] 根据已获得每一超参数组中各个超参数的取值和每一超参数组的权值,利用贝叶斯寻优算法估算各个超参数的一个取值,获得由估算得到的取值构成的新超参数组,并累计超参数组的个数;

[0159] 当累计后超参数组的个数未达到第二预设数量时,以新超参数组中的取值作为所述预设分类算法中超参数的取值,通过所述预设的分类算法对所述第一数据集和第二数据集中的数据进行分类处理,得到新超参数组对应的第一置信度和第二置信度;根据新超参数组对应的第一置信度和第二置信度,计算新超参数组的权值;返回所述根据已获得每一超参数组中各个超参数的取值和每一超参数组的权值,利用贝叶斯寻优算法估算各个超参数的一个取值,获得由估算得到的取值构成的新超参数组,并累计超参数组的个数的步骤;

[0160] 当累计后超参数组的个数达到第二预设数量时,将目标超参数组包含的取值确定为所述预设分类算法中各个超参数的取值,其中,所述目标超参数组为权值最大的超参数组。

[0161] 需要说明的是,通过上述计算机程序产品实现超参数确定方法的其他实施例,与前述方法实施例部提供的实施例相同,这里不再赘述。

[0162] 本发明实施例提供的各个方案中,通过运行上述包含指令的计算机程序产品,可以根据所获得的每一超参数组中各个超参数的取值和所获得的每一超参数组的权值,通过贝叶斯寻优算法估算各个超参数的一个取值,获得由估算得到的取值构成的新超参数组,并累计超参数组的个数,当累计后超参数组的数量达到第二预设数量时,将权值最大的超参数组包含的取值确定为预设分类算法中各个超参数的取值。应用本发明实施例提供的技术方案无需程序员对超参数进行反复调整,因而,能够降低程序员的工作量,并且通过自动确定超参数的方式能够提高超参数的确定效率。

[0163] 需要说明的是,在本文中,诸如第一和第二等之类的关系术语仅仅用来将一个实体或者操作与另一个实体或操作区分开来,而不一定要求或者暗示这些实体或操作之间存在任何这种实际的关系或者顺序。而且,术语“包括”、“包含”或者其任何其他变体意在涵盖非排他性的包含,从而使得包括一系列要素的过程、方法、物品或者设备不仅包括那些要素,而且还包括没有明确列出的其他要素,或者是还包括为这种过程、方法、物品或者设备所固有的要素。在没有更多限制的情况下,由语句“包括一个……”限定的要素,并不排除在包括所述要素的过程、方法、物品或者设备中还存在另外的相同要素。

[0164] 本说明书中的各个实施例均采用相关的方式描述,各个实施例之间相同相似的部分互相参见即可,每个实施例重点说明的都是与其他实施例的不同之处。尤其,对于系统实施例而言,由于其基本相似于方法实施例,所以描述的比较简单,相关之处参见方法实施例的部分说明即可。

[0165] 以上所述仅为本发明的较佳实施例而已,并非用于限定本发明的保护范围。凡在本发明的精神和原则之内所作的任何修改、等同替换、改进等,均包含在本发明的保护范围内。

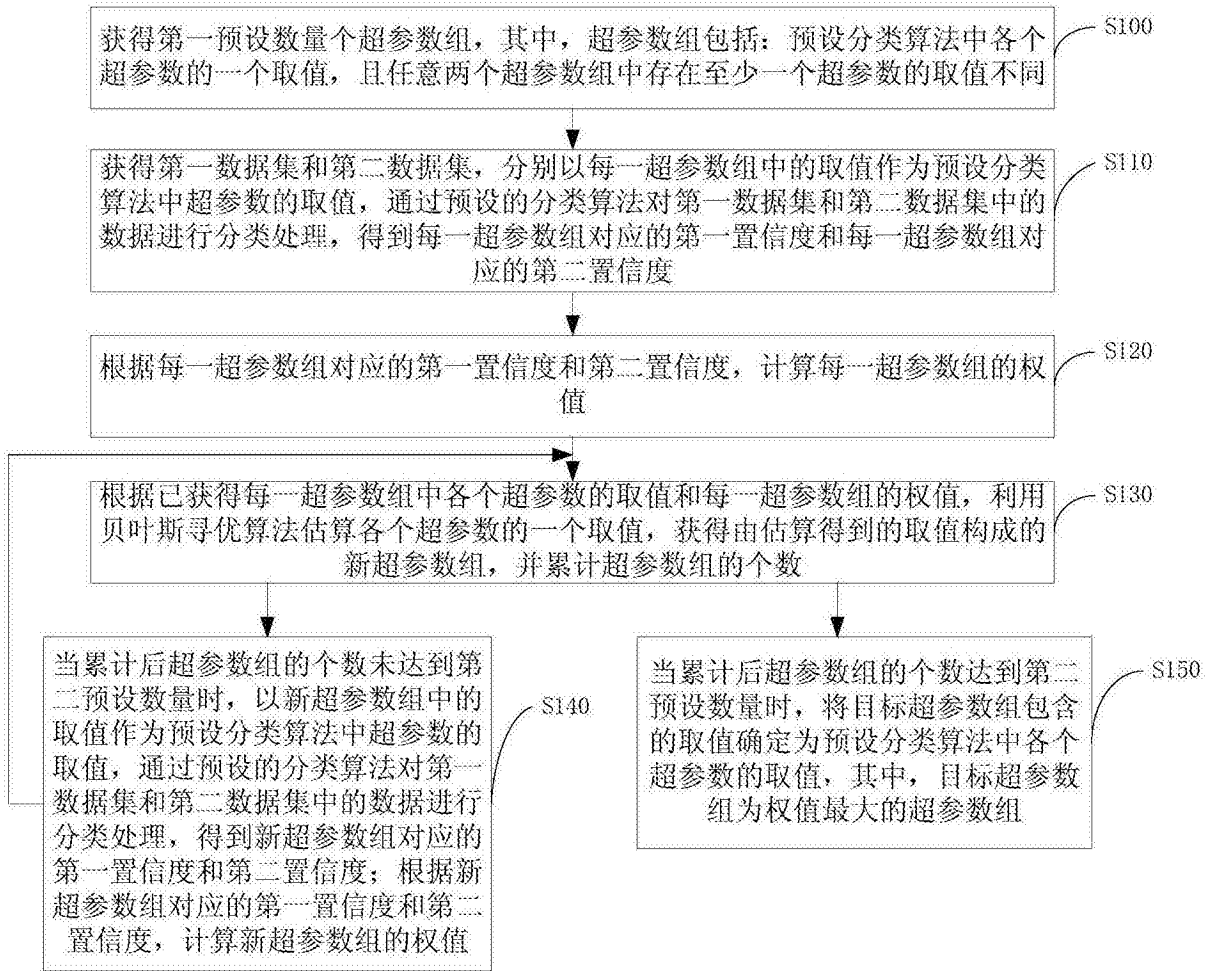


图1

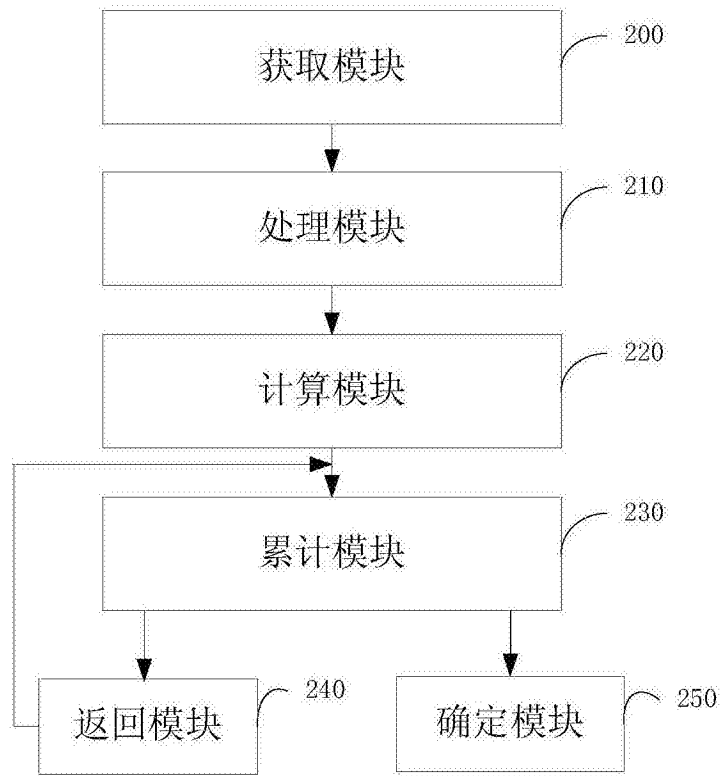


图2

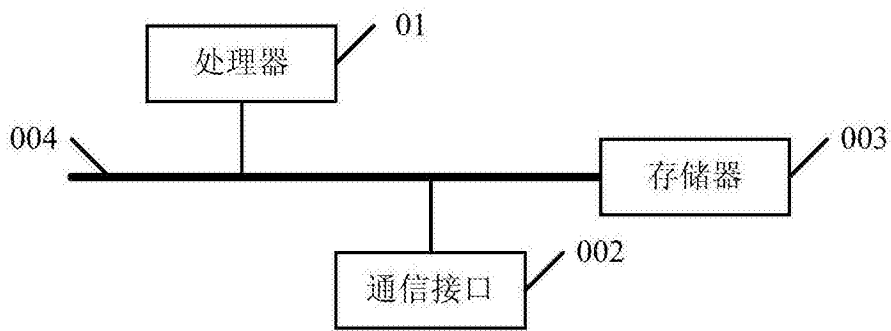


图3