



# (12) 发明专利申请

(10) 申请公布号 CN 115525921 A

(43) 申请公布日 2022. 12. 27

(21) 申请号 202210053611.7

(22) 申请日 2022.01.18

(71) 申请人 富算科技(上海)有限公司  
地址 200135 上海市浦东新区自由贸易试  
验区浦东大道1200号2层A区

(72) 发明人 尤志强 卞阳 赵东

(74) 专利代理机构 上海弼兴律师事务所 31283  
专利代理师 林嵩 罗朗

(51) Int. Cl.  
G06F 21/62 (2013.01)  
G06N 20/00 (2019.01)

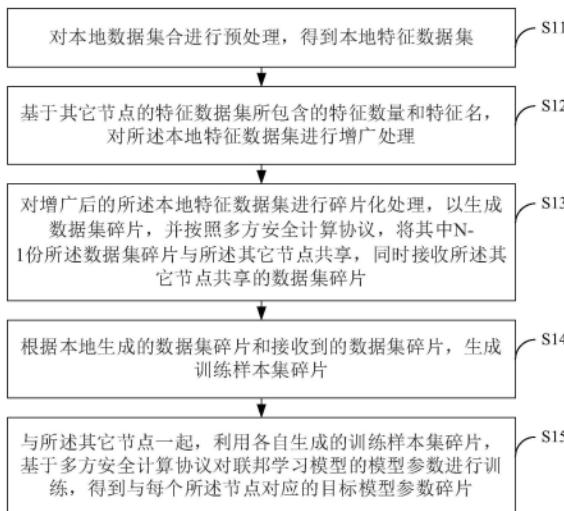
权利要求书3页 说明书13页 附图10页

## (54) 发明名称

基于MPC的联邦学习模型训练和预测方法、系统、设备及介质

## (57) 摘要

本发明提供一种基于MPC的联邦学习模型训练和预测方法、系统、设备及介质,该方法包括:对本地数据集合进行预处理,得到本地特征数据集;基于其它节点特征数据集的特征数量和特征名,对本地特征数据集进行增广;对增广后的本地特征数据集进行碎片化处理,以生成数据集碎片,并按照多方安全计算协议,将其中N-1份数据集碎片与所述其它节点共享,同时接收所述其它节点共享的数据集碎片;根据本地生成的数据集碎片和接收到的数据集碎片,生成训练样本集碎片;与其它节点一起,利用各自生成的训练样本集碎片,基于多方安全计算协议对联邦学习模型的模型参数进行训练,得到与每个节点对应的目标模型参数碎片。本发明能够实现联邦学习模型训练和预测的全流程安全保护。



1. 一种基于MPC的联邦学习模型训练方法,适用于合作执行纵向联邦学习任务的N个节点中的任意节点,其中N为大于1的整数,其特征在于,所述方法包括:

对本地数据集合进行预处理,得到本地特征数据集;

基于其它节点的特征数据集所包含的特征数量和特征名,对所述本地特征数据集进行增广处理;

对增广后的所述本地特征数据集进行碎片化处理,以生成数据集碎片,并按照多方安全计算协议,将其中N-1份所述数据集碎片与所述其它节点共享,同时接收所述其它节点共享的数据集碎片;

根据本地生成的数据集碎片和接收到的数据集碎片,生成训练样本集碎片;

与所述其它节点一起,利用各自生成的训练样本集碎片,基于多方安全计算协议对联邦学习模型的模型参数进行训练,得到与每个所述节点对应的目标模型参数碎片。

2. 如权利要求1所述的联邦学习模型训练方法,其特征在于,所述对本地数据集合进行预处理,得到本地特征数据集,包括:

根据发起联邦学习任务的节点所配置的预处理规则,对所述本地数据集合进行预处理,得到所述本地特征数据集。

3. 如权利要求2所述的联邦学习模型训练方法,其特征在于,所述预处理规则包括缺失值处理规则、异常值处理规则、特征转换规则、和/或标准化处理规则。

4. 如权利要求1所述的联邦学习模型训练方法,其特征在于,所述根据本地生成的数据集碎片和接收到的数据集碎片,生成训练样本集碎片,包括:

对本地生成的数据集碎片和接收到的数据集碎片进行拼接和聚合处理,以生成所述训练样本集碎片。

5. 如权利要求1所述的联邦学习模型训练方法,其特征在于,所述与所述其它节点一起,利用各自生成的训练样本集碎片,基于多方安全计算协议对联邦学习模型的模型参数进行训练,得到与每个所述节点对应的目标模型参数碎片,包括与所述其它节点一起执行以下步骤:

获取各自对应的模型参数碎片的初始值,所述初始值通过对所述联邦学习模型的初始模型参数进行碎片化处理得到;

基于多方安全计算协议,利用各自对应的所述模型参数碎片的初始值对各自生成的所述训练样本集碎片进行处理,得到与每个所述节点对应的预测结果碎片;

基于多方安全计算协议,根据各自对应的所述预测结果碎片与相应的标准标签碎片,计算与每个所述节点对应的损失值碎片;

基于多方安全计算协议,根据各自计算得到的所述损失值碎片对各自对应的所述模型参数碎片进行训练,直至满足预设的训练终止条件,得到与每个所述节点对应的所述目标模型参数碎片。

6. 如权利要求5所述的联邦学习模型训练方法,其特征在于,所述与所述其它节点一起基于多方安全计算协议,根据各自计算得到的所述损失值碎片对各自对应的所述模型参数碎片进行训练,包括:

与所述其它节点一起,基于多方安全计算协议获取各自对应的梯度值碎片,并根据各自对应的所述梯度值碎片对各自对应的所述模型参数碎片进行更新。

7. 一种基于MPC的联邦学习模型预测方法,适用于合作执行纵向联邦学习任务的N个节点中的任意节点,其中N为大于1的整数,其特征在于,所述方法包括:

对本地预测数据进行预处理,得到本地预测特征数据;

基于其它节点的预测特征数据所包含的特征数量和特征名,对所述本地预测特征数据进行增广处理;

对增广后的所述本地预测特征数据进行碎片化处理,以生成特征数据碎片,并按照多方安全计算协议,将其中N-1份所述特征数据碎片与所述其它节点共享,同时接收所述其它节点共享的特征数据碎片;

根据本地生成的特征数据碎片和接收到的特征数据碎片,生成目标预测数据碎片;

与所述其它节点一起,利用前述权利要求1-6中任一项所述的模型训练方法训练得到的与各所述节点对应的所述目标模型参数碎片,基于多方安全计算协议对各自生成的所述目标预测数据碎片进行处理,得到与每个所述节点对应的处理结果,并将所述处理结果共享给结果拥有方,以供所述结果拥有方根据各所述处理结果得到目标结果。

8. 一种基于MPC的联邦学习模型训练系统,所述系统适用于合作执行纵向联邦学习任务的N个节点,其中N为大于1的整数,其特征在于,所述系统包括:

第一预处理模块,用于分别对本地数据集进行预处理,得到本地特征数据集;

第一增广模块,用于基于其它节点的特征数据集所包含的特征数量和特征名,对所述本地特征数据集进行增广处理;

第一碎片化分发模块,用于对增广后的所述本地特征数据集进行碎片化处理,以生成数据集碎片,并按照多方安全计算协议,将其中N-1份所述数据集碎片与所述其它节点共享,同时接收所述其它节点共享的数据集碎片;

第一融合模块,用于根据本地生成的数据集碎片和接收到的数据集碎片,生成训练样本集碎片;

训练模块,用于与所述其它节点一起,利用各自生成的训练样本集碎片,基于多方安全计算协议对联邦学习模型的模型参数进行训练,得到与所述N个节点对应的目标模型参数碎片。

9. 一种基于MPC的联邦学习模型预测方系统,适用于合作执行纵向联邦学习任务的N个节点中的任意节点,其中N为大于1的整数,其特征在于,所述系统包括:

第二预处理模块,用于对本地预测数据进行预处理,得到本地预测特征数据;

第二增广模块,用于基于其它节点的预测特征数据所包含的特征数量和特征名,对所述本地预测特征数据进行增广处理;

第二碎片化分发模块,用于对所述本地预测特征数据进行碎片化处理,以生成特征数据碎片,并按照多方安全计算协议,将其中N-1份所述特征数据碎片与所述其它节点共享,同时接收所述其它节点共享的特征数据碎片;

第二融合模块,用于根据本地生成的特征数据碎片和接收到的特征数据碎片,生成目标预测数据碎片;

预测模块,用于与所述其它节点一起,利用前述权利要求8所述的模型训练系统训练得到的与各所述节点对应的所述目标模型参数碎片,基于多方安全计算协议对各自生成的所述目标预测数据碎片进行处理,得到与每个所述节点对应的处理结果,并将所述处理结果

共享给结果拥有方,以供所述结果拥有方根据各所述处理结果得到目标结果。

10.一种电子设备,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程序,其特征在于,所述处理器执行所述计算机程序时实现如权利要求1至7中任一项所述的方法。

11.一种计算机可读存储介质,其上存储有计算机程序,其特征在于,所述计算机程序被处理器执行时实现如权利要求1至7中任一项所述的方法。

## 基于MPC的联邦学习模型训练和预测方法、系统、设备及介质

### 技术领域

[0001] 本发明涉及联邦学习领域,尤其涉及一种基于MPC的联邦学习模型训练和预测方法、系统、设备及介质。

### 背景技术

[0002] 伴随着金融科技,尤其是互联网金融科技的快速发展,已经有越来越多的技术应用于金融领域,其中,联邦学习技术基于对用户隐私和数据的安全保障,正逐渐受到越来越多的重视。

[0003] 联邦学习(federated learning)是指,通过联合不同的参与者(participant,或者party,也称为数据所有者(data owner)、或者客户(client))进行机器学习建模的方法。在联邦学习中,参与者不需要向其它参与者和协调者(coordinator,也称为服务器(server),参数服务器(parameter server),或者聚合服务器(aggregation server))暴露自己所拥有的数据,因而联邦学习可以很好的保护用户隐私和保障数据安全,并可以解决数据孤岛问题。

[0004] 然而,在现有的纵向联邦学习中,因为不同参与者之间拥有的是不相同数据特征的数据,因此在各参与者每进行一次本地的模型训练之后,都需要各参与者交换各自的中间计算结果,特别是需要交换关于梯度信息的中间计算结果,因而无法实现全流程安全保护。

### 发明内容

[0005] 为了解决现有技术联邦学习模型训练方法无法实现全流程安全保护的技术问题,本发明提供一种基于MPC的联邦学习模型训练和预测方法、系统、设备及介质。

[0006] 为了实现上述目的,本发明采用以下技术方案:

[0007] 第一方面,本发明提供一种基于MPC的联邦学习模型训练方法,适用于合作执行纵向联邦学习任务的N个节点中的任意节点,其中N为大于1的整数,所述方法包括:

[0008] 对本地数据集合进行预处理,得到本地特征数据集;

[0009] 基于其它节点的特征数据集所包含的特征数量和特征名,对所述本地特征数据集进行增广处理;

[0010] 对增广后的所述本地特征数据集进行碎片化处理,以生成数据集碎片,并按照多方安全计算协议,将其中N-1份所述数据集碎片与所述其它节点共享,同时接收所述其它节点共享的数据集碎片;

[0011] 根据本地生成的数据集碎片和接收到的数据集碎片,生成训练样本集碎片;

[0012] 与所述其它节点一起,利用各自生成的训练样本集碎片,基于多方安全计算协议对联邦学习模型的模型参数进行训练,得到与每个所述节点对应的目标模型参数碎片。

[0013] 优选地,所述对本地数据集合进行预处理,得到本地特征数据集,包括:

[0014] 根据发起联邦学习任务的节点所配置的预处理规则,对所述本地数据集合进行预

处理,得到所述本地特征数据集。

[0015] 优选地,所述预处理规则包括缺失值处理规则、异常值处理规则、特征转换规则、和/或标准化处理规则。

[0016] 优选地,所述根据本地生成的数据集碎片和接收到的数据集碎片,生成训练样本集碎片,包括:

[0017] 对本地生成的数据集碎片和接收到的数据集碎片进行拼接和聚合处理,以生成所述训练样本集碎片。

[0018] 优选地,所述与所述其它节点一起,利用各自生成的训练样本集碎片,基于多方安全计算协议对联邦学习模型的模型参数进行训练,得到与每个所述节点对应的目标模型参数碎片,包括与所述其它节点一起执行以下步骤:

[0019] 获取各自对应的模型参数碎片的初始值,所述初始值通过对所述联邦学习模型的初始模型参数进行碎片化处理得到;

[0020] 基于多方安全计算协议,利用各自对应的所述模型参数碎片的初始值对各自生成的所述训练样本集碎片进行处理,得到与每个所述节点对应的预测结果碎片;

[0021] 基于多方安全计算协议,根据各自对应的所述预测结果碎片与相应的标准标签碎片,计算与每个所述节点对应的损失值碎片;

[0022] 基于多方安全计算协议,根据各自计算得到的所述损失值碎片对各自对应的所述模型参数碎片进行训练,直至满足预设的训练终止条件,得到与每个所述节点对应的所述目标模型参数碎片。

[0023] 优选地,所述与所述其它节点一起基于多方安全计算协议,根据各自计算得到的所述损失值碎片对各自对应的所述模型参数碎片进行训练,包括:

[0024] 与所述其它节点一起,基于多方安全计算协议获取各自对应的梯度值碎片,并根据各自对应的所述梯度值碎片对各自对应的所述模型参数碎片进行更新。

[0025] 第二方面,本发明提供一种基于MPC的联邦学习模型预测方法,适用于合作执行纵向联邦学习任务的N个节点中的任意节点,其中N为大于1的整数,所述方法包括:

[0026] 对本地预测数据进行预处理,得到本地预测特征数据;

[0027] 基于其它节点的预测特征数据所包含的特征数量和特征名,对所述本地预测特征数据进行增广处理;

[0028] 对增广后的所述本地预测特征数据进行碎片化处理,以生成特征数据碎片,并按照多方安全计算协议,将其中N-1份所述特征数据碎片与所述其它节点共享,同时接收所述其它节点共享的特征数据碎片;

[0029] 根据本地生成的特征数据碎片和接收到的特征数据碎片,生成目标预测数据碎片;

[0030] 与所述其它节点一起,利用前述模型训练方法训练得到的与各所述节点对应的所述目标模型参数碎片,基于多方安全计算协议对各自生成的所述目标预测数据碎片进行处理,得到与每个所述节点对应的处理结果,并将所述处理结果共享给结果拥有方,以供所述结果拥有方根据各所述处理结果得到目标结果。

[0031] 第三方面,本发明提供一种基于MPC的联邦学习模型训练系统,所述系统适用于合作执行纵向联邦学习任务的N个节点,其中N为大于1的整数,所述系统包括:

- [0032] 第一预处理模块,用于分别对本地数据集合进行预处理,得到本地特征数据集;
- [0033] 第一增广模块,用于基于其它节点的特征数据集所包含的特征数量和特征名,对所述本地特征数据集进行增广处理;
- [0034] 第一碎片化分发模块,用于对增广后的所述本地特征数据集进行碎片化处理,以生成数据集碎片,并按照多方安全计算协议,将其中N-1份所述数据集碎片与所述其它节点共享,同时接收所述其它节点共享的数据集碎片;
- [0035] 第一融合模块,用于根据本地生成的数据集碎片和接收到的数据集碎片,生成训练样本集碎片;
- [0036] 训练模块,用于与所述其它节点一起,利用各自生成的训练样本集碎片,基于多方安全计算协议对联邦学习模型的模型参数进行训练,得到与所述N个节点对应的目标模型参数碎片。
- [0037] 优选地,所述第一预处理模块根据发起联邦学习任务的节点所配置的预处理规则,对所述本地数据集合进行预处理,得到所述本地特征数据集。
- [0038] 优选地,所述预处理规则包括缺失值处理规则、异常值处理规则、特征转换规则、和/或标准化处理规则。
- [0039] 优选地,所述第一融合模块对本地生成的数据集碎片和接收到的数据集碎片进行拼接和聚合处理,以生成所述训练样本集碎片。
- [0040] 优选地,所述训练模块用于与所述其它节点一起执行以下步骤:
- [0041] 获取各自对应的模型参数碎片的初始值,所述初始值通过对所述联邦学习模型的初始模型参数进行碎片化处理得到;
- [0042] 基于多方安全计算协议,利用各自对应的所述模型参数碎片的初始值对各自生成的所述训练样本集碎片进行处理,得到与每个所述节点对应的预测结果碎片;
- [0043] 基于多方安全计算协议,根据各自对应的所述预测结果碎片与相应的标准标签碎片,计算与每个所述节点对应的损失值碎片;
- [0044] 基于多方安全计算协议,根据各自计算得到的所述损失值碎片对各自对应的所述模型参数碎片进行训练,直至满足预设的训练终止条件,得到与每个所述节点对应的所述目标模型参数碎片。
- [0045] 第四方面,本发明提供一种基于MPC的联邦学习模型预测方系统,适用于合作执行纵向联邦学习任务的N个节点中的任意节点,其中N为大于1的整数,所述系统包括:
- [0046] 第二预处理模块,用于对本地预测数据进行预处理,得到本地预测特征数据;
- [0047] 第二增广模块,用于基于其它节点的预测特征数据所包含的特征数量和特征名,对所述本地预测特征数据进行增广处理;
- [0048] 第二碎片化分发模块,用于对所述本地预测特征数据进行碎片化处理,以生成特征数据碎片,并按照多方安全计算协议,将其中N-1份所述特征数据碎片与所述其它节点共享,同时接收所述其它节点共享的特征数据碎片;
- [0049] 第二融合模块,用于根据本地生成的特征数据碎片和接收到的特征数据碎片,生成目标预测数据碎片;
- [0050] 预测模块,用于与所述其它节点一起,利用前述模型训练系统训练得到的与各所述节点对应的所述目标模型参数碎片,基于多方安全计算协议对各自生成的所述目标预测

数据碎片进行处理,得到与每个所述节点对应的处理结果,并将所述处理结果共享给结果拥有方,以供所述结果拥有方根据各所述处理结果得到目标结果。

[0051] 第五方面,本发明提供一种电子设备,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程序,所述处理器执行所述计算机程序时实现上述任一项所述的方法。

[0052] 第六方面,本发明提供一种计算机可读存储介质,其上存储有计算机程序,其所述程序被处理器执行时实现上述任一项所述的方法。

[0053] 通过采用上述技术方案,本发明具有如下有益效果:

[0054] 本发明中的任意节点通过首先对本地数据集合进行预处理,得到本地特征数据集;而后基于其它节点的特征数据集所包含的特征数量和特征名,对所述本地特征数据集进行增广处理;再而后对增广后的所述本地特征数据集进行碎片化处理,以生成数据集碎片,并按照多方安全计算协议,将其中N-1份所述数据集碎片与所述其它节点共享,同时接收所述其它节点共享的数据集碎片;再而后根据本地生成的数据集碎片和接收到的数据集碎片,生成训练样本集碎片;最后与所述其它节点一起,利用各自生成的训练样本集碎片,基于多方安全计算协议对联邦学习模型的模型参数进行训练,得到与每个所述节点对应的目标模型参数碎片。可见,本发明中的训练方法只需要在各节点进行计算,无需发送中间结果到相关参与方,能够实现全流程安全保护。

## 附图说明

[0055] 图1为本发明实施例1的基于MPC的联邦学习模型训练方法的应用场景示意图;

[0056] 图2为本发明实施例1的基于MPC的联邦学习模型训练方法的流程图;

[0057] 图3为本发明实施例1中步骤S1的流程图;

[0058] 图4为本发明实施例1中步骤S12-S14的整体流程图;

[0059] 图5为本发明实施例1中发起方进行碎片化分发处理的流程图;

[0060] 图6为本发明实施例1中参与方进行碎片化分发处理的流程图;

[0061] 图7为本发明实施例1中原始数据状态到数据碎片化状态的流程图;

[0062] 图8为本发明实施例1中步骤S15的流程图;

[0063] 图9为本发明实施例2的基于MPC的联邦学习模型预测方法的流程图;

[0064] 图10为本发明实施例2的基于MPC的联邦学习模型预测方法的原理图;

[0065] 图11为本发明实施例3的基于MPC的联邦学习模型训练系统的结构框图;

[0066] 图12为本发明实施例4的基于MPC的联邦学习模型预测系统的结构框图;

[0067] 图13为本发明实施例5中电子设备的硬件架构图。

## 具体实施方式

[0068] 为了使本发明的目的、技术方案及优点更加清楚明白,以下结合附图及实施例,对本发明进行进一步详细说明。应当理解,此处所描述的具体实施例仅用以解释本发明,并不用于限定本发明。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。

[0069] 在本发明使用的术语是仅仅出于描述特定实施例的目的,而非旨在限制本公开。



在本公开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,本文中使用的术语“和/或”是指并包含一个或多个相关联的列出项目的任意或所有可能组合。

[0070] 多方安全计算(MPC)是指在无可信第三方情况下,通过多方共同参与,安全地完成某种协同计算。即在一个分布式的网络中,每个参与者都各自持有秘密输入,希望共同完成对某个函数的计算,但要求每个参与者除计算结果外均不能得到其他参与实体的任何输入信息。也就是说多方安全计算技术可以实现数据的可用不可见,获取数据使用价值,同时不泄露原始数据内容,实现数据安全与隐私保护。

[0071] 纵向联邦学习是指在参与者的数据特征重叠较小,而用户重叠较多的情况下,取出参与者用户相同而用户数据特征不同的那部分用户及特征数据进行联合机器学习训练。例如,当有两个参与方时,纵向联邦学习的场景如图1所示。

[0072] 实施例1

[0073] 本实施例提供一种基于MPC的联邦学习模型训练方法,适用于合作执行纵向联邦学习任务的N个节点中的任意节点,其中N为大于1的整数。

[0074] 如图2所示,本实施例的方法具体包括以下步骤:

[0075] S11,对本地数据集合进行预处理,得到本地特征数据集。

[0076] 下面结合图3,以两个节点合作执行纵向联邦学习任务为例说明本步骤的具体实现过程:

[0077] 首先,发起方(即发起联邦学习任务的节点)向参与方发送联邦学习合作任务请求,并选择发起方侧的数据集合D1作为发起方侧的数据集合,选择参与方侧的数据集合D2作为参与方侧的数据集合,在该过程中,可以选择对应的数据集特征。其中,在图3所示的示例中,以带标签Y的数据方作为发起方(当然也可以将不带标签的数据方作为发起方),数据集D1包含特征列Xa及标签列Y,数据集D2包含特征列Xb。若发起方A与参与方B建立了合作关系,则A、B能相互查看对方的数据集摘要列表:包括数据集名称、特征名、字段类型、特征数量等概要信息。

[0078] 而后,参与方审核联邦学习合作任务请求,并将审核结果(通过或者失败)发送给发起方。如果审核结果为通过,则发起方与参与方建立联邦学习合作关系,否则任务终止。

[0079] 当发起方与参与方建立联邦学习合作关系后,发起方配置统一的预处理规则并发送给参与方,该预处理规则包含缺失值处理规则、异常值处理规则、特征转换规则、和/或标准化处理规则等。其中,缺失值处理规则包括设置缺失值填补默认值或者使用众数等填补策略;异常值处理规则包括删除特征值超出正常范围以外的样本;特征转换规则涉及采用数值转换操作,如采用平方或者对数等非线性转换动作,或者设置文本转数值规则等;标准化处理规则包括Z-score或者Min-Max等标准化处理规则。

[0080] 最后,发起方与参与方分别根据前述预处理规则,执行相应的特征预处理,分别得到数值化的特征数据集D1'和D2'。

[0081] S12,所述任意节点基于其它节点的特征数据集所包含的特征数量和特征名,对所述本地特征数据集进行增广处理。下面以图3得到的D1'和D2'为例,结合图4详细介绍本实施例的增广处理过程:

[0082] 假设D1'具有三个特征,分别为Xa1,Xa2,Xa3;D2'具有两个特征,分别为Xb1,Xb2,

则发起方将对其所对应的本地特征数据集D1'中的特征进行增广处理,具体通过采用特征值0补足D1'不具备的特征Xb1、Xb2;同理,参与方将按照相同规则对其所对应的本地特征数据集D2'中的特征进行增广处理,即通过采用特征值0补足D2'不具备的特征Xa1、Xa2、Xa3。

[0083] 例如,假设D1'包括4个样本(一个样本表示一个用户),D2'包括3个样本,且原始D1'如表1所示,原始D2'如表2所示,则对D1'增广后的D1''如表3所示,对D2'增广后的D2''如表4所示。

[0084] 表1

[0085]

uid	Xa1	Xa2	Xa3
124360	1.3	5.2	3
328492	2.5	3.3	-2
572683	-1	0.5	0.2
930913	0.9	0.12	1

[0086] 表2

[0087]

uid	Xb1	Xb2
748329	0.89	1.41
328492	2.3	1.9
930913	-1.2	-0.1

[0088] 表3

[0089]

uid	Xa1	Xa2	Xa3	Xb1	Xb2
124360	1.3	5.2	3	0	0
328492	2.5	3.3	-2	0	0
572683	-1	0.5	0.2	0	0
930913	0.9	0.12	1	0	0

[0091] 表4

[0092]

uid	Xa1	Xa2	Xa3	Xb1	Xb2
748329	0	0	0	0.89	1.41
328492	0	0	0	2.3	1.9
930913	0	0	0	-1.2	-0.1

[0093] 根据表3和表4可以看出,增广处理后得到的D1''和D2''具有同样多的特征数,特征数都为5。

[0094] S13,对增广后的所述本地特征数据集进行碎片化处理,以生成数据集碎片,并将其中N-1份所述数据集碎片与所述其它节点共享,同时接收所述其它节点共享的数据集碎片。

[0095] 具体地,首先基于MPC协议算子,在发起方侧对D1''样本矩阵、在参与方侧对D2''样本矩阵进行碎片化处理,碎片化之后得到的数据集碎片D1''和D2''分别如表5和表6所示。

[0096] 表5

[0097]

uid	Xa1	Xa2	Xa3	Xb1	Xb2
[124360]	[1.3]	[5.2]	[3]	[0]	[0]

[328492]	[2.5]	[3.3]	[-2]	[0]	[0]
[572683]	[-1]	[0.5]	[0.2]	[0]	[0]
[930913]	[0.9]	[0.12]	[1]	[0]	[0]

[0098] 表6

[0099]

uid	Xa1	Xa2	Xa3	Xb1	Xb2
[748329]	[0]	[0]	[0]	[0.89]	[1.41]
[328492]	[0]	[0]	[0]	[2.3]	[1.9]
[930913]	[0]	[0]	[0]	[-1.2]	[-0.1]

[0100] 在表5和表6中，“[ ]”表示的是碎片状态的数值，比如[1.3]其实表示的是两份碎片，如2.7和-1.4，发起方与参与方将各持有一份。在本实施例中，uid(即用户id)必须数值化，如果是字符串可以先hash为数值型，再进行碎片化，比如uid[124360]表示的是两份碎片，如235378和-111018，发起方与参与方将各持有一份碎片。

[0101] 图5和图6以两方为例，展示了发起方侧和参与方侧的数据碎片化方式。其中，Ra和Rb分别是在发起方侧和参与方侧各自生成的随机数，随机数生成的规则是在指定范围内采用uniformrandom的方式随机生成一个值，指定范围视具体的精度选择确定，64bit精度下范围是 $[-2^{63}, 2^{63}]$ ，128bit精度下为 $[-2^{127}, 2^{127}]$ ，甚至扩展到256bit、512bit等范围，精度越高，随机数范围越大。

[0102] 通过图5和图6展示的算术运算，可以将原始的Xa拆成两份碎片，如图7所示，Xa-Ra由发起方自己保留，碎片Ra则共享给参与方。同理，参与方将Xb拆分成两份碎片，Xb-Rb共享给发起方，Rb由参与方自己保留。经过碎片化处理分发，可以看到发起方和参与方都持有两份碎片，但无法推出对方原始的值。同理，Y标签数据在发起方端做同样的碎片化处理，得到碎片状态的标签数据Y’。

[0103] S14，根据本地生成的数据集碎片和接收到的数据集碎片，生成训练样本集碎片。本步骤的具体实现过程如下：

[0104] S141，分别将本地生成的数据集碎片和接收到的数据集碎片进行concat(拼接)操作。

[0105] 具体地，再次参阅图4，各节点在对增广后的本地特征数据集进行碎片化分发处理之后，通过MPC协议的concat隐私算子，在发起方侧和参与方侧分别对碎片化的D1”和碎片化的D2”分别各自拼接在一起，得到碎片化状态下的所有样本拼接矩阵，这样将得到如表7的示的样本量为7的样本矩阵D12”。

[0106] 表7

[0107]

uid	Xa1	Xa2	Xa3	Xb1	Xb2
[124360]	[1.3]	[5.2]	[3]	[0]	[0]
[328492]	[2.5]	[3.3]	[-2]	[0]	[0]
[572683]	[-1]	[0.5]	[0.2]	[0]	[0]
[930913]	[0.9]	[0.12]	[1]	[0]	[0]
[748329]	[0]	[0]	[0]	[0.89]	[1.41]
[328492]	[0]	[0]	[0]	[2.3]	[1.9]
[930913]	[0]	[0]	[0]	[-1.2]	[-0.1]

[0108] S142,对拼接后的数据集碎片D12”中的uid列进行group(聚合)操作,得到所述训练样本集碎片。其中,聚合操作是指找出并保留uid存在重复的样本,即D1’与D2’存在的重叠的样本。

[0109] 在本实施例中,可以采用MPC的比较算子来实现该逻辑,比较uid是否相等,如果存在相等情形,或者[uid]的数量为2,则保留该样本,并对特征进行合并处理;否则丢弃该样本。这样,基于MPC协议,可以达到样本对齐的目标。

[0110] 本实施例对group的实现逻辑不做限定,可以有多种优化实现的方式,比如可以先通过对uid采取模糊方式做分桶,或者基于排序后做顺序比较等,在这里仅描述基于比较算子通过比较大小可以实现该逻辑,不做具体的限定说明。通过group操作后可以获得uid存在重叠的样本,保留该样本,并进行特征合并,特征合并是指将原先以0填充的特征改为使用本地真实的样本碎片化特征填充。对表7进行聚合处理后得到的训练样本集碎片将如表8所示:

[0111] 表8

uid	Xa1	Xa2	Xa3	Xb1	Xb2
[328492]	[2.5]	[3.3]	[-2]	[2.3]	[1.9]
[930913]	[0.9]	[0.12]	[1]	[-1.2]	[-0.1]

[0113] S15,与其它节点一起,利用各自生成的训练样本集碎片,基于多方安全计算协议对联邦学习模型的模型参数进行训练,得到与所述N个节点对应的目标模型参数碎片。

[0114] 在本实施例中,整体训练框架包括应用层和算子层,其中应用层处理模型的创建、参数设置、主流程逻辑的计算、pipeline的构建。应用层中涉及运算算子,如加法、减法、乘法、矩阵乘法等算子,使用MPC协议算子库提供的封装接口。通过调用接口,将具体的运算交给底层算子库中的函数执行,也就是说算子层为真正执行计算的主体,计算结束将结果发回应用层。这种架构类似于tensorflow(底层为C代码执行逻辑,上层为python封装接口),这样的架构,有助于降低编程难度,上手更迅速,对整体模型的改造也更小,而执行逻辑下沉到底层,应用层不需要关心具体的算子实现逻辑。算子层可以采用更高效的编程语言实现,如C、rust或者go。这种架构更适合工程化。

[0115] 下面结合图8,对步骤S15的具体实现过程进行详细描述:

[0116] S151,各节点分别获取各自对应的模型参数碎片的初始值,所述初始值由发起联邦学习任务的节点对所述联邦学习模型的初始模型参数进行碎片化处理得到。

[0117] 具体地,首先由发起方进行初始化,定义待训练的联邦学习模型,比如定义该模型为Logisticregression或者CNN、DNN等常用的机器学习模型。而后,发起方设置模型所需的超参数,如学习率、迭代次数、损失终止条件、梯度变化率、惩罚系数以及优化方法的选择。同时,发起方初始化生成模型需要学习的模型参数的初始值w(不同模型对生成的参数需求不同,此处以模型权重系数w为例),并对w进行碎片化处理,发起方保留w\_1碎片,将w\_2(即w-w\_1)碎片发送给参与方持有。

[0118] S152,各节点基于多方安全计算协议,利用各自对应的所述模型参数碎片的初始值对各自生成的所述训练样本集碎片进行处理,得到与每个所述节点对应的预测结果碎片。

[0119] 具体地,各节点基于各自对应的训练样本集碎片计算当前模型的预测结果,当然

该预测结果为碎片化形式。其中,发起方执行 $w_1$ 与训练样本集碎片 $D12_1$ 的计算,参与方执行 $w_2$ 与训练样本集碎片 $D12_2$ 的计算,比如LR模型下,发起方计算的是 $Y'_1_{tmp}=w_1 * X12_1 + b_1$ ,参与方计算的是 $Y'_2_{tmp}=w_2 * X12_2 + b_2$ 。然后再通过MPC的sigmoid算子进行计算得到 $Y'_1$ 与 $Y'_2$ 。如果是DNN,则使用MPC的relu算子。 $Y'_1$ 与 $Y'_2$ 是模型预测值 $Y'$ 的两份碎片。为了避免误解,说明一下,这里的 $Y'$ 不是单一个值,而是量级对应样本量的多个值,样本量多少, $Y'$ 就有多少,表示形式是矩阵类型。 $Y'_1$ 与 $Y'_2$ 分别由发起方与参与方各自持有,不共享。这一步也称为碎片化前向预测值计算。

[0120] S153,各节点基于多方安全计算协议,根据各自对应的所述预测结果碎片与相应的标准标签碎片,计算与每个所述节点对应的损失值碎片。

[0121] 具体地,有了碎片化状态的预测结果 $Y'$ ,就可以与碎片化的真实的标准标签 $Y$ ,计算碎片化的损失值LOSS。

[0122] S154,各节点基于多方安全计算协议,根据各自计算得到的所述损失值碎片对各自对应的所述模型参数碎片进行训练,直至满足预设的训练终止条件,得到与每个所述节点对应的所述目标模型参数碎片。

[0123] 在本实施例中,发起方执行碎片化的LOSS值与模型初始化定义的模型误差终止条件值进行比较,如果满足条件,则终止模型训练,反之则继续执行。底层mpc算子库支持密文碎片状态下的值与明文的比较,得到结果True或者False。另外,如果迭代次数达到初始化设定的次数,也会终止模型训练。

[0124] 在本实施例中,各节点基于多方安全计算协议,根据各自计算得到的所述损失值碎片对各自对应的所述模型参数碎片进行训练,包括:

[0125] 首先,各节点基于多方安全计算协议获取各自对应的梯度值碎片。具体地,梯度值的计算公式有多种,常见的有 $(Y-Y')$ 或者 $(Y-Y') * X$ 等,基于不同模型采取的具体梯度计算公式进行具体计算。一般用到的算子为减法、乘法或者矩阵乘法等。通过计算得到梯度值的碎片化状态结果。

[0126] 而后,各节点基于多方安全计算协议,根据各自对应的所述梯度值碎片对各自对应的所述模型参数碎片进行更新。具体地,可以基于梯度值碎片,结合学习率、惩罚系数等超参数,对模型参数碎片进行更新。更新的具体逻辑常见梯度下降法或者牛顿迭代法。一般涉及的MPC算子为乘法/矩阵乘法、加法、减法、开方等。

[0127] 各节点得到各自对应的所述目标模型参数碎片后,存储到各自的数据库进行持久化,比如hdfs、mysql、fastdfs等,并通过模型id进行标识。

[0128] 需要说明的是,本实施例的模型训练可以不局限在LR、梯度决策树、CNN这些模型,其他模型同样可以扩展,实践发现基于MPC的实现,对原始明文的模型改动代价较小。

[0129] 在本实施例中,合执行联邦学习任务的任意节点通过首先对本地数据集进行预处理,得到本地特征数据集;而后基于其它节点的特征数据集所包含的特征数量和特征名,对所述本地特征数据集进行增广处理;再而后对增广后的所述本地特征数据集进行碎片化处理,以生成数据集碎片,并按照多方安全计算协议,将其中 $N-1$ 份所述数据集碎片与所述其它节点共享,同时接收所述其它节点共享的数据集碎片;再而后根据本地生成的数据集碎片和接收到的数据集碎片,生成训练样本集碎片;最后与所述其它节点一起,利用各自生成的训练样本集碎片,基于多方安全计算协议对联邦学习模型的模型参数进行训练,得到

与每个所述节点对应的目标模型参数碎片。可见,本发明中的训练方法只需要在各节点进行计算,无需发送中间结果到相关参与方,能够实现全流程安全保护。

#### [0130] 实施例2

[0131] 本实施例提供一种基于MPC的联邦学习模型预测方法,适用于合作执行纵向联邦学习任务的N个节点中的任意节点,其中N为大于1的整数,所述方法基于实施例1的训练方法实现,如图9-图10所示,该方法具体包括以下步骤:

[0132] S21,对本地预测数据(包括测试数据)进行预处理,得到本地预测特征数据。

[0133] S22,基于其它节点的预测特征数据所包含的特征数量和特征名,对所述本地预测特征数据进行增广处理。

[0134] S23,对增广后的所述本地预测特征数据进行碎片化处理,以生成特征数据碎片,并按照多方安全计算协议,将其中N-1份所述特征数据碎片与所述其它节点共享,同时接收所述其它节点共享的特征数据碎片。

[0135] S24,根据本地生成的特征数据碎片和接收到的特征数据碎片,生成目标预测数据碎片。

[0136] 其中,步骤S21-S24的具体实现过程可对应参考步骤S11-S14,其中,步骤S21采用的预处理规则需要与训练阶段保持一致,也就是说,两者预处理的config配置文件需要使用同一份。在训练阶段进行预处理时需要保存该配置信息,在执行预测的时候,加载该预处理规则进行处理。

[0137] S25,与所述其它节点一起,利用前述模型训练方法训练得到的与各所述节点对应的所述目标模型参数碎片,基于多方安全计算协议对各自生成的所述目标预测数据碎片进行处理,得到与每个所述节点对应的处理结果,并将所述处理结果共享给结果拥有方,以供所述结果拥有方根据各所述处理结果得到目标结果。

[0138] 具体地,首先,发起方侧下发需要加载的模型id(具体获取的逻辑由具体产品设计设定,也可以通过发起方从模型列表选取),发起方向参与方发起协同预测请求,并将该模型id发送给参与方。发起方与参与方各自根据模型id,从持久化数据库中加载该对应模型的参数碎片,分别得到模型参数碎片 $w_1$ 和 $w_2$ 。

[0139] 而后,发起方与参与方利用各自生成的目标预测数据碎片 $D_{12\_1}$ 和 $D_{12\_2}$ 执行碎片化前向预测值计算,计算方式同训练阶段的前向预测值计算逻辑。发起方侧计算得到预测值碎片 $Y'_1$ ,参与方获得预测值碎片 $Y'_2$ 。

[0140] 最后,只有结果拥有方才有权利拿到明文结果。在本实施例中,假设发起方为结果拥有方,那么其将会向参与方发起请求,获取模型处理结果碎片 $Y'_2$ ,参与方同意后将该碎片发送给结果拥有方,结果拥有方执行明文reveal解密,计算 $Y' = Y'_1 + Y'_2$ 获得最终的明文状态的目标结果。

[0141] 可见,本实施例的预测方法在数据不出门及全流程碎片密文的安全保护限度下,能够实现高效的隐私模型预测。

#### [0142] 实施例3

[0143] 本实施例提供一种基于MPC的联邦学习模型训练系统,所述系统适用于合作执行纵向联邦学习任务的N个节点,其中N为大于1的整数,如图11所示,所述系统包括:

[0144] 第一预处理模块11,用于分别对本地数据集合进行预处理,得到本地特征数据集;

[0145] 第一增广模块12,用于基于其它节点的特征数据集所包含的特征数量和特征名,对所述本地特征数据集进行增广处理;

[0146] 第一碎片化分发模块13,用于对增广后的所述本地特征数据集进行碎片化处理,以生成数据集碎片,并按照多方安全计算协议,将其中N-1份所述数据集碎片与所述其它节点共享,同时接收所述其它节点共享的数据集碎片;

[0147] 第一融合模块14,用于根据本地生成的数据集碎片和接收到的数据集碎片,生成训练样本集碎片;

[0148] 训练模块15,用于与所述其它节点一起,利用各自生成的训练样本集碎片,基于多方安全计算协议对联邦学习模型的模型参数进行训练,得到与所述N个节点对应的目标模型参数碎片。

[0149] 优选地,所述第一预处理模块根据发起联邦学习任务的节点所配置的预处理规则,对所述本地数据集进行预处理,得到所述本地特征数据集。

[0150] 优选地,所述预处理规则包括缺失值处理规则、异常值处理规则、特征转换规则、和/或标准化处理规则。

[0151] 优选地,所述第一融合模块对本地生成的数据集碎片和接收到的数据集碎片进行拼接和聚合处理,以生成所述训练样本集碎片。

[0152] 优选地,所述训练模块用于与所述其它节点一起执行以下步骤:

[0153] 获取各自对应的模型参数碎片的初始值,所述初始值通过对所述联邦学习模型的初始模型参数进行碎片化处理得到;

[0154] 基于多方安全计算协议,利用各自对应的所述模型参数碎片的初始值对各自生成的所述训练样本集碎片进行处理,得到与每个所述节点对应的预测结果碎片;

[0155] 基于多方安全计算协议,根据各自对应的所述预测结果碎片与相应的标准标签碎片,计算与每个所述节点对应的损失值碎片;

[0156] 基于多方安全计算协议,根据各自计算得到的所述损失值碎片对各自对应的所述模型参数碎片进行训练,直至满足预设的训练终止条件,得到与每个所述节点对应的所述目标模型参数碎片。

[0157] 本实施例的训练系统在数据不出门及全流程碎片密文的安全保护限度下,能够实现高效的隐私模型训练。

[0158] 实施例4

[0159] 本实施例提供一种基于MPC的联邦学习模型预测方系统,适用于合作执行纵向联邦学习任务的N个节点中的任意节点,其中N为大于1的整数,如图12所示,所述系统包括:

[0160] 第二预处理模块21,用于对本地预测数据进行预处理,得到本地预测特征数据;

[0161] 第二增广模块22,用于基于其它节点的预测特征数据所包含的特征数量和特征名,对所述本地预测特征数据进行增广处理;

[0162] 第二碎片化分发模块23,用于对所述本地预测特征数据进行碎片化处理,以生成特征数据碎片,并按照多方安全计算协议,将其中N-1份所述特征数据碎片与所述其它节点共享,同时接收所述其它节点共享的特征数据碎片;

[0163] 第二融合模块24,用于根据本地生成的特征数据碎片和接收到的特征数据碎片,生成目标预测数据碎片;

[0164] 预测模块25,用于与所述其它节点一起,利用前述模型训练系统训练得到的与各所述节点对应的所述目标模型参数碎片,基于多方安全计算协议对各自生成的所述目标预测数据碎片进行处理,得到与每个所述节点对应的处理结果,并将所述处理结果共享给结果拥有方,以供所述结果拥有方根据各所述处理结果得到目标结果。

[0165] 本实施例的预测系统在数据不出门及全流程碎片密文的安全保护限度下,能够实现高效的隐私模型预测。

[0166] 对于本系统实施例而言,由于其基本对应于方法实施例,所以相关之处参见方法实施例的部分说明即可。以上所描述的系统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发明方案的目的。本领域普通技术人员在不付出创造性劳动的情况下,即可以理解并实施。

[0167] 实施例5

[0168] 本实施例提供一种电子设备,电子设备可以通过计算设备的形式表现(例如可以为服务器设备),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程序,其中处理器执行计算机程序时可以实现实施例1或2提供的方法。

[0169] 图13示出了本实施例的硬件结构示意图,如图13所示,电子设备9具体包括:

[0170] 至少一个处理器91、至少一个存储器92以及用于连接不同系统组件(包括处理器91和存储器92)的总线93,其中:

[0171] 总线93包括数据总线、地址总线和控制总线。

[0172] 存储器92包括易失性存储器,例如随机存取存储器(RAM)921和/或高速缓存存储器922,还可以进一步包括只读存储器(ROM)923。

[0173] 存储器92还包括具有一组(至少一个)程序模块924的程序/实用工具925,这样的程序模块924包括但不限于:操作系统、一个或者多个应用程序、其它程序模块以及程序数据,这些示例中的每一个或某种组合中可能包括网络环境的实现。

[0174] 处理器91通过运行存储在存储器92中的计算机程序,从而执行各种功能应用以及数据处理,例如本发明实施例1提供的方法。

[0175] 电子设备9进一步可以与一个或多个外部设备94(例如键盘、指向设备等)通信。这种通信可以通过输入/输出(I/O)接口95进行。并且,电子设备9还可以通过网络适配器96与一个或者多个网络(例如局域网(LAN),广域网(WAN)和/或公共网络,例如因特网)通信。网络适配器96通过总线93与电子设备9的其它模块通信。应当明白,尽管图中未示出,可以结合电子设备9使用其它硬件和/或软件模块,包括但不限于:微代码、设备驱动器、冗余处理器、外部磁盘驱动阵列、RAID(磁盘阵列)系统、磁带驱动器以及数据备份存储系统等。

[0176] 应当注意,尽管在上文详细描述中提及了电子设备的若干单元/模块或子单元/模块,但是这种划分仅仅是示例性的并非强制性的。实际上,根据本申请的实施方式,上文描述的两个或更多单元/模块的特征和功能可以在一个单元/模块中具体化。反之,上文描述的一个单元/模块的特征和功能可以进一步划分为由多个单元/模块来具体化。

[0177] 实施例6

[0178] 本实施例提供了一种计算机可读存储介质,其上存储有计算机程序,所述程序被处理器执行时实现实施例1的方法的步骤。



[0179] 其中,可读存储介质可以采用的更具体可以包括但不限于:便携式盘、硬盘、随机存取存储器、只读存储器、可擦拭可编程只读存储器、光存储器件、磁存储器件或上述的任意合适的组合。

[0180] 在可能的实施方式中,本发明还可以实现为一种程序产品的形式,其包括程序代码,当所述程序产品在终端设备上运行时,所述程序代码用于使所述终端设备执行实现实施例1的方法的步骤。

[0181] 其中,可以以一种或多种程序设计语言的任意组合来编写用于执行本发明的程序代码,所述程序代码可以完全地在用户设备上执行、部分地在用户设备上执行、作为一个独立的软件包执行、部分在用户设备上部分在远程设备上执行或完全在远程设备上执行。

[0182] 虽然以上描述了本发明的具体实施方式,但是本领域的技术人员应当理解,这仅是举例说明,本发明的保护范围是由所附权利要求书限定的。本领域的技术人员在不背离本发明的原理和实质的前提下,可以对这些实施方式做出多种变更或修改,但这些变更和修改均落入本发明的保护范围。

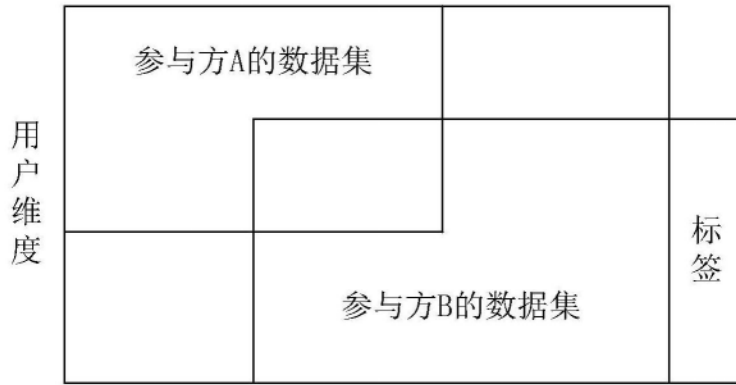


图1

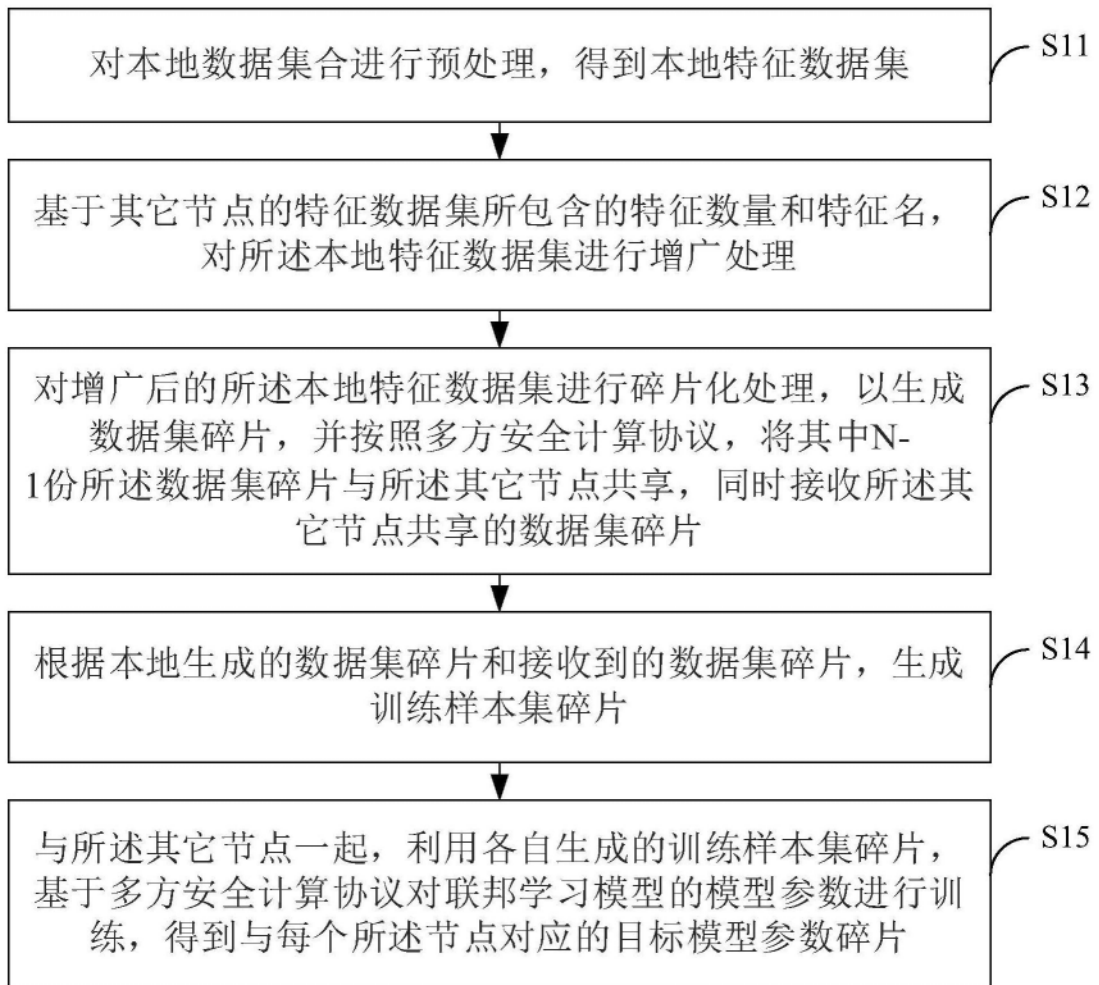


图2

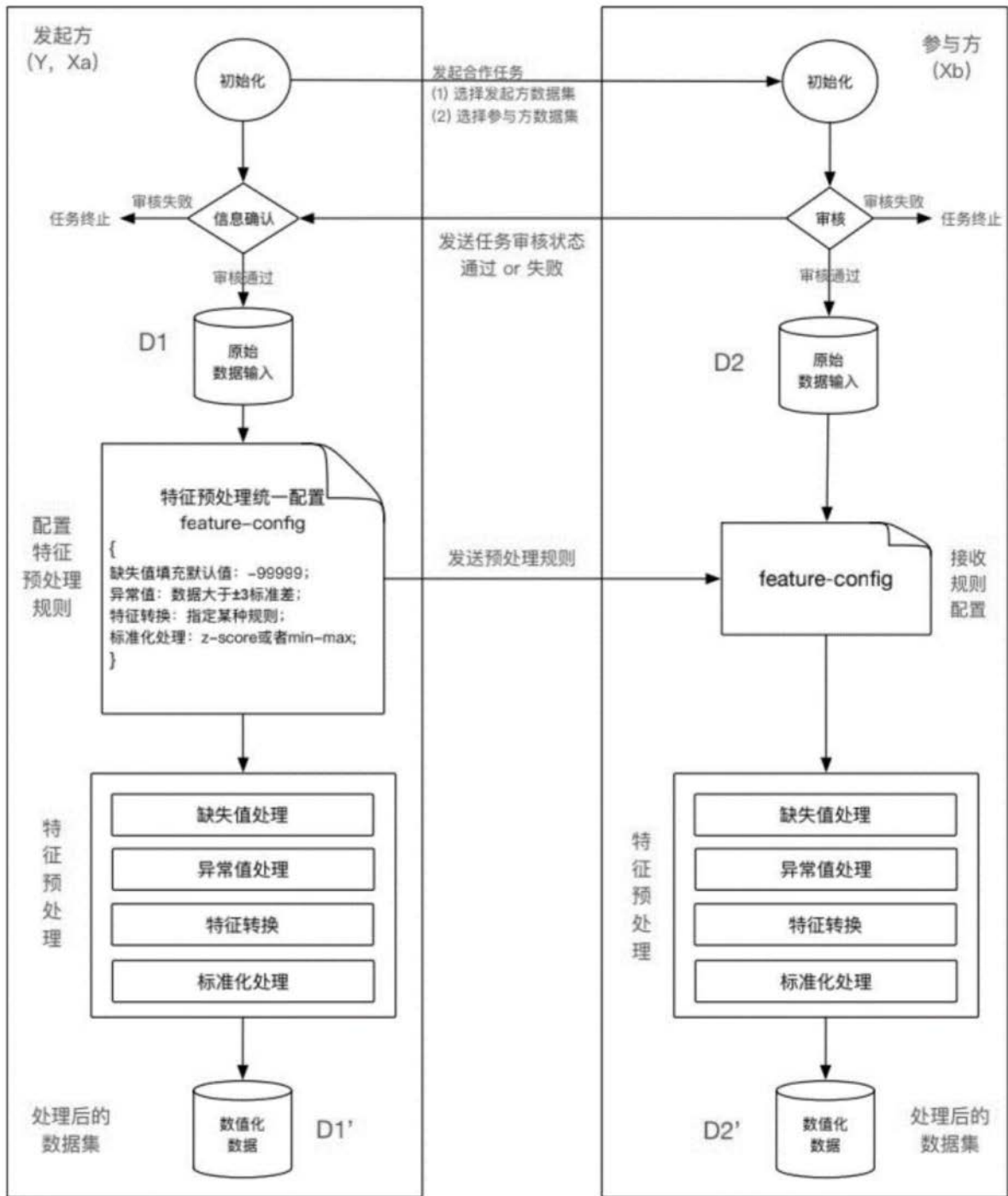


图3

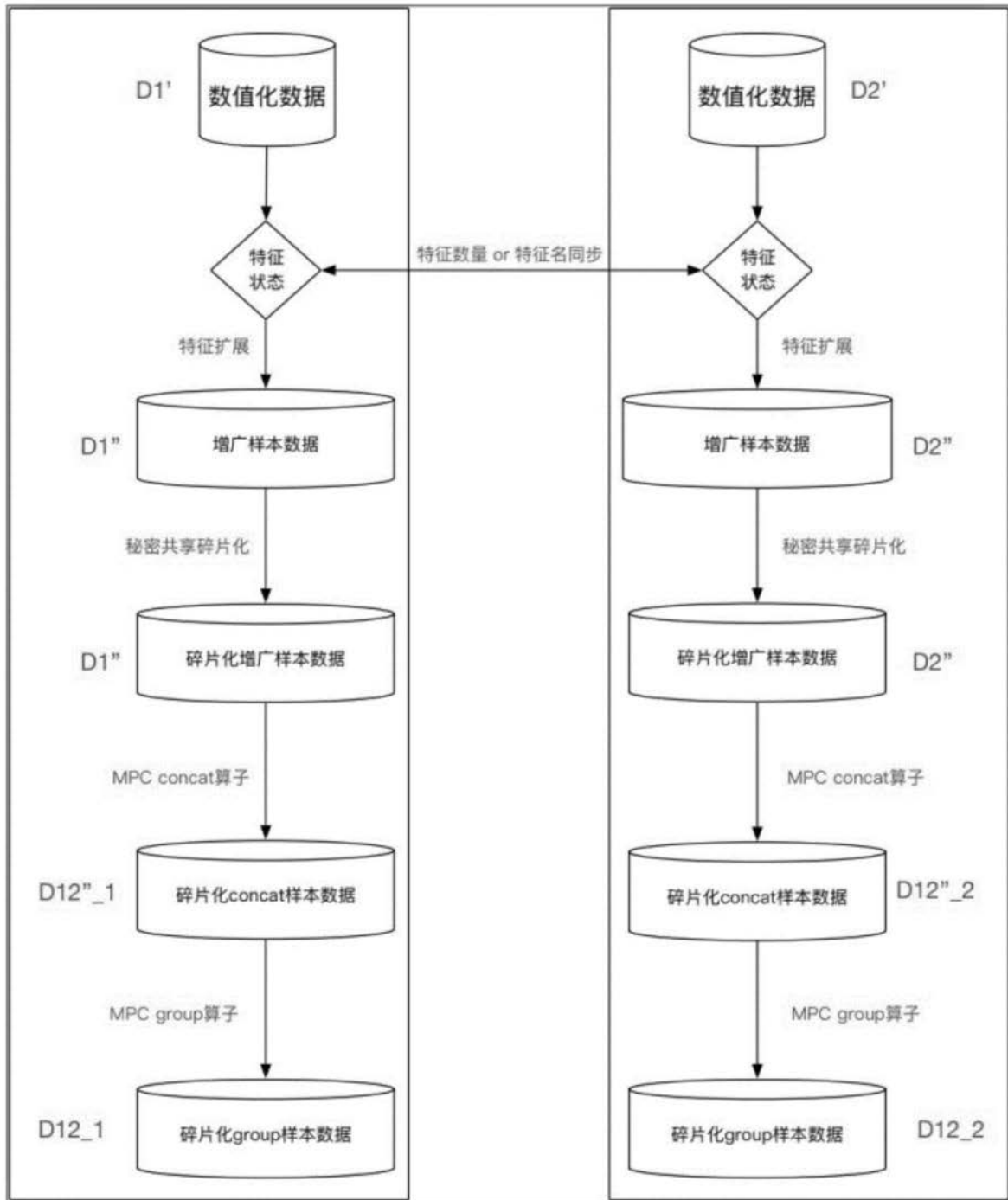


图4

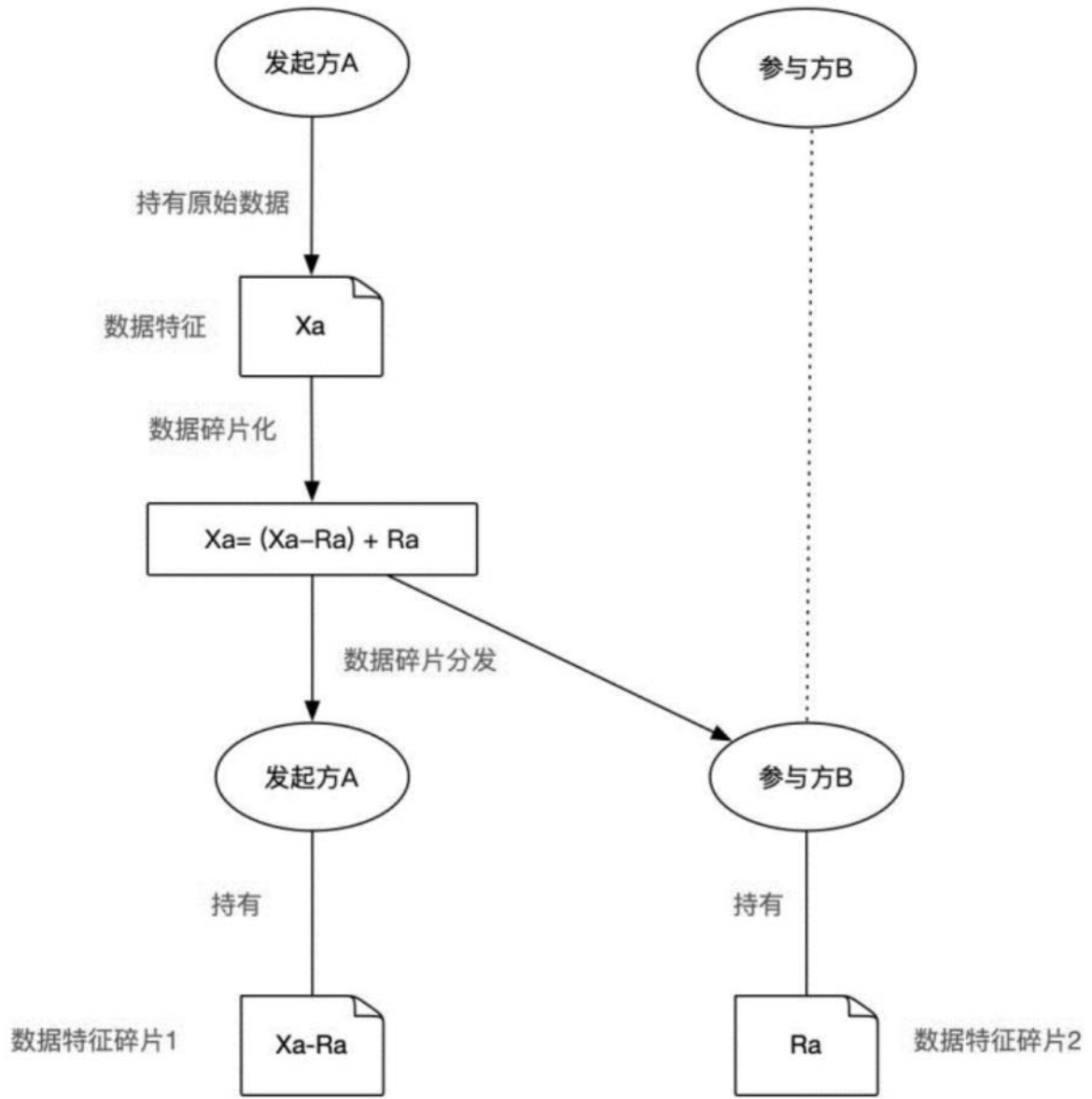


图5

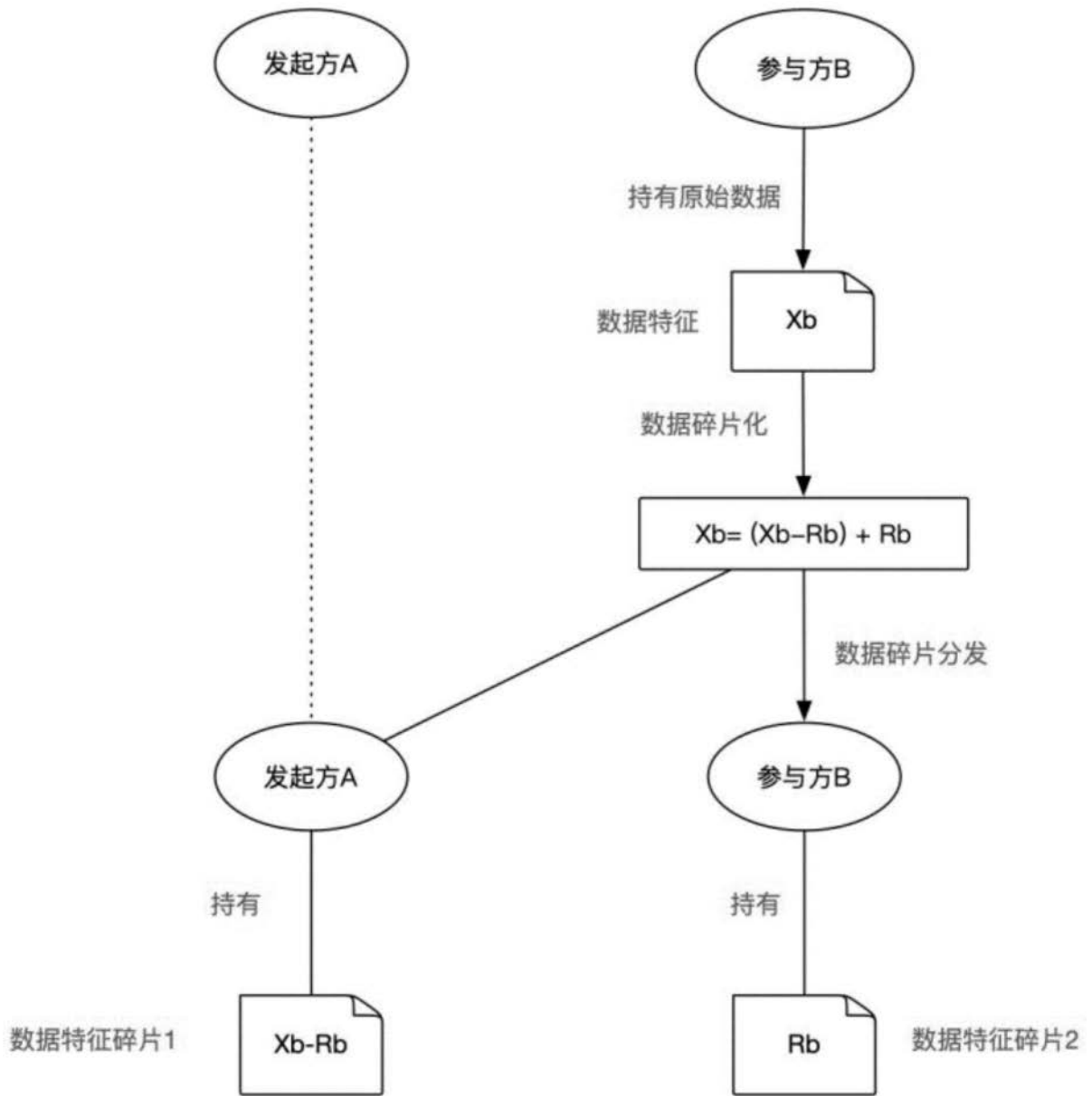


图6

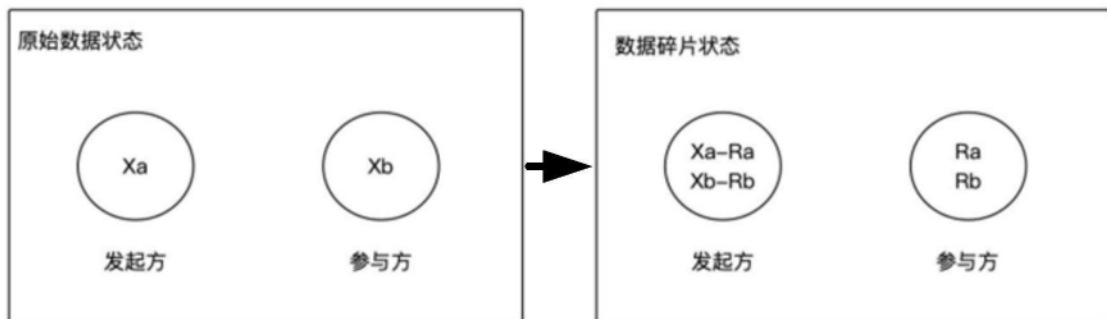


图7

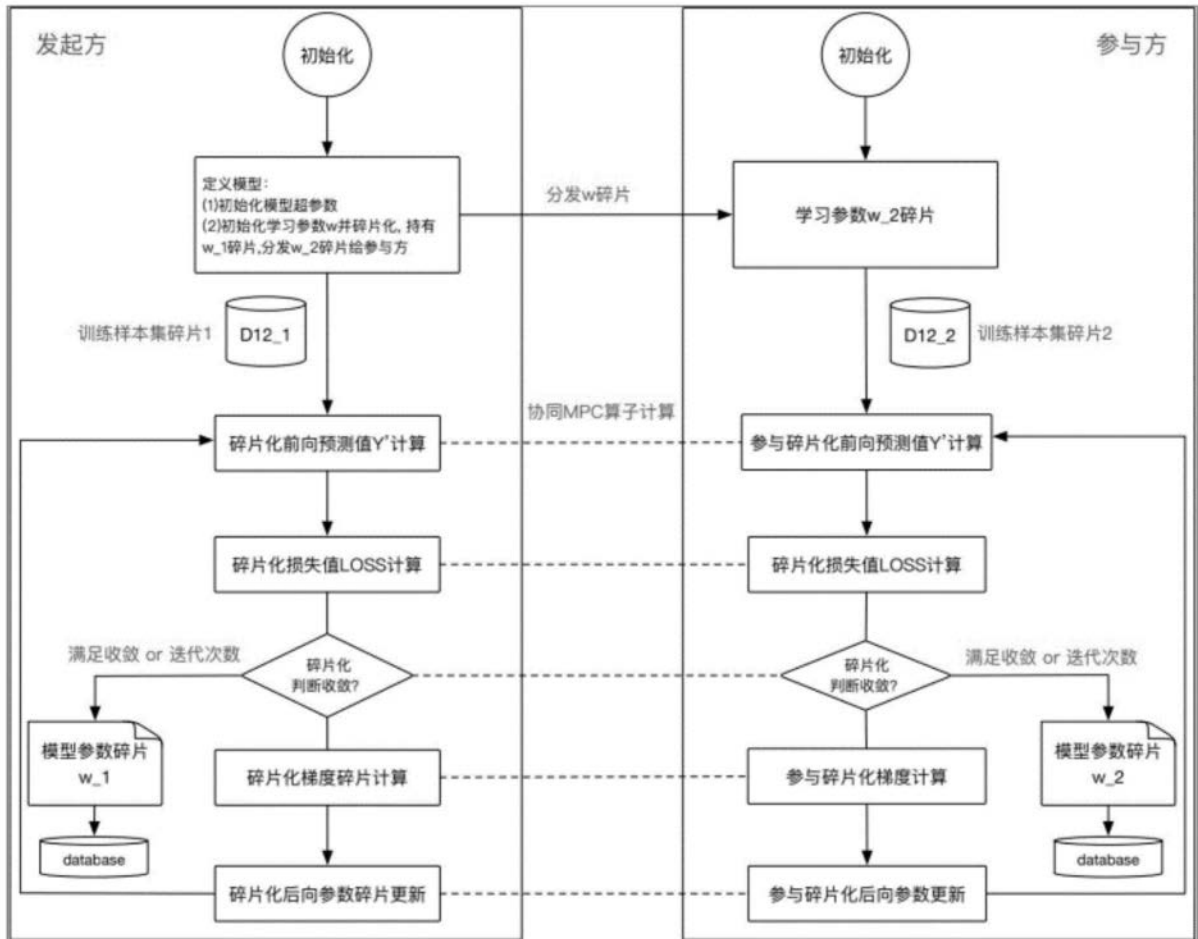


图8

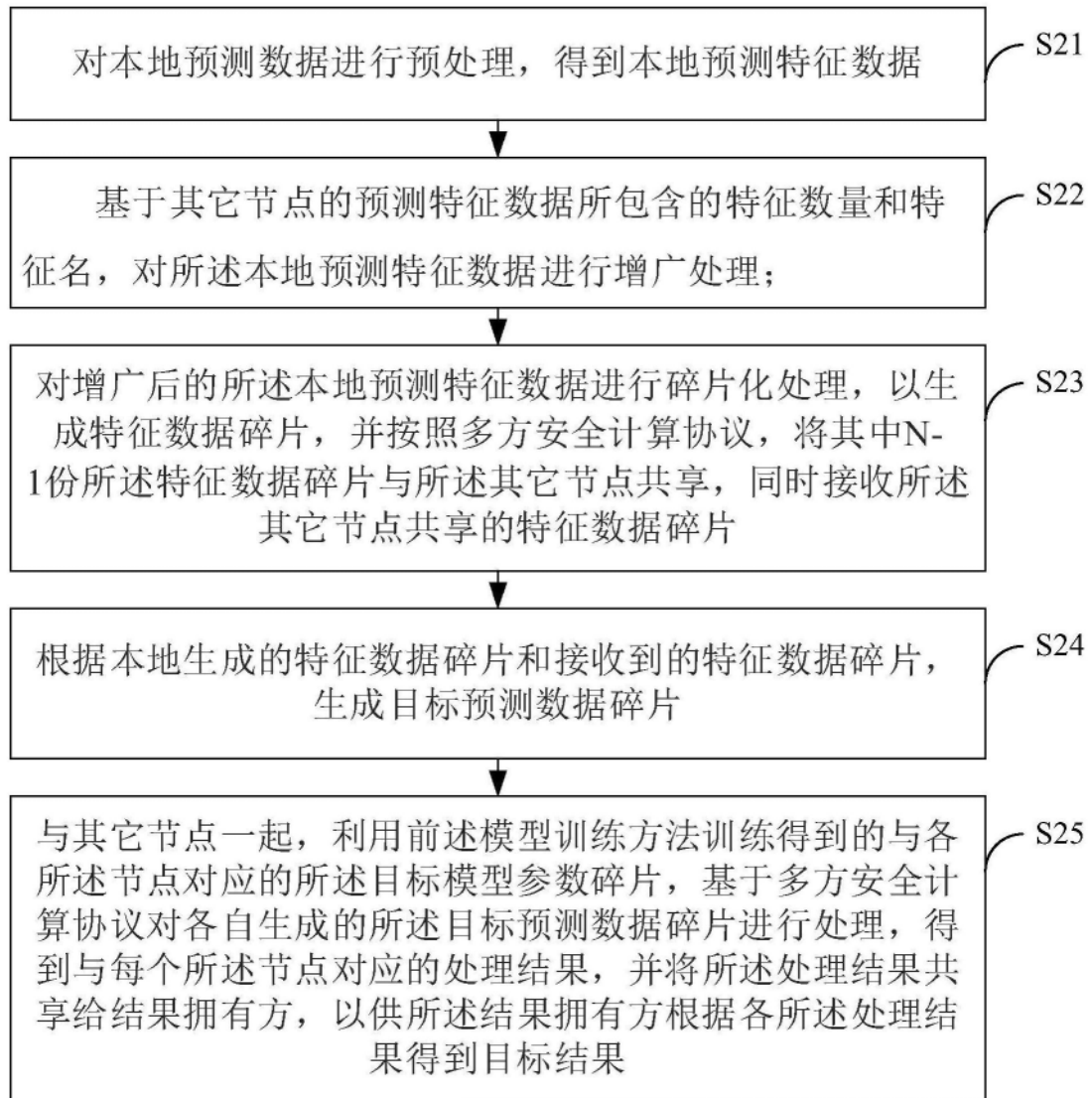


图9



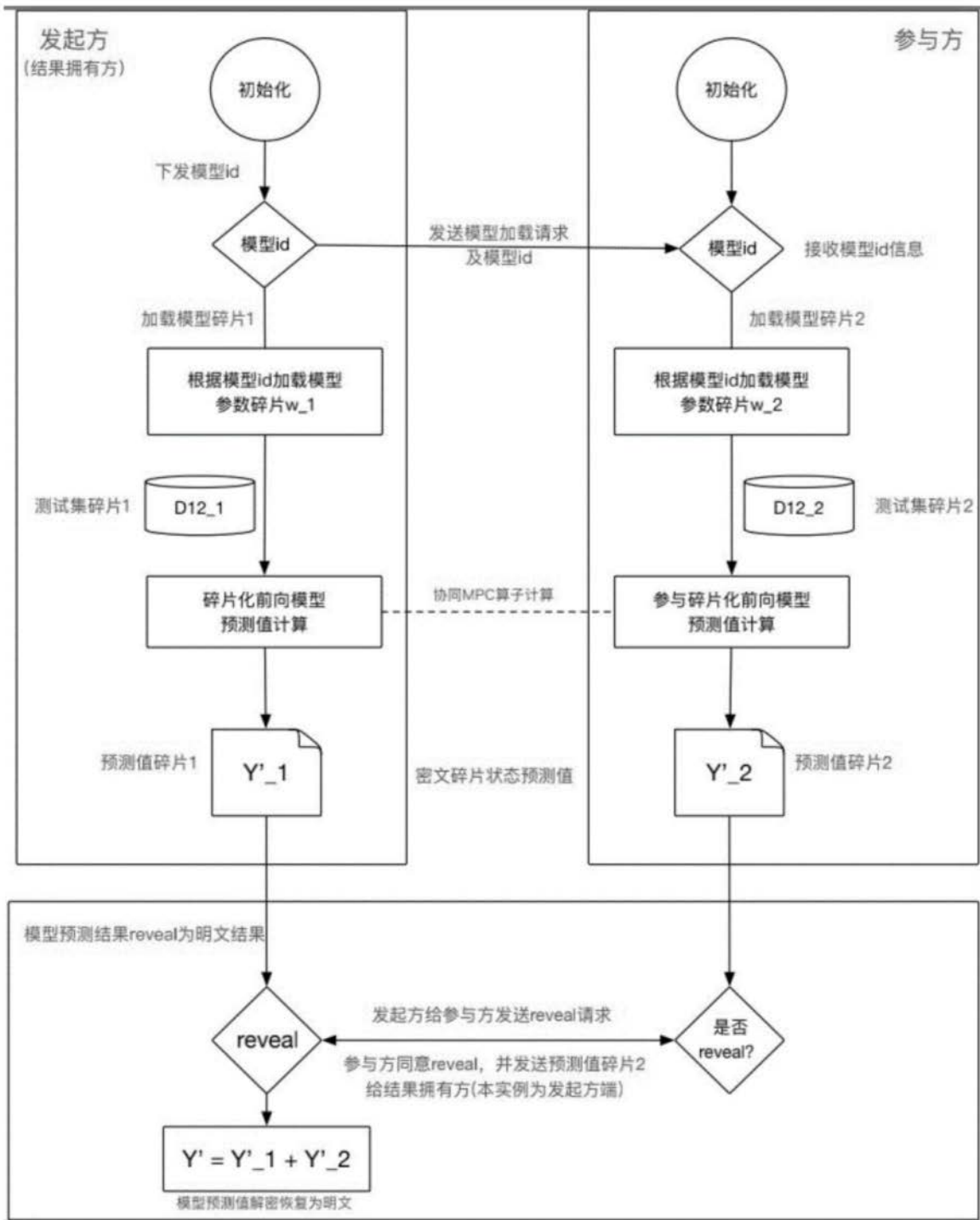


图10

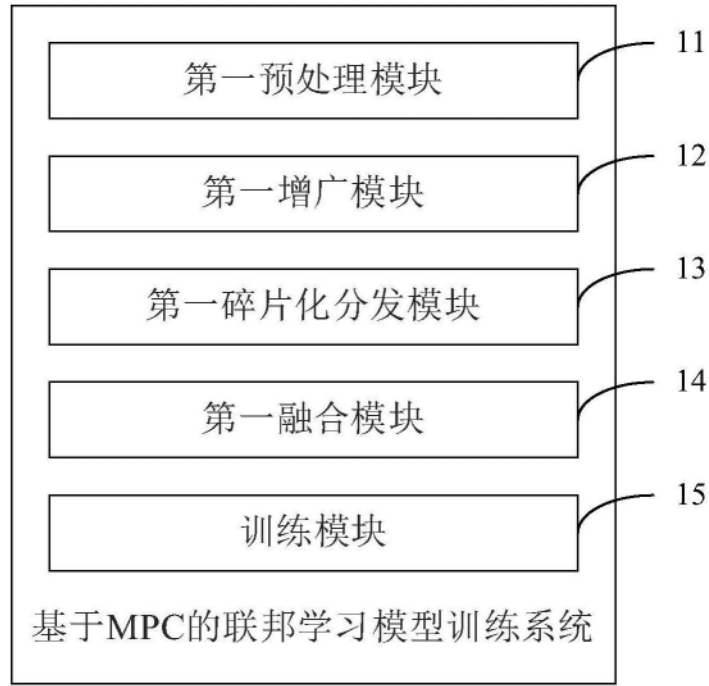


图11

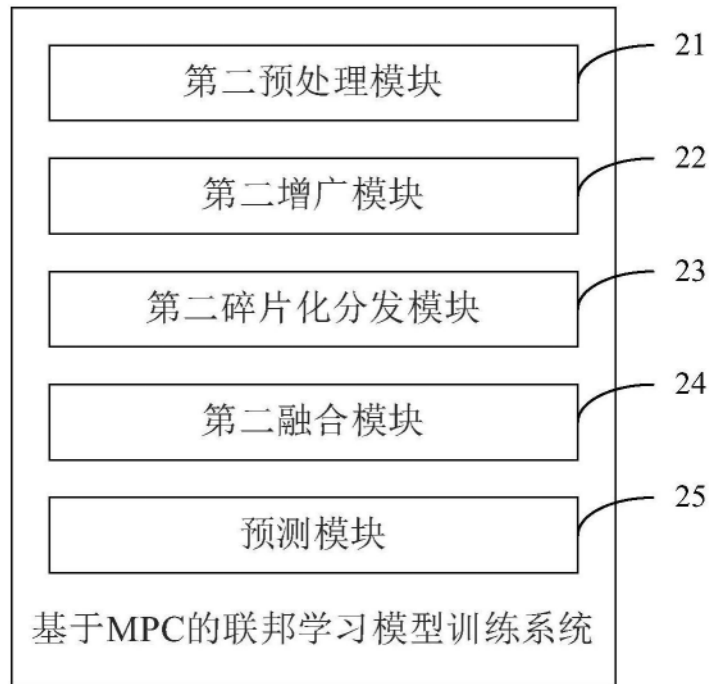


图12

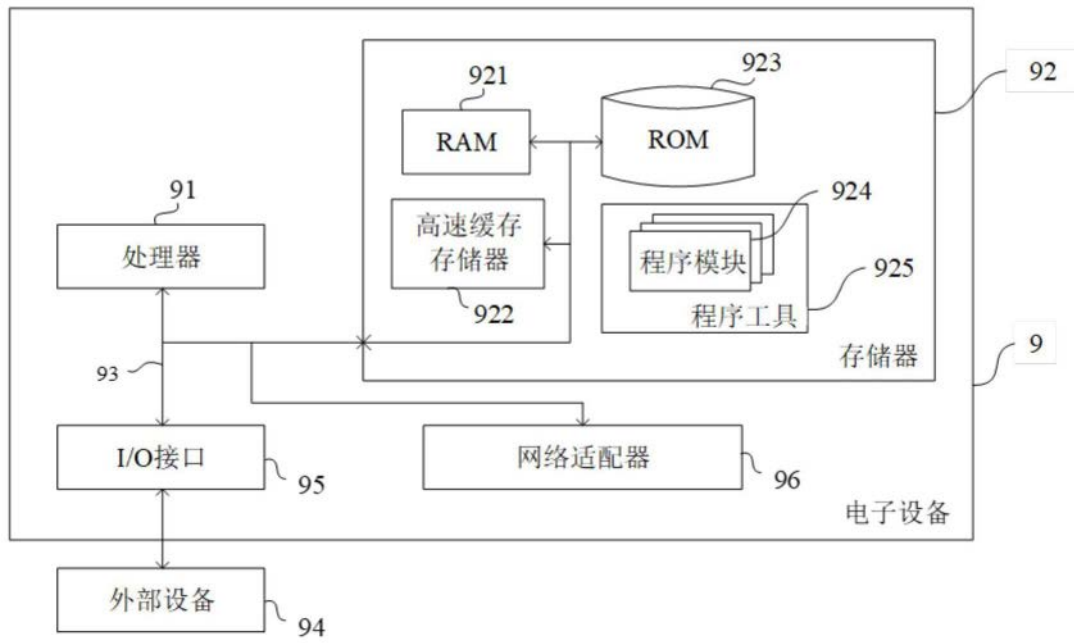


图13