



## (12) 发明专利申请

(10) 申请公布号 CN 113435537 A

(43) 申请公布日 2021.09.24

(21) 申请号 202110806104.1

(22) 申请日 2021.07.16

(71) 申请人 同盾控股有限公司

地址 311121 浙江省杭州市余杭区五常街  
道文一西路998号18幢704室

申请人 同盾科技有限公司

(72) 发明人 周一竞 孟丹 李宏宇 李晓林

(74) 专利代理机构 北京律智知识产权代理有限公司 11438

代理人 王辉 阚梓瑄

(51) Int.Cl.

G06K 9/62 (2006.01)

G06N 20/00 (2019.01)

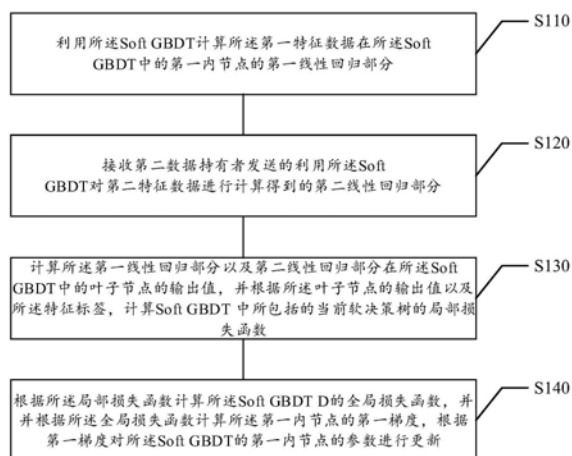
权利要求书3页 说明书18页 附图6页

### (54) 发明名称

基于Soft GBDT的跨特征联邦学习方法、预测方法

### (57) 摘要

本公开是关于一种基于Soft GBDT的跨特征联邦学习方法、数据预测方法及装置,涉及机器学习技术领域,该基于Soft GBDT的跨特征联邦学习方法包括:利用Soft GBDT计算第一特征数据在Soft GBDT中的第一内节点的第一线性回归部分;计算第一线性回归部分以及第二线性回归部分在Soft GBDT中的叶子节点的输出值,并计算Soft GBDT中所包括的当前软决策树的局部损失函数;根据局部损失函数计算Soft GBDT的全局损失函数,并根据全局损失函数计算第一内节点的第一梯度,以对Soft GBDT的第一内节点参数进行更新。本公开有利于加快模型训练速度,提高模型训练效率。



1. 一种基于Soft GBDT的跨特征联邦学习方法,其特征在于,配置于多方机器学习中提供第一特征数据以及特征标签的第一数据持有者,用于对Soft GBDT进行训练,所述基于Soft GBDT的跨特征联邦学习方法包括:

利用所述Soft GBDT计算所述第一特征数据在所述Soft GBDT中的第一内节点的第一线性回归部分;

接收第二数据持有者发送的利用所述Soft GBDT对第二特征数据进行计算得到的第二线性回归部分;

计算所述第一线性回归部分以及第二线性回归部分在所述Soft GBDT中的叶子节点的输出值,并根据所述叶子节点的输出值以及所述特征标签,计算所述Soft GBDT中所包括的当前软决策树的局部损失函数;

根据所述局部损失函数计算所述Soft GBDT的全局损失函数,并根据所述全局损失函数计算所述第一内节点的第一梯度,根据所述第一梯度对所述Soft GBDT的第一内节点参数进行更新。

2. 根据权利要求1所述的基于Soft GBDT的跨特征联邦学习方法,其特征在于,计算所述第一线性回归部分以及第二线性回归部分在所述Soft GBDT中的叶子节点的输出值,包括:

对所述第一线性回归部分以及第二线性回归部分进行求和运算得到第一和运算结果;

利用所述Soft GBDT的叶子节点所在的归一化层对所述第一和运算结果进行归一化处理,得到所述第一和运算结果在所述当前软决策树所在叶子节点的输出值。

3. 根据权利要求2所述的基于Soft GBDT的跨特征联邦学习方法,其特征在于,根据所述叶子节点的输出值以及所述特征标签,计算所述Soft GBDT中所包括的当前软决策树的局部损失函数,包括:

计算排在当前软决策树前面的所有软决策树在叶子节点的输出值的总和;

对排在当前软决策树前面的所有软决策树在叶子节点的输出值的总和进行归一化运算,得到所述第一和运算结果中所包括的样本特征在排在当前软决策树前面的所有软决策树中的第一预测结果;

根据所述第一预测结果、所述样本特征在当前软决策树中的叶子节点的输出值以及与所述样本特征对应的特征标签,构建所述当前软决策树的局部损失函数。

4. 根据权利要求3所述的基于Soft GBDT的跨特征联邦学习方法,其特征在于,根据所述全局损失函数计算所述第一内节点的第一梯度,包括:

根据所述全局损失函数对所述第一内节点参数进行一阶求导,并根据所述局部损失函数对所述叶子节点参数进行二阶求导;

根据所述全局损失函数的一阶导数、局部损失函数的二阶导数以及所述软决策树的输出值,计算所述软决策树的第一内节点的第一梯度。

5. 根据权利要求4所述的基于Soft GBDT的跨特征联邦学习方法,其特征在于,所述基于Soft GBDT的跨特征联邦学习方法还包括:

根据所述全局损失函数对所述叶子节点参数进行一阶求导,得到所述软决策树中所包括的叶子节点的叶子梯度;

根据所述叶子节点的叶子梯度,对所述软决策树中的叶子节点参数进行更新。

6.一种基于Soft GBDT的跨特征联邦学习方法,其特征在于,配置于多方机器学习中提供第二特征数据的第二数据持有者,用于对Soft GBDT进行训练,所述基于Soft GBDT的跨特征联邦学习方法包括:

利用所述Soft GBDT计算所述第二特征数据在所述Soft GBDT中的内节点的第二线性回归部分,并将所述第二线性回归部分发送至第一数据持有者;

接收所述第一数据持有者发送的加密后的第一梯度;其中,所述第一梯度是根据第一线性回归部分以及第二线性回归部分计算得到的,所述第一线性回归部分是所述第一数据持有者利用所述Soft GBDT对与所述第二特征数据具有相同数据生产者的第一特征数据计算得到的,所述第一梯度用于对所述第一数据持有者侧的所述Soft GBDT的第一内节点参数进行更新;

根据所述加密后的第一梯度以及第二特征数据计算第二梯度,并接收所述第一数据持有者发送的解密后的第二梯度;

利用解密后的第二梯度对所述第二数据持有者侧的Soft GBDT的第二内节点参数进行更新。

7.一种基于Soft GBDT的跨特征联邦预测方法,其特征在于,配置于多方机器学习中提供第一待预测数据的第一数据持有者,用于根据对Soft GBDT进行训练得到的数据预测模型进行数据预测,所述基于Soft GBDT的跨特征联邦预测方法包括:

利用所述数据预测模型计算所述第一待预测数据在所述数据预测模型中的内节点的第三线性回归部分;其中,所述数据预测模型是根据权利要求1-6任一项所述基于Soft GBDT的跨特征联邦学习方法对Soft GBDT训练得到的;

接收第二数据持有者发送的利用所述数据预测模型对第二待预测数据进行计算得到的第四线性回归部分;

对所述第三线性回归部分以及第四线性回归部分进行求和运算,并利用所述数据预测模型的叶子节点所在的归一化层对第二和运算结果进行归一化处理,得到数据预测结果。

8.根据权利要求7所述的基于Soft GBDT的跨特征联邦预测方法,其特征在于,利用所述数据预测模型的叶子节点所在的归一化层对第二和运算结果进行归一化处理,得到数据预测结果,包括:

利用所述数据预测模型的叶子节点所在的归一化层计算所述第二和运算结果在每一个叶子节点上的分支概率;

根据所述分支概率计算所述数据预测结果。

9.一种基于Soft GBDT的跨特征联邦学习装置,其特征在于,配置于多方机器学习中提供第一特征数据以及特征标签的第一数据持有者,用于对Soft GBDT进行训练,所述基于Soft GBDT的跨特征联邦学习装置包括:

第一计算模块,用于利用所述Soft GBDT计算所述第一特征数据在所述Soft GBDT中的第一内节点的第一线性回归部分;

第一接收模块,用于接收第二数据持有者发送的利用所述Soft GBDT对第二特征数据进行计算得到的第二线性回归部分;其中,所述第一特征数据以及第二特征数据的数据生产者相同;

第二计算模块,用于计算所述第一线性回归部分以及第二线性回归部分在所述Soft

GBDT中的叶子节点的输出值,并根据所述叶子节点的输出值以及所述特征标签,计算所述Soft GBDT中所包括的当前软决策树的局部损失函数;

第一参数更新模块,用于根据所述局部损失函数计算所述Soft GBDT的全局损失函数,并根据所述全局损失函数计算所述第一内节点的第一梯度,以根据所述第一梯度对所述Soft GBDT的第一内节点参数进行更新。

10.一种基于Soft GBDT的跨特征联邦预测装置,其特征在于,配置于多方机器学习中提供第一待预测数据的第一数据持有者,用于根据对Soft GBDT进行训练得到的数据预测模型进行数据预测,所述基于Soft GBDT的跨特征联邦预测方法包括:

第五计算模块,用于利用所述数据预测模型计算所述第一待预测数据在所述数据预测模型中的内节点的第三线性回归部分;其中,所述数据预测模型是根据权利要求1-6任一项所述基于Soft GBDT的跨特征联邦学习方法对Soft GBDT训练得到的;

第三接收模块,用于接收第二数据持有者发送的利用所述数据预测模型对第二待预测数据进行计算得到的第四线性回归部分;

数据预测模块,用于对所述第三线性回归部分以及第四线性回归部分进行求和运算,并利用所述数据预测模型的叶子节点所在的归一化层对第二和运算结果进行归一化处理,得到数据预测结果。

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权利要求1-6任一项所述的基于Soft GBDT的跨特征联邦学习方法,或者7-8任一项所述的基于Soft GBDT的跨特征联邦预测方法。

12.一种电子设备,其特征在于,包括:

处理器;以及

存储器,用于存储所述处理器的可执行指令;

其中,所述处理器配置为经由执行所述可执行指令来执行权利要求1-6任一项所述的基于Soft GBDT的跨特征联邦学习方法,或者7-8任一项所述的基于Soft GBDT的跨特征联邦预测方法。

## 基于Soft GBDT的跨特征联邦学习方法、预测方法

### 技术领域

[0001] 本公开实施例涉及机器学习技术领域,具体而言,涉及一种基于Soft GBDT的跨特征联邦学习方法、基于Soft GBDT的跨特征联邦学习装置、基于Soft GBDT的跨特征联邦预测方法、基于Soft GBDT的跨特征联邦预测装置、计算机可读存储介质以及电子设备。

### 背景技术

[0002] 隐私保护机器学习中的跨特征场景,隶属于知识联邦理论体系中的跨特征联邦,是指多机构所使用的训练或推理样本是一致的,而特征不同,只有一个机构持有标签,训练和推理可以在多方协助下完成。

[0003] 在当前的基于跨特征场景的机器学习中,基于决策树的模型方法依旧是主流的以及实用的模型训练方法。

[0004] 但是,基于Boosting的决策树模型,例如XGboost,由于其内部结构是串联的,无法进行独立的计算,进而使得在联邦化计算的场景下,模型的训练速度慢,训练效率较低。

[0005] 因此,需要提供一种新的基于Soft GBDT的跨特征联邦学习方法、预测方法及装置。

[0006] 需要说明的是,在上述背景技术部分发明的信息仅用于加强对本公开的背景的理解,因此,可以包括不构成对本领域普通技术人员已知的现有技术的信息。

### 发明内容

[0007] 本公开的目的在于提供一种基于Soft GBDT的跨特征联邦学习方法、基于Soft GBDT的跨特征联邦学习装置、基于Soft GBDT的跨特征联邦预测方法、基于Soft GBDT的跨特征联邦预测装置、计算机可读存储介质以及电子设备,进而至少在一定程度上克服由于相关技术的限制和缺陷导致的模型的训练速度慢,训练效率低的问题。

[0008] 根据本公开的一个方面,提供一种基于Soft GBDT的跨特征联邦学习方法,配置于多方机器学习中提供第一特征数据以及特征标签的第一数据持有者,用于对Soft GBDT进行训练,所述基于Soft GBDT的跨特征联邦学习方法包括:

[0009] 利用所述Soft GBDT计算所述第一特征数据在所述Soft GBDT中的第一内节点的第一线性回归部分;

[0010] 接收第二数据持有者发送的利用所述Soft GBDT对第二特征数据进行计算得到的第二线性回归部分;

[0011] 计算所述第一线性回归部分以及第二线性回归部分在所述Soft GBDT中的叶子节点的输出值,并根据所述叶子节点的输出值以及所述特征标签,计算所述Soft GBDT中所包括的当前软决策树的局部损失函数;

[0012] 根据所述局部损失函数计算所述Soft GBDT的全局损失函数,并根据所述全局损失函数计算所述第一内节点的第一梯度,根据所述第一梯度对所述Soft GBDT的第一内节点的参数进行更新。

[0013] 在本公开的一种示例性实施例中,计算所述第一线性回归部分以及第二线性回归部分在所述Soft GBDT中的叶子节点的输出值,包括:

[0014] 对所述第一线性回归部分以及第二线性回归部分进行求和运算得到第一和运算结果;

[0015] 利用所述Soft GBDT的叶子节点所在的归一化层对所述第一和运算结果进行归一化处理,得到所述第一和运算结果在所述当前软决策树所在叶子节点的输出值。

[0016] 在本公开的一种示例性实施例中,根据所述叶子节点的输出值以及所述特征标签,计算所述Soft GBDT中所包括的当前软决策树的局部损失函数,包括:

[0017] 计算排在当前软决策树前面的所有软决策树在叶子节点的输出值的总和;

[0018] 对排在当前软决策树前面的所有软决策树在叶子节点的输出值的总和进行归一化运算,得到所述第一和运算结果中所包括的样本特征在排在当前软决策树前面的所有软决策树中的第一预测结果;

[0019] 根据所述第一预测结果、所述样本特征在当前软决策树中的叶子节点的输出值以及与所述样本特征对应的特征标签,构建所述当前软决策树的局部损失函数。

[0020] 在本公开的一种示例性实施例中,根据所述全局损失函数计算所述第一内节点的第一梯度,包括:

[0021] 根据所述全局损失函数对所述第一内节点的参数进行一阶求导,并根据所述局部损失函数对所述叶子节点的参数进行二阶求导;

[0022] 根据所述全局损失函数的一阶导数、局部损失函数的二阶导数以及所述软决策树的输出值,计算所述软决策树的第一内节点的第一梯度。

[0023] 在本公开的一种示例性实施例中,所述基于Soft GBDT的跨特征联邦学习方法还包括:

[0024] 根据所述全局损失函数对所述叶子节点的参数进行一阶求导,得到所述软决策树中所包括的叶子节点的叶子梯度;

[0025] 根据所述叶子节点的叶子梯度,对所述软决策树中的叶子节点的参数进行更新。

[0026] 根据本公开的一个方面,提供一种基于Soft GBDT的跨特征联邦学习方法,配置于多方机器学习中提供第二特征数据第二数据持有者,用于对Soft GBDT进行训练,所述基于Soft GBDT的跨特征联邦学习方法包括:

[0027] 利用所述Soft GBDT计算所述第二特征数据在所述Soft GBDT中的内节点的第二线性回归部分,并将所述第二线性回归部分发送至第一数据持有者;

[0028] 接收所述第一数据持有者发送的加密后的第一梯度;其中,所述第一梯度是根据第一线性回归部分以及第二线性回归部分计算得到的,所述第一线性回归部分是所述第一数据持有者利用所述Soft GBDT对与所述第二特征数据具有相同数据生产者的第一特征数据计算得到的,所述第一梯度用于对所述第一数据持有者侧的所述Soft GBDT的第一内节点的参数进行更新;

[0029] 根据所述加密后的第一梯度以及第二特征数据计算第二梯度,并接收所述第一数据持有者发送的解密后的第二梯度;

[0030] 利用解密后的第二梯度对所述第二数据持有者侧的Soft GBDT的第二内节点的参数进行更新。



[0031] 根据本公开的一个方面,提供一种基于Soft GBDT的跨特征联邦预测方法,配置于多方机器学习中提供第一待预测数据的第一数据持有者,用于根据对Soft GBDT进行训练得到的数据预测模型进行数据预测,所述基于Soft GBDT的跨特征联邦预测方法包括:

[0032] 利用所述数据预测模型计算所述第一待预测数据在所述数据预测模型中的内节点的第三线性回归部分;其中,所述数据预测模型是根据上述任意一项所述基于Soft GBDT的跨特征联邦学习方法对Soft GBDT训练得到的;

[0033] 接收第二数据持有者发送的利用所述数据预测模型对第二待预测数据进行计算得到的第四线性回归部分;

[0034] 对所述第三线性回归部分以及第四线性回归部分进行求和运算,并利用所述数据预测模型的叶子节点所在的归一化层对第二和运算结果进行归一化处理,得到数据预测结果。

[0035] 在本公开的一种示例性实施例中,利用所述数据预测模型的叶子节点所在的归一化层对第二和运算结果进行归一化处理,得到数据预测结果,包括:

[0036] 利用所述数据预测模型的叶子节点所在的归一化层计算所述第二和运算结果在每一个叶子节点上的分支概率;

[0037] 根据所述分支概率计算所述数据预测结果。

[0038] 根据本公开的一个方面,提供一种基于Soft GBDT的跨特征联邦学习装置,配置于多方机器学习中提供第一特征数据以及特征标签的第一数据持有者,用于对Soft GBDT进行训练,所述基于Soft GBDT的跨特征联邦学习装置包括:

[0039] 第一计算模块,用于利用所述Soft GBDT计算所述第一特征数据在所述Soft GBDT中的第一内节点的第一线性回归部分;

[0040] 第一接收模块,用于接收第二数据持有者发送的利用所述Soft GBDT对第二特征数据进行计算得到的第二线性回归部分;其中,所述第一特征数据以及第二特征数据的数据生产者相同;

[0041] 第二计算模块,用于计算所述第一线性回归部分以及第二线性回归部分在所述Soft GBDT中的叶子节点的输出值,并根据所述叶子节点的输出值以及所述特征标签,计算所述Soft GBDT中所包括的当前软决策树的局部损失函数;

[0042] 第一参数更新模块,用于根据所述局部损失函数计算所述Soft GBDT的全局损失函数,并根据所述全局损失函数计算所述第一内节点的第一梯度,以根据所述第一梯度对所述Soft GBDT的第一内节点参数进行更新。

[0043] 根据本公开的一个方面,提供一种基于Soft GBDT的跨特征联邦预测装置,配置于多方机器学习中提供第一待预测数据的第一数据持有者,用于根据对Soft GBDT进行训练得到的数据预测模型进行数据预测,所述基于Soft GBDT的跨特征联邦预测方法包括:

[0044] 第五计算模块,用于利用所述数据预测模型计算所述第一待预测数据在所述数据预测模型中的内节点的第三线性回归部分;其中,所述数据预测模型是根据上述任意一项所述基于Soft GBDT的跨特征联邦学习方法对Soft GBDT训练得到的;

[0045] 第三接收模块,用于接收第二数据持有者发送的利用所述数据预测模型对第二待预测数据进行计算得到的第四线性回归部分;

[0046] 数据预测模块,用于对所述第三线性回归部分以及第四线性回归部分进行求和运

算,并利用所述数据预测模型的叶子节点所在的归一化层对第二和运算结果进行归一化处理,得到数据预测结果。

[0047] 根据本公开的一个方面,提供一种计算机可读存储介质,其上存储有计算机程序,所述计算机程序被处理器执行时实现上述任意一项所述的基于Soft GBDT的跨特征联邦学习方法,或者上述任一项所述的基于Soft GBDT的跨特征联邦预测方法。

[0048] 根据本公开的一个方面,提供一种电子设备,包括:

[0049] 处理器;以及

[0050] 存储器,用于存储所述处理器的可执行指令;

[0051] 其中,所述处理器配置为经由执行所述可执行指令来执行上述任意一项所述的基于Soft GBDT的跨特征联邦学习方法,或者上述任一项所述的基于Soft GBDT的跨特征联邦预测方法。

[0052] 本公开实施例提供的一种基于Soft GBDT的跨特征联邦学习方法,一方面,由于Soft GBDT内部所包括的每一棵树是独立的,其可以进行独立的计算,进而可以避免由于无法进行独立的计算,使得在联邦化计算的场景下,模型的训练效率较低的问题;另一方面,由于可以通过计算第一线性回归部分以及第二线性回归部分在Soft GBDT中的叶子节点的输出值,并根据叶子节点的输出值以及特征标签,计算当前软决策树的局部损失函数;再根据局部损失函数计算Soft GBDT的全局损失函数,并根据全局损失函数计算第一内节点的第一梯度,以根据第一梯度对Soft GBDT的第一内节点参数进行更新,第二数据持有者也可以根据第一梯度对第二内节点参数进行更新,进而使得第二数据持有者在不泄露本方所持有的第二特征数据的前提下,即可根据全局损失函数对本方的第二内节点参数进行更新,确保了第二特征树的安全性;再一方面,由于整体的训练过程是分布在多方进行的,进而避免由于需要集中在某一方进行训练导致的数据处理负担较重,训练效率较低的问题,进一步的提高了模型训练效率。

[0053] 应当理解的是,以上的一般描述和后文的细节描述仅是示例性和解释性的,并不能限制本公开。

## 附图说明

[0054] 此处的附图被并入说明书中并构成本说明书的一部分,示出了符合本公开的实施例,并与说明书一起用于解释本公开的原理。显而易见地,下面描述中的附图仅仅是本公开的一些实施例,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。

[0055] 图1示意性示出根据本公开示例实施例的一种在第一数据持有者侧进行的基于Soft GBDT的跨特征联邦学习方法的流程图。

[0056] 图2示意性示出根据本公开示例实施例的一种软决策树的结构示例图。

[0057] 图3示意性示出根据本公开示例实施例的一种模型训练系统的框图。

[0058] 图4示意性示出根据本公开示例实施例的一种根据所述叶子节点的输出值以及所述特征标签,计算所述Soft GBDT中所包括的当前软决策树的局部损失函数的方法流程图。

[0059] 图5示意性示出根据本公开示例实施例的一种在第二数据持有者侧进行的基于Soft GBDT的跨特征联邦学习方法的流程图。



[0060] 图6示意性示出根据本公开示例实施例的一种在第一数据持有者侧以及第二数据持有者侧进行跨特征的基于Soft GBDT的跨特征联邦学习方法的流程图。

[0061] 图7示意性示出根据本公开示例实施例的一种基于Soft GBDT的跨特征联邦预测方法的流程图。

[0062] 图8示意性示出根据本公开示例实施例的一种基于Soft GBDT的跨特征联邦学习装置的框图。

[0063] 图9示意性示出根据本公开示例实施例的另一种基于Soft GBDT的跨特征联邦学习装置的框图。

[0064] 图10示意性示出根据本公开示例实施例的一种基于Soft GBDT的跨特征联邦预测装置的框图。

[0065] 图11示意性示出根据本公开示例实施例的一种用于实现上述基于Soft GBDT的跨特征联邦学习方法、基于Soft GBDT的跨特征联邦预测方法的电子设备示意图。

### 具体实施方式

[0066] 现在将参考附图更全面地描述示例实施方式。然而,示例实施方式能够以多种形式实施,且不应被理解为限于在此阐述的范例;相反,提供这些实施方式使得本公开将更加全面和完整,并将示例实施方式的构思全面地传达给本领域的技术人员。所描述的特征、结构或特性可以以任何合适的方式结合在一个或更多实施方式中。在下面的描述中,提供许多具体细节从而给出对本公开的实施方式的充分理解。然而,本领域技术人员将意识到,可以实践本公开的技术方案而省略所述特定细节中的一个或更多,或者可以采用其它的方法、组元、装置、步骤等。在其它情况下,不详细示出或描述公知技术方案以避免喧宾夺主而使得本公开的各方面变得模糊。

[0067] 此外,附图仅为本公开的示意性图解,并非一定是按比例绘制。图中相同的附图标记表示相同或类似的部分,因而将省略对它们的重复描述。附图中所示的一些方框图是功能实体,不一定必须与物理或逻辑上独立的实体相对应。可以采用软件形式来实现这些功能实体,或在一个或多个硬件模块或集成电路中实现这些功能实体,或在不同网络和/或处理器装置和/或微控制器装置中实现这些功能实体。

[0068] 本公开示例实施例是为了确保数据的保密性和安全性的情况下,能够充分利用多方数据特性进行联邦建模的目的,涉及到“知识联邦”,知识联邦旨在确保各个训练方服务器的数据不离开本地的情况下交换数据中的“知识”,从而建立一个充分利用各训练方服务器的本地数据的模型,达到“数据可用不可见,知识共创可共享”的目的。

[0069] 根据各个训练方服务器的数据分布的特点,知识联邦可分为跨特征联邦、跨样本联邦以及复合型联邦,其中,跨特征联邦是指不同训练方服务器中有很多共同的用户样本,但样本特征数据分布不同,可能只有一方是有标签数据;跨样本联邦是指每个训练方服务器的样本数据具有相同的特征分布,但各方的样本数据是独立的,而且每个参与方服务器都有自己样本对应的标签数据;复合型联邦是指同时涉及跨样本联邦和跨特征联邦,只有一小部分的样本或特征是参与各方的交集,其他数据无论是特征分布或样本分布都是不相同的。

[0070] 本示例实施方式中首先提供了一种基于Soft GBDT的跨特征联邦学习方法,用于

对Soft GBDT进行训练;该方法可以运行于多方机器学习中提供第一特征数据以及特征标签的第一数据持有者所在的服务器、服务器集群或云服务器等;当然,本领域技术人员也可以根据需求在其他平台运行本公开的方法,本示例性实施例中对此不做特殊限定。具体的,参考图1所示,该基于Soft GBDT的跨特征联邦学习方法可以包括以下步骤:

[0071] 步骤S110.利用所述Soft GBDT计算所述第一特征数据在所述Soft GBDT中的第一内节点的第一线性回归部分;

[0072] 步骤S120.接收第二数据持有者发送的利用所述Soft GBDT对第二特征数据进行计算得到的第二线性回归部分;其中,所述第一特征数据以及第二特征数据的数据生产者相同;

[0073] 步骤S130.计算所述第一线性回归部分以及第二线性回归部分在所述Soft GBDT中的叶子节点的输出值,并根据所述叶子节点的输出值以及所述特征标签,计算所述Soft GBDT中所包括的当前软决策树的局部损失函数;

[0074] 步骤S140.根据所述局部损失函数计算所述Soft GBDT的全局损失函数,并根据所述全局损失函数计算所述第一内节点的第一梯度,以根据所述第一梯度对所述Soft GBDT的第一内节点参数进行更新。

[0075] 上述基于Soft GBDT的跨特征联邦学习方法中,一方面,由于Soft GBDT内部所包括的每一棵树是独立的,其可以进行独立的计算,进而可以避免由于无法进行独立的计算,使得在联邦化计算的场景下,模型的训练效率较低的问题;另一方面,由于可以通过计算第一线性回归部分以及第二线性回归部分在Soft GBDT中的叶子节点的输出值,并根据叶子节点的输出值以及特征标签,计算当前软决策树的局部损失函数;再根据局部损失函数计算Soft GBDT的全局损失函数,并根据全局损失函数计算第一内节点的第一梯度,以根据第一梯度对Soft GBDT的第一内节点参数进行更新,第二数据持有者也可以根据第一梯度对第二内节点参数进行更新,进而使得第二数据持有者在不泄露本方所持有的第二特征数据的前提下,即可根据全局损失函数对本方的第二内节点参数进行更新,确保了第二特征树的安全性;再一方面,由于整体的训练过程是分布在多方进行的,进而避免由于需要集中在某一方进行训练进而导致的数据处理负担较重,训练速度慢,效率低的问题,进一步的提高了模型训练效率。

[0076] 以下,将结合附图对本公开示例实施例基于Soft GBDT的跨特征联邦学习方法进行详细的解释以及说明。

[0077] 首先,对本公开示例实施例所涉及到的软决策树进行解释以及说明。

[0078] 软决策树,区别于普通的决策树,每层用逻辑回归来决策左右分支概率,最后一层的叶子节点用输出固定分布,最终结果可取概率最大路径对应的分类分布,或是所有路径的加权平均(权重即每条路径的概率),具体的结构示例图可以参考图2所示。

[0079] 具体的,在软决策树中,每个样本不再是严格落在一个分支上,而是以一定的概率落在两个或者多个分支上;软决策树的每个节点可以单独计算,所以可以并行化。和GBDT类似,soft GBDT由多个基础学习器组成,这里基础学习器是软决策树;对每棵树在上一棵树的基础上的残差进行局部损失计算,多棵树串联组成了soft GBDT;同时,由于每棵树的计算是独立的可以并行计算,且每棵树的损失函数和GBDT的计算是一致的,最后会使用一个全局的损失函数,进而可以提高最终的计算效率。

[0080] 同时,在训练过程中,可以通过最小化全局损失函数 $L$ 来训练整个模型;虽然,每棵树的结果拟合的是上一颗树的残差,但由于软决策树的节点和树之间的独立性,只需要在最后综合所有树的局部损失函数进行求和,使得所有局部损失函数的和的结果最小化来进行训练,这样可以充分利用并行计算来大大提升模型训练的速度。

[0081] 其次,对本公开示例实施例的应用场景进行解释以及说明。

[0082] 基于上述基于软决策树的梯度提升学习器,考虑跨特征的联邦方案,本方案适用于任何联邦学习的分类问题。具体的,在跨特征场景中,所有的参与方拥有相同的样本但是特征不相交,也即所有参与方所持有的特征数据的数据生产者是不同的,例如,同一用户,在第一数据持有者侧所具有的特征为基本用户信息,在第二数据持有者侧所具有的特征为银行交易信息等等,本示例对此不做特殊限制;同时,所有参与方使用相同的模型结构,该模型结构可以包括 $K$ 个base learner(基础学习器,也即软决策树),每个base learner的大小即深度`tree_depth`一致(或者每个参与方对应的树的大小一致)。

[0083] 进一步的,对本公开示例实施例中所涉及到的模型训练系统进行解释以及说明。具体的本公开示例实施例的基于Soft GBDT的跨特征联邦学习方法,是对不少于两个训练方服务器所拥有的样本数据进行跨样本联邦训练(例如第一数据持有者以及第二数据持有者),各训练方服务器得到联邦Soft GBDT。

[0084] 如图3所示,该模型训练系统可以包括提供第一特征数据以及特征标签的第一数据持有者310以及提供第二特征数据的第二数据持有者320;第一数据持有者、第二数据持有者之间互相通过网络进行连接。

[0085] 其中,第一数据持有者可以用于:利用所述Soft GBDT计算所述第一特征数据在所述Soft GBDT中的第一内节点的第一线性回归部分;接收第二数据持有者发送的利用所述Soft GBDT对第二特征数据进行计算得到的第二线性回归部分;其中,所述第一特征数据以及第二特征数据的数据生产者不同;计算所述第一线性回归部分以及第二线性回归部分在所述Soft GBDT中的叶子节点的输出值,并根据所述叶子节点的输出值以及所述特征标签,计算所述Soft GBDT中所包括的当前软决策树的局部损失函数;根据所述局部损失函数计算所述Soft GBDT的全局损失函数,并根据所述全局损失函数计算所述第一内节点的第一梯度,以根据所述第一梯度对所述Soft GBDT的第一内节点参数进行更新。

[0086] 第二数据持有者可以用于:利用所述Soft GBDT计算所述第二特征数据在所述Soft GBDT中的内节点的第二线性回归部分,并将所述第二线性回归部分发送至第一数据持有者;接收所述第一数据持有者发送的加密后的第一梯度;其中,所述第一梯度是根据第一线性回归部分以及第二线性回归部分计算得到的,所述第一线性回归部分是所述第一数据持有者利用所述Soft GBDT对与所述第二特征数据具有相同数据生产者的第一特征数据计算得到的,所述第一梯度用于对所述第一数据持有者侧的所述Soft GBDT的第一内节点参数进行更新;根据所述加密后的第一梯度以及第二特征数据计算第二梯度,并接收所述第一数据持有者发送的解密后的第二梯度;利用解密后的第二梯度对所述第二数据持有者侧的Soft GBDT的第二内节点参数进行更新。

[0087] 以下,结合上述软决策树、具体的应用场景以及模型训练系统,对本公开示例实施例基于Soft GBDT的跨特征联邦学习方法中所涉及到的各步骤进行详细的解释以及说明。

[0088] 在本公开示例实施例的一种基于Soft GBDT的跨特征联邦学习方法中:

[0089] 在步骤S110中,利用所述Soft GBDT计算所述第一特征数据在所述Soft GBDT中的第一内节点的第一线性回归部分。

[0090] 具体的,参考图2所示,第一数据持有者可以通过软决策树的第一内节点计算第一特征数据的第一线性回归部分(Logistic Regression)。具体的,该第一内节点可以包括多层,比如第零层、第一层、第二层等等,具体的层数可以根据实际需要自行设定,本示例对此不做特殊限制。进一步的,该第一特征数据例如可以包括用户的基本属性信息,例如姓名、年龄、职业、学历、婚育情况、家庭住址等等,当然也可以包括其他用户信息,本示例对此不做特殊限制。

[0091] 在步骤S120中,接收第二数据持有者发送的利用所述Soft GBDT对第二特征数据进行计算得到的第二线性回归部分;其中,所述第一特征数据以及第二特征数据的数据生产者相同。

[0092] 具体的,为了实现跨特征的联邦计算,还需要考虑第二数据持有者侧(例如银行或者其他第三方平台)持有的第二特征数据以及与第二特征数据对应的第二线性回归部分;其中,第二特征数据可以包括用户的交易信息、消费信息、信用等级等等,也可以包括其他用户信息,本示例对此不做特殊限制。此处需要补充说明的是,在跨特征的联邦计算过程中,不同的数据持有者所具有的用户样本(也即数据生产者)是相同的,但是各数据持有者所持有的样本特征数据是不同的。

[0093] 在步骤S130中,计算所述第一线性回归部分以及第二线性回归部分在所述Soft GBDT中的叶子节点的输出值,并根据所述叶子节点的输出值以及所述特征标签,计算所述Soft GBDT中所包括的当前软决策树的局部损失函数。

[0094] 在本示例实施例中,首先,计算所述第一线性回归部分以及第二线性回归部分在所述Soft GBDT中的叶子节点的输出值。具体的,可以包括:首先,对所述第一线性回归部分以及第二线性回归部分进行求和运算得到第一和运算结果;利用所述Soft GBDT的叶子节点所在的归一化层对所述第一和运算结果进行归一化处理,得到所述第一和运算结果在所述当前软决策树所在叶子节点的输出值。

[0095] 详细来说,计算所述第一线性回归部分以及第二线性回归部分的和,具体可以如下公式(1)所示:

[0096]  $w = u + v$ ; 公式(1)

[0097] 其中, $w$ 为第一和运算结果, $u$ 为第一线性回归部分, $v$ 为第二线性回归部分。

[0098] 其次,叶子节点的输出值。具体可以如下公式(2)所示:

[0099]  $p = \text{sigmoid}(w) = \text{sigmoid}(u + v)$ ; 公式(2)

[0100] 其次,根据所述叶子节点的输出值以及所述特征标签,计算所述Soft GBDT中所包括的当前软决策树的局部损失函数。具体的,参考图4所示,可以包括以下步骤:

[0101] 步骤S410,计算排在当前软决策树前面的所有软决策树在叶子节点的输出值的总和;

[0102] 步骤S420,对排在当前软决策树前面的所有软决策树在叶子节点的输出值的总和进行归一化运算,得到所述第一和运算结果中所包括的样本特征在排在当前软决策树前面的所有软决策树中的第一预测结果;

[0103] 步骤S430,根据所述第一预测结果、所述样本特征在当前软决策树中的叶子节点

的输出值以及与所述样本特征对应的特征标签,构建所述当前软决策树的局部损失函数。

[0104] 以下,将对步骤S410-步骤S430进行解释以及说明。首先,假设共有M棵软决策树,数据集共有K个类别,每棵软决策树的每一个叶子节点有K个输出  $\Phi_{mk}$ ,  $k=1, \dots, K$ 。若每个参与方(第一数据持有者以及第二数据持有者)的特征个数为num\_fea,则每个内节点的参数为w, w为大小为num\_fea\*1的向量。

[0105] 其次,为求得第一数据持有者以及第二数据持有者的模型参数的梯度,进行模型的更新训练,需要根据模型的损失函数对各个参数进行求导,得到导数的表达式,进而可以更新参数。因此,令  $l_m$  为第m棵树的局部损失函数,  $r_m^i$  为第  $x_i$  个样本在前第m-1棵树上的残差,  $o_m^i$  是指第m棵树的输出结果,  $r_m^i$  以及  $o_m^i$  的维度均为  $K*1$ 。同时,  $l_m$  表示的是用第m棵树的输出来拟合前面m-1棵树的残差。  $\Phi_s$  表示第s棵树的叶子节点的参数,其仅在第一数据持有者侧拥有。

$$[0106] \quad l_m = \sum_{x_i \in B} \|r_m^i - o_m^i\|_2^2; \quad \text{公式 (1)}$$

[0107] 然后,令loss表示交叉熵损失函数,那么第  $x_i$  个样本在前第m-1棵树上的残差  $r_m^i$  可以表示为:

$$[0108] \quad r_{mk}^i = -\frac{\partial \text{loss}(\Phi_m^i, y^i)}{\partial \Phi_{mk}^i}; \quad \text{公式 (2)}$$

[0109] 其中,  $\Phi_m^i$  表示前m-1棵树的输出结果的加和,那么有:

$$[0110] \quad \Phi_m^i = (\Phi_{m1}^i, \Phi_{m2}^i, \dots, \Phi_{mk}^i)^T = \sum_{n=0}^{m-1} o_n^i = \sum_{n=0}^{m-1} (o_{n1}^i, o_{n2}^i, \dots, o_{nK}^i); \quad \text{公式 (3)}$$

$$[0111] \quad \begin{cases} \text{loss}(\Phi_m^i, y^i) = -\sum_k y_k^i \log(p_{k,m-1}^i) \\ \frac{\partial \text{loss}}{\partial \Phi_{mk}^i} = y_k^i - p_{k,m-1}^i \end{cases}; \quad \text{公式 (4)}$$

[0112] 其中,  $p_{k,m-1}$  表示前m-1棵树对于样本  $x_i$  的预测为类别k的概率,是前m-1棵树输出结果的softmax概率,所以有:

$$[0113] \quad p_{k,m-1} = \frac{\exp(\Phi_{mk}^i)}{\sum_k \exp(\Phi_{mk}^i)}; \quad \text{公式 (5)}$$

[0114] 进一步的,结合公式(3)可以得到:

$$[0115] \quad p_{k,m-1} = \frac{\exp(\Phi_{mk}^i)}{\sum_k \exp(\Phi_{mk}^i)} = \frac{\exp\left(\sum_{n=0}^{m-1} (o_{n1}^i, o_{n2}^i, \dots, o_{nk}^i)\right)}{\sum_k \exp\left(\sum_{n=0}^{m-1} (o_{n1}^i, o_{n2}^i, \dots, o_{nk}^i)\right)}; \quad \text{公式 (6)}$$

[0116] 至此可以得知:

$$[0117] \quad l_m = \sum_{x_i \in B} \|r_m^i - o_m^i\|_2^2 = \sum_{x_i \in B} \sum_{k=1}^K (r_{mk}^i - o_{mk}^i)^2 = \sum_{x_i \in B} \sum_{k=1}^K (y_k^i - p_{k,m-1}^i - o_{mk}^i)^2; \quad \text{公式 (7)}$$

[0118] 在步骤S140中,根据所述局部损失函数计算所述Soft GBDT的全局损失函数,并根据所述全局损失函数计算所述第一内节点的第一梯度,以根据所述第一梯度对所述Soft GBDT的第一内节点的参数进行更新。

[0119] 在本示例实施例中,首先,根据所述局部损失函数计算所述Soft GBDT的全局损失函数,具体可以如下公式(7)所示:

$$[0120] \quad L = \sum_{m=1}^M l_m = \sum_{m=1}^M \sum_{x_i \in B} \|r_m^i - o_m^i\|_2^2 = \sum_{m=1}^M \sum_{x_i \in B} \sum_{k=1}^K (r_{mk}^i - o_{mk}^i)^2; \quad \text{公式(7)}$$

[0121] 然后,根据所述全局损失函数计算所述第一内节点的第一梯度,具体的可以包括:首先,根据所述全局损失函数对所述第一内节点的参数进行一阶求导,并根据所述局部损失函数对所述叶子节点的参数进行二阶求导;其次,根据所述全局损失函数的一阶导数、局部损失函数的二阶导数以及所述软决策树的输出值,计算所述软决策树的第一内节点的第一梯度。详细而言:

$$[0122] \quad \frac{\partial L}{\partial w_s^j} = \sum_{m \geq s} \sum_{x_i \in B} \sum_{k=1}^K 2(r_{mk}^i - o_{mk}^i) \left[ \frac{\partial r_{mk}^i}{\partial w_s^j} - \frac{\partial o_{mk}^i}{\partial w_s^j} \right]; \quad \text{公式(8)}$$

[0123] 其中,  $w_s^j$  表示第s棵树的第j个第一内节点的模型参数;

[0124] 进一步的,对于每棵软决策树的第一内节点的模型参数  $w_s^j$  用链式法则求导,可以得到:

$$[0125] \quad \frac{\partial r_{mk}^i}{\partial w_s^j} = - \sum_{k'} \frac{\partial^2 \text{loss}(\Phi_m^i, y^i)}{\partial \Phi_{mk}^i \partial \Phi_{mk'}^i} \frac{\partial \Phi_{mk'}^i}{\partial w_s^j}, s \leq m-1; \quad \text{公式(9)}$$

[0126] 同时,每棵树有多层,每一个内节点的权重矩阵(也即内节点的参数)是  $w_m^j$  (第m棵树的第j个内节点),最后一层的logits输出为  $\phi_m^j$ , 大小为  $K \times 1$ ,  $o_{mk}^i$  即为第m棵树的第k个分类的输出为叶子节点上第k个分类的加权和,这个加权就是每个叶子节点上的路径的概率,具体如下:

$$[0127] \quad o_{mk}^i = \sum_{j=1}^{m_j} \phi_{mk}^j P_m^j(x_i), k=1,2,\dots,K; \quad \text{公式(10)}$$

[0128] 其中,  $m_j$  为所有的路径条数,路径条数可以根据软决策树的深度计算得到;  $P_m^j(x_i)$  为每条路径的概率, j表示第j条路径, l表示第l个内节点,计算方法如下:

$$[0129] \quad P_m^j(x_i) = 1;$$

$$[0130] \quad \text{如果不是叶子节点: } P_m^l(x_i) = \text{sigmoid} \left[ \left( w_m^l \right)^T x_i \right]$$

$$[0131] \quad \text{如果是左侧叶子节点: } P_m^j(x_i) * = P_m^l(x_i);$$

$$[0132] \quad \text{如果是右侧叶子节点: } P_m^j(x_i) * = 1 - P_m^l(x_i);$$

[0133] 另外,对公式(13)用链式法则求导可以得到:

$$[0134] \quad \frac{\partial o_{mk}^i}{\partial w_m^j} = \sum_{j \in J_1} \phi_{mk}^j P_m^j (1 - p_m^j) x_i + \sum_{j \in J_2} -\phi_{mk}^j P_m^j p_m^j x_i; \quad \text{公式(11)}$$



[0135] 其中,  $J_1$  以及  $J_2$  分别表示左边以及右边经过 node\_1 的路径。记:

$$[0136] \quad u_{mk}^i = \sum_{j \in J_1} j + \sum_{j \in J_2} -\phi_{mk}^j P_m^j p_m^j; \quad \text{公式(12)}$$

[0137] 由公式 (11) 以及公式 (12) 可以得知:

$$[0138] \quad \frac{\partial o_{mk}^i}{\partial w_m^j} = u_{mk}^i x_i; \quad \text{公式 (13)}$$

$$[0139] \quad \frac{\partial o_{mk}^i}{\partial \phi_m} = P_m^i; \quad \text{公式 (14)}$$

[0140] 其中,  $P_m^i$  表示叶子节点  $\phi_m$  对应的路径的概率。

[0141] 同时, 由公式 (3) 可以得知:

$$[0142] \quad \Phi_{mk}^i = \sum_{n=0}^{m-1} o_{nk}^i;$$

$$[0143] \quad \frac{\partial \Phi_{mk}^i}{\partial w_s^j} = \frac{\partial o_{sk}^i}{\partial w_s^j};$$

[0144] 因此有:

$$[0145] \quad \frac{\partial L}{\partial w_s^j} = \sum_{x_i \in B} \left[ \sum_{m \geq s} \sum_{k=1}^K 2(r_{mk}^i - o_{mk}^i) \left[ -\sum_{k'} \frac{\partial^2 \text{loss}(\Phi_m^i, y^i)}{\partial \Phi_{mk}^i \partial \Phi_{mk'}^i} - \delta_{kk'}(ms) \right] u_{mk}^i \right] x_{a,b}^i; \quad \text{公式(15)}$$

[0146] 其中,  $\delta_{kk'}(ms) = 1$ , 当且仅当  $m=s$  且  $k=k'$  时等于 1, 否则为 0,  $x_{a,b}^i$  表示每个第一和计算结果, 也即样本特征。记:

$$[0147] \quad v_i = \sum_{m \geq s} \sum_{k=1}^K 2(r_{mk}^i - o_{mk}^i) \left[ -\sum_{k'} \frac{\partial^2 l(\Phi_m^i, y^i)}{\partial \Phi_{mk}^i \partial \Phi_{mk'}^i} - \delta_{kk'}(ms) \right] u_{mk}^i; \quad \text{公式 (16)}$$

$$[0148] \quad \text{则有: } \frac{\partial L}{\partial w_s^j} = \sum_{x_i \in B} v_i x_{a,b}^i; \quad \text{公式 (17)}$$

[0149] 至此, 即可根据公式 (17) 对第一梯度和/或第二梯度进行更新。此处需要补充说明的是, 当对第一梯度进行更新时, 可以通过  $v_i$  (第一梯度) 以及第一特征数据进行相乘, 即可根据乘积结果对第一内节点的参数进行更新。

[0150] 进一步的, 由于第一数据持有者侧还具有叶子节点的参数, 因此, 为了对叶子节点的参数进行更新, 该基于 Soft GBDT 的跨特征联邦学习方法还包括: 根据所述全局损失函数对所述叶子节点的参数进行一阶求导, 得到所述软决策树中所包括的叶子节点的叶子梯度; 根据所述叶子节点的叶子梯度, 对所述软决策树中的叶子节点的参数进行更新。详细而言:

[0151] 首先, 对于每棵树的叶子节点参数用链式法则求导, 可以得到:

$$[0152] \quad \frac{\partial r_{mk}^i}{\partial \phi_s} = -\sum_{k'} \frac{\partial^2 \text{loss}(\Phi_m^i, y^i)}{\partial \Phi_{mk}^i \partial \Phi_{mk'}^i} \frac{\partial \Phi_{mk'}^i}{\partial \phi_s} = -\sum_{k'} \frac{\partial^2 \text{loss}(\Phi_m^i, y^i)}{\partial \Phi_{mk}^i \partial \Phi_{mk'}^i} \frac{\partial o_{sk'}^i}{\partial \phi_s}; \quad \text{公式 (18)}$$

[0153] 其中, 公式 (18) 表示前  $m-1$  棵树的残差对于前  $m-1$  棵树的叶子节点的输出求导, 只

在第一数据持有者侧计算;因此,叶子节点的叶子梯度的计算方法可以如下公式(19)以及公式(20)所示:

$$[0154] \quad \frac{\partial L}{\partial \phi_s} = \sum_{m \geq s} \sum_{x_i \in B} \sum_{k=1}^K 2(r_{mk}^i - o_{mk}^i) \left[ -\sum_{k'} \frac{\partial^2 \text{loss}(\Phi_m^i, y^i)}{\partial \Phi_{mk}^i \partial \Phi_{mk'}^i} - \delta_{kk'}(ms) \right] \frac{\partial o_{sk'}^i}{\partial \phi_s}; \quad \text{公式(19)}$$

[0155] 同时,由公式(14)可以得出:

$$[0156] \quad \frac{\partial L}{\partial \phi_s} = \sum_{m \geq s} \sum_{x_i \in B} \sum_{k=1}^K 2(r_{mk}^i - o_{mk}^i) \left[ -\sum_{k'} \frac{\partial^2 \text{loss}(\Phi_m^i, y^i)}{\partial \Phi_{mk}^i \partial \Phi_{mk'}^i} - \delta_{kk'}(ms) \right] P_m^i. \quad \text{公式(20)}$$

[0157] 至此,已经完全完成了第一数据持有者侧的Soft GBDT的第一内节点的参数以及叶子节点的参数的更新。

[0158] 进一步的,为了完成本公开示例实施例所提出的跨特征联邦模型训练,本公开示例实施例还提供了另一种基于Soft GBDT的跨特征联邦学习方法,配置于多方机器学习中提供第二特征数据第二数据持有者,用于对Soft GBDT进行训练。参考图5所示,该基于Soft GBDT的跨特征联邦学习方法可以包括以下步骤:

[0159] 步骤S510,利用所述Soft GBDT计算所述第二特征数据在所述Soft GBDT中的内节点的第二线性回归部分,并将所述第二线性回归部分发送至第一数据持有者;

[0160] 步骤S520,接收所述第一数据持有者发送的加密后的第一梯度;其中,所述第一梯度是根据第一线性回归部分以及第二线性回归部分计算得到的,所述第一线性回归部分是所述第一数据持有者利用所述Soft GBDT对与所述第二特征数据具有相同数据生产者的第一特征数据计算得到的,所述第一梯度用于对所述第一数据持有者侧的所述Soft GBDT的第一内节点的参数进行更新;

[0161] 步骤S530,根据所述加密后的第一梯度以及第二特征数据计算第二梯度,并接收所述第一数据持有者发送的解密后的第二梯度;

[0162] 步骤S540,利用解密后的第二梯度对所述第二数据持有者侧的Soft GBDT的第二内节点的参数进行更新。

[0163] 以下,将对步骤S510-S540进行解释以及说明。首先,第二线性回归部分的计算过程与前文的第一线性回归部分的计算方法类似,此处不再赘述;其次,当第一数据持有者计算得到第一梯度以后,可以通过自身生成的公钥对该第一梯度进行加密,进而将加密后的第一梯度( $[v^1]$ )发送至第二数据持有者,第二数据持有者根据加密后的第一梯度以及第二特征数据计算第二梯度,具体计算方法可以如下公式(21)所示:

$$[0164] \quad [grad\_b] = \frac{\partial L}{\partial w_s^i} = v^i * x_b^i; \quad \text{公式(21)}$$

[0165] 当得到第二梯度以后,则将该第二梯度发送至第一数据持有者,第一数据持有者利用自身持有的私钥对第二梯度解密后,将解密完成的第二梯度发送至第二数据持有者,第二数据持有者即可根据解密后的第二梯度,对自身持有的第二内节点的参数进行更新。通过该方法,使得第二数据持有者在不泄露本方所持有的第二特征数据的前提下,即可根据全局损失函数对本方的第二内节点的参数进行更新,确保了第二特征数据的安全性;并且,由于整体的训练过程是分布在多方进行的,进而避免由于需要在某一方进行训练进而

导致的数据处理负担较重,训练速度慢,训练效率较低的问题,进一步的提高了模型训练效率。

[0166] 以下,结合图6对本公开示例实施例的基于Soft GBDT的跨特征联邦学习方法进行进一步的解释以及说明。参考图6所示,该基于Soft GBDT的跨特征联邦学习方法可以包括:

[0167] S601,第二数据持有者计算本方每棵树的每个节点的第二线性回归部分 $v$ ,然后传给第一数据持有者;

[0168] S602,第一数据持有者接收第二线性回归部分,计算本方对于树的节点的第一线性回归部分 $u$ ,然后计算每个节点的概率值 $p = \text{sigmoid}(u+v)$ ;

[0169] S603,第一数据持有者根据概率值计算损失函数和本方每个第一内节点的梯度以及叶子节点的梯度,并更新梯度;

[0170] S604,第一数据持有者把加密后的梯度传输给第二数据持有者;

[0171] S605,第二数据持有者根据加密后的梯度计算本方每个第二内节点的参数的梯度;

[0172] S606,第二数据持有者将第二内节点的梯度发送给第一数据持有者;

[0173] S607,第一数据持有者对第二内节点的梯度进行解密,并将解密后的第二内节点的梯度发送至第二数据持有者;

[0174] S608,第二数据持有者根据解密后的第二内节点的梯度对第二内节点的参数进行更新。

[0175] 至此,即完成了整个的跨特征联邦模型的训练过程。

[0176] 本公开示例实施例还提供了一种基于Soft GBDT的跨特征联邦预测方法,配置于多方机器学习中提供第一待预测数据的第一数据持有者,用于根据对Soft GBDT进行训练得到的数据预测模型进行数据预测。参考图7所示,该基于Soft GBDT的跨特征联邦预测方法可以包括以下步骤:

[0177] 步骤S710,利用所述数据预测模型计算所述第一待预测数据在所述数据预测模型中的内节点的第三线性回归部分;其中,所述数据预测模型是根据前述基于Soft GBDT的跨特征联邦学习方法对Soft GBDT训练得到的;

[0178] 步骤S720,接收第二数据持有者发送的利用所述数据预测模型对第二待预测数据进行计算得到的第四线性回归部分;其中,所述第一待预测数据以及第二待预测数据的数据生产者相同;

[0179] 步骤S730,对所述第三线性回归部分以及第四线性回归部分进行求和运算,并利用所述数据预测模型的叶子节点所在的归一化层对第二和运算结果进行归一化处理,得到数据预测结果。

[0180] 具体的,利用所述数据预测模型的叶子节点所在的归一化层对第二和运算结果进行归一化处理,得到数据预测结果,可以包括:首先,利用所述数据预测模型的叶子节点所在的归一化层计算所述第二和运算结果在每一个叶子节点上的分支概率;其次,根据所述分支概率计算所述数据预测结果。

[0181] 以下,将对步骤S710-步骤S730进行解释以及说明。首先,第一数据持有者以及第二数据持有者分别计算每棵树的每个节点的第三线性回归部分以及第四线性回归部分,然后第二数据持有者将第四线性回归部分发送给第一数据持有者,当第一数据持有者接收到

第四线性回归部分以后,即可计算每个节点的分支概率,然后再根据分支概率计算数据预测模型最终的预测结果,并输出该预测结果。通过该方法,可以使得第一数据持有者以及第二数据持有者对用户的风险进行预测,并且,由于加入多方所持有的用户的特征,进而可以提高预测结果的准确率,进而进一步的避免由于用户风险较高带来的经济损失。

[0182] 本公开示例实施例还提供了一种基于Soft GBDT的跨特征联邦学习装置,配置于多方机器学习中提供第一特征数据以及特征标签的第一数据持有者,用于对Soft GBDT进行训练。参考图8所示,该基于Soft GBDT的跨特征联邦学习装置可以包括第一计算模块810、第一接收模块820、第二计算模块830以及第一参数更新模块840。其中:

[0183] 第一计算模块810可以用于利用所述Soft GBDT计算所述第一特征数据在所述Soft GBDT中的第一内节点的第一线性回归部分;

[0184] 第一接收模块820可以用于接收第二数据持有者发送的利用所述Soft GBDT对第二特征数据进行计算得到的第二线性回归部分;其中,所述第一特征数据以及第二特征数据的数据生产者相同;

[0185] 第二计算模块830可以用于计算所述第一线性回归部分以及第二线性回归部分在所述Soft GBDT中的叶子节点的输出值,并根据所述叶子节点的输出值以及所述特征标签,计算所述Soft GBDT中所包括的当前软决策树的局部损失函数;

[0186] 第一参数更新模块840可以用于根据所述局部损失函数计算所述Soft GBDT的全局损失函数,并根据所述全局损失函数计算所述第一内节点的第一梯度,以根据所述第一梯度对所述Soft GBDT的第一内节点参数进行更新。

[0187] 在本公开的一种示例性实施例中,计算所述第一线性回归部分以及第二线性回归部分在所述Soft GBDT中的叶子节点的输出值,包括:

[0188] 对所述第一线性回归部分以及第二线性回归部分进行求和运算得到第一和运算结果;

[0189] 利用所述Soft GBDT的叶子节点所在的归一化层对所述第一和运算结果进行归一化处理,得到所述第一和运算结果在所述当前软决策树所在叶子节点的输出值。

[0190] 在本公开的一种示例性实施例中,根据所述叶子节点的输出值以及所述特征标签,计算所述Soft GBDT中所包括的当前软决策树的局部损失函数,包括:

[0191] 计算排在当前软决策树前面的所有软决策树在叶子节点的输出值的总和;

[0192] 对排在当前软决策树前面的所有软决策树在叶子节点的输出值的总和进行归一化运算,得到所述第一和运算结果中所包括的样本特征在排在当前软决策树前面的所有软决策树中的第一预测结果;

[0193] 根据所述第一预测结果、所述样本特征在当前软决策树中的叶子节点的输出值以及与所述样本特征对应的特征标签,构建所述当前软决策树的局部损失函数。

[0194] 在本公开的一种示例性实施例中,根据所述全局损失函数计算所述第一内节点的第一梯度,包括:

[0195] 根据所述全局损失函数对所述第一内节点参数进行一阶求导,并根据所述局部损失函数对所述叶子节点参数进行二阶求导;

[0196] 根据所述全局损失函数的一阶导数、局部损失函数的二阶导数以及所述软决策树的输出值,计算所述软决策树的第一内节点的第一梯度。

[0197] 在本公开的一种示例性实施例中,所述基于Soft GBDT的跨特征联邦学习装置还包括:

[0198] 第一参数求导模块,可以用于根据所述全局损失函数对所述叶子节点的参数进行一阶求导,得到所述软决策树中所包括的叶子节点的叶子梯度;

[0199] 第二参数更新模块,可以用于根据所述叶子节点的叶子梯度,对所述软决策树中的叶子节点的参数进行更新。

[0200] 本公开示例实施例还提供了另一种基于Soft GBDT的跨特征联邦学习装置,配置于多方机器学习中提供第二特征数据第二数据持有者,用于对Soft GBDT进行训练。参考图9所示,该基于Soft GBDT的跨特征联邦学习装置可以包括第三计算模块910、第二接收模块920、第四计算模块930以及第三参数更新模块940。其中:

[0201] 第三计算模块910可以用于利用所述Soft GBDT计算所述第二特征数据在所述Soft GBDT中的内节点的第二线性回归部分,并将所述第二线性回归部分发送至第一数据持有者;

[0202] 第二接收模块920可以用于接收所述第一数据持有者发送的加密后的第一梯度;其中,所述第一梯度是根据第一线性回归部分以及第二线性回归部分计算得到的,所述第一线性回归部分是所述第一数据持有者利用所述Soft GBDT对与所述第二特征数据具有相同数据生产者的第一特征数据计算得到的,所述第一梯度用于对所述第一数据持有者侧的所述Soft GBDT的第一内节点的参数进行更新;

[0203] 第四计算模块930可以用于根据所述加密后的第一梯度以及第二特征数据计算第二梯度,并接收所述第一数据持有者发送的解密后的第二梯度;

[0204] 第三参数更新模块940可以用于利用解密后的第二梯度对所述第二数据持有者侧的Soft GBDT的第二内节点的参数进行更新。

[0205] 本公开示例实施例还提供了一种基于Soft GBDT的跨特征联邦预测装置,配置于多方机器学习中提供第一待预测数据的第一数据持有者,用于根据对Soft GBDT进行训练得到的数据预测模型进行数据预测。参考图10所示,所述基于Soft GBDT的跨特征联邦预测装置可以包括第五计算模块1010、第三接收模块1020以及数据预测模块1030。其中:

[0206] 第五计算模块1010,可以用于利用所述数据预测模型计算所述第一待预测数据在所述数据预测模型中的内节点的第三线性回归部分;其中,所述数据预测模型是根据前述所述基于Soft GBDT的跨特征联邦学习装置对Soft GBDT训练得到的;

[0207] 第三接收模块1020,可以用于接收第二数据持有者发送的利用所述数据预测模型对第二待预测数据进行计算得到的第四线性回归部分;其中,所述第一待预测数据以及第二待预测数据的数据生产者相同;

[0208] 数据预测模块1030,可以用于对所述第三线性回归部分以及第四线性回归部分进行求和运算,并利用所述数据预测模型的叶子节点所在的归一化层对第二和运算结果进行归一化处理,得到数据预测结果。

[0209] 在本公开的一种示例性实施例中,利用所述数据预测模型的叶子节点所在的归一化层对第二和运算结果进行归一化处理,得到数据预测结果,包括:

[0210] 利用所述数据预测模型的叶子节点所在的归一化层计算所述第二和运算结果在每一个叶子节点上的分支概率;

[0211] 上述基于Soft GBDT的跨特征联邦学习装置以及基于Soft GBDT的跨特征联邦预测装置中各模块的具体细节已经在对应的基于Soft GBDT的跨特征联邦学习方法以及基于Soft GBDT的跨特征联邦预测方法中进行了详细的描述,因此此处不再赘述。

[0212] 应当注意,尽管在上文详细描述中提及了用于动作执行的设备的若干模块或者单元,但是这种划分并非强制性的。实际上,根据本公开的实施方式,上文描述的两个或更多模块或者单元的特征和功能可以在一个模块或者单元中具体化。反之,上文描述的一个模块或者单元的特征和功能可以进一步划分为由多个模块或者单元来具体化。

[0213] 此外,尽管在附图中以特定顺序描述了本公开中方法的各个步骤,但是,这并非要求或者暗示必须按照该特定顺序来执行这些步骤,或是必须执行全部所示的步骤才能实现期望的结果。附加的或备选的,可以省略某些步骤,将多个步骤合并为一个步骤执行,以及/或者将一个步骤分解为多个步骤执行等。

[0214] 在本公开的示例性实施例中,还提供了一种能够实现上述方法的电子设备。

[0215] 所属技术领域的技术人员能够理解,本公开的各个方面可以实现为系统、方法或程序产品。因此,本公开的各个方面可以具体实现为以下形式,即:完全的硬件实施方式、完全的软件实施方式(包括固件、微代码等),或硬件和软件方面结合的实施方式,这里可以统称为“电路”、“模块”或“系统”。

[0216] 下面参照图11来描述根据本公开的这种实施方式的电子设备1100。图11显示的电子设备1100仅仅是一个示例,不应对本公开实施例的功能和使用范围带来任何限制。

[0217] 如图11所示,电子设备1100以通用计算设备的形式表现。电子设备1100的组件可以包括但不限于:上述至少一个处理单元1110、上述至少一个存储单元1120、连接不同系统组件(包括存储单元1120和处理单元1110)的总线1130以及显示单元1140。

[0218] 其中,所述存储单元存储有程序代码,所述程序代码可以被所述处理单元1110执行,使得所述处理单元1110执行本说明书上述“示例性方法”部分中描述的根据本公开各种示例性实施方式的步骤。例如,所述处理单元1110可以执行如图1中所示的步骤S110:利用所述Soft GBDT计算所述第一特征数据在所述Soft GBDT中的第一内节点的第一线性回归部分;步骤S120:接收第二数据持有者发送的利用所述Soft GBDT对第二特征数据进行计算得到的第二线性回归部分;步骤S130:计算所述第一线性回归部分以及第二线性回归部分在所述Soft GBDT中的叶子节点的输出值,并根据所述叶子节点的输出值以及所述特征标签,计算所述Soft GBDT中所包括的当前软决策树的局部损失函数;步骤S140:根据所述局部损失函数计算所述Soft GBDT的全局损失函数,并根据所述全局损失函数计算所述第一内节点的第一梯度,以根据所述第一梯度对所述Soft GBDT的第一内节点的参数进行更新。

[0219] 所述处理单元1110可以执行如图5中所示的步骤S510:利用所述Soft GBDT计算所述第二特征数据在所述Soft GBDT中的内节点的第二线性回归部分,并将所述第二线性回归部分发送至第一数据持有者;步骤S520:接收所述第一数据持有者发送的加密后的第一梯度;其中,所述第一梯度是根据第一线性回归部分以及第二线性回归部分计算得到的,所述第一线性回归部分是所述第一数据持有者利用所述Soft GBDT对与所述第二特征数据具有相同数据生产者的第一特征数据计算得到的,所述第一梯度用于对所述第一数据持有者侧的所述Soft GBDT的第一内节点的参数进行更新;步骤S530:根据所述加密后的第一梯度以及第二特征数据计算第二梯度,并接收所述第一数据持有者发送的解密后的第二梯度;



步骤S540:利用解密后的第二梯度对所述第二数据持有者侧的Soft GBDT的第二内节点参数进行更新。

[0220] 所述处理单元1110可以执行如图6中所示的步骤S610:利用所述数据预测模型计算所述第一待预测数据在所述数据预测模型中的内节点的第三线性回归部分;其中,所述数据预测模型是根据前述所述基于Soft GBDT的跨特征联邦学习方法对Soft GBDT训练得到的;步骤S620:接收第二数据持有者发送的利用所述数据预测模型对第二待预测数据进行计算得到的第三线性回归部分;步骤S630:对所述第三线性回归部分以及第四线性回归部分进行求和运算,并利用所述数据预测模型的叶子节点所在的归一化层对第二和运算结果进行归一化处理,得到数据预测结果。

[0221] 存储单元1120可以包括易失性存储单元形式的可读介质,例如随机存取存储单元(RAM) 11201和/或高速缓存存储单元11202,还可以进一步包括只读存储单元(ROM) 11203。

[0222] 存储单元1120还可以包括具有一组(至少一个)程序模块11205的程序/实用工具11204,这样的程序模块11205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。

[0223] 总线1130可以为表示几类总线结构中的一种或多种,包括存储单元总线或者存储单元控制器、外围总线、图形加速端口、处理单元或者使用多种总线结构中的任意总线结构的局域总线。

[0224] 电子设备1100也可以与一个或多个外部设备1200(例如键盘、指向设备、蓝牙设备等)通信,还可与一个或者多个使得用户能与该电子设备1100交互的设备通信,和/或与使得该电子设备1100能与一个或多个其它计算设备进行通信的任何设备(例如路由器、调制解调器等等)通信。这种通信可以通过输入/输出(I/O)接口1150进行。并且,电子设备1100还可以通过网络适配器1160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。如图所示,网络适配器1160通过总线1130与电子设备1100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1100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、冗余处理单元、外部磁盘驱动阵列、RAID系统、磁带驱动器以及数据备份存储系统等。

[0225] 通过以上的实施方式的描述,本领域的技术人员易于理解,这里描述的示例实施方式可以通过软件实现,也可以通过软件结合必要的硬件的方式来实现。因此,根据本公开实施方式的技术方案可以以软件产品的形式体现出来,该软件产品可以存储在一个非易失性存储介质(可以是CD-ROM,U盘,移动硬盘等)中或网络上,包括若干指令以使得一台计算设备(可以是个人计算机、服务器、终端装置、或者网络设备等)执行根据本公开实施方式的方法。

[0226] 在本公开的示例性实施例中,还提供了一种计算机可读存储介质,其上存储有能够实现本说明书上述方法的程序产品。在一些可能的实施方式中,本公开的各个方面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本说明书上述“示例性方法”部分中描述的根据本公开各种示例性实施方式步骤。

[0227] 根据本公开的实施方式的用于实现上述方法的程序产品,其可以采用便携式紧凑盘只读存储器(CD-ROM)并包括程序代码,并可以在终端设备,例如个人电脑上运行。然而,

本公开的程序产品不限于此,在本文件中,可读存储介质可以是任何包含或存储程序的有形介质,该程序可以被指令执行系统、装置或者器件使用或者与其结合使用。

[0228] 所述程序产品可以采用一个或多个可读介质的任意组合。可读介质可以是可读信号介质或者可读存储介质。可读存储介质例如可以为但不限于电、磁、光、电磁、红外线、或半导体的系统、装置或器件,或者任意以上的组合。可读存储介质的更具体的例子(非穷举的列表)包括:具有一个或多个导线的电连接、便携式盘、硬盘、随机存取存储器(RAM)、只读存储器(ROM)、可擦式可编程只读存储器(EPROM或闪存)、光纤、便携式紧凑盘只读存储器(CD-ROM)、光存储器件、磁存储器件、或者上述的任意合适的组合。

[0229] 计算机可读信号介质可以包括在基带中或者作为载波一部分传播的数据信号,其中承载了可读程序代码。这种传播的数据信号可以采用多种形式,包括但不限于电磁信号、光信号或上述的任意合适的组合。可读信号介质还可以是可读存储介质以外的任何可读介质,该可读介质可以发送、传播或者传输用于由指令执行系统、装置或者器件使用或者与其结合使用的程序。

[0230] 可读介质上包含的程序代码可以用任何适当的介质传输,包括但不限于无线、有线、光缆、RF等等,或者上述的任意合适的组合。

[0231] 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公开操作的程序代码,所述程序设计语言包括面向对象的程序设计语言—诸如Java、C++等,还包括常规的过程式程序设计语言—诸如“C”语言或类似的设计语言。程序代码可以完全地在用户计算设备上执行、部分地在用户设备上执行、作为一个独立的软件包执行、部分在用户计算设备上部分在远程计算设备上执行、或者完全在远程计算设备或服务器上执行。在涉及远程计算设备的情形中,远程计算设备可以通过任意种类的网络,包括局域网(LAN)或广域网(WAN),连接到用户计算设备,或者,可以连接到外部计算设备(例如利用因特网服务提供商来通过因特网连接)。

[0232] 此外,上述附图仅是根据本公开示例性实施例的方法所包括的处理的示意性说明,而不是限制目的。易于理解,上述附图所示的处理并不表明或限制这些处理的时间顺序。另外,也易于理解,这些处理可以是例如在多个模块中同步或异步执行的。

[0233] 本领域技术人员在考虑说明书及实践这里发明的发明后,将容易想到本公开的其他实施例。本申请旨在涵盖本公开的任何变型、用途或者适应性变化,这些变型、用途或者适应性变化遵循本公开的一般性原理并包括本公开未发明的本技术领域中的公知常识或惯用技术手段。说明书和实施例仅被视为示例性的,本公开的真正范围和精神由权利要求指出。

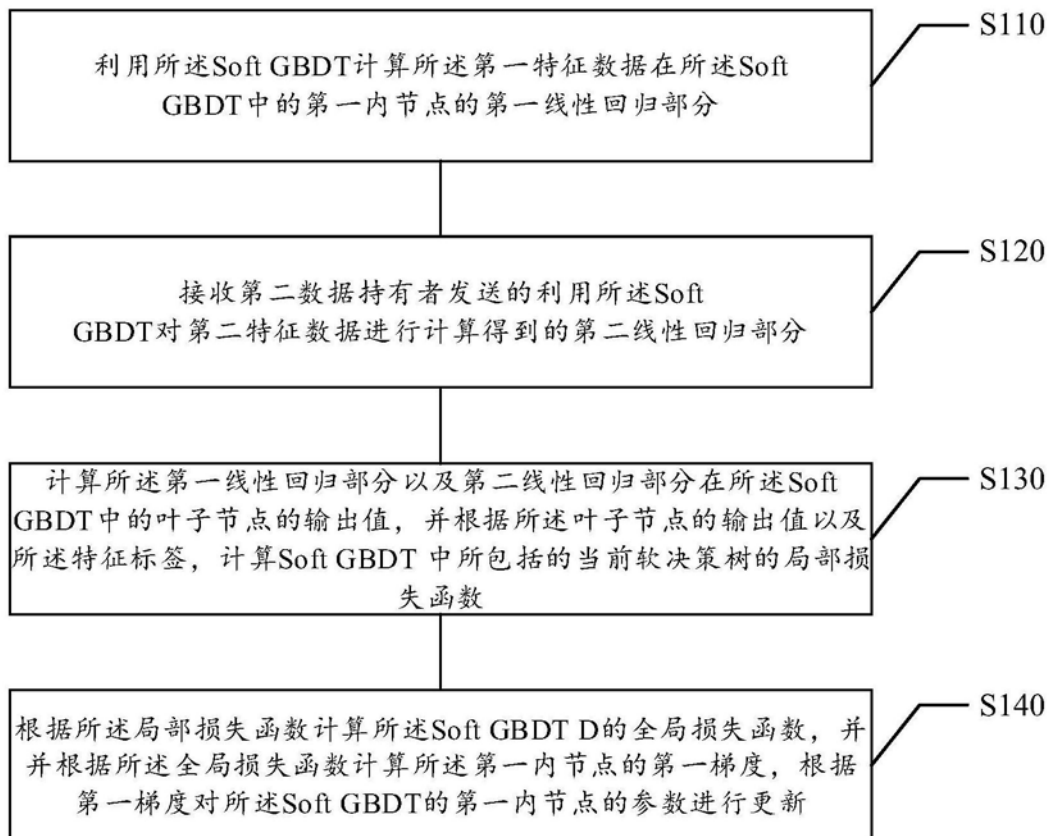


图1

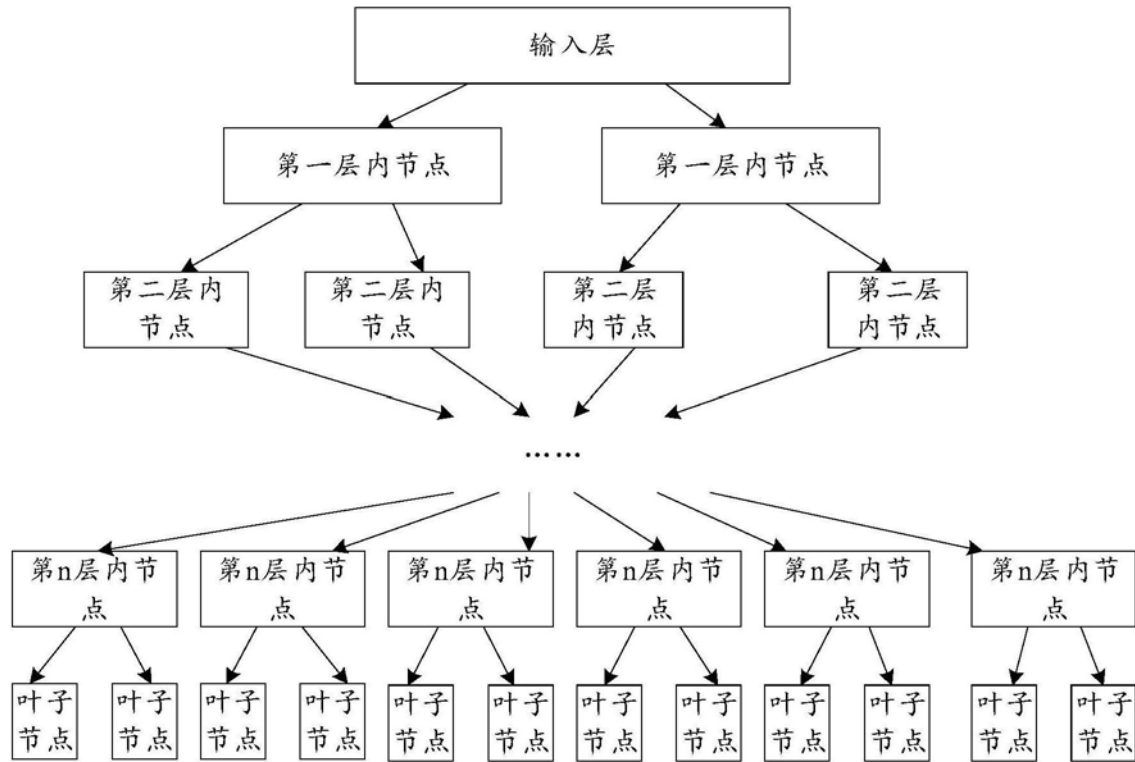


图2

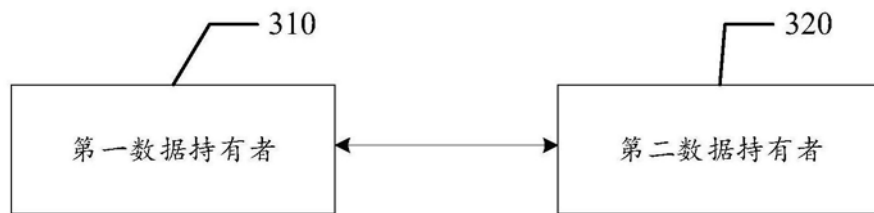


图3

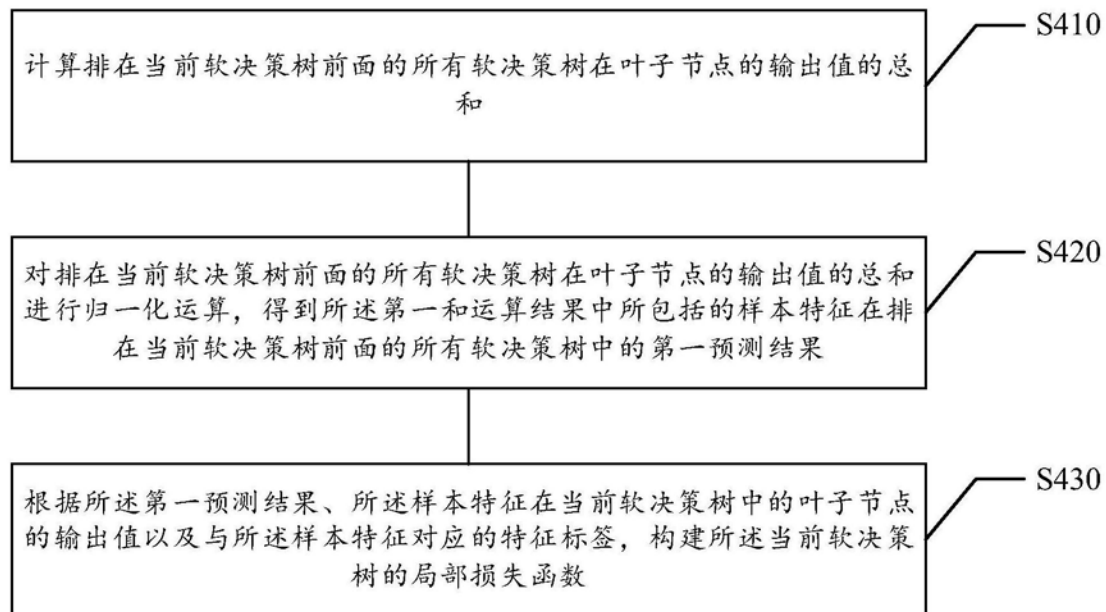


图4

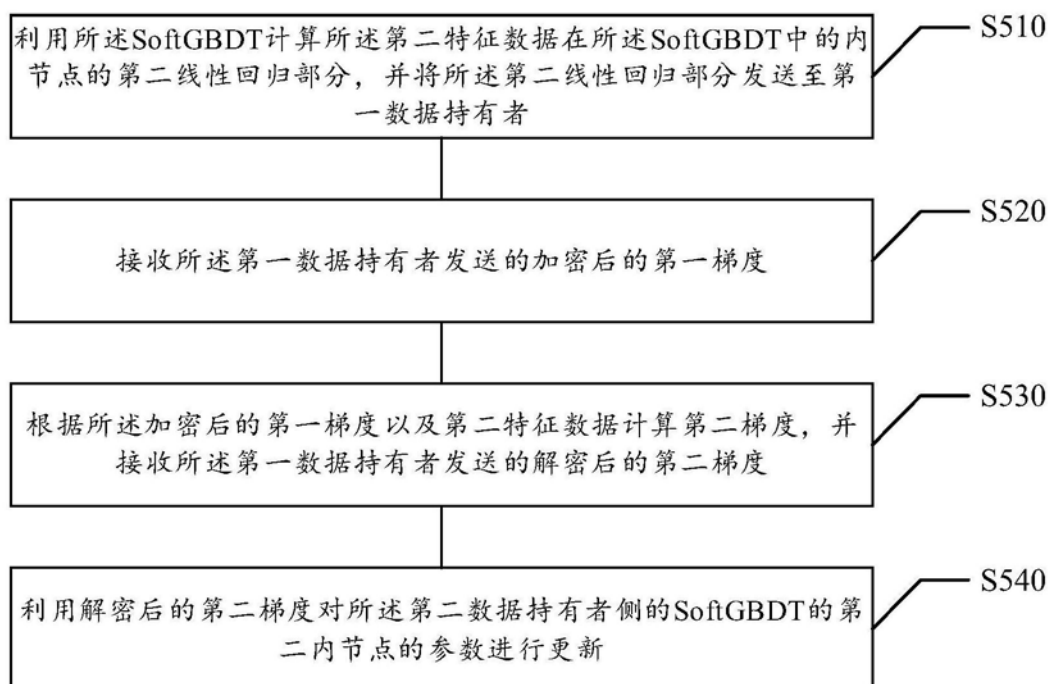


图5

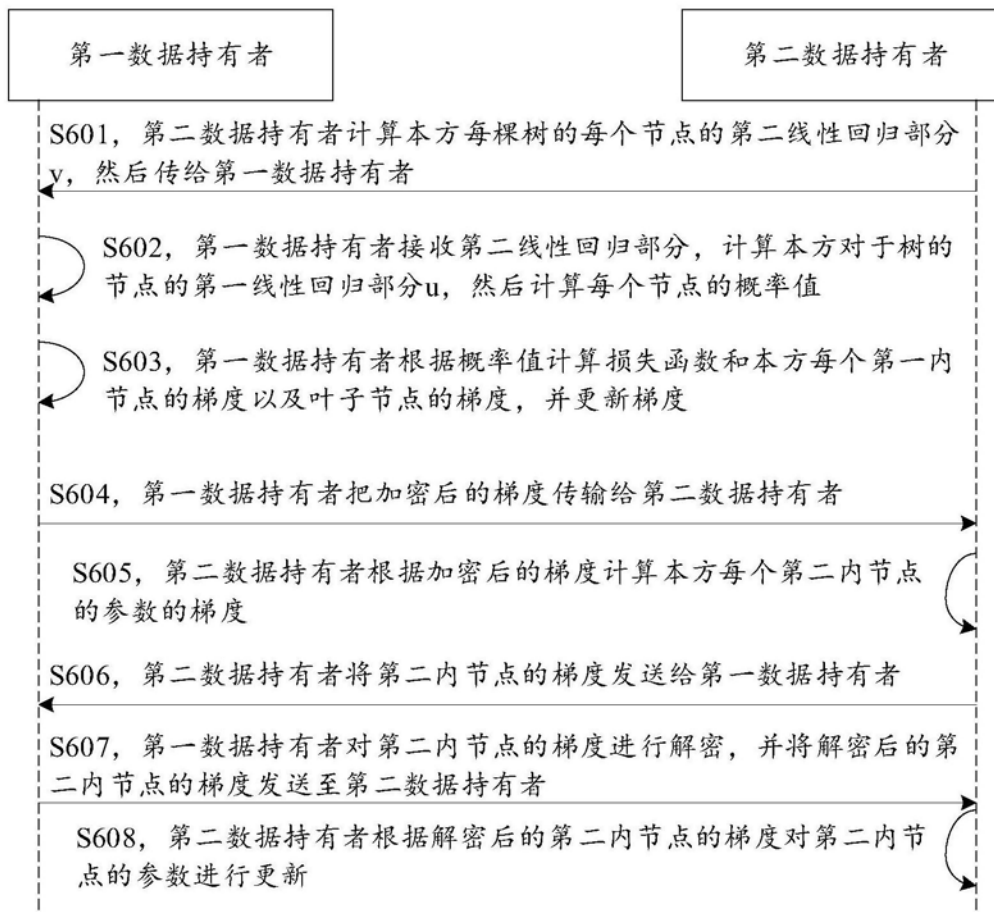


图6

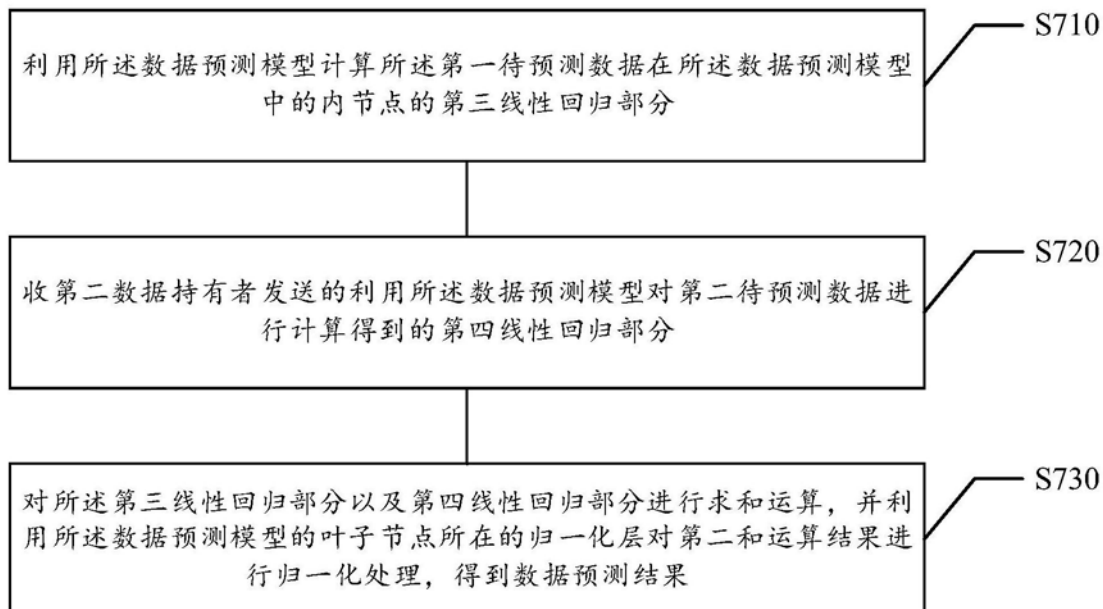


图7



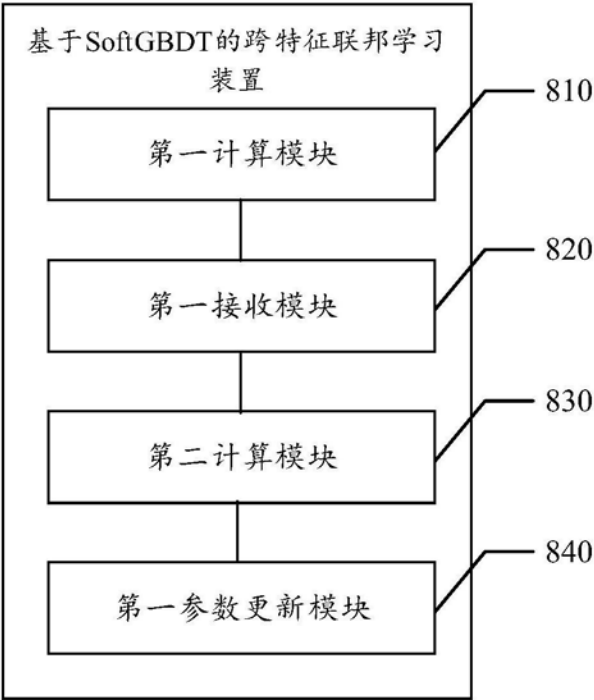


图8

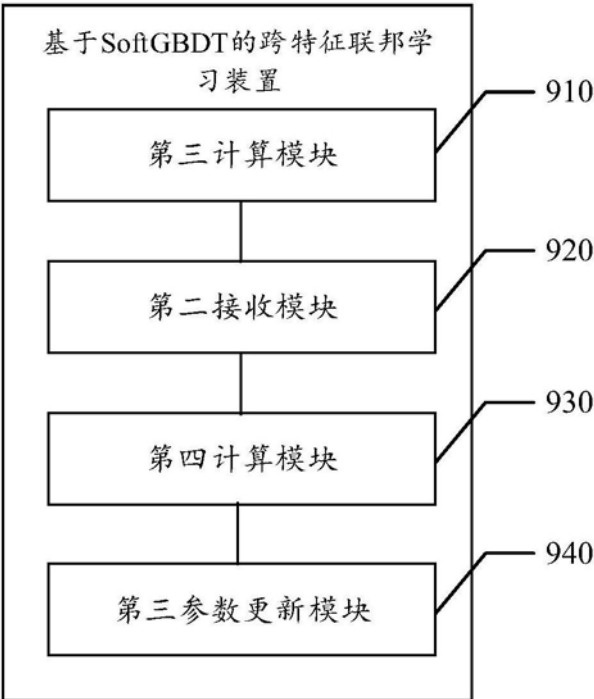


图9

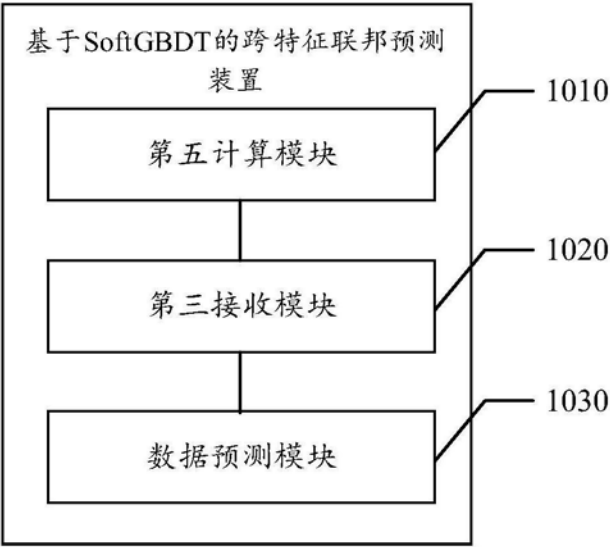


图10

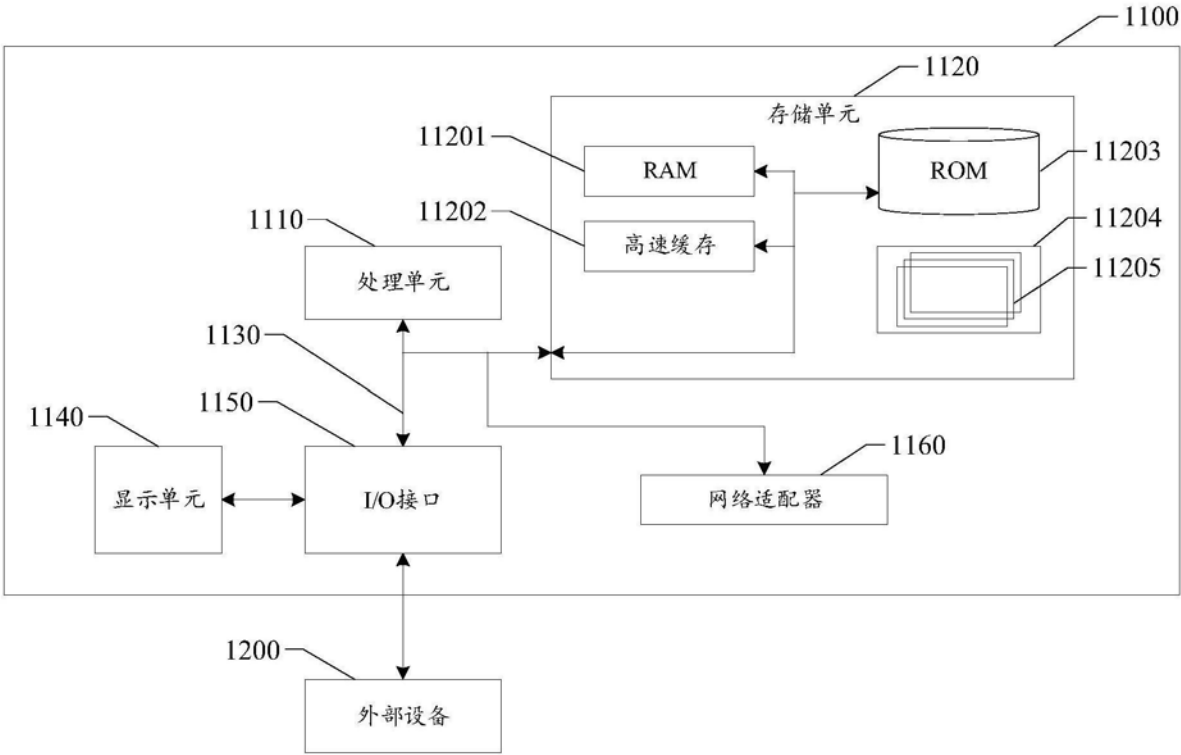


图11