



(12) 发明专利申请

(10) 申请公布号 CN 112580817 A

(43) 申请公布日 2021. 03. 30

(21) 申请号 202011038278.X

(22) 申请日 2020.09.28

(30) 优先权数据

16/587,713 2019.09.30 US

(71) 申请人 脸谱公司

地址 美国加利福尼亚州

(72) 发明人 贾宏钟 杰伊·帕瑞克

(74) 专利代理机构 北京安信方达知识产权代理有限公司 11262

代理人 周靖 杨明钊

(51) Int. Cl.

G06N 20/00 (2019.01)

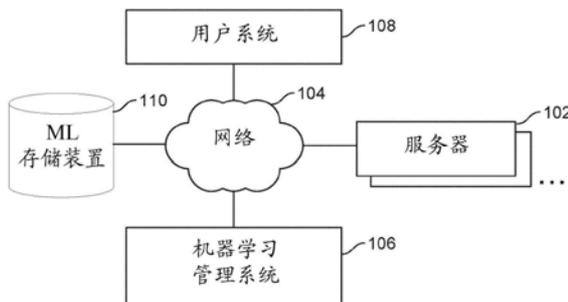
权利要求书2页 说明书8页 附图5页

(54) 发明名称

管理机器学习特征

(57) 摘要

本申请涉及管理机器学习特征。机器学习模型被训练。为机器学习模型的多个机器学习特征中的每个机器学习特征确定特征重要性度量。基于特征重要性度量,管理机器学习模型的多个机器学习特征中的一个或多个机器学习特征。



1. 一种方法,包括:
训练机器学习模型;
为所述机器学习模型的多个机器学习特征中的每个机器学习特征确定特征重要性度量;和
基于所述特征重要性度量,管理所述机器学习模型的所述多个机器学习特征中的一个或更多个机器学习特征。
2. 根据权利要求1所述的方法,其中管理所述一个或更多个机器学习特征包括基于所述特征重要性度量生成所述机器学习模型的新版本。
3. 根据权利要求1所述的方法,其中管理所述一个或更多个机器学习特征包括基于确定所述一个或更多个机器学习特征中的一个机器学习特征的相关特征重要性度量不满足阈值来确定移除所述一个机器学习特征。
4. 根据权利要求1所述的方法,其中管理所述一个或更多个机器学习特征包括重新训练所述机器学习模型,以从所述机器学习模型中移除所述一个或更多个机器学习特征中的至少一个特征。
5. 根据权利要求4所述的方法,其中管理所述一个或更多个机器学习特征包括自动删除所移除的至少一个特征的存储数据。
6. 根据权利要求4所述的方法,其中管理所述一个或更多个机器学习特征包括自动使所移除的至少一个特征的数据不再被收集。
7. 根据权利要求1所述的方法,其中管理所述一个或更多个机器学习特征包括修改所述一个或更多个机器学习特征中的至少一个特征,并且使用所述一个或更多个机器学习特征中的经修改的至少一个特征来重新训练所述机器学习模型。
8. 根据权利要求1所述的方法,其中管理所述一个或更多个机器学习特征包括基于所述一个或更多个机器学习特征中的至少一个而生成新特征,并使用所述新特征重新训练所述机器学习模型。
9. 根据权利要求1所述的方法,其中确定所述机器学习特征中的选择的机器学习特征的所述特征重要性度量包括将针对选择的测试数据集的所述机器学习模型的基本性能与针对所述选择的机器学习特征的所述选择的测试数据集的修改版本的所述机器学习模型的新性能进行比较。
10. 根据权利要求9所述的方法,其中所述选择的机器学习特征的所述选择的测试数据集的修改版本至少部分地通过以下方式生成:使用基于机器学习模型的属性而选择的修改方法来修改对应于所述选择的机器学习特征的值。
11. 根据权利要求9所述的方法,其中所述选择的机器学习特征的所述选择的测试数据集的修改版本至少部分地通过使对应于所述选择的机器学习特征的值随机化来生成。
12. 根据权利要求1所述的方法,其中每个特征重要性度量包括表示相应机器学习特征对所述机器学习模型的推理结果的贡献量的数值。
13. 根据权利要求1所述的方法,其中每个特征重要性度量包括表示相应机器学习特征相比于所述机器学习特征中的其他机器学习特征的排名顺序的数值。
14. 根据权利要求1所述的方法,其中管理所述机器学习模型的所述一个或更多个机器学习特征包括一起管理多个不同机器学习模型的机器学习特征。

15. 根据权利要求14所述的方法,其中一起管理所述多个不同机器学习模型的机器学习特征包括确定所述多个不同机器学习模型的机器学习特征的排序。

16. 根据权利要求14所述的方法,其中一起管理所述多个不同机器学习模型的机器学习特征包括识别所述多个不同机器学习模型之间的特征共享。

17. 根据权利要求16所述的方法,其中识别所述多个不同机器学习模型之间的特征共享包括识别相同共享特征的多个不同特征重要性度量。

18. 根据权利要求16所述的方法,其中一起管理所述多个不同机器学习模型的机器学习特征包括基于识别出的在多个不同机器学习模型中的对所述不同机器学习模型中的一个机器学习模型的选择的机器学习特征的共享来确定是否移除所述选择的机器学习特征。

19. 一种系统,包括:

处理器,所述处理器被配置为:

训练机器学习模型;

为所述机器学习模型的多个机器学习特征中的每个机器学习特征确定特征重要性度量;和

基于所述特征重要性度量,管理所述机器学习模型的所述多个机器学习特征中的一个或多个机器学习特征;和

存储器,所述存储器耦合到所述处理器并被配置为向所述处理器提供指令。

20. 一种计算机程序产品,所述计算机程序产品包含在有形的计算机可读存储介质中,并且包括用于以下操作的计算机指令:

训练机器学习模型;

为所述机器学习模型的多个机器学习特征中的每个机器学习特征确定特征重要性度量;和

基于所述特征重要性度量,管理所述机器学习模型的所述多个机器学习特征中的一个或多个机器学习特征。

管理机器学习特征

[0001] 发明背景

[0002] 机器学习允许预测和决策基于从训练数据中自动学习的模式。使用训练数据构建的机器学习模型的精度受到可用训练数据的种类和数量的严重影响。用作构建模型的输入的训练数据的方面通常被称为模型的机器学习特征。为了试图提高精度，一直在努力增加模型中所利用的特征数量。发现拥有数万或数十万个特征的模型并不罕见。然而，随着特征数量的增加，训练和利用模型所需的存储和处理的量也在增加。因此，训练的效率、稳定性和可靠性都变得难以管理。

[0003] 附图简述

[0004] 在以下详细描述和附图中公开了本发明的各种实施例。

[0005] 图1是示出用于管理机器学习特征的系统环境的实施例的框图。

[0006] 图2是示出用于训练和部署机器学习模型的过程的实施例的流程图。

[0007] 图3是示出用于训练机器学习模型并管理其特征的过程的实施例的流程图。

[0008] 图4是示出用于确定特征重要性度量的过程的实施例的流程图。

[0009] 图5是示出用于自动管理机器学习特征的过程的实施例的流程图。

[0010] 详细描述

[0011] 本发明可以以多种方式实现，包括作为过程；装置；系统；物质的组成；体现在计算机可读存储介质上的计算机程序产品；和/或处理器，例如被配置为执行存储在耦合到处理器的存储器上和/或由该存储器提供的指令的处理器。在本说明书中，这些实现或者本发明可以采取的任何其他形式可以被称为技术。通常，在本发明的范围内，可以改变所公开的过程的步骤顺序。除非另有说明，否则被描述为被配置为执行任务的诸如处理器或存储器的组件可以被实现为在给定时间被临时配置为执行任务的通用组件或者被制造为执行任务的特定组件。如本文所使用的，术语“处理器”指的是被配置成处理数据（例如计算机程序指令）的一个或更多个设备、电路和/或处理核心。

[0012] 下面提供了本发明的一个或更多个实施例的详细描述连同说明本发明原理的附图。结合这些实施例描述了本发明，但是本发明不限于任何实施例。本发明的范围仅由权利要求限定，并且本发明包括许多替代、修改和等同物。为了提供对本发明的全面理解，在以下描述中阐述了许多具体细节。这些细节是出于示例的目的而提供的，并且本发明可以根据权利要求来被实施，而不需要这些具体细节中的一些或全部。为了清楚起见，没有详细描述与本发明相关的技术领域中的已知的技术材料，以便不会不必要地模糊本发明。

[0013] 虽然机器学习特征数量的增加通常会导致机器学习模型性能的提高，但是由于新特征而导致的性能提高可能太小而不足以证明相关收集、训练和处理的成本是合理的。更糟糕的是，它甚至可能产生伴随特征过拟合 (feature overfitting) 的相反的效果。

[0014] 在一些实施例中，训练机器学习模型。例如，使用收集的对应于机器学习模型的特征的数据，训练机器学习模型。确定机器学习模型的多个 (a plurality of) 机器学习特征中的每个机器学习特征的重要性度量。例如，在机器学习训练过程期间评估特征，以确定特征的值和模型的性能之间是否存在相关性。基于重要性度量，管理机器学习模型的多个机

器学习特征中的一个或多个机器学习特征。例如,如果确定特征与模型精度的相关性不足,则该特征将被移除。这自动导致不再收集该特征的数据,并且当模型被重新训练时,该特征被移除。这允许生成计算效率更高的模型(例如,部署模型所需的处理和存储更少),以及减少需要收集和存储的数据量。在各种实施例中,特征的管理可以考虑多个不同机器学习模型的机器学习特征的重要性度量以及不同机器学习模型之间的特征共享。例如,虽然某一特征对于一个模型来说并不重要,但是对于另一个模型来说是重要的,该特征可以被保留。

[0015] 图1是示出用于管理机器学习特征的系统环境的实施例的框图。

[0016] 服务器102可以包括一个或多个计算、存储、web、应用和/或其他处理服务器。服务器102可以位于一个或多个不同的数据中心。可以指示服务器102收集训练数据、执行机器学习训练和/或使用经训练的机器学习模型执行推理。由服务器102执行的处理工作的至少一部分可以包括托管和/或执行与终端用户(end-user)请求相关联的处理(例如,为社交网络服务向终端用户的网页或应用提供所请求的数据/内容)。在一些实施例中,一些服务器102在低流量时期或未被充分利用时被用于机器学习。例如,一些服务器102被临时重新用于应对机器学习训练。如果再次需要这些服务器来应对与终端用户请求相关的处理,则可以将这些服务器中的一个或多个返回到可用服务器的池中。

[0017] 在一些实施例中,在机器学习模型的机器学习过程之后或期间,服务器102确定模型的一个或多个机器学习特征的重要性的量度(measure)。例如,在机器学习训练过程中评估特征,以确定特征的值和模型的性能之间是否存在相关性。重要性的量度可以被提供给机器学习管理系统106。

[0018] 机器学习管理系统106包括被配置成协调和管理机器学习和机器学习特征的一个或多个服务器/计算机。例如,机器学习管理系统106在服务器102中的选定的服务器上启动机器学习训练。机器学习管理系统106也可以管理机器学习特征。例如,机器学习管理系统106管理为机器学习训练收集哪些数据和收集的数据的量,以及由各种机器学习模型使用哪些机器学习特征。在一个示例中,如果确定特征与模型精度的相关性不足,则将从模型中移除该特征,并且不再收集该特征的数据。在各种实施例中,特征的管理可以考虑多个不同机器学习模型的机器学习特征的重要性度量以及不同机器学习模型之间的特征共享。例如,虽然某一特征对于一个模型来说并不重要,但是对于另一个模型来说是重要的,该特征可以被保留。

[0019] 在一些实施例中,机器学习存储装置(storage)110存储机器学习模型和相关联数据的储存库(repository)。例如,对于每个执行的机器学习模型构建/训练,存储训练进度、训练数据、训练状态、参数、元数据、结果、结果模型和/或其他相关信息。该储存库是可搜索的,并且允许用户识别匹配搜索查询的机器学习模型。在一些实施例中,机器学习存储装置110存储与机器学习模型的机器学习特征相关联的数据。例如,为各种机器学习模型存储特征和相关联度量(例如,每个特征的重要性的量度)的列表。机器学习管理系统106可以使用来自存储装置110的数据来管理各种机器学习模型的特征。

[0020] 在一些实施例中,用户利用用户系统108来访问由一个或多个服务器102提供的服务。例如,由服务器102提供的社交网络服务由用户使用用户系统108来访问。服务器102还可以使用一个或多个经过训练的机器学习模型来执行推理,以提供用户经由用户系统

108访问的服务。用户系统108的示例包括个人计算机、膝上型计算机、平板计算机、移动设备、显示设备、用户输入设备和任何其他计算设备。

[0021] 尽管为了简化图表,只显示了有限数量的组件实例,但是图1中所示的任何组件的附加实例都可能存在。图1中未显示的组件也可能存在。所示的组件通过网络104相互通信。网络104的示例包括下列项中的一个或更多个:直接或间接物理通信连接、移动通信网络、互联网、内联网、局域网、广域网、存储区域网以及将两个或更多个系统、组件或存储设备连接在一起的任何其他形式。

[0022] 图2是示出用于训练和部署机器学习模型的过程的实施例的流程图。图2的过程可以至少部分地在服务器102和/或机器学习管理系统106中的一个或更多个上实现。

[0023] 在202,接收要收集的数据的规范。在一些实施例中,指定了要记录/存储的数据的类型和数量。例如,在可与所提供的数字服务相关联地存储的大量且几乎无穷无尽的数据中,要被捕获(例如,由图1的服务器102捕获)和被存储(例如,存储在服务器102的存储装置或存储装置110中)的数据的规范(例如,哪些数据、何时捕获数据、要捕获的数据量、数据的保留期等)被接收。在一些实施例中,收集的数据的至少一部分将被用于训练机器学习模型,并且要收集的数据的标识指定训练机器学习模型所需的训练数据。例如,接收期望被训练的机器学习模型的规范,并且该规范识别将被用来训练机器学习模型的训练数据。使用该规范,自动识别为产生期望的训练数据而收集的数据。因为收集、处理和存储数据会消耗大量的存储和计算资源,所以希望只收集将是有用的数据。如果需要或要求发生变化,希望能够动态修改要收集的数据。通过能够指定要收集的数据,它以允许计算机资源利用的效率和优化的方式提供了对数据收集的动态控制。

[0024] 在204,收集指定的数据。收集指定的数据包括记录/存储在202中识别的数据。例如,当要收集的数据被检测、提供和/或生成时(例如,在终端用户利用社交网络服务期间),该数据被存储在存储装置中。所收集的数据可以以允许有效识别和检索以供以后使用(例如,训练机器学习模型)的方式被标记和/或组织。在一些实施例中,数据由提供和/或执行与终端用户访问的服务相关联的处理的一个或更多个生产服务器(例如,图1的一个或更多个服务器102)收集。

[0025] 在206,为待训练的机器学习模型选择机器学习特征。每个机器学习特征是机器学习模型所基于的属性(property)、特性(characteristic)、数据成分、变量或性质(attribute)。例如,训练数据中的机器学习特征充当用于训练机器学习模型的训练数据的方面,并且生产数据的机器学习特征充当生产部署的机器学习模型的输入。在一个示例中,在训练数据表中,行表示每个数据记录,而列表示数据记录的不同数据字段,被用作训练机器学习模型的输入的所选列是模型的机器学习特征。执行机器学习特征选择可以包括接收将被用于被指定要被训练的机器学习模型的机器学习特征的规范。这些机器学习特征是在204中收集的用于训练机器学习模型的数据的方面。例如,在204中收集的数据的可用特征中选择机器学习特征(例如,自动选择和/或手动选择)。在一些实施例中,所接收的要收集的数据的规范至少部分基于要训练的一个或更多个机器学习模型的已识别的机器学习特征。

[0026] 在一些实施例中,至少一个选定的机器学习特征被标记为生成的机器学习特征,和/或至少一个选定的机器学习特征被标记为基本机器学习特征。基本机器学习特征可在

收集的数据和/或提供的训练数据(例如,在204中收集的数据的至少一部分)中直接获得,并且生成的机器学习特征基于一个或多个基本机器学习特征或另一个生成的机器学习特征而生成。例如,基本机器学习特征的值被处理、修改和/或与另一特征组合以生成生成的机器学习特征。标记哪个特征是基本机器学习特征以及哪个特征是生成的机器学习特征以通知未来的数据收集决策可以是重要的。例如,如果在未来模型中不再使用一基本机器学习特征,则如果没有其他被使用的生成的机器学习特征基于该基本机器学习特征,那么该基本机器学习特征的数据可以被识别为不再收集。然而,如果生成的机器学习特征将在未来模型中不再被使用,则相关联的基本机器学习特征可以被识别并被单独分析以确定它是否仍然需要被收集(例如,如果基本机器学习特征没有被直接使用并且没有其他生成的机器学习特征基于该基本机器学习特征,则该基本机器学习特征可以被识别为不再被收集)。

[0027] 在208,基于所选择的机器学习特征来训练机器学习模型。例如,对应于所选择的机器学习特征的训练数据的部分被用来训练机器学习模型(例如,训练卷积神经网络)。在一些实施例中,使用图1的一个或多个服务器102来训练机器学习模型。

[0028] 在一些实施例中,在训练过程中,评估一个或多个所选择的特征以确定每个特征的重要性。例如,为一个或多个机器学习特征中的每一个计算重要性度量,该重要性度量识别特征对机器学习模型的精度/性能的贡献的量度。使用该评估,可以管理一个或多个机器学习模型的特征。例如,可以丢弃不太重要的特征,并且可以在没有丢弃的特征的情况下重新训练机器学习模型,并且丢弃的特征的数据可以被识别为不再被收集。这允许生成计算效率更高的模型(例如,部署模型所需的处理和存储更少),并且减少了需要收集和存储的数据量。在各种实施例中,特征的管理可以考虑跨多个不同机器学习模型的机器学习特征的重要性度量。例如,虽然一个特征对于一个模型来说并不重要,但是对于另一个模型来说是重要的,该特征可以被保留。

[0029] 在210,部署经训练的机器学习模型。例如,经训练的机器学习模型被部署用于服务器中的生产使用,以执行与提供给终端用户的服务(例如,社交网络服务)相关联的推理工作。可以在部署期间收集额外的训练数据,并将其用于重新训练机器学习模型。例如,可以重复图2的过程。在图2的过程的后续迭代期间,可以基于在模型训练期间在208中确定的机器学习特征的重要性来修改在202中要收集的数据的规范。例如,被确定为不满足重要性阈值的机器学习特征将不再被利用,并且这些被丢弃的特征的数据将不再被识别为要被收集。

[0030] 图3是示出用于训练机器学习模型并管理其特征的过程的实施例的流程图。图3的过程可以至少部分地在服务器102和/或机器学习管理系统106中的一个或多个上实现。在一些实施例中,图3的过程的至少一部分被包括在图2的208中。

[0031] 在302,训练机器学习模型。在一些实施例中,机器学习模型是基于所选的机器学习特征在图2的208中训练的机器学习模型。例如,对应于所选机器学习特征的训练数据的部分被用来训练机器学习模型。在一些实施例中,使用图1的一个或多个服务器102来训练机器学习模型。

[0032] 在304,确定机器学习模型的多个机器学习特征的重要性度量。每个重要性度量可以包括指示相应机器学习特征的重要性的值。例如,为模型的每个机器学习特征计算重要

性度量值,该重要性度量值是特征对机器学习模型的精度/性能的贡献的度量。在一些实施例中,确定重要性度量包括确定模型的机器学习特征之间的相对重要性。例如,确定基于特征对模型的结果的重要性的排序的特征列表。

[0033] 特定特征的重要性度量可以通过将模型针对测试数据集的基本性能与模型针对测试数据集的修改版本的新性能进行比较(例如,确定二者之间的差异)来确定,该测试数据集的修改版本具有指定特征的交替值(例如,翻转值(flipped value)、随机化值、零值、移除值等)。如果修改后的性能测试数据集与原始性能测试数据集相比,性能差异较大且更差,则特定特征是重要的;然而,如果修改后的性能测试数据集与原始性能测试数据集相比,性能差异很小或更优,则特定特征就不那么重要了。因为特征重要性评估可能是一个计算成本非常高的过程,所以可以将测试数据集选择成在大小方面受到限制,以降低特征重要性评估的计算成本,同时保持特征重要性评估的精度。例如,基于可用测试数据集的数量,只有满足大小/数量标准的所有可用测试数据集的一部分被选择用于特征重要性评估(例如,仅利用最后N天的数据,仅加载x%的数据,等等)。

[0034] 在306,至少部分基于重要性度量,管理一个或多个机器学习特征。

[0035] 在各种实施例中,管理机器学习特征包括执行以下一项或更多项:从模型中移除/丢弃特征、修改特征、删除对应于特征的数据、使对应于特征的数据不再被收集、为对应于特征的数据选择存储层(tier)、或基于管理的特征生成新的特征。例如,如果机器学习特征的重要性度量低于阈值,则机器学习特征将被从模型中移除(例如,模型被重新训练以移除特征),并且对应于该机器学习特征的现有收集数据被自动删除,并且对应于该机器学习特征的未来数据被自动地不再收集(例如,在图2的202中被指定为不再收集)。该阈值可以基于下列项中的一个或多个来动态确定:其他特征的重要性度量(例如,基于平均值)、其他特征的总数、用于存储特征数据的存储量、可用的存储资源量、可用的处理资源量或利用该特征的机器学习模型的重要性(例如,排名、类别等)。在另一个示例中,如果机器学习特征的重要性度量低于阈值,则机器学习特征被修改或变换,以试图产生该特征的一个版本来提高模型的性能(例如,使用修改的特征重新训练模型)。在另一示例中,如果机器学习特征的重要性度量高于阈值,则机器学习特征被用于生成新特征(例如,与其他机器学习特征相结合以生成用于重新训练模型的新特征)。

[0036] 在一些实施例中,对应于机器学习特征的存储数据的存储层是基于其重要性度量来确定的。例如,对应于重要性度量在第一值范围内(例如,高重要性)的特征的数据被存储在高性能存储装置(例如,固态驱动器)中;对应于重要性度量在第二值范围内(例如,中等重要性)的特征的数据被存储在中等性能存储装置(例如,硬盘驱动器)中;并且对应于重要性度量在第三值范围内(例如,低重要性)的特征的数据被存储在低性能存储装置(例如,冷存储装置(cold storage))中。

[0037] 在各种实施例中,特征的管理可以考虑多个不同机器学习模型的其他机器学习特征的重要性度量。例如,如果一个特征对于一个模型不重要,但是对于使用相同特征的另一个模型重要,则可以保留该特征和/或其数据。

[0038] 在308,基于对一个或多个机器学习特征的管理(如果适用的话)来生成机器学习模型的新版本。例如,如果在306中对特征的管理已经导致机器学习特征被移除,则机器学习模型的新版本被利用不包括移除的特征的训练数据生成/重新训练以生成新版本。因

此,可以丢弃不太重要的特征,并且可以在没有丢弃的特征的情况下重新训练机器学习模型,从而允许删除丢弃的特征的收集的数据,并且不再收集和存储未来的相关数据。这允许生成计算效率更高的模型(例如,部署模型所需的处理和存储更少),并且减少了需要收集和存储的数据量。在另一个示例中,如果在306中对特征的管理已经导致机器学习特征被添加、变换或修改,则机器学习模型的新版本利用包括新的/变换的/修改的特征的训练数据来生成/重新训练。利用这个新的/变换的/修改的特征,机器学习模型的新版本可以更有效地执行和/或更精确。

[0039] 图4是示出用于确定特征重要性度量的过程的实施例的流程图。图4的过程可以至少部分地在服务器102和/或机器学习管理系统106中的一个或多个上实现。在一些实施例中,图4的过程的至少一部分被包括在图2的208和/或图3的304中。

[0040] 在402,提供测试数据集作为机器学习模型的输入,以确定模型的初始基本性能。因为特征重要性评估可以是一个计算成本非常高的过程,所以可以将测试数据集选择成在大小上受限制,以降低特征重要性评估的计算成本,同时保持特征重要性评估的精度。例如,基于可用测试数据集的数量,所有可用测试数据集中只有满足大小/数量标准的一部分被选择用于特征重要性评估(例如,仅利用最后N天的数据,仅加载x%的数据,等等)。在一些实施例中,机器学习模型是在图2的208和/或图3的302中训练的模型。例如,已经被训练的机器学习模型将被分析以确定模型的特征的相对重要性。为了建立初始基线性能,将测试数据集的每个条目作为模型的输入来提供,并将模型的相应推断结果与相应的已知正确结果进行比较,以确定模型结果的精度/性能度量(例如,识别模型结果与相应的已知正确结果的接近度的精度百分比)。在一些实施例中,跨越测试数据集的多个条目的模型精度的统计量度(例如,平均值)被计算为模型的初始基本性能。

[0041] 在404,选择待评估的模型的至少一部分机器学习特征。例如,模型的每个特征要被评估,并且逐个选择每个特征进行评估,直到已经评估了模型的所有特征。在一些实施例中,只评估模型的特征的一部分。例如,被识别为对应于高于阈值总大小的存储数据大小的特征被选择用于评估(例如,仅测试需要大量存储空间来存储的特征)。在另一示例中,随机选择用于评估的特征(例如,抽查特征,因为评估所有特征可能消耗太多计算资源)。

[0042] 对每个选择的特征重复步骤406-410。

[0043] 在406,修改原始测试数据集中对应于所选特征的值。例如,原始测试数据集中对应于所选特征的值(例如,对应于所选特征的值列)被修改以创建修改的测试数据集。修改这些值包括用以下项中的一个或多个值替换这些值:随机生成的值、零值、翻转值、移除的值或任何其他更改/修改的值。模型精度受特征变化影响的程度可以指示特征对模型精度的重要性和贡献。这些修改测试数据集的不同方式的计算成本因各种精度权衡而不同。在各种机器学习用例中,可以基于模型的属性(例如,基于模型类型、模型用例、精度要求等)来在许多方法中选择某个测试数据集修改方法。例如,对于期望高结果精度的医学机器学习用例,选择具有高质量特征重要性评估的测试数据集修改方法(例如,利用随机化值而不仅仅是从测试数据集中移除值)。

[0044] 在408,经修改的测试数据集被提供作为机器学习模型的输入,以确定机器学习模型针对所选特征的经修改的测试数据集的新性能。将经修改的测试数据集的每个条目被提供作为模型的输入,并将模型的相应推断结果与相应的已知正确结果进行比较,以确定模

型结果的精度/性能度量(例如,识别模型结果与相应的已知正确结果的接近度的精度百分比)。在一些实施例中,跨经修改的测试数据集的多个条目的模型精度的统计量度(例如,平均值)被计算为该模型针对所选特征的经修改的测试数据集的新性能。

[0045] 在410,基于模型的初始基本性能和针对所选特征的经修改的测试数据集的新性能之间的比较来确定特征重要性度量。在一些实施例中,特征重要性度量至少部分地通过计算模型的初始基本性能和针对所选特征的经修改的测试数据集的新性能之间的差异来计算。较大的差异表示特征对模型结果的重要性/贡献较大,而较小的差异表示特征对模型结果的重要性/贡献较小。在一个示例中,特征重要性度量是原始差值。在另一个示例中,特征重要性度量是百分比变化值(例如,针对所选特征的机器学习模型的初始基本性能和修改性能之间的百分比变化)。在另一个示例中,特征重要性度量是排名顺序位置值。每个所选特征的排名顺序位置值可以通过以下操作来确定:确定模型的初始基本性能和针对每个所选特征的经修改的测试数据集的新性能之间的每个差值的大小,并且当所选特征的所有差值从最大到最小的大小排名时,识别所选特征的排名顺序位置。

[0046] 在一些实施例中,如果针对所选特征的经修改的测试数据集的模型的新性能优于初始基本性能,则特征重要性度量被设置为对应于低特征重要性的值(例如,零值)。例如,如果表示初始基本性能的值减去表示针对所选特征的经修改的测试数据集的模型的新性能的值是负数,则特征重要性度量被设置为零值(例如,更大的特征重要性度量值表示更大的重要性)。

[0047] 图5是示出用于自动管理机器学习特征的过程的实施例的流程图。图5的过程可以至少部分地在服务器102和/或机器学习管理系统106中的一个或多个上实现。在一些实施例中,图5的过程的至少一部分被包括在图2的208和/或图3的306中。

[0048] 在502,为多个不同机器学习模型的机器学习特征确定特征重要性度量。例如,对于共享相同存储和/或计算资源的多个不同机器学习模型中的每个不同机器学习模型,执行图4的过程以确定每个机器学习模型的特征重要性度量。各种不同模型的这些不同特征重要性度量可以被收集并存储在中央储存库中(例如,存储在图1的管理系统106的存储装置中和/或存储装置110中)。

[0049] 在一些实施例中,仅对多个不同机器学习模型的机器学习特征的选定部分确定特征重要性度量。例如,被识别为对应于高于阈值总大小的存储数据大小的特征被选择用于评估(例如,仅测试需要大量存储空间来存储的特征)。在另一示例中,随机选择用于评估的特征(例如,抽查特征,因为评估所有特征可能消耗太多计算资源)。

[0050] 在504,跨多个不同模型的特征共享被识别。例如,为了管理由系统支持的机器学习模型的特征和特征的相应数据,希望知道特征和它们的数据在所有不同的机器学习模型中的何处被利用,因为这些模型共享相同的存储和/或处理资源池。在一些实施例中,对于每个独特特征,识别利用特定特征(例如,利用特征的收集/存储数据)的模型列表以及利用特定特征的每个模型的特定特征的相关特征重要性度量。例如,存储关于模型特征的信息的储存库被搜索以识别利用每个特征的模型。

[0051] 在506,基于所识别的特征共享和相关联的重要性度量,不同机器学习模型的机器学习特征被一起管理。例如,不是针对每个机器学习模型在单独的级别上独立地管理特征,而是跨共享相同的存储和/或计算资源池的多个不同的机器学习模型一起管理特征。在各

种实施例中,每个管理机器学习特征包括执行以下一项或更多项:从一个或多个模型中移除/丢弃特征、修改特征、删除对应于特征的数据、使对应于特征的数据不再被收集、为对应于特征的数据选择存储层、或基于被管理的特征生成新的特征。在一些实施例中,在确定不同机器模型的重要特征的超集(super set)时生成特征谱系图(lineage graph)。不在此超集中的特征可被移除,导致已移除的特征的相关数据不再被记录,也不再被发送到后端存储层进行存储。

[0052] 在一个示例中,如果特定模型的特定机器学习特征的重要性度量低于特定阈值,则该特定机器学习特征将从该模型中移除(例如,重新训练模型以移除特征)。然而,该特定特征的数据可能不能被删除,因为该特定特征的数据可能仍然被另一个模型使用。只有当没有其他模型要利用该特征的数据时,对应于特定机器学习特征的现有数据才会被自动删除,并且对应于特定机器学习特征的未来数据也不会被自动收集。

[0053] 在一些实施例中,对应于机器学习特征的数据的存储层基于其在所有使用它的模型中的重要性度量来确定(例如,通过将不同模型的特征的特征重要性度量加在一起,在特征的不同特征重要性度量中选择最高的特征重要性度量,对特征的不同特征重要性度量进行平均,等等,来确定总的重要性度量)。例如,对应于具有第一值范围内(例如,高重要性)的总重要性度量的特征的数据被存储在高性能存储装置(例如,固态驱动器)中;对应于具有在第二值范围内(例如,中等重要性)的总重要性度量的特征的数据被存储在中等性能存储装置(例如,硬盘驱动器)中;并且对应于具有在第三值范围内(例如,低重要性)的总重要性度量的特征的数据被存储在低性能存储装置(例如,冷存储装置)中。

[0054] 在一些实施例中,基于所识别的特征共享和相关联的重要性度量来确定不同模型中的特征之间的相对重要性。例如,对于每个独特的特征,通过计算不同模型中的任何一个所利用的相同特征的不同特征重要性度量的统计量度来确定总重要性度量(例如,将不同的特征重要性度量加在一起,在不同的特征重要性度量中选择最高的特征重要性度量,对不同的特征重要性度量进行平均,等等)。然后,可以对总特征重要性度量进行排名和排序,以识别根据它们的相对重要性(例如,对模型结果的贡献量)进行排序的独特特征的列表。在一些实施例中,总重要性度量和/或排名考虑了使用相应特征的模型的重要性类别和用于存储相应特征的数据的存储/资源量。如果需要额外的存储和/或处理资源,可以选择从所有模型中移除列表底部的特征,并且可以删除和不再收集与移除的特征相对应的数据。

[0055] 在一些实施例中,如果特征的总重要性度量低于阈值,则该机器学习特征将从所有模型中移除(例如,模型被重新训练以移除该特征),并且对应于该机器学习特征的现有数据被自动删除,并且对应于该机器学习特征的未来数据不再被自动收集。该阈值可以基于以下一项或更多项来动态确定:不同模型中其他特征的总重要性度量、不同模型中其他特征的总数、用于存储特征数据的存储量、可用存储资源量、可用处理资源量或利用该特征的不同机器学习模型的重要性(例如,排名、类别等)。

[0056] 尽管为了清楚理解的目的已经详细描述了前述实施例,但是本发明不限于所提供的细节。有许多实现本发明的替代方式。所公开的实施例是说明性的,而不是限制性的。

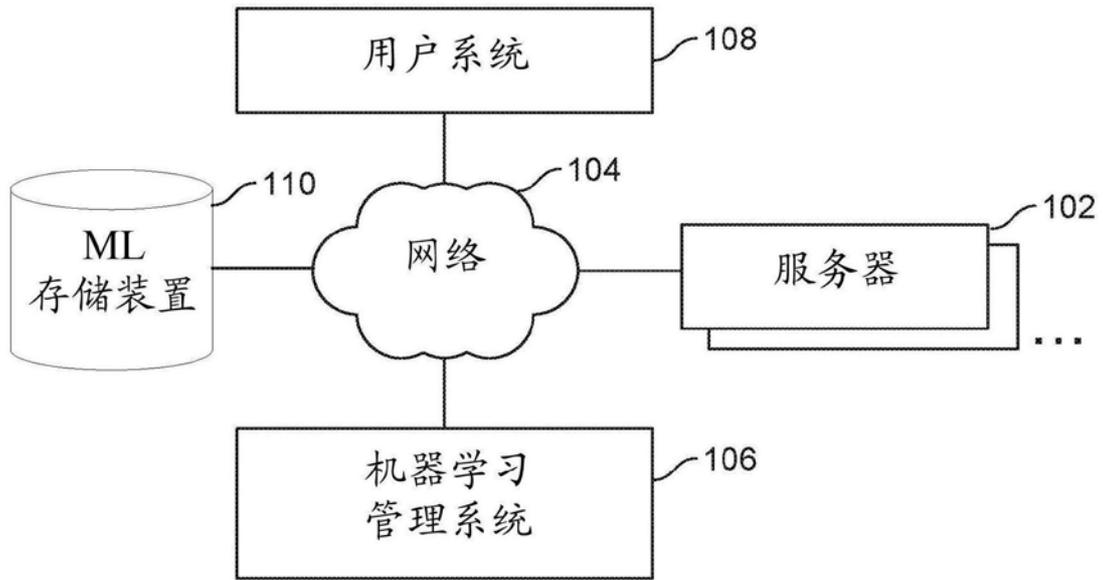


图1

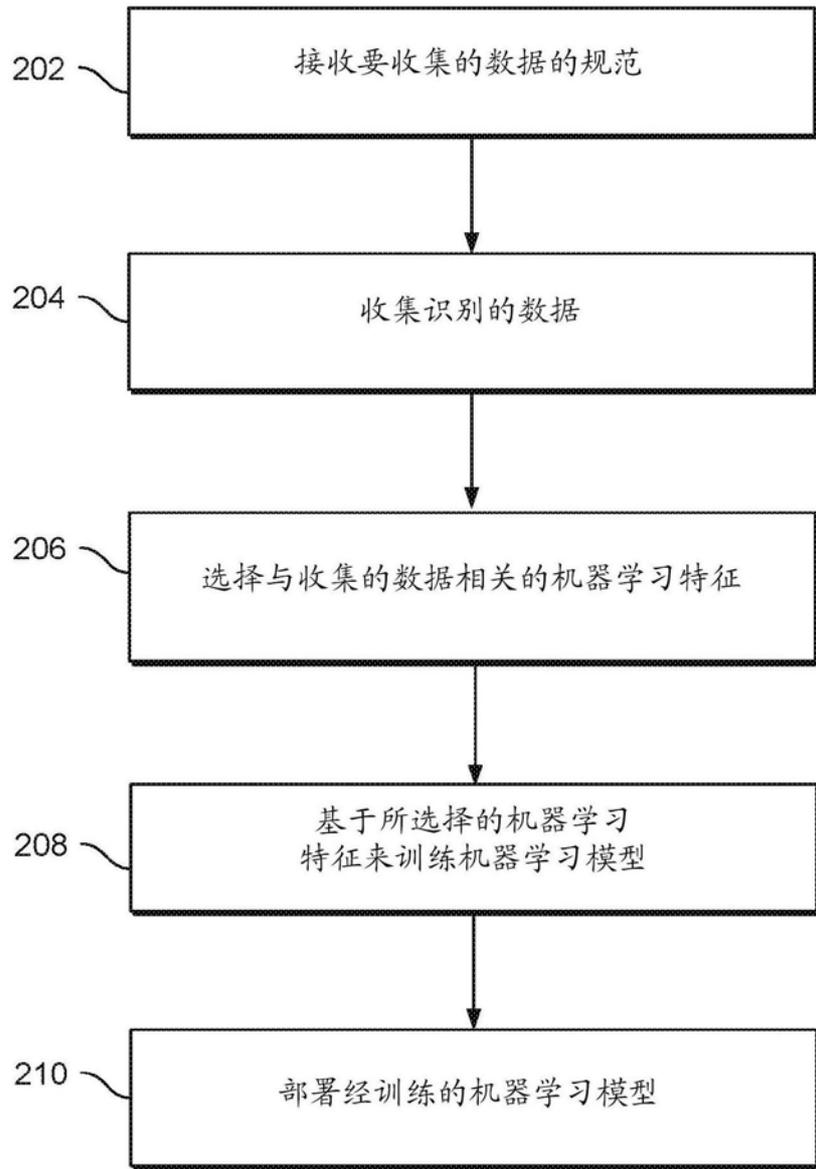


图2

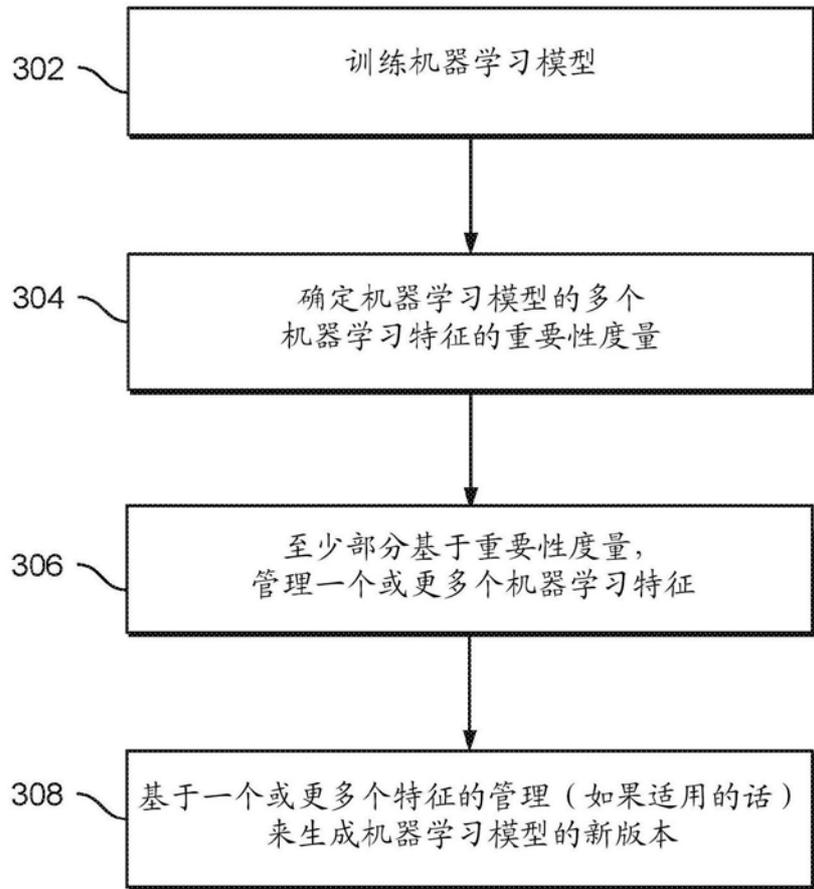


图3

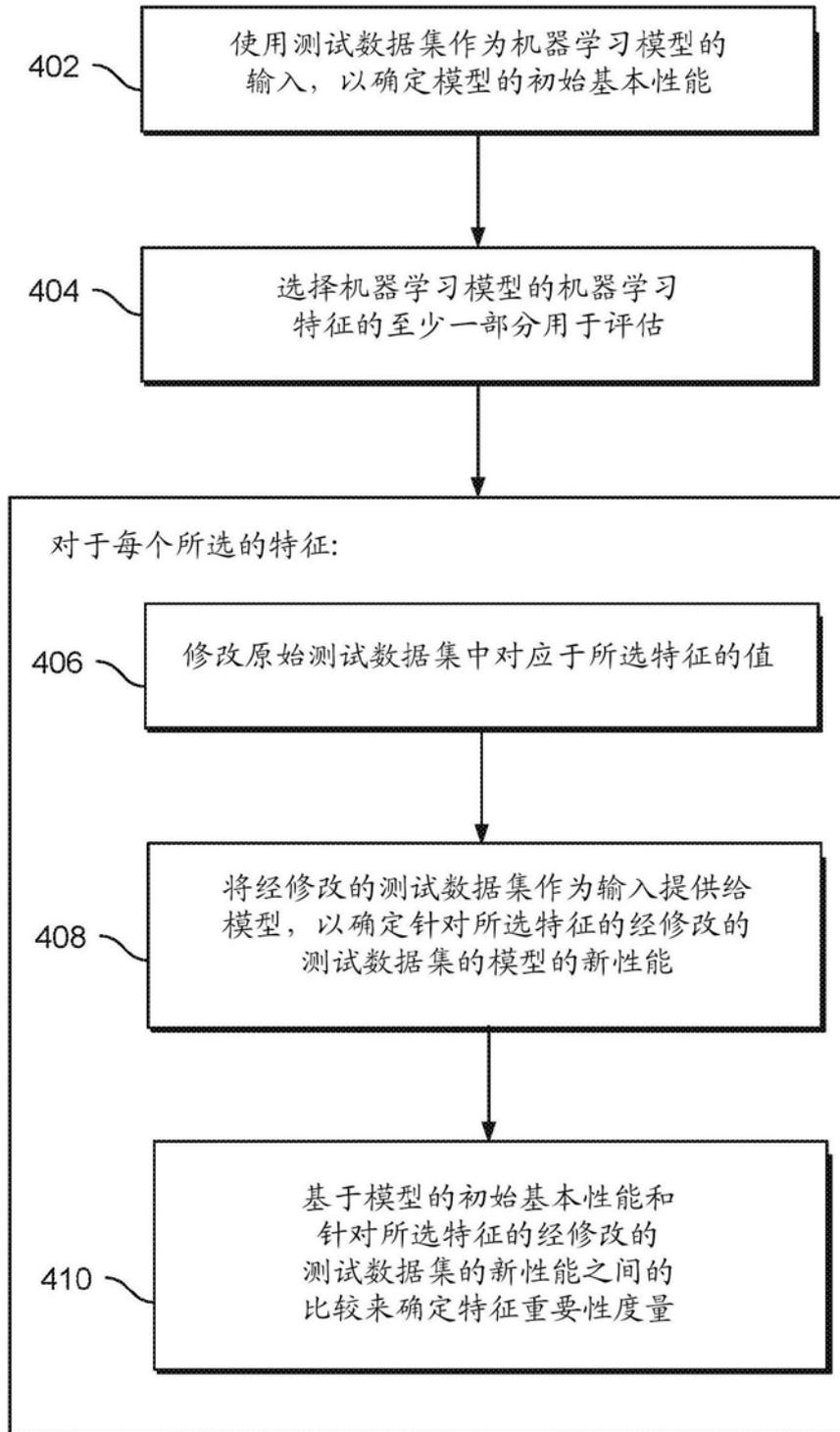


图4

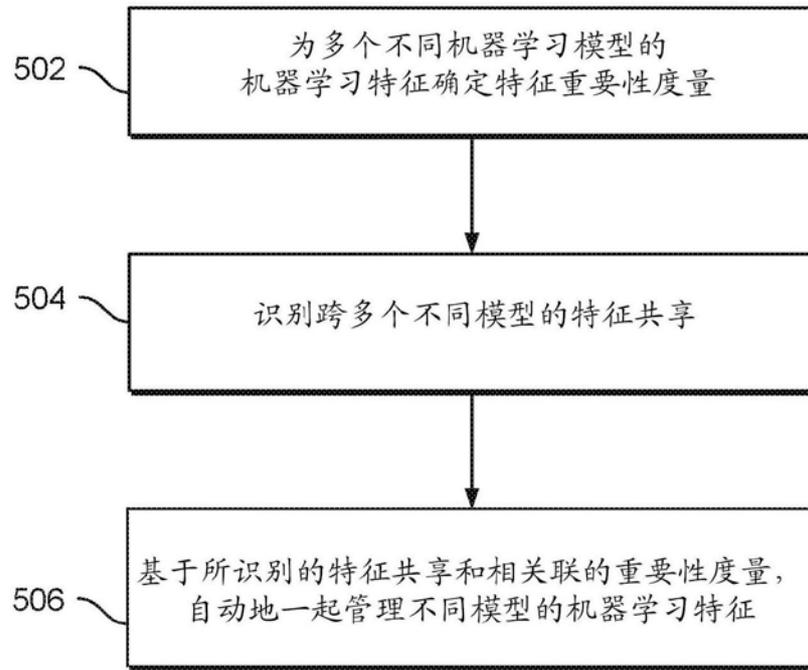


图5