

(19)日本国特許庁(JP)

(12)公表特許公報(A)

(11)公表番号

特表2024-520312

(P2024-520312A)

(43)公表日 令和6年5月24日(2024.5.24)

(51)国際特許分類 F I  
 G 0 6 N 3/08 (2023.01) G 0 6 N 3/08

審査請求 未請求 予備審査請求 未請求 (全30頁)

<p>(21)出願番号 特願2023-570346(P2023-570346)</p> <p>(86)(22)出願日 令和4年6月1日(2022.6.1)</p> <p>(85)翻訳文提出日 令和5年11月14日(2023.11.14)</p> <p>(86)国際出願番号 PCT/IB2022/055104</p> <p>(87)国際公開番号 WO2022/259089</p> <p>(87)国際公開日 令和4年12月15日(2022.12.15)</p> <p>(31)優先権主張番号 17/303,732</p> <p>(32)優先日 令和3年6月7日(2021.6.7)</p> <p>(33)優先権主張国・地域又は機関                  米国(US)</p> <p>(81)指定国・地域 AP(BW,GH,GM,KE,LR,LS,MW,MZ,NA,                  ,RW,SD,SL,ST,SZ,TZ,UG,ZM,ZW),EA(                  AM,AZ,BY,KG,KZ,RU,TJ,TM),EP(AL,A                  T,BE,BG,CH,CY,CZ,DE,DK,EE,ES,FI,FR                  ,GB,GR,HR,HU,IE,IS,IT,LT,LU,LV,MC,                  最終頁に続く</p>	<p>(71)出願人 390009531                  インターナショナル・ビジネス・マシー                  ンズ・コーポレーション                  INTERNATIONAL BUSI                  NESS MACHINES CORPO                  RATION                  アメリカ合衆国10504 ニューヨー                  ク州 アーモンク ニュー オーチャード                  ロード                  New Orchard Road, A                  rmonk, New York 105                  04, United States of                  America</p> <p>(74)代理人 100112690                  弁理士 太佐 種一                  最終頁に続く</p>
--	---

(54)【発明の名称】 人工知能モジュール訓練中のバイアス低減

(57)【要約】

本明細書において、1つまたは複数の選ばれた変数を含む入力データ・セットを受け取ることに応じて分析結果を提供するように訓練される、調節可能なパラメータをもつ人工知能モデルを訓練する方法が開示される。この方法は、訓練分析結果と対になる訓練入力データの多数のグループを含む訓練データ・セットを受け取ることと、訓練入力データの多数のグループを人工知能モデルに入力することに依りて、人工知能モデルから試行分析結果を受け取ることと、試行分析結果と訓練分析結果との間の比較を記述する正確度メトリックを計算することと、1つまたは複数の選ばれた変数を試行分析結果と比較することによって公平性スコア・メトリックを計算することと、公平性スコア・メトリックおよび正確度メトリックから組み合わせメトリックを計算することと、少なくとも組み合わせメトリックを受け取る訓練アルゴリズムを使用して、調節可能なパラメータを修正することを含む。

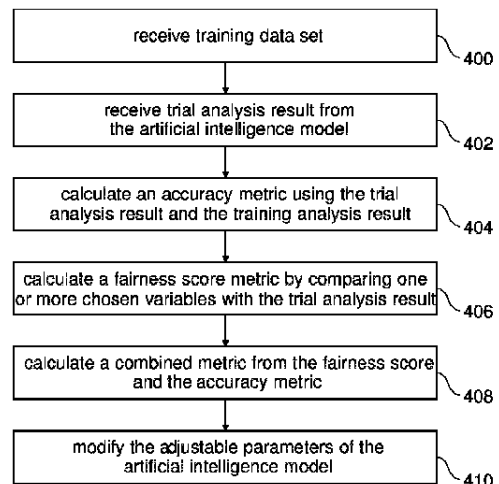


Fig. 4

## 【特許請求の範囲】

## 【請求項 1】

人工知能モデルを訓練する方法であって、前記人工知能モデルが、調節可能なパラメータを有し、前記人工知能モデルが、入力データ・セットを受け取ることに応じて、分析結果を提供するように訓練され、前記入力データ・セットが、1つまたは複数の選ばれた変数を含み、前記方法が、

前記人工知能モデルを訓練するための訓練データ・セットを受け取ることであり、前記訓練データ・セットが、訓練分析結果と対になる訓練入力データの多数のグループを含む、前記受け取ることと、

前記訓練入力データの多数のグループを前記入力データ・セットとして前記人工知能モデルに入力することに応じて、前記人工知能モデルから試行分析結果を受け取ることと、

前記試行分析結果と前記訓練分析結果との間の比較を記述する正確度メトリックを計算することと、

前記1つまたは複数の選ばれた変数を前記試行分析結果と比較することによって公平性スコア・メトリックを計算することと、

前記公平性スコア・メトリックおよび前記正確度メトリックから組み合わせメトリックを計算することと、

少なくとも前記組み合わせメトリックを入力として受け取る訓練アルゴリズムを使用して、前記人工知能モデルの前記調節可能なパラメータを修正することと

を含む、方法。

## 【請求項 2】

前記方法が、

多数の訓練された人工知能モデルを受け取ることであり、前記多数の訓練された人工知能モデルが前記人工知能モデルからなる、前記受け取ることと、

前記多数の人工知能モデルをテストするためのテスト・データ・セットを受け取ることとであり、前記テスト・データ・セットが、テスト分析結果と対になるテスト入力データの多数のグループを含む、前記受け取ることと、

前記テスト入力データの多数のグループを前記入力データ・セットとして入力することに応じて、前記多数の人工知能モデルの各々から緩和分析結果を受け取ることと、

前記多数の訓練された人工知能モデルの各々に対する前記緩和分析結果と、前記テスト分析結果との間の比較を記述する、前記多数の訓練された人工知能モデルの各々に対する正確度スコアを計算することと、

前記1つまたは複数の選ばれた変数を前記試行分析結果と比較することによって、前記多数の訓練された人工知能モデルの各々に対する公平性評価メトリックを計算することと

、

前記多数の訓練された人工知能モデルの各々に対して前記公平性評価メトリックと前記正確度スコアとを組み合わせることによって、前記多数の訓練された人工知能モデルの各々に対する公平性重み付きランキングを計算することと

によって、前記多数の訓練された人工知能モデルの各々に対する前記公平性重み付きランキングを提供することをさらに含む、請求項 1 に記載の方法。

## 【請求項 3】

前記公平性評価メトリックが、前記1つまたは複数の選ばれた変数の1つまたは複数の選ばれた値と、前記試行分析結果との間の相関を記述する、請求項 2 に記載の方法。

## 【請求項 4】

前記多数の訓練された人工知能モデルが、異なるタイプのものである、請求項 2 または 3 に記載の方法。

## 【請求項 5】

前記多数の訓練された人工知能モデルの各々が、独立して、以下のもの、すなわち、ニューラル・ネットワーク、分類器ニューラル・ネットワーク、畳み込みニューラル・ネットワーク、ベイジアン・ニューラル・ネットワーク、ベイジアン・ネットワーク、ベイズ

10

20

30

40

50

・ネットワーク、単純ベイズ分類器、信念ネットワーク、または決定ネットワーク、決定木、サポート・ベクトル機械、回帰分析、および遺伝的アルゴリズムのうちの任意の1つである、先行する請求項2ないし4のいずれか一項に記載の方法。

【請求項6】

前記公平性重み付きランキングが、以下のもの、すなわち、前記公平性評価メトリックと前記正確度スコアの最小2乗組み合わせ、前記公平性評価メトリックと前記正確度スコアの重み付き最小2乗組み合わせ、前記公平性評価メトリックと前記正確度スコアの線形組み合わせ、前記公平性評価メトリックと前記正確度スコアの重み付き組み合わせ、および前記公平性評価メトリックと前記正確度スコアの多項式組み合わせのうちの任意の1つを含む、請求項1ないし5のいずれか一項に記載の方法。

10

【請求項7】

前記組み合わせメトリックが、前記正確度スコアにスケール係数を乗算し、それを所定の冪乗したものであり、前記スケール係数が、前記公平性評価メトリックの関数である、請求項2ないし5のいずれか一項に記載の方法。

【請求項8】

前記スケール係数が、前記公平性評価メトリックの逆数である、請求項7に記載の方法。

【請求項9】

前記公平性スコア・メトリックが、前記1つまたは複数の選ばれた変数の1つまたは複数の選ばれた値と、前記試行分析結果との間の相関を記述する、請求項1ないし8のいずれか一項に記載の方法。

20

【請求項10】

前記組み合わせメトリックが、以下のもの、すなわち、前記公平性スコア・メトリックとテスト・メトリックの最小2乗組み合わせ、前記公平性スコア・メトリックと前記テスト・メトリックの重み付き最小2乗組み合わせ、前記公平性スコア・メトリックと前記テスト・メトリックの線形組み合わせ、前記公平性スコア・メトリックと前記テスト・メトリックの重み付き組み合わせ、および前記公平性スコア・メトリックと前記テスト・メトリックの多項式組み合わせのうちの任意の1つを含む、請求項1ないし9のいずれか一項に記載の方法。

【請求項11】

前記組み合わせメトリックが、以下のもの、すなわち、前記公平性スコア・メトリックに対する制約、前記テスト・メトリックに対する制約、前記公平性スコア・メトリックに対する最大許容値、および前記テスト・メトリックに対する最大許容値のうちの任意の1つを含む、請求項9または10に記載の方法。

30

【請求項12】

前記人工知能モデルが、以下のもの、すなわち、ニューラル・ネットワーク、分類器ニューラル・ネットワーク、畳み込みニューラル・ネットワーク、ベイジアン・ニューラル・ネットワーク、ベイジアン・ネットワーク、ベイズ・ネットワーク、単純ベイズ分類器、信念ネットワーク、または決定ネットワーク、決定木、サポート・ベクトル機械、回帰分析、および遺伝的アルゴリズムのうちの任意の1つである、請求項1ないし11のいずれか一項に記載の方法。

40

【請求項13】

前記人工知能モデルが畳み込みニューラル・ネットワークであり、前記訓練アルゴリズムが深層学習アルゴリズムである、請求項1ないし12のいずれか一項に記載の方法。

【請求項14】

コンピュータ可読プログラム・コードが具現化されたコンピュータ可読ストレージ媒体を含むコンピュータ・プログラム製品であって、前記コンピュータ可読のプログラム・コードが、請求項1ないし13に記載の前記方法を実施するように構成される、コンピュータ・プログラム製品。

【請求項15】

50

コンピュータ・システムを制御するように構成されたプロセッサと、  
機械実行命令を格納するメモリであって、前記命令の実行により、前記プロセッサが、  
人工知能モデルを訓練するための訓練データ・セットを受け取ることであり、前記人工知能モデルが、調節可能なパラメータを有し、前記人工知能モデルが、入力データ・セットを受け取ることに応じて、分析結果を提供するように訓練され、前記入力データ・セットが、1つまたは複数の選ばれた変数を含み、前記訓練データ・セットが、訓練分析結果と対になる訓練入力データの多数のグループを含む、前記受け取ることと、

前記訓練入力データの多数のグループを前記入力データ・セットとして前記人工知能モデルに入力することに応じて、前記人工知能モデルから試行分析結果を受け取ることと

10

前記試行分析結果と前記訓練分析結果との間の比較を記述する正確度メトリックを計算することと、

前記1つまたは複数の選ばれた変数を前記試行分析結果と比較することによって計算される公平性スコア・メトリックを計算することと、

前記公平性スコア・メトリックおよび前記正確度メトリックから組み合わせメトリックを計算することと、

少なくとも前記組み合わせメトリックを入力として受け取る訓練アルゴリズムを使用して、前記人工知能モデルの前記調節可能なパラメータを修正することと  
を行う、メモリと

を含むコンピュータ・システム。

20

#### 【請求項16】

前記命令の実行により、さらに、前記プロセッサが、

多数の訓練された人工知能モデルを受け取ることであり、前記多数の訓練された人工知能モデルが前記人工知能モデルからなる、前記受け取ることと、

前記多数の人工知能モデルをテストするためのテスト・データ・セットを受け取ることとであり、前記テスト・データ・セットが、テスト分析結果と対になるテスト入力データの多数のグループを含む、前記受け取ることと、

前記テスト入力データの多数のグループを前記入力データ・セットとして入力することに応じて、前記多数の人工知能モデルの各々から緩和分析結果を受け取ることと、

前記多数の訓練された人工知能モデルの各々に対する前記緩和分析結果と、前記テスト分析結果との間の比較を記述する、前記多数の訓練された人工知能モデルの各々に対する正確度スコアを計算することと、

30

前記1つまたは複数の選ばれた変数を前記試行分析結果と比較することによって、前記多数の訓練された人工知能モデルの各々に対する公平性評価メトリックを計算することと

、  
前記多数の訓練された人工知能モデルの各々に対して前記公平性評価メトリックと前記正確度スコアとを組み合わせることによって、前記多数の訓練された人工知能モデルの各々に対する公平性重み付きランキングを計算することと  
を行う、請求項15に記載のコンピュータ・システム。

#### 【請求項17】

前記人工知能モデルが、以下のもの、すなわち、ニューラル・ネットワーク、分類器ニューラル・ネットワーク、畳み込みニューラル・ネットワーク、ベイジアン・ニューラル・ネットワーク、ベイジアン・ネットワーク、ベイズ・ネットワーク、単純ベイズ分類器、信念ネットワーク、または決定ネットワーク、決定木、サポート・ベクトル機械、回帰分析、および遺伝的アルゴリズムのうちの任意の1つである、請求項15または16に記載のコンピュータ・システム。

40

#### 【請求項18】

前記人工知能モデルが畳み込みニューラル・ネットワークであり、前記訓練アルゴリズムが深層学習アルゴリズムである、請求項15ないし17のいずれか一項に記載のコンピュータ・システム。

50

## 【請求項 19】

コンピュータ・プログラム製品であって、前記コンピュータ・プログラム製品が、請求項 1 ないし 13 に記載の前記方法に従って訓練された人工知能モデルを格納したコンピュータ可読ストレージ媒体を含む、コンピュータ・プログラム製品。

## 【請求項 20】

データ処理システムで実行されるアプリケーション・プログラムがアクセスするためのデータを格納するためのメモリであって、請求項 1 ないし 13 に記載の前記方法に従って訓練される人工知能モデルを含む、メモリ。

## 【請求項 21】

多数の訓練された人工知能モデルの各々に公平性重み付きランキングを提供する方法であって、前記方法が、

前記多数の訓練された人工知能モデルを受け取ることと、

前記多数の人工知能モデルをテストするためのテスト・データ・セットを受け取ることであり、前記テスト・データ・セットが、テスト分析結果と対になるテスト入力データの多数のグループを含む、前記受け取ることと、

前記テスト入力データの多数のグループを入力データ・セットとして入力することに応じて、前記多数の人工知能モデルの各々から緩和分析結果を受け取ることと、

前記多数の訓練された人工知能モデルの各々に対する前記緩和分析結果と、前記テスト分析結果との間の比較を記述する、前記多数の訓練された人工知能モデルの各々に対する正確度スコアを計算することと、

1 つまたは複数の選ばれた値を試行分析結果と比較することによって、前記多数の訓練された人工知能モデルの各々に対する公平性評価メトリックを計算することと、

前記多数の訓練された人工知能モデルの各々に対して前記公平性評価メトリックと前記正確度スコアとを組み合わせることによって、前記多数の訓練された人工知能モデルの各々に対する前記公平性重み付きランキングを計算することとを含む、方法。

## 【発明の詳細な説明】

## 【技術分野】

## 【0001】

本発明は、人工知能モデルの訓練に関する。

## 【背景技術】

## 【0002】

人工知能（AI）モジュールおよびアルゴリズムの自動訓練は、非常に普及しており、それらを訓練するために必要とされる人間の労力の低減を可能にする。しかしながら、好結果の訓練は、現在、人工知能モジュールが系統的なバイアスまたは偏見を含む結果を生成しないように、訓練を非常に注意深く実行することに依拠する。現在、訓練データが系統的なバイアスを含む場合、訓練された人工知能モジュールもそうなることになる。

## 【発明の概要】

## 【0003】

1 つの態様では、本発明は、人工知能モデルを訓練する方法を提供する。人工知能モデルは、調節可能なパラメータを有する。調節可能なパラメータは、人工知能モデルの実行および動作に影響を与える。それゆえに、人工知能モデルは、調節可能なパラメータを修正または調節することによって訓練することができる。人工知能モデルは、入力データ・セットを受け取ることに応じて、分析結果を提供するように訓練される。入力データ・セットは、1 つまたは複数の選ばれた変数を含む。

## 【0004】

この方法は、人工知能モデルを訓練するための訓練データ・セットを受け取ることを含む。訓練データ・セットは、訓練分析結果と対になる訓練入力データの多数のグループを含む。訓練入力データは、人工知能モデルへの入力としての試行として使用されるデータとすることができる。次いで、人工知能モデルの出力は、訓練分析結果と比較することが

できる。この方法は、訓練入力データの多数のグループを入力データとして人工知能モデルに入力することに応じて、人工知能モデルから試行分析結果を受け取ることさらに含む。このステップにおいて、訓練入力データが人工知能モデルに入力され、それに応じて、試行分析結果を受け取られる。この方法は、前記試行分析結果と前記訓練分析結果との間の比較を記述する正確度 (accuracy) メトリックを計算することをさらに含む。人工知能モデルから得られる結果である試行分析結果は、訓練分析結果と比較され、正確度メトリックは、試行分析結果が訓練分析結果にどれだけ近いまたは正確であるかを評価する尺度または値を提供する。

【0005】

この方法は、1つまたは複数の選ばれた変数を試行分析結果と比較することによって公平性スコア・メトリックを計算することをさらに含む。人工知能の公平性尺度または公平性スコアは、特定の変数、またはこの場合には1つまたは複数の選ばれた変数が、人工知能モデルの出力にどれだけ影響を与えるかの尺度を指す。

10

【0006】

この方法は、公平性スコア・メトリックおよび正確度メトリックから組み合わせメトリックを計算することをさらに含む。この方法は、少なくとも前記組み合わせメトリックを入力として受け取る訓練アルゴリズムを使用して、人工知能モデルの調節可能なパラメータを修正することをさらに含む。

【0007】

本発明のさらなる態様によれば、本発明は、プロセッサと、機械実行可能命令を格納するメモリとを含むコンピュータ・システムを提供する。機械実行可能命令の実行により、プロセッサは、一実施形態による方法を実施する。

20

【0008】

本発明のさらなる態様によれば、本発明は、コンピュータ可読プログラム・コードが具現化されたコンピュータ可読ストレージ媒体を含むコンピュータ・プログラム製品を提供する。コンピュータ可読プログラム・コードは、一実施形態による方法を実施するように構成される。

【0009】

本発明のさらなる態様によれば、本発明は、コンピュータ・プログラム製品を提供する。コンピュータ・プログラム製品は、この方法の一実施形態に従って訓練された人工知能モデルを格納したコンピュータ可読ストレージ媒体を含む。

30

【0010】

本発明のさらなる態様によれば、本発明は、データ処理システムで実行されるアプリケーション・プログラムがアクセスするためのデータを格納するメモリを提供する。これは、この方法の一実施形態に従って訓練された人工知能モデルを含む。

【0011】

以下において、本発明の実施形態が、単に例として、図面を参照してより詳細に説明される。

【図面の簡単な説明】

【0012】

40

【図1】コンピュータ・システムの一例を示す図である。

【図2】図1のコンピュータ・システムが接続される例示的なコンピューティング環境を示す図である。

【図3】コンピュータ・システムのさらなる例を示す図である。

【図4】図3のコンピュータ・システムを使用する方法を示す流れ図である。

【図5】コンピュータ・システムのさらなる例を示す図である。

【図6】図5のコンピュータ・システムを使用する方法を示す流れ図である。

【発明を実施するための形態】

【0013】

本発明の様々な実施形態の説明は、例証の目的のために提示されるが、網羅的であるこ

50

と、または開示される実施形態に限定されることを意図するものではない。説明される実施形態の範囲および思想から逸脱することなく、多くの変更および変形が当業者には明らかであろう。本明細書で使用される用語は、実施形態の原理、実際の適用、もしくは市場で見いだされる技術に対する技術的改善を最も良く説明するために、または本明細書で開示される実施形態を他の当業者が理解できるようにするために選ばれた。

【0014】

実施形態は、1つまたは複数の選ばれた変数に対して望ましくないバイアスを低減する手段を提供することができるので有益であり得る。これは、例えば、訓練データ・セットが望ましくないバイアスまたは偏見を含むにもかかわらず、低減されたバイアスで人工知能モジュールを訓練することを可能にすることができる。

10

【0015】

例えば、人工知能モデルは、機械の保守が実行されるべきかどうかおよびいつ実行されるべきかを評価するように訓練される。人工知能モデルを訓練するために使用されるデータには以前の経験および個人的な好み起因するバイアスがある場合がある。

【0016】

通常、人工知能モデルが訓練されるとき、正確度メトリックのみが、調整可能パラメータを評価し、次いで、修正するために使用される。組み合わせメトリックは、正確な結果を提供するための人工知能モデルの必要性和、いわゆる公平な結果を提供することの必要性和をバランスさせる手段を提供することができる。それは、特定の変数における、またはこの場合、1つまたは複数の選ばれた変数における望ましくないバイアスを排除しようと試みることである。

20

【0017】

例えば、訓練アルゴリズムの入力としてただ単に正確度メトリックを使用する代わりに、組み合わせメトリックが代わりに使用される。すぐ上で説明したように、これは、1つまたは複数の選ばれた変数における望ましくないバイアスを除去する手段を提供することができる。ニューラル・ネットワークの例では、正確度メトリックは損失関数とすることができる。ニューラル・ネットワークの場合、正確度メトリックの結果の代わりに、組み合わせメトリックを、逆伝播アルゴリズムへの入力として使用することができる。ニューラル・ネットワークでは、組み合わせメトリックは、公平性スコア・メトリックの値を通常の損失関数または従来損失関数と組み合わせる修正された損失関数になることになる。

30

【0018】

別の実施形態では、この方法は、多数の訓練された人工知能モデルを最初に受け取ることによって、多数の訓練された人工知能モデルの各々に対して公平性重み付きランキングを提供することをさらに含む。多数の訓練された人工知能モデルは、人工知能モデルからなる。多数の訓練された人工知能モデルの各々に対する公平性重み付きランキングは、例えば、多数の訓練された人工知能モデルの各々が1つまたは複数の選ばれた変数においてどれだけのバイアスを有するかを識別するランキングとすることができる。

【0019】

この方法は、前記多数のインテリジェンス・モデルをテストするためにテスト・データ・セットを受け取ることさらに含む。テスト・データ・セットは、テスト分析結果と対になるテスト入力データの多数のグループを含む。テスト・データ・セットは、本質的に、多数の訓練される人工知能モデルの各々に入力するために使用される試行データである。特定のテスト・データ・セットでは、本質的に、人工知能モデルのうちの1つの正しいまたは所望の出力を提供するためにラベル付けされたグラウンド・トゥールズまたはデータであるテスト分析結果が存在する。

40

【0020】

この方法は、前記テスト入力データの多数のグループを前記入力データ・セットとして入力することに応じて、前記多数の人工知能モデルの各々から緩和分析結果を受け取ることさらに含む。緩和分析結果は、多数の人工知能モデルの試行の結果であると考えら

50

とができる。

【0021】

この方法は、多数の訓練された人工知能モデルの各々に対する緩和分析結果と、テスト分析結果との間の比較を記述する、多数の訓練された人工知能モデルの各々に対する正確度スコアを計算することをさらに含む。

【0022】

この方法は、1つまたは複数の選ばれた変数を試行分析結果と比較することによって、多数の訓練された人工知能モデルの各々に対する公平性評価メトリックを計算することをさらに含む。正確度スコアは、多数の訓練された人工知能モデルの各々がどれだけ正確であるかの尺度である。公平性評価メトリックは、多数の訓練された人工知能モデルの各々 10  
に対して、1つまたは複数の選ばれた変数にどれだけ望ましくないバイアスがあるかの尺度を提供する。

【0023】

この方法は、次いで多数の訓練された人工知能モデルの各々に対して公平性評価メトリックと正確度スコアとを組み合わせることによって、多数の訓練された人工知能モデルの各々に対する公平性重み付きランキングを計算することを含む。そのため、正確度スコアを使用することによって多数の訓練された人工知能モデルをランキングする代わりに、正確度スコアと公平性評価メトリックの組み合わせが代わりに使用される。これは、モデルがどれだけ正確であるかの値だけでなく、どれだけ望ましくないバイアスが様々な人工知能モデルに存在するかの値も提供する。次いで、公平性重み付きランキングは、最良の人工知能モデルの自動選択に役立つことができ、またはユーザに表示されてもよく、ユーザは、公平性重み付きランキングに基づいてどのモデルを使用するかを選択を決定してもよい。 20

【0024】

別の実施形態では、公平性評価メトリックは、前記1つまたは複数の選ばれた変数の1つまたは複数の選ばれた値と、試行分析結果との間の相関を記述する。例えば、公平性評価メトリックは、1つまたは複数の選ばれた変数の特定の値が差別されているかどうかを調べるために計算することができる。前記の例を使用して、特定のジェンダを選ぶことができ、この特定のジェンダが、訓練された人工知能モデルにバイアスをもたらすかどうかを調べる 30  
ことができる。これは、公平性評価メトリックが、訓練された人工知能モデルにおける特定のバイアスをチェックするために使用され得るので有益である。

【0025】

別の実施形態では、多数の訓練された人工知能モデルは異なるタイプである。例えば、多数の訓練された人工知能モデルは、異なるニューラル・ネットワーク・トポロジを使用することができる。他の例では、異なるタイプは、人工知能の完全に異なる実施態様でさえあり得る。1つの例は、あるモデルがニューラル・ネットワークであり、他のモデルがベイズ決定モデルである場合であろう。この実施形態は、最良の人工知能トポロジまたはモデル・タイプあるいはその両方を選択することを可能にすることができるので有益であり得る。 40

【0026】

別の実施形態では、多数の訓練された人工知能モデルのうちの1つはニューラル・ネットワークである。

【0027】

別の実施形態では、多数の訓練された人工知能モデルのうちの1つは、分類器ニューラル・ネットワークである。

【0028】

別の実施形態では、多数の訓練された人工知能モデルのうちの1つは、畳み込みニューラル・ネットワークである。

【0029】

別の実施形態では、多数の訓練された人工知能モデルのうちの1つは、ベイジアン・ニ 50

ューラル・ネットワークである。

【0030】

別の実施形態では、多数の訓練された人工知能モデルのうちの一つは、ベイジアン・ネットワークである。

【0031】

別の実施形態では、多数の訓練された人工知能モデルのうちの一つは、ベイズ・ネットワークである。

【0032】

別の実施形態では、多数の訓練された人工知能モデルのうちの一つは、単純ベイズ分類器である。

【0033】

別の実施形態では、多数の訓練された人工知能モデルのうちの一つは、信念ネットワークである。

【0034】

別の実施形態では、多数の訓練された人工知能モデルのうちの一つは、決定ネットワークである。

【0035】

別の実施形態では、多数の訓練された人工知能モデルのうちの一つは、決定木である。

【0036】

別の実施形態では、多数の訓練された人工知能モデルのうちの一つは、サポート・ベクトル機械である。

【0037】

別の実施形態では、多数の訓練された人工知能モデルのうちの一つは、回帰分析である。

【0038】

別の実施形態では、多数の訓練された人工知能モデルのうちの一つは、遺伝的アルゴリズムである。

【0039】

別の実施形態では、公平性重み付きランキングは、公平性評価メトリックと正確度スコアの最小2乗組み合わせを含む。

【0040】

別の実施形態では、公平性重み付きランキングは、公平性評価メトリックと正確度スコアの重み付き最小2乗組み合わせを含む。例えば、公平性評価メトリックを2乗し、次いで、第1の係数を乗算することができ、次いで、正確度スコアを2乗し、第2の係数を乗算し、次いで、2つのものを加算する。

【0041】

別の実施形態では、公平性重み付きランキングは、評価メトリックと正確度スコアの線形結合を含む。

【0042】

別の実施形態では、公平性重み付きランキングは、公平性評価メトリックと正確度スコアの重み付き組み合わせを含む。

【0043】

別の実施形態では、公平性重み付きランキングは、公平性評価メトリックと正確度スコアの多項式組み合わせを含む。例えば、多項式は、様々な係数を用いて選ぶことができ、次いで、公平性評価メトリックおよび正確度スコアは、各々、異なる組み合わせで多項式に入れることができる。

【0044】

別の実施形態では、組み合わせメトリックは、正確度スコアにスケール係数を乗算し、次いで、それを所定の冪乗したものである。スケール係数は、公平性評価メトリックの関数である。この実施形態は、公平性と正確度の良好な組み合わせ尺度を提供する

10

20

30

40

50

ことが示されているので有益であり得る。

【0045】

別の実施形態では、スケーリング係数は、公平性評価メトリックの逆数である。

【0046】

別の実施形態では、公平性スコア・メトリックは、前記1つまたは複数の選ばれた変数の1つまたは複数の選ばれた値と、試行分析結果との間の相関を記述する。公平性スコア・メトリックは、訓練中に人工知能モデルを評価するために使用される。この実施形態では、1つまたは複数の選ばれた値のうち特定の値を選択することができ、これが差別されているかどうかまたは望ましくないバイアスを有しているかどうかを評価することができる。例えば、特定のジェンダに対する差別を避けるようにモデルを訓練することができる。

10

【0047】

別の実施形態では、組み合わせメトリックは、公平性スコア・メトリックとテスト・メトリックの最小2乗組み合わせを含む。組み合わせメトリックは、公平性スコア・メトリックとテスト・メトリックの重み付き最小2乗組み合わせを含む。

【0048】

別の実施形態では、組み合わせメトリックは、公平性スコア・メトリックとテスト・メトリックの線形組み合わせを含む。

【0049】

別の実施形態では、組み合わせメトリックは、公平性スコア・メトリックとテスト・メトリックの重み付き組み合わせを含む。

20

【0050】

別の実施形態では、組み合わせメトリックは、公平性スコア・メトリックとテスト・メトリックの多項式組み合わせを含む。

【0051】

別の実施形態では、組み合わせメトリックは、公平性スコア・メトリックに対する制約を含む。例えば、制約は、公平性スコア・メトリックがどれだけ大きくなるのが許容されるかに対して限定される可能性がある。これは、特定の変数に対してどれだけのバイアスがあるかに対して限定を有する訓練された人工知能モデルを提供することができる。

【0052】

別の実施形態では、組み合わせメトリックは、テスト・メトリックに対する制約を含む。これは、例えば、訓練に許容できる最小正確度が存在するように訓練を限定するために使用することができるので、有用であり得る。これは、公平であるだけでなく正確でもあるモデルを構築するのに役立つことができる。

30

【0053】

別の実施形態では、組み合わせメトリックは、公平性スコア・メトリックの最大許容値を含む。

【0054】

別の実施形態では、組み合わせメトリックは、テスト・メトリックの最大許容値を含む。

40

【0055】

別の実施形態では、人工知能モデルは、ニューラル・ネットワークである。

【0056】

別の実施形態では、人工知能モデルは、分類器ニューラル・ネットワークである。

【0057】

別の実施形態では、人工知能モデルは、畳み込みニューラル・ネットワークである。

【0058】

別の実施形態では、人工知能モデルは、ベイジアン・ニューラル・ネットワークである。

【0059】

50

別の実施形態では、人工知能モデルは、ベイジアン・ネットワークである。

【0060】

別の実施形態では、人工知能モデルは、ベイズ・ネットワークである。

【0061】

別の実施形態では、人工知能モデルは、単純ベイズ分類器である。

【0062】

別の実施形態では、人工知能モデルは、信念ネットワークである。

【0063】

別の実施形態では、人工知能モデルは、決定ネットワークである。

【0064】

別の実施形態では、人工知能モデルは、決定木である。

【0065】

別の実施形態では、人工知能モデルは、サポート・ベクトル機械である。

【0066】

別の実施形態では、人工知能モデルは、回帰分析である。

【0067】

別の実施形態では、人工知能モデルは、遺伝的アルゴリズムである。

【0068】

別の実施形態では、人工知能モデルは、畳み込みニューラル・ネットワークである。訓練アルゴリズムは、深層学習アルゴリズムである。例えば、訓練アルゴリズムは、組み合わせメトリックを損失関数として使用する逆伝播アルゴリズムとすることができる。

10

20

【0069】

本発明の実施形態は、コンピュータ・システム、クライアント、またはサーバとも呼ばれることがあるコンピューティング・デバイスを使用して実施することができる。次に図1を参照すると、コンピュータ・システムの一例の概略図が示される。コンピュータ・システム10は、適切なコンピュータ・システムの単なる1つの例であり、本明細書に記載の本発明の実施形態の使用または機能の範囲に関していかなる限定も示唆するように意図されていない。それにもかかわらず、コンピュータ・システム10は、上述に記載の機能のいずれかを実施または実行あるいはその両方を行うことができる。

【0070】

コンピュータ・システム10内には、非常に多くの他の汎用または専用コンピューティング・システム環境または構成により動作可能なコンピュータ・システム/サーバ12がある。コンピュータ・システム/サーバ12とともに使用するのに好適であり得るよく知られているコンピューティング・システム、環境、または構成、あるいはその組み合わせの例には、限定はしないが、パーソナル・コンピュータ・システム、サーバ・コンピュータ・システム、シン・クライアント、シック・クライアント、ハンドヘルドまたはラップトップ・デバイス、マルチプロセッサ・システム、マイクロプロセッサ・ベース・システム、セット・トップ・ボックス、プログラマブル家庭用電化製品、ネットワークPC、ミニコンピュータ・システム、メインフレーム・コンピュータ・システム、および上述のシステムまたはデバイスのいずれかを含む分散コンピューティング環境、などが含まれる。

30

40

【0071】

コンピュータ・システム/サーバ12は、コンピュータ・システムによって実行されるプログラム・モジュールなどのコンピュータ・システム実行可能命令の一般的な文脈で説明することができる。一般に、プログラム・モジュールは、特定のタスクを実行するか、または特定の抽象データ型を実施するルーチン、プログラム、オブジェクト、コンポーネント、論理、データ構造、などを含むことができる。コンピュータ・システム/サーバ12は、通信ネットワークを介してリンクされるリモート処理デバイスによってタスクが実行される分散コンピューティング環境において実践され得る。分散コンピューティング環境では、プログラム・モジュールは、メモリ・ストレージ・デバイスを含むローカルおよびリモートの両方のコンピュータ・システム・ストレージ媒体に配置することができる。

50

## 【 0 0 7 2 】

図 1 に示されるように、コンピュータ・システム 1 0 内のコンピュータ・システム / サーバ 1 2 は、汎用コンピューティング・デバイスの形態で示される。コンピュータ・システム / サーバ 1 2 の構成要素は、限定はしないが、1 つまたは複数のプロセッサまたは処理ユニット 1 6 と、システム・メモリ 2 8 と、システム・メモリ 2 8 を含む様々なシステム構成要素をプロセッサ 1 6 に結合させるバス 1 8 とを含むことができる。バス 1 8 は、メモリ・バスまたはメモリ・コントローラ、周辺バス、アクセラレーテッド・グラフィック・ポート、および様々なバス・アーキテクチャのうちのいずれかを使用するプロセッサまたはローカル・バスを含むいくつかのタイプのバス構造のいずれかの 1 つまたは複数を表す。例として、限定ではなく、そのようなアーキテクチャには、産業標準アーキテクチャ ( I S A ) バス、マイクロ・チャンネル・アーキテクチャ ( M C A ) バス、エンハンスド I S A ( E I S A ) バス、ビデオ・エレクトロニクス標準協会 ( V E S A ) ローカル・バス、および周辺構成要素相互接続 ( P C I ) バスが含まれる。

10

## 【 0 0 7 3 】

コンピュータ・システム / サーバ 1 2 は、一般に、様々なコンピュータ・システム可読媒体を含む。そのような媒体は、コンピュータ・システム / サーバ 1 2 によってアクセス可能な任意の利用可能な媒体とすることができ、揮発性および不揮発性媒体と取り外し可能および取り外し不可媒体の両方を含む。

## 【 0 0 7 4 】

システム・メモリ 2 8 は、ランダム・アクセス・メモリ ( R A M ) 3 0 またはキャッシュ・メモリ 3 2 あるいはその両方などの揮発性メモリの形態のコンピュータ・システム可読媒体を含むことができる。コンピュータ・システム / サーバ 1 2 は、他の取り外し可能 / 取り外し不可、揮発性 / 不揮発性のコンピュータ・システム・ストレージ媒体をさらに含むことができる。単なる例として、ストレージ・システム 3 4 は、取り外し不可、不揮発性磁気媒体 ( 図示せず、一般に、「ハード・ドライブ」と呼ばれる ) から読み出し、それに書き込むために設けることができる。図示されていないが、取り外し可能、不揮発性磁気ディスク ( 例えば、「フロッピー ( R ) ・ディスク」 ) から読み出し、それに書き込むための磁気ディスク・ドライブと、C D - R O M、D V D - R O M、または他の光学媒体などの取り外し可能、不揮発性光ディスクから読み出し、またはそれに書き込むための光ディスク・ドライブとが設けられてもよい。そのような場合には、各々は、1 つまたは複数のデータ媒体インターフェースによってバス 1 8 に接続され得る。以下でさらに図示および説明されるように、メモリ 2 8 は、本発明の実施形態の機能を実行するように構成されたプログラム・モジュールのセット ( 例えば、少なくとも 1 つ ) を有する少なくとも 1 つのプログラム製品を含むことができる。

20

30

## 【 0 0 7 5 】

プログラム・モジュール 4 2 のセット ( 少なくとも 1 つ ) を有するプログラム / ユーティリティ 4 0 は、例として、限定ではなく、オペレーティング・システム、1 つまたは複数のアプリケーション・プログラム、他のプログラム・モジュール、およびプログラム・データと同様にメモリ 2 8 に格納することができる。オペレーティング・システム、1 つまたは複数のアプリケーション・プログラム、他のプログラム・モジュール、およびプログラム・データの各々、またはそれらの何らかの組み合わせは、ネットワーキング環境の実施態様を含むことができる。プログラム・モジュール 4 2 は、通常、本明細書に記載の本発明の実施形態の機能または技法あるいはその両方を実行する。

40

## 【 0 0 7 6 】

コンピュータ・システム / サーバ 1 2 はまた、キーボード、ポインティング・デバイス、ディスプレイ 2 4、などの 1 つまたは複数の外部デバイス 1 4、ユーザがコンピュータ・システム / サーバ 1 2 と対話することを可能にする 1 つまたは複数のデバイス、またはコンピュータ・システム / サーバ 1 2 が 1 つまたは複数の他のコンピューティング・デバイスと通信することを可能にする任意のデバイス ( 例えば、ネットワーク・カード、モデム、など ) あるいはその組み合わせと通信することができる。そのような通信は、入力 /

50

出力（I/O）インターフェース 22 を介して行うことができる。さらに、コンピュータ・システム / サーバ 12 は、ネットワーク・アダプタ 20 を介して、ローカル・エリア・ネットワーク（LAN）、汎用ワイド・エリア・ネットワーク（WAN）、またはパブリック・ネットワーク（例えば、インターネット）、あるいはその組み合わせなどの 1 つまたは複数のネットワークと通信することができる。図示のように、ネットワーク・アダプタ 20 は、バス 18 を介してコンピュータ・システム / サーバ 12 の他の構成要素と通信する。図示されていないが、他のハードウェア構成要素またはソフトウェア・コンポーネントあるいはその両方を、コンピュータ・システム / サーバ 12 とともに使用することができることを理解されたい。例は、限定はしないが、マイクロコード、デバイス・ドライバ、冗長処理ユニット、外部ディスク・ドライブ・アレイ、RAID システム、テープ・ドライブ、およびデータ・アーカイブ・ストレージ・システム、などを含む。

10

#### 【0077】

図 1 に示されるコンピュータ・システム 10 などのコンピュータ・システムは、人工知能モジュールの訓練などの本明細書に開示される動作を実行するために使用することができる。そのようなコンピュータ・システムは、人工知能モジュールを訓練するための訓練データ・セットなどの処理されるべきデータをローカル・インターフェースを通して受け取ることができるネットワーク接続性のないスタンド・アロン・コンピュータとすることができる。しかしながら、そのような動作は、通信ネットワークまたはコンピューティング・ネットワークあるいはその両方などのネットワークに接続されるコンピュータ・システムを使用して、同様に実行されてもよい。

20

#### 【0078】

図 2 は、例示的なコンピューティング環境を示し、コンピュータ・システム 10 などのコンピュータ・システムは、例えばネットワーク・アダプタ 20 を使用してネットワーク 200 に接続される。限定はしないが、ネットワーク 200 は、インターネット、ローカル・エリア・ネットワーク（LAN）、モバイル通信ネットワークなどの無線ネットワーク、などのような通信ネットワークとすることができる。ネットワーク 200 は、クラウド・コンピューティング・ネットワークなどのコンピューティング・ネットワークを含むことができる。コンピュータ・システム 10 は、人工知能モデルを訓練するための訓練データ・セットなどの処理されるべきデータをネットワーク 200 から受け取ることができ、または訓練データ・セットを使用して訓練された後の訓練された人工知能モジュールなどのコンピューティング結果を、ネットワーク 200 を介してコンピュータ・システム 10 に接続された別のコンピューティング・デバイスに提供することができ、あるいはその両方である。

30

#### 【0079】

コンピュータ・システム 10 は、ネットワーク 200 を介して受け取った要求に応じて、本明細書に記載の動作を、完全にまたは部分的に、実行することができる。特に、コンピュータ・システム 10 は、ネットワーク 200 を介してコンピュータ・システム 10 に接続され得る 1 つまたは複数のさらなるコンピュータ・システムとともに、そのような動作を分散計算で実行することができる。その目的のために、コンピューティング・システム 10 または任意のさらなる関連するコンピュータ・システムあるいはその両方は、ネットワーク 200 を使用して、専用メモリまたは共有メモリなどのさらなるコンピューティング・リソースにアクセスすることができる。

40

#### 【0080】

図 3 は、コンピュータ・システム 10 の理想化されたものを示す。コンピュータ 10 の処理ユニット 16 またはプロセッサ、ならびにネットワーク・アダプタ 20 および I/O インターフェース 22 が示される。メモリ 28 は、処理ユニット 16 がアクセスすることができる様々なタイプのメモリを表す。処理ユニットは、機械実行可能命令 300 を含むものとして示されている。機械実行可能命令 300 は、プログラム・モジュール 42 のうちの 1 つに相当する。メモリ 28 の様々な内容は、様々な場所、例えば、RAM 30、キャッシュ 32、または永続メモリなどに格納され得る。メモリ 28 は、さらに、調整可能

50

パラメータを有する人工知能モデル 302 を含むものとして示されている。

【0081】

人工知能モデルは、入力データ・セットを受け取ることに応じて、分析結果を提供するように訓練することができる。メモリ 28 は、さらに、人工知能モデル 302 を訓練するために使用される訓練データ・セット 304 を含むものとして示されている。訓練データ・セット 304 は、訓練入力データ 306 の多数のグループのうちの一つのグループと、訓練入力データの各々にとって利用可能であり得る訓練分析結果 308 とに分割することができる。訓練入力データ 306 は、人工知能モデル 302 に入力され、試行分析結果 310 を提供することができる。これは、メモリ 28 に格納されるものとして示されている。

【0082】

メモリ 28 は、さらに、正確度メトリック 312 を含むものとして示されている。正確度メトリック 312 は、試行分析結果 310 と訓練分析結果 308 との間で計算された。メモリ 28 は、さらに、入力データ・セットの一つまたは複数の選ばれた変数を試行分析結果 310 と比較することによって計算された公平性スコア・メトリック 314 を含むものとして示されている。メモリ 28 は、さらに、公平性スコア・メトリック 314 と正確度メトリック 312 を組み合わせることによって計算された組み合わせメトリック 316 を含むものとして示されている。次いで、組み合わせメトリック 316 は、人工知能モデル 302 の調整可能パラメータを調節するために、訓練アルゴリズム 318 とともに使用される。

【0083】

図 4 は、図 3 のコンピュータ 10 を動作させる方法を示す流れ図を示す。最初に、ステップ 400 において、訓練データ・セット 304 が受け取られる。次に、ステップ 402 において、試行分析結果 310 が、入力訓練データ 306 の多数のグループを入力データ・セットとして人工知能モデル 302 に入力することによって、人工知能モデル 302 から受け取られる。次に、ステップ 404 において、正確度メトリック 312 が計算され、それは、試行分析結果 310 と訓練分析結果 308 との間の比較を記述する。次いで、ステップ 406 において、公平性スコア・メトリック 314 が、一つまたは複数の選ばれた変数を試行分析結果 310 と比較することによって計算される。次に、ステップ 408 において、組み合わせメトリック 316 が、公平性スコア・メトリック 314 および正確度メトリック 312 から計算される。最後に、ステップ 410 において、人工知能モデル 302 の調節可能なパラメータが、少なくとも組み合わせメトリック 316 を入力として受け取る訓練アルゴリズム 318 を使用して修正される。

【0084】

図 5 は、コンピュータ 10 のさらなる図を示す。図 3 において示されたコンピュータ 10 の特徴は、図 5 に示される特徴と組み合わせることができる。

【0085】

メモリ 28 は、機械実行可能命令 300 を含むものとして示されている。メモリは、さらに、多数の訓練人工知能モデル 500 を含むものとして示されている。図 3 に示された人工知能モデル 302 は、多分、多数の訓練された人工知能モデル 500 のうちのひとつとすることができるであろう。メモリ 28 は、さらに、テスト入力データ 504 とテスト分析結果 506 とを含むテスト・データ・セット 502 を含むものとして示されている。テスト・データ・セット 502 は、多数の訓練された人工知能モデル 500 をテストおよび評価するために使用される。テスト入力データ 504 は、入力として使用され、様々な人工知能モデルの出力が、テスト分析結果と比較される。

【0086】

メモリ 28 は、さらに、緩和分析結果を含むものとして示されている。緩和分析結果 508 は、テスト入力データが様々な人工知能モデルに入力されるときに、様々な人工知能モデルによって返される結果である。メモリ 28 は、さらに、正確度スコア 510 を含むものとして示されている。正確度スコア 510 は、緩和分析結果 508 がテスト分析結果 506 に対してどれだけ正確であるかを評価するスコアである。メモリ 28 は、さらに、

10

20

30

40

50

1つまたは複数の選ばれた変数を緩和分析結果508と比較することによって、多数の訓練された人工知能モデル500の各々に対して計算された公平性評価メトリック512を含むものとして示されている。メモリ28は、さらに、公平性重み付きランキング514を含むものとして示されている。公平性重み付きランキング514は、正確度スコア510と公平性評価メトリック512の組み合わせである。

#### 【0087】

図6は、図5のコンピュータ・システム10を動作させる方法を示す流れ図を示す。図6に示される流れ図は、図4に示された流れ図と組み合わせられてもよい。例えば、様々な人工モデルの訓練が図4に示された方法を使用して実行された後、多数の訓練された人工知能モデルは、図6に示される方法を使用して比較することができる。

#### 【0088】

最初に、ステップ600において、多数の訓練された人工知能モデル500が受け取られる。次に、ステップ602において、テスト・データ・セット502が受け取られる。次に、ステップ604において、緩和分析結果508が、様々な訓練された人工知能モデル500にテスト入力データ504を入力することによって受け取られる。次に、ステップ606において、正確度スコア510が、多数の訓練された人工知能モデル500の各々に対して、特定の知能モデル500の緩和分析結果508と、テスト分析結果506とを比較することによって計算される。次に、ステップ608において、公平性評価メトリック512が、多数の訓練された人工知能モデル500の各々に対して、1つまたは複数の選ばれた変数を緩和分析結果508と比較することによって計算される。最後に、ステップ610において、公平性重み付きランキング514が、多数の訓練された人工知能モデル500の各々に対して、公平性評価メトリック512と正確度スコア510とを組み合わせることによって計算される。

#### 【0089】

自動機械学習手法は、今日では非常に普及している。それは、手動のデータ科学者の作業を自動化し、モデル開発プロセスを高速化することを可能にする。残念なことに、最良のモデルを見つけるには、かなりの量の時間およびリソースを必要とする可能性がある。自動機械学習プロセスの目標は、最も正確なモデルを見つけることである。

#### 【0090】

モデルが公平であることを確認することは、今日では、関連する可能性がある別の側面である。モデル公平性を評価し、緩和を可能にするように構成された専用の監視システムまたはライブラリがある。

#### 【0091】

実施形態は、バイアスチェックおよび緩和手順を自動機械学習プロセスに注入することができる。その手順は、スコアラ概念に基づく。

#### 【0092】

例示のシステムは、多分、2つのモジュール、すなわち、検出モジュール（人工知能モジュールの調節可能なパラメータを修正するために組み合わせメトリックを計算するために使用される）と、緩和モジュール（多数の訓練された人工知能モジュールに公平性評価メトリックを提供するための）とに基づくことができるであろう。モジュールは、別々にまたは一緒に使用することができる。

#### 【0093】

##### 1. 検出モジュール

#### 【0094】

検出モジュールは、公平性計算スコアラ（公平性評価メトリック）によって正規スコアラ・リストを拡張することに基づくことができる。スコアラ関数（本明細書では正確度スコアラと呼ぶ）は、機械学習モデル（人工知能モデル）を評価するために使用される。サンプル・スコアラは、正確度、Brierスコアラ損失、平均精度、バランス正確度、f1スコアラ、などを含む。自動ML（autoML）プロセスの各段階中に、選択されたスコアラが探索プロセスを最適化するために使用され、その結果、最良のスコアラ値をもつモデ

10

20

30

40

50

ルが見いだされる。スコアラは探索プロセスを最適化するために使用され、その結果、最良のスコアラ値をもつモデルが見いだされる。スコアラは、モデルの性能（正確度）を記述する機械学習スコアラである。これらは、本明細書では「`ml__scorers`」と呼ばれる。

#### 【0095】

このモジュールには、公平性メトリック・スコアラ（公平性スコア・メトリック）をプロセスに追加することによるスコアラの拡張リストが存在する。言い換えれば、新しいタイプのスコアラが、既存のMLアーキテクチャに注入されている。これは、本明細書では、「`fairness__scorer`」または公平性スコア・メトリックと呼ばれる。`ml__scorer`が計算されるたびに、「`fairness__scorer`」が（スコアラ・リストに追加されてから）、同様に実行され得る。

10

#### 【0096】

その結果、新しいメトリクスが、ユーザに返されることが可能であり、正確度、精度、および再現度などの機械学習メトリクスの次に、公平性スコア・メトリックが計算される。公平性スコア・メトリックは、本明細書では`disparate__impact`と呼ばれ、「`fairness__metrics`」カテゴリの下で計算される。

#### 【0097】

`disparate__impact`を計算するために、データ・セット内のあり得る困難（`adversity`）またはバイアスに関する情報が提供されてもよい。あり得るバイアスまたは偏見に関するこの情報は、本明細書では「`fairness__info`」と呼ばれる。`fairness__info`の例および説明が以下で説明される。この情報は、パラメータとして`AutoML`システムに渡され、公平性スコア・メトリックは、その情報に基づいて検出モジュールの各段階で計算される。公平性情報をもつ擬似コードにおけるシステムの例示的な呼出しは、以下のように示される。

20

```
automl= AutoMLSystem(scorer= 'accuracy',
learning_type= 'classification',
positive_label= "No Risk",
fairness_info= fairness_info
)
```

```
automl.fit(training_data,training_labels)
```

30

#### 【0098】

「`accuracy`」は、使用される正確度メトリックのタイプを指す。「`training_data`」は訓練入力データに対応し、「`training_labels`」は訓練分析結果に対応する。以下の公平性情報の保護された属性は、1つまたは複数の選ばれた変数に対応する。「`protected_attributes`」の「`privileged_groups`」は、1つまたは複数の選ばれた変数の1つまたは複数の選ばれた値に対応する。

#### 【0099】

公平性情報の例：

#### 【0100】

- 分類： 「`privileged_groups`」の分類の以下の例は、バイアスされることがある1つまたは複数の選ばれた変数の特定の値であり得る。

40

```
fairness_info= {
  "protected_attributes": [
    {"feature": "Gender", "privileged_groups": ['male']},
    {"feature": "Age", "privileged_groups": [[0.0, 40.0]]},
  ],
  "favorable_labels": ["No Risk"]}

```

#### 【0101】

- 回帰

50

```

fairness_info = {
  "favorable_labels": [[-100000.0, 100]],
  "protected_attributes": [
    {"feature": "B", "privileged_groups": [[0.0, 40.0]]},
  ]
}

```

ここで、

- `protected_attribute` (アイテムのディクショナリ) - 公平性が望ましい特徴名および特権グループのサブセット。

- `favorable_labels` (アレイ) - 好ましい (すなわち、「肯定的である」と考えられるラベル値。利用可能なタイプ: スtring、数、数のアレイ

10

【0102】

例示的なメトリクス出力:

パイプライン0のスコア: 異種の影響: 0.81、正確度および異種の影響: 0.71

パイプライン1のスコア: 異種の影響: 0.84、正確度および異種の影響: 0.77

パイプライン2のスコア: 異種の影響: 0.67、正確度および異種の影響: 0.82

パイプライン3のスコア: 異種の影響: 0.66、正確度および異種の影響: 0.84

【0103】

上述において、「異種の影響」は「公平性スコア・メトリック」であり、「正確度および異種の影響」は「組み合わせメトリック」である。

【0104】

20

2. 緩和モジュール

【0105】

緩和モジュールは、再び、スコアラ手法に基づく。ここで、いわゆる組み合わせスコアラが、もう一度導入される。いくつかの重みに基づいてML (正確度スコア) と公平性メトリック (公平性評価メトリック) の両方を組み合わせた組み合わせスコアラは、本明細書では公平性重み付きランキングまたは本明細書では「`accuracy_and_disparate_impact_scorer`」とも呼ばれる。次に、そのようなスコアラは、ランキング・スコアラとして設定され、最適化プロセスで使用される。これは、計算されたスコア値 (公平性重み付きランキング) に従って最良のモデルを見つけることに関与するプロセスである。緩和モジュールでは、それは、組み合わせされた値である。それは、`autoML`システムの各段階で計算されるが、加えて、モデル選択ステップ中のモデル・ランキングのために使用される (多数の訓練された人工知能モデルの各々に公平性重み付きランキングを提供することによって)。組み合わせスコアラのうちの一つは、公平性スコアラ (公平性評価メトリック) であり、検出モジュールの公平性スコア・メトリックと類似しており、それは、すべての`autoML`システム・ステップで計算することができ、提供された`fairness_info`も使用する。組み合わせスコアラの最終値は、異種の影響の比に依存する。

30

【0106】

公平性メトリック (公平性評価メトリック) がNaN (ゼロ除算によって生じるものなどの非数 (`not a number`)) である場合、公平性情報は、データ・セット・サンプル (例えば、`k`-分割交差検証からのサンプル) に適していないので、組み合わせメトリックからの第2のメトリック、例えば正確度が返される。

40

【0107】

異種の影響の比 (公平性評価メトリック) が0.0に等しいとき、組み合わせメトリック (公平性重み付きランキング) の最終値は0.0である。

【0108】

そうでない場合、組み合わせメトリックは、以下の式を使用して、両方のメトリックの混合として計算される。

【0109】

正確度 (正確度スコア) および異種の影響 (公平性評価メトリック) = 正確度 \* (スケ

50

ーリング係数) ^ (スケーリング硬度)

【0110】

ここで、

【0111】

スケーリング係数は、0.9に設定されたパラメータである異種の影響の閾値(この閾値を超える値は公平と考えられる)と、以下で説明するパラメータである対称影響値(symmetric impact value)とに依存する。異種の影響が0と0.9の間にある場合、対称影響は、異種の影響に等しい。異種の影響が1.0より大きい場合、対称影響は、以下の式を使用して計算される。

【0112】

$scaling\_factor = (\text{対称影響}) / (\text{異種の影響の閾値})$

【0113】

スケーリング硬度は、4.0に設定されたパラメータである。

【0114】

例示的な緩和モジュールにおいて利用可能な2つの組み合わせスコアラ(公平性重み付きランキングを計算するための)がある。

【0115】

- 回帰: `r2_and_disparage_impact`

【0116】

- 分類: `accuracy_and_disparate_impact`

【0117】

緩和モジュールの例示的な呼出し:

```
automl = AutoMLSystem(scorer = 'accuracy_and_disparate_impact',
learning_type= 'classification'
positive_label= "No Risk",
fairness_info= fairness_info
)

automl.fit(training_data,training_labels)
```

【0118】

公平性情報の例:

- 分類

```
Fairness_info ={
  "protected_attributes": [
{"feature": "GENDER", "privileged_groups": ['F']},
{"feature": "BP", "privileged_groups": ["LOW", "NORMAL"]}
],
  "favorable_labels": ["drugA", "durgC"]}
```

【0119】

例示的なメトリクス出力:

パイプライン0のスコア: 異種の影響: 0.60、正確度および異種の影響: 0.64

パイプライン1のスコア: 異種の影響: 0.66、正確度および異種の影響: 0.68

パイプライン2のスコア: 異種の影響: 0.71、正確度および異種の影響: 0.77

パイプライン3のスコア: 異種の影響: 0.70、正確度および異種の影響: 0.81

【0120】

上述において、「異種の影響」は「公平性評価メトリック」であり、「正確度および異種の影響」は「公平性重み付きランキング」である。

【0121】

モデル・ランキングはまた、解釈を容易にするために、両方のメトリック(分離された)、すなわち、正確度のような機械学習メトリックと異種の影響のような公平性メトリックとを使用して行うことができる。それは、エンド・ユーザへの有用な提示と、選択され

10

20

30

40

50

たメトリックに基づいてランキングまたはソートあるいはその両方を行う能力とを可能にする。

【0122】

その選択はまた、いくつかの閾値に基づくフィルタリングに容易に拡張することができる。ユーザは、例えば、最良の公平性パイプラインを提供するが、精度が0.8以上であるという制約を設定する。

【0123】

本発明は、任意の可能な技術的詳細レベルの統合におけるシステム、方法、またはコンピュータ・プログラム製品、あるいはその組み合わせであり得る。コンピュータ・プログラム製品は、プロセッサに本発明の態様を実行させるためのコンピュータ可読プログラム命令を有する1つのコンピュータ可読ストレージ媒体（または複数の媒体）を含む。

10

【0124】

コンピュータ可読ストレージ媒体は、命令実行デバイスによる使用のための命令を保持および格納することができる有形のデバイスとすることができる。コンピュータ可読ストレージ媒体は、例えば、限定はしないが、電子ストレージ・デバイス、磁気ストレージ・デバイス、光学ストレージ・デバイス、電磁気ストレージ・デバイス、半導体ストレージ・デバイス、または前述のものの任意の適切な組み合わせとすることができる。コンピュータ可読ストレージ媒体のより具体的な例の非網羅的なリストは、以下のもの、すなわち、ポータブル・コンピュータ・ディスク、ハード・ディスク、ランダム・アクセス・メモリ（RAM）、読み取り専用メモリ（ROM）、消去可能プログラマブル読み取り専用メモリ（EPROMまたはフラッシュ・メモリ）、スタティック・ランダム・アクセス・メモリ（SRAM）、ポータブル・コンパクト・ディスク読み取り専用メモリ（CD-ROM）、デジタル・バーサタイル・ディスク（DVD）、メモリ・スティック、フロッピー（R）・ディスク、命令が記録されたパンチカードまたは溝内の隆起構造などの機械的にコード化されたデバイス、および前述のものの任意の適切な組み合わせを含む。本明細書で使用されるコンピュータ可読ストレージ媒体は、電波もしくは他の自由に伝播する電磁波、導波路もしくは他の伝送媒体を通して伝搬する電磁波（例えば、光ファイバ・ケーブルを通過する光パルス）、またはワイヤを通して伝送される電気信号などのそれ自体一過性信号であると解釈されるべきではない。

20

【0125】

本明細書に記載のコンピュータ可読プログラム命令は、コンピュータ可読ストレージ媒体からそれぞれのコンピューティング/処理デバイスに、あるいはネットワーク、例えば、インターネット、ローカル・エリア・ネットワーク、ワイド・エリア・ネットワーク、または無線ネットワーク、あるいはその組み合わせを介して外部コンピュータまたは外部ストレージ・デバイスにダウンロードされ得る。ネットワークは、銅伝送ケーブル、光伝送ファイバ、無線伝送、ルータ、ファイアウォール、スイッチ、ゲートウェイ・コンピュータ、またはエッジ・サーバ、あるいはその組み合わせを含むことができる。各コンピューティング/処理デバイスのネットワーク・アダプタ・カードまたはネットワーク・インターフェースは、コンピュータ可読プログラム命令をネットワークから受け取り、そのコンピュータ可読プログラム命令をそれぞれのコンピューティング/処理デバイス内のコンピュータ可読ストレージ媒体に格納するために転送する。

30

40

【0126】

本発明の動作を実行するためのコンピュータ可読プログラム命令は、アセンブラ命令、命令セット・アーキテクチャ（ISA）命令、機械命令、機械依存命令、マイクロコード、ファームウェア命令、状態設定データ、集積回路のための構成データ、あるいはSmalltalk（R）、C++などのようなオブジェクト指向プログラミング言語および「C」プログラミング言語または同様のプログラミング言語などの手続き型プログラミング言語を含む1つまたは複数のプログラミング言語の任意の組み合わせで書かれたソース・コードまたはオブジェクトコードのいずれかとすることができる。コンピュータ可読プログラム命令は、全面的にユーザのコンピュータで、部分的にユーザのコンピュータで、

50

スタンドアロン・ソフトウェア・パッケージとして、部分的にユーザのコンピュータでおよび部分的にリモート・コンピュータで、または全面的にリモート・コンピュータもしくはサーバで実行することができる。後者のシナリオでは、リモート・コンピュータは、ローカル・エリア・ネットワーク（LAN）もしくはワイド・エリア・ネットワーク（WAN）を含む任意のタイプのネットワークを通してユーザのコンピュータに接続されてもよく、または接続が外部コンピュータに対して行われてもよい（例えば、インターネット・サービス・プロバイダを使用してインターネットを通して）。いくつかの実施形態では、例えば、プログラマブル論理回路、フィールド・プログラマブル・ゲート・アレイ（FPGA）、またはプログラマブル論理アレイ（PLA）を含む電子回路は、本発明の態様を実行するために、コンピュータ可読プログラム命令の状態情報を利用して電子回路を個人専用にすることによってコンピュータ可読プログラム命令を実行することができる。

10

## 【0127】

本発明の態様は、本発明の実施形態による方法、装置（システム）、およびコンピュータ・プログラム製品の流れ図またはブロック図あるいはその両方を参照して本明細書に記載される。流れ図またはブロック図あるいはその両方の各ブロック、および流れ図またはブロック図あるいはその両方におけるブロックの組み合わせは、コンピュータ可読プログラム命令によって実現され得ることが理解されよう。

## 【0128】

これらのコンピュータ可読プログラム命令は、汎用コンピュータ、専用コンピュータ、または他のプログラマブル・データ処理装置のプロセッサに提供されて、コンピュータまたは他のプログラマブル・データ処理装置のプロセッサを介して実行される命令が流れ図またはブロック図あるいはその両方の1つまたは複数のブロックにおいて指定された機能/動作を実施するための手段を作り出すような機械を生成することができる。これらのコンピュータ可読プログラム命令はまた、命令が格納されたコンピュータ可読ストレージ媒体が流れ図またはブロック図あるいはその両方の1つまたは複数のブロックにおいて指定された機能/動作の態様を実施する命令を含む製品を構成するように、コンピュータ、プログラマブル・データ処理装置、または他のデバイス、あるいはその組み合わせに、特定のやり方で機能するように指示することができるコンピュータ可読ストレージ媒体に格納されてもよい。

20

## 【0129】

コンピュータ可読プログラム命令はまた、コンピュータ、他のプログラマブル・データ処理装置、または他のデバイスにロードされて、一連の動作ステップをコンピュータ、他のプログラマブル装置、または他のデバイスで実行させて、コンピュータ実施プロセスを作り出し、その結果、コンピュータ、他のプログラマブル装置、または他のデバイスで実行される命令が流れ図またはブロック図あるいはその両方の1つまたは複数のブロックにおいて指定された機能/動作を実施することができる。

30

## 【0130】

図における流れ図およびブロック図は、本発明の様々な実施形態によるシステム、方法、およびコンピュータ・プログラム製品の可能な実施態様のアーキテクチャ、機能、および動作を示す。これに関しては、流れ図またはブロック図の各ブロックは、指定された論理機能を実施するための1つまたは複数の実行可能命令を含む命令のモジュール、セグメント、または一部を表すことができる。いくつかの代替実施形態では、ブロックに記された機能は、図に記された順序から外れて行われてもよい。例えば、連続して示された2つのブロックは、実際には、実質的に同時に実行されてもよく、またはブロックは、時には、関連する機能に応じて逆の順序で実行されてもよい。ブロック図または流れ図あるいはその両方の各ブロック、およびブロック図または流れ図あるいはその両方のブロックの組み合わせは、指定された機能または動作を実行するかあるいは専用ハードウェア命令とコンピュータ命令の組み合わせを実行する専用ハードウェア・ベース・システムで実施され得ることに留意されたい。

40

## 【0131】

50

様々な例は、多分、以下の番号付けされた条項における以下の特徴のうちの1つまたは複数によって記載され得るであろう。

【0132】

条項1。

人工知能モデルを訓練する方法であって、人工知能モデルが、調節可能なパラメータを有し、前記人工知能モデルが、入力データ・セットを受け取ることに応じて、分析結果を提供するように訓練され、前記入力データ・セットが、1つまたは複数の選ばれた変数を含み、前記方法が、

前記人工知能モデルを訓練するための訓練データ・セットを受け取ることであり、前記訓練データ・セットが、訓練分析結果と対になる訓練入力データの多数のグループを含む、受け取ることと、

前記訓練入力データの多数のグループを前記入力データ・セットとして前記人工知能モデルに入力することに応じて、前記人工知能モデルから試行分析結果を受け取ることと、

前記試行分析結果と前記訓練分析結果との間の比較を記述する正確度メトリックを計算することと、

前記1つまたは複数の選ばれた値を前記試行分析結果と比較することによって公平性スコア・メトリックを計算することと、

前記公平性スコア・メトリックおよび前記正確度メトリックから組み合わせメトリックを計算することと、

少なくとも前記組み合わせメトリックを入力として受け取る訓練アルゴリズムを使用して、人工知能モデルの調節可能なパラメータを修正することとを含む、方法。

【0133】

条項2。

前記方法は、

前記多数の訓練された人工知能モデルを受け取ることであり、前記多数の訓練された人工知能モデルが前記人工知能モデルからなる、受け取ることと、

前記多数の人工知能モデルをテストするためのテスト・データ・セットを受け取ることであり、前記テスト・データ・セットが、テスト分析結果と対になるテスト入力データの多数のグループを含む、受け取ることと、

前記テスト入力データの多数のグループを前記入力データ・セットとして入力することに応じて、前記多数の人工知能モデルの各々から緩和分析結果を受け取ることと、

前記多数の訓練された人工知能モデルの各々に対する前記緩和分析結果と、前記テスト分析結果との間の比較を記述する、前記多数の訓練された人工知能モデルの各々に対する正確度スコアを計算することと、

前記1つまたは複数の選ばれた値を前記試行分析結果と比較することによって、前記多数の訓練された人工知能モデルの各々に対する公平性評価メトリックを計算することと、

多数の訓練された人工知能モデルの各々に対して前記公平性評価メトリックと正確度スコアとを組み合わせることによって、前記多数の訓練された人工知能モデルの各々に対する前記公平性重み付きランキングを計算することと

によって、多数の訓練された人工知能モデルの各々に対する公平性重み付きランキングを提供することをさらに含む、条項1に記載の方法。

【0134】

条項3。

前記公平性評価メトリックが、前記1つまたは複数の選ばれた変数の1つまたは複数の選ばれた値と、前記試行分析結果との間の相関を記述する、条項2に記載の方法。

【0135】

条項4。

前記多数の訓練された人工知能モデルが、異なるタイプのものである、条項2または3に記載の方法。

10

20

30

40

50

## 【0136】

条項 5。

前記多数の訓練された人工知能モデルの各々が、独立して、以下のもの、すなわち、ニューラル・ネットワーク、分類器ニューラル・ネットワーク、畳み込みニューラル・ネットワーク、ベイジアン・ニューラル・ネットワーク、ベイジアン・ネットワーク、ベイズ・ネットワーク、単純ベイズ分類器、信念ネットワーク、または決定ネットワーク、決定木、サポート・ベクトル機械、回帰分析、および遺伝的アルゴリズムのうちの任意の1つである、条項 2、3、または 4 に記載の方法。

## 【0137】

条項 6。

前記公平性重み付きランキングが、以下のもの、すなわち、前記公平性評価メトリックと前記正確度スコアの最小 2 乗組み合わせ、前記公平性評価メトリックと前記正確度スコアの重み付き最小 2 乗組み合わせ、前記公平性評価メトリックと前記正確度スコアの線形組み合わせ、前記公平性評価メトリックと前記正確度スコアの重み付き組み合わせ、および前記公平性評価メトリックと前記正確度スコアの多項式組み合わせのうちの任意の1つを含む、条項 1 ないし 5 のいずれか一項に記載の方法。

10

## 【0138】

条項 7。

前記組み合わせメトリックが、前記正確度スコアにスケール係数を乗算し、それを所定の冪乗したものであり、前記スケール係数が、前記公平性評価メトリックの関数である、条項 2 ないし 5 のいずれか一項に記載の方法。

20

## 【0139】

条項 8。

前記スケール係数が、前記公平性評価メトリックの逆数である、条項 7 に記載の方法。

## 【0140】

条項 9。

前記公平性スコア・メトリックが、前記 1 つまたは複数の選ばれた変数の 1 つまたは複数の選ばれた値と、前記試行分析結果との間の相関を記述する、条項 1 ないし 8 のいずれか一項に記載の方法。

30

## 【0141】

条項 10。

前記組み合わせメトリックが、以下のもの、すなわち、前記公平性スコア・メトリックと前記テスト・メトリックの最小 2 乗組み合わせ、前記公平性スコア・メトリックと前記テスト・メトリックの重み付き最小 2 乗組み合わせ、前記公平性スコア・メトリックと前記テスト・メトリックの線形組み合わせ、前記公平性スコア・メトリックと前記テスト・メトリックの重み付き組み合わせ、および前記公平性スコア・メトリックと前記テスト・メトリックの多項式組み合わせのうちの任意の1つを含む、条項 1 ないし 9 のいずれか一項に記載の方法。

## 【0142】

条項 11。

前記組み合わせメトリックが、以下のもの、すなわち、前記公平性スコア・メトリックに対する制約、前記テスト・メトリックに対する制約、前記公平性スコア・メトリックに対する最大許容値、および前記テスト・メトリックに対する最大許容値のうちの任意の1つを含む、条項 9 または 10 に記載の方法。

40

## 【0143】

条項 12。

前記人工知能モデルが、以下のもの、すなわち、ニューラル・ネットワーク、分類器ニューラル・ネットワーク、畳み込みニューラル・ネットワーク、ベイジアン・ニューラル・ネットワーク、ベイジアン・ネットワーク、ベイズ・ネットワーク、単純ベイズ分類器

50

、信念ネットワーク、または決定ネットワーク、決定木、サポート・ベクトル機械、回帰分析、および遺伝的アルゴリズムのうちの任意の1つである、条項1ないし11のいずれか一項に記載の方法。

【0144】

条項13。

前記人工知能モデルが畳み込みニューラル・ネットワークであり、前記訓練アルゴリズムが深層学習アルゴリズムである、条項1ないし12のいずれか一項に記載の方法。

【0145】

条項14。

コンピュータ可読プログラム・コードが具現化されたコンピュータ可読ストレージ媒体を含むコンピュータ・プログラム製品であって、前記コンピュータ可読のプログラム・コードが、条項1ないし13のいずれか一項に記載の方法を実施するように構成される、コンピュータ・プログラム製品。 10

【0146】

条項15。

コンピュータ・システムを制御するように構成されたプロセッサと、  
機械実行命令を格納するメモリであり、前記命令の実行により、前記プロセッサが、  
人工知能モデルを訓練するための訓練データ・セットを受け取ることであり、人工知能モデルが、調節可能なパラメータを有し、前記人工知能モデルが、入力データ・セットを受け取ることに応じて、分析結果を提供するように訓練され、前記入力データ・セットが、1つまたは複数の選ばれた変数を含み、前記訓練データ・セットが、訓練分析結果と対になる訓練入力データの多数のグループを含む、受け取ることと、 20

前記訓練入力データの多数のグループを前記入力データ・セットとして前記人工知能モデルに入力することに応じて、前記人工知能モデルから試行分析結果を受け取ることと、  
前記試行分析結果と前記訓練分析結果との間の比較を記述する正確度メトリックを計算することと、

前記1つまたは複数の選ばれた値を前記試行分析結果と比較することによって計算される公平性スコア・メトリックを計算することと、

前記公平性スコア・メトリックおよび前記正確度メトリックから組み合わせメトリックを計算することと、 30

少なくとも前記組み合わせメトリックを入力として受け取る訓練アルゴリズムを使用して、人工知能モデルの調節可能なパラメータを修正することと  
を行う、メモリと  
を含むコンピュータ・システム。

【0147】

条項16。

命令の実行により、さらに、前記プロセッサが、

前記多数の訓練された人工知能モデルを受け取ることであり、前記多数の訓練された人工知能モデルが前記人工知能モデルからなる、受け取ることと、 40

前記多数の人工知能モデルをテストするためのテスト・データ・セットを受け取ることであり、前記テスト・データ・セットが、テスト分析結果と対になるテスト入力データの多数のグループを含む、受け取ることと、

前記テスト入力データの多数のグループを前記入力データ・セットとして入力することに応じて、前記多数の人工知能モデルの各々から緩和分析結果を受け取ることと、

前記多数の訓練された人工知能モデルの各々に対する前記緩和分析結果と、前記テスト分析結果との間の比較を記述する、前記多数の訓練された人工知能モデルの各々に対する正確度スコアを計算することと、

前記1つまたは複数の選ばれた値を前記試行分析結果と比較することによって、前記多数の訓練された人工知能モデルの各々に対する公平性評価メトリックを計算することと、 50

多数の訓練された人工知能モデルの各々に対して前記公平性評価メトリックと正確度スコアとを組み合わせることによって、前記多数の訓練された人工知能モデルの各々に対する前記公平性重み付きランキングを計算することと  
を行う、条項 15 に記載のコンピュータ・システム。

【0148】

条項 17。

人工知能モデルが、以下のもの、すなわち、ニューラル・ネットワーク、分類器ニューラル・ネットワーク、畳み込みニューラル・ネットワーク、ベイジアン・ニューラル・ネットワーク、ベイジアン・ネットワーク、ベイズ・ネットワーク、単純ベイズ分類器、信念ネットワーク、または決定ネットワーク、決定木、サポート・ベクトル機械、回帰分析、および遺伝的アルゴリズムのうちの任意の 1 つである、条項 15 ないし 16 のいずれか一項に記載のコンピュータ・システム。

【0149】

条項 18。

前記人工知能モデルが畳み込みニューラル・ネットワークであり、前記訓練アルゴリズムが深層学習アルゴリズムである、条項 15 ないし 17 のいずれか一項に記載のコンピュータ・システム。

【0150】

条項 19。

コンピュータ・プログラム製品であって、前記コンピュータ・プログラム製品が、条項 1 ないし 12 のいずれか一項に記載の方法に従って訓練された人工知能モデルを格納したコンピュータ可読ストレージ媒体を含む、コンピュータ・プログラム製品。

【0151】

条項 20。

データ処理システムで実行されるアプリケーション・プログラムがアクセスするためのデータを格納するためのメモリであって、条項 1 ないし 12 のいずれか一項に記載の方法に従って訓練される人工知能モデルを含む、メモリ。

【0152】

条項 21。

多数の訓練された人工知能モデルの各々に公平性重み付きランキングを提供する方法であって、この方法が、

前記多数の訓練された人工知能モデルを受け取ることと、

前記多数の人工知能モデルをテストするためのテスト・データ・セットを受け取ることであり、前記テスト・データ・セットが、テスト分析結果と対になるテスト入力データの多数のグループを含む、受け取ることと、

前記テスト入力データの多数のグループを前記入力データ・セットとして入力することに応じて、前記多数の人工知能モデルの各々から緩和分析結果を受け取ることと、

前記多数の訓練された人工知能モデルの各々に対する前記緩和分析結果と、前記テスト分析結果との間の比較を記述する、前記多数の訓練された人工知能モデルの各々に対する正確度スコアを計算することと、

前記 1 つまたは複数の選ばれた値を前記試行分析結果と比較することによって、前記多数の訓練された人工知能モデルの各々に対する公平性評価メトリックを計算することと、

多数の訓練された人工知能モデルの各々に対して前記公平性評価メトリックと正確度スコアとを組み合わせることによって、前記多数の訓練された人工知能モデルの各々に対する前記公平性重み付きランキングを計算することと  
を含む、方法。

【0153】

本発明の様々な実施形態の説明は、例証の目的のために提示されたが、網羅的であること、または開示された実施形態に限定されることを意図するものではない。説明された実施形態の範囲および思想から逸脱することなく、多くの変更および変形が当業者には明ら

10

20

30

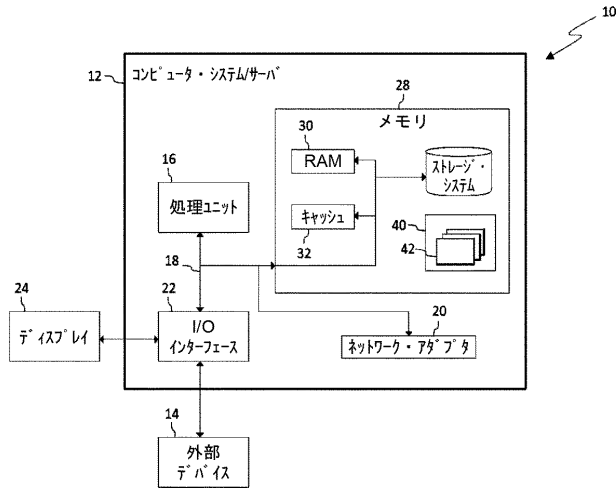
40

50

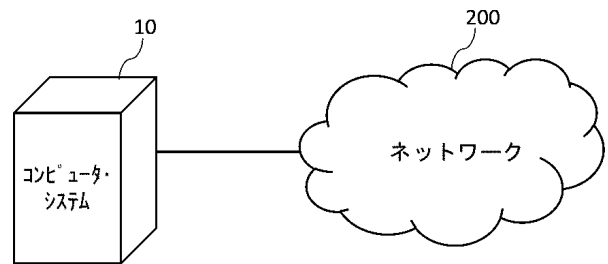
かであろう。本明細書で使用される用語は、実施形態の原理、実際の適用、もしくは市場で見いだされる技術に対する技術的改善を最もよく説明するために、または他の当業者が本明細書に開示される実施形態を理解できるようにするために選ばれた。

【図面】

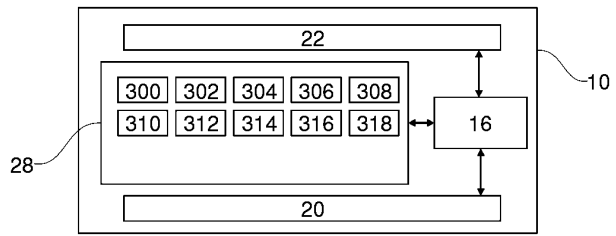
【図 1】



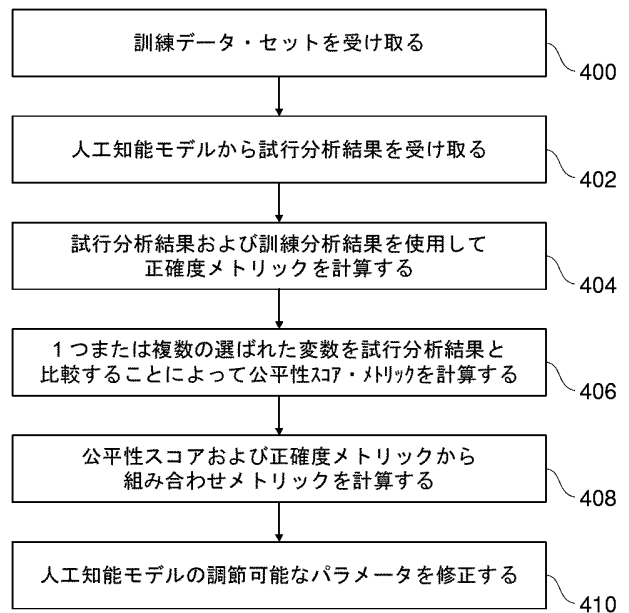
【図 2】



【図 3】



【図 4】



10

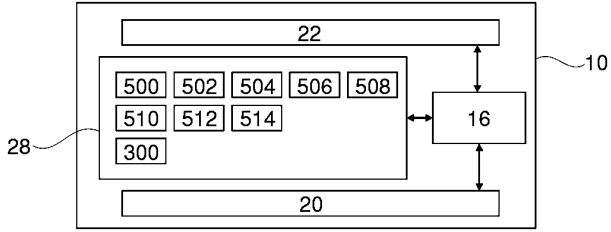
20

30

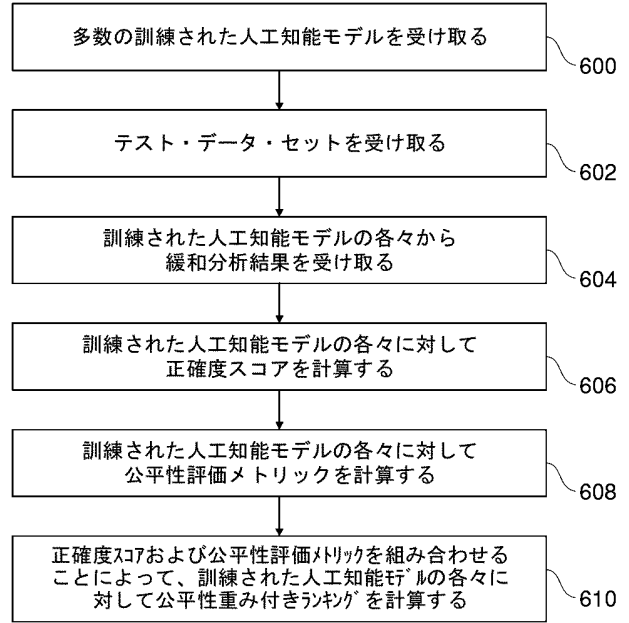
40

50

【 図 5 】



【 図 6 】



10

20

30

40

50

【 国際調査報告 】

INTERNATIONAL SEARCH REPORT

International application No.  
**PCT/IB2022/055104**

<b>A. CLASSIFICATION OF SUBJECT MATTER</b> G06N 3/08(2006.01)i; G06N 20/00(2019.01)i; G06N 7/00(2006.01)i  According to International Patent Classification (IPC) or to both national classification and IPC		
<b>B. FIELDS SEARCHED</b>  Minimum documentation searched (classification system followed by classification symbols) G06N3/-,G06N20/-,G06N7/-  Documentation searched other than minimum documentation to the extent that such documents are included in the fields searched  Electronic data base consulted during the international search (name of data base and, where practicable, search terms used) CNTXT, CNABS, DWPI, CNKI, WPABS, ENTXT, ENTXTC:AI,artificial w intelligence, model, deep 3d learning, neural 3d network, class, classifier, attribute, race, ethnicity, age, sex, dataset, data w set,input, output, variable, adjustable w parameter+, training, fair, accuracy, fairness, metric+		
<b>C. DOCUMENTS CONSIDERED TO BE RELEVANT</b>		
Category*	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
X	US 2020302524 A1 (ZESTFINANCE INC.) 24 September 2020 (2020-09-24) description, paragraphs [0042], [0048] to [0136] and figures 1A to 7	1-21
A	US 2020320428 A1 (IBM) 08 October 2020 (2020-10-08) the whole document	1-21
A	CN 110782004 A (CHAOCANSHU TECHNOLOGY SHENZHEN CO., LTD.) 11 February 2020 (2020-02-11) the whole document	1-21
A	CN 112541579 A (BEIJING BEIMING SHUKE INFORMATION TECHNOLOGY CO., LTD.) 23 March 2021 (2021-03-23) the whole document	1-21
A	US 2020372406 A1 (ORACLE INTERNATIONAL CORP.) 26 November 2020 (2020-11-26) the whole document	1-21
A	US 2021158204 A1 (IBM) 27 May 2021 (2021-05-27) the whole document	1-21
A	US 2020342307 A1 (IBM) 29 October 2020 (2020-10-29) the whole document	1-21
<input checked="" type="checkbox"/> Further documents are listed in the continuation of Box C. <input checked="" type="checkbox"/> See patent family annex.		
* Special categories of cited documents: "A" document defining the general state of the art which is not considered to be of particular relevance "E" earlier application or patent but published on or after the international filing date "L" document which may throw doubts on priority claim(s) or which is cited to establish the publication date of another citation or other special reason (as specified) "O" document referring to an oral disclosure, use, exhibition or other means "P" document published prior to the international filing date but later than the priority date claimed	"T" later document published after the international filing date or priority date and not in conflict with the application but cited to understand the principle or theory underlying the invention "X" document of particular relevance; the claimed invention cannot be considered novel or cannot be considered to involve an inventive step when the document is taken alone "Y" document of particular relevance; the claimed invention cannot be considered to involve an inventive step when the document is combined with one or more other such documents, such combination being obvious to a person skilled in the art "&" document member of the same patent family	
Date of the actual completion of the international search <b>15 July 2022</b>	Date of mailing of the international search report <b>29 August 2022</b>	
Name and mailing address of the ISA/CN <b>National Intellectual Property Administration, PRC 6, Xitucheng Rd., Jimen Bridge, Haidian District, Beijing 100088, China</b> Facsimile No. (86-10)62019451	Authorized officer <b>HUANG, Bin</b> Telephone No. 86-(10)-53962532	

Form PCT/ISA/210 (second sheet) (January 2015)

10

20

30

40

50

**INTERNATIONAL SEARCH REPORT**

International application No.  
**PCT/IB2022/055104**

<b>C. DOCUMENTS CONSIDERED TO BE RELEVANT</b>		
<b>Category*</b>	<b>Citation of document, with indication, where appropriate, of the relevant passages</b>	<b>Relevant to claim No.</b>
A	US 2020226489 A1 (ADOBE INC.) 16 July 2020 (2020-07-16) the whole document	1-21

10

20

30

40

50

**INTERNATIONAL SEARCH REPORT**  
**Information on patent family members**

International application No.  
**PCT/IB2022/055104**

Patent document cited in search report			Publication date (day/month/year)	Patent family member(s)			Publication date (day/month/year)
US	2020302524	A1	24 September 2020	WO	2020191057	A1	24 September 2020
				US	2021133870	A1	06 May 2021
				CA	3134043	A1	24 September 2020
				EP	3942384	A1	26 January 2022
				JP	2022525702	A	18 May 2022
US	2020320428	A1	08 October 2020	WO	2020208444	A1	15 October 2020
				CN	113692594	A	23 November 2021
				JP	2022527536	A	02 June 2022
				DE	11202000537	T5	21 October 2021
GB	2597406	A	26 January 2022				
CN	110782004	A	11 February 2020	None			
CN	112541579	A	23 March 2021	None			
US	2020372406	A1	26 November 2020	None			
US	2021158204	A1	27 May 2021	None			
US	2020342307	A1	29 October 2020	None			
US	2020226489	A1	16 July 2020	None			

10

20

30

40

## フロントページの続き

MK,MT,NL,NO,PL,PT,RO,RS,SE,SI,SK,SM,TR),OA(BF,BJ,CF,CG,CI,CM,GA,GN,GQ,GW,KM,ML,MR,N  
E,SN,TD,TG),AE,AG,AL,AM,AO,AT,AU,AZ,BA,BB,BG,BH,BN,BR,BW,BY,BZ,CA,CH,CL,CN,CO,CR,CU,  
CZ,DE,DJ,DK,DM,DO,DZ,EC,EE,EG,ES,FI,GB,GD,GE,GH,GM,GT,HN,HR,HU,ID,IL,IN,IQ,IR,IS,IT, JM,  
JO,JP,KE,KG,KH,KN,KP,KR,KW,KZ,LA,LC,LK,LR,LS,LU,LY,MA,MD,ME,MG,MK,MN,MW,MX,MY,M  
Z,NA,NG,NI,NO,NZ,OM,PA,PE,PG,PH,PL,PT,QA,RO,RS,RU,RW,SA,SC,SD,SE,SG,SK,SL,ST,SV,SY,TH,  
TJ,TM,TN,TR,TT,TZ,UA,UG,US,UZ,VC,VN,WS,ZA,ZM,ZW

(74)代理人 100120710

弁理士 片岡 忠彦

(72)発明者 クミエロウスキ、ルーカス

ポーランド 30 - 150 クラクフ アレヤ・アルミイ・クラヨベイ 18

(72)発明者 クハルチク、シモン

ポーランド 30 - 150 クラクフ アレヤ・アルミイ・クラヨベイ 18

(72)発明者 ヒルツェル、マーティン

アメリカ合衆国 10598 ニューヨーク州ヨークタウン・ハイツ キッチャワン・ロード 1101

(72)発明者 ラジャク、ドロタ

ポーランド 30 - 150 クラクフ アレヤ・アルミイ・クラヨベイ 18