

(19)日本国特許庁(JP)

(12)特 許 公 報(B2)

(11)特許番号

特許第7289086号  
(P7289086)

(45)発行日 令和5年6月9日(2023.6.9)

(24)登録日 令和5年6月1日(2023.6.1)

(51)Int. Cl.		F I
G 0 6 N 20/00	(2019.01)	G 0 6 N 20/00
G 0 6 Q 10/04	(2023.01)	G 0 6 Q 10/04

請求項の数 19 (全 24 頁)

(21)出願番号	特願2021-528927(P2021-528927)	(73)特許権者	390009531
(86)(22)出願日	令和1年11月28日(2019.11.28)		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
(65)公表番号	特表2022-510142(P2022-510142A)		INTERNATIONAL BUSINESS MACHINES CORPORATION
(43)公表日	令和4年1月26日(2022.1.26)		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
(86)国際出願番号	PCT/IB2019/060288		New Orchard Road, Armonk, New York 10504, United States of America
(87)国際公開番号	W02020/121104		
(87)国際公開日	令和2年6月18日(2020.6.18)		
審査請求日	令和4年4月18日(2022.4.18)	(74)代理人	100112690
(31)優先権主張番号	16/214,703		弁理士 太佐 種一
(32)優先日	平成30年12月10日(2018.12.10)		
(33)優先権主張国・地域又は機関	米国(US)		

最終頁に続く

(54)【発明の名称】 インスタンスレベルおよびグループレベルの予測メトリックの事後改善

(57)【特許請求の範囲】

【請求項1】

インスタンスレベルおよびグループレベルの予測メトリックの事後改善のための後処理コンピュータ実装方法であって、

グループ・バイアスへの制約付きで、所定の個人バイアス閾値よりも大きい個人バイアスを有するサンプルを検出することを学習するバイアス検出器を訓練することと、

実行時に、前記所定の個人バイアス閾値よりも大きい個人バイアスを有する実行時サンプル内のバイアスがかかったサンプルを選択するために、前記バイアス検出器を前記実行時サンプルに適用することと、

前記実行時に、前記バイアスがかかったサンプルに対するバイアス除去された予測を提案することと、

を含む、後処理コンピュータ実装方法。

【請求項2】

前記適用することおよび前記提案することは、後処理で動作し、個人バイアスを有する前記バイアスがかかったサンプルを是正の対象として、個人およびグループの両方の公平性メトリックに基づいて前記バイアスがかかったサンプルのバイアスを変更するようにする、請求項1に記載の後処理コンピュータ実装方法。

【請求項3】

前記バイアス検出器は、

前記サンプルを有する訓練セット内の保護属性を摂動させることと、

10

20

前記訓練セット内の各サンプル・ポイントについて、複数回の摂動後に平均個人バイアスを求めることと、

有利なクラスに関する前記個人バイアスの差を求めることと、

非特権グループの場合、前記個人バイアスを前記平均個人バイアスとして設定することと、

前記有利なクラスに関する前記個人バイアスの前記差の降順で前記訓練セット内のサンプルをソートすることと、

によって訓練される、請求項 1 または 2 に記載の後処理コンピュータ実装方法。

【請求項 4】

前記バイアス除去された予測を前記提案することは、

訓練セット内の保護属性を摂動させることと、

前記摂動の結果を顧客モデルにかけることと、

前記顧客モデルにかけられた前記摂動の結果に対する最も尤度の高い予測を採取することと、

この結果が有利なクラスに属する場合、非特権グループ・メンバーに対して前記摂動の結果に対する前記最も尤度の高い予測に変更することと、

によって実行される、請求項 1 ~ 3 のいずれか 1 項に記載の後処理コンピュータ実装方法。

【請求項 5】

前記バイアス検出器によって予測された非特権グループの中から検出された個人バイアス・サンプルが、前記提案することによって優先的に訂正され、

前記提案することは、摂動させたサンプルを顧客モデルにかけ、最も尤度の高い予測を採取することによって、前記バイアスを正すための訂正を実行し、

ユーザは、前記バイアスの元の値、または前記提案されたバイアス除去された予測のいずれを選ぶかを決定する、

請求項 1 ~ 4 のいずれか 1 項に記載の後処理コンピュータ実装方法。

【請求項 6】

前記摂動は保護属性にわたって実行され、集約結果が決定される、請求項 5 に記載の後処理コンピュータ実装方法。

【請求項 7】

複数のクラスの中から結果を選ぶことは、

摂動後の各クラスの集約予測を確認することと、

前記摂動後の最も尤度の高い予測結果を見つけることと、

の一方によって行われる、請求項 5 または 6 に記載の後処理コンピュータ実装方法。

【請求項 8】

前記訓練中に、前記バイアス検出器は、

非特権グループ・サンプルのペイロード・データ内の保護属性を摂動させ、

前記摂動させたデータと元のデータとの有利な結果の確率の差を求めることにより、非特権グループ・サンプルの個人バイアス・スコアを計算する個人バイアス・チェックを実装することと、

前記所定の個人バイアス閾値よりも大きい前記個人バイアスを有する前記非特権グループ・サンプルにフラグを立てることと、

前記フラグが立てられたサンプルとフラグが立てられていないサンプルとを区別するように前記バイアス検出器を訓練することと、

によって訓練される、請求項 1 ~ 7 のいずれか 1 項に記載の後処理コンピュータ実装方法。

【請求項 9】

前記実行時に、

前記適用することは、前記実行時に前記バイアス検出器を各非特権グループ・サンプルに適用して、前記個人バイアスの尤度を計算し、

10

20

30

40

50

保護属性を摂動させ、摂動後の結果をチェックすることにより、前記個人バイアスがかかったサンプルをさらにテストし、

前記摂動後の前記結果が元の結果と異なる場合、個人バイアスがかかったサンプルをバイアス除去された予測としてアービタに提案し、前記アービタは、前記バイアスの元の値と、バイアス除去された予測とから選ぶことができる、

請求項 1 ~ 8 のいずれか 1 項に記載の後処理コンピュータ実装方法。

【請求項 1 0】

クラウド・コンピューティング環境で具現化される、請求項 1 ~ 9 のいずれか 1 項に記載の後処理コンピュータ実装方法。

【請求項 1 1】

インスタンスレベルおよびグループレベルの予測メトリックの事後改善のための後処理コンピュータ・プログラムであって、前記コンピュータ・プログラムは、コンピュータに

グループ・バイアスへの制約付きで、所定の個人バイアス閾値よりも大きい個人バイアスを有するサンプルを検出することを学習するバイアス検出器を訓練することと、

実行時に、前記所定の個人バイアス閾値よりも大きい個人バイアスを有する実行時サンプル内のバイアスがかかったサンプルを選択するために、前記バイアス検出器を前記実行時サンプルに適用することと、

前記実行時に、前記バイアスがかかったサンプルに対するバイアス除去された予測を提案することと、

を実行させる、後処理コンピュータ・プログラム。

【請求項 1 2】

インスタンスレベルおよびグループレベルの予測メトリックの事後改善のための後処理システムであって、

プロセッサと、  
メモリと、

を備え、前記メモリは前記プロセッサに、

グループ・バイアスへの制約付きで、所定の個人バイアス閾値よりも大きい個人バイアスを有するサンプルを検出することを学習するバイアス検出器を訓練することと、

実行時に、前記所定の個人バイアス閾値よりも大きい個人バイアスを有する実行時サンプル内のバイアスがかかったサンプルを選択するために、前記バイアス検出器を前記実行時サンプルに適用することと、

前記実行時に、前記バイアスがかかったサンプルに対するバイアス除去された予測を提案することと、

を実行させるための命令を記憶する、後処理システム。

【請求項 1 3】

インスタンスレベルおよびグループレベルの予測メトリックの事後改善のための後処理コンピュータ実装方法であって、

所定の個人閾値よりも大きい値を有する実行時サンプル内のサンプルを選択するために、訓練された検出器を前記実行時サンプルに適用することと、

前記サンプルに対する変更された予測を提案することと、

を含む、後処理コンピュータ実装方法。

【請求項 1 4】

前記適用することおよび前記提案することは、個人バイアスを有する前記サンプルを是正の対象として、個人およびグループの両方の公平性メトリックに基づいて前記サンプルの前記値を変更するようにする後処理で動作する、請求項 1 3 に記載の後処理コンピュータ実装方法。

【請求項 1 5】

前記訓練された検出器は、完全なブラックボックスとして受け取られる、請求項 1 3 または 1 4 に記載の後処理コンピュータ実装方法。

10

20

30

40

50

**【請求項 16】**

前記適用することは、検証セットにグラウンド・トゥース・クラス・ラベルを必要としない純粋な実行時の手法を含む、請求項 13 ~ 15 のいずれか 1 項に記載の後処理コンピュータ実装方法。

**【請求項 17】**

インスタンスレベルおよびグループレベルの予測メトリックの事後改善のための後処理コンピュータ実装方法であって、

個人バイアスについて、非特権グループの各サンプルをテストすることと、

前記サンプルが所定の個人バイアス閾値より大きいバイアスを有する場合、特権グループの結果を前記サンプルに割り当てることと、

を含む、後処理コンピュータ実装方法。

10

**【請求項 18】**

前記非特権グループは、前記非特権グループの分類に応じて前記個人バイアスに基づいて前記結果を受け取る、請求項 17 に記載の後処理コンピュータ実装方法。

**【請求項 19】**

前記特権グループ内のサンプルの前記結果は、前記割り当てることによって変更されない、請求項 17 または 18 に記載の後処理コンピュータ実装方法。

**【発明の詳細な説明】****【技術分野】****【0001】**

20

本発明は、一般に後処理方法に関連し、より詳細には、限定としてではなく、差別的効果 (disparate impact) のグループ公平性尺度を改善することを目的としたバイアス軽減アルゴリズムにおいてデータサンプルに優先順位を付けるために使用される個人バイアス検出器を介して後処理で個人およびグループの両方の公平性を高めるためのシステム、方法、および記録媒体に関する。

**【0002】**

従来、バイアスがかかった分類器の予測の公平性を高めるための後処理手法は、グループの公平性のみに対処している。公平性、無差別、および望ましくないバイアスは、人間の意思決定において常に懸念事項であったが、過去の人間の意思決定が、雇用、信用、および刑事司法などの危険度の高い用途における機械学習モデル用の訓練データとして現在使用されているので、ますます注目されている。バイアスを軽減しなければ、そのような意思決定に関して訓練されたモデルは、人間のバイアスを永続させ、増減させるので、安全でなく、信頼できない。最近では、検出、推定理論、および情報理論を使用して基本原理が定義された機械学習およびデータ・マイニングの文献では、アルゴリズムの公平性に関する活動が活発になっている。

30

**【0003】**

従来、意思決定における公平性には、グループの公平性および個人の公平性という 2 つの主要な概念がある。グループの公平性は、その最も広い意味において、母集団を保護属性によって定義されたグループに分割し、ある統計的尺度がグループ間で等しくなることを求める。異なる統計的尺度に関する多くの異なるグループの公平性の概念が存在し、そのような概念の 1 つは差別的効果である。個人の公平性は、その最も広い意味において、類似の個人が同様に扱われることを求める。グループの公平性のチェックは、統計的メトリックの簡単な計算であるが、個人の公平性のチェックは、多くの値を有する多くの保護属性が存在する場合により計算が複雑になり、モデルを使用してサンプルをスコアリングするにはコストがかかる。最近では、不平等指数に基づくグループおよび個人の両方の公平性に関する統一されたメトリックが提案されている。

40

**【0004】**

機械学習パイプラインには、望ましくないバイアスを軽減するための 3 つの可能な介入ポイント、すなわち、訓練データ、学習手順、および出力予測が含まれ、3 つの対応するクラスの、すなわち、前処理、処理中、および後処理のバイアス軽減アルゴリズムがある

50

。後処理アルゴリズムの利点は、訓練プロセスへのアクセスを必要としないので、実行時環境に適していることである。さらに、後処理アルゴリズムはブラックボックス方式で動作し、すなわち、モデルの内部、その派生物 ( d e r i v a t i v e s ) などにアクセスする必要がないので、任意の機械学習モデルに適用可能である。

#### 【 0 0 0 5 】

従来、バイアス軽減アルゴリズムの大多数はグループの公平性に対処しているが、個人の公平性に対処しているものは少数である。一部の前処理アルゴリズムは、グループおよび個人の両方の公平性に対処している。既存の後処理アルゴリズムは全て、グループの公平性のためだけのものである。

#### 【 0 0 0 6 】

したがって、グループおよび個人の両方の公平性を考慮した後処理バイアス軽減技法が当技術分野で必要とされている。さらに、バイアス軽減アルゴリズムの訓練中に、検証サンプルにグラウンド・トゥルース・クラス・ラベルを必要としない技法が当技術分野で必要とされている。

#### 【 0 0 0 7 】

したがって、効果的な説明方法を生成するように構成することができる攻略法が当技術分野で必要とされている。

#### 【 発明の概要 】

#### 【 0 0 0 8 】

後処理アルゴリズムの一般的な方法論は、サンプルのサブセットを取得し、グループ公平性要件を満たすように、予測されたラベルを適切に変更することである。後処理に関する興味深い観察は、メトリックが期待値であるので、グループ公平性要件を達成するように任意のサンプルを改変できることである。当技術分野における問題を考慮して、本発明者らは、個人の公平性の問題を有するまたは有し得るサンプルを選び、これにより、グループおよび個人の両方の公平性に一緒に対処することができる。当技術分野における問題を解決する本発明の手法の出発点は、個人バイアス検出器であり、これは、保護属性を変化させ、他の全ての特徴を一定のままにした場合に、モデル予測が変化するサンプルを見つけるものである。この技法で多くの効率化が実現されるが、それでも計算コストは高くなる。検出器を継続的に実行できないという制限を克服するために、本発明は、小さなポイント・セットで個人の公平性をチェックし、新しいサンプルに適用される分類器を訓練することによって、それらから汎化する。個人バイアスを有し得るサンプルは、予測されたラベルの変更が検討されるサンプルである。そうすることにより、本発明は、不確定性に焦点を当てることから個人バイアスに焦点を当てることへと当技術分野を改善する。

#### 【 0 0 0 9 】

例示的な実施形態では、本発明は、インスタンスレベルおよびグループレベルの予測メトリックの事後改善のための後処理コンピュータ実装方法であって、グループ・バイアスへの制約付きで、所定の個人バイアス閾値よりも大きい個人バイアスを有するサンプルを検出することを学習するバイアス検出器を訓練することと、所定の個人バイアス閾値よりも大きいバイアスを有する実行時サンプル内のバイアスがかかったサンプルを選択するために、バイアス検出器を実行時サンプルに適用することと、バイアスがかかったサンプルに対するバイアス除去された予測を提案することと、を含む、後処理コンピュータ実装方法を提供することができる。

#### 【 0 0 1 0 】

1 つまたは複数の他の例示的な実施形態は、コンピュータ・プログラム製品およびシステムを含む。

#### 【 0 0 1 1 】

当技術分野への本貢献をよりよく理解できるように、本発明の他の詳細および実施形態を以下に説明する。それにもかかわらず、本発明は、その適用において、説明に記載しているかまたは図面に示しているそのような詳細、表現、用語、図解、もしくは配置またはそれらの組み合わせに限定されない。むしろ、本発明は、記載したものに加えた実施形態

10

20

30

40

50

が可能であり、様々な方法で実施および実行することが可能であり、限定的なものとは見なされるべきではない。

#### 【 0 0 1 2 】

したがって、当業者は、本開示の基礎となる概念が、本発明のいくつかの目的を実行するための他の構造、方法、およびシステムの設計の基礎として容易に利用され得ることを理解するであろう。そのため、特許請求の範囲が、本発明の主旨および範囲から逸脱しない限り、そのような同等の構造を含むものと見なされることが重要である。

#### 【 0 0 1 3 】

本発明の態様は、図面を参照して、本発明の例示的な実施形態の以下の詳細な説明からよりよく理解されよう。

#### 【図面の簡単な説明】

#### 【 0 0 1 4 】

【図 1】後処理方法 1 0 0 の高レベルのフロー・チャートを例示的に示す図である。

【図 2】本発明の一実施形態による方法 1 0 0 の訓練フェーズを例示的に示す図である。

【図 3】本発明の一実施形態による方法 1 0 0 の実行時フェーズを例示的に示す図である。

。【図 4】本発明の一実施形態による方法 1 0 0 の第 1 のアルゴリズムを例示的に示す図である。

【図 5】本発明の一実施形態による方法 1 0 0 の第 1 の結果を例示的に示す図である。

【図 6】本発明の一実施形態による方法 1 0 0 の第 2 の結果を例示的に示す図である。

【図 7】本発明の一実施形態による方法 1 0 0 の第 3 の結果を例示的に示す図である。

【図 8】本発明の一実施形態による方法 1 0 0 の第 4 の結果を例示的に示す図である。

【図 9】本発明の一実施形態によるクラウド・コンピューティング・ノード 1 0 を示す図である。

【図 1 0】本発明の一実施形態によるクラウド・コンピューティング環境 5 0 を示す図である。

【図 1 1】本発明の一実施形態による抽象化モデル・レイヤを示す図である。

#### 【発明を実施するための形態】

#### 【 0 0 1 5 】

#### 【数 1】

$$\hat{y}$$

#### 【 0 0 1 6 】

は以降  $y$  ハットと記載する。

#### 【 0 0 1 7 】

#### 【数 2】

$$\check{y}$$

#### 【 0 0 1 8 】

は以降  $y$  キャロン ( c a r o n ) と記載する。

#### 【 0 0 1 9 】

#### 【数 3】

$$\mathfrak{b}$$

#### 【 0 0 2 0 】

は以降  $b$  ドットと記載する。

#### 【 0 0 2 1 】

以下、図 1 ~ 図 1 1 を参照して本発明を説明し、図全体を通して、同様の参照番号は同様の部分を指す。一般的な慣行に準じて、図面の様々な特徴は必ずしも縮尺通りではないことを強調しておく。それどころか、様々な特徴の寸法は、明確にするために任意に拡大

10

20

30

40

50

または縮小する場合がある。

#### 【0022】

ここで図1に示す例を参照すると、後処理方法100は、任意の確率的機械学習モデル（またはモデルの融合）に適用可能なブラックボックス手法である、実行時環境に適した様々なステップを含み、また、個人およびグループの両方の公平性メトリックを改善する。

#### 【0023】

図9に示すように、本発明の一実施形態によるコンピュータ・システム12の1つまたは複数のコンピュータは、図1のステップを実行するための命令がストレージ・システムに記憶されたメモリ28を含むことができる。

10

#### 【0024】

1つまたは複数の実施形態（たとえば、図9～図11を参照）は、クラウド環境50（たとえば、図10を参照）に実装され得るが、それにもかかわらず、本発明はクラウド環境以外に実装できることは理解されよう。

#### 【0025】

図1～図8全体に関連して、本発明では、特徴 $X$ 、カテゴリー的保護属性 $D$ 、およびカテゴリー・ラベル $Y$ を有する教師あり分類問題を考える。訓練サンプルのセット $\{(x_1, d_1, y_1), \dots, (x_n, d_n, y_n)\}$ が与えられると、本発明は分類器（すなわち、ステップ101） $\hat{y}$ ハット： $X \times D \rightarrow Y$ を学習する。

#### 【0026】

「分類器」および「バイアス検出器」は同じ意味で使用していることに留意されたい。

20

#### 【0027】

説明を簡単にするために、スカラーのバイナリ保護属性（すなわち、 $D = \{0, 1\}$ ）と、バイナリ分類問題（すなわち、 $Y = \{0, 1\}$ ）とを考える。値 $d = 1$ は、「特権グループ」（たとえば、刑事司法適用における米国の第1のグループ）に対応するように設定され、 $d = 0$ は、「非特権グループ」（たとえば、第1のグループと比較して否定的な扱いを受ける、またはビジネスの場で融資を受けられたり受けられなかったりする、刑事司法適用における米国の第2のグループ）に対応するように設定される。値 $y = 1$ は「有利な結果」（たとえば、融資を受けられたり、逮捕されたりしないこと）に対応するように設定され、 $y = 0$ は「不利な結果」（たとえば、融資を受けられなかったり、逮捕されたりすること）に対応するように設定される。

30

#### 【0028】

コンテキストに基づいて、連続出力スコア $\hat{y}_s$ ハット： $[0, 1]$ を有する確率的バイナリ分類器も、 $\{0, 1\}$ への閾値と共に使用される。一方、本発明は、複数の保護属性およびマルチクラスの結果の場合に拡張可能である。実際に、複数のクラスが、複数の「有利」または「不利」な結果を有し得る。

#### 【0029】

個人バイアスの定義の1つは次のとおりである。 $\hat{y}_s$ ハット（ $x_i, d = 0$ ） $\hat{y}_s$ ハット（ $x_i, d = 1$ ）の場合、サンプル $i$ は個人バイアスを有する。 $b_i = I[\hat{y}_s$ ハット（ $x_i, d = 0$ ） $\hat{y}_s$ ハット（ $x_i, d = 1$ ）]とし、ここで、 $I[\cdot]$ は指示関数である。個人バイアス・スコア $b_{s,i} = \hat{y}_s$ ハット（ $x_i, d = 1$ ） $\hat{y}_s$ ハット（ $x_i, d = 0$ ）は、 $b_i$ のソフト・バージョンである。個人バイアス要約統計量を計算するために、テスト・サンプルにわたって $b_i$ の平均が取られる。

40

#### 【0030】

「差別的効果」として知られているグループの公平性の1つの概念は、次のように定義される。式（4）が1未満であるか、または $(1 - \text{式（4）}) - 1$ より大きい場合、「差別的効果」が存在し、ここで、式（4）の一般的な値は0.2である。

#### 【0031】

【数 4】

$$\frac{E[\hat{y}(X, D) \mid D = 0]}{E[\hat{y}(X, D) \mid D = 1]} \quad (4)$$

【0032】

個人バイアス検出用のテスト生成に関して、個人バイアス検出における2つの特有の問題がある。1つ目は、個人バイアスのケースの有無を判定することである。2つ目は、全サンプルの個人バイアス・ステータスを決定することである。

【0033】

本発明はこれに対し、1つ目の問題に関して、任意のブラック・ボックス分類器（たとえば、顧客モデル）の決定空間を体系的に探索して、バイアスがかかっている可能性が高いテスト・サンプルを生成する技法の改良を含む。

【0034】

この方法は、2種類の検索を使用する。1つ目はグローバル検索であり、これにより、様々な領域が網羅されるように決定空間が探索される。2つ目はローカル検索であり、これにより、発見済みの個人バイアスがかかったサンプルの非保護特徴の値を巧みに摂動（*p e r t u r b*）させることによって、テスト・ケースが生成される。重要なアイデアの一例は、動的シンボリック実行、すなわち、プログラム・パス内の制約を無効化して検索制約を生成し、制約ソルバを使用して新しい検索パスを見つける、プログラム用の既存の体系的なテスト・ケース生成技法を使用することである。このアルゴリズムは、多数の属性および属性値を含む状況でサンプルのバッチに使用された場合に、計算の観点から2つ目の特有の問題を解決するのに役立つ。

【0035】

しかしながら、許容可能なグループの公平性を達成するために、分類器  $y$  ハット<sub>*i*</sub> のラベル出力を他のラベル  $y$  キャロン<sub>*i*</sub> に変更する様々な後処理方法が適用され得る。棄却オプション分類（*R O C : r e j e c t o p t i o n c l a s s i f i c a t i o n*）法は、あるマージン・パラメータ  $\epsilon$  について、 $|y \text{ ハット}_i - 0.5| < \epsilon$  となる不確定サンプルを考え（ $0.5$  が分類閾値であると仮定する）、 $d_i = 0$  のサンプルには  $y \text{ キャロン}_i = 1$  を割り当て、 $d_i = 1$  のサンプルには  $y \text{ キャロン}_i = 0$  を割り当てる。いわゆる「棄却オプション」バンド外の特定のサンプルについては、 $y \text{ キャロン}_i = y \text{ ハット}_i$  である。 $\epsilon$  の値は、差別的効果に関する要件を達成するように最適化され得る。

【0036】

本発明は、公平性後処理アルゴリズムを含む。差別的効果のようなグループ公平性メトリックに関する後処理における重要な観察は、それらが期待値として定義されているので、個人のサンプルが交換可能であることである。本発明は、個人バイアスを有し得る  $X$  の部分からサンプルを選出する。本発明はこれを、まず個人バイアス検出を行い（たとえば、訓練されたバイアス検出器を使用）、次に後処理バイアス軽減アルゴリズムを使用する（たとえば、非特権グループの結果を有利なクラスの結果に変更する）ことによって行う。

【0037】

個人バイアス検出について、訓練データセット・パーティションで訓練済みの分類器  $y$  ハットを考える。本発明は、上記で提供した個人バイアスの定義を、それに従うラベルを有さない検証パーティションで評価することができる。これらの検証サンプルには個人バイアスを有するものと有さないものがある。 $X$  における個人バイアスのある程度の一貫性または平滑性の仮定の下で、本発明は、（たとえば、ステップ101において）この検証セットから個人バイアスの分類器または検出器（たとえば、バイアス検出器）を学習することができ、これは個人バイアスが未知である未見のサンプルへと汎化する。

【0038】

ここで、スコア出力を提供する任意の分類または異常検出アルゴリズムを使用し得る。



本発明では例示的に、経験的結果にロジスティック回帰を使用する。

【0039】

正式には、非特権グループ ( $d_j = 0$ ) に属する検証セット・サンプル ( $x_j, d_j$ )、 $j = 1, \dots, m$  の  $d_j$  を摂動させることによって、本発明は個人バイアス・スコア  $b_{s,j}$  を求める。方法 100 はさらに、データセット  $\{(x_1, d_1), \dots, (x_m, d_m)\}$  を構築し、これを使用して個人バイアス検出器  $b$  ドット ( $\cdot$ ) を訓練する。最大の個人バイアス・スコアを有するサンプルの場合、 $d_j$  は 1 であり、残りの場合は 0 である。この割り当ては、検証セット全体に対する差別的効果の制約に基づいて選ばれた個人バイアス・スコアに対する閾値 (たとえば、「所定の個人バイアス閾値」) によって決定される。これは、差別的効果の要件に基づいてマージン・パラメータが調整されるアルゴリズムに似ている。

10

【0040】

たとえば、図 2 は、バイアス検出器の訓練段階を例示的に示している。訓練では、ペイロード・データ 250 (たとえば、サンプル) を受け取り、これらは、摂動 240 を使用する (すなわち、検証セット・サンプルの  $d_j$  を摂動させる) 個人バイアス・チェッカ 201 にかけて、個人バイアス・スコアが求められる。チェッカ 201 は、結果に対する制約である顧客モデル 220 と協議する (たとえば、誰が承認されるべきか、または承認されないべきか)。顧客モデルは、見込み客 (銀行、政府、企業など) から受け取った例示的なブラックボックス・データ・セットである。データセット  $\{(x_1, d_1), \dots, (x_m, d_m)\}$  を使用して、汎化された個人バイアス・チェッカ 201 からの個人バイアス・スコアの入力を用いて、バイアス・サンプル検出器 202 を訓練する。すなわち、個人バイアス・スコアは、摂動 240 によって訓練セット内の保護属性 (複数可) を摂動させ、各サンプル・ポイントについて、場合によっては複数回の摂動後に平均スコアを求め、有利なクラスに関するスコアの差を求め、ここで、特権グループの場合、元のスコアおよび平均摂動後スコアを求め、「非特権グループ」の場合、元のスコアに対する平均摂動後スコアを求め、スコア差の降順で (個人バイアスの高いものから低いものへ) サンプルをソートすることによって求められる。

20

【0041】

上述のように、グループ公平性制約 230 を使用してバイアス・サンプル検出器 202 を訓練し、ここで、最大の個人バイアス・スコアを有するサンプルの場合、 $d_j$  は 1 であり、残りの場合は 0 である。この割り当ては、検証セット全体に対する差別的効果の制約に基づいて選ばれた個人バイアス・スコアに対する閾値 (たとえば、「所定の個人バイアス閾値」) によって決定される。訓練されたバイアス・サンプル検出器 202 は、ソートされたサンプルからバイアスがかかったサンプル・ポイントの最適なセットを識別し、上から開始して各ポイントをバイアス除去し、グループ・バイアスをチェックすることを許容レベルに達するまで行い、バイアス除去された標本 (example) に +1 のラベルを付け、バイアス除去されなかった標本に「-1」のラベルを付ける。本発明では、非特権グループの標本のみにはラベルを付ける。特権グループの標本は考慮しない。これらのラベル付けされた標本から、バイアスがかかったサンプルを検出する分類器 (最近傍、ロジスティック回帰、ランダム・フォレストなど) が訓練される。この分類器は、ソフト・バイアス尤度を提供可能でなければならない。持続検出器 (persist detector) 260 は、バイアス検出器モデルを将来の使用のためにメモリに保存する。

30

40

【0042】

これにより、訓練の際に、任意のサンプルを同等に改変することにより、グループ公平性メトリックを改善できるので、個人バイアスを有し得るサンプルに焦点を当てることができる。これは、計算コストの高い個人バイアス・チェッカから汎化して、個人バイアスを有し得る新しいサンプルを識別するモデルを作成し、まず、これらのサンプルを改変してグループ公平性メトリック要件を達成するようにする。

【0043】

なお、訓練された個人バイアス検出器は不要であるという主張もあり得、全てのサンプ

50

ルについて、実行時に来たときに単に  $b_i$  を計算すればよく、その理由は、そうするためには、ブラックボックス分類器モデルを使用したスコアリングしか必要ないためである。これは真実であり得るが、以下の注意点がある。第一に、 $d_i$  がスカラーかつバイナリであることが仮定されており、多くの場合はそうではない。そのため、 $b_i$  を計算するには複数回のモデル評価が必要になり得、これは特に、スコアリングされるサンプルごとに、モデルを配備してバイアスを是正するエンティティによって支払われる一定の金額がかかることが想定される産業用途では、法外になり得る。第二に、本発明は、グループ公平性制約に基づいてバイナリの  $b_i$  値を計算し、これによって、最大の個人バイアス・スコアを有する標本のみに対してバイアス除去が検討され、過剰な支払いがなくなる。全ての標本が  $b_i = 1$  で等しくバイアスがかかっていると見なされる場合、このレベルの制御は不可能である。

10

#### 【0044】

個人バイアス検出器  $b$  ドットが検証セットで訓練されると（すなわち、ステップ 101）、実行時にバイアス軽減アルゴリズムが適用されて、次のようにサンプルがテストされる（すなわち、ステップ 102）。非特権グループ（ $d_i = 0$ ）からの各サンプルは、個人バイアスのテストを受け、個人バイアスを有し得る場合（すなわち、 $b \text{ ドット } i = 1$ ）、このサンプルには、それが有利なクラスにあれば受け取ったはずの結果が割り当てられる（すなわち、 $y \text{ キャロン } i = y \text{ ハット } (x_{k,i})$ ）。発明と同様の人間の感受性をエンコードするために、特権グループからのサンプルを含め、他の全てのサンプルは未変更のままにする。図 4 に示すアルゴリズムは、上記の説明を例示的にまとめたものである。

20

#### 【0045】

図 3 に示すように、実行時環境（すなわち、ステップ 102 ~ 103）には、訓練済みの顧客モデルに加え、本発明のバイアス検出器モデルが既に存在する。定期的な間隔で、個人バイアス・チェックがテスト・サンプルに対して実行され、特徴空間全体に汎化される。新しいラベルのないサンプルが入ってくると、モデルはそれをスコアリングし、汎化された個人バイアス・チェックは、それが個人バイアスを有するか否かを予測する。これらの予測は後処理バイアス軽減アルゴリズムに供給され、これは、個人およびグループの両方の公平性メトリックを達成するように（たとえば、摂動および提案 303 と、提案されたバイアス除去された予測 350 とを介して）スコアを改変する方法を決定する。最後のアービタ 340 は、元のモデル予測または改変された予測を最終出力予測 360 として選ぶオプションをユーザに与える。

30

#### 【0046】

すなわち、実行時のバイアス・サンプル検出器およびバイアス除去は、テスト標本 330 に対して、バイアス検出器 301 の適用を通じて、それらのソフト・バイアス尤度に従って降順で適用される。個人バイアス尤度が閾値よりも大きい場合（302 で判定される）、本発明はそれらの標本について 303 以降を続ける。本発明では、グループ・バイアスをチェックしない。このため、提案手法は、実行時に 1 つの標本にさえ適用することができる。

#### 【0047】

これは、「非特権グループ」の標本に対してのみ実行される。302 で識別されたサンプルから、上から開始して各ポイントをバイアス除去し、グループ・バイアスをチェックすることを、許容レベルに達するまで行う。バイアス除去された予測 350 は後処理で、各サンプル・ポイントに対するバイアス除去手順によって実行され、この手順は、訓練セット内の保護属性（複数可）を摂動させ（たとえば、図 3 の 303）、摂動させた標本を顧客モデル 320 にかけ、摂動させたデータに対する最も尤度の高い予測を修正すべき提案値として採取することによって行われる。

40

#### 【0048】

これにより、検出器によって（「非特権グループ」の中で）最大の個人バイアスを有すると予測されたサンプルが優先的に訂正され、提案される訂正は、摂動させた標本を顧客モデルにかけ、最も尤度の高い予測を採取することを含み、アービタは、元の予測または

50

提案されたバイアス除去された予測のいずれを選択するかを決定することができる。また、本発明は非特権グループのみに焦点を当てているので、特権グループについては、本発明は個人バイアスを全く計算しない。

【 0 0 4 9 】

一実施形態では、摂動が保護属性にわたって実行され、集約スコア / 結果が決定されることを除いて、動作は上記と同様である。他の実施形態では、複数のクラスの場合、本発明がバイナリ・クラスではなく複数のクラスの中から結果を選ぶことを除いて、動作は同様である。これはいくつかの方法で、すなわち、摂動後の各クラスの集約予測スコアを確認し、摂動後の最も尤度の高い予測結果を見つけること、または任意のユーザ定義の方法で、実行することができる。「クラス」は「特権または非特権」であり得る。

10

【 0 0 5 0 】

アルゴリズムの公平性はビジネスおよび社会にとって重要なトピックであり、できるだけ多くの公平性の側面に対処する新しいバイアス軽減アルゴリズムを開発することは重要である。本明細書に開示する本発明は、個人およびグループの両方の公平性メトリックを改善するために、個人バイアスを有するサンプルを是正の対象とする新しい後処理アルゴリズムを含み、実験によりいくつかの実世界のデータセットで分類精度をそれほど損なうことなくそうすることを示した。機械学習業界は、モデル構築とモデル配備とが分離されるパラダイムに向かって進んでいる。これには、配備者が訓練済みのモデルの内部にアクセスするための能力が限られていることが含まれる。そのため、後処理アルゴリズム、特に分類器を完全なブラックボックスとして扱うことができるアルゴリズムが必要である。以前の技術と比較して、方法 1 0 0 は、個人およびグループの両方の公平性に取り組むだけでなく、検証セットにグラウンド・トゥルース・クラス・ラベルを必要としないので、純粋な実行時の手法にもなる。

20

【 0 0 5 1 】

「クラス」は通常、結果を表すことに留意されたい。たとえば、融資が承認された場合は「有利」と見なされ、融資が拒否された場合は「不利」と見なされる。「グループ」は、保護されたグループなどの人口統計群を表す。たとえば、1 つ目の人口統計群を「特権グループ」とすることができ、2 つ目の人口統計群を「非特権グループ」とすることができる。

【 0 0 5 2 】

本発明の上記の説明を考慮すると、本発明は、機械学習モデルからの予測に関するインスタンスレベルおよびグループレベルの両方のメトリックを事後的に改善し得る。グループレベルのメトリックの一例は、グループ・バイアスである。グループレベルおよびインスタンスレベルのメトリックは、たとえば、ある程度の効用を各個人に提供するシステムを含むことができ、ユーザは、適格な個人に対するあるレベルの保証に加え、集約されたグループ・レベルでのある程度の保証を確保することを望む。この場合、本発明は提案手法の変形を使用することができる。一般的に、グループ・レベルでの平均的な保証によって、適格な個人ごとの保証は確保されない。

30

【 0 0 5 3 】

本発明は、グループ・バイアスへの制約付きで、高い個人バイアスを有するサンプルを検出することを学習するバイアス検出器を訓練することによって、インスタンスレベルおよびグループレベルのメトリックを改善する。次に、本発明は、バイアス検出器を新規の実行時サンプルに適用して、バイアスを有し得るサンプルを選び、それらに対するバイアス除去された予測を提案する。

40

【 0 0 5 4 】

バイアス検出器の訓練段階の間に、本発明は、非特権グループ・サンプルのペイロード・データの保護属性を摂動させ、摂動させたデータと元のデータとの有利な結果の確率の差を求めることによって、それらの個人バイアス・スコアを計算する個人バイアス・チェッカを実装する。複数回の摂動の場合、有利な結果の平均確率が使用される。高い個人バイアスを有する非特権グループ・サンプルには「バイアスがかかったサンプル」のフラグ

50

が立てられ、グループ・バイアス制約が満たされるまでこの処理が繰り返される。バイアス検出器は、フラグが立てられたサンプルとフラグが立てられていないサンプルとを区別するように訓練される（たとえば、図2を参照）。

#### 【0055】

実行時/適用段階の間に、本発明は、実行時に各非特権グループ・サンプルにバイアス検出器を適用し、個人バイアスの尤度を計算する。最も尤度の高いバイアスがかかったサンプルが、さらなるテストのために選ばれる。さらなるテストは、保護属性を摂動させ、摂動後の予測をチェックすることを含む。複数回の摂動の場合、複数回の摂動の結果が単一の結果に集約される。摂動後の結果が元の結果と異なる場合、これはバイアス除去された予測としてアービタに提案され、アービタは、さらなる考慮事項に基づいて、元の予測と、バイアス除去された予測とから選ぶことができる（たとえば、図3を参照）。

10

#### 【0056】

本発明は、保護属性に対して複数回の摂動を実行することができる。このため、本発明は、数値的保護属性およびカテゴリー的保護属性の両方を扱うことができる。複数の保護属性を使用することもできる。また、本発明は、選択された結果が有利と見なされ、残りが不利と見なされる、一般的なカテゴリー的な結果を扱うことができる。

#### 【0057】

##### 実験結果

上記の方法100は、UCI Adult（1994年の米国国勢調査データベースに基づく収入データセット；45,222サンプル；有利な結果：\$50,000を超える収入；保護属性：性別、人種）、UCI Statlog German Credit（クレジット・スコアリング・データセット；1,000サンプル；有利な結果：低リスク；保護属性：性別、年齢）、およびProPublica COMPAS（刑務所の常習性データセット；6,167サンプル；有利な結果：再犯しない；保護属性：性別、人種）の3つの標準データセットで評価する。3つのデータセットのそれぞれは、2つのバイナリ保護属性を有し、我々はこれらを2つの異なる問題と見なすので、全部で6つの問題が生じる。方法100の結果は、棄却オプション分類（ROC）技法、および等化オッズ後処理（EOP: equalized odds post-processing）技法と比較される。

20

#### 【0058】

これらのデータセットは、例示的なAI Fairness 360ツールキットを使用して処理およびロードされ、60%の訓練パーティション、20%の検証パーティション、および20%のテスト・パーティションにランダムに分割される。これらの実験は、データセットの25個のそのようなランダムなパーティションを用いて実行され、以下の実験結果にエラー・バーを設けることが可能になる。訓練パーティションを使用して、実験では、ブラックボックス分類器として、 $L_2$ 正則化ロジスティック回帰と、ランダム・フォレストとの両方を組み込む。ランダム・フォレストの場合、実験では、ツリーの数を100個に設定し、リーフ・ノードあたりの最小サンプル数を20個に設定する。

30

#### 【0059】

3つのバイアス軽減手法全てのパラメータは、データセットの検証パーティションで最適化される。ROCおよびEOPの両方の手法では、検証セットにグラウンド・トゥールース・クラス・ラベルが必要であるが、純粋な実行時の方法である方法100では必要ない。ROCおよび方法100は、範囲（0.8, 1.25）で、すなわち、 $\epsilon = 0.2$ で、「差別的効果」を達成するように最適化される。「差別的効果」ではなく等化オッズ用に設計されたEOPは、差別的効果の範囲に対して最適化することはできない。

40

#### 【0060】

以下のサブセクションでは、実験は最初に、提案方法100で使用する個人バイアス検出器の有効性を実証し、その後、分類精度、差別的効果、および個人の公平性について3つのアルゴリズムを比較する。

#### 【0061】

50

図 6 ~ 図 8 において、「IGD」ラベルは方法 100 の結果に対応することに留意されたい。

#### 【0062】

個人バイアスの汎化に関する検証結果の図を提供する図 5 を参照すると、実験では、未見のテスト・データに対する個人バイアス検出器の汎化性能を検証する。個人バイアス検出器は、非特権グループ・サンプル ( $d = 0$ ) のみで使用されるので、その性能尺度はこのサブセットのみについて計算される。バイアス検出器のグラウンド・トゥールズ・ラベルは、テスト・データ内の全ての「非特権グループ」のサンプルに対して個人バイアス・スコア ( $b_{s,k}$ ) を実際に計算し、「差別的効果」の制約に基づいてグラウンド・トゥールズ・バイアス・ラベル ( $\hat{b}_k$ ) を識別することによって、得られる。これらのラベルは、バイアス検出器 ( $b^k$ ) によって予測されたラベルと比較され、バランス分類精度が計算される。

10

#### 【0063】

バイアス検出器のこの性能を、ブラックボックス分類器がロジスティック回帰である場合の全てのデータセットおよび保護属性の組み合わせについて、図 5 に示す。全ての精度値は 0.85 を超えており、これは目下の目的に対する明確な有効性を示している。検出器は、ブラックボックス分類器がランダム・フォレストの場合も同様に機能し、最小精度は約 0.80 である。

#### 【0064】

公平性の比較のために、実験では、EOP、ROC、および方法 100 を比較するために、(a) 個人バイアス、(b) 差別的効果、および (c) バランス分類精度、の 3 つの尺度を使用する。これらの尺度は、後処理された予測  $y$  キャロンを使用して計算される。個人バイアス尺度は、上記の段落 0038 で説明した要約統計量であり、差別的効果尺度は式 (4) で定義され、バランス分類精度は、真のラベル  $y$  に対する予測  $y$  キャロンに関して求められた真陽性率および真陰性率の平均である。実験では、元の (オリジナルの) 予測  $y$  ハットに関するこれらの尺度も求める。図 6 ~ 図 8 に示すように、方法 100 は、元の分類器に近い精度を維持しながら、両方の公平性尺度を一貫して改善する唯一の方法である。全ての結果は、ブラックボックス分類器としてのロジスティック回帰に関して示しているが、ランダム・フォレストでも同様の結果が観察される (スペースの制約のため省略している)。

20

30

#### 【0065】

実験に基づいて、個人バイアスでは、方法 100 は German データセットおよび COMPAS データセットで最高性能を発揮している。ROC 法は、Adult データセットに対して最高性能を発揮しているが、バランス精度を低下させるという犠牲を払っている。EOP 法および ROC 法は、個人バイアスを増加させることがあるが、方法 100 には全く当てはまらない。方法 100 はまた、元の予測に対して差別的効果を一貫して改善しており、ただし、6 つのケースのうち 5 つで ROC 法が上回っている。ROC 手法の強力な性能は、個人バイアスも最適化していないためと考えられる。EOP 法は、差別的効果に関して性能が低く、これはおそらく、オッズを等化するように設計されていたために、「差別的効果」が改善される場合と、常に改善されない場合とがあるためである。

40

#### 【0066】

また、方法 100 は、検証パーティションでグラウンド・トゥールズ・ラベルが使用されていないにもかかわらず、元の予測と比較してバランス分類器精度を維持するのに最適である。

#### 【0067】

クラウド・コンピューティング環境を使用した例示的な態様

この詳細な説明はクラウド・コンピューティング環境における本発明の例示的な実施形態を含むが、本明細書に列挙した教示の実装形態はそのようなクラウド・コンピューティング環境に限定されないことを理解されたい。むしろ、本発明の実施形態は、現在知られているまたは今後開発される他の任意のタイプのコンピューティング環境と共に実装する

50

ことが可能である。

【 0 0 6 8 】

クラウド・コンピューティングは、最小限の管理労力またはサービスのプロバイダとのやり取りによって迅速に供給および公開することができる、設定可能なコンピューティング・リソース（たとえば、ネットワーク、ネットワーク帯域幅、サーバ、処理、メモリ、ストレージ、アプリケーション、仮想マシン、およびサービス）の共有プールへの便利かつオンデマンドのネットワーク・アクセスを可能にするためのサービス提供のモデルである。このクラウド・モデルは、少なくとも5つの特徴と、少なくとも3つのサービス・モデルと、少なくとも4つの配備モデルとを含み得る。

【 0 0 6 9 】

特徴は以下の通りである。

オンデマンド・セルフ・サービス：クラウド・コンシューマは、サービスのプロバイダとの人的な対話を必要とせずに、必要に応じて自動的に、サーバ時間およびネットワーク・ストレージなどのコンピューティング能力を一方的に供給することができる。

ブロード・ネットワーク・アクセス：能力はネットワークを介して利用することができ、異種のシンまたはシック・クライアント・プラットフォーム（たとえば、携帯電話、ラップトップ、およびPDA）による使用を促進する標準的なメカニズムを介してアクセスされる。

リソース・プーリング：プロバイダのコンピューティング・リソースをプールして、様々な物理リソースおよび仮想リソースが需要に応じて動的に割り当ておよび再割り当てされるマルチ・テナント・モデルを使用して複数のコンシューマにサービス提供する。一般にコンシューマは、提供されるリソースの正確な位置に対して何もできず、知っているわけでもないが、より高い抽象化レベル（たとえば、国、州、またはデータセンターなど）における位置を特定可能であり得るという点で位置非依存の感覚がある。

迅速な弾力性：能力を迅速かつ弾力的に、場合によっては自動的に供給して素早くスケール・アウトし、迅速に解放して素早くスケール・インすることができる。コンシューマにとって、供給可能な能力は無制限であるように見えることが多く、任意の時間に任意の数量で購入することができる。

測定されるサービス：クラウド・システムは、サービスのタイプ（たとえば、ストレージ、処理、帯域幅、およびアクティブ・ユーザ・アカウント）に適したある抽象化レベルでの計量機能を活用して、リソースの使用を自動的に制御し最適化する。リソース使用量を監視、制御、および報告して、利用されるサービスのプロバイダおよびコンシューマの両方に透明性を提供することができる。

【 0 0 7 0 】

サービス・モデルは以下の通りである。

ソフトウェア・アズ・ア・サービス（SaaS：Software as a Service）：コンシューマに提供される能力は、クラウド・インフラストラクチャ上で動作するプロバイダのアプリケーションを使用することである。アプリケーションは、ウェブ・ブラウザ（たとえば、ウェブベースの電子メール）などのシン・クライアント・インターフェースを介して様々なクライアント回路からアクセス可能である。コンシューマは、限定されたユーザ固有のアプリケーション構成設定を可能性のある例外として、ネットワーク、サーバ、オペレーティング・システム、ストレージ、さらには個々のアプリケーション機能を含む、基盤となるクラウド・インフラストラクチャを管理も制御もしない。

プラットフォーム・アズ・ア・サービス（PaaS：Platform as a Service）：コンシューマに提供される能力は、プロバイダによってサポートされるプログラミング言語およびツールを使用して作成された、コンシューマが作成または取得したアプリケーションをクラウド・インフラストラクチャ上に配備することである。コンシューマは、ネットワーク、サーバ、オペレーティング・システム、またはストレージを含む、基盤となるクラウド・インフラストラクチャを管理も制御もしないが、配備されたアプリケーションおよび場合によってはアプリケーション・ホスティング環境構成を制御

10

20

30

40

50

する。

インフラストラクチャ・アズ・ア・サービス (IaaS: Infrastructure as a Service): コンシューマに提供される能力は、オペレーティング・システムおよびアプリケーションを含むことができる任意のソフトウェアをコンシューマが配備して動作させることが可能な、処理、ストレージ、ネットワーク、および他の基本的なコンピューティング・リソースを供給することである。コンシューマは、基盤となるクラウド・インフラストラクチャを管理も制御もしないが、オペレーティング・システム、ストレージ、配備されたアプリケーションを制御し、場合によっては選択したネットワーク・コンポーネント (たとえば、ホスト・ファイアウォール) を限定的に制御する。

10

#### 【0071】

配備モデルは以下の通りである。

プライベート・クラウド: クラウド・インフラストラクチャは組織専用運用される。これは組織または第三者によって管理され得、構内または構外に存在し得る。

コミュニティ・クラウド: クラウド・インフラストラクチャはいくつかの組織によって共有され、共通の懸念 (たとえば、ミッション、セキュリティ要件、ポリシー、およびコンプライアンスの考慮事項など) を有する特定のコミュニティをサポートする。これは組織または第三者によって管理され得、構内または構外に存在し得る。

パブリック・クラウド: クラウド・インフラストラクチャは、一般大衆または大規模な業界団体に対して利用可能にされ、クラウド・サービスを販売する組織によって所有される。

20

ハイブリッド・クラウド: クラウド・インフラストラクチャは、固有のエンティティのままであるが、データおよびアプリケーションの移植性を可能にする標準化技術または独自技術 (たとえば、クラウド間の負荷分散のためのクラウド・パースティング) によって結合された2つ以上のクラウド (プライベート、コミュニティ、またはパブリック) を合成したものである。

#### 【0072】

クラウド・コンピューティング環境は、ステートレス性、低結合性、モジュール性、および意味論的相互運用性に重点を置いたサービス指向型である。クラウド・コンピューティングの中核にあるのは、相互接続されたノードのネットワークを含むインフラストラクチャである。

30

#### 【0073】

ここで図9を参照すると、クラウド・コンピューティング・ノードの一例の概略図が示されている。クラウド・コンピューティング・ノード10は、適切なノードの一例に過ぎず、本明細書に記載の本発明の実施形態の使用または機能性の範囲に関するいかなる制限も示唆することを意図していない。いずれにしても、クラウド・コンピューティング・ノード10は、本明細書に記載の機能性のいずれかを実装されること、もしくは実行すること、またはその両方が可能である。

#### 【0074】

クラウド・コンピューティング・ノード10はコンピュータ・システム/サーバ12として図示しているが、他の多くの汎用または専用のコンピューティング・システム環境または構成で動作可能であることを理解されたい。コンピュータ・システム/サーバ12での使用に適し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、もしくは構成、またはそれらの組み合わせの例には、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、ハンドヘルドもしくはラップトップ回路、マルチプロセッサ・システム、マイクロプロセッサベースのシステム、セット・トップ・ボックス、プログラム可能な家庭用電化製品、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、および上記のシステムもしくは回路のいずれかを含む分散型クラウド・コンピューティング環境などが含まれるが、これらに限定されない。

40

50

## 【0075】

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コンテキストで記述され得る。一般に、プログラム・モジュールには、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装するルーチン、プログラム、オブジェクト、コンポーネント、ロジック、データ構造などが含まれ得る。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たりモート処理回路によってタスクが実行される分散型クラウド・コンピューティング環境で実施され得る。分散型クラウド・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・ストレージ回路を含むローカルおよびリモート両方のコンピュータ・システム記憶媒体に配置され得る。

10

## 【0076】

図9を再度参照すると、コンピュータ・システム/サーバ12は、汎用コンピューティング回路の形態で示されている。コンピュータ・システム/サーバ12のコンポーネントは、1つまたは複数のプロセッサまたは処理ユニット16と、システム・メモリ28と、システム・メモリ28を含む様々なシステム・コンポーネントをプロセッサ16に結合するバス18と、を含み得るが、これらに限定されない。

## 【0077】

バス18は、メモリ・バスまたはメモリ・コントローラ、ペリフェラル・バス、加速グラフィックス・ポート、および様々なバス・アーキテクチャ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うちのいずれかの1つまたは複数を表す。限定ではなく例として、そのようなアーキテクチャには、業界標準アーキテクチャ(ISA: Industry Standard Architecture)バス、マイクロ・チャンネル・アーキテクチャ(MCA: Micro Channel Architecture)バス、拡張ISA(EISA: Enhanced ISA)バス、ビデオ・エレクトロニクス規格協会(VESA: Video Electronics Standards Association)ローカル・バス、および周辺機器相互接続(PCI: Peripheral Component Interconnects)バスが含まれる。

20

## 【0078】

コンピュータ・システム/サーバ12は、典型的には、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12によってアクセス可能な任意の利用可能な媒体であり得、揮発性および不揮発性の媒体、取り外し可能および取り外し不可能な媒体の両方を含む。

30

## 【0079】

システム・メモリ28は、ランダム・アクセス・メモリ(RAM: random access memory)30もしくはキャッシュ・メモリ32またはその両方などの、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12は、他の取り外し可能/取り外し不可能な揮発性/不揮発性のコンピュータ・システム記憶媒体をさらに含み得る。単なる例として、取り外し不可能な不揮発性の磁気媒体(図示せず、典型的には「ハード・ドライブ」と呼ばれるもの)に読み書きするためのストレージ・システム34を設けることができる。図示していないが、取り外し可能な不揮発性の磁気ディスク(たとえば、「フレキシブル・ディスク」)に読み書きするための磁気ディスク・ドライブと、CD-ROM、DVD-ROM、または他の光学媒体などの取り外し可能な不揮発性の光学ディスクに読み書きするための光学ディスク・ドライブと、を設けることができる。そのような例では、それぞれを、1つまたは複数のデータ・メディア・インターフェースによってバス18に接続することができる。以下でさらに図示および説明するように、メモリ28は、本発明の実施形態の機能を実行するように構成されるプログラム・モジュールのセット(たとえば、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み得る。

40

## 【0080】

50



プログラム・モジュール 42 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 40 は、限定ではなく例として、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 28 に記憶され得る。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・データまたはそれらの何らかの組み合わせのそれぞれは、ネットワーク環境の実装形態を含み得る。プログラム・モジュール 42 は、一般に、本明細書に記載の本発明の実施形態の機能もしくは方法論またはその両方を実行する。

#### 【0081】

コンピュータ・システム / サーバ 12 はまた、キーボード、ポインティング回路、ディスプレイ 24 などの 1 つまたは複数の外部デバイス 14、ユーザがコンピュータ・システム / サーバ 12 とやりとりすることを可能にする 1 つまたは複数の回路、ならびに / あるいはコンピュータ・システム / サーバ 12 が 1 つまたは複数の他のコンピューティング回路と通信することを可能にする任意の回路（たとえば、ネットワーク・カード、モデムなど）と通信し得る。そのような通信は、入力 / 出力（I / O : Input / Output）インターフェース 22 を介して行うことができる。またさらに、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 : local area network）、一般的なワイド・エリア・ネットワーク（WAN : wide area network）、もしくはパブリック・ネットワーク（たとえば、インターネット）、またはそれら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 の他のコンポーネントと通信する。図示していないが、他のハードウェアもしくはソフトウェアまたはその両方のコンポーネントを、コンピュータ・システム / サーバ 12 と併用できることを理解されたい。例には、マイクロコード、回路ドライバ、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・システムなどが含まれるが、これらに限定されない。

#### 【0082】

ここで図 10 を参照すると、例示的なクラウド・コンピューティング環境 50 が示されている。図示のように、クラウド・コンピューティング環境 50 は 1 つまたは複数のクラウド・コンピューティング・ノード 10 を含み、これらを使用して、たとえば、パーソナル・デジタル・アシスタント（PDA : personal digital assistant）もしくは携帯電話 54 A、デスクトップ・コンピュータ 54 B、ラップトップ・コンピュータ 54 C、または自動車コンピュータ・システム 54 N あるいはそれらの組み合わせなどの、クラウド・コンシューマによって使用されるローカル・コンピューティング回路が通信し得る。ノード 10 は互いと通信し得る。これらは、たとえば、上述のプライベート、コミュニティ、パブリック、もしくはハイブリッド・クラウド、またはそれらの組み合わせなどの 1 つまたは複数のネットワークにおいて、物理的または仮想的にグループ化され得る（図示せず）。これにより、クラウド・コンピューティング環境 50 は、クラウド・コンシューマがローカル・コンピューティング回路上にリソースを維持する必要がない、インフラストラクチャ・アズ・ア・サービス、プラットフォーム・アズ・ア・サービス、またはソフトウェア・アズ・ア・サービス、あるいはそれらの組み合わせを提供することが可能になる。図 10 に示すコンピューティング回路 54 A ~ N のタイプは例示的なものにすぎないことが意図されており、コンピューティング・ノード 10 およびクラウド・コンピューティング環境 50 は、任意のタイプのネットワークまたはネットワーク・アドレス指定可能接続（たとえば、Web ブラウザを使用）あるいはその両方を介して任意のタイプのコンピュータ化回路と通信できることを理解されたい。

#### 【0083】

ここで図 11 を参照すると、クラウド・コンピューティング環境 50（図 10）によって提供される機能的抽象化レイヤの例示的なセットが示されている。図 11 に示すコンポ

10

20

30

40

50

ーネット、レイヤ、および機能は例示的なものにすぎないことが意図されており、本発明の実施形態はこれらに限定されないことを事前に理解されたい。図示のように、以下のレイヤおよび対応する機能が提供される。

【 0 0 8 4 】

ハードウェアおよびソフトウェア・レイヤ 6 0 は、ハードウェア・コンポーネントおよびソフトウェア・コンポーネントを含む。ハードウェア・コンポーネントの例には、メインフレーム 6 1、RISC ( 縮小命令セット・コンピュータ : Reduced Instruction Set Computer ) アーキテクチャ・ベースのサーバ 6 2、サーバ 6 3、ブレード・サーバ 6 4、ストレージ回路 6 5、ならびにネットワークおよびネットワークング・コンポーネント 6 6 が含まれる。一部の実施形態では、ソフトウェア・コンポーネントは、ネットワーク・アプリケーション・サーバ・ソフトウェア 6 7 およびデータベース・ソフトウェア 6 8 を含む。

10

【 0 0 8 5 】

仮想化レイヤ 7 0 は、仮想エンティティの以下の例 : 仮想サーバ 7 1、仮想ストレージ 7 2、仮想プライベート・ネットワークを含む仮想ネットワーク 7 3、仮想アプリケーションおよびオペレーティング・システム 7 4、ならびに仮想クライアント 7 5 の提供元になり得る抽象化レイヤを提供する。

【 0 0 8 6 】

一例では、管理レイヤ 8 0 は、下記の機能を提供し得る。リソース供給 8 1 は、クラウド・コンピューティング環境内でタスクを実行するために利用されるコンピューティング・リソースおよび他のリソースの動的調達を提供する。計量および価格設定 8 2 は、クラウド・コンピューティング環境内でリソースが利用されるときのコスト追跡と、これらのリソースの消費に対する課金または請求とを提供する。一例では、これらのリソースはアプリケーション・ソフトウェア・ライセンスを含み得る。セキュリティは、クラウド・コンシューマおよびタスクの識別情報検証だけでなく、データおよび他のリソースに対する保護も提供する。ユーザ・ポータル 8 3 は、コンシューマおよびシステム管理者にクラウド・コンピューティング環境へのアクセスを提供する。サービス・レベル管理 8 4 は、要求されたサービス・レベルが満たされるような、クラウド・コンピューティング・リソースの割り当ておよび管理を提供する。サービス・レベル合意 ( S L A ) の計画および履行 8 5 は、S L A に従って将来要求されると予想されるクラウド・コンピューティング・リソースの事前手配および調達を提供する。

20

30

【 0 0 8 7 】

ワークロード・レイヤ 9 0 は、クラウド・コンピューティング環境が利用され得る機能性の例を提供する。このレイヤから提供され得るワークロードおよび機能の例は、マッピングおよびナビゲーション 9 1、ソフトウェア開発およびライフサイクル管理 9 2、仮想教室教育配信 9 3、データ分析処理 9 4、取引処理 9 5、および、より詳細には、本発明に関連して、方法 1 0 0、を含む。

【 0 0 8 8 】

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、もしくはコンピュータ・プログラム製品またはそれら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、本発明の態様をプロセッサに実行させるためのコンピュータ可読プログラム命令をその上に有するコンピュータ可読記憶媒体 ( または複数の媒体 ) を含み得る。

40

【 0 0 8 9 】

コンピュータ可読記憶媒体は、命令実行デバイスによる使用のために命令を保持および記憶できる有形のデバイスとすることができる。コンピュータ可読記憶媒体は、たとえば、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、電磁ストレージ・デバイス、半導体ストレージ・デバイス、またはこれらの任意の適切な組み合わせであり得る。コンピュータ可読記憶媒体のより具体的な例の非網羅的なリストには、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 ( R A M )、読み取り専用メモリ ( R O M : r e a d - o n

50

ly memory)、消去可能プログラム可能読み取り専用メモリ(EPR0M:erasable programmable read-only memory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:static random access memory)、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:portable compact disc read-only memory)、デジタル多用途ディスク(DVD:digital versatile disk)、メモリー・スティック(登録商標)、フレキシブル・ディスク、命令が記録されたパンチ・カードまたは溝の隆起構造などの機械的にコード化されたデバイス、およびこれらの任意の適切な組み合わせが含まれる。コンピュータ可読記憶媒体は、本明細書で使用する場合、たとえば、電波または他の自由に伝搬する電磁波、導波管もしくは他の伝送媒体を伝搬する電磁波(たとえば、光ファイバ・ケーブルを通過する光パルス)、または有線で伝送される電気信号など、一過性の信号自体であると解釈されるべきではない。

10

#### 【0090】

本明細書に記載のコンピュータ可読プログラム命令は、コンピュータ可読記憶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、たとえば、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、もしくは無線ネットワーク、またはそれらの組み合わせなどのネットワーク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することができる。ネットワークは、銅線伝送ケーブル、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、もしくはエッジ・サーバ、またはそれらの組み合わせを含み得る。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、ネットワークからコンピュータ可読プログラム命令を受信し、コンピュータ可読プログラム命令を転送して、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記憶媒体に記憶する。

20

#### 【0091】

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:instruction-set-architecture)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、集積回路の構成データ、あるいは、Smalltalk(登録商標)、C++などの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類似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・コードであり得る。コンピュータ可読プログラム命令は、完全にユーザのコンピュータ上で、部分的にユーザのコンピュータ上で、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上かつ部分的にリモート・コンピュータ上で、あるいは完全にリモート・コンピュータまたはサーバ上で実行し得る。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)またはワイド・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介してユーザのコンピュータに接続され得、または(たとえば、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介して)外部コンピュータに接続され得る。一部の実施形態では、たとえば、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:field-programmable gate array)、またはプログラマブル・ロジック・アレイ(PLA:programmable logic array)を含む電子回路は、本発明の態様を実行するために、電子回路を個人向けにするための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、コンピュータ可読プログラム命令を実行し得る。

30

40

#### 【0092】

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・プログラム製品のフローチャート図もしくはブロック図またはその両方を参照して本

50

明細書で説明している。フローチャート図もしくはブロック図またはその両方の各ブロック、およびフローチャート図もしくはブロック図またはその両方におけるブロックの組み合わせが、コンピュータ可読プログラム命令によって実装できることは理解されよう。

【0093】

これらのコンピュータ可読プログラム命令を、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ム可能データ処理装置のプロセッサに提供して、コンピュータまたは他のプログラム可能データ処理装置のプロセッサを介して実行されたそれらの命令が、フローチャートもしくはブロック図またはその両方の1つまたは複数のブロックにおいて指定された機能／行為を実装するための手段を生成するようなマシンを生成し得る。また、これら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記憶媒体に記憶され得、それによって、コンピュータ、プログラム可能データ処理装置、もしくは他のデバイス、またはそれらの組み合わせに特定の方法で機能するように指示することができ、その結果、命令が記憶されたコンピュータ可読記憶媒体は、フローチャートもしくはブロック図またはその両方の1つまたは複数のブロックにおいて指定された機能／行為の態様を実装する命令を含む製造品を含む。

10

【0094】

また、コンピュータ可読プログラム命令をコンピュータ、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して、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上で一連の動作ステップを実行させることによって、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上で実行されたそれらの命令が、フローチャートもしくはブロック図またはその両方の1つまたは複数のブロックにおいて指定された機能／行為を実装するようなコンピュータ実装処理を生成し得る。

20

【0095】

図中のフローチャート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装形態のアーキテクチャ、機能性、および動作を示している。これに関して、フローチャートまたはブロック図の各ブロックは、指定された論理的機能（複数可）を実装するための1つまたは複数の実行可能命令を含むモジュール、セグメント、または命令の一部を表し得る。一部の代替的実装形態では、ブロックに示す機能は、図に示す順序以外で行われ得る。たとえば、関与する機能性に応じて、連続して示す2つのブロックは、実際には実質的に同時に実行され得、またはそれらのブロックは、場合により逆の順序で実行され得る。ブロック図もしくはフローチャート図またはその両方の各ブロック、およびブロック図もしくはフローチャート図またはその両方におけるブロックの組み合わせは、指定された機能もしくは行為を実行する専用のハードウェア・ベースのシステムによって実装されるか、または専用ハードウェアおよびコンピュータ命令の組み合わせを実行することができることにも留意されたい。

30

【0096】

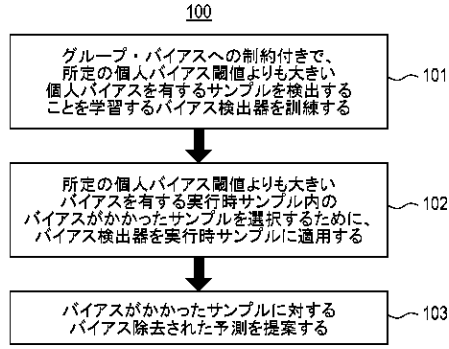
本発明の様々な実施形態の説明は、例示の目的で提示しているが、網羅的であることも、開示した実施形態に限定されることも意図したものではない。開示した実施形態の範囲および主旨から逸脱することなく、多くの修正および変形が当業者には明らかであろう。本明細書で使用している用語は、実施形態の原理、市場で見られる技術に対する実際の適用または技術的改善を最もよく説明するために、または当業者が本明細書に開示した実施形態を理解できるようにするために選択している。

40

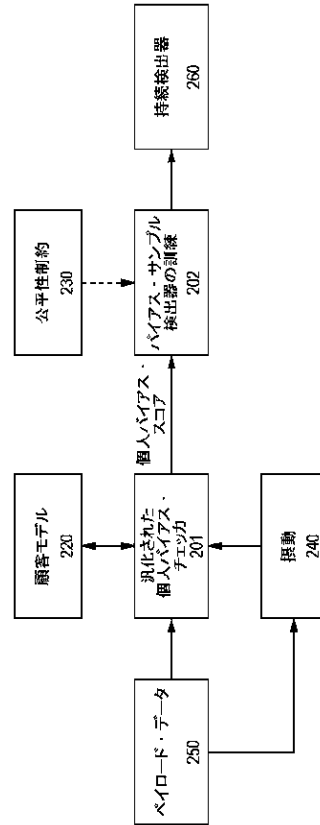
【0097】

さらに、出願人の意図は、全ての請求項の要素の均等物を包含することであり、本出願のいかなる請求項への補正も、補正した請求項の任意の要素または特徴の均等物に対する利益または権利の放棄として解釈されるべきではない。

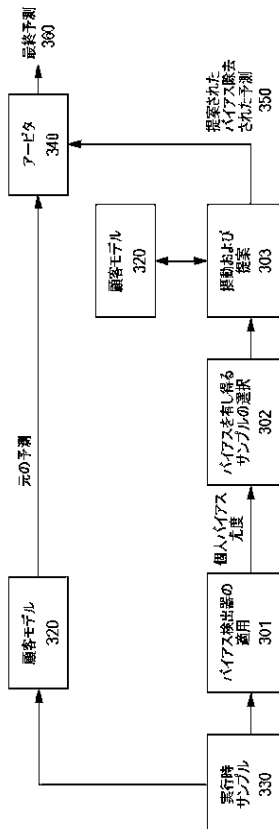
【図 1】



【図 2】



【図 3】

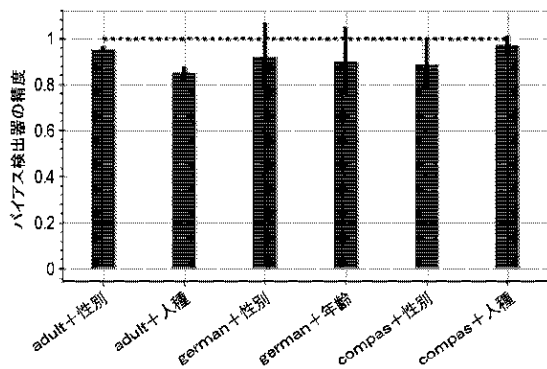


【図 4】

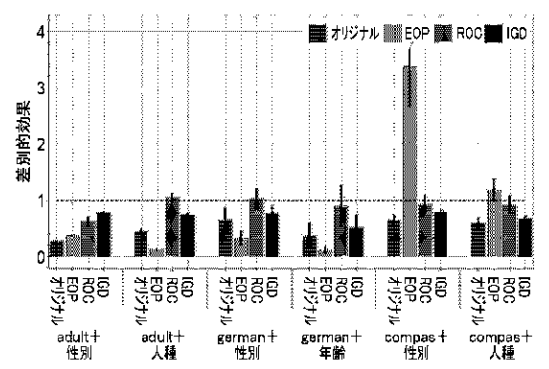
**Algorithm 1** Individual+Group Debiasing (IGD) Post-Processing

Given classifier  $\hat{y}$  trained on training set  $\{(\mathbf{x}_i, d_i, y_i)\}$ , and  
 Given validation set  $\{\mathbf{x}_j \mid d_j = 0\}$ , compute individual bias  
 scores  $\{b_{S,j} \mid d_j = 0\}$ .  
**if**  $b_{S,j} > \tau$  **then**  
    $\beta_j \leftarrow 1$   
**else**  
    $\beta_j \leftarrow 0$   
**end if**  
 Construct auxiliary dataset  $\{(\mathbf{x}_j, \beta_j) \mid d_j = 0\}$ .  
 Train individual bias detector  $\hat{b}$  on auxiliary dataset.  
**for all** run-time test samples  $(\mathbf{x}_k, d_k)$  **do**  
    $\hat{y}_k \leftarrow \hat{y}(\mathbf{x}_k, d_k)$   
   **if**  $d_k == 0$  **then**  
      $\hat{b}_k \leftarrow \hat{b}(\mathbf{x}_k)$   
     **if**  $\hat{b}_k == 1$  **then**  
        $\hat{y}_k \leftarrow \hat{y}(\mathbf{x}_k, 1)$   
     **else**  
        $\hat{y}_k \leftarrow \hat{y}_k$   
     **end if**  
   **else**  
      $\hat{y}_k \leftarrow \hat{y}_k$   
   **end if**  
**end for**

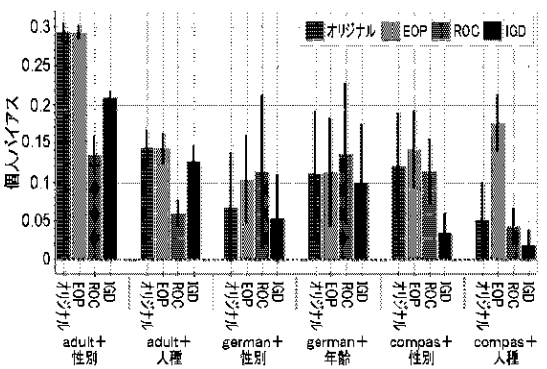
【図 5】



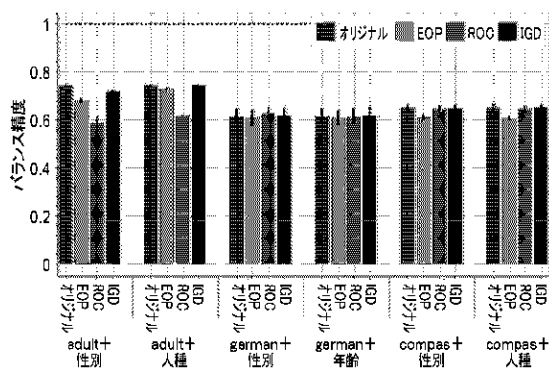
【図 7】



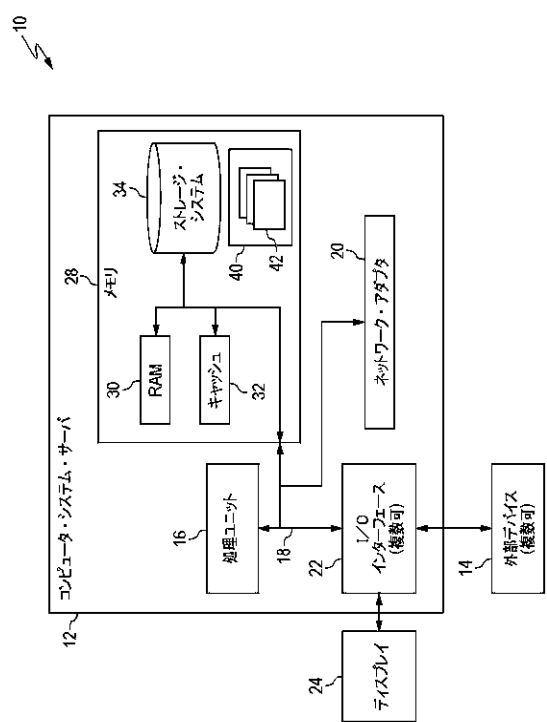
【図 6】



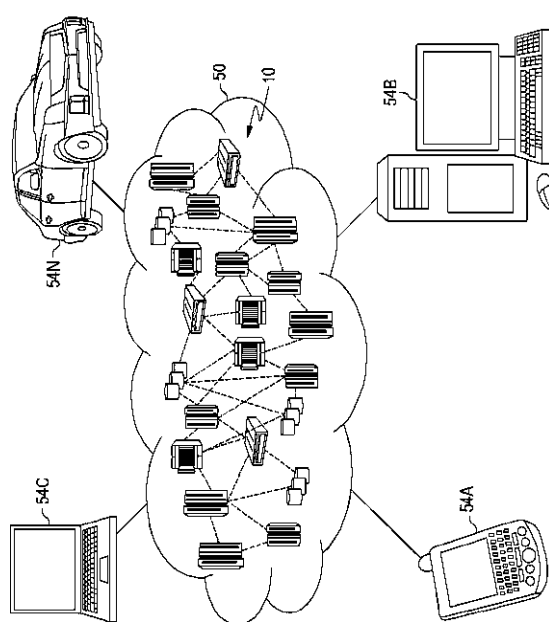
【図 8】



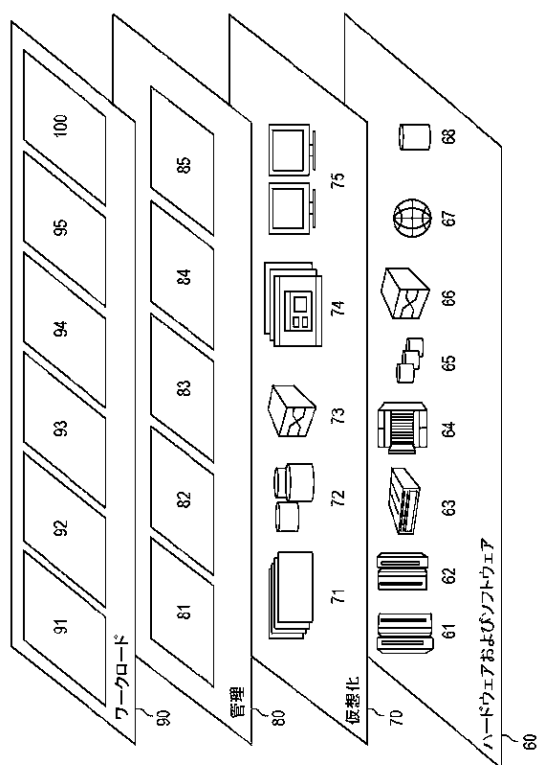
【図 9】



【図 10】



【図 11】



## フロントページの続き

- (72)発明者 バイド、マニシュ、アナン  
インド 5 0 0 0 8 1 テランガーナ州 ハイデラバード マダプール ハイテクシティ マインド  
スペース
- (72)発明者 ロヒア、プラネイ、クマール  
インド 5 6 0 0 4 5 カルナータカ州 バンガロール ラケナハリ・アンド・ナガワラ・ビレッジ  
アウター・リング・ロード マニヤタ・エンバシー・ビー
- (72)発明者 ナテサン ラママーシー、カーティケヤン  
アメリカ合衆国 1 0 5 9 8 ニューヨーク州 ヨークタウン・ハイツ ピーオー・ボックス 2 1 8  
キッチャワン・ロード 1 1 0 1
- (72)発明者 ブリ、ルチル  
アメリカ合衆国 1 0 5 9 8 ニューヨーク州 ヨークタウン・ハイツ ピーオー・ボックス 2 1 8  
キッチャワン・ロード 1 1 0 1
- (72)発明者 サハ、ディプティカルヤン  
インド 5 6 0 0 4 5 カルナータカ州 バンガロール ラケナハリ・アンド・ナガワラ・ビレッジ  
アウター・リング・ロード マニヤタ・エンバシー・ビー
- (72)発明者 ヴァーシュニー、クシュ、ラージ  
アメリカ合衆国 1 0 5 9 8 ニューヨーク州 ヨークタウン・ハイツ ピーオー・ボックス 2 1 8  
キッチャワン・ロード 1 1 0 1

審査官 武田 広太郎

- (56)参考文献 特開 2 0 1 4 - 2 1 1 7 6 1 ( J P , A )  
KAMIRAN, Faisal , Decision Theory for A 式外 Discrimination aware Classification , Proce  
edinds of 12th International Conf. on Data Mining , 2013年01月17日 , web.lums.edu.pk/ ak  
arim/pub/decision theory icdm2012.pdf

## (58)調査した分野(Int.Cl. , D B 名)

G 0 6 N 2 0 / 0 0  
G 0 6 Q 1 0 / 0 4