

(12) 特許協力条約に基づいて公開された国際出願

(19) 世界知的所有権機関
国際事務局

(43) 国際公開日
2024年8月15日(15.08.2024)



(10) 国際公開番号

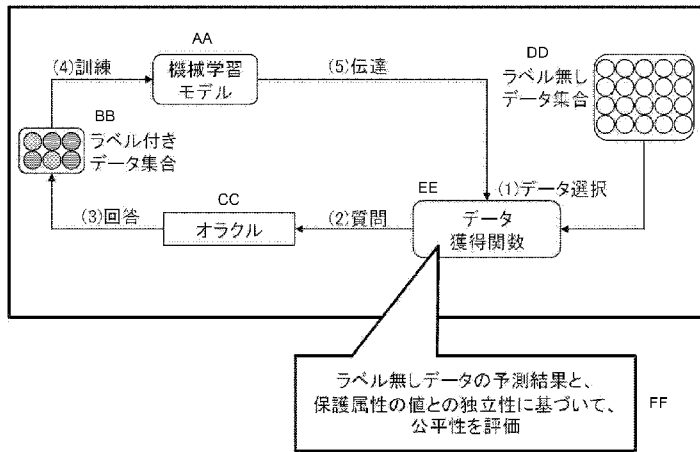
WO 2024/166331 A1

- (51) 国際特許分類:
G06N 3/091 (2023.01) G06N 20/00 (2019.01)
- (21) 国際出願番号: PCT/JP2023/004458
- (22) 国際出願日: 2023年2月9日(09.02.2023)
- (25) 国際出願の言語: 日本語
- (26) 国際公開の言語: 日本語
- (71) 出願人: 富士通株式会社 (FUJITSU LIMITED)
[JP/JP]; 〒2118588 神奈川県川崎市中原区上小田中4丁目1番1号 Kanagawa (JP).
- (72) 発明者: 園田 亮介 (SONODA, Ryosuke);
〒2118588 神奈川県川崎市中原区上小田中4丁目1番1号 富士通株式会社内 Kanagawa (JP).
- (74) 代理人: 弁理士法人太陽国際特許事務所(TAIYO, NAKAJIMA & KATO); 〒1600022 東京都新宿区新宿4丁目3番17号 Tokyo (JP).
- (81) 指定国(表示のない限り、全ての種類の国内保護が可能): AE, AG, AL, AM, AO, AT, AU, AZ, BA, BB, BG, BH, BN, BR, BW, BY, BZ, CA, CH, CL, CN, CO, CR, CU, CV, 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, MG, MK, MN, MW, MX, MY, MZ, 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.

(54) Title: MACHINE LEARNING PROGRAM, METHOD, AND DEVICE

(54) 発明の名称: 機械学習プログラム、方法、及び装置

[図3]



- (1) Data selection
(2) Question
(3) Answer
(4) Training
(5) Transfer
AA Machine learning model
BB Labeled dataset
CC Oracle
DD Unlabeled dataset
EE Data acquisition function
FF Evaluate fairness based on independence between prediction result for unlabeled data and value of protected attribute

(57) Abstract: Provided is a machine learning device that: calculates the independence between a prediction result, which is obtained when inputting each of a plurality of pieces of unlabeled data into a machine learning model, and the value of a first attribute of each of the plurality of pieces of data; selects a first piece of data from the plurality of pieces of data on the basis of the independence; obtains a label for the first piece of data; and executes training of the machine learning model on the basis of the first piece of data and the label.

(84) 指定国(表示のない限り、全ての種類の広域保護が可能): ARIPO (BW, CV, GH, GM, KE, LR, LS, MW, MZ, NA, RW, SC, SD, SL, ST, SZ, TZ, UG, ZM, ZW), ユーラシア (AM, AZ, BY, KG, KZ, RU, TJ, TM), ヨーロッパ (AL, AT, BE, BG, CH, CY, CZ, DE, DK, EE, ES, FI, FR, GB, GR, HR, HU, IE, IS, IT, LT, LU, LV, MC, ME, MK, MT, NL, NO, PL, PT, RO, RS, SE, SI, SK, SM, TR), OAPI (BF, BJ, CF, CG, CI, CM, GA, GN, GQ, GW, KM, ML, MR, NE, SN, TD, TG).

添付公開書類:

一 国際調査報告 (条約第21条(3))

(57) 要約: 機械学習装置は、ラベル付けされていない複数のデータのそれぞれを機械学習モデルに入力した場合の予測結果と、複数のデータのそれぞれの第1の属性の値との独立性を算出し、独立性に基づいて、複数のデータから第1のデータを選択し、第1のデータのラベルを取得し、第1のデータとラベルとに基づいて機械学習モデルの訓練を実行する。

明 細 書

発明の名称：機械学習プログラム、方法、及び装置

技術分野

[0001] 開示の技術は、機械学習プログラム、機械学習方法、及び機械学習装置に関する。

背景技術

[0002] 従来、公平性に配慮した機械学習モデルに関する技術が提案されている。例えば、分類器を学習するための訓練データと、訓練データに含まれる変数間の因果関係を表す因果グラフとを入力する学習装置が提案されている。この学習装置は、入力された訓練データと因果グラフとを用いて、所定の変数間の因果効果の平均が所定の範囲内にあり、かつ、因果効果の分散が所定の値以下である制約付き最適化問題を解くことにより分類器を学習する。

[0003] また、例えば、少数属性のデータを人工的に増やすことによって、各入力データに対して公平な判断を行うための学習用データを生成する情報処理装置が提案されている。この情報処理装置は、機械学習モデルの学習に使用する第1の学習用データを保持し、学習用データの偏りに関する情報を取得し、偏りに関する情報に基づいて、学習用データに含まれるデータを用いて第2の学習用データを生成する。そして、この情報処理装置は、第1の学習用データ及び第2の学習用データを用いて機械学習モデルを学習する。

[0004] また、例えば、ラベル付けされていないデータに、ラベルバイアスの量に従ってラベル付けするシステムが提案されている。このシステムは、選択バイアスの量及び特徴の希少性間の不一致に従って入力データをサンプリングし、サンプリングされてラベル付けされたデータと、追加のラベル付けされていないデータとを使用して分類器を訓練する。

[0005] また、公平性に配慮した能動学習により機械学習モデルの訓練を実行する手法として、能動学習と半教師あり学習とを統合した手法が提案されている。この手法は、ラベル付けされていない最も価値のあるデータを選択し、ラ

ベル付けのためにエキスパートシステムに送信する。そして、この手法は、ラベル付けされていないデータとラベル付けされているデータとの間の接続を確立し、ラベル付けされていないデータの固有の情報を使用してモデルを改善し、疑似ラベルをそれらのサンプルに割り当てる。

先行技術文献

特許文献

- [0006] 特許文献1：国際公開第2021/084609号公報
特許文献2：国際公開第2022/123907号公報
特許文献3：米国特許出願公開第2020/0372406号明細書

非特許文献

- [0007] 非特許文献1：Quan Ren, Hongbing Zhang, Dailu Zhang, Xiang Zhao, Lizhi Yan, Jianwen Rui, Fanxin Zeng, Xinyi Zhu, "A framework of active learning and semi-supervised learning for lithology identification based on improved naive Bayes," Expert Systems with Applications, Volume 202, 15 September 2022, 117278.

発明の概要

発明が解決しようとする課題

- [0008] しかしながら、従来の公平性に配慮した能動学習により機械学習モデルの訓練を実行する学習方法（以下、「公平能動学習」という）では、ラベル付けの対象となるデータを選択する処理においても機械学習モデルの訓練が必要であり、処理負荷が高いという問題がある。
- [0009] 一つの側面として、開示の技術は、公平能動学習の処理負荷を低減することを目的とする。

課題を解決するための手段

- [0010] 一つの態様として、開示の技術は、ラベル付けされていない複数のデータのそれぞれを機械学習モデルに入力した場合の予測結果と、前記複数のデータのそれぞれの第1の属性の値との独立性を算出する。そして、開示の技術

は、前記独立性に基づいて、前記複数のデータから第1のデータを選択し、前記第1のデータのラベルを取得し、前記第1のデータと前記ラベルとに基づいて前記機械学習モデルの訓練を実行する。

発明の効果

[0011] 一つの側面として、開示の技術は、公平能動学習の処理負荷を低減することができる、という効果を有する。

図面の簡単な説明

[0012] [図1]能動学習を説明するための図である。

[図2]従来の公平能動学習を説明するための図である。

[図3]本実施形態における公平能動学習を説明するための図である。

[図4]本実施形態に係る機械学習装置の機能ブロック図である。

[図5]精度改善度合いの一例を説明するための図である。

[図6]機械学習装置として機能するコンピュータの概略構成を示すブロック図である。

[図7]機械学習処理の一例を示すフローチャートである。

[図8]機械学習モデルの予測精度、公平性、及び公平能動学習の実行時間について、本手法と比較手法との比較を示す図である。

発明を実施するための形態

[0013] 以下、図面を参照して、開示の技術に係る実施形態の一例を説明する。

[0014] 実施形態の詳細を説明する前に、機械学習モデルの公平能動学習及びその課題について説明する。

[0015] まず、機械学習モデルの訓練とは、データのラベル（結果変数）と特徴量（説明変数）との関係を学習し、その関係を近似するパラメータを特定することである。また、機械学習における公平性とは、機械学習モデルによる予測結果に、意思決定における個人や集団の先天的又は後天的な特性（以下、「保護属性」という）に基づく偏りや差別がないことである。機械学習モデルの社会実装では、ローン審査における性差、顔認識における人種差、病気診断における年齢差等、集団に基づく公平性及び格差の指標の改善が求めら

れることが多い。そのため、公平な予測結果が得られるように機械学習モデルを訓練することが重要である。また、機械学習モデルの公平性と予測精度とはトレードオフの関係にあるため、これらのバランスをとることも重要である。

[0016] 教師あり学習により機械学習モデルを訓練する場合、正解を示すラベルが付与されたデータ（以下、「ラベル付きデータ」という）が必要である。ラベル付きデータは、1つ以上のラベルでタグ付けされたデータである。ラベルは通常、人間、その他情報源であるオラクルによってタグ付けされる。十分な数のラベル付きデータがない場合、機械学習モデルの公平性及び予測精度を十分に改善することができない。しかし、ラベル付きデータは、ラベル付けがされていないデータ（以下、「ラベル無しデータ」という）よりも収集コストが高い。

[0017] また、能動学習とは、質問によって機械学習モデルを改善する対話的な機械学習手法である。具体的には、能動学習を実行する情報処理装置は、図1に示すように、（1）データ選択、（2）質問、（3）回答、（4）訓練、及び（5）伝達の処理を実行する。

[0018] より具体的には、情報処理装置は、「（1）データ選択」の処理として、ラベル無しデータ集合に含まれる各データについて、データ獲得関数を算出し、データ獲得関数に基づいて、機械学習モデルの訓練に有用なデータを優先して選択する。データ獲得関数は、現在の機械学習モデルのパラメータ等の情報を用いて、各データに対する機械学習モデルによる予測の曖昧さ、ラベル無しデータ集合に対する各データの代表性等を表す指標である。また、情報処理装置は、「（2）質問」の処理として、選択したデータのラベルをオラクルへ問い合わせる。また、情報処理装置は、「（3）回答」の処理として、オラクルからの回答であるラベルを取得し、取得したラベルを選択したデータに付与してラベル付きデータとし、ラベル付きデータ集合に追加する。また、情報処理装置は、「（4）訓練」の処理として、ラベル付きデータ集合を用いて機械学習モデルを訓練する。また、情報処理装置は、「（5

) 伝達」の処理として、訓練後の機械学習モデルのパラメータ等の情報を（
1）データ選択の処理へ伝達する。

[0019] 能動学習では、上記（1）～（5）の処理が繰り返し実行されることにより、ラベル無しデータのうち、機械学習モデルの訓練に有用なデータから優先的にラベル付けが行われるため、効果的にラベル付きデータを収集することができる。

[0020] なお、図1において、丸は各データを表し、白丸はラベル無しデータ、ハッチング付きの丸はラベル付きデータを表し、ハッチングの相違はラベルの相違を表す。以下の図2及び図3においても同様である。

[0021] 通常の能動学習では、（1）データ選択において、機械学習モデルの予測精度の改善に役立つデータが選択され、公平性が考慮されない。そのため、ラベル付きデータが追加され、機械学習モデルの訓練が進むほど、公平性の悪化を招く場合がある。例えば、顔表情認識のための機械学習モデルにおいて、ラベル無しデータ集合から特定の人種のデータばかりが選択されてしまう場合が考えられる。

[0022] そこで、従来の公平能動学習では、（1）データ選択の処理で、公平性と予測精度とのトレードオフを考慮してデータを選択する。具体的には、図2に示すように、従来の公平能動学習では、ラベル無しデータ集合の一部をラベル無し検証データ集合とし、残りをラベル無し候補データ集合とする。従来の公平能動学習を実行する情報処理装置は、図2に示す（A）仮質問、（B）仮回答、（C）訓練、（D）評価、及び（E）選択の処理を実行することで、各候補データについて、検証データにおける機械学習モデルの不公平性改善度合いを推定する。

[0023] より具体的には、情報処理装置は、「（A）仮質問」の処理として、データに対する仮ラベルを出力するラベル付けモデルに各候補データを入力する。また、情報処理装置は、「（B）仮回答」の処理として、ラベル付けモデルから出力された仮ラベルを取得し、取得した仮ラベルを各候補データに付与して、仮ラベル付き候補データ集合とする。また、情報処理装置は、「（

C) 訓練」の処理として、仮ラベル付き候補データ集合を用いて機械学習モデルを訓練する。また、情報処理装置は、「(D) 評価」の処理として、ラベル無し検証データ集合を用いて、訓練の前後での機械学習モデルの不公平性の差分を考慮して、各仮ラベル付き候補データの公平性を評価する。また、情報処理装置は、ラベル無し検証データ集合を用いて機械学習モデルの予測精度を評価する。また、情報処理装置は、「(E) 選択」の処理として、公平性と予測精度とのトレードオフを考慮した指標に基づいて、最も値の良い候補データを選択する。

[0024] 従来の公平能動学習では、「(1) データ選択」において、機械学習モデルの訓練が実行されるため、機械学習モデルの決定境界を考慮して、予測精度が改善するようなデータが選択される。また、検証データを用いた評価により、公平性も改善するようなデータが選択される。しかし、機械学習モデルの訓練は処理負荷が高く、効率的にデータ選択を実行することができないという課題がある。特に、機械学習モデルが深層学習モデル、非線形モデル等の複雑なモデルの場合には、現実的な実行時間で公平能動学習を適用することが困難である。

[0025] そこで、本実施形態では、図3に示すように、ラベル無しデータに対するモデルの予測結果と、保護属性の値との独立性に基づいて、機械学習モデルの訓練を要することなく、各候補データの公平性を評価する。

[0026] また、上記の課題を解決するために容易に思いつく手段として、各候補データに対する機械学習モデルの予測結果に基づいて、各候補データの不公平性度合いを推定し、ラベル無し候補データ集合内から不公平性度合いが低いデータを選択することが考えられる。しかし、この場合、選択されるデータは、検証データへの影響が考慮されていないため、外れ値又は同じようなデータが選択され易い等、データの代表性が低くなる。そこで、さらに、本実施形態では、ラベル無しデータに対するモデルの予測結果として、候補データが与えられた場合の検証データの予測結果を用いる。以下、本実施形態に係る機械学習装置について説明する。

- [0027] 図4に示すように、機械学習装置10には、ラベル付きデータ集合20及びラベル無しデータ集合22が入力される。ラベル付きデータ集合20に含まれるラベル付きデータの数は、ラベル無しデータ集合22に含まれるラベル無しデータの数に比べて非常に少ないものとする。機械学習装置10は、ラベル無しデータ集合22から、公平性に配慮したデータを選択して機械学習モデル24の訓練を実行する。すなわち、機械学習装置10は、公平能動学習を実行する。
- [0028] 機械学習装置10は、機能的には、図4に示すように、制御部11を含む。制御部11は、さらに、訓練部12と、算出部14と、選択部16と、取得部18とを含む。また、機械学習装置10の所定の記憶領域には、機械学習モデル24が記憶される。
- [0029] 訓練部12は、ラベル付きデータ集合20に含まれる複数のラベル付きデータを訓練データとして用いて、機械学習モデル24の訓練を実行する。また、後述するように、本実施形態では、取得部18により、ラベル付きデータ集合20に新たなラベル付きデータが追加される。ラベル付きデータ集合20に新たなラベル付きデータが追加された場合には、訓練部12は、当初のラベル付きデータと、追加されたラベル付きデータとを用いて、機械学習モデル24の訓練を実行する。
- [0030] 算出部14は、ラベル無しデータ集合22に含まれる複数のラベル無しデータのそれぞれを機械学習モデル24に入力した場合の予測結果（以下、「予測ラベル」という）と、複数のデータのそれぞれの保護属性の値との独立性を算出する。保護属性は、開示の技術の「第1の属性」の一例である。算出部14は、独立性として、予測ラベルと保護属性の値との相互情報量を算出する。相互情報量は、2つの変数が依存しているか否かを定量的に表す指標であり、2つの変数がお互いに完全に独立している場合は、相互情報量は0となる。すなわち、予測ラベルと保護属性の値との相互情報量が0の場合、機械学習モデル24は完全に公平であると言える。したがって、この相互情報量が最も小さいラベル無しデータが最も公平なデータであると言える。

[0031] 本実施形態では、算出部14は、各候補データで条件付けられた、検証データについての予測ラベルと保護属性の値との相互情報量に基づいて、機械学習モデル24の不公平性改善度合いを算出する。具体的には、算出部14は、ラベル無しデータ集合22に含まれる複数のラベル無しデータの一部を検証データ、検証データ以外を候補データとして設定する。算出部14は、各検証データの予測ラベルと保護属性の値との独立性であって、各候補データの予測ラベルと保護属性の値との独立性に条件付けられた独立性を算出する。

[0032] より具体的には、算出部14は、検証データ v の予測ラベル Y_v と、検証データ v における保護属性の値 S_v との相互情報量 I の、候補データ u が与えられる前後での差分を、候補データ u の不公平性改善度合い F_u として、下記(1)式により算出する。

[0033] [数1]

$$F_u = \sum_v \{I(Y_v; S_v) - I(Y_v; S_v | Y_u, S_u)\} \quad (1)$$

[0034] (1)式において、 Y_u は、候補データ u の予測ラベル、 S_u は、候補データ u における保護属性の値である。(1)式右辺の Σ 内の第1項は、候補データ u が与えられる前の検証データ v の相互情報量であり、第2項は、候補データ u が与えられた後の検証データ v の相互情報量である。(1)式の場合、不公平性改善度合い F_u の値が大きい候補データ u ほど、不公平性の改善度合いが高いことを表す。

[0035] なお、算出部14は、(1)式右辺の Σ 内の第2項について、下記のように算出する。まず、算出部14は、第2項を下記(2)式のように変換する。(2)式において、 $H(X)$ は、 X のエントロピーである。

[0036] [数2]

$$\begin{aligned} & I(Y_v; S_v | Y_u, S_u) \\ &= H(Y_v | Y_u, S_u) - H(Y_v | S_v, Y_u, S_u) \\ &= \{H(Y_v, Y_u, S_u) - H(Y_u, S_u)\} \\ & \quad - \{H(Y_v, S_v, Y_u, S_u) - H(S_v, Y_u, S_u)\} \end{aligned} \quad (2)$$

[0037] 次に、算出部 14 は、モンテカルロドロップアウトにより、各エントロピーに対応する確率分布を近似する。算出部 14 は、機械学習モデル 24 のパラメータと、機械学習モデル 24 の予測ラベルの分布とが条件付き独立であると仮定し、下記 (3) 式により、 Y_i の確率 $P(Y_i)$ を算出する。ここでは、 $Y_i = \{Y_v, S_v, Y_u, S_u\}$ である。

[0038] [数3]

$$P(Y_i) = P(Y_i|\theta) \approx \frac{1}{M} \sum_{m=1}^M P(Y_i|\theta) \quad (3)$$

[0039] (3) 式において、 θ は、機械学習モデル 24 のパラメータであり、 M は、モンテカルロサンプリング回数である。算出部 14 は、(3) 式の確率分布を用いて、(2) 式の相互情報量を算出する。

[0040] また、算出部 14 は、各候補データを機械学習モデル 24 に入力した場合の予測結果の不確実性に基づく、各候補データによる機械学習モデル 24 の予測精度の改善度合いを、各候補データについて算出する。例えば、機械学習モデル 24 の決定境界付近に位置するデータは、機械学習モデル 24 にとって判断が難しいデータと言えるため、このようなデータを訓練データとして用いることで、機械学習モデル 24 の予測精度の改善が図られる。例えば、機械学習モデル 24 において、図 5 に示すように、特徴空間で決定境界が規定されているとする。なお、図 5 の例では、線形モデルの 2 値分類の例を示しており、各丸は各データの特徴量を表す。この場合、決定境界付近のデータ（例えば、図 5 中の網点の丸で示すデータ）は、どちらのラベルに属するか不安定であり、機械学習モデル 24 の精度改善に有用である。そこで、算出部 14 は、候補データが機械学習モデルの決定境界に近いほど高くなる精度改善度合いを算出する。例えば、算出部 14 は、下記 (4) 式に示すように、候補データ u の予測ラベル Y_u の不確かさを示すエントロピーを、精度改善度合い A_u として算出する。(4) 式において、 Y は、機械学習モデル 24 の予測ラベルの集合である。

[0041] [数4]

$$A_u = H(Y_u|u) = - \sum_{Y_u \in Y} P(Y_u|u) \log P(Y_u|u) \quad (4)$$

[0042] また、算出部14は、例えば、下記(5)式に示すように、不公平性改善度合い F_u と、精度改善度合い A_u と、不公平性改善度合い F_u と精度改善度合い A_u とのトレードオフを表す係数 α とで表される、各候補データ u についての評価値 E_u を算出する。

$$E_u = \alpha \times F_u + (1 - \alpha) \times A_u \quad (5)$$

α は、0~1の値(例えば、0.6)であり、不公平性改善度合い F_u と精度改善度合い A_u とのどちらをどの程度優先させるかを規定する係数である。

[0043] 選択部16は、算出部14により算出された各候補データ u についての評価値 E_u に基づいて、複数の候補データ u から、ラベル付けの対象となる対象データを選択する。対象データは、開示の技術の「第1のデータ」の一例である。例えば、選択部16は、評価値 E_u が最も高い候補データ u を選択してもよいし、評価値 E_u が所定値以上の候補データ u を選択してもよいし、評価値 E_u が上位所定個の候補データ u を選択してもよい。

[0044] 取得部18は、選択部16により選択された対象データのラベルを、人間、その他情報源であるオラクルへ問い合わせ、オラクルからの回答であるラベルを取得する。取得部18は、取得したラベルを対象データに付与してラベル付きデータとし、ラベル付きデータ集合20に追加する。これにより、上述したように、訓練部12が、追加されたラベル付きデータも用いて、機械学習モデル24の訓練を実行する。

[0045] 機械学習装置10は、例えば図6に示すコンピュータ40で実現されてよい。コンピュータ40は、CPU (Central Processing Unit) 41と、GPU (Graphics Processing Unit) 42と、一時記憶領域としてのメモリ43と、不揮発性の記憶装置44とを備える。また、コンピュータ40は、入力装置、表示装置等の入出力装置45と、記憶媒体49に対するデータの読み

込み及び書き込みを制御するR/W (Read/Write) 装置46とを備える。また、コンピュータ40は、インターネット等のネットワークに接続される通信I/F (Interface) 47を備える。CPU41、GPU42、メモリ43、記憶装置44、入出力装置45、R/W装置46、及び通信I/F47は、バス48を介して互いに接続される。

[0046] 記憶装置44は、例えば、HDD (Hard Disk Drive)、SSD (Solid State Drive)、フラッシュメモリ等である。記憶媒体としての記憶装置44には、コンピュータ40を、機械学習装置10として機能させるための機械学習プログラム50が記憶される。機械学習プログラム50は、訓練プロセス制御命令52と、算出プロセス制御命令54と、選択プロセス制御命令56と、取得プロセス制御命令58とを有する。また、記憶装置44は、機械学習モデル24を構成する情報が記憶される情報記憶領域60を有する。

[0047] CPU41は、機械学習プログラム50を記憶装置44から読み出してメモリ43に展開し、機械学習プログラム50が有する制御命令を順次実行する。CPU41は、訓練プロセス制御命令52を実行することで、図4に示す訓練部12として動作する。また、CPU41は、算出プロセス制御命令54を実行することで、図4に示す算出部14として動作する。また、CPU41は、選択プロセス制御命令56を実行することで、図4に示す選択部16として動作する。また、CPU41は、取得プロセス制御命令58を実行することで、図4に示す取得部18として動作する。また、CPU41は、情報記憶領域60から情報を読み出して、機械学習モデル24をメモリ43に展開する。これにより、機械学習プログラム50を実行したコンピュータ40が、機械学習装置10として機能することになる。なお、プログラムを実行するCPU41はハードウェアである。また、プログラムの一部は、GPU62により実行されてもよい。

[0048] なお、機械学習プログラム50により実現される機能は、例えば半導体集積回路、より詳しくはASIC (Application Specific Integrated Circuit)、FPGA (Field-Programmable Gate Array) 等で実現されてもよい。

- [0049] 次に、本実施形態に係る機械学習装置 10 の動作について説明する。機械学習モデル 24 に対する公平能動学習が指示されると、機械学習装置 10 が、図 7 に示す機械学習処理を実行する。なお、機械学習処理は、開示の技術の機械学習方法の一例である。
- [0050] ステップ S 10 で、訓練部 12 が、ラベル付きデータ集合 20 を取得して、ラベル付きデータを訓練データとして用いて、機械学習モデル 24 の訓練を実行する。次に、ステップ S 12 で、訓練部 12 が、公平能動学習の終了条件を満たすか否かを判定する。終了条件は、例えば、ラベル付きデータ集合 20 に新たに追加されたデータ数が所定数を越えた場合としてよい。終了条件を満たさない場合には、ステップ S 14 へ移行する。
- [0051] ステップ S 14 では、算出部 14 が、ラベル無しデータ集合 22 に含まれる複数のラベル無しデータの一部を検証データ、残りを候補データとして設定する。次に、ステップ S 16 で、算出部 14 が、例えば (1) 式に示す、検証データ v の予測ラベル Y_v と、検証データ v における保護属性の値 S_v との相互情報量 I の、候補データ u が与えられる前後での差分を、各候補データ u の不公平性改善度合い F_u として算出する。
- [0052] 次に、ステップ S 18 で、算出部 14 が、例えば (4) 式に示す、各候補データ u の予測ラベル Y_u の不確かさを示すエントロピーを、精度改善度合い A_u として算出する。次に、ステップ S 20 で、算出部 14 が、例えば (5) 式に示す、不公平性改善度合い F_u と、精度改善度合い A_u と、不公平性改善度合い F_u と精度改善度合い A_u とのトレードオフを表す係数 α とで表される、各候補データ u についての評価値 E_u を算出する。
- [0053] 次に、ステップ S 22 で、選択部 16 が、各候補データ u についての評価値 E_u に基づいて、複数の候補データ u から、ラベル付けの対象となる対象データを選択する。次に、ステップ S 24 で、取得部 18 が、対象データのラベルをオラクルへ問い合わせ、オラクルからの回答であるラベルを取得する。次に、ステップ S 26 で、取得部 18 が、取得したラベルを対象データに付与してラベル付きデータとし、ラベル付きデータ集合 20 に追加すると共

に、対象データとなった候補データをラベル無しデータ集合 22 から削除し、ステップ S10 に戻る。

[0054] ステップ S10 に戻ると、訓練部 12 が、追加されたラベル付きデータも用いて、機械学習モデル 24 の訓練を実行する。次に、ステップ S12 で、公平能動学習の終了条件を満たすと判定されると、ステップ S28 へ移行する。ステップ S28 では、訓練部 12 が、公平能動学習による訓練済みの機械学習モデルを出力し、機械学習処理は終了する。

[0055] 以上説明したように、本実施形態に係る機械学習装置は、複数のラベル無しデータのそれぞれを機械学習モデルに入力した場合の予測結果と、複数のラベル無しデータのそれぞれの保護属性の値との独立性を算出する。また、機械学習装置は、算出した独立性に基づいて、複数のラベル無しデータから対象データを選択し、オラクルへ問い合わせで対象データのラベルを取得し、対象データと取得したラベルとに基づいて機械学習モデルの訓練を実行する。すなわち、本実施形態に係る機械学習装置は、ラベル無しデータの公平性を情報理論的アプローチに基づいて評価し、機械学習モデルの再学習を行うことなく、ラベル付けの対象となるデータを選択する。これにより、本実施形態に係る機械学習装置は、公平能動学習の処理負荷を低減することができる。

[0056] また、本実施形態に係る機械学習装置は、各候補データについて、予測結果と保護属性の値との独立性として相互情報量を用いた不公平性改善度合いを算出すると共に、候補データの不確かさに基づいて、機械学習モデルの精度改善度合いを算出する。そして、機械学習装置は、不公平性改善度合いと精度改善度合いとのトレードオフを考慮した評価値に基づいて対象データを選択する。これにより、本実施形態に係る機械学習装置は、処理負荷を低減しつつ、不公平性改善度合いと精度改善度合いとのトレードオフを最適化するようなデータを対象データとして選択することができる。

[0057] さらに、本実施形態に係る機械学習装置は、ラベル無しデータの一部を検証データ、残りを候補データとし、検証データの予測結果と保護属性の値と

の相互情報量の、候補データを与える前後での差分を不公平性改善度合いとして算出する。このように、検証データへの影響を考慮して候補データの不公平性改善度合いが算出されることで、候補データの中から代表性の高いデータが選択され易くなる。

- [0058] 図8に、機械学習モデルの予測精度、公平性、及び公平能動学習の実行時間について、本実施形態の手法（以下、「本手法」という）と比較手法との比較を概略的に示す。ここでの比較手法は、非特許文献1に記載の手法のように、ラベル付けの対象となるデータを選択する際に、機械学習モデルの再学習が必要な手法である。本手法は、予測精度及び公平性のトレードオフは比較手法と同等である。また、本手法は、比較手法に比べ、実行時間が大幅に削減される。なお、公平能動学習の理論的な計算コストは、比較手法が「 $O((T + N_v) C N_u)$ 」、本手法が「 $O(N_u N_v C^2 M)$ 」である。なお、 T は機械学習モデルの訓練の計算コスト、 N_v は検証データの数、 N_u は候補データの数、 C はラベルの数、 M はモンテカルロサンプリング回数である。
- [0059] なお、上記実施形態では、ラベル無しデータの一部を検証データとして用いて公平性を評価する場合について説明したが、これに限定されない。ラベル無しデータ集合に含まれる全てのデータを候補データとしてもよい。この場合、公平性の指標としては、例えば、(1)式右辺の Σ 内の第2項を除いた式を用いるようにすればよい。また、上記実施形態では、公平性と予測精度とのトレードオフを考慮した評価値を用いてデータを選択する場合について説明したが、公平性のみを考慮してデータを選択するようにしてもよい。
- [0060] また、上記実施形態では、ラベル無しデータについての機械学習モデルの予測結果と保護属性の値との独立性を示す指標として相互情報量を用いる場合について説明した。相互情報量は、予測結果と保護属性の値との独立性を数式的に定義できればよく、例えば、カルバックライブラー情報量、イエンセン・シャノン情報量、共分散、Demographic parity difference、Disparate impact ratio等を適用してよい。
- [0061] また、上記実施形態では、機械学習プログラムが記憶装置に予め記憶（イ

ンストール) されているが、これに限定されない。開示の技術に係るプログラムは、CD-ROM、DVD-ROM、USBメモリ等の記憶媒体に記憶された形態で提供されてもよい。

符号の説明

- [0062] 10 機械学習装置
- 11 制御部
- 12 訓練部
- 14 算出部
- 16 選択部
- 18 取得部
- 20 ラベル付きデータ集合
- 22 ラベル無しデータ集合
- 24 機械学習モデル
- 40 コンピュータ
- 41 CPU
- 42 GPU
- 43 メモリ
- 44 記憶装置
- 45 入出力装置
- 46 R/W装置
- 47 通信I/F
- 48 バス
- 49 記憶媒体
- 50 機械学習プログラム
- 52 訓練プロセス制御命令
- 54 算出プロセス制御命令
- 56 選択プロセス制御命令
- 58 取得プロセス制御命令

60 情報記憶領域

請求の範囲

- [請求項1] ラベル付けされていない複数のデータのそれぞれを機械学習モデルに入力した場合の予測結果と、前記複数のデータのそれぞれの第1の属性の値との独立性を算出し、
- 前記独立性に基づいて、前記複数のデータから第1のデータを選択し、
- 前記第1のデータのラベルを取得し、
- 前記第1のデータと前記ラベルとに基づいて前記機械学習モデルの訓練を実行する、
- 処理をコンピュータに実行させることを特徴とする機械学習プログラム。
- [請求項2] 前記独立性は、前記予測結果と前記第1の属性の値との相互情報量である請求項1に記載の機械学習プログラム。
- [請求項3] 前記第1のデータを選択する処理は、前記複数のデータのそれぞれを前記機械学習モデルに入力した場合の予測結果の不確実性に基づく、前記複数のデータのそれぞれによる前記機械学習モデルの予測精度の改善度合いと、前記独立性とに基づいて選択することを含む、
- ことを特徴とする請求項1又は請求項2に記載の機械学習プログラム。
- [請求項4] 前記改善度合いは、前記第1のデータが前記機械学習モデルの決定境界に近いほど高くなる、
- ことを特徴とする請求項3に記載の機械学習プログラム。
- [請求項5] 前記第1のデータを選択する処理は、前記改善度合いと、前記独立性と、前記改善度合いと前記独立性とのトレードオフを表す係数とで表される指標が所定値以上、又は、前記指標が上位所定個のデータを選択することを含む、
- ことを特徴とする請求項3に記載の機械学習プログラム。
- [請求項6] 前記独立性を算出する処理は、前記複数のデータの一部を検証デー

タ、前記複数のデータのうち前記検証データ以外を候補データとし、前記検証データのそれぞれを前記機械学習モデルに入力した場合の予測結果と、前記検証データのそれぞれの前記第1の属性の値との独立性であって、前記候補データのそれぞれを前記機械学習モデルに入力した場合の予測結果と、前記候補データのそれぞれの前記第1の属性の値との独立性に条件付けられた独立性を算出することを含み、

前記第1のデータを選択する処理は、前記候補データから前記第1のデータを選択することを含む、

ことを特徴とする請求項1又は請求項2に記載の機械学習プログラム。

[請求項7] ラベル付けされていない複数のデータのそれぞれを機械学習モデルに入力した場合の予測結果と、前記複数のデータのそれぞれの第1の属性の値との独立性を算出し、

前記独立性に基づいて、前記複数のデータから第1のデータを選択し、

前記第1のデータのラベルを取得し、

前記第1のデータと前記ラベルとに基づいて前記機械学習モデルの訓練を実行する、

処理をコンピュータが実行することを特徴とする機械学習方法。

[請求項8] 前記独立性は、前記予測結果と前記第1の属性の値との相互情報量である請求項7に記載の機械学習方法。

[請求項9] 前記第1のデータを選択する処理は、前記複数のデータのそれぞれを前記機械学習モデルに入力した場合の予測結果の不確実性に基づく、前記複数のデータのそれぞれによる前記機械学習モデルの予測精度の改善度合いと、前記独立性とに基づいて選択することを含む、

ことを特徴とする請求項7又は請求項8に記載の機械学習方法。

[請求項10] 前記改善度合いは、前記第1のデータが前記機械学習モデルの決定境界に近いほど高くなる、

ことを特徴とする請求項 9 に記載の機械学習方法。

[請求項11]

前記第 1 のデータを選択する処理は、前記改善度合いと、前記独立性と、前記改善度合いと前記独立性とのトレードオフを表す係数とで表される指標が所定値以上、又は、前記指標が上位所定個のデータを選択することを含む、

ことを特徴とする請求項 9 に記載の機械学習方法。

[請求項12]

前記独立性を算出する処理は、前記複数のデータの一部を検証データ、前記複数のデータのうち前記検証データ以外を候補データとし、前記検証データのそれぞれを前記機械学習モデルに入力した場合の予測結果と、前記検証データのそれぞれの前記第 1 の属性の値との独立性であって、前記候補データのそれぞれを前記機械学習モデルに入力した場合の予測結果と、前記候補データのそれぞれの前記第 1 の属性の値との独立性に条件付けられた独立性を算出することを含み、

前記第 1 のデータを選択する処理は、前記候補データから前記第 1 のデータを選択することを含む、

ことを特徴とする請求項 7 又は請求項 8 に記載の機械学習方法。

[請求項13]

ラベル付けされていない複数のデータのそれぞれを機械学習モデルに入力した場合の予測結果と、前記複数のデータのそれぞれの第 1 の属性の値との独立性を算出し、

前記独立性に基づいて、前記複数のデータから第 1 のデータを選択し、

前記第 1 のデータのラベルを取得し、

前記第 1 のデータと前記ラベルとに基づいて前記機械学習モデルの訓練を実行する、

処理を実行する制御部を含むことを特徴とする機械学習装置。

[請求項14]

前記独立性は、前記予測結果と前記第 1 の属性の値との相互情報量である請求項 13 に記載の機械学習装置。

[請求項15]

前記第 1 のデータを選択する処理は、前記複数のデータのそれぞれ

を前記機械学習モデルに入力した場合の予測結果の不確実性に基づく、前記複数のデータのそれぞれによる前記機械学習モデルの予測精度の改善度合いと、前記独立性とに基づいて選択することを含む、

ことを特徴とする請求項 13 又は請求項 14 に記載の機械学習装置。

[請求項16] 前記改善度合いは、前記第 1 のデータが前記機械学習モデルの決定境界に近いほど高くなる、

ことを特徴とする請求項 15 に記載の機械学習装置。

[請求項17] 前記第 1 のデータを選択する処理は、前記改善度合いと、前記独立性と、前記改善度合いと前記独立性とのトレードオフを表す係数とで表される指標が所定値以上、又は、前記指標が上位所定個のデータを選択することを含む、

ことを特徴とする請求項 15 に記載の機械学習装置。

[請求項18] 前記独立性を算出する処理は、前記複数のデータの一部を検証データ、前記複数のデータのうち前記検証データ以外を候補データとし、前記検証データのそれぞれを前記機械学習モデルに入力した場合の予測結果と、前記検証データのそれぞれの前記第 1 の属性の値との独立性であって、前記候補データのそれぞれを前記機械学習モデルに入力した場合の予測結果と、前記候補データのそれぞれの前記第 1 の属性の値との独立性に条件付けられた独立性を算出することを含み、

前記第 1 のデータを選択する処理は、前記候補データから前記第 1 のデータを選択することを含む、

ことを特徴とする請求項 13 又は請求項 14 に記載の機械学習装置。

[請求項19] ラベル付けされていない複数のデータのそれぞれを機械学習モデルに入力した場合の予測結果と、前記複数のデータのそれぞれの第 1 の属性の値との独立性を算出し、

前記独立性に基づいて、前記複数のデータから第 1 のデータを選択

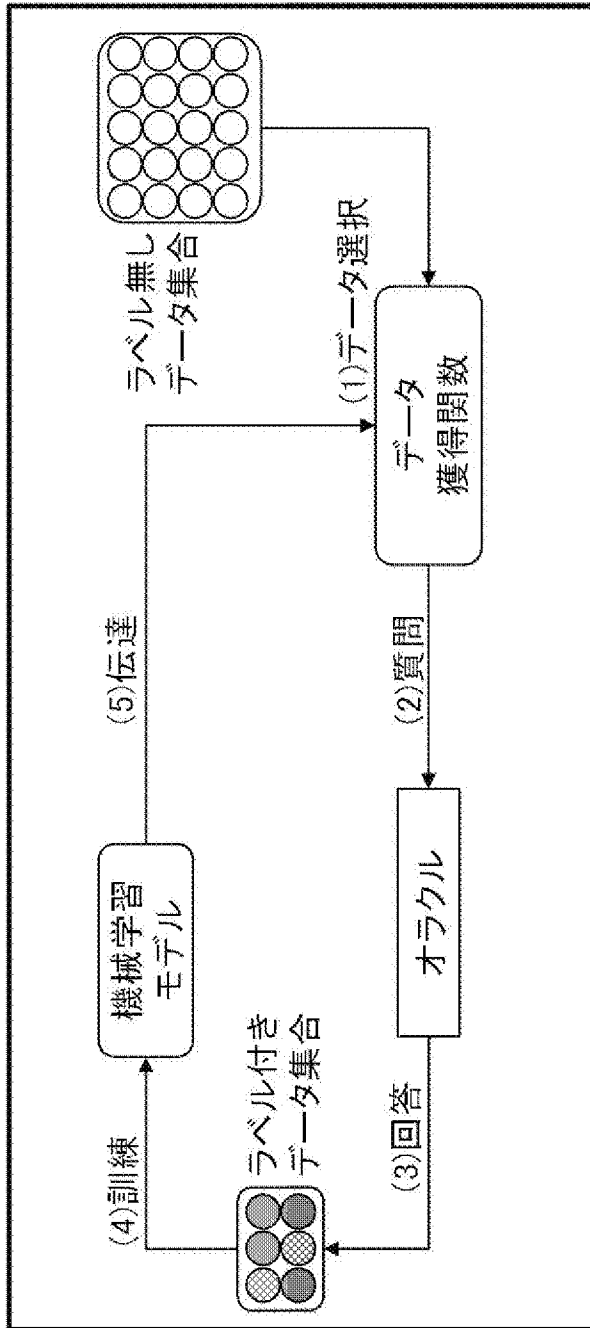
し、

前記第 1 のデータのラベルを取得し、

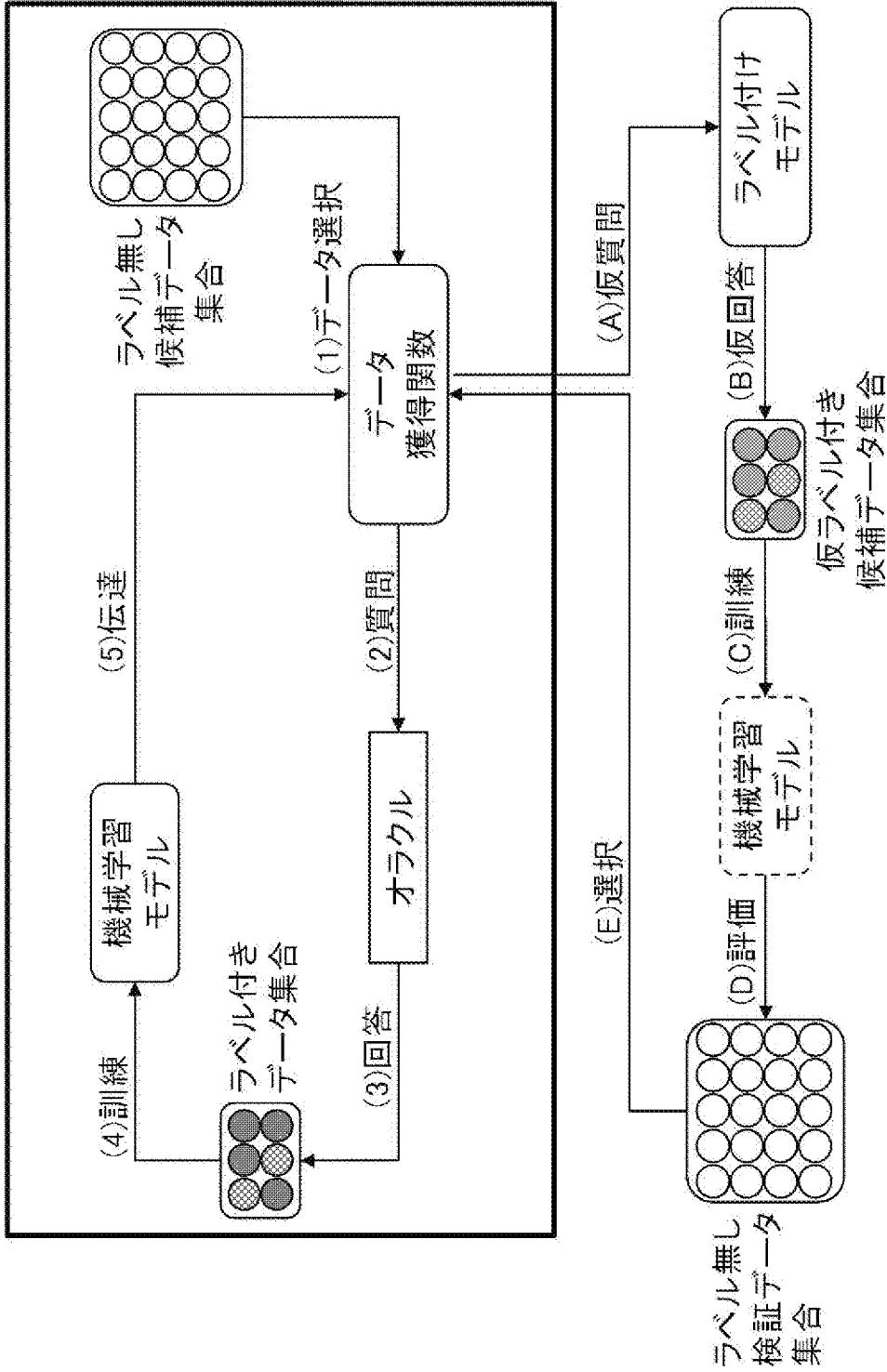
前記第 1 のデータと前記ラベルとに基づいて前記機械学習モデルの
訓練を実行する、

処理をコンピュータに実行させることを特徴とする機械学習プログラ
ムを記憶した非一時的記憶媒体。

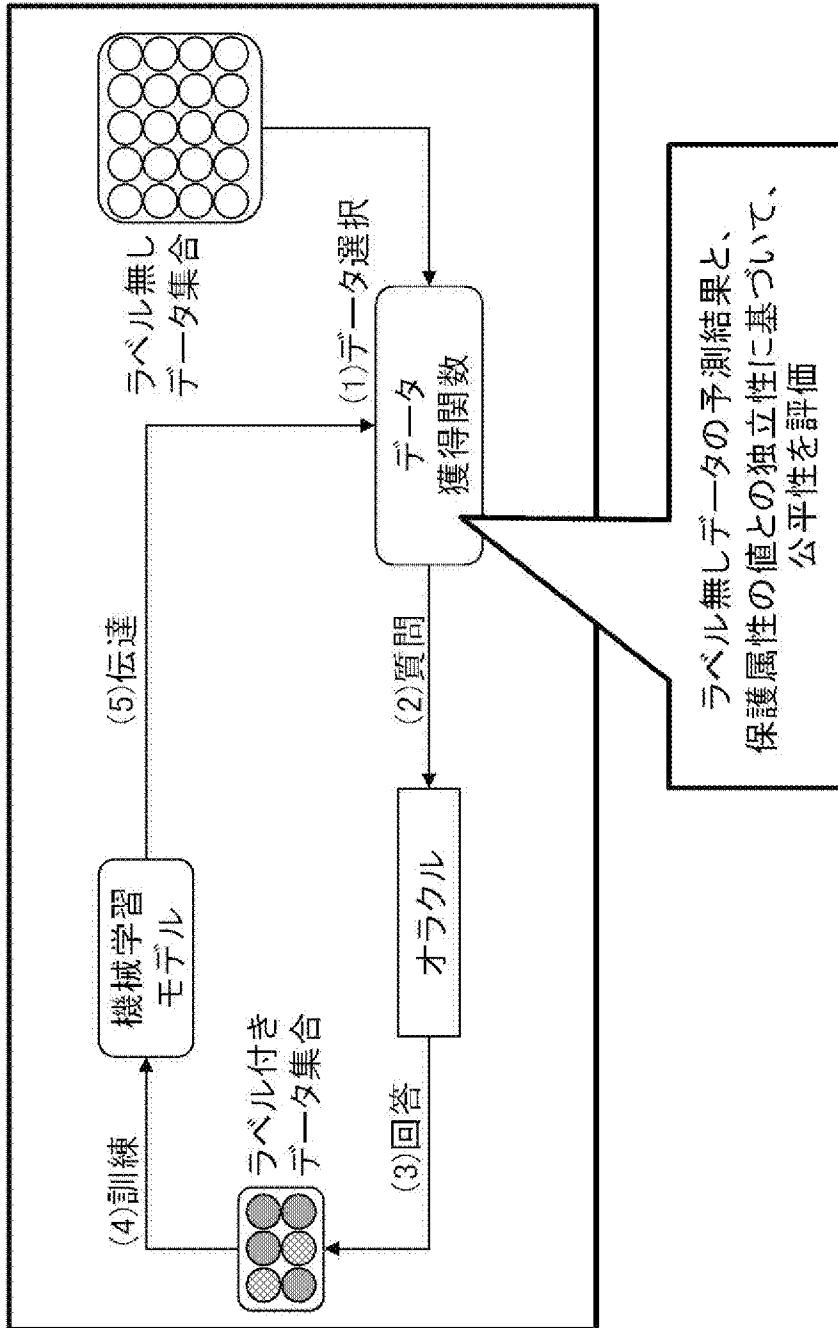
[図1]



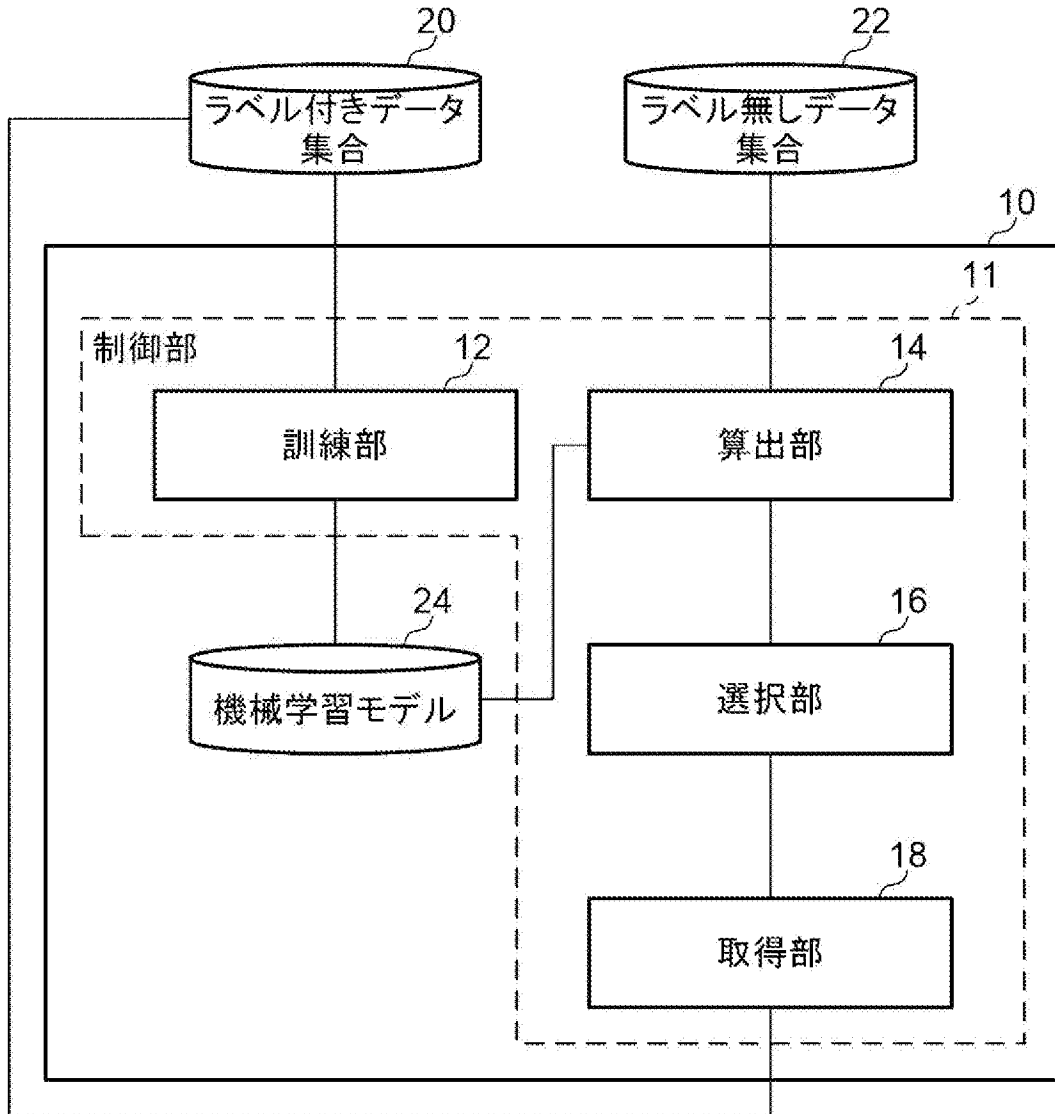
[図2]



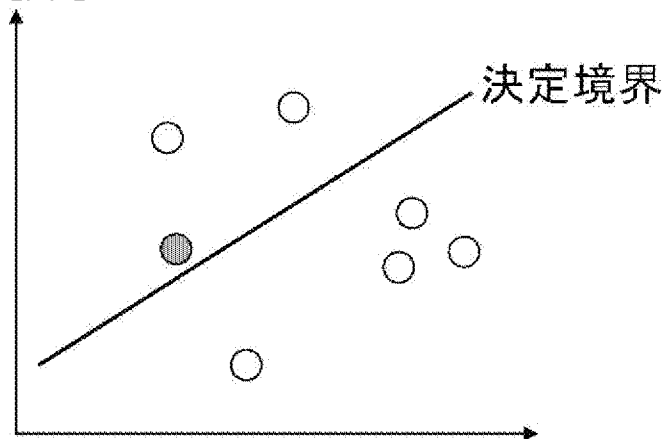
[図3]



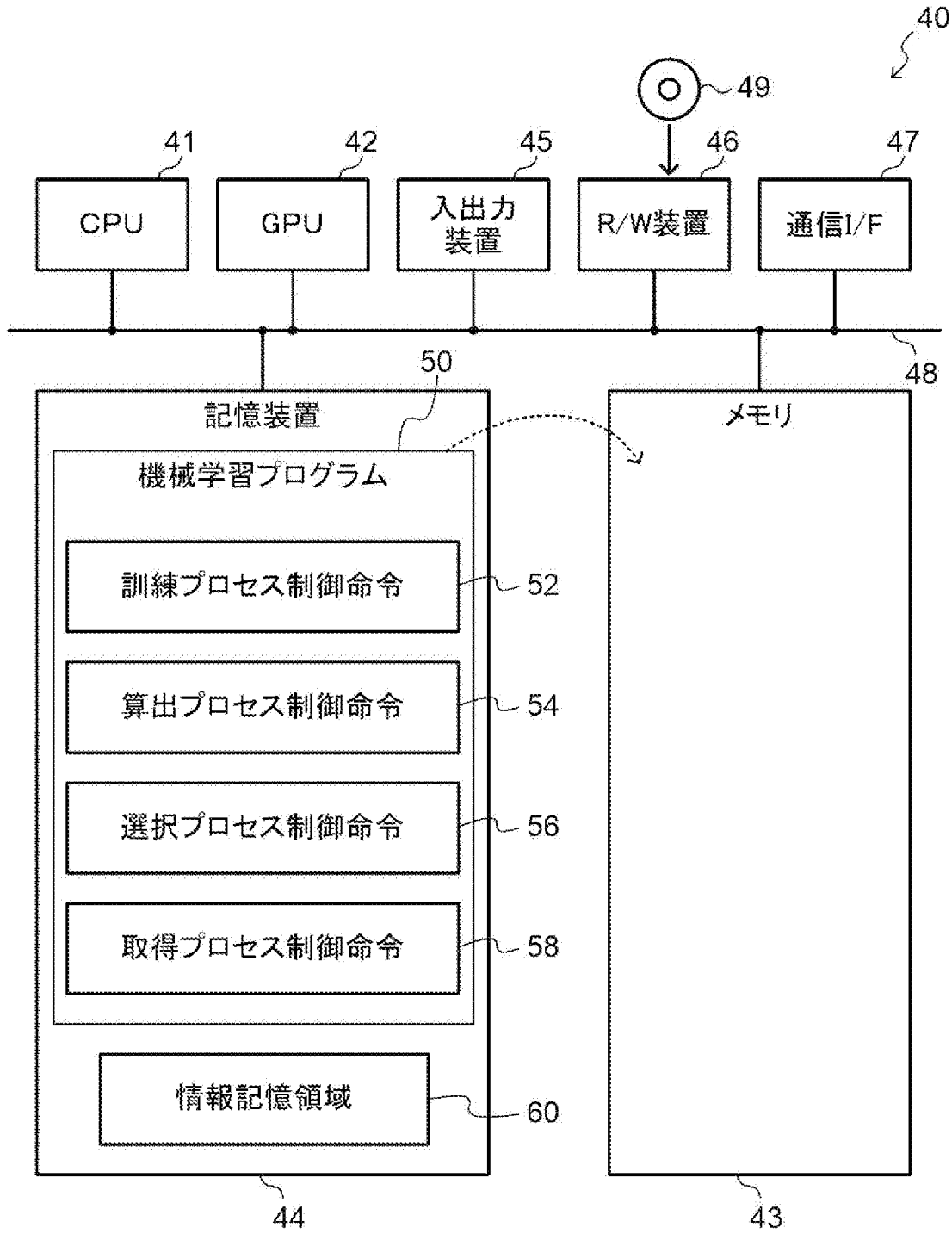
[図4]



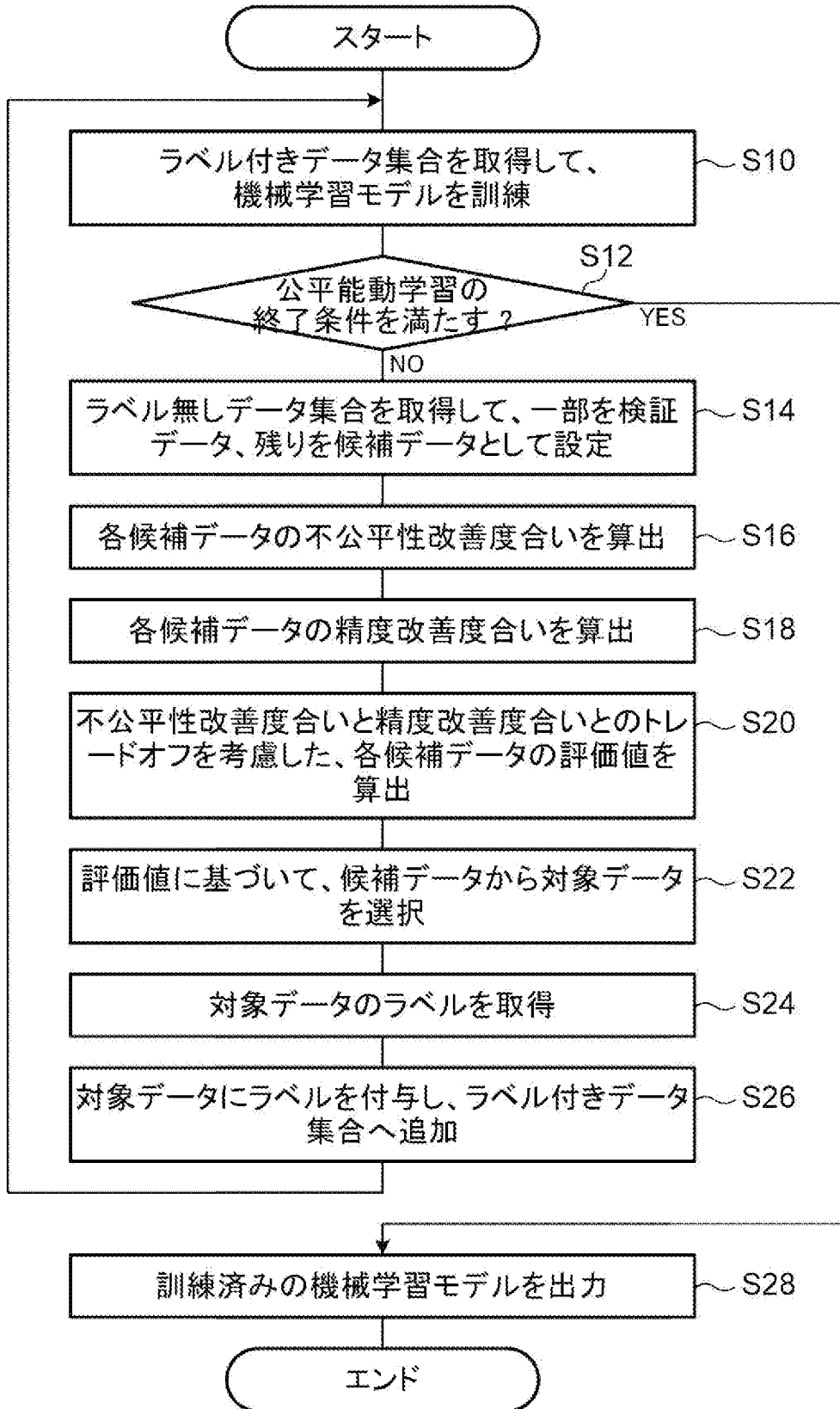
[図5]

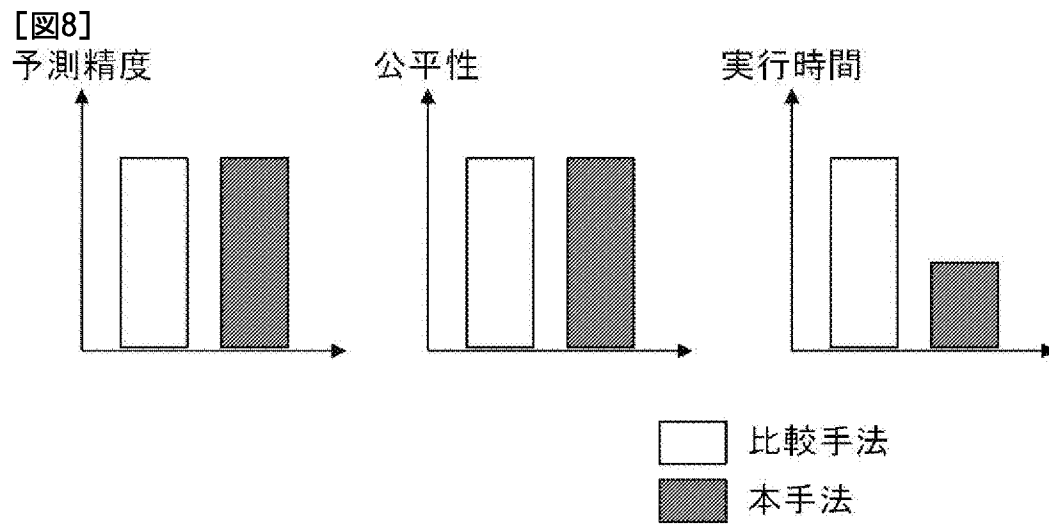


[図6]



[図7]





INTERNATIONAL SEARCH REPORT

International application No.

PCT/JP2023/004458

A. CLASSIFICATION OF SUBJECT MATTER		
<i>G06N 3/091</i> (2023.01)i; <i>G06N 20/00</i> (2019.01)i FI: G06N3/091; G06N20/00 130		
According to International Patent Classification (IPC) or to both national classification and IPC		
B. FIELDS SEARCHED		
Minimum documentation searched (classification system followed by classification symbols) G06N3/00-99/00		
Documentation searched other than minimum documentation to the extent that such documents are included in the fields searched Published examined utility model applications of Japan 1922-1996 Published unexamined utility model applications of Japan 1971-2023 Registered utility model specifications of Japan 1996-2023 Published registered utility model applications of Japan 1994-2023		
Electronic data base consulted during the international search (name of data base and, where practicable, search terms used)		
C. DOCUMENTS CONSIDERED TO BE RELEVANT		
Category*	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
A	ANAHIDEH, Hadis et al., Fair Active Learning, arXiv [online], 31 March 2021, [retrieved on 20 April 2023], Internet: <URL: https://arxiv.org/abs/2001.01796v5 > entire text, all drawings	1-19
A	US 2021/0035014 A1 (INTERNATIONAL BUSINESS MACHINES CORPORATION) 04 February 2021 (2021-02-04) entire text, all drawings	1-19
A	WO 2022/254626 A1 (FUJITSU LIMITED) 08 December 2022 (2022-12-08) entire text, all drawings	1-19
<input 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 20 April 2023		Date of mailing of the international search report 09 May 2023
Name and mailing address of the ISA/JP Japan Patent Office (ISA/JP) 3-4-3 Kasumigaseki, Chiyoda-ku, Tokyo 100-8915 Japan		Authorized officer Telephone No.

INTERNATIONAL SEARCH REPORT
Information on patent family members

International application No.

PCT/JP2023/004458

Patent document cited in search report	Publication date (day/month/year)	Patent family member(s)	Publication date (day/month/year)
US 2021/0035014 A1	04 February 2021	(Family: none)	
WO 2022/254626 A1	08 December 2022	(Family: none)	

A. 発明の属する分野の分類（国際特許分類（IPC）） G06N 3/091(2023.01)i; G06N 20/00(2019.01)i FI: G06N3/091; G06N20/00 130		
B. 調査を行った分野 調査を行った最小限資料（国際特許分類（IPC）） G06N3/00-99/00 最小限資料以外の資料で調査を行った分野に含まれるもの 日本国実用新案公報 1922-1996年 日本国公開実用新案公報 1971-2023年 日本国実用新案登録公報 1996-2023年 日本国登録実用新案公報 1994-2023年		
国際調査で使用した電子データベース（データベースの名称、調査に使用した用語）		
C. 関連すると認められる文献		
引用文献の カテゴリー*	引用文献名 及び一部の箇所が関連するときは、その関連する箇所の表示	関連する 請求項の番号
A	ANAHI DEH, Hadis et al., Fair Active Learning, arXiv [online], 2021.03.31, [検索日 2023.04.20], インターネット: <URL:https://arxiv.org/ abs/2001.01796v5> 全文、全図	1-19
A	US 2021/0035014 A1 (INTERNATIONAL BUSINESS MACHINES CORPORATION) 04.02.2021 (2021-02-04) 全文、全図	1-19
A	WO 2022/254626 A1 (富士通株式会社) 08.12.2022 (2022-12-08) 全文、全図	1-19
<input type="checkbox"/> C欄の続きにも文献が列挙されている。 <input checked="" type="checkbox"/> パテントファミリーに関する別紙を参照。		
* 引用文献のカテゴリー “A” 特に関連のある文献ではなく、一般的技術水準を示すもの “E” 国際出願日前の出願または特許であるが、国際出願日以後に 公表されたもの “L” 優先権主張に疑義を提起する文献又は他の文献の発行日若しく は他の特別な理由を確立するために引用する文献（理由を 付す） “O” 口頭による開示、使用、展示等に言及する文献 “P” 国際出願日前で、かつ優先権の主張の基礎となる出願の日の 後に公表された文献 “T” 国際出願日又は優先日後に公表された文献であって出願と抵 触するものではなく、発明の原理又は理論の理解のために引 用するもの “X” 特に関連のある文献であって、当該文献のみで発明の新規性 又は進歩性がないと考えられるもの “Y” 特に関連のある文献であって、当該文献と他の1以上の文献 との、当業者にとって自明である組合せによって進歩性がな いと考えられるもの “&” 同一パテントファミリー文献		
国際調査を完了した日	国際調査報告の発送日	
20.04.2023	09.05.2023	
名称及びあて先 日本国特許庁(ISA/JP) 〒100-8915 日本国 東京都千代田区霞が関三丁目4番3号	権限のある職員（特許庁審査官） 山本 俊介 5B 5087 電話番号 03-3581-1101 内線 3545	

国際調査報告
パテントファミリーに関する情報

国際出願番号

PCT/JP2023/004458

引用文献	公表日	パテントファミリー文献	公表日
US 2021/0035014 A1	04.02.2021	(ファミリーなし)	
WO 2022/254626 A1	08.12.2022	(ファミリーなし)	