

(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
2021年11月18日(18.11.2021)

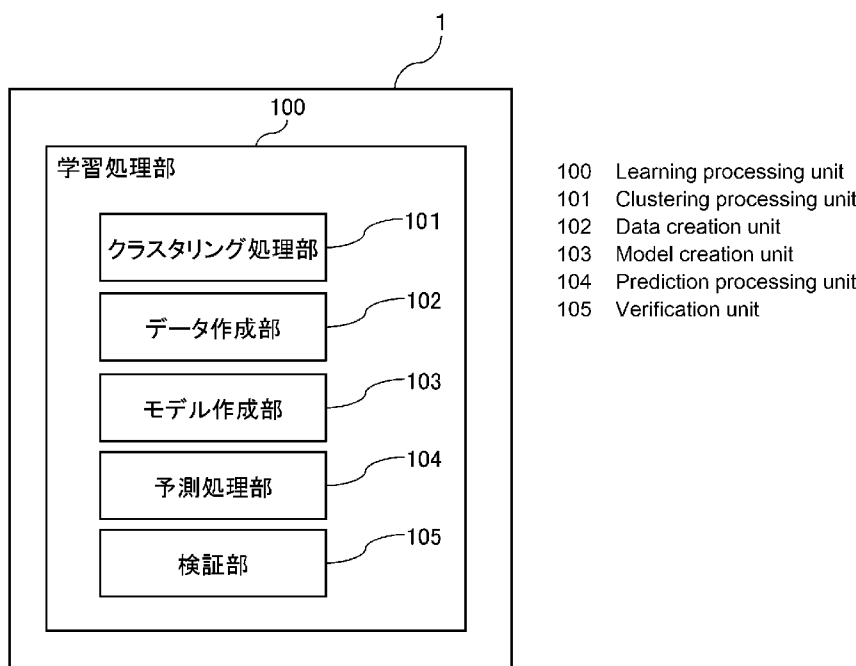


(10) 国際公開番号
WO 2021/229630 A1

- (51) 国際特許分類:
G06N 20/00 (2019.01)
- (21) 国際出願番号: PCT/JP2020/018777
- (22) 国際出願日: 2020年5月11日(11.05.2020)
- (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: 松尾 達 (**MATSUO, Tatsuru**); 〒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: 真田 有, 外 (**SANADA, Tamotsu et al.**); 〒1800004 東京都武蔵野市吉祥寺本町
- 1 丁目 10 番 31 号 NMF 吉祥寺本町ビル8階 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, CZ, DE, DJ, DK, DM, DO, DZ, EC, EE, EG, ES, FI, GB, GD, GE, GH, GM, GT, HN, HR, HU, ID, IL, IN, IR, IS, 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, 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.
- (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, GH, GM, KE, LR, LS,

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

(54) 発明の名称: 機械学習プログラム、機械学習方法および機械学習装置



(57) **Abstract:** This invention makes it possible to avoid overtraining by causing a computer to execute processing for: clustering a plurality of data pieces; generating a model based on machine learning using data pieces sorted into a first group by the clustering; and using data pieces sorted into a second group by the clustering to verify the output accuracy of the generated model.



WO 2021/229630 A1

MW, MZ, NA, RW, 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, 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：複数のデータをクラスタリングし、クラスタリングで第1のグループに分類されたデータを用いた機械学習によってモデルを生成し、クラスタリングで第2のグループに分類されたデータを用いて、生成されたモデルの出力精度を検証する、処理をコンピュータに実行させることで、過学習を回避できるようにする。

明 細 書

発明の名称：

機械学習プログラム，機械学習方法および機械学習装置

技術分野

[0001] 本発明は、機械学習技術に関する。

背景技術

[0002] 機械学習の手法の一つとして、入力データと出力データとを備える教師データ（正解付きデータ）を用いて入出力関係を学習させる教師あり学習が知られている。

[0003] また、一般的に、教師あり機械学習においては、複数の教師データのうちの一部をモデル作成用データ（訓練データ）として使用することでモデル（機械学習モデル）を作成し、複数の教師データのうち残りの一部をモデル検証用データ（評価データ）として使用することで過学習の判断を行なうことが知られている。

[0004] ここで、過学習（overfitting）とは、モデルが訓練データだけに最適化されてしまい汎用性がない状態に陥ることをいい、モデル作成用データについては高精度に予測できるが、それ以外のデータについては予測が低精度となる。

[0005] 上述した教師データの一部をモデル検証用データとして用いた過学習の判断手法においては、作成したモデルを用いてモデル作成用データを予測した時の予測精度とモデル検証用データを予測した時の予測精度が大幅に異なる場合に過学習の状態と判断される。

先行技術文献

特許文献

[0006] 特許文献1：特開2019-66993号公報

発明の概要

発明が解決しようとする課題

[0007] しかしながら、教師データ取得時に偏りがある場合、入力データ空間全体で見たときに、例外的なクラスタが存在する可能性がある。

[0008] なお、ここでいう「偏り」とは、入力に関するものであり、偶然では起こりえないほど入力が類似したデータ群（クラスタ）が教師データに含まれている状況を指す。取得可能な教師データに制限がある等の事情により、このような偏りが生じ得る。

[0009] このような例外的なクラスタのデータまで正しく予測するモデルを作成すると過学習になりやすいが、上述の通りモデル作成／検証用データのいずれも高精度に予測できてしまうため、過学習となったことが検出されない場合がある。

[0010] 1つの側面では、過学習を抑制することを目的とする。

課題を解決するための手段

[0011] このため、この機械学習プログラムは、複数のデータをクラスタリングし、前記クラスタリングで第1のグループに分類されたデータを用いた機械学習によってモデルを生成し、前記クラスタリングで第2のグループに分類されたデータを用いて、生成された前記モデルの出力精度を検証する、処理をコンピュータに実行させる。

発明の効果

[0012] 一実施形態によれば、過学習を抑制することができる。

図面の簡単な説明

[0013] [図1]実施形態の一例としての計算機システムのハードウェア構成を例示する図である。

[図2]実施形態の一例としての計算機システムの機能構成を例示する図である。

[図3]ニューラルネットワークの概要を示す図である。

[図4]実施形態の一例としての計算機システムのクラスタリング処理部によるクラスタリング手法を説明するための図である。

[図5]実施形態の一例としての計算機システムのデータ作成部による処理を説

明するための図である。

[図6]実施形態の一例としての計算機システムのモデル作成部による処理を説明するための図である。

[図7]実施形態の一例としての計算機システムの予測処理部による処理を説明するための図である。

[図8]実施形態の一例としての計算機システムの検証部による処理を説明するための図である。

[図9]実施形態の一例としての計算機システムにおける処理を説明するためのフローチャートである。

[図10]実施形態の一例としての計算機システムの学習処理部により行なわれる二値分類を説明するための図である。

[図11]図10におけるモデル作成用クラスタを抽出して示す図である。

[図12]図10におけるモデル検証用クラスタを抽出して示す図である。

[図13]機械学習方法における過学習について説明するための図である。

発明を実施するための形態

[0014] 図13は機械学習方法における過学習について説明するための図であり、入力データ空間に配置された教師データを例示する。この図13においては、多数の微小点が配置された入力データ空間を例示している。微小点のそれぞれは教師データを表しており、それぞれ入力データに応じた位置にプロットされている。

[0015] また、この入力データ空間においては、複数の教師データ（微小点）が局所的に集まることで複数の小規模のクラスタ（データ群）が形成されている。図13中においては、教師データの集合によるクラスタに符号aまたは符号bが付されている。

[0016] これらの符号aまたは符号bは教師データの出力を表しており、符号aが付されたクラスタを構成する教師データの出力はそれぞれaであり、符号bが付されたクラスタを構成する教師データの出力はそれぞれbである。すなわち、図13に示す例においては、aまたはbを予測する二値分類を表して

いる。

[0017] この図13に示す例において、太破線は、モデル作成用データを用いて全て正解できる高精度なモデルを作成した場合の予測の境界を示す。当該モデルにおいては、この太破線の左側に位置する教師データの出力をbと予測し、この太破線の右側に位置する教師データの出力をaと予測する。

[0018] ここで、この図13に例示する入力データ空間上の教師データには、モデル作成用のデータとモデル検証用のデータとが混在しており、各クラスを構成する教師データにも、モデル作成用のデータとモデル検証用のデータとが混在している。その場合、図13における太破線を境界として予測するモデルはモデル作成用データおよびモデル検証用データのいずれに対しても高精度に予測できる。

[0019] しかしながら、教師データ取得時に偏りがある場合、入力データ空間全体で見たときに、例外的なクラスが存在する可能性がある。

[0020] このような例外的なクラスデータのデータまで正しく予測するモデルを作成すると過学習になりやすいが、上述の通りモデル作成／検証用データのいずれも高精度に予測できてしまうため、過学習となったことが検出されない場合がある。例えば、図13中において、符号P1を付して示す四角点線で囲んだクラスが例外的なクラスであった場合、図13に示す入力空間においては、太破線を境界として予測するモデルは過学習の状態にあり、同図中に一点鎖線を境界として予測するモデルの方が望ましい。機械学習においては、このような過学習を抑制することが望まれている。

[0021] 以下、図面を参照して本機械学習プログラム、機械学習方法および機械学習装置にかかる実施の形態を説明する。ただし、以下に示す実施形態はあくまでも例示に過ぎず、実施形態で明示しない種々の変形例や技術の適用を排除する意図はない。すなわち、本実施形態を、その趣旨を逸脱しない範囲で種々変形（実施形態および各変形例を組み合わせる等）して実施することができる。また、各図は、図中に示す構成要素のみを備えるという趣旨ではなく、他の機能等を含むことができる。

[0022] 図1は実施形態の一例としての計算機システム1のハードウェア構成を例示する図である。計算機システム1は、機械学習装置であって、例えば、ニューラルネットワークを実現する。計算機システム1は、図1に示すように、CPU (Central Processing Unit) 10、メモリ11およびアクセラレータ12を備える。これらのCPU10、メモリ11およびアクセラレータ12は、通信バス13を介して相互に通信可能に接続されている。通信バス13は、本計算機システム1内のデータ通信を行なう。

[0023] メモリ11は、ROM (Read Only Memory) およびRAM (Random Access Memory) を含む記憶メモリである。メモリ11のROMには、後述するCPU10によって実行されるプログラムやこのプログラム用のデータ類が書き込まれている。メモリ11上のソフトウェアプログラムは、CPU10に適宜読み込まれて実行される。また、メモリ11のRAMは、一次記憶メモリあるいはワーキングメモリとして利用される。メモリ11のRAMには、教師データ (モデル作成用データ、モデル検証用データ) やモデルを構成する情報およびモデルを用いた予測結果等も格納される。アクセラレータ12は、例えば、行列演算などのニューラルネットワークの計算に必要な演算処理を実行する。

[0024] CPU10は、種々の制御や演算を行なう処理装置 (プロセッサ) であり、実装されたプログラムに基づき、計算機システム1全体を制御する。そして、このCPU10がメモリ11等に格納された機械学習プログラム (図示省略) を実行することで、後述する学習処理部100 (図2参照) としての機能を実現する。計算機システム1は、機械学習プログラムを実行することにより機械学習装置として機能する。

[0025] なお、学習処理部100としての機能を実現するためのプログラム (機械学習プログラム) は、例えばフレキシブルディスク、CD (CD-ROM, CD-R, CD-RW等)、DVD (DVD-ROM, DVD-RAM, DVD-R, DVD+R, DVD-RW, DVD+RW, HD DVD等)、ブルーレイディスク、磁気ディスク、光ディスク、光磁気ディスク等の、コ

コンピュータ読取可能な記録媒体に記録された形態で提供される。そして、コンピュータ（計算機システム1）はその記録媒体からプログラムを読み取って内部記憶装置または外部記憶装置に転送し格納して用いる。また、そのプログラムを、例えば磁気ディスク、光ディスク、光磁気ディスク等の記憶装置（記録媒体）に記録しておき、その記憶装置から通信経路を介してコンピュータに提供するようにしてもよい。

[0026] 学習処理部100としての機能を実現する際には、内部記憶装置（本実施形態ではメモリ11のRAMやROM）に格納されたプログラムがコンピュータのマイクロプロセッサ（本実施形態ではCPU10）によって実行される。このとき、記録媒体に記録されたプログラムをコンピュータが読み取って実行するようにしてもよい。

[0027] 図2は実施形態の一例としての計算機システム1の機能構成を例示する図である。計算機システム1は、図2に示すように、学習処理部100としての機能を備える。学習処理部100は、例えば、ニューラルネットワークにおける深層学習を実施する。

[0028] ニューラルネットワークは、ハードウェア回路であってもよいし、CPU10等によりコンピュータプログラム上で仮想的に構築される階層間を接続するソフトウェアによる仮想的なネットワークであってもよい。

[0029] 図3にニューラルネットワークの概要を示す。図3に示すニューラルネットワークは、入力層と出力層との間に複数の隠れ層を含むディープニューラルネットワークである。隠れ層は、例えば、畳み込み層、プーリング層または全結合層等である。図3中において、各層に示す丸印は、所定の計算をそれぞれ実行するノードを示す。

[0030] ニューラルネットワークは、入力データを入力層に入力し、畳み込み層やプーリング層などで構成される隠れ層にて所定の計算を順次実行することで、演算により得られる情報を入力側から出力側に順次伝えるフォワード方向の処理（順伝播処理）を実行する。フォワード方向の処理の実行後、出力層から出力される出力データと正解データとから得られる誤差関数の値を

小さくするために、フォワード方向の処理で使用するパラメータを決定するバックワード方向の処理（逆伝播処理）を実行する。そして、逆伝播処理の結果に基づいて重み等の変数を更新する更新処理が実行される。

[0031] 図2に示すように、学習処理部100は、クラスタリング処理部101、データ作成部102、モデル作成部103、予測処理部104および検証部105を備える。

[0032] クラスタリング処理部101は、複数の教師データに対して、偏りが認識できるようにクラスタリングを行なうことで、複数のクラスタ（データ群）を作成する。教師データは、予め図示しない記憶装置に格納されてもよく、本計算機システム1の外部から入力されてもよい。クラスタリング処理部101は、複数の教師データに対して階層型クラスタリングを行なう。

[0033] 図4は実施形態の一例としての計算機システム1のクラスタリング処理部101によるクラスタリング手法を説明するための図である。この図4においては、階層型クラスタリングにおけるデンドログラム（樹形図）を例示している。

[0034] 階層型クラスタリングにおいては、複数の入力データに対して、データ間の距離に応じて結合（グルーピング、マージ）することを繰り返し行なうことでクラスタリングを実現する。

[0035] 本計算機システム1において、クラスタリング処理部101は、最遠隣法により、クラスタリングを実現する。なお、最遠隣法におけるデータ間の距離は、例えば、ユーグリッド距離を用いてもよく、適宜変更して実施することができる。

[0036] また、階層型クラスタリングにおいては、例えば、システム管理者等が、同一クラスタとするためのデータ間の距離を閾値として設定することができる。この閾値を設定することで、クラスタリング処理部101は、データ間の距離が閾値未満となるデータどうしを同一のクラスタとなるようにクラスタリングする。閾値は、クラスタのマージ停止条件に相当し、例えば、システム管理者等が任意に設定してもよい。図4においては、符号D0～D9で

表されるデータに対して階層型クラスタリングを行なう例を示しており、閾値=5が設定されている。

[0037] 隣接する入力データ間の距離が近いものから順に結合（グルーピング，マージ）することで、例えば、データD3，D4が一つのクラスタC1を形成している。同様に、データD8，D5，D7がクラスタC2を、データD2，D1，D6がクラスタC5をそれぞれ形成している。データD0，D9はいずれも他のデータからの距離が遠いものであるため、それぞれ単独で独立したクラスタC3，C4を形成する。

[0038] これらのクラスタC1～C5は、各クラスタ内におけるデータ間の距離が閾値（図4に示す例では5）未満であることが保証されており、データ空間内におけるデータの偏りを実現する。

クラスタリング処理部101は、このような階層型クラスタリング手法を用いて教師データに偏りが認識されるクラスタリングを実現する。

[0039] また、クラスタのマージ停止条件（閾値）は、教師データ取得時の偏りによるとみなせる入力データ間の距離とすることが望ましい。この閾値は、例えば、対象データに対するドメイン知識を持つ人がデータの素性に基づいて任意に設定してもよい。

[0040] データ作成部102は、モデル作成用データ（教師データ）およびモデル検証用データを作成する。モデル作成用データは、後述するモデル作成部103が機械学習のモデルを作成するために用いる教師データである。モデル検証用データは、後述する検証部105が作成されたモデルの検証を行なうために用いる教師データである。

[0041] 以下、モデル学習用データを用いてモデル作成を行なう過程を学習フェーズ（第1フェーズ）という場合があり、モデル作成用データを用いてモデルの検証を行なう過程を検証フェーズ（第2フェーズ）という場合がある。図5は実施形態の一例としての計算機システム1のデータ作成部102による処理を説明するための図である。

[0042] データ作成部102は、クラスタリング処理部101により作成された複

数のクラスタを、モデル作成用クラスタとモデル検証用クラスタとに分類する。なお、モデル作成用クラスタおよびモデル検証用クラスタの各数は、適宜変更して実施することができる。例えば、複数のクラスタをモデル作成用クラスタまたはモデル検証用クラスタへランダムに振り分けることで分類してもよく、適宜変更して分類を実施することができる。なお、複数のクラスタのモデル作成用クラスタまたはモデル検証用クラスタへの分類は、クラスタリング処理部101が行なってもよく、適宜変更して実施することができる。

[0043] 本計算機システム1においては、異なるクラスタのデータを使用して機械学習と検証とを実行する。すなわち、複数のクラスタのうち第1のクラスタ（第1のグループ）のデータを用いて機械学習のモデルを作成し、第2のクラスタ（第2のグループ）のデータを用いてモデルの出力精度の検証を行なう。

モデル作成用クラスタは、機械学習によってモデルを生成するために用いられるデータによる第1のグループであってもよい。また、モデル検証用クラスタは、生成されたモデルの出力精度を検証するために用いられるデータによる第2のグループであってもよい。

[0044] データ作成部102は、複数のモデル作成用クラスタからデータを均等にサンプリング（抽出）して、モデル作成用データを作成する。複数のモデル作成用クラスタからデータを均等にサンプリングする理由は、複数のモデル作成用クラスタ間においてデータ数の偏りがある恐れがあるからである。データ作成部102は、複数のモデル作成用クラスタから異なるサンプリングを行なうことで、複数のモデル作成用データを作成する。

[0045] 同様に、データ作成部102は、複数のモデル検証用クラスタからデータを均等にサンプリング（抽出）して、モデル検証用データを作成する。複数のモデル検証用クラスタからデータを均等にサンプリングする理由は、複数のモデル検証用クラスタ間においてもデータ数の偏りがあるおそれがあるからである。データ作成部102は、複数のモデル検証用クラスタから異なる

サンプリングを行なうことで、複数のモデル検証用データを作成する。

[0046] 複数のモデル作成用クラスタ、複数のモデル検証用クラスタ、複数のモデル作成用データおよび複数のモデル検証用データは、それぞれメモリ 11 の所定の記憶領域に格納してもよく、また、図示しない記憶装置に格納してもよい。

[0047] モデル作成部 103 は、モデル作成用データ（教師データ）を用いた機械学習によってモデル（学習モデル）を作成する。モデルは、入力値を受け取り、何かしらの評価・判定をして出力値を出力する。モデルの出力を予測結果といってもよい。なお、モデルの作成は既知の手法を用いて実現することができ、モデル作成部 103 によるモデル作成手法の説明は省略する。また、モデル作成部 103 は、機械学習に複数のモデル検証用データを用いることで、これらのモデル作成用データに応じた複数のモデルを作成する。モデル検証用データは、クラスタリングで第 3 のグループに分類されたデータに相当する。図 6 は実施形態の一例としての計算機システム 1 のモデル作成部 103 による処理を説明するための図である。

[0048] 図 6 に示す例においては、2 つのモデル作成用データ # 1, # 2 が示されている。モデル作成部 103 は、モデル作成用データ # 1 を用いて教師あり学習（機械学習）を行なうことでモデル # 1 を作成し、モデル作成用データ # 2 を用いて教師あり学習（機械学習）を行なうことでモデル # 2 を作成する。作成されたモデル # 1, # 2 には、モデル作成用データやモデル検証用データが入力される。モデル作成用データ # 1 は、第 1 のグループに分類されたデータのうちの第 1 のデータに相当する。モデル作成用データ # 2 は、第 1 のグループに分類されたデータのうちの第 2 のデータに相当する。

[0049] 予測処理部 104 は、モデル作成部 103 が作成した複数のモデルを用いて、予測対象データをこれらのモデルに入力した場合の出力の予測を行なう。予測処理部 104 は、予測対象データをモデル作成部 103 が作成した複数のモデルのそれぞれに入力し、各モデルの出力（予測結果）をアンサンブル（統合、集計）する。予測処理部 104 は、このアンサンブル結果を最終

的な出力（予測結果）とする。予測処理部104は、複数のモデルの各出力を統合（アンサンブル）して一の出力を生成するアンサンブル処理部に相当する。

予測対象データには、第1フェーズにおいてはモデル作成用データが用いられ、第2フェーズにおいてはモデル検証用データが用いられる。すなわち、予測処理部104は、第1フェーズにおいては、モデル作成用データを複数のモデルのそれぞれに入力し、各モデルの出力をアンサンブルした結果を最終的な出力（予測結果）とする。

[0050] また、予測処理部104は、第2フェーズにおいては、モデル検証用データを複数のモデルのそれぞれに入力し、各モデルの出力をアンサンブルした結果を最終的な出力（予測結果）とする。

図7は実施形態の一例としての計算機システム1の予測処理部104による処理を説明するための図である。この図7に示す例においては、2つのモデル#1、#2に予測対象データ、すなわち、モデル作成用データもしくはモデル検証用データが入力されている。各モデル#1、#2からそれぞれ出力される予測結果がアンサンブルされ、予測結果（予測対象データの予測結果）が出力される。

[0051] この図7に示す例において、予測対象データが第2のグループに分類されたデータに含まれる第3のデータに相当する。予測処理部104は、予測対象データ（第3のデータ）のモデル#1への入力に応じて当該モデル#1が出力した第1の結果と、予測対象データ（第3のデータ）のモデル#2への入力に応じて当該モデル#2が出力した第2の結果とに基づいて、第1の出力精度を算出する。

[0052] なお、複数のモデル出力のアンサンブルは平均値の演算等の既知の手法を用いて実現することができ、予測処理部104によるモデル出力のアンサンブル手法の説明は省略する。

[0053] 検証部105は、データ作成部102によって作成されたモデル検証用データを用いて、モデル作成部103が作成したモデルの検証を行なう。図8

は実施形態の一例としての計算機システム1の検証部105による処理を説明するための図である。検証部105は、データ作成部102によって作成されたモデル検証用データを用いて、モデル作成部103が作成したモデルの検証を行なう。

[0054] 検証部105は、データ作成部102によって作成された複数のモデル検証用データを、モデル作成部103によって作成された複数のモデルのそれぞれに入力させる。検証部105は、例えば、予測処理部104の機能を用いて、モデル検証用データ（予測対象データ）をモデル作成部103が作成した複数のモデルのそれぞれに入力し、各モデルの出力（予測結果）をアンサンブル（集計）する。予測処理部104は、このアンサンブル結果を最終的な出力（予測結果）とする。

[0055] 図8に示す例においては、モデル検証用データ#1がモデル#1、#2にそれぞれ入力され、各モデル#1、#2からそれぞれ出力される予測結果がアンサンブルされ、予測結果（モデル検証用データの予測結果）#1が出力されている。また、モデル検証用データ#2がモデル#1、#2にそれぞれ入力され、各モデル#1、#2からそれぞれ出力される予測結果がアンサンブルされ、予測結果（モデル検証用データの予測結果）#2が出力されている。

[0056] 検証部105は、予測結果#1を、モデル検証用データ#1の出力データと比較することで正答率（精度）を算出する。また、検証部105は、予測結果#2を、モデル検証用データ#2の出力データと比較することで正答率（精度）を算出する。検証部105は、これらの精度（正答率）の平均を算出することで、モデル検証用クラスタの精度を決定する。

[0057] すなわち、検証部105は、各モデル検証用データに対する予測精度の平均を算出して、モデル検証用クラスタについての最終的（全体的）な予測精度を取得する。

[0058] 例えば、検証部105は、モデル検証用データに基づいて出力された予測結果の精度と、モデル作成用データに基づいて出力された予測結果の精度と

の差が許容閾値内であるかを判断してもよい。すなわち、検証部105は、モデル検証用データに基づいて出力された予測結果の精度と、モデル作成用データに基づいて出力された予測結果の精度とが同レベルの精度であるかを判断してもよい。また、検証部105は、モデル検証用データに基づいて出力された予測結果の精度が所定の閾値以上であるかを判断してもよい。

[0059] 図8に示す例において、モデル検証用データ#1は、第2のグループに分類されたデータに含まれる第3のデータに相当する。モデル検証用データ#2は、第2のグループに分類されたデータに含まれる第4のデータに相当する。

[0060] 検証部105は、モデル検証用データ（第3のデータ）#1のモデル#1への入力に応じて当該モデル#1が出力した第1の結果と、モデル検証用データ（第3のデータ）#1のモデル#2への入力に応じて当該モデル#2が出力した第2の結果とに基づいて予測結果#1（第1の出力精度）を算出する。

[0061] また、検証部105は、モデル検証用データ（第4のデータ）#2のモデル#1への入力に応じて当該モデル#1が出力した第3の結果と、モデル検証用データ（第4のデータ）#2のモデル#2への入力に応じて当該モデル#2が出力した第4の結果とに基づいて予測結果#2（第2の出力精度）を算出する。検証部105は、これらの予測結果#1（第1の出力精度）と予測結果#2（第2の出力精度）とに基づいて予測精度の検証を行なう。

[0062] 上述の如く構成された実施形態の一例としての計算機システム1における処理を、図9に示すフローチャート（ステップS1～S4）に従って説明する。

[0063] ステップS1において、クラスタリング処理部101が、予め用意された教師データに対して階層型クラスタリングを行なうことで、偏りが認識できる複数のクラスタを作成する。データ作成部102は、クラスタリング処理部101が作成した複数のクラスタを、モデル作成用クラスタとモデル検証用クラスタとに分ける。

[0064] そして、データ作成部102は、複数のモデル作成用クラスタからデータを均等にサンプリングしてモデル作成用データを作成する。この際、データ作成部102は、複数のモデル作成用クラスタから異なるサンプリングを複数行なうことで、複数のモデル作成用データを作成する。

[0065] また、データ作成部102は、複数のモデル検証用クラスタからデータを均等にサンプリングしてモデル検証用データを作成する。この際、データ作成部102は、複数のモデル検証用クラスタから異なるサンプリングを複数行なうことで、複数のモデル検証用データを作成する。

ステップS2において、モデル作成部103は、機械学習にモデル作成用データ（教師データ）を用いてモデルを作成する。

[0066] ステップS3において、予測処理部104は、モデル作成部103が作成した複数のモデルを用いて、予測対象データをこれらのモデルに入力した場合の出力の予測を行なう。

[0067] ステップS4において、検証部105は、データ作成部102によって作成されたモデル検証用データを用いて、モデル作成部103が作成したモデルの検証を行なう。

[0068] このように、実施形態の一例としての計算機システム1によれば、クラスタリング処理部101が作成した一のクラスタを、データ作成部102がモデル作成用データもしくはモデル検証用データのいずれかに割り当てる。これにより、入力データ空間全体で見たときに例外的なクラスタが存在していても、同一クラスタ内のデータはモデル作成データとモデル検証用データのいずれか一方にしか含まれない。そのため、モデル作成データの予測精度とモデル検証用データの予測精度とが同時に高くなることはない。このように、同一クラスタ内のデータがモデル作成データとモデル検証用データとに分かれることがないため、過学習を回避することができる。

[0069] 図10は実施形態の一例としての計算機システム1の学習処理部100により行なわれる二値分類を説明するための図であり、入力データ空間に配置された教師データを例示する。この図10においては、多数の微小点が配置

された入力データ空間を例示している。微小点のそれぞれは教師データを表しており、それぞれ入力データに応じた位置にプロットされている。

[0070] また、この入力データ空間においては、破線の丸で囲まれた教師データの集合はモデル作成用クラスタを示し、実線の丸で囲まれた教師データの集合はモデル検証用クラスタを示す。

[0071] また、この図10においては、各クラスタに符号aまたは符号bが付されている。これらの符号aまたは符号bは教師データの出力を表しており、符号aが付されたクラスタを構成する教師データの出力はそれぞれaであり、符号bが付されたクラスタを構成する教師データの出力はそれぞれbである。すなわち、図10に示す例においては、aまたはbを予測する二値分類を表している。

[0072] この図10に示す例において、モデル作成用クラスタからサンプリングしたデータで高精度なモデルを作成すると、符号 α を付した太破線を境界として予測するモデルとなる。

[0073] 図11は図10におけるモデル作成用クラスタを抽出して示す図である。この図11に示すように、モデル作成用クラスタに関しては、符号 α を付した太破線の左側には全ての出力bが配置され、その右側には全ての出力aが配置されている。すなわち、モデル作成用クラスタからサンプリングしたデータに対する予測精度が高いことがわかる。

[0074] 図12は図10におけるモデル検証用クラスタを抽出して示す図である。この図12に示すように、モデル検証用クラスタに関しては、符号 α を付した太破線の左側には出力bとともに出力aも配置されており、図11に示したモデル作成用クラスタからサンプリングしたデータに比べて予測精度が低いことがわかる。すなわち、過学習していると判断できる。図10に示す例においては、符号 β を付した一点鎖線を境界として予測するモデルが過学習のない好適なモデルとなる。

[0075] クラスタリング処理部101が階層型クラスタリングを行なうことで、複数の教師データに対して偏りが認識できるようにクラスタリングを行なうこ

とができる。

- [0076] モデル作成部103が、機械学習に、モデル作成用クラスタに備えられる複数のクラスタデータ群のそれぞれから抽出（サンプリング）して生成したデータ（モデル作成用データ）を用いる。複数のクラスタから均等にサンプリングすることで取得したモデル作成用データを用いることで、モデルの出力精度を向上させることができる。
- [0077] 検証部105が、複数のモデル検証用データをそれぞれモデルに適用することで、複数のクラスタの各データを検証に反映させることができ、検出精度を向上させることができる。
- [0078] 開示の技術は上述した実施形態に限定されるものではなく、本実施形態の趣旨を逸脱しない範囲で種々変形して実施することができる。本実施形態の各構成および各処理は、必要に応じて取捨選択することができ、あるいは適宜組み合わせてもよい。
- [0079] 例えば、上述した実施形態においては、第1フェーズにおいて、データ作成部102が複数のモデル作成用データを作成し、モデル作成部103がこれらの複数のモデル作成用データを用いて複数のモデルを作成しているが、これに限定されるものではない。モデル作成部103は、全てのモデル作成用クラスタのデータを用いて一つのモデルを作成してもよい。
- [0080] なおこの場合、第2フェーズにおいては、上述した実施形態と同様に、複数のモデル検証用データを作成し、これらの複数のモデル検証用データをモデルにそれぞれ適用することが望ましい。そして、予測処理部104は、これらの複数の入力データに基づいて出力された複数の予測結果を用いて精度を求めることが望ましい。
- [0081] 検証を行なう際に、複数のクラスタのデータを一つにまとめた場合には、データ数が多いクラスタの精度が優先されてしまい検出精度が低下するおそれがある。そこで、複数のモデル検証用データをそれぞれモデルに適用することで、複数のクラスタの各データを検証に反映させることができ、検出精度を向上させることができる。

[0082] 上述した実施形態においては、機械学習をニューラルネットワークに適用した例を示しているが、これに限定されるものではなく、種々変形して実施することができる。また、上述した開示により本実施形態を当業者によって実施・製造することが可能である。

符号の説明

- [0083]
- 1 計算機システム
 - 10 CPU
 - 11 メモリ
 - 12 アクセラレータ
 - 13 通信バス
 - 100 学習処理部
 - 101 クラスタリング処理部
 - 102 データ作成部
 - 103 モデル作成部
 - 104 予測処理部
 - 105 検証部

請求の範囲

- [請求項1] 複数のデータをクラスタリングし、
前記クラスタリングで第1のグループに分類されたデータを用いた機械学習によってモデルを生成し、
前記クラスタリングで第2のグループに分類されたデータを用いて、生成された前記モデルの出力精度を検証する、
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] 前記クラスタリングは、階層型クラスタリングである、
ことを特徴とする請求項1に記載の機械学習プログラム。
- [請求項3] 前記モデルを生成する処理は、前記クラスタリングで第3のグループに分類されたデータを用いた機械学習を含み、
前記検証する処理は、前記クラスタリングで第4のグループに分類されたデータを用いて実行される、
ことを特徴とする請求項1または2に記載の機械学習プログラム。
- [請求項4] 前記モデルは、前記第1のグループに分類されたデータのうち第1のデータを用いた機械学習により生成され、
前記第1のグループに分類されたデータのうち第2のデータを用いた機械学習により他のモデルを生成する、
処理を前記コンピュータに実行させることを特徴とする請求項1乃至3のいずれか1項に記載の機械学習プログラム。
- [請求項5] 前記検証する処理は、前記第2のグループに分類されたデータに含まれる第3のデータの前記モデルへの入力に応じて前記モデルが出力した第1の結果と、前記第3のデータの前記他のモデルへの入力に応じて前記他のモデルが出力した第2の結果とに基づいて、第1の出力精度を算出する処理を含む、
ことを特徴とする請求項4に記載の機械学習プログラム。
- [請求項6] 前記検証する処理は、前記第2のグループに分類されたデータに含

まれる第4のデータの前記モデルへの入力に応じて前記モデルが出力した第3の結果と、前記第4のデータの前記他のモデルへの入力に応じて前記他のモデルが出力した第4の結果とに基づいて算出された第2の出力精度と、前記第1の出力精度とに基づいて実行される、ことを特徴とする請求項5項に記載の機械学習プログラム。

[請求項7] 複数のデータをクラスタリングする処理と、
前記クラスタリングで第1のグループに分類されたデータを用いた機械学習によってモデルを生成する処理と、
前記クラスタリングで第2のグループに分類されたデータを用いて、生成された前記モデルの出力精度を検証する処理と、
を行なうことを特徴とする機械学習方法。

[請求項8] 前記クラスタリングは、階層型クラスタリングである、
ことを特徴とする、請求項7に記載の機械学習方法。

[請求項9] 前記モデルを生成する処理は、前記クラスタリングで第3のグループに分類されたデータを用いた機械学習を含み、
前記検証する処理は、前記クラスタリングで第4のグループに分類されたデータを用いて実行される、
ことを特徴とする、請求項7または8に記載の機械学習方法。

[請求項10] 前記モデルは、前記第1のグループに分類されたデータのうち第1のデータを用いた機械学習により生成され、
前記第1のグループに分類されたデータのうち第2のデータを用いた機械学習により他のモデルを生成する、
ことを特徴とする、請求項7乃至9のいずれか1項に記載の機械学習方法。

[請求項11] 前記検証する処理は、前記第2のグループに分類されたデータに含まれる第3のデータの前記モデルへの入力に応じて前記モデルが出力した第1の結果と、前記第3のデータの前記他のモデルへの入力に応じて前記他のモデルが出力した第2の結果とに基づいて、第1の出力

精度を算出する処理を含む、

ことを特徴とする、請求項10に記載の機械学習方法。

[請求項12] 前記検証する処理は、前記第2のグループに分類されたデータに含まれる第4のデータの前記モデルへの入力に応じて前記モデルが出力した第3の結果と、前記第4のデータの前記他のモデルへの入力に応じて前記他のモデルが出力した第4の結果とに基づいて算出された第2の出力精度と、前記第1の出力精度とに基づいて実行される、

ことを特徴とする、請求項11に記載の機械学習方法。

[請求項13] 複数のデータをクラスタリングし、
前記クラスタリングで第1のグループに分類されたデータを用いた機械学習によってモデルを生成し、

前記クラスタリングで第2のグループに分類されたデータを用いて、生成された前記モデルの出力精度を検証する、

処理部を有することを特徴とする機械学習装置。

[請求項14] 前記クラスタリングは、階層型クラスタリングである、
ことを特徴とする請求項13に記載の機械学習装置。

[請求項15] 前記モデルを生成する処理は、前記クラスタリングで第3のグループに分類されたデータを用いた機械学習を含み、
前記検証する処理は、前記クラスタリングで第4のグループに分類されたデータを用いて実行される、
ことを特徴とする請求項13または14に記載の機械学習装置。

[請求項16] 前記モデルは、前記第1のグループに分類されたデータのうち第1のデータを用いた機械学習により生成され、
前記第1のグループに分類されたデータのうち第2のデータを用いた機械学習により他のモデルを生成する、

ことを特徴とする請求項13乃至15のいずれか1項に記載の機械学習装置。

[請求項17] 前記検証する処理は、前記第2のグループに分類されたデータに含

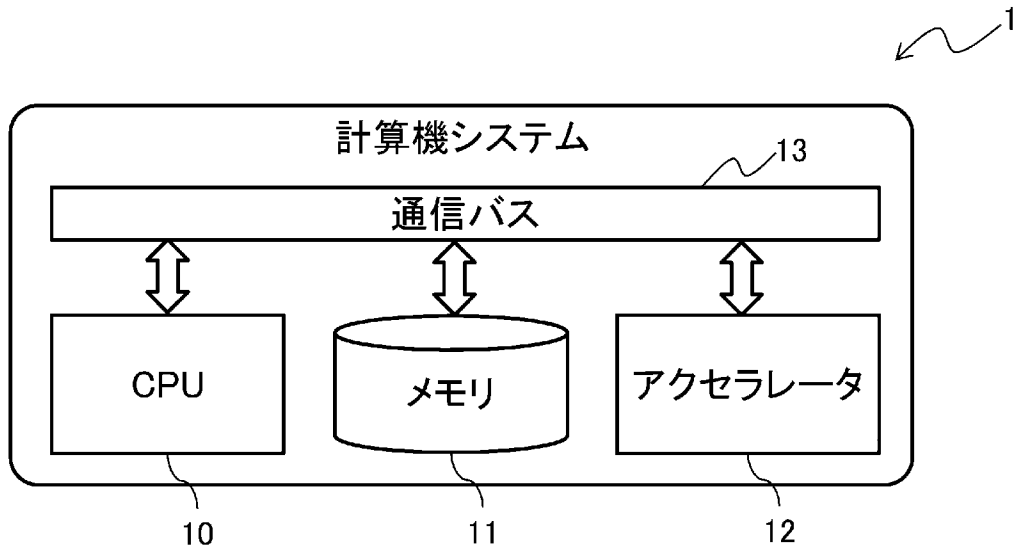
まれる第3のデータの前記モデルへの入力に応じて前記モデルが出力した第1の結果と、前記第3のデータの前記他のモデルへの入力に応じて前記他のモデルが出力した第2の結果とに基づいて、第1の出力精度を算出する処理を含む、

ことを特徴とする請求項16に記載の機械学習装置。

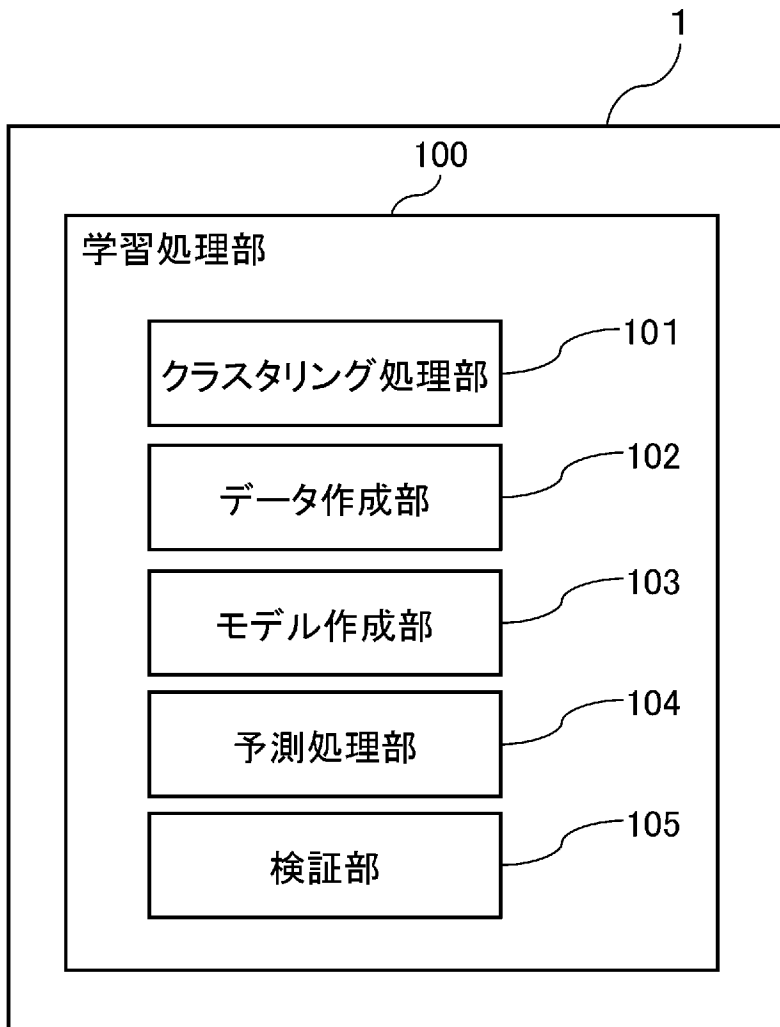
[請求項18]

前記検証する処理は、前記第2のグループに分類されたデータに含まれる第4のデータの前記モデルへの入力に応じて前記モデルが出力した第3の結果と、前記第4のデータの前記他のモデルへの入力に応じて前記他のモデルが出力した第4の結果とに基づいて算出された第2の出力精度と、前記第1の出力精度とに基づいて実行される、ことを特徴とする請求項17項に記載の機械学習装置。

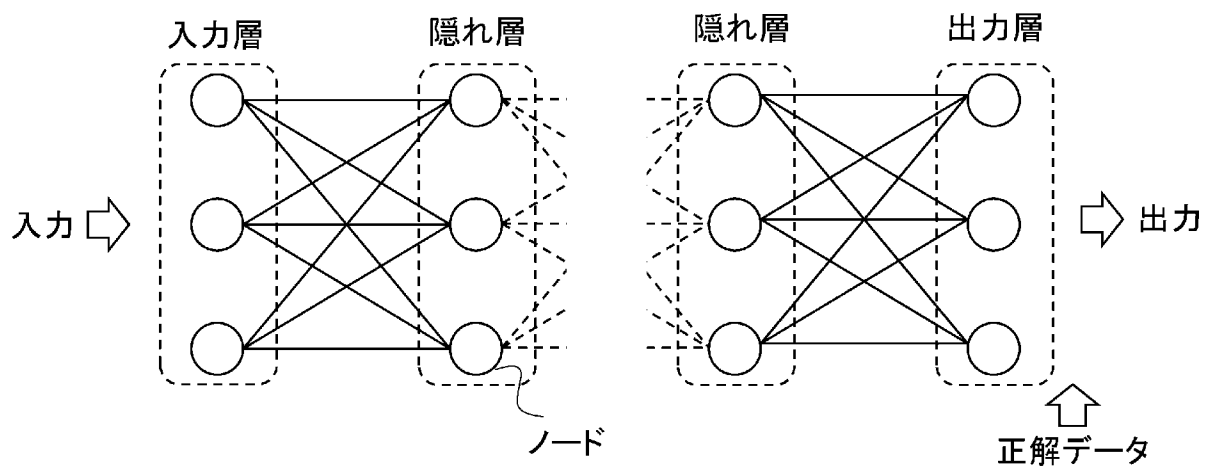
[図1]



[図2]

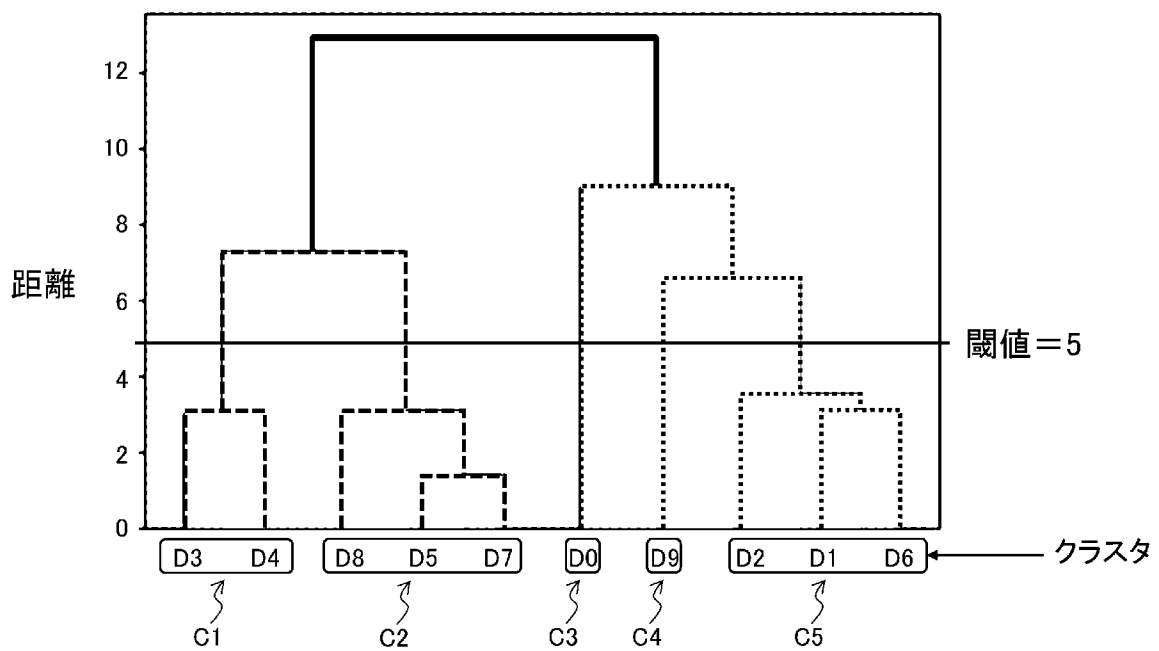


[図3]

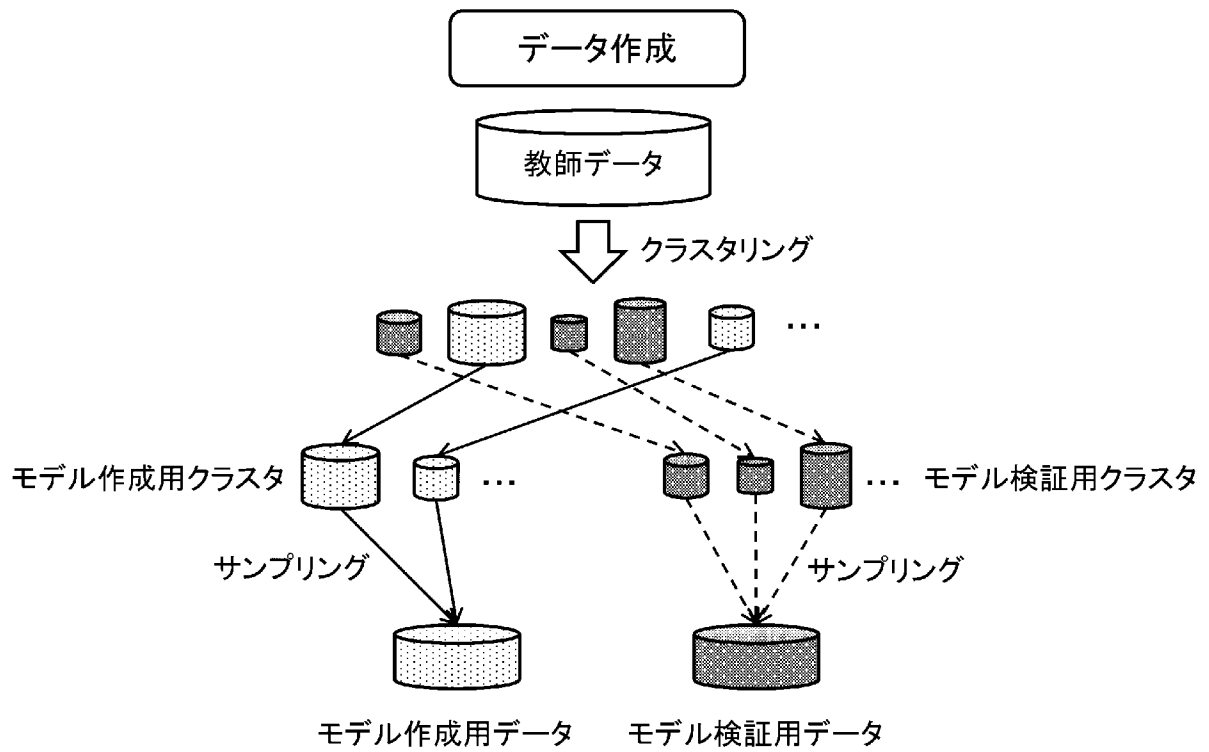


[図4]

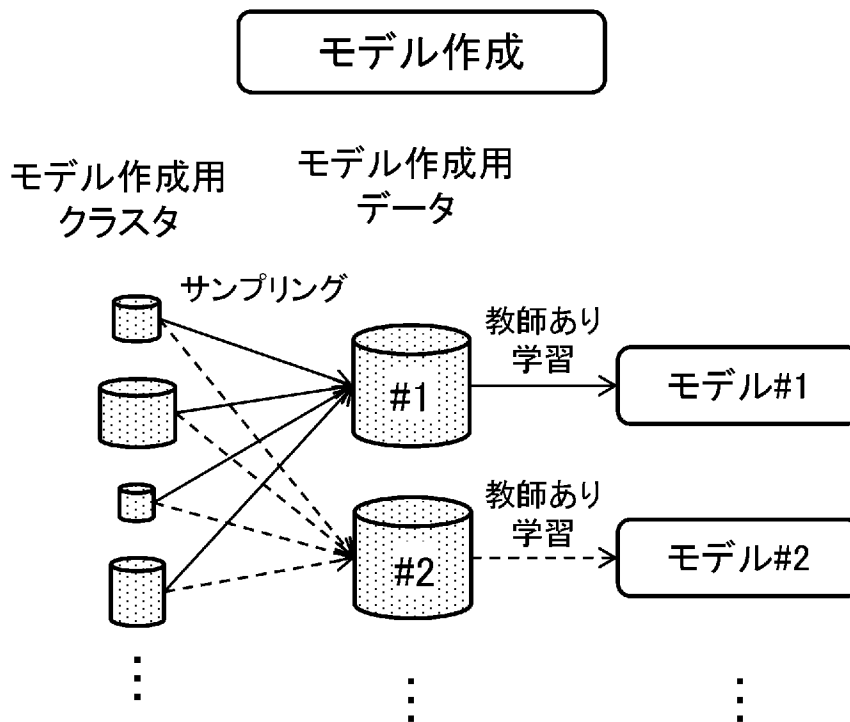
階層型クラスタリング



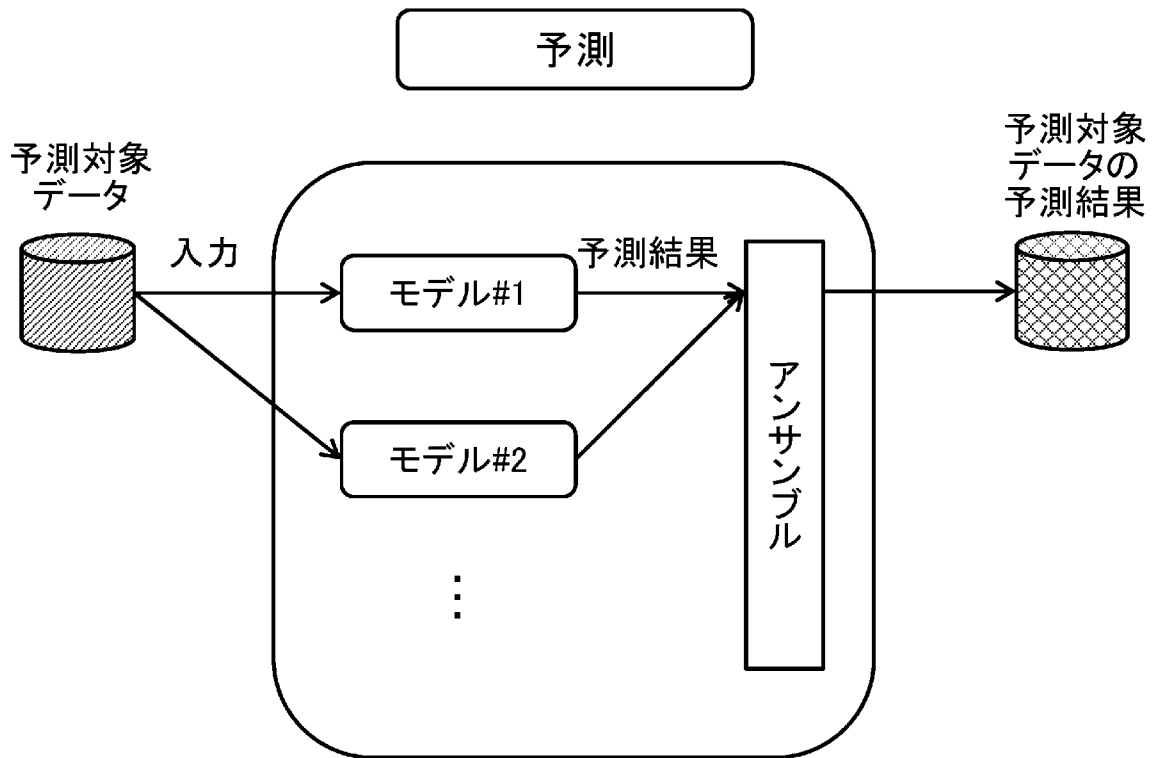
[図5]



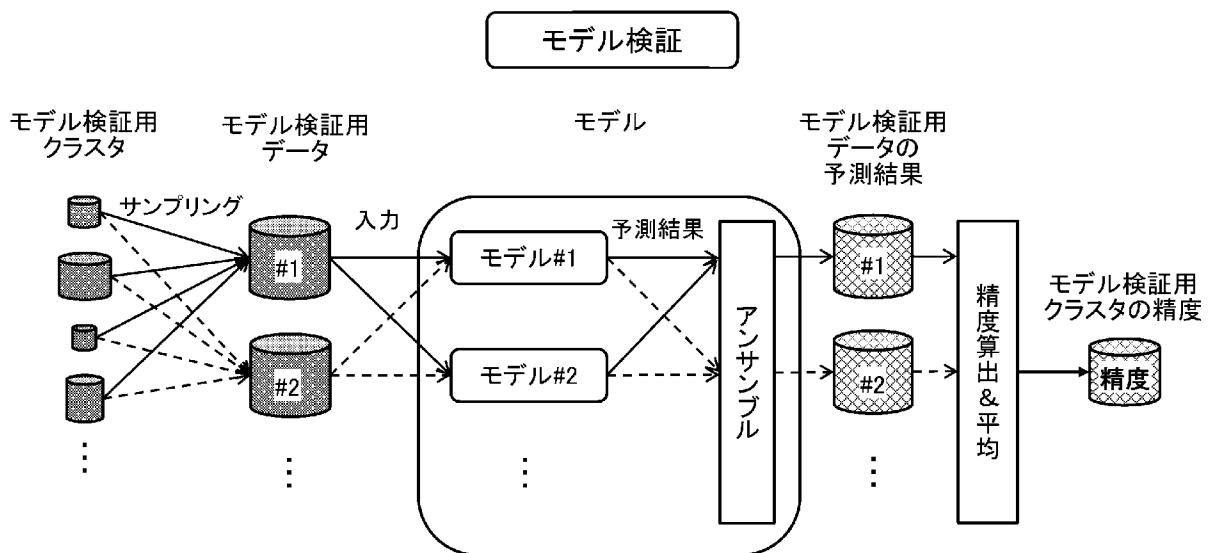
[図6]



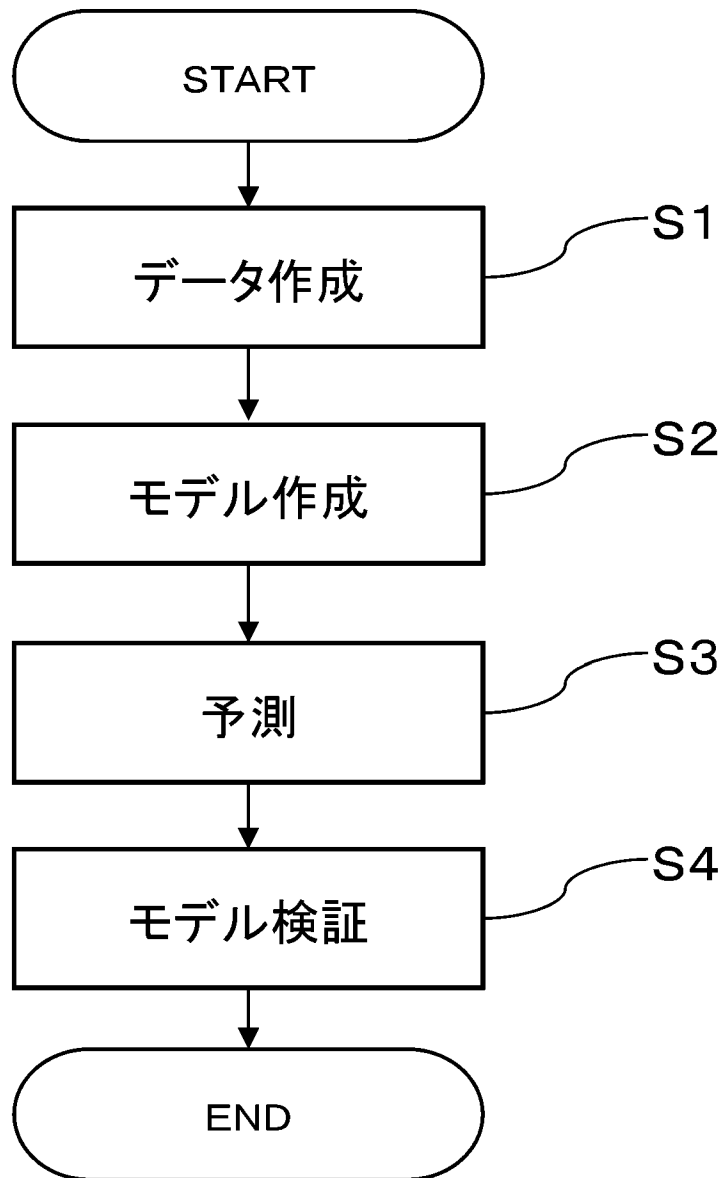
[図7]



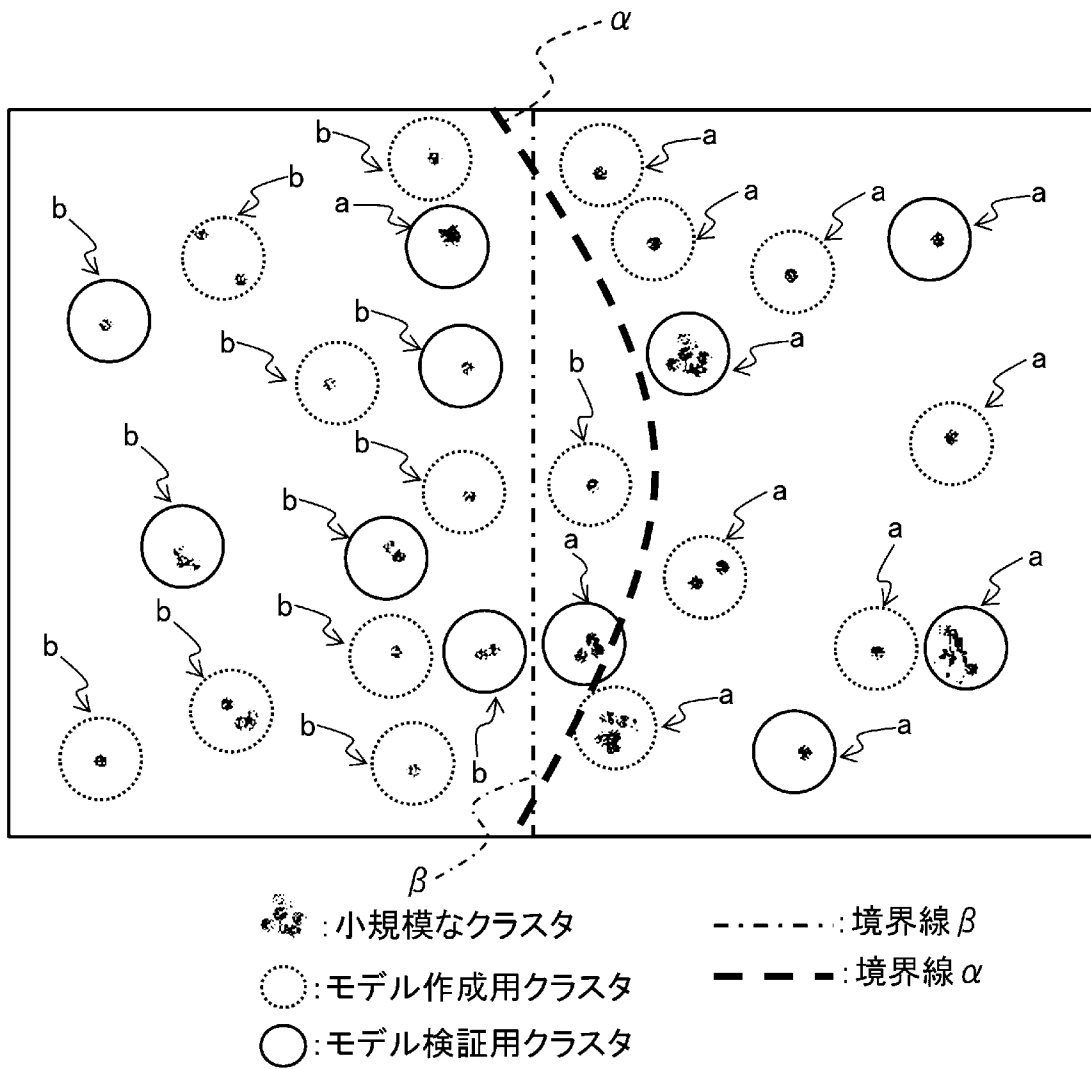
[図8]



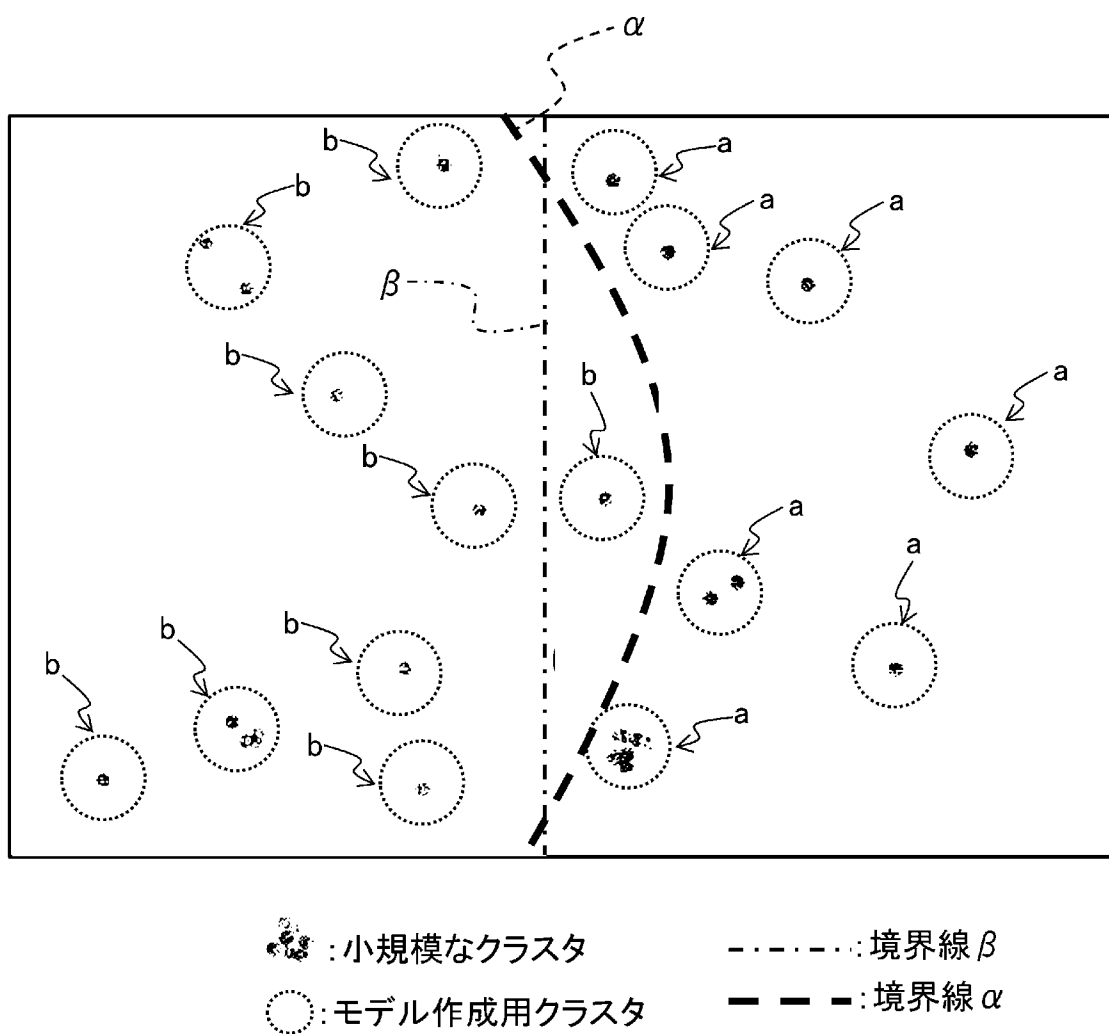
[図9]



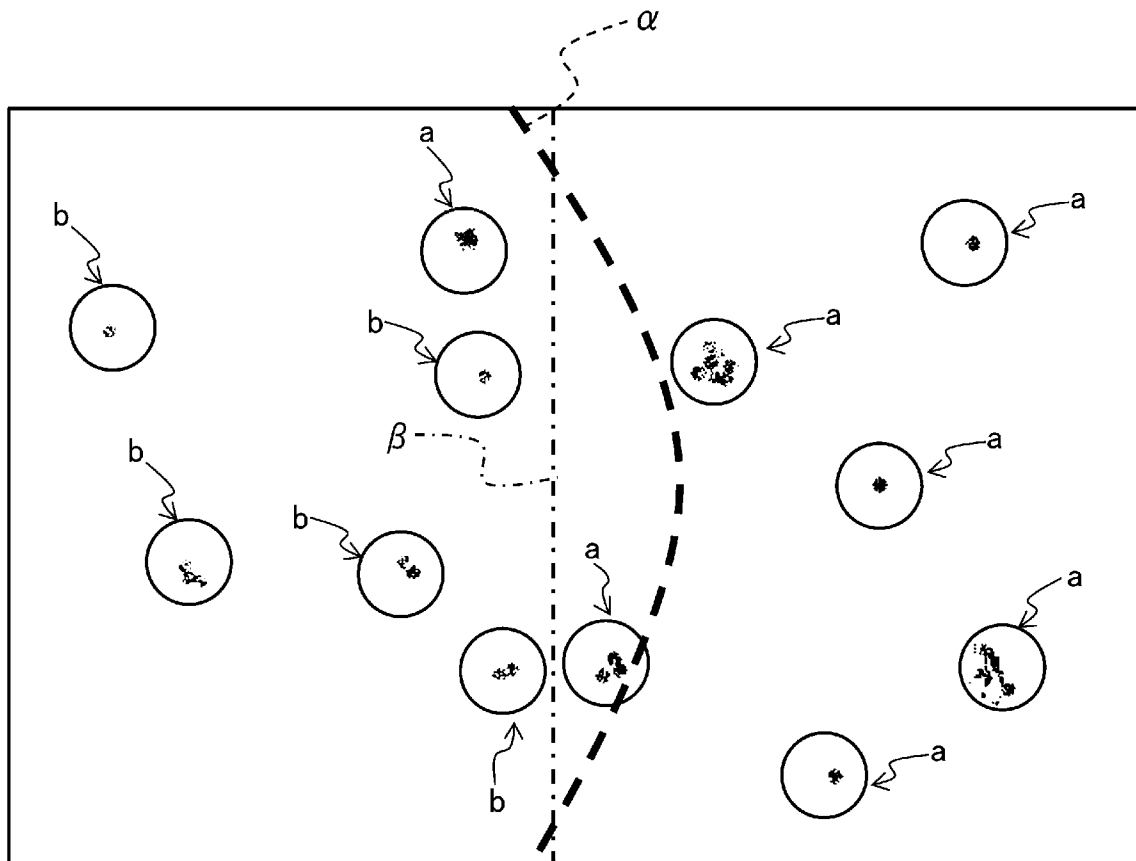
[図10]





[図11]





[図12]



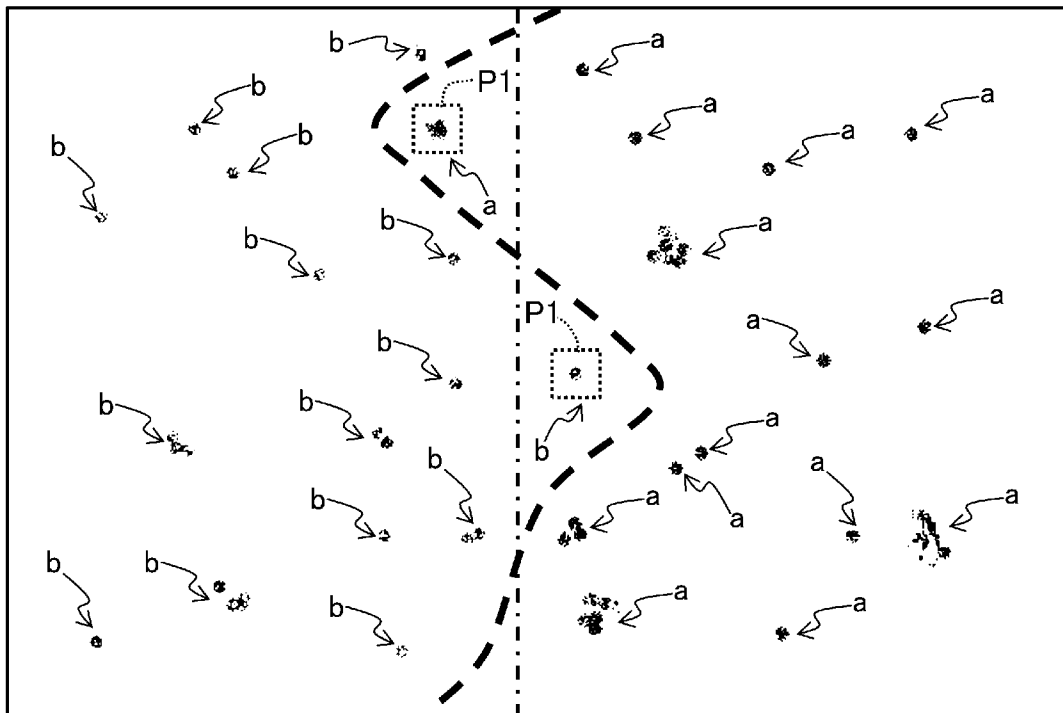
 : 小規模なクラスター


 : モデル検証用クラスター


 : 境界線 β

 : 境界線 α

[図13]



 : 小規模なクラスタ

 : 例外的なクラスタ

----- : 境界線

----- : 境界線

INTERNATIONAL SEARCH REPORT

International application No.

PCT/JP2020/018777

A. CLASSIFICATION OF SUBJECT MATTER

Int. Cl. G06N20/00 (2019.01) i
FI: G06N20/00

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)

Int. Cl. G06N20/00

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 1971-2020
Registered utility model specifications of Japan 1996-2020
Published registered utility model applications of Japan 1994-2020

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.
X	JP 2019-159538 A (OMRON CORP.) 19 September 2019, paragraphs [0006], [0019], [0020], [0033], [0034]	1-3, 7-9, 13-15
A		4-6, 10-12, 16-18
A	JP 2016-133895 A (CANON INC.) 25 July 2016, paragraphs [0034]-[0049], fig. 5	4-6, 10-12, 16-18
A	JP 2018-190125 A (NIPPON TELEGRAPH AND TELEPHONE CORP.) 29 November 2018, paragraphs [0026]-[0032]	4-6, 10-12, 16-18
A	JP 2019-45929 A (CANON INC.) 22 March 2019, paragraphs [0003]-[0007]	1-18

Further documents are listed in the continuation of Box C.

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
13.07.2020

Date of mailing of the international search report
21.07.2020

Name and mailing address of the ISA/
Japan Patent Office
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/JP2020/018777

Patent Documents referred to in the Report	Publication Date	Patent Family	Publication Date
JP 2019-159538 A	19.09.2019	(Family: none)	
JP 2016-133895 A	25.07.2016	US 2016/0210535 A1 paragraphs [0048]- [0061], fig. 5	
JP 2018-190125 A	29.11.2018	(Family: none)	
JP 2019-45929 A	22.03.2019	(Family: none)	

A. 発明の属する分野の分類（国際特許分類（IPC）） G06N 20/00(2019.01)i FI: G06N20/00		
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）） G06N20/00 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 - 2020年 日本国実用新案登録公報 1996 - 2020年 日本国登録実用新案公報 1994 - 2020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
X	JP 2019-159538 A (オムロン株式会社) 19.09.2019 (2019 - 09 - 19) [0006],[0019]-[0020],[0033]-[0034]	1-3,7-9,13-15
A		4-6,10-12,16-18
A	JP 2016-133895 A (キヤノン株式会社) 25.07.2016 (2016 - 07 - 25) [0034]-[0049],[図5]	4-6,10-12,16-18
A	JP 2018-190125 A (日本電信電話株式会社) 29.11.2018 (2018 - 11 - 29) [0026]-[0032]	4-6,10-12,16-18
A	JP 2019-45929 A (キヤノン株式会社) 22.03.2019 (2019 - 03 - 22) [0003]-[0007]	1-18
<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、一般的な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 13.07.2020	国際調査報告の発送日 21.07.2020	
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） 渡部 博樹 5B 3868 電話番号 03-3581-1101 内線 3545	

国際調査報告
特許ファミリーに関する情報

国際出願番号

PCT/JP2020/018777

引用文献	公表日	特許ファミリー文献	公表日
JP 2019-159538 A	19.09.2019	(ファミリーなし)	
JP 2016-133895 A	25.07.2016	US 2016/0210535 A1 [0048]-[0061], Fig. 5	
JP 2018-190125 A	29.11.2018	(ファミリーなし)	
JP 2019-45929 A	22.03.2019	(ファミリーなし)	