

(19)日本国特許庁(JP)

(12)特許公報(B1)

(11)特許番号
特許第7112802号
(P7112802)

(45)発行日 令和4年8月4日(2022.8.4)

(24)登録日 令和4年7月27日(2022.7.27)

(51)国際特許分類 F I
G 0 6 N 20/00 (2019.01) G 0 6 N 20/00 1 3 0
G 0 6 N 3/04 (2006.01) G 0 6 N 3/04

請求項の数 9 (全20頁)

(21)出願番号	特願2022-73380(P2022-73380)	(73)特許権者	520040304 窪田 望 東京都港区港南1丁目9番36号 NTT DATA 品川ビル13階 株式会社クリ エーターズネクスト内
(22)出願日	令和4年4月27日(2022.4.27)	(74)代理人	100079108 弁理士 稲葉 良幸
審査請求日	令和4年4月27日(2022.4.27)	(74)代理人	100109346 弁理士 大貫 敏史
早期審査対象出願		(74)代理人	100117189 弁理士 江口 昭彦
		(74)代理人	100134120 弁理士 内藤 和彦
		(72)発明者	窪田 望 東京都港区港南1丁目9番36号 NTT 最終頁に続く

(54)【発明の名称】 学習モデルの軽量化

(57)【特許請求の範囲】

【請求項1】

情報処理装置に含まれる1又は複数のプロセッサが、
 所定の学習データを取得すること、
 ニューラルネットワークを用いる所定の学習モデルに対して、蒸留処理された第1学習
 モデル、枝刈り処理された第2学習モデル、及び量子化処理された第3学習モデルの少な
 くとも2つのモデルを含む各モデルがそれぞれ重み付けされた重み学習モデルに、所定の
 データを入力して機械学習を行うこと、
 前記各モデルそれぞれの重みを変更された重み学習モデルごとに、前記所定の学習デー
 タを入力して前記機械学習が行われた場合の学習結果を取得すること、
 変更された各重みが重み付けされた各重み学習モデルと、前記各重み学習モデルで学習
 されたときの各学習結果とを含む学習データを用いて、教師あり学習を行うこと、
 前記教師あり学習により、任意の学習データを入力する場合に、各重みの集合ごとに学
 習結果を予測する予測モデルを生成すること、
 を実行する、情報処理方法。

【請求項2】

前記1又は複数のプロセッサは、
 任意の学習データを前記予測モデルに入力し、前記各重みの集合ごとに、前記重み学習
 モデルを実行した場合の学習結果を予測することを実行する、請求項1に記載の情報処理
 方法。

【請求項 3】

前記 1 又は複数のプロセッサが、
前記任意の学習データを前記所定の学習モデルに入力した場合の学習結果と、前記予測モデルにより予測された学習結果とが、軽量化に関する所定条件を満たすか否かを判定すること、

前記所定条件の判定結果に基づいて、前記各重みの有効性を判定すること、
をさらに実行する請求項 2 に記載の情報処理方法。

【請求項 4】

前記 1 又は複数のプロセッサが、
前記軽量化に関する所定条件に関するユーザ操作を受け付けること、
前記ユーザ操作に基づいて、前記軽量化に関する所定条件を設定すること、
をさらに実行する請求項 3 に記載の情報処理方法。

10

【請求項 5】

前記プロセッサは、
前記学習結果に含まれる学習精度を第 1 変数、前記学習結果に含まれるモデルサイズに関する値を第 2 変数とし、前記第 1 変数及び前記第 2 変数と、前記各重みとを対応付ける関係情報を生成すること、
を実行する請求項 1 に記載の情報処理方法。

【請求項 6】

前記プロセッサは、
前記第 1 変数の第 1 値及び前記第 2 変数の第 2 値を取得すること、
前記関係情報に基づいて、前記第 1 値及び前記第 2 値に対応する各重みを特定すること、
を実行する請求項 5 に記載の情報処理方法。

20

【請求項 7】

前記重み学習モデルは、前記第 1 学習モデル、前記第 2 学習モデル、及び前記第 3 学習モデルそれぞれに重みが付与されて線形結合されたモデルを含む、請求項 1 に記載の情報処理方法。

【請求項 8】

メモリと、1 又は複数のプロセッサとを備える情報処理装置であって、
前記メモリは、
ニューラルネットワークを用いる所定の学習モデルと、
前記所定の学習モデルに対して、蒸留処理された第 1 学習モデル、枝刈り処理された第 2 学習モデル、及び量子化処理された第 3 学習モデルの少なくとも 2 つのモデルを含む各モデルがそれぞれ重み付けされた重み学習モデルと、を記憶し、
前記 1 又は複数のプロセッサは、
所定の学習データを取得すること、
前記重み学習モデルに、前記所定の学習データを入力して機械学習を行うこと、
前記各モデルそれぞれの重みを変更された重み学習モデルごとに、前記所定の学習データを入力して前記機械学習が行われた場合の学習結果を取得すること、
変更された各重みが重み付けされた各重み学習モデルと、前記各重み学習モデルで学習されたときの各学習結果とを含む学習データを用いて、教師あり学習を行うこと、
前記教師あり学習により、任意の学習データを入力する場合に、各重みの集合ごとに学習結果を予測する予測モデルを生成すること、
を実行する、情報処理装置。

30

40

【請求項 9】

情報処理装置に含まれる 1 又は複数のプロセッサに、
所定の学習データを取得すること、
ニューラルネットワークを用いる所定の学習モデルに対して、蒸留処理された第 1 学習モデル、枝刈り処理された第 2 学習モデル、及び量子化処理された第 3 学習モデルの少なくとも 2 つのモデルを含む各モデルがそれぞれ重み付けされた重み学習モデルに、所定の

50

学習データを入力して機械学習を行うこと、

前記各モデルそれぞれの重みを変更された重み学習モデルごとに、前記所定の学習データを入力して前記機械学習が行われた場合の学習結果を取得すること、

変更された各重みが重み付けされた各重み学習モデルと、前記各重み学習モデルで学習されたときの各学習結果とを含む学習データを用いて、教師あり学習を行うこと、

前記教師あり学習により、任意の学習データを入力する場合に、各重みの集合ごとに学習結果を予測する予測モデルを生成すること、

を実行させる、プログラムを記録したコンピュータ読み取り可能な非一時的な記憶媒体。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、学習モデルの軽量化に関する情報処理方法、プログラム及び情報処理装置に関する。

【背景技術】

【0002】

近年、学習モデルを軽量化する研究が行われている。例えば、下記特許文献1には、パラメータの量子化を用いて学習モデルの軽量化を行う技術が記載されている。

【先行技術文献】

【特許文献】

【0003】

【文献】特開2019-133628号公報

【発明の概要】

【発明が解決しようとする課題】

【0004】

ここで、学習モデルを軽量化するには、枝刈り(Pruning)、量子化(Quantization)、蒸留(Distillation)などの手法がある。これらの軽量化手法の少なくとも1つを学習モデルに適用し、エンジニアが適宜パラメータ調整等を行って軽量化が行われる。

【0005】

しかしながら、学習データや学習モデルなどの様々な条件に応じて適切な軽量化手法は異なるものと思われるが、エンジニアが適宜パラメータを調整して決定された軽量化手法が最適なものであるとは限らなかった。

【0006】

そこで、本発明の目的の1つは、学習モデルに対する軽量化手法を、より適切にすることを可能にする情報処理方法、プログラム及び情報処理装置を提供する。

【課題を解決するための手段】

【0007】

本発明の一態様に係る情報処理方法は、情報処理装置に含まれる1又は複数のプロセッサが、所定の学習データを取得すること、ニューラルネットワークを用いる所定の学習モデルに対して、蒸留処理された第1学習モデル、枝刈り処理された第2学習モデル、及び量子化処理された第3学習モデルの少なくとも2つのモデルを含む各モデルがそれぞれ重み付けされた重み学習モデルに、所定のデータを入力して機械学習を行うこと、前記各モデルそれぞれの重みを変更された重み学習モデルごとに、前記所定の学習データを入力して前記機械学習が行われた場合の学習結果を取得すること、変更された各重みが重み付けされた各重み学習モデルと、前記各重み学習モデルで学習されたときの各学習結果とを含む学習データを用いて、教師あり学習を行うこと、前記教師あり学習により、任意の学習データを入力する場合に、各重みの集合ごとに学習結果を予測する予測モデルを生成すること、を実行する。

【発明の効果】

【0008】

本発明によれば、学習モデルに対する軽量化手法を、より適切にすることを可能にする

10

20

30

40

50

情報処理方法、プログラム及び情報処理装置を提供することができる。

【図面の簡単な説明】

【0009】

【図1】実施形態に係るシステム構成の一例を示す図である。

【図2】実施形態に係る情報処理装置の物理的構成の一例を示す図である。

【図3】実施形態に係る情報処理装置の処理ブロックの一例を示す図である。

【図4】学習済みモデルの蒸留を説明するための図である。

【図5】学習済みモデルの枝刈りを説明するための図である。

【図6】実施形態に係る情報処理装置の処理ブロックの一例を示す図である。

【図7】実施形態に係る関係情報の一例を示す図である。

10

【図8】実施形態に係る関係情報の表示例を示す図である。

【図9】実施形態に係る予測モデルの生成に関する処理の一例を示すフローチャートである。

【図10】実施形態に係るユーザが利用する情報処理装置における処理の一例を示すフローチャートである。

【発明を実施するための形態】

【0010】

添付図面を参照して、本発明の実施形態について説明する。なお、各図において、同一の符号を付したものは、同一又は同様の構成を有する。

【0011】

20

[実施形態]

<システム構成>

図1は、実施形態に係るシステム構成の一例を示す図である。図1に示す例では、サーバ10と、各情報処理装置20A、20B、20C、20Dとが、ネットワークを介してデータ送受信可能なように接続される。情報処理装置を個別に区別しない場合は情報処理装置20とも表記する。

【0012】

サーバ10は、データを収集、分析可能な情報処理装置であり、1つ又は複数の情報処理装置から構成されてもよい。情報処理装置20は、スマートフォン、パーソナルコンピュータ、タブレット端末、サーバ、コネクティッドカーなど、機械学習を実行可能な情報処理装置である。なお、情報処理装置20は、脳波をセンシングする侵襲型又は非侵襲型の電極に直接的又は間接的に接続されており、脳波データを解析、送受信可能な装置でもよい。

30

【0013】

図1に示すシステムでは、サーバ10は、例えば、所定の学習データにおいて学習済みの学習モデルに、様々な軽量化手法（軽量化アルゴリズム）を適用する。様々な軽量化手法は、既存の軽量化手法1つを適用したり、任意の軽量化手法の組み合わせを適用したりすることを含む。このとき、サーバ10は、所定のデータセット、所定の学習モデル及び所定の軽量化手法のときの学習結果を関連付けて記憶する。

【0014】

40

次に、サーバ10は、任意のデータセットと、任意の軽量化手法と、これらの学習結果（例えば学習精度）とを訓練データとして、学習結果が適切な軽量化手法を特定する予測モデルを学習して生成する。学習結果の適切さは、例えば、学習精度や、モデルサイズの圧縮率などにより決定される。

【0015】

これにより、学習済みの学習モデルに対し、軽量化をより適切に行うことが可能になる。また、サーバ10は、各軽量化手法を重み付けして線形結合したモデルを用いて、各軽量化手法の適用割合を定める各重みを適切に調整してもよい。

【0016】

<ハードウェア構成>

50

図2は、実施形態に係る情報処理装置10の物理的構成の一例を示す図である。情報処理装置10は、演算部に相当する1又は複数のCPU(Central Processing Unit)10aと、記憶部に相当するRAM(Random Access Memory)10bと、記憶部に相当するROM(Read only Memory)10cと、通信部10dと、入力部10eと、表示部10fと、を有する。これらの各構成は、バスを介して相互にデータ送受信可能に接続される。

【0017】

実施形態では、情報処理装置10が一台のコンピュータで構成される場合について説明するが、情報処理装置10は、複数のコンピュータ又は複数の演算部が組み合わされて実現されてもよい。また、図2で示す構成は一例であり、情報処理装置10はこれら以外の構成を有してもよいし、これらの構成のうち一部を有さなくてもよい。

10

【0018】

CPU10aは、RAM10b又はROM10cに記憶されたプログラムの実行に関する制御やデータの演算、加工を行う制御部である。CPU10aは、より適切な軽量化手法を調べるための学習モデルを用いて学習を行うプログラム(学習プログラム)や、任意のデータを入力したときに適切な軽量化手法を出力する予測モデルを生成するための学習を行うプログラム(予測プログラム)を実行する演算部である。CPU10aは、入力部10eや通信部10dから種々のデータを受け取り、データの演算結果を表示部10fに表示したり、RAM10bに格納したりする。

【0019】

RAM10bは、記憶部のうちデータの書き換えが可能なものであり、例えば半導体記憶素子で構成されてよい。RAM10bは、CPU10aが実行するプログラム、様々な軽量化手法に関する軽量化データ(例えば軽量化アルゴリズム)、適切な軽量化手法を予測する予測モデル、学習対象のデータに関する情報と、このデータに対応する適切な軽量化手法との対応関係を示す関係情報などのデータを記憶してもよい。なお、これらは例示であって、RAM10bには、これら以外のデータが記憶されていてもよいし、これらの一部が記憶されていなくてもよい。

20

【0020】

ROM10cは、記憶部のうちデータの読み出しが可能なものであり、例えば半導体記憶素子で構成されてよい。ROM10cは、例えば学習プログラムや、書き換えが行われないデータを記憶してよい。

30

【0021】

通信部10dは、情報処理装置10を他の機器に接続するインターフェースである。通信部10dは、インターネット等の通信ネットワークに接続されてよい。

【0022】

入力部10eは、ユーザからデータの入力を受け付けるものであり、例えば、キーボード及びタッチパネルを含んでよい。

【0023】

表示部10fは、CPU10aによる演算結果を視覚的に表示するものであり、例えば、LCD(Liquid Crystal Display)により構成されてよい。表示部10fが演算結果を表示することは、XAI(eXplainable AI:説明可能なAI)に貢献し得る。表示部10fは、例えば、学習結果や、学習モデルに関する情報を表示してもよい。

40

【0024】

学習プログラムは、RAM10bやROM10c等のコンピュータによって読み取り可能な記憶媒体に記憶されて提供されてもよいし、通信部10dにより接続される通信ネットワークを介して提供されてもよい。情報処理装置10では、CPU10aが学習プログラムを実行することにより、後述する図3を用いて説明する様々な動作が実現される。なお、これらの物理的な構成は例示であって、必ずしも独立した構成でなくてもよい。例えば、情報処理装置10は、CPU10aとRAM10bやROM10cが一体化したLSI(Large-Scale Integration)を備えていてもよい。また、情報処理装置10は、GP

50

U (Graphical Processing Unit) や A S I C (Application Specific Integrated Circuit) を備えていてもよい。

【 0 0 2 5 】

なお、情報処理装置 2 0 の構成は、図 2 に示す情報処理装置 1 0 の構成と同様であるため、その説明を省略する。また、情報処理装置 1 0 と情報処理装置 2 0 とは、データ処理を行う基本的な構成である CPU 1 0 a や RAM 1 0 b 等を有していればよく、入力部 1 0 e や表示部 1 0 f は設けられなくてもよい。また、入力部 1 0 e や表示部 1 0 f は、外部からインターフェースを用いて接続されてもよい。

【 0 0 2 6 】

< 処理構成 >

図 3 は、実施形態に係る情報処理装置 1 0 の処理ブロックの一例を示す図である。情報処理装置 1 0 は、取得部 1 0 1、第 1 学習部 1 0 2、変更部 1 0 3、第 2 学習部 1 0 4、予測部 1 0 5、判定部 1 0 6、設定部 1 0 7、関連付け部 1 0 8、特定部 1 0 9、表示制御部 1 1 0、出力部 1 1 1、及び記憶部 1 1 2 を備える。例えば、図 3 に示す第 1 学習部 1 0 2、変更部 1 0 3、第 2 学習部 1 0 4、予測部 1 0 5、判定部 1 0 6、設定部 1 0 7、関連付け部 1 0 8、特定部 1 0 9、表示制御部 1 1 0 は、例えば CPU 1 0 a などにより実行されて実現され、取得部 1 0 1 及び出力部 1 1 1 は、例えば通信部 1 0 d などにより実現され、記憶部 1 1 2 は、RAM 1 0 b 及び / 又は ROM 1 0 c などにより実現され得る。

【 0 0 2 7 】

取得部 1 0 1 は、所定の学習データを取得する。例えば、取得部 1 0 1 は、所定の学習データとして、画像データ、系列データ、テキストデータなどの公知のデータセットを取得してもよい。なお、取得部 1 0 1 は、記憶部 1 1 2 に記憶されたデータを取得してもよいし、他の情報処理装置により送信されたデータを取得してもよい。

【 0 0 2 8 】

第 1 学習部 1 0 2 は、所定の問題を解くため、ニューラルネットワークを用いる所定の学習モデル 1 0 2 a に対して、蒸留処理された第 1 学習モデル、枝刈り処理された第 2 学習モデル、及び量子化処理された第 3 学習モデルの少なくとも 2 つのモデルを含む各モデルがそれぞれ重み付けされた重み学習モデルに、所定の学習データを入力して機械学習を行う。

【 0 0 2 9 】

ここで、学習済みの学習モデル 1 0 2 a の軽量化手法の例として、蒸留 (Distillation)、枝刈り (Pruning)、及び量子化 (Quantization) の各アルゴリズムについて、以下に簡単に説明する。

【 0 0 3 0 】

図 4 は、学習済みモデルの蒸留を説明するための図である。図 4 に示す蒸留は、学習済みモデル M 1 1 の予測結果を教師データとして、より小さいモデル M 1 2 を学習することで軽量化を行う。このとき、この小さいモデル M 1 2 は、大きいモデル M 1 1 と同程度の精度を持つ場合がある。

【 0 0 3 1 】

例えば、蒸留において、学習済みモデル M 1 1 は、Teacher モデル、小さいモデル M 1 2 は、Student モデルと呼ばれる。Student モデルは、エンジニアが適宜設計する。

【 0 0 3 2 】

図 4 に示す例では、分類器を例とした学習データについて説明する。モデル M 1 1 の Teacher モデルは、0 と 1 とで表現され、1 が正解である教師データを用いて学習を行なう。これに対し、モデル M 1 2 の Student モデルは、Teacher モデルが出力した値 (例 : A = 0 . 7、B = 0 . 3) を教師データとして学習する。実施形態では、1 つの学習モデル M 1 1 に対して複数の異なる蒸留後のモデル M 1 2 が用意されてもよい。

10

20

30

40

50

【 0 0 3 3 】

図 5 は、学習済みモデルの枝刈りを説明するための図である。図 5 に示す枝刈りは、学習済みモデル M 2 1 の重みやノードを削除することで、軽量化が行われたモデル M 2 2 が生成される。これにより計算回数、メモリ使用量の削減を行うことが可能になる。

【 0 0 3 4 】

枝刈りの手法は、ノード間の接続において重みの小さいところを対象に削除が行われてもよい。例えば、枝刈りは、蒸留と違い別途モデルを設計する必要はないが、パラメータの削除が行われるため、再学習を行い、学習精度を維持するとよい。例えば、学習への影響が小さい枝（エッジ）、例えば重みが所定値以下の枝をカットし、軽量化が行われてもよい。

10

【 0 0 3 5 】

量子化は、モデルに含まれるパラメータを少ないビット数で表現する。これにより、ネットワークの構造を変えずにモデルを小さくすることが可能になる。例えば、重みパラメータを 6 個持つ簡単なネットワークを例にした場合、3 2 ビット精度の場合は合計 1 9 2 ビットを必要とするが、8 ビット精度の制約にすると合計 4 8 ビットで表現することになり、軽量化が行われていることになる。

【 0 0 3 6 】

図 3 に戻り、例えば、第 1 学習部 1 0 2 は、学習済みの学習モデル 1 0 2 a に対して、第 1 モデル、第 2 モデル、及び第 3 モデルのうち少なくとも 2 つの軽量化モデルが選択され、各モデルに付与される重みとして、デフォルトの重みを設定する。

20

【 0 0 3 7 】

第 1 モデル、第 2 モデル、及び第 3 モデルは、学習済みのモデルのカテゴリごとに予め設定されていてもよいし、学習済みモデルごとに、所定の基準に従って自動で生成されてもよい。例えば、第 1 学習部 1 0 2 は、蒸留の場合、学習済みモデルに適した蒸留後のモデルを機械学習により決定してもよく、枝刈りの場合、重みが所定値以下の枝をカットして枝刈り後のモデルを生成してもよく、量子化の場合、所定値ビット精度の制約（量子化）にしてもよい。また、1 つの学習済みモデルに対し、複数の第 1 モデル、複数の第 2 モデル、複数の第 3 モデルが設定され、それぞれのモデルに重みが付与されてもよい。

【 0 0 3 8 】

所定の問題は、例えば画像データ、系列データ及びテキストデータの少なくともいずれかについて、分類、生成及び最適化の少なくともいずれかを行う問題を含む。ここで、画像データは、静止画のデータと、動画のデータとを含む。系列データは、音声データや株価のデータを含む。

30

【 0 0 3 9 】

また、所定の学習モデル 1 0 2 a は、ニューラルネットワークを含む学習済みの学習モデルであり、例えば、画像認識モデル、系列データ解析モデル、ロボットの制御モデル、強化学習モデル、音声認識モデル、音声生成モデル、画像生成モデル、自然言語処理モデル等の少なくとも 1 つを含む。また、具体例としては、所定の学習モデル 1 0 2 a は、C N N (Convolutional Neural Network)、R N N (Recurrent Neural Network)、D N N (Deep Neural Network)、L S T M (Long Short-Term Memory)、双方向 L S T M、D Q N (Deep Q-Network)、V A E (Variational AutoEncoder)、G A N s (Generative Adversarial Networks)、f l o w - b a s e d 生成モデル等のいずれかでもよい。

40

【 0 0 4 0 】

変更部 1 0 3 は、所定の学習データ及び/又は重み学習モデルの各重みを変更する。例えば、変更部 1 0 3 は、複数の学習データの中から、第 1 学習部 1 2 に入力される所定の学習データを 1 つずつ順に変更する。また、変更部 1 0 3 は、ある 1 つの重み学習モデルに対して全ての所定の学習データが入力されて学習が行われた場合、重み学習モデルの別の各重みを利用するため、複数の各重みの集合（セット）の中から 1 つのセットを選択し、用意された全てのセットで学習を行い、学習結果を取得してもよい。

50

【 0 0 4 1 】

また、第 1 学習部 1 0 2 は、所定の学習データを重み学習モデルに入力し、適切な学習結果が出力されるように、重み学習モデルのハイパーパラメータ等の学習を行う。このとき、第 1 学習部 1 0 2 は、ハイパーパラメータが更新（調整）される際に、重み学習モデルの各モデルに付与される各重みも所定の方法により調整する。

【 0 0 4 2 】

例えば、各重みの調整については、あらかじめ設定される初期値から逐次的に各重みが調整されるとよい。このとき、各重みが全て加算して 1 になるように調整され、以前に行った調整と異なる調整が行われればいずれの調整方法が用いられてもよい。例えば、第 1 学習部 1 0 2 は、各重みを順に所定値ずつ変更していき、全ての組み合わせについて変更する。例えば、第 1 学習部 1 0 2 は、重み w_k に対して初期値から所定値ずつ減算し、重み w_{k+1} に対して初期値から所定値ずつ加算し、どちらかの重みが 0 以下になると、 k に 1 を加算して、各初期値からの変更を繰り返す。また、各重みが全て加算して 1 になる条件は設けなくてもよく、この場合、 Softmax 関数などを用いて、各重みを加算して 1 になるように最後に調整されればよい。

10

【 0 0 4 3 】

これにより、所定の学習データと所定の各重みのセットとの任意の組み合わせに対して学習させることが可能になる。例えば、変更部 1 0 3 は、所定の学習データと所定の各重みのセットとの全ての組み合わせが学習されるように、所定の学習データ及び / 又は所定の各重みのセットを 1 つずつ順に変更してもよいし、所定の条件が満たされるまで所定の学習データ及び / 又は所定の各重みのセットを 1 つずつ順に変更してもよい。所定の条件は、例えば、学習精度やモデルサイズの圧縮率などにより設定されてもよい。

20

【 0 0 4 4 】

取得部 1 0 1 又は第 1 学習部 1 0 2 は、各モデルそれぞれの重みが変更された重み学習モデルごとに、所定の学習データを入力して機械学習が行われた場合の学習結果を取得する。例えば、取得部 1 0 1 又は第 1 学習部 1 0 2 は、様々な組み合わせの所定の学習データ及び / 又は所定の各重みのセットを用いて学習された学習結果を取得する。

【 0 0 4 5 】

ここで、重み学習モデルについて具体例を用いて説明する。例えば、第 1 学習部 1 0 2 は、第 1 モデル、第 2 モデル、第 3 モデルにそれぞれ重み w_1 , w_2 , w_3 を付与して線形結合した重み学習モデルを利用してもよい。この場合の重み学習の関数 $M(x)$ として、式 (1) が挙げられるが、一例にすぎない。

30

$$M_1(x) = w_1 m_1(x) + w_2 m_2(x) + w_3 m_3(x) \quad \dots \text{式 (1)}$$

w_n : 重み (各重みの集合 (セット) を W と表記する)

$m_n(x)$: 第 n モデル

x : 学習データ

【 0 0 4 6 】

変更部 1 0 3 は、例えば、 $w_1 + w_2 + w_3 = 1$ となるように、各重みを所定基準に従って、1 つずつ順に変更する。第 1 学習部 1 0 2 は、変更後の各重みに対する学習結果を取得し、各重みのセットに対して学習結果を関連付けておく。学習結果は、学習精度と、軽量化による効果を示すモデルサイズの圧縮率である。モデルサイズの圧縮率とは、例えば、軽量化後の学習済みモデルのパラメータ数の、軽量化前の学習済みモデルのパラメータ数に対する割合である。

40

【 0 0 4 7 】

また、変更部 1 0 3 が所定の学習データを変更すると、第 1 学習部 1 0 2 は、変更後の学習データに対して、上述されたように各重みのセットでの重み学習モデルを学習し、学習結果を取得する。これにより、任意の学習データ、任意の各重みのセット、これらの場合の学習結果を含む訓練データが生成される。

【 0 0 4 8 】

第 2 学習部 1 0 4 は、変更された各重みが重み付けされた各重み学習モデルと、各重み

50

学習モデルで学習されたときの各学習結果とを含む学習データを用いて、教師あり学習を行う。例えば、第2学習部104は、任意の学習データ及び任意の各重みのセットを用いて学習された際の学習結果（例えば学習性能及び/又はモデルサイズの圧縮率）を正解ラベルとする訓練データを用いて、教師あり学習を行う。

【0049】

また、第2学習部104は、教師あり学習により、任意の学習データを入力する場合に、各重みのセットごとに学習結果を予測する予測モデル104aを生成する。例えば、第2学習部104は、任意の学習データを入力すると、この学習データに対する各軽量化手法の各重みのセットごとに、学習精度やモデルサイズの圧縮率を出力する予測モデルを生成する。

10

【0050】

以上の構成により、様々な学習データや、様々な軽量化手法により軽量化した各学習モデルを用いた学習結果を訓練データとして教師あり学習を行うことにより、各重みのセットごとに、学習結果を予測する予測モデルを生成することができる。その結果、第2学習部104により生成された予測モデルを用いることで、軽量化手法をより適切にすることが可能になる。

【0051】

予測部105は、任意の学習データを予測モデル104aに入力し、各モデルそれぞれの重みのセットごとに、重み学習モデルを実行した場合の学習結果を予測する。例えば、予測部105は、学習データとして画像のデータセットを入力した場合、特定の重みセット W_n (w_{1n} , w_{2n} , w_{3n}) ごとに、学習精度とモデルサイズに関する値（例えば圧縮率）とを予測する。

20

【0052】

これにより、任意のデータ（例、データセット）に対して、各軽量化手法をどれくらい適用するかの各重みのセットごとに、学習結果が予測されるため、この学習結果に基づいて、より適切な各重みを選択することなどが可能になる。

【0053】

判定部106は、任意の学習データを所定の学習モデル102aに入力した場合の学習結果と、予測モデル104aにより予測された学習結果とが、軽量化に関する所定条件を満たすか否かを判定する。例えば、判定部106は、軽量化前の学習済みの学習モデル102aに学習データAを入力したときの学習精度A1と、予測モデル104aにより予測された学習精度B1との第1差分値が、第1閾値内であるか否かを判定する。この第1差分値が小さければと小さいほど、学習モデルを軽量化後でも学習精度を維持することができており、学習精度B1のときの各重みは適切な軽量化手法となる。

30

【0054】

また、判定部106は、軽量化前の学習済みの学習モデル102aの圧縮率A2（=1）と、予測モデル104aにより予測された圧縮率B2との第2差分値が、第2閾値以上であるか否かを判定する。この第2差分値が大きければ大きいほど、学習モデルがより軽量化できていることを示す。

【0055】

判定部106は、軽量化に関する判定結果に基づいて、各重みの有効性を判定する。例えば、判定部106は、第1差分値と第2差分値とに基づいて、圧縮率B2が大きく、学習精度B1が軽量化前の精度を維持できている各重みに対し、有効な軽量化手法であると判定する。具体例としては、判定部106は、第1差分値が第1閾値以下、第2差分値が第2閾値以上の各重みを有効な軽量化手法、それ以外の各重みを有効ではない軽量化手法と判定してもよい。

40

【0056】

これにより、モデルサイズに関する値（例えば圧縮率）と学習精度とに基づき、それぞれの予測値を参考に、適切な各重みを選定することができる。例えば、判定部106は、最も学習精度が高い各重みを選定してもよいし、圧縮率が第2閾値以上のもので、学習精

50

度が最も高い各重みを選定してもよい。

【0057】

設定部107は、軽量化に関する所定条件に関するユーザ操作を受け付ける。例えば、設定部107は、表示部10fに表示された条件入力画面から、ユーザが入力部10eを操作して軽量化に関する所定条件を入力した場合、この入力操作を受け付ける。

【0058】

設定部107は、受け付けたユーザ操作に基づいて、軽量化に関する所定条件を判定部106の判定条件に設定する。例えば、設定部107は、ユーザの入力操作に基づいて、学習性能に関する第1閾値、及び/又は、モデルサイズに関する第2閾値を設定することを可能としてもよい。

10

【0059】

これにより、ユーザが所望する条件を用いて、有効な軽量化手法を特定することができるようになる。

【0060】

関連付け部108は、学習結果に含まれる学習精度を第1変数、学習結果に含まれるモデルサイズに関する値(例えば圧縮率)を第2変数とし、第1変数及び第2変数と、各重みとを対応付ける関係情報を生成する。例えば、関連付け部108は、縦軸を第1変数、横軸を第2変数とする場合に、それぞれの変数の交点に各重みWを対応付けたマトリックスを生成してもよい。また、関連付け部108は、各情報処理装置20から取得された学習精度や圧縮率に基づいて、第1変数及び第2変数と、各重みWとを対応付ける関係情報(実測関係情報)を生成してもよい。

20

【0061】

以上の処理により、第1変数又は第2変数に変更された場合に、対応する各重みWを迅速に特定することが可能になる。また、第1変数と第2変数とは、適宜変更されてもよい。例えば、第1変数として学習精度、第2変数として各重みWを適用し、特定される情報モデルサイズに関する値でもよい。

【0062】

また、取得部101は、第1変数の第1値及び第2変数の第2値を取得してもよい。例えば、取得部101は、ユーザから指定される第1変数の第1値及び第2変数の第2値を取得する。第1値又は第2値はユーザにより適宜指定される。

30

【0063】

この場合、特定部109は、関連付け部108により生成された関係情報に基づいて、第1変数の第1値及び第2変数の第2値に対応する各重みWを特定する。例えば、特定部109は、関係情報を用いて、変更される第1変数の値、又は第2変数の値に対応する各重みWを特定する。

【0064】

表示制御部110は、特定部109により特定された各重みWを表示装置(表示部10f)に表示制御する。また、表示制御部110は、第1変数及び第2変数を変更可能にしたマトリックスをGUI(Graphical User Interface)で表してもよい(例えば、後述する図8等)。

40

【0065】

以上の処理により、ユーザにより指定された第1変数又は第2変数に応じて特定される各重みWを、ユーザに対して可視化することが可能になる。ユーザは、第1変数又は第2変数を変更することで、所望の各重みWを特定し、学習済みモデルの軽量化に適用することができる。

【0066】

出力部111は、第2学習部104により予測された各重みWを、他の情報処理装置20に出力してもよい。例えば、出力部111は、所定の学習データを送信した情報処理装置20であって、適切な各重みWの取得を要求した情報処理装置20に対し、所定の学習データに対応する適切な各重みWを出力してもよい。また、出力部111は、予測された

50

各重みWを記憶部112に出力してもよい。

【0067】

記憶部112は、学習に関するデータを記憶する。記憶部112は、所定のデータセット112aや、軽量化手法112bに関するデータ、上述した関係情報112c、訓練データ、学習途中のデータ、学習結果に関する情報などを記憶する。

【0068】

図6は、実施形態に係る情報処理装置20の処理ブロックの一例を示す図である。情報処理装置20は、取得部201、学習部202、出力部203、及び記憶部204を備える。情報処理装置20は、汎用のコンピュータで構成されてもよい。

【0069】

取得部201は、他の情報処理装置（例えばサーバ10）により、分散学習の指示とともに、所定の重み学習モデルに関する情報や所定のデータセットに関する情報を取得してもよい。所定の重み学習モデルに関する情報は、各重みを示す情報や、重み学習モデル自体を示す情報でもよい。所定のデータセットに関する情報は、データセット自体でもよく、所定のデータセットが格納された格納先を示す情報でもよい。

【0070】

学習部202は、所定の重み学習モデル202aに学習対象の所定のデータセットを入力して学習を行う。学習部202は、学習後の学習結果をサーバ10にフィードバックするように制御する。学習結果は、例えば、学習性能などを含み、モデルサイズに関する情報をさらに含んでもよい。学習部202は、学習対象のデータセットの種類、及び/又は、解くべき問題に応じて、学習モデル202aを選択してもよい。

【0071】

また、所定の重み学習モデル202aは、ニューラルネットワークを含む学習モデルであり、例えば、画像認識モデル、系列データ解析モデル、ロボットの制御モデル、強化学習モデル、音声認識モデル、音声生成モデル、画像生成モデル、自然言語処理モデル等の少なくとも1つをベースに、各軽量化手法が重み付けされたモデルを含む。また、具体例としては、所定の重み学習モデル202aのベースは、CNN（Convolutional Neural Network）、RNN（Recurrent Neural Network）、DNN（Deep Neural Network）、LSTM（Long Short-Term Memory）、双方向LSTM、DQN（Deep Q-Network）、VAE（Variational AutoEncoder）、GANs（Generative Adversarial Networks）、flow-based生成モデル等のいずれかでもよい。

【0072】

出力部203は、分散学習の学習結果に関する情報を他の情報処理装置に出力する。例えば、出力部203は、学習部202による学習結果に関する情報をサーバ10に出力する。例えば、分散学習の学習結果に関する情報は、上述したように、学習性能を含み、モデルサイズに関する情報をさらに含んでもよい。

【0073】

記憶部204は、学習部202に関するデータを記憶する。記憶部204は、所定のデータセット204aや、サーバ10から取得したデータ、学習途中のデータ、学習結果に関する情報などを記憶する。

【0074】

これにより、情報処理装置20は、他の情報処理装置（例えばサーバ10）からの指示により、所定のデータセットに対して、所定の重み学習モデルを適用した分散学習を実行し、学習結果をサーバ10にフィードバックすることが可能になる。

【0075】

また、出力部203は、所定のデータに関する情報を他の情報処理装置（例えばサーバ10）に出力する。出力部203は、所定のデータ（例えば学習対象のデータセット）を出力してもよいし、所定のデータの特徴情報を出力してもよい。

【0076】

取得部201は、他の情報処理装置から、所定のデータに対応する各重みWを取得して

10

20

30

40

50

もよい。取得される各重み W は、他の情報処理装置が予測モデルを利用して予測した、所定のデータに適切な各重みである。

【0077】

学習部202は、取得された各重みを重み学習モデル202aに適用する。このとき、重み学習モデル202aは、上述した学習に用いられた重み学習モデル22aに各重みを適用してもよい。また、重み学習モデル202aは、他の情報処理装置10から取得される学習モデルでもよいし、自装置で管理する学習モデルでもよい。

【0078】

学習部202は、各重みが適用された重み学習モデル202aに、所定のデータを入力して学習結果を取得する。この学習結果は、所定のデータに適した各重みを用いて学習した結果である。学習部202は、学習性能を保ちつつ適切に軽量化された学習モデルを使用することができる。

<データ例>

図7は、実施形態に係る関係情報の一例を示す図である。図7に示す例では、関係情報は、各第1変数(例、 P_{11})及び各第2変数(例、 P_{21})に対応する各重み(例、 W_1)を含む。第1変数 P_{1n} は、例えば学習精度であり、第2変数 P_{2n} は例えばモデルサイズの圧縮率であり、変数としてはいずれかの変数だけでもよい。各重み $W(P_{1n}, P_{2m})$ は、第1変数 P_{1n} 及び第2変数 P_{2n} の場合の重みである。

【0079】

図7に示す関係情報について、サーバ10は、所定の分散インスタンス数とハイパーパラメータの組み合わせで分散学習を行わせた情報処理装置20から、又は、自装置の教師あり学習の結果から、学習精度(第1変数)と、圧縮率(第2変数)とを取得する。サーバ10は、取得された学習精度と圧縮率に、各重み W を対応付ける。サーバ10は、教師あり学習により実測された学習精度と圧縮率とを取得するたびにを行うことで、図7に示す関係情報を生成することが可能になる。また、関係情報は、予測部105により予測された結果に基づいて、任意のデータセットに対する予測関係情報が生成されてもよい。

【0080】

<ユーザインタフェースの例>

図8は、実施形態に係る関係情報の表示例を示す図である。図8に示す例では、関係情報に含まれる第1変数と第2変数とをスライダーを用いて変更可能にする。ユーザが第1変数又は第2変数に対してスライダーを用いて移動させることで、例えば、移動後の第1変数(P_{1n})又は第2変数(P_{2m})に対応する各重み W のセット $W(P_{1n}, P_{2m})$ が、対応する点に関連付けて表示される。

【0081】

また、ユーザは、第1変数及び第2変数の二次元のグラフ上に所定の点を指定することで、指定された点に対応する学習精度と、圧縮率との組み合わせが表示されるようにしてもよい。

【0082】

これにより、サーバ10は、第1変数と第2変数との組み合わせに対応する、適切な各重み W を表示可能になる。また、視覚的に対応関係をユーザに示しながら、これから分散学習が行われる任意のデータセットに対して適切な分散インスタンス数やハイパーパラメータを選択させるユーザインタフェースを提供することが可能になる。

【0083】

<動作>

図9は、実施形態に係る予測モデルの生成に関する処理の一例を示すフローチャートである。図9に示す処理は、情報処理装置10により実行される。

【0084】

ステップS102において、情報処理装置10の取得部101は、所定の学習データを取得する。所定の学習データは、記憶部112のデータセット112aから選択されてもよいし、他の装置からネットワークを介して受信された所定のデータでもよいし、ユーザ

10

20

30

40

50

操作に応じて入力された所定のデータを取得してもよい。

【0085】

ステップS104において、情報処理装置10の第1学習部102は、ニューラルネットワークを用いる所定の学習モデルに対して、蒸留処理された第1学習モデル、枝刈り処理された第2学習モデル、及び量子化処理された第3学習モデルの少なくとも2つのモデルを含む各モデルがそれぞれ重み付けされた重み学習モデルに、所定のデータを入力して機械学習を行う。

【0086】

ステップS106において、情報処理装置10の第2学習部104は、各モデルそれぞれの重みを変更された重み学習モデルごとに、所定の学習データを入力して機械学習が行われた場合の学習結果を取得する。

10

【0087】

ステップS108において、情報処理装置10の第2学習部104は、変更された各重みが重み付けされた各重み学習モデルと、各重み学習モデルで学習されたときの各学習結果とを含む学習データを用いて、教師あり学習を行う。

【0088】

ステップS110において、情報処理装置10の第2学習部104は、教師あり学習により、任意の学習データを入力する場合に、各重みの組み合わせごとに学習結果を予測する予測モデルを生成する。

【0089】

以上の処理により、生成された予測モデルを利用することで、ニューラルネットワークを用いる学習済みモデルを、学習精度を保ちつつ、より適切に軽量化を行うことを可能にする。

20

【0090】

図10は、実施形態に係るユーザが利用する情報処理装置20における処理の一例を示すフローチャートである。ステップS202において、情報処理装置20の出力部203は、学習対象の所定の学習データに関する情報を他の情報処理装置(例えばサーバ10)に出力する。

【0091】

ステップS204において、情報処理装置20の取得部201は、他の情報処理装置(例えばサーバ10)から、所定の学習データに対応する各重みを示す情報を取得する。

30

【0092】

ステップS206において、情報処理装置20の学習部202は、取得された各重みを所定の重み学習モデル202aに適用する。

【0093】

ステップS208において、情報処理装置20の学習部202は、各重みが適用された学習モデル202aに、所定の学習データを入力して学習結果を取得する。

【0094】

これにより、エッジ側の情報処理装置であっても、学習対象のデータに対して、適切な軽量化を行った学習モデルを用いて学習を行うことで、学習精度を保つことができる。

40

【0095】

以上説明した実施形態は、本発明の理解を容易にするためのものであり、本発明を限定して解釈するためのものではない。実施形態が備える各要素並びにその配置、材料、条件、形状及びサイズ等は、例示したものに限定されるわけではなく適宜変更することができる。また、第1学習部102を備える装置と、第2学習部104を備える装置とは別のコンピュータでもよい。この場合、生成された第1学習部102により学習された学習結果が、ネットワークを介して、第2学習部104を備える装置に送信されてもよい。

【0096】

また、情報処理装置10は、変更部103を必ずしも設けなくてもよい。例えば、情報処理装置10は、任意の学習対象のデータと任意の重みのセットとの組の各学習性能を取

50

得して第2学習部104による学習を行ってもよい。

【符号の説明】

【0097】

10...情報処理装置、10a...CPU、10b...RAM、10c...ROM、10d...通信部、10e...入力部、10f...表示部、101...取得部、102...第1学習部、102a...学習モデル、103...変更部、104...第2学習部、104a...予測モデル、105...予測部、106...判定部、107...設定部、108...関連付け部、109...特定部、110...表示制御部、111...出力部、112...記憶部、112a...データセット、112b...軽量化手法、112c...関係情報、201...取得部、202...学習部、202a...学習モデル、203...出力部、204...記憶部、204a...データセット

10

20

30

40

50

【要約】

【課題】学習モデルに対する軽量化手法を、より適切にする。

【解決手段】情報処理方法は、情報処理装置に含まれる1又は複数のプロセッサが、所定の学習データを取得すること、ニューラルネットワークを用いる所定の学習モデルに対して、各軽量化モデルの少なくとも2つのモデルを含む各モデルがそれぞれ重み付けされた重み学習モデルに、所定のデータを入力して機械学習を行うこと、重み学習モデルごとに、所定の学習データを入力して機械学習が行われた場合の学習結果を取得すること、各重み学習モデルと、各重み学習モデルで学習されたときの各学習結果とを含む学習データを用いて、教師あり学習を行うこと、教師あり学習により、任意の学習データを入力する場合に、各重みの集合ごとに学習結果を予測する予測モデルを生成すること、を実行する。

10

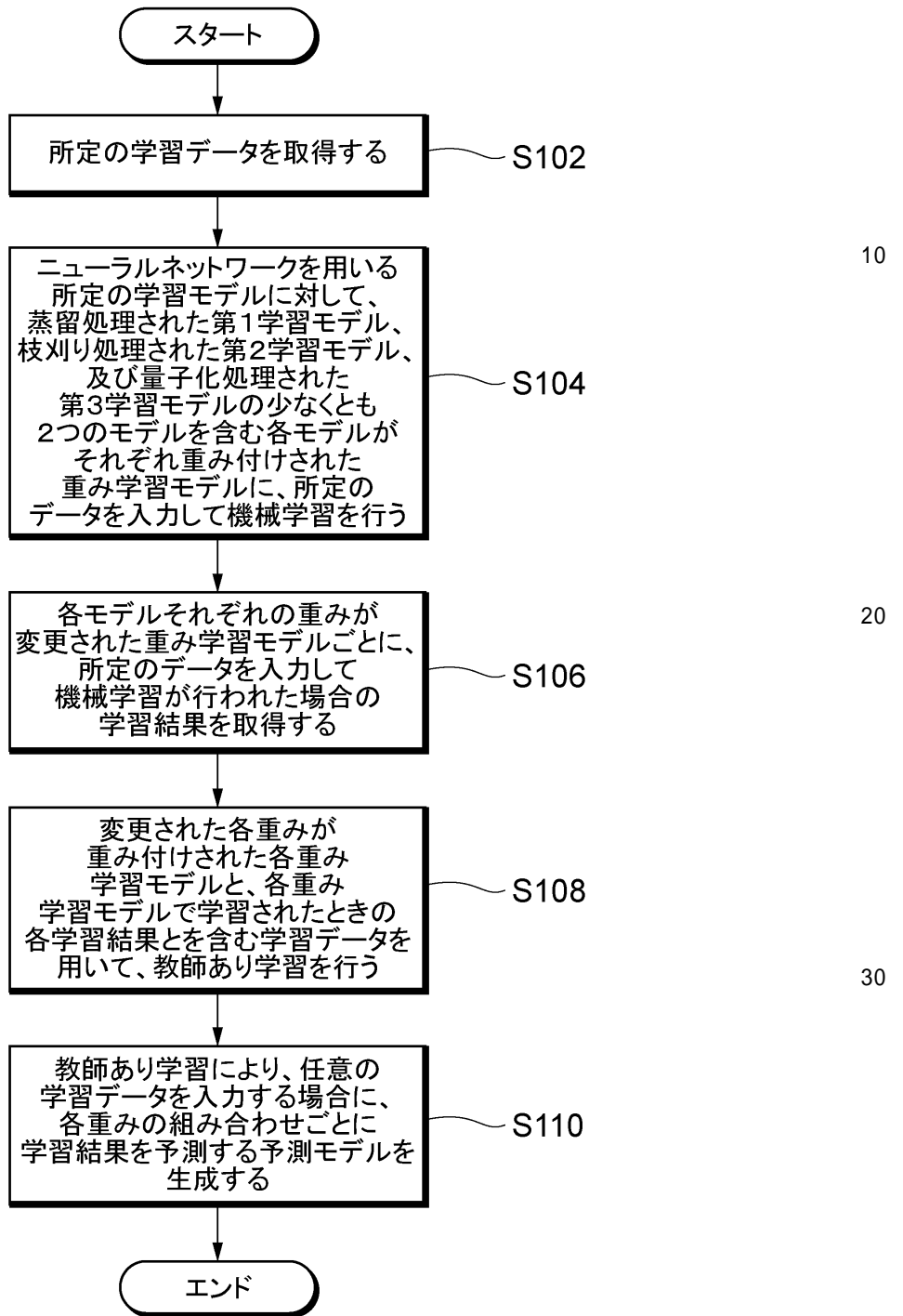
【選択図】図9

20

30

40

50



10

20

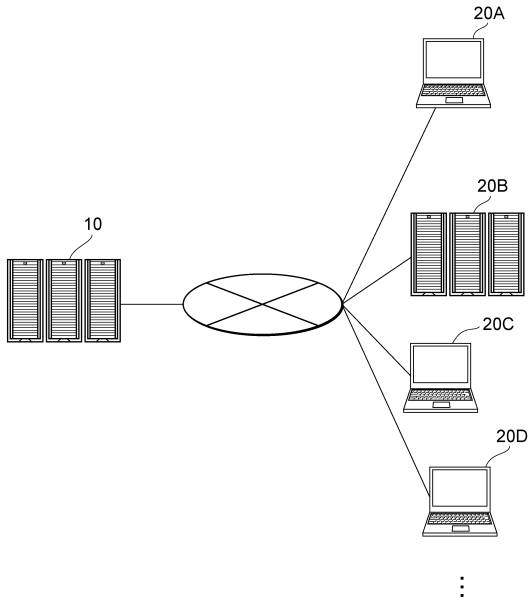
30

40

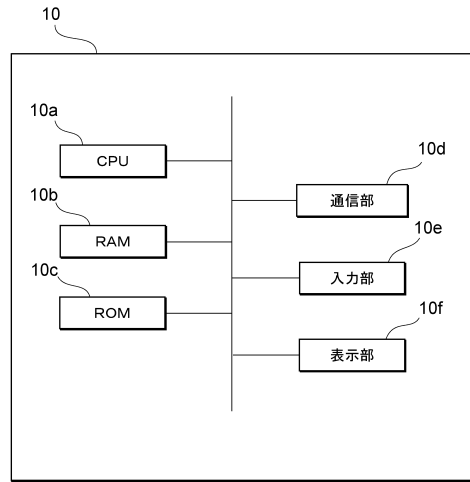
50

【図面】

【図 1】



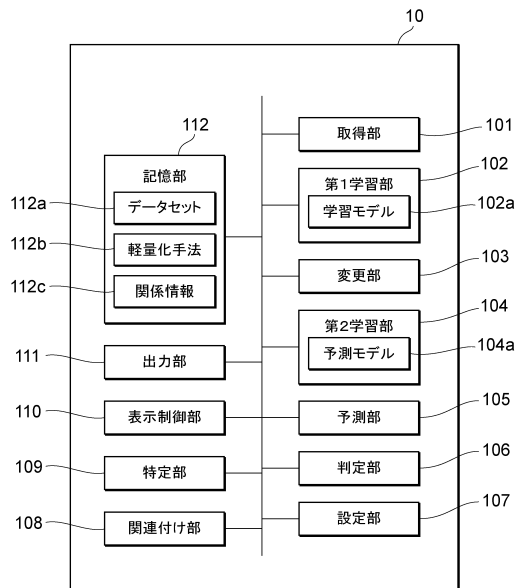
【図 2】



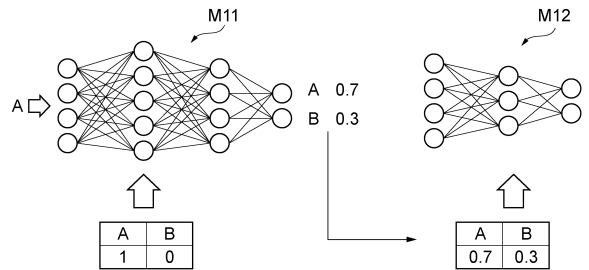
10

20

【図 3】



【図 4】

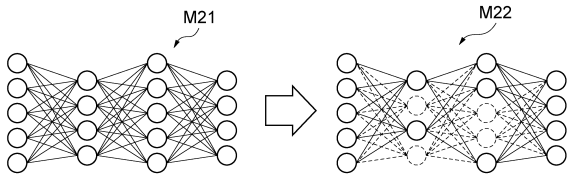


30

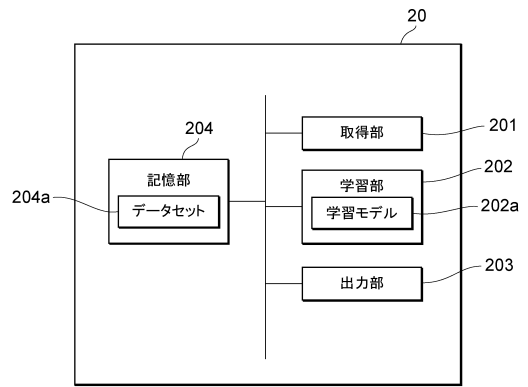
40

50

【図 5】



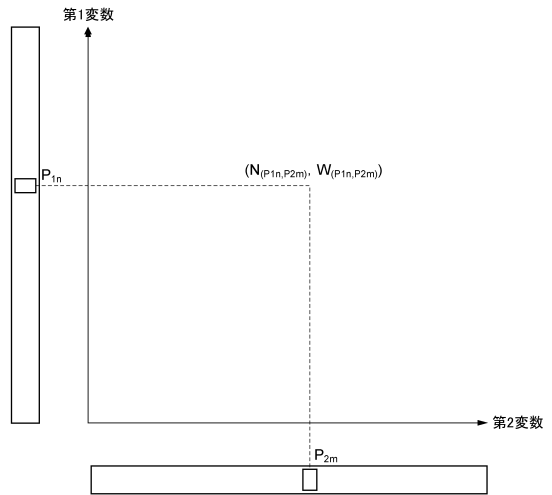
【図 6】



【図 7】

関係情報		
第1変数	第2変数	重み
P_{11}	P_{21}	$W_1(w_1, w_2, w_3)$
...

【図 8】



10

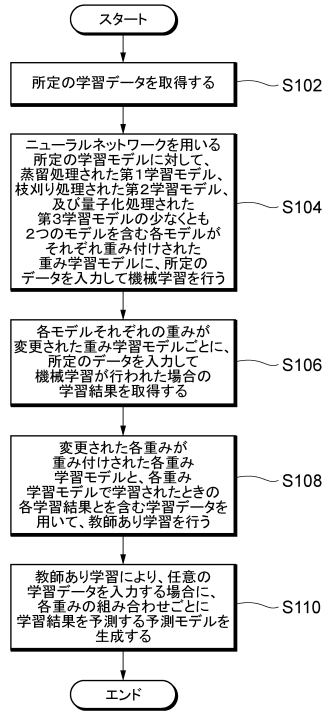
20

30

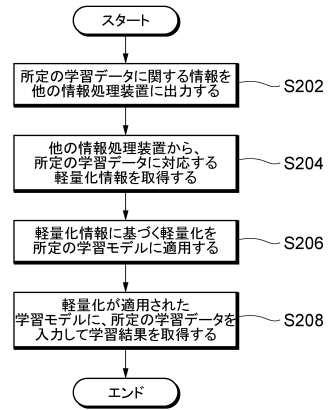
40

50

【 図 9 】



【 図 10 】



10

20

30

40

50

フロントページの続き

DATA 品川ビル13階 株式会社クリエイターズネクスト内

審査官 中村 信也

- (56)参考文献 国際公開第2022/023022(WO, A1)
米国特許出願公開第2020/0125956(US, A1)
米国特許出願公開第2020/0311552(US, A1)
- (58)調査した分野 (Int.Cl., DB名)
- | | |
|------|--------------|
| G06N | 3/00 - 3/12 |
| G06N | 7/08 - 99/00 |
| G06N | 5/00 - 7/06 |