

(19)日本国特許庁(JP)

(12)公開特許公報(A)

(11)公開番号

特開2023-44336

(P2023-44336A)

(43)公開日 令和5年3月30日(2023.3.30)

(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 3/08 (2023.01)

G 0 6 N 3/08

審査請求 未請求 請求項の数 14 O L (全19頁)

(21)出願番号 特願2021-152315(P2021-152315)

(22)出願日 令和3年9月17日(2021.9.17)

(71)出願人 000000295

沖電気工業株式会社

東京都港区虎ノ門一丁目7番12号

(74)代理人 100140958

弁理士 伊藤 学

(74)代理人 100137888

弁理士 大山 夏子

(74)代理人 100190942

弁理士 風間 竜司

(72)発明者 国定 恭史

東京都港区虎ノ門一丁目7番12号 沖

電気工業株式会社内

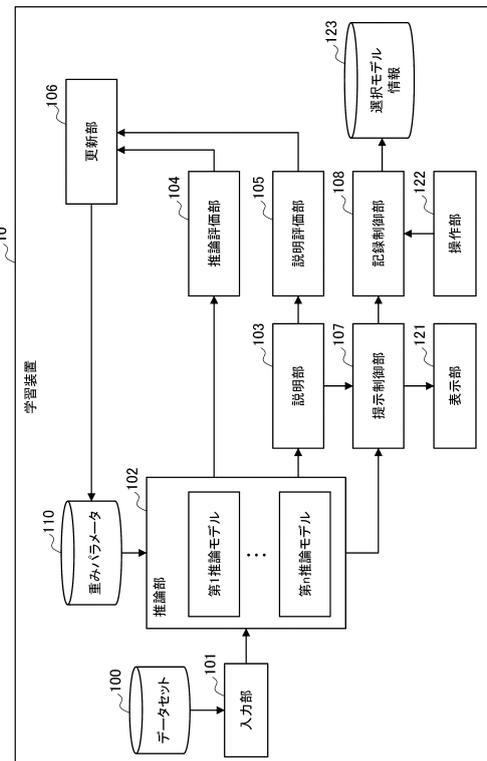
(54)【発明の名称】 学習装置、学習方法およびプログラム

(57)【要約】

【課題】人的コストを抑制しながら、解釈性および精度の高いモデルを得ることが可能な技術が提供されることが望まれる。

【解決手段】第1の入力データと前記第1の入力データの正解値とを取得する入力部と、前記第1の入力データと複数の推論モデルとに基づいて、複数の推論モデルそれぞれに対応する第1の推論値を出力する推論部と、前記第1の推論値に対する前記第1の入力データの寄与の大きさを示す前記複数の推論モデルそれぞれに対応する第1の説明情報を出力する説明部と、前記正解値と前記第1の推論値とに基づいて推論評価結果を得る推論評価部と、前記複数の推論モデルそれぞれに対応する第1の説明情報同士の一貫度に基づいて説明評価結果を得る説明評価部と、前記推論評価結果と前記説明評価結果とに基づいて、前記複数の推論モデルの第1の重みパラメータの更新を行う更新部と、を備える、学習装置が提供される。

【選択図】図1



【特許請求の範囲】

【請求項 1】

第 1 の入力データと前記第 1 の入力データの正解値とを取得する入力部と、
 前記第 1 の入力データと複数の推論モデルとに基づいて、複数の推論モデルそれぞれに対応する第 1 の推論値を出力する推論部と、
 前記第 1 の推論値に対する前記第 1 の入力データの寄与の大きさを示す前記複数の推論モデルそれぞれに対応する第 1 の説明情報を出力する説明部と、
 前記正解値と前記第 1 の推論値とに基づいて推論評価結果を得る推論評価部と、
 前記複数の推論モデルそれぞれに対応する第 1 の説明情報同士の一貫性に基づいて説明評価結果を得る説明評価部と、
 前記推論評価結果と前記説明評価結果とに基づいて、前記複数の推論モデルの第 1 の重みパラメータの更新を行う更新部と、
 を備える、学習装置。

10

【請求項 2】

前記入力部は、第 2 の入力データを取得し、
 前記推論部は、前記第 2 の入力データと前記第 1 の重みパラメータの更新後の複数の推定モデルである複数の学習済みモデルとに基づいて、前記複数の学習済みモデルそれぞれに対応する第 2 の推論値を出力し、
 前記説明部は、前記第 2 の推論値に対する前記第 2 の入力データの寄与の大きさを示す前記複数の学習済みモデルそれぞれに対応する第 2 の説明情報を出力し、
 前記学習装置は、前記第 2 の説明情報のユーザへの提示を制御する提示制御部を備える、
 請求項 1 に記載の学習装置。

20

【請求項 3】

前記提示制御部は、前記第 2 の推論値および前記第 2 の説明情報の前記ユーザへの提示を制御する、
 請求項 2 に記載の学習装置。

【請求項 4】

前記学習装置は、
 前記複数の学習済みモデルから前記ユーザによって選択された 1 または複数の学習済みモデルを示す情報の記録を制御する記録制御部を備える、
 請求項 2 または 3 に記載の学習装置。

30

【請求項 5】

前記説明評価結果は、前記複数の推論モデルそれぞれに対応する第 1 の説明情報同士の一貫性が大きいほど小さい値を取る、
 請求項 1 ~ 4 のいずれか一項に記載の学習装置。

【請求項 6】

前記説明評価部は、前記複数の推論モデルそれぞれに対応する第 1 の説明情報を正規化したベクトルの内積に基づいて前記説明評価結果を得る、
 請求項 5 に記載の学習装置。

40

【請求項 7】

前記説明評価部は、前記複数の推論モデルごとに、前記第 1 の説明情報の二値化を行ってマスクを生成するとともに、自身以外の推論モデルに対応する前記第 1 の説明情報から生成したマスクと自身の推論モデルに対応する前記第 1 の説明情報との積を計算し、前記複数の推論モデルごとの前記積の和に基づいて、前記説明評価結果を得る、
 請求項 5 に記載の学習装置。

【請求項 8】

前記説明部は、誤差逆伝播が可能な関数を含む、
 請求項 1 ~ 7 のいずれか一項に記載の学習装置。

【請求項 9】

50

前記説明部は、第 2 の重みパラメータを有し、
前記更新部は、誤差逆伝播法によって前記第 2 の重みパラメータの更新を行う、
請求項 8 に記載の学習装置。

【請求項 10】

前記複数の推論モデルの少なくとも一つは、ニューラルネットワークを含む、
請求項 1 ~ 9 のいずれか一項に記載の学習装置。

【請求項 11】

前記更新部は、前記推論評価結果と前記説明評価結果との加算結果に基づいて、前記第 1 の重みパラメータの更新を行う、
請求項 1 ~ 10 のいずれか一項に記載の学習装置。

10

【請求項 12】

前記第 1 の説明情報は、前記第 1 の推論値に対する前記第 1 の入力データの寄与の大きさを示すヒートマップである、
請求項 1 ~ 11 のいずれか一項に記載の学習装置。

【請求項 13】

第 1 の入力データと前記第 1 の入力データの正解値とを取得することと、
前記第 1 の入力データと複数の推論モデルとに基づいて、複数の推論モデルそれぞれに対応する第 1 の推論値を出力することと、
前記第 1 の推論値に対する前記第 1 の入力データの寄与の大きさを示す前記複数の推論モデルそれぞれに対応する第 1 の説明情報を出力することと、
前記正解値と前記第 1 の推論値とに基づいて推論評価結果を得ることと、
前記複数の推論モデルそれぞれに対応する第 1 の説明情報同士の一貫性に基づいて説明評価結果を得ることと、
前記推論評価結果と前記説明評価結果とに基づいて、前記複数の推論モデルの第 1 の重みパラメータの更新を行うことと、
を含む、学習方法。

20

【請求項 14】

コンピュータを、
第 1 の入力データと前記第 1 の入力データの正解値とを取得する入力部と、
前記第 1 の入力データと複数の推論モデルとに基づいて、複数の推論モデルそれぞれに対応する第 1 の推論値を出力する推論部と、
前記第 1 の推論値に対する前記第 1 の入力データの寄与の大きさを示す前記複数の推論モデルそれぞれに対応する第 1 の説明情報を出力する説明部と、
前記正解値と前記第 1 の推論値とに基づいて推論評価結果を得る推論評価部と、
前記複数の推論モデルそれぞれに対応する第 1 の説明情報同士の一貫性に基づいて説明評価結果を得る説明評価部と、
前記推論評価結果と前記説明評価結果とに基づいて、前記複数の推論モデルの第 1 の重みパラメータの更新を行う更新部と、
を備える、学習装置として機能させるプログラム。

30

【発明の詳細な説明】

40

【技術分野】

【0001】

本発明は、学習装置、学習方法およびプログラムに関する。

【背景技術】

【0002】

ニューラルネットワーク（以下、「NN」とも表記する。）は、画像認識などにおいて高い性能を有する。しかし、一般的に NN は、膨大なパラメータと複雑なモデルとによって構成されており、NN のパラメータと NN からの出力結果との関係を解釈することが難しい。かかる課題を解決するため、解釈性の高い NN を得る手法が幾つか提案されている。なお、「解釈性が高い」は、「人間の感覚との一致度が高い」とも換言され得る。

50

【 0 0 0 3 】

例えば、NNのモデルが判断のために注目すべき領域を示したヒートマップのラベルを手によって付しておき、そのヒートマップと一致するようにモデルを学習させることによって人にも解釈しやすいモデルを得る手法が知られている（例えば、非特許文献1参照）。また、モデルから得られたヒートマップの解釈性が低い場合には、そのヒートマップと一致しないようにモデルを再学習させることによって、より解釈性の高いモデルを得ることもできる。

【 0 0 0 4 】

また、入力データのうちNNが判断を行うための注目領域を抽出する機構をネットワーク内に導入することによって、NNの精度を向上させる手法も知られている（例えば、非特許文献2参照）。かかる手法によって得られた注目領域を人間が修正し、修正した注目領域とNNの注目領域が一致するようにNNを再学習させることによって、NNの解釈性および精度を向上させることができる。

【 先行技術文献 】

【 非特許文献 】

【 0 0 0 5 】

【非特許文献1】Andrew Ross、他2名、"Right for the Right Reasons: Training Differentiable Models by Constraining their Explanations"、[online]、[令和3年9月8日検索]、インターネット<<https://arxiv.org/abs/1703.03717>>

【非特許文献2】Masahiro Mitsuhashi、他6名、"Embedding Human Knowledge into Deep Neural Network via Attention Map"、[online]、[令和3年9月8日検索]、インターネット<<https://arxiv.org/abs/1905.03540>>

【非特許文献3】"Grad-CAM: Visual Explanations from Deep Networks via Gradient-based Localization"、[online]、[令和3年9月8日検索]、インターネット<<https://arxiv.org/abs/1610.02391v3>>

【 発明の概要 】

【 発明が解決しようとする課題 】

【 0 0 0 6 】

しかしながら、非特許文献1および非特許文献2に記載された、人手によってヒートマップのラベルを用意する手法は、ラベル付けのための人的コストが大きい。

【 0 0 0 7 】

一方、ラベル付けを必要としない手法としては、非特許文献1に記載の学習済みモデルのヒートマップと一致しないようにモデルを再学習させる手法が挙げられる。しかし、かかる手法では、再学習により精度が低下してしまう可能性が高いという点が課題として挙げられる。さらに、かかる手法では、全てのデータに対して一様にヒートマップの一致度が低下してしまうため、個々のデータに対してはかえってヒートマップの解釈性を低下させてしまう場合があるという点が課題として挙げられる。

【 0 0 0 8 】

そこで、人的コストを抑制しながら、解釈性および精度の高いモデルを得ることが可能な技術が提供されることが望まれる。

【 課題を解決するための手段 】

【 0 0 0 9 】

上記問題を解決するために、本発明のある観点によれば、第1の入力データと前記第1の入力データの正解値とを取得する入力部と、前記第1の入力データと複数の推論モデルとに基づいて、複数の推論モデルそれぞれに対応する第1の推論値を出力する推論部と、前記第1の推論値に対する前記第1の入力データの寄与の大きさを示す前記複数の推論モデルそれぞれに対応する第1の説明情報を出力する説明部と、前記正解値と前記第1の推論値とに基づいて推論評価結果を得る推論評価部と、前記複数の推論モデルそれぞれに対応する第1の説明情報同士の一貫性に基づいて説明評価結果を得る説明評価部と、前記推論評価結果と前記説明評価結果とに基づいて、前記複数の推論モデルの第1の重みパラメ

10

20

30

40

50

ータの更新を行う更新部と、を備える、学習装置が提供される。

【0010】

前記入力部は、第2の入力データを取得し、前記推論部は、前記第2の入力データと前記第1の重みパラメータの更新後の複数の推定モデルである複数の学習済みモデルとに基づいて、前記複数の学習済みモデルそれぞれに対応する第2の推論値を出力し、前記説明部は、前記第2の推論値に対する前記第2の入力データの寄与の大きさを示す前記複数の学習済みモデルそれぞれに対応する第2の説明情報を出力し、前記学習装置は、前記第2の説明情報のユーザへの提示を制御する提示制御部を備えてもよい。

【0011】

前記提示制御部は、前記第2の推論値および前記第2の説明情報の前記ユーザへの提示を制御してもよい。 10

【0012】

前記学習装置は、前記複数の学習済みモデルから前記ユーザによって選択された1または複数の学習済みモデルを示す情報の記録を制御する記録制御部を備えてもよい。

【0013】

前記説明評価結果は、前記複数の推論モデルそれぞれに対応する第1の説明情報同士の一貫性が大きいほど小さい値を取ってもよい。

【0014】

前記説明評価部は、前記複数の推論モデルそれぞれに対応する第1の説明情報を正規化したベクトルの内積に基づいて前記説明評価結果を得てもよい。 20

【0015】

前記説明評価部は、前記複数の推論モデルごとに、前記第1の説明情報の二値化を行ってマスクを生成するとともに、自身以外の推論モデルに対応する前記第1の説明情報から生成したマスクと自身の推論モデルに対応する前記第1の説明情報との積を計算し、前記複数の推論モデルごとの前記積の和に基づいて、前記説明評価結果を得てもよい。

【0016】

前記説明部は、誤差逆伝播が可能な関数を含んでもよい。

【0017】

前記説明部は、第2の重みパラメータを有し、前記更新部は、誤差逆伝播法によって前記第2の重みパラメータの更新を行ってもよい。 30

【0018】

前記複数の推論モデルの少なくとも一つは、ニューラルネットワークを含んでもよい。なお、ニューラルネットワークは、機械学習アルゴリズムの一例に過ぎない。したがって、ニューラルネットワークの代わりに他の機械学習アルゴリズムが用いられてもよい。

【0019】

前記更新部は、前記推論評価結果と前記説明評価結果との加算結果に基づいて、前記第1の重みパラメータの更新を行ってもよい。

【0020】

前記第1の説明情報は、前記第1の推論値に対する前記第1の入力データの寄与の大きさを示すヒートマップであってもよい。 40

【0021】

また、本発明の別の観点によれば、第1の入力データと前記第1の入力データの正解値とを取得することと、前記第1の入力データと複数の推論モデルとに基づいて、複数の推論モデルそれぞれに対応する第1の推論値を出力することと、前記第1の推論値に対する前記第1の入力データの寄与の大きさを示す前記複数の推論モデルそれぞれに対応する第1の説明情報を出力することと、前記正解値と前記第1の推論値とに基づいて推論評価結果を得ることと、前記複数の推論モデルそれぞれに対応する第1の説明情報同士の一貫性に基づいて説明評価結果を得ることと、前記推論評価結果と前記説明評価結果とに基づいて、前記複数の推論モデルの第1の重みパラメータの更新を行うことと、を含む、学習方法が提供される。 50

【 0 0 2 2 】

また、本発明の別の観点によれば、コンピュータを、第 1 の入力データと前記第 1 の入力データの正解値とを取得する入力部と、前記第 1 の入力データと複数の推論モデルとに基づいて、複数の推論モデルそれぞれに対応する第 1 の推論値を出力する推論部と、前記第 1 の推論値に対する前記第 1 の入力データの寄与の大きさを示す前記複数の推論モデルそれぞれに対応する第 1 の説明情報を出力する説明部と、前記正解値と前記第 1 の推論値とに基づいて推論評価結果を得る推論評価部と、前記複数の推論モデルそれぞれに対応する第 1 の説明情報同士の一貫性に基づいて説明評価結果を得る説明評価部と、前記推論評価結果と前記説明評価結果とに基づいて、前記複数の推論モデルの第 1 の重みパラメータの更新を行う更新部と、を備える、学習装置として機能させるプログラムが提供される。

10

【 発明の効果 】

【 0 0 2 3 】

以上説明したように本発明によれば、人的コストを抑制しながら、解釈性および精度の高いモデルを得ることが可能な技術が提供される。

【 図面の簡単な説明 】

【 0 0 2 4 】

【 図 1 】本発明の実施形態に係る学習装置の機能構成例を示す図である。

【 図 2 】ヒートマップを二値化したマスクと他のヒートマップとの掛け合わせによって説明評価結果を得る手法について説明するための図である。

【 図 3 】同実施形態に係る学習装置の学習段階における動作例を示すフローチャートである。

20

【 図 4 】同実施形態に係る学習装置のテスト段階における動作例を示すフローチャートである。

【 図 5 】学習装置の例としての情報処理装置のハードウェア構成を示す図である。

【 発明を実施するための形態 】

【 0 0 2 5 】

以下に添付図面を参照しながら、本発明の好適な実施の形態について詳細に説明する。なお、本明細書及び図面において、実質的に同一の機能構成を有する構成要素については、同一の符号を付することにより重複説明を省略する。

【 0 0 2 6 】

30

また、本明細書および図面において、実質的に同一の機能構成を有する複数の構成要素を、同一の符号の後に異なる数字を付して区別する場合がある。ただし、実質的に同一の機能構成を有する複数の構成要素等の各々を特に区別する必要がない場合、同一符号のみを付する。また、異なる実施形態の類似する構成要素については、同一の符号の後に異なるアルファベットを付して区別する場合がある。ただし、異なる実施形態の類似する構成要素等の各々を特に区別する必要がない場合、同一符号のみを付する。

【 0 0 2 7 】

(0 . 実施形態の概要)

本発明の実施形態の概要について説明する。本発明の実施形態では、入力データ（学習用データ）と正解値との組み合わせに基づいてニューラルネットワークの学習を行う学習装置について説明する。しかし、ニューラルネットワークは、機械学習アルゴリズムの一例に過ぎない。したがって、ニューラルネットワークの代わりに他の機械学習アルゴリズムが用いられてもよい。例えば、機械学習アルゴリズムの他の一例として、SVM (Support Vector Machine) などが用いられてもよい。

40

【 0 0 2 8 】

(1 . 実施形態の詳細)

本発明の実施形態について詳細に説明する。

【 0 0 2 9 】

(1 . 1 . 学習装置の構成例)

図 1 は、本発明の実施形態に係る学習装置 10 の機能構成例を示す図である。図 1 に示

50

されるように、本発明の実施形態に係る学習装置 10 は、入力部 101 と、推論部 102 と、説明部 103 と、推論評価部 104 と、説明評価部 105 と、更新部 106 と、提示制御部 107 と、記録制御部 108 と、表示部 121 と、操作部 122 とを備える。

【0030】

本発明の実施形態では、推論部 102 が、 n 個 (n は 1 より大きい整数) の推論モデル、すなわち、「第 1 推論モデル」から「第 n 推論モデル」までを含む場合を主に想定する。また、本発明の実施形態では、第 1 推論モデルから第 n 推論モデルまでのそれぞれが、ニューラルネットワークを含んで構成される場合を主に想定する。以下では、ニューラルネットワークを「NN」とも表記する。

【0031】

第 1 推論モデルから第 n 推論モデルまでのそれぞれに含まれる NN は、重みパラメータ 110 (第 1 の重みパラメータ) を使用する。このとき、第 1 推論モデルから第 n 推論モデルまでのそれぞれに含まれる NN は、共通の構造を有し、使用する重みパラメータ 110 (第 1 の重みパラメータ) が異なってもよい。あるいは、第 1 推論モデルから第 n 推論モデルまでのそれぞれに含まれる NN は、別々の構造を有していてもよい。

【0032】

なお、第 1 推論モデルから第 n 推論モデルまでの少なくとも一つが、NN を含んでもよい。例えば、第 1 推論モデルから第 n 推論モデルまでの一部が NN を含んでもよく、第 1 推論モデルから第 n 推論モデルまでの他の一部は、NN の代わりに他の機械学習アルゴリズムを含んでもよい。

【0033】

さらに、本発明の実施形態では、説明部 103 が NN を含んで構成される場合を主に想定する。説明部 103 に含まれる NN は、重みパラメータ (第 2 の重みパラメータ) を使用する。

【0034】

データセット 100、第 1 推論モデルから第 n 推論モデルまでの重みパラメータ 110 (第 1 の重みパラメータ) および説明部 103 が有する重みパラメータ (第 2 の重みパラメータ) は、図示しない記憶部によって記憶される。かかる記憶部は、RAM (Random Access Memory)、ハードディスクドライブまたはフラッシュメモリなどのメモリによって構成されてよい。

【0035】

入力部 101 と、推論部 102 と、説明部 103 と、推論評価部 104 と、説明評価部 105 と、更新部 106 と、提示制御部 107 と、記録制御部 108 とは、CPU (Central Processing Unit) または GPU (Graphics Processing Unit) などの演算装置を含み、ROM (Read Only Memory) により記憶されているプログラムが演算装置により RAM に展開されて実行されることにより、その機能が実現され得る。このとき、当該プログラムを記録した、コンピュータに読み取り可能な記録媒体も提供され得る。あるいは、これらのブロックは、専用のハードウェアにより構成されていてもよいし、複数のハードウェアの組み合わせにより構成されてもよい。演算装置による演算に必要なデータは、図示しない記憶部によって適宜記憶される。

【0036】

初期状態において、第 1 推論モデルから第 n 推論モデルまでの重みパラメータ 110 および説明部 103 が有する重みパラメータそれぞれには、初期値が設定されている。例えば、これらに設定される初期値は、ランダムな値であってよいが、どのような値であっててもよい。例えば、これらに設定される初期値は、あらかじめ学習によって得られた学習済みの値であっててもよい。

【0037】

(データセット 100)

データセット 100 は、学習段階において使用される複数の入力データ (第 1 の入力デ

10

20

30

40

50

ータ)と当該複数の入力データそれぞれの正解値とを含む。学習段階において使用される複数の入力データは、学習用データに該当し得る。さらに、データセット100は、テスト段階において使用される複数の入力データ(第2の入力データ)を含む。テスト段階において使用される複数の入力データは、テスト用データに該当し得る。

【0038】

なお、テスト用データは、学習用データと別のデータとして用意されていることが主に想定される。しかし、テスト用データは、学習用データの一部を含んでもよい。

【0039】

また、本発明の実施形態では、入力データが画像データである場合(特に、静止画像データである場合)を主に想定する。しかし、入力データの種類は特に限定されず、画像データ以外も入力データとして用いられ得る。例えば、入力データは、複数のフレームを含んだ動画データであってもよいし、音響データであってもよい。

【0040】

(入力部101)

入力部101は、学習段階において、データセット100から学習段階において使用される入力データおよび正解値の組み合わせを順次に取得する。入力部101は、学習段階において使用される入力データおよび正解値の組み合わせを順次に推論部102に出力する。また、入力部101は、テスト段階において、データセット100からテストにおいて使用される入力データを順次に取得する。入力部101は、テスト段階において使用される入力データを順次に推論部102に出力する。

【0041】

なお、例えば、入力部101は、データセット100から学習段階において使用される入力データおよび正解値の組み合わせを全部取得して出力し終わった場合には、最初から当該組み合わせを取得し直して再度出力する動作を所定の回数繰り返してよい。かかる場合には、入力部101よりも後段のブロックにおいても、再度の入力に基づいて順次に各自の処理が繰り返し実行されてよい。一方、例えば、入力部101は、データセット100からテスト段階において使用される入力データを全部取得して出力し終わった場合には、入力データの取得を終了してよい。

【0042】

(推論部102)

推論部102は、学習段階において、入力部101から入力された入力データと第1推論モデルから第n推論モデルまでとに基づいて、第1推論モデルから第n推論モデルまでのそれぞれに対応する推論値(第1の推論値)を得る。同様に、推論部102は、テスト段階において、入力部101から入力された入力データと第1推論モデルから第n推論モデルまでとに基づいて、第1推論モデルから第n推論モデルまでのそれぞれに対応する推論値(第2の推論値)を得る。

【0043】

第1推論モデルから第n推論モデルまでが使用する重みパラメータ110は図示しない記憶部によって記憶されている。したがって、推論部102は、図示しない記憶部から重みパラメータ110を取得し、取得した重みパラメータ110と入力部101から入力された入力データとに基づいて、第1推論モデルから第n推論モデルまでによる推論を行う。

【0044】

なお、本明細書においては、NNへの入力に基づいてNNからの出力を得ることを広く「推論」と言う。

【0045】

一例として、i番目の推論モデルを示す関数を F_i (iは1~nまでの整数)とし、i番目の推論モデルへの入力をxとすると、i番目の推論モデルからの出力は $F_i(x)$ と表現され得る。

【0046】

10

20

30

40

50

なお、後にも説明するように、説明部 103 が用いる説明手法（すなわち、説明情報の生成手法）には、推論値の他に第 1 推論モデルから第 n 推論モデルまでのそれぞれから出力される特徴量（中間特徴量）などの情報を必要とする説明手法が存在する場合があります。かかる場合には、推論部 102 は、推論値とともに、第 1 推論モデルから第 n 推論モデルまでのそれぞれの間層から出力される特徴量を説明部 103 に出力してよい。

【0047】

第 1 推論モデルから第 n 推論モデルまでの具体的な構成は、特に限定されない。しかし、第 1 推論モデルから第 n 推論モデルまでのそれぞれの出力の形式は、入力データに対応する正解値の形式と合わせて設定されているのがよい。例えば、正解値が分類問題のクラスである場合、第 1 推論モデルから第 n 推論モデルまでのそれぞれの出力は、クラス数分の長さを持つ one-hot ベクトルであるとよい。

10

【0048】

推論部 102 は、学習段階において、第 1 推論モデルから第 n 推論モデルまでのそれぞれに対応する推論値を、説明部 103 および推論評価部 104 それぞれに出力する。一方、推論部 102 は、テスト段階において、第 1 推論モデルから第 n 推論モデルまでのそれぞれに対応する推論値を、説明部 103 および提示制御部 107 それぞれに出力する。

【0049】

（説明部 103）

説明部 103 は、第 1 推論モデルから第 n 推論モデルまでのそれぞれについて、推論部 102 から入力された推論値の判断根拠を説明する説明情報を生成する。

20

【0050】

ここで、説明情報は、推論部 102 から入力された推論値に対する入力データの寄与の大きさを示す情報である。以下では、説明情報が推論値に対する入力データの寄与の大きさを領域（例えば、画像を構成するピクセルなど）または変数ごとに示すヒートマップである場合について主に説明する。ヒートマップによれば、入力データのうち判断に寄与した重要な領域または変数が示され得る。

【0051】

入力データが画像データなどである場合には、ヒートマップは 2 次元ベクトルによって表現され得る。あるいは、入力データが表形式データなどである場合には、ヒートマップは 1 次元ベクトルによって表現され得る。

30

【0052】

ヒートマップはどのように生成されてもよい。例えば、説明部 103 は、推論部 102 から入力された推論値に基づいて、ヒートマップを生成してもよい。あるいは、上記したように、推論部 102 から説明部 103 に推論値だけでなく特徴量も入力される場合があります。かかる場合には、説明部 103 は、推論部 102 から入力された推論値と特徴量とに基づいて、ヒートマップを生成してもよい。

【0053】

例えば、説明部 103 は、誤差逆伝播が可能な関数を含んでいてもよい。このとき、後に説明するように、更新部 106 によって説明部 103 が有する重みパラメータが誤差逆伝播法によって更新され得る。すなわち、説明部 103 は、誤差逆伝播法による更新後の重みパラメータによってヒートマップを生成してもよい。

40

【0054】

誤差逆伝播法による更新後の重みパラメータによってヒートマップを生成する説明手法としては、非特許文献 3 に記載された、いわゆる Grad-CAM などが適用され得る。Grad-CAM は、NN への入力のうち推論値への寄与度が高い領域を示すヒートマップを出力する説明手法である。その他にも、Vanilla Gradient、Smooth Grad といった各種の説明手法が適用され得る。

【0055】

上記したように、i 番目の推論モデルに対応する推論値は $F_i(x)$ と表現され得るため、一例として、ヒートマップの生成処理を示す関数を G とすると、説明部 103 によ

50

て生成される i 番目の推論モデルに対応するヒートマップ $T_i(x)$ は、以下の式 (1) のように表現され得る。

【0056】

$$T_i(x) = G(F_i(x)) \cdots (1)$$

【0057】

説明部 103 は、学習段階において、生成した n 個のヒートマップ (第 1 の説明情報) を説明評価部 105 に出力する。一方、説明部 103 は、テスト段階において、生成した n 個のヒートマップ (第 2 の説明情報) を提示制御部 107 に出力する。

【0058】

(推論評価部 104)

推論評価部 104 は、推論部 102 から入力された第 1 推論モデルから第 n 推論モデルまでのそれぞれに対応する推論値と入力部 101 によって取得された正解値とに基づいて、推論評価結果を得る。より詳細に、推論評価部 104 は、第 1 推論モデルから第 n 推論モデルまでのそれぞれに対応する推論値と入力部 101 によって取得された正解値とを比較することによって、推論評価結果を得る。

【0059】

本発明の実施形態では、推論評価部 104 が、推論値と正解値とに応じた損失関数の第 1 推論モデルから第 n 推論モデルまでについての和を推論評価結果の例としての損失関数 L_1 として算出する場合を想定する。ここで、推論値と正解値とに応じた損失関数は特定の関数に限定されず、一般的なニューラルネットワークにおいて用いられる損失関数と同様の損失関数が用いられてよい。例えば、推論値と正解値とに応じた損失関数は、正解値と推論値との差分に基づくクロスエントロピー誤差であってもよい。

【0060】

推論評価部 104 は、推論評価結果を更新部 106 に出力する。

【0061】

(説明評価部 105)

説明評価部 105 は、説明部 103 から入力された第 1 推論モデルから第 n 推論モデルまでのそれぞれに対応するヒートマップに基づいて説明評価結果を得る。より詳細に、説明評価部 105 は、第 1 推論モデルから第 n 推論モデルまでのそれぞれに対応するヒートマップ同士を比較する。そして、説明評価部 105 は、比較結果としての n 個のヒートマップ同士の一致度に基づいて、説明評価結果を得る。

【0062】

本発明の実施形態では、 n 個のヒートマップ同士の一致度が大きいほど説明評価結果が小さい値を取る損失関数である場合を主に想定する。なお、 n 個のヒートマップ同士の一致度は、 n 個のヒートマップ同士がどの程度乖離しているかを示す乖離度と換言されてもよい。かかる場合には、 n 個のヒートマップ同士の乖離度が小さいほど説明評価結果が小さい値を取る損失関数であってよい。

【0063】

n 個のヒートマップから説明評価結果を得る手法は限定されない。ここでは、説明評価結果を得る手法として、ヒートマップを二値化したマスクと他のヒートマップとの掛け合わせによって説明評価結果を得る手法、および、正規化されたヒートマップ同士の内積によって説明評価結果を得る手法について順に説明する。

【0064】

図 2 は、ヒートマップを二値化したマスクと他のヒートマップとの掛け合わせによって説明評価結果を得る手法について説明するための図である。図 2 に示された例では、説明を簡便にするため、 $n = 2$ である場合、すなわち、推論部 102 が、第 1 推論モデルおよび第 2 推論モデルを有する場合を想定する。

【0065】

図 2 を参照すると、第 1 推論モデルからは、推論値とヒートマップ H_1 とが出力されている。一方、第 2 推論モデルからは、推論値とヒートマップ H_2 とが出力されている。図

10

20

30

40

50

2では、ヒートマップH1およびヒートマップH2において、入力データのうち推論値への寄与が大きい領域ほど濃い色によって示されている。

【0066】

説明評価部105は、ヒートマップH1の二値化を行ってマスクM1を生成するとともに、ヒートマップH2の二値化を行ってマスクM2を生成する。なお、二値化は、閾値c以上である要素（例えば、ヒートマップを構成するピクセル）の値を1とし、閾値cよりも小さい要素の値を0とすることによって実行され得る。図2においては、二値のうち1が黒によって示され、0が白によって示されている。

【0067】

説明評価部105は、第1推論モデルから出力されたヒートマップH1と、第2の推論モデルから出力されたヒートマップH2から生成したマスクM2との積を、要素ごとに計算する。同様に、説明評価部105は、第2推論モデルから出力されたヒートマップH2と、第1の推論モデルから出力されたヒートマップH1から生成したマスクM1との積を、要素ごとに計算する。これによって、各要素に対応する積の集合が推論モデルごとに得られる。

10

【0068】

説明評価部105は、各要素に対応する積を全部の推論モデルについて足し合わせることによって積の和を計算する。そして、説明評価部105は、このようにして計算した積の和を全要素について足し合わせることによって合計値を計算する。説明評価部105は、この合計値を説明評価結果の例としての損失関数L2とする。

20

【0069】

図2を参照しながらn=2である場合について説明した。nを1より大きい任意の整数であるとして説明すると、以下の通りである。

【0070】

すなわち、説明評価部105は、 $i = 1 \sim n$ について、ヒートマップ $T_i(x)$ の各要素の値を二値化したマスク $M_i(x)$ を生成する。次に、説明評価部105は、推論モデルごとに、自身の推論モデルから出力されたヒートマップ $T_i(x)$ と、自身以外の推論モデルに対応するヒートマップから生成したマスク $M_1(x) \sim M_{i-1}(x)$ 、 $M_{i+1}(x) \sim M_n(x)$ の和との積を要素ごとに計算する。

【0071】

説明評価部105は、各要素に対応する積を第1推論モデルから第n推論モデルまでについて足し合わせることによって積の和を計算する。そして、説明評価部105は、このようにして計算した積の和に基づいて、説明評価結果を得る。より詳細に、説明評価部105は、積の和を全要素について足し合わせることによって合計値を計算する。説明評価部105は、この合計値を説明評価結果の例としての損失関数L2とする。

30

【0072】

この損失関数L2は、各ヒートマップにおいて、自身以外のヒートマップにおいて閾値以上の値を持つ領域の合計値である。この損失関数L2の値を小さくするように学習が行われることによって、ヒートマップの一致度が小さいn個の推論モデルが得られる。なお、このときの損失関数L2は、以下の式(2)のように表現され得る。式(2)において、eは、要素番号を示す。ここで、ヒートマップ $T_i(x)$ は、ヒートマップ $T_i(x)$ の大きさ $|T_i(x)|$ で割るなどして正規化してもよい。また、ヒートマップ $T_i(x)$ にはsigmoidなどの活性化関数をかけてもよい。

40

【0073】

【数1】

$$L2(x) = \sum_e \left(\sum_{i=1}^n \left(T_i(x) \cdot \sum_{j=1, j \neq i}^n M_j(x) \right) \right) \dots (2)$$

50

【 0 0 7 4 】

図 2 を参照しながら、ヒートマップを二値化したマスクと他のヒートマップとの掛け合わせによって説明評価結果を得る手法について説明した。続いて、正規化されたヒートマップ同士の内積によって説明評価結果を得る手法について説明する。

【 0 0 7 5 】

説明評価部 1 0 5 は、 $i = 1 \sim n$ について、ヒートマップ $T_i(x)$ をヒートマップ $T_i(x)$ の大きさ $|T_i(x)|$ で割ることによって正規化して、 $i = 1 \sim n$ についての正規化したベクトルを生成する。そして、説明評価部 1 0 5 は、 $i = 1 \sim n$ についての正規化したベクトルの内積に基づいて説明評価結果を得る。より詳細に、説明評価部 1 0 5 は、内積を全要素について足し合わせることによって合計値を計算する。説明評価部 1 0 5 は、この合計値を説明評価結果の例としての損失関数 L_2 とする。

10

【 0 0 7 6 】

正規化したベクトルの内積が大きいほど、この損失関数 L_2 は、大きい値となる。正規化したベクトルの内積が大きいことは、ヒートマップ同士の一致度が高いことを意味する。したがって、この損失関数 L_2 の値を小さくするように学習が行われることによって、ヒートマップの一致度が小さい n 個の推論モデルが得られる。なお、このときの損失関数 L_2 は、以下の式 (3) のように表現され得る。式 (3) において、 e は、要素番号を示す。

【 0 0 7 7 】

【 数 2 】

20

$$L2(x) = \prod_{i=1}^n \frac{Ti(x)}{|Ti(x)|} \quad \dots (3)$$

【 0 0 7 8 】

説明評価部 1 0 5 は、説明評価結果を更新部 1 0 6 に出力する。

【 0 0 7 9 】

(更新部 1 0 6)

更新部 1 0 6 は、推論評価部 1 0 4 から入力された推論評価結果と、説明評価部 1 0 5 から入力された説明評価結果とに基づいて、第 1 推論モデルから第 n 推論モデルまでのそれぞれが使用する重みパラメータ 1 1 0 の更新を行う。これによって、第 1 推論モデルから第 n 推論モデルまでのそれぞれから出力される推論値が正解値に近づくように、かつ、説明部 1 0 3 から出力される n 個のヒートマップ同士の一致度が小さくなるように、重みパラメータ 1 1 0 が更新され得る。重みパラメータ 1 1 0 は、誤差逆伝播法 (バックプロパゲーション) によって更新されてよい。

30

【 0 0 8 0 】

例えば、更新部 1 0 6 は、推論評価部 1 0 4 から入力された推論評価結果と、説明評価部 1 0 5 から入力された説明評価結果とを加算し、加算結果に基づいて、重みパラメータ 1 1 0 の更新を行えばよい。このとき、更新部 1 0 6 は、計算した加算結果を誤差として、誤差逆伝播法 (バックプロパゲーション) によって重みパラメータ 1 1 0 を更新すればよい。上記のように、推論評価結果が損失関数 L_1 と表現され、説明評価結果が損失関数 L_2 と表現される場合、加算結果は、 $L_1 + L_2$ である。

40

【 0 0 8 1 】

さらに、更新部 1 0 6 は、説明部 1 0 3 が有する重みパラメータを更新してよい。より詳細に、説明部 1 0 3 が、誤差逆伝播が可能な関数を含む場合、更新部 1 0 6 は、推論評価結果と説明評価結果とに基づいて、誤差逆伝播法 (バックプロパゲーション) によって、説明部 1 0 3 が有する重みパラメータを更新してよい。

【 0 0 8 2 】

50

なお、学習の終了条件（すなわち、重みパラメータ更新の終了条件）は特に限定されず、第1推論モデルから第n推論モデルまでの学習がある程度行われたことを示す条件であればよい。具体的に、学習の終了条件は、損失関数 $L_1 + L_2$ の値が閾値よりも小さいという条件を含んでもよい。あるいは、学習の終了条件は、損失関数 $L_1 + L_2$ の値の変化が閾値よりも小さいという条件（損失関数 $L_1 + L_2$ の値が収束状態になったという条件）を含んでもよい。あるいは、学習の終了条件は、重みパラメータの更新が所定の回数行われたという条件を含んでもよい。あるいは、推論評価部104によって正解値と推論値とに基づいて精度（例えば、正答率など）が算出される場合、学習の終了条件は、精度が所定の割合（例えば、90%など）を超えるという条件を含んでもよい。

【0083】

（提示制御部107）

提示制御部107は、テスト段階において、推論部102から入力された第1推論モデルから第n推論モデルまでのそれぞれに対応する推論値と、説明部103から入力されたn個のヒートマップとが、ユーザに提示されるように制御する。より詳細に、提示制御部107は、第1推論モデルから第n推論モデルまでのそれぞれに対応する推論値と、n個のヒートマップとが表示されるように表示部121を制御する。なお、n個のヒートマップは表示されるが、第1推論モデルから第n推論モデルまでのそれぞれに対応する推論値は表示されない形態も想定され得る。

【0084】

（表示部121）

表示部121は、ディスプレイによって構成され、提示制御部107による制御に従って各種情報の表示を行う機能を有する。例えば、表示部121は、n個の推論値とn個のヒートマップとを表示することが可能である。ここで、表示部121の形態は特に限定されない。例えば、表示部121は、液晶ディスプレイ（LCD）装置であってもよいし、OLED（Organic Light Emitting Diode）装置であってもよいし、ランプなどの表示装置であってもよい。

【0085】

（操作部122）

操作部122は、ユーザによる操作を受け付ける。例えば、ユーザがn個の推論値とn個のヒートマップとを参照しながら、n個の推論モデルから解釈性の高い1または複数の推論モデル（以下、「選択モデル」とも言う。）を見つけたとする。このとき、ユーザは、選択モデルを示す情報（以下、「選択モデル情報」とも言う。）を操作部122に入力し、操作部122は、選択モデル情報123を受け付ける。例えば、選択モデル情報123は、選択モデルを示す番号であってもよい。

【0086】

なお、本発明の実施形態では、操作部122がマウスおよびキーボードである場合を主に想定する。しかし、操作部122の形態は特に限定されない。例えば、操作部122は、タッチパネルであってもよいし、他の入力装置であってもよい。

【0087】

（記録制御部108）

記録制御部108は、操作部122によってユーザから受け付けられた選択モデル情報123の記録を制御する。より詳細に、記録制御部108は、操作部122によってユーザから受け付けられた選択モデル情報123を図示しない記憶部に記憶させる。選択モデル情報123は、図示しない記憶部から後に取得され、選択モデル情報123によって示される選択モデルが、解釈性の高い学習済みモデルとして用いられ得る。

【0088】

なお、テストの終了条件は特に限定されず、ユーザにとって十分な回数のテストが行われたことを示す条件であればよい。具体的に、テストの終了条件は、テスト段階においてユーザによって推論結果の確認が所定の回数以上行われたという条件を含んでもよい。

【0089】

10

20

30

40

50

以上、本発明の実施形態に係る学習装置 10 の構成例について説明した。

【0090】

(1.2. 学習段階における動作)

図3を参照しながら、本発明の実施形態に係る学習装置10の学習段階における動作の流れについて説明する。図3は、本発明の実施形態に係る学習装置10の学習段階における動作例を示すフローチャートである。

【0091】

まず、図3に示されたように、入力部101は、データセット100から入力データ(すなわち、学習用データ)および正解値の組み合わせを取得する。さらに、推論部102は、n個の推論モデルそれぞれに対応する重みパラメータ110を取得する(S11)。推論部102は、入力部101によって取得された入力データとn個の推論モデルとに基づいて推論を行い(S12)、推論によって得られたn個の推論値を推論評価部104および説明部103それぞれに出力する。

10

【0092】

説明部103は、推論部102から入力されたn個の推論値に基づいて、n個の推論値それぞれの判断根拠を説明するヒートマップを生成する(S13)。説明部103は、生成したn個のヒートマップを説明評価部105に出力する。

【0093】

推論評価部104は、入力部101によって取得された正解値に基づいて、推論部102から入力されたn個の推論値を評価して推論評価結果を得る。より詳細に、推論評価部104は、正解値とn個の推論値とに応じた損失関数を推論評価結果として算出する。推論評価部104は、算出した推論評価結果を更新部106に出力する。

20

【0094】

説明評価部105は、説明部103から入力されたn個のヒートマップの一致度に基づいて、説明評価結果を得る。より詳細に、説明評価部105は、説明部103から入力されたn個のヒートマップ同士の一貫性に応じた損失関数を説明評価結果として算出する。説明評価部105は、算出した説明評価結果を更新部106に出力する(S14)。

【0095】

更新部106は、推論評価部104から入力された推論評価結果と、説明評価部105から入力された説明評価結果とに基づいて、第1推論モデルから第n推論モデルまでのそれぞれに対応する重みパラメータ110の更新を行う(S15)。より詳細に、更新部106は、推論評価結果と説明評価結果とに基づいて、誤差逆伝播法によって、重みパラメータ110を更新する。さらに、更新部106は、推論評価結果と説明評価結果とに基づく誤差逆伝播法によって説明部103が有する重みパラメータの更新を行う。

30

【0096】

更新部106は、入力データに基づく重みパラメータの更新が終わるたびに、学習の終了条件が満たされたか否かを判断する(S16)。学習の終了条件が満たされていないと判断した場合には(S16において「NO」)、S11に動作が移行され、入力部101によって次の入力データが取得され、推論部102、説明部103、推論評価部104、説明評価部105および更新部106それぞれによって、当該次の入力データに基づく各自の処理が再度実行される。一方、更新部106によって、学習の終了条件が満たされたと判断された場合には(S16において「YES」)、学習が終了される。

40

【0097】

以上、本発明の実施形態に係る学習装置10の学習段階における動作の流れについて説明した。

【0098】

(1.3. テスト段階における動作)

図4を参照しながら、本発明の実施形態に係る学習装置10のテスト段階における動作の流れについて説明する。図4は、本発明の実施形態に係る学習装置10のテスト段階における動作例を示すフローチャートである。

50

【 0 0 9 9 】

まず、図 4 に示されたように、入力部 1 0 1 は、データセット 1 0 0 から入力データ（すなわち、テスト用データ）および正解値の組み合わせを取得する。さらに、推論部 1 0 2 は、n 個の推論モデルそれぞれに対応する重みパラメータ 1 1 0 を取得する（S 2 1）。推論部 1 0 2 は、入力部 1 0 1 によって取得された入力データと n 個の推論モデルとに基づいて推論を行い（S 2 2）、推論によって得られた n 個の推論値を説明部 1 0 3 および提示制御部 1 0 7 それぞれに出力する。

【 0 1 0 0 】

説明部 1 0 3 は、推論部 1 0 2 から入力された n 個の推論値に基づいて、n 個の推論値それぞれの判断根拠を説明するヒートマップを生成する（S 2 3）。説明部 1 0 3 は、生成した n 個のヒートマップを提示制御部 1 0 7 に出力する。

10

【 0 1 0 1 】

提示制御部 1 0 7 は、推論部 1 0 2 から入力された n 個の推論値と、説明部 1 0 3 から入力された n 個のヒートマップとがユーザに提示されるように表示部 1 2 1 を制御する。表示部 1 2 1 は、提示制御部 1 0 7 による制御に従って、n 個の推論値と、n 個のヒートマップとを表示する（S 2 4）。

【 0 1 0 2 】

操作部 1 2 2 は、n 個の推論モデルから解釈性が高いと判断された 1 または複数の推論モデルを示す情報（選択モデル情報 1 2 3）をユーザから受け付ける。記録制御部 1 0 8 は、操作部 1 2 2 によってユーザから受け付けられた選択モデル情報 1 2 3 の記録を制御する（S 2 5）。図示しない記憶部は、記録制御部 1 0 8 による制御に従って、選択モデル情報 1 2 3 を記憶する。

20

【 0 1 0 3 】

記録制御部 1 0 8 は、入力データに基づく選択モデル情報 1 2 3 の記録制御が終わるたびに、テストの終了条件が満たされたか否かを判断する（S 2 6）。テストの終了条件が満たされていないと判断した場合には（S 2 6 において「NO」）、S 2 1 に動作が移行され、入力部 1 0 1 によって次の入力データが取得され、推論部 1 0 2、説明部 1 0 3、提示制御部 1 0 7 および記録制御部 1 0 8 それぞれによって、当該次の入力データに基づく各自の処理が再度実行される。一方、記録制御部 1 0 8 によって、テストの終了条件が満たされたと判断された場合には（S 2 6 において「YES」）、テストが終了される。

30

【 0 1 0 4 】

以上、本発明の実施形態に係る学習装置 1 0 のテスト段階における動作の流れについて説明した。

【 0 1 0 5 】

(1 . 4 . 実施形態の効果)

以上に説明したように、本発明の実施形態によれば、第 1 推論モデルから第 n 推論モデルまでのそれぞれから出力される推論値が正解値に近づくように、かつ、説明情報として出力される n 個のヒートマップ同士の一貫性が小さくなるように、学習が行われ得る。これによって、互いに異なる複数のヒートマップを出力する推論モデルを得ることができる。これによって、ユーザは、n 個のモデルの中からより解釈性の高いヒートマップを出力するモデルを選んで使用することができる。

40

【 0 1 0 6 】

以上、本発明の実施形態が奏する効果について説明した。

【 0 1 0 7 】

(2 . ハードウェア構成例)

続いて、本発明の実施形態に係る学習装置 1 0 のハードウェア構成例について説明する。以下では、本発明の実施形態に係る学習装置 1 0 のハードウェア構成例として、情報処理装置 9 0 0 のハードウェア構成例について説明する。なお、以下に説明する情報処理装置 9 0 0 のハードウェア構成例は、学習装置 1 0 のハードウェア構成の一例に過ぎない。したがって、学習装置 1 0 のハードウェア構成は、以下に説明する情報処理装置 9 0 0 の

50

ハードウェア構成から不要な構成が削除されてもよいし、新たな構成が追加されてもよい。

【0108】

図5は、本発明の実施形態に係る学習装置10の例としての情報処理装置900のハードウェア構成を示す図である。情報処理装置900は、CPU(Central Processing Unit)901と、ROM(Read Only Memory)902と、RAM(Random Access Memory)903と、ホストバス904と、ブリッジ905と、外部バス906と、インタフェース907と、入力装置908と、出力装置909と、ストレージ装置910と、通信装置911と、を備える。

【0109】

CPU901は、演算処理装置および制御装置として機能し、各種プログラムに従って情報処理装置900内の動作全般を制御する。また、CPU901は、マイクロプロセッサであってもよい。ROM902は、CPU901が使用するプログラムや演算パラメータ等を記憶する。RAM903は、CPU901の実行において使用するプログラムや、その実行において適宜変化するパラメータ等を一時記憶する。これらはCPUバス等から構成されるホストバス904により相互に接続されている。

【0110】

ホストバス904は、ブリッジ905を介して、PCI(Peripheral Component Interconnect/Interface)バス等の外部バス906に接続されている。なお、必ずしもホストバス904、ブリッジ905および外部バス906を分離構成する必要はなく、1つのバスにこれらの機能を実装してもよい。

【0111】

入力装置908は、マウス、キーボード、タッチパネル、ボタン、マイクロフォン、スイッチおよびレバー等ユーザが情報を入力するための入力手段と、ユーザによる入力に基づいて入力信号を生成し、CPU901に出力する入力制御回路等から構成されている。情報処理装置900を操作するユーザは、この入力装置908を操作することにより、情報処理装置900に対して各種のデータを入力したり処理動作を指示したりすることができる。

【0112】

出力装置909は、例えば、CRT(Cathode Ray Tube)ディスプレイ装置、液晶ディスプレイ(LCD)装置、OLED(Organic Light Emitting Diode)装置、ランプ等の表示装置およびスピーカ等の音声出力装置を含む。

【0113】

ストレージ装置910は、データ格納用の装置である。ストレージ装置910は、記憶媒体、記憶媒体にデータを記録する記録装置、記憶媒体からデータを読み出す読出し装置および記憶媒体に記録されたデータを削除する削除装置等を含んでもよい。ストレージ装置910は、例えば、HDD(Hard Disk Drive)で構成される。このストレージ装置910は、ハードディスクを駆動し、CPU901が実行するプログラムや各種データを格納する。

【0114】

通信装置911は、例えば、ネットワークに接続するための通信デバイス等で構成された通信インタフェースである。また、通信装置911は、無線通信または有線通信のどちらに対応してもよい。

【0115】

以上、本発明の実施形態に係る学習装置10のハードウェア構成例について説明した。

【0116】

(3.まとめ)

以上、添付図面を参照しながら本発明の好適な実施形態について詳細に説明したが、本発明はかかる例に限定されない。本発明の属する技術の分野における通常の知識を有する

10

20

30

40

50

者であれば、特許請求の範囲に記載された技術的思想の範疇内において、各種の変更例または修正例に想到し得ることは明らかであり、これらについても、当然に本発明の技術的範囲に属するものと了解される。

【0117】

例えば、上記した例では、学習装置10がn個の推論モデルを同時に学習する場合を主に想定している。しかし、学習装置10は、n個の推論モデルの全部を同時に学習しなくてもよい。例えば、n個の推論モデルの一部として、学習済みの推論モデルが使用されてもよい。このとき、学習済みの推論モデルの重みパラメータは、更新されずに一定の値に固定され得る。

【0118】

また、上記した例では、説明部103におけるヒートマップの生成手法の種類が、1種類である場合を主に想定している。しかし、説明部103におけるヒートマップの生成手法の種類は複数であってもよい。このとき、説明部103は、ヒートマップ同士の一貫性に基づく損失の複数種類のヒートマップ生成手法についての合計値を説明評価結果の例として更新部106に出力してもよい。

【符号の説明】

【0119】

- 10 学習装置
- 100 データセット
- 101 入力部
- 102 推論部
- 103 説明部
- 104 推論評価部
- 105 説明評価部
- 106 更新部
- 107 提示制御部
- 108 記録制御部
- 110 重みパラメータ
- 121 表示部
- 122 操作部
- 123 選択モデル情報

10

20

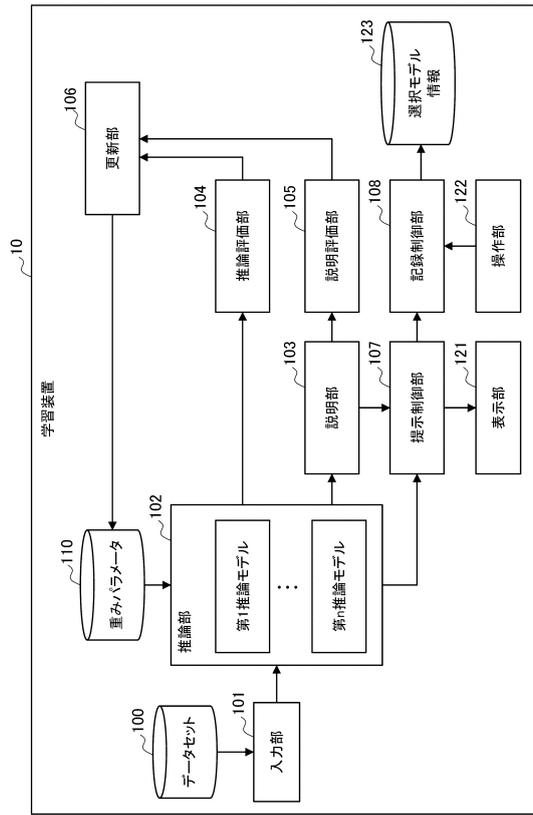
30

40

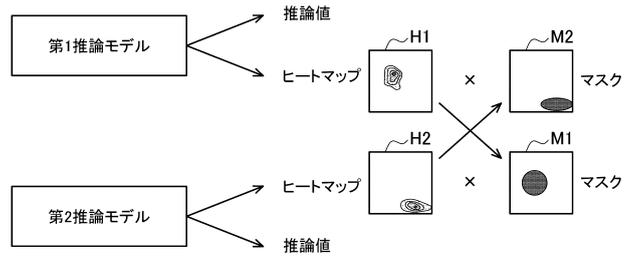
50

【 図 面 】

【 図 1 】



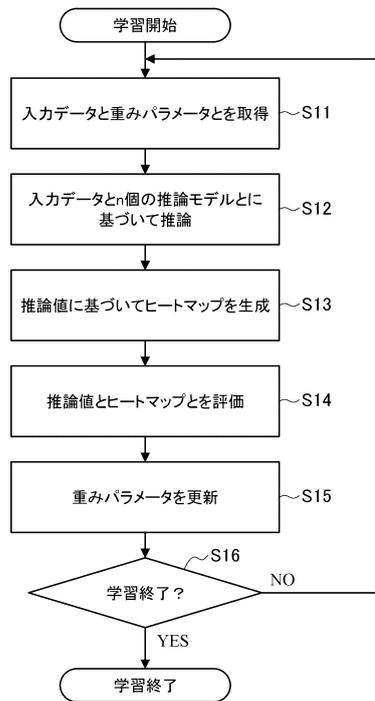
【 図 2 】



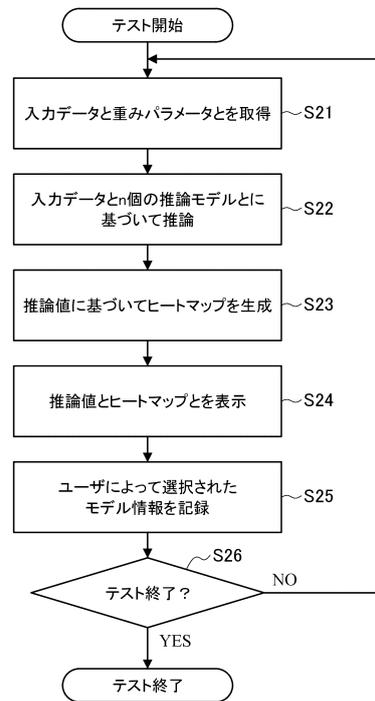
10

20

【 図 3 】



【 図 4 】

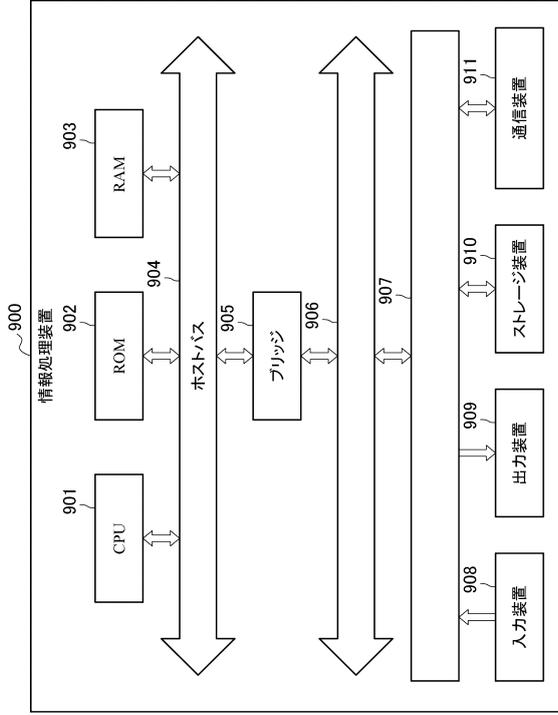


30

40

50

【図 5】



10

20

30

40

50