

【公報種別】特許法第17条の2の規定による補正の掲載
 【部門区分】第6部門第3区分
 【発行日】令和6年2月19日(2024.2.19)

【国際公開番号】WO2022/249392
 【出願番号】特願2023-523860(P2023-523860)
 【国際特許分類】
 G 0 6 N 2 0 / 0 0 (2 0 1 9 . 0 1)
 【 F I 】
 G 0 6 N 2 0 / 0 0

10

【手続補正書】
 【提出日】令和5年11月17日(2023.11.17)
 【手続補正1】
 【補正対象書類名】特許請求の範囲
 【補正対象項目名】全文
 【補正方法】変更
 【補正の内容】
 【特許請求の範囲】

【請求項1】

20

複数の訓練用例を取得する取得手段と、
 用例を入力として予測結果を出力する機械学習モデルを複数含む機械学習モデル群を、
 前記複数の訓練用例を用いて訓練する訓練手段と、
 前記複数の訓練用例のうち、訓練後の前記機械学習モデル群を用いて得られる複数の予測結果にばらつきがある訓練用例を選択する選択手段と、
 前記複数の訓練用例のうち、前記選択した訓練用例を含む2つ以上の訓練用例を合成して人工用例を生成する生成手段と、
 を備えた情報処理装置。

【請求項2】

前記生成手段は、前記選択した訓練用例と、特徴量空間において前記選択した訓練用例の近傍に存在する用例とを合成して前記人工用例を生成する、請求項1に記載の情報処理装置。

30

【請求項3】

前記選択手段は、前記複数の訓練用例のうち、2つ以上の前記複数の予測結果にばらつきがある訓練用例を選択し、
 前記生成手段は、2つ以上の前記選択した訓練用例を合成して前記人工用例を生成する、請求項1に記載の情報処理装置。

【請求項4】

前記生成手段は、
 前記選択した訓練用例と、特徴量空間において前記選択した訓練用例の近傍に存在する用例とを合成する第1生成処理と、
 2つ以上の前記選択した訓練用例を合成して前記人工用例を生成する第2生成処理と、
 の何れかを実行することにより前記人工用例を生成する、請求項1に記載の情報処理装置。

40

【請求項5】

前記人工用例を前記複数の訓練用例に追加して、前記取得手段、前記訓練手段、前記選択手段、及び前記生成手段を再度機能させる、請求項1から4の何れか1項に記載の情報処理装置。

【請求項6】

50

前記生成手段は、

複数の前記人工用例を生成し、

複数の前記人工用例のうち類似条件を満たす2つの人工用例を1つの人工用例に統合する、請求項1から4の何れか1項に記載の情報処理装置。

【請求項7】

前記生成手段は、前記人工用例のうち、訓練後の前記機械学習モデル群を用いて得られる複数の予測結果にばらつきがある人工用例を出力する、請求項1から6の何れか1項に記載の情報処理装置。

【請求項8】

前記機械学習モデル群は、前記人工用例を用いて訓練する訓練対象の機械学習モデルを含む、請求項1から7の何れか1項に記載の情報処理装置。

10

【請求項9】

複数の訓練用例を取得すること、

用例を入力として予測結果を出力する機械学習モデルを複数含む機械学習モデル群を、前記複数の訓練用例を用いて訓練すること、

前記複数の訓練用例のうち、訓練後の前記機械学習モデル群を用いて得られる複数の予測結果にばらつきがある訓練用例を選択すること、及び、

前記複数の訓練用例のうち、前記選択した訓練用例を含む2つ以上の訓練用例を合成して人工用例を生成すること、を含む情報処理方法。

20

【請求項10】

コンピュータを情報処理装置として機能させるためのプログラムであって、前記コンピュータを、

複数の訓練用例を取得する取得手段と、

用例を入力として予測結果を出力する機械学習モデルを複数含む機械学習モデル群を、前記複数の訓練用例を用いて訓練する訓練手段と、

前記複数の訓練用例のうち、訓練後の前記機械学習モデル群を用いて得られる複数の予測結果にばらつきがある訓練用例を選択する選択手段と、

前記複数の訓練用例のうち、前記選択した訓練用例を含む2つ以上の訓練用例を合成して人工用例を生成する生成手段と、として機能させるプログラム。

30

【手続補正2】

【補正対象書類名】明細書

【補正対象項目名】全文

【補正方法】変更

【補正の内容】

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、機械学習に用いる用例を生成する技術に関する。

40

【背景技術】

【0002】

機械学習モデルの推論の精度は、その機械学習モデルを構築する際に用いた訓練用例の数や内容に依存することが知られている。機械学習モデルの推論精度を向上させるために、事前に用意された訓練用例から人工用例を生成することにより、訓練用例を増強する技術が知られている。例えば非特許文献1には、サポートベクタマシン(Support Vector Machine)の決定境界に最も近いマイノリティクラスのインスタンス(訓練用例)とその近傍のマイノリティクラスのインスタンスとを合成して、マイノリティクラスの仮想インスタンスを生成することが記載されている。

【先行技術文献】

50

【非特許文献】

【0003】

【非特許文献1】Seyda Ertekin, "Adaptive Oversampling for Imbalanced Data Classification", Information Sciences and Systems 2013", proceedings of the 28th International Symposium on Computer and Information Sciences (ISCIS), pp. 261 - 269), 2013

【発明の概要】

【発明が解決しようとする課題】

【0004】

しかしながら、非特許文献1に記載の技術では、決定境界付近に仮想インスタンス（人工用例）が生成され、決定境界付近以外で訓練用例が不足している領域に人工用例が生成されない、という問題があった。

【0005】

本発明の一態様は、上記の問題に鑑みてなされたものであり、その目的の一例は、機械学習の訓練に用いる訓練用例が不足している領域に人工用例を生成可能な技術を提供することである。

【課題を解決するための手段】

【0006】

本発明の一側面に係る情報処理装置は、複数の訓練用例を取得する取得手段と、用例を入力として予測結果を出力する機械学習モデルを複数含む機械学習モデル群を、前記複数の訓練用例を用いて訓練する訓練手段と、前記複数の訓練用例のうち、訓練後の前記機械学習モデル群を用いて得られる複数の予測結果にばらつきがある訓練用例を選択する選択手段と、前記複数の訓練用例のうち、前記選択した訓練用例を含む2つ以上の訓練用例を合成して人工用例を生成する生成手段と、を備える。

【0007】

本発明の一側面に係る情報処理方法は、複数の訓練用例を取得すること、用例を入力として予測結果を出力する機械学習モデルを複数含む機械学習モデル群を、前記複数の訓練用例を用いて訓練すること、前記複数の訓練用例のうち、訓練後の前記機械学習モデル群を用いて得られる複数の予測結果にばらつきがある訓練用例を選択すること、及び前記複数の訓練用例のうち、前記選択した訓練用例を含む2つ以上の訓練用例を合成して人工用例を生成すること、を含む。

【0008】

本発明の一側面に係るプログラムは、コンピュータを情報処理装置として機能させるためのプログラムであって、前記コンピュータを、複数の訓練用例を取得する取得手段と、用例を入力として予測結果を出力する機械学習モデルを複数含む機械学習モデル群を、前記複数の訓練用例を用いて訓練する訓練手段と、前記複数の訓練用例のうち、訓練後の前記機械学習モデル群を用いて得られる複数の予測結果にばらつきがある訓練用例を選択する選択手段と、前記複数の訓練用例のうち、前記選択した訓練用例を含む2つ以上の訓練用例を合成して人工用例を生成する生成手段と、として機能させる。

【発明の効果】

【0009】

本発明の一態様によれば、機械学習の訓練に用いる訓練用例が不足している領域に人工用例を生成することができる。

【図面の簡単な説明】

【0010】

【図1】本発明の例示的实施形態1に係る情報処理装置の構成を示すブロック図である。

【図2】本発明の例示的实施形態1に係る情報処理方法の流れを示すフロー図である。

【図3】本発明の例示的实施形態1に係る情報処理方法の具体例を模式的に示す図である。

【図4】本発明の例示的实施形態1により生成される人工用例を模式的に説明する図である。

る。

【図 5】本発明の例示的实施形態 2 に係る情報処理装置の構成を示すブロック図である。

【図 6】本発明の例示的实施形態 2 に係る情報処理方法の流れを示すフロー図である。

【図 7】本発明の例示的实施形態 2 に係る情報処理方法の具体例を模式的に示す図である。

【図 8】本発明の例示的实施形態 2 に係る第 1 生成処理の流れを示すフロー図である。

【図 9】本発明の例示的实施形態 2 に係る第 2 生成処理の流れを示すフロー図である。

【図 10】本発明の例示的实施形態 2 に係る第 3 生成処理の流れを示すフロー図である。

【図 11】本発明の例示的实施形態 3 に係る情報処理方法の流れを示すフロー図である。

【図 12】本発明の例示的实施形態 4 に係る情報処理方法の流れを示すフロー図である。

10

【図 13】本発明の例示的实施形態 5 に係る情報処理方法を模式的に説明する図である。

【図 14】非特許文献 1 に記載の技術が生成する人工用例を模式的に説明する図である。

【図 15】本発明の例示的实施形態 1 ~ 5 に係る情報処理装置として機能するコンピュータの構成を示すブロック図である。

【発明を実施するための形態】

【0011】

〔例示的实施形態 1〕

本発明の第 1 の例示的实施形態について、図面を参照して詳細に説明する。本例示的实施形態は、後述する例示的实施形態の基本となる形態である。

【0012】

20

< 情報処理装置の構成 >

本例示的实施形態に係る情報処理装置 10 の構成について、図 1 を参照して説明する。図 1 は、情報処理装置 10 の構成を示すブロック図である。情報処理装置 10 は、複数の訓練用例から、機械学習モデル群を用いて人工用例を生成する装置である。

【0013】

情報処理装置 10 は、図 1 に示すように、取得部 11 と、訓練部 12 と、選択部 13 と、生成部 14 とを含む。取得部 11 は、請求の範囲に記載した取得手段を実現する構成の一例である。訓練部 12 は、請求の範囲に記載した訓練手段を実現する構成の一例である。選択部 13 は、請求の範囲に記載した選択手段を実現する構成の一例である。生成部 14 は、請求の範囲に記載した生成手段を実現する構成の一例である。

30

【0014】

取得部 11 は、複数の訓練用例を取得する。訓練部 12 は、用例を入力として予測結果を出力する機械学習モデルを複数含む機械学習モデル群を、複数の訓練用例を用いて訓練する。選択部 13 は、複数の訓練用例のうち、訓練後の機械学習モデル群を用いて得られる複数の予測結果にばらつきがある訓練用例を選択する。生成部 14 は、複数の訓練用例のうち、選択した訓練用例を含む 2 つ以上の訓練用例を合成して人工用例を生成する。

【0015】

(機械学習モデル群)

機械学習モデル群は、複数の機械学習モデルを含む。各機械学習モデルは、用例を入力として予測結果を出力する。予測結果は、例えば、複数のラベルの各々が予測される予測確率を含むものであってもよい。この場合、最も予測確率が高いラベルを、予測結果と記載する場合もある。機械学習モデルは、一例として、決定木、ニューラルネットワーク、ランダムフォレスト、又はサポートベクターマシン等の機械学習アルゴリズムを用いて生成されたモデルである。ただし、各機械学習モデルの生成に用いられる機械学習アルゴリズムは、これらに限られない。複数の機械学習モデルは、全て同一の機械学習アルゴリズムを用いて生成されたモデルであってもよい。また、複数の機械学習モデルのうち少なくとも 2 つが、互いに異なる機械学習アルゴリズムを用いて生成されたモデルであってもよい。機械学習モデル群は、例えば情報処理装置 10 のメモリに記憶されていてもよいし、情報処理装置 10 と通信可能に接続された他の装置に記憶されていてもよい。

40

【0016】

50

また、機械学習モデル群は、必ずしも全てが、「生成した人工用例を用いて訓練する訓練対象の機械学習モデル」でなくてもよい。換言すると、機械学習モデル群は、訓練対象である機械学習モデルの一部又は全部を含んでいてもよい。また、機械学習モデル群は、訓練対象である機械学習モデルを含んでいなくてもよい。訓練対象である機械学習モデルの数は複数であってもよく、また単数であってもよい。

【 0 0 1 7 】

(用例、訓練用例、人工用例)

用例は、各機械学習モデルに入力される情報であり、特徴量を含む。換言すると、用例は、特徴量空間に存在する。訓練用例は、機械学習モデル群の訓練に用いる用例である。訓練用例は、観測により得られる用例であってもよいし、人工的に生成された人工用例であつてもよい。

10

【 0 0 1 8 】

(予測結果にばらつきがある訓練用例)

複数の予測結果にばらつきがある訓練用例とは、ばらつきの評価結果が「ばらつきが大きい」ことを示す訓練用例である。例えば、ばらつきの評価とは、複数の予測結果のばらつきが大きいか否かを評価することである。具体例として、ばらつきの評価は、投票結果のエントロピーに基づく評価であってもよい。投票結果のエントロピーについては、後述の例示的实施形態2で詳細を説明する。また、ばらつきの評価は、複数の予測結果のうち同一のラベルを示す予測結果の割合に基づく評価であってもよい。ただし、ばらつきの評価は、上述したものに限られない。以降、「複数の予測結果のばらつきが大きいと評価した訓練用例」を、「予測結果にばらつきがある訓練用例」とも記載する。また、「複数の予測結果のばらつきが小さくないと評価した訓練用例」を、「予測結果のばらつきが小さい訓練用例」とも記載する。

20

【 0 0 1 9 】

< 情報処理方法の流れ >

本例示的实施形態に係る情報処理方法 S 1 0 の流れについて、図 2 を参照して説明する。図 2 は、情報処理方法 S 1 0 の流れを示すフロー図である。図 2 に示すように、情報処理方法 S 1 0 は、ステップ S 1 0 1 ~ S 1 0 4 を含む。

【 0 0 2 0 】

(ステップ S 1 0 1)

ステップ S 1 0 1 (取得処理)において、取得部 1 1 は、複数の訓練用例を取得する。例えば、取得部 1 1 は、複数の訓練用例をメモリから読み込むことにより取得してもよい。また、例えば、取得部 1 1 は、複数の訓練用例を、入力装置から取得してもよいし、ネットワークを介して接続された装置から取得してもよい。本ステップで取得する複数の訓練用例は、観測用例及び人工用例の一方又は両方を含んでいる。

30

【 0 0 2 1 】

(ステップ S 1 0 2)

ステップ S 1 0 2 (訓練処理)において、訓練部 1 2 は、ステップ S 1 0 1 で取得した複数の訓練用例を用いて、機械学習モデル群を訓練する。ここで、機械学習モデル群の各々の訓練に用いる訓練用例は、ステップ S 1 0 1 で取得した複数の訓練用例の一部であってもよいし全部であってもよい。

40

【 0 0 2 2 】

(ステップ S 1 0 3)

ステップ S 1 0 3 (選択処理)において、選択部 1 3 は、複数の訓練用例のうち、予測結果にばらつきがある訓練用例を選択する。選択部 1 3 は、そのような訓練用例を 1 つ選択してもよいし、複数選択してもよい。

【 0 0 2 3 】

具体的には、選択部 1 3 は、複数の訓練用例のうち評価対象の訓練用例を、訓練後の各機械学習モデルに入力し、各機械学習モデルから出力される予測結果を取得する。これにより、選択部 1 3 は、評価対象の訓練用例について複数の予測結果を得る。また、選択部

50

13は、得られた複数の予測結果のばらつきを評価する。選択部13は、複数の予測結果のばらつきが大きいと評価した場合、当該訓練用例を「予測結果にばらつきがある訓練用例」として選択する。

【0024】

なお、選択部13は、複数の訓練用例の一部又は全部を、ばらつきの評価対象とする。例えば、複数の訓練用例の一部を用いて機械学習モデル群の各々を訓練した場合、選択部13は、複数の訓練用例の他の一部（すなわち、機械学習モデル群の訓練に用いなかった訓練用例）をそれぞれ評価対象としてもよい。

【0025】

（ステップS104）

ステップS104（生成処理）において、生成部14は、複数の訓練用例のうち、ステップS103で選択した訓練用例を含む2つ以上の訓練用例を合成して人工用例を生成する。例えば、生成部14は、選択した訓練用例と、特徴量空間においてその近傍に存在する他の訓練用例とを合成してもよい。また、例えば、生成部14は、ステップS103において複数の訓練用例を選択した場合、選択した複数の訓練用例同士を合成してもよい。また、生成部14は、2つの訓練用例を合成して1つの人工用例を生成してもよいし、3以上の訓練用例を合成して1つの人工用例を生成してもよい。また、生成部14は、本ステップにおいて、1つの人工用例を生成してもよいし、複数の人工用例を生成してもよい。

10

【0026】

（合成処理の具体例）

2つの訓練用例を合成して1つの人工用例を生成する場合、生成部14が行う合成処理は、一例として、以下の（1）式で表される。

20

【数1】

$$\hat{x}_v = \lambda x_i + (1 - \lambda)x_j \quad \cdots (1)$$

式（1）において、 \hat{x}_v は人工用例を表し、 x_i は、選択部13が選択した訓練用例を表す。 x_j は、選択部13が選択した他の訓練用例であってもよいし、選択しなかった他の訓練用例であってもよい。選択しなかった他の訓練用例である場合、 x_j は、 x_i の近傍に存在する訓練用例である。 λ は、0から1を満たす重み係数である。生成部14は、一例として、係数 λ の値を、ランダム関数により発生させた乱数を用いて決定する。なお、生成部14が行う合成処理は上述した手法に限られず、生成部14は他の手法により複数の訓練用例を合成してもよい。

30

【0027】

<情報処理方法の具体例>

情報処理方法S10の具体例について、図3を参照して説明する。図3は、情報処理方法S10の具体例を模式的に示す図である。

【0028】

本具体例では、ステップS101において取得部11が取得する複数の訓練用例Tは、訓練用例t1, t2, t3, ...を含む。ステップS102において訓練部12が訓練する機械学習モデル群は、機械学習モデルm1, m2, m3, ...を含む。機械学習モデルm1, m2, m3, ...のそれぞれは、用例が入力されると、予測結果として「A」及び「B」の何れかのラベルを出力する。機械学習モデルm1, m2, m3, ...のそれぞれは、訓練用例Tの一部又は全部を用いて訓練される。ステップS103において、選択部13は、評価対象の訓練用例t1～t10について、複数の予測結果のばらつきを評価する。図3では、実線の円は、予測結果にばらつきがある訓練用例を示し、破線の円は、予測結果のばらつきが小さい訓練用例を示す。

40

【0029】

具体的には、訓練用例t1, t2, t5は、機械学習モデルm1, m2, m3, ...から

50

得られる複数の予測結果が全て「A」であり、この例では、予測結果のばらつきが大きくなると評価される。訓練用例 t 6、t 9、t 10 は、機械学習モデル m 1, m 2, m 3, ... から得られる複数の予測結果が全て「B」であり、この例では、予測結果のばらつきが大きくなると評価される。訓練用例 t 3、t 4 は、機械学習モデル m 1, m 2, m 3, ... から得られる複数の予測結果のうち2つが「A」であり、1つが「B」であり、この例では、予測結果のばらつきが大きいと評価される。訓練用例 t 7、t 8 は、機械学習モデル m 1, m 2, m 3, ... から得られる複数の予測結果のうち2つが「B」であり、1つが「A」であり、この例では、予測結果のばらつきが大きいと評価される。

【0030】

したがって、選択部 13 は、予測結果にばらつきがある訓練用例 t 3、t 4、t 7、t 8 を選択する。ステップ S 104 において、生成部 14 は、予想結果にばらつきがある訓練用例 t 3 とその近傍の訓練用例 t 5 とを合成することにより、人工用例 t 51 を生成する。また、生成部 14 は、予想結果にばらつきがある訓練用例 t 4 とその近傍の訓練用例 t 1 とを合成することにより、人工用例 t 52 を生成する。また、生成部 14 は、複数の予想結果にばらつきがある訓練用例 t 7、t 8 同士を合成することにより、人工用例 t 53 を生成する。図 3 では、二重線の円は、人工用例を示している。

【0031】

< 本例示的实施形態の効果 >

本例示的实施形態は、複数の訓練用例を取得し、用例を入力として予測結果を出力する機械学習モデルを複数含む機械学習モデル群を、複数の訓練用例を用いて訓練し、複数の訓練用例のうち、訓練後の機械学習モデル群を用いて得られる複数の予測結果にばらつきがある訓練用例を選択し、複数の訓練用例のうち、選択した訓練用例を含む2つ以上の訓練用例を合成して人工用例を生成する、との構成を採用している。

【0032】

ここで、複数の予想結果にばらつきがある訓練用例は、特徴量空間において訓練用例が不足している領域にあると考えられる。そのような訓練用例を含む複数の訓練用例を合成して得られる人工用例は、訓練用例が不足している領域に生成される可能性が高い。したがって、本例示的实施形態は、訓練用例が不足している領域に人工用例を生成可能である。

【0033】

このような本例示的实施形態の効果について、図 4、及び図 14 を参照して説明する。図 4 は、本例示的实施形態が生成する人工用例を模式的に説明する図である。図 14 は、非特許文献 1 に記載の技術が生成する人工用例を模式的に説明する図である。図 4、及び図 14 において、実線の円は、予測結果のばらつきが小さい訓練用例を示し、破線の円は、予測結果にばらつきがある訓練用例を示し、二重線の円は、人工用例を示す。領域 R 1、R 2、R 3 は、特徴量空間における領域を示す。領域 R 1、R 2、R 3 は、予測結果にばらつきがある人工用例を含んでおり、人工用例が不足している領域である。

【0034】

図 14 に示すように、非特許文献 1 に記載の技術は、サポートベクタマシンによる決定境界 B の近傍である領域 R 1 に人工用例を生成する。しかしながら、非特許文献 1 に記載の技術は、決定境界 B の近傍以外で訓練用例が不足している領域 R 2、R 3 に人工用例を生成することが難しい。

【0035】

これに対して、図 4 に示すように、本例示的实施形態は、予測結果にばらつきがある訓練用例を含む複数の訓練用例を合成して人工用例を生成する。したがって、本例示的实施形態は、訓練用例が不足している領域 R 1、R 2、R 3 に、人工用例を生成することができる。また、本例示的实施形態は、一部の領域 R 1 に偏って人工用例を生成することを抑制できる。

【0036】

また、本例示的实施形態は、予測結果にばらつきがある訓練用例を選択するために、機

械学習モデル群を用いる、との構成を採用している。

【0037】

これにより、本例示的实施形態は、非特許文献1に記載の技術のように決定境界の近傍に人工用例を生成する場合と比較して、偏った領域に人工用例を生成することを抑制することができる。

【0038】

また、これにより、本例示的实施形態は、予測結果にばらつきがある訓練用例を選択するために予測確率を用いる場合と比較して、より訓練用例が不足している領域に人工用例を生成することができる。これは、例えば、機械学習モデル群に決定木が含まれる場合、決定木は、予測確率1で間違った予測をする可能性があるためである。

10

【0039】

〔例示的实施形態2〕

本発明の第2の例示的实施形態について、図面を参照して詳細に説明する。なお、例示的实施形態1にて説明した構成要素と同じ機能を有する構成要素については、同じ符号を付し、その説明を適宜省略する。

【0040】

<情報処理装置の構成>

本例示的实施形態に係る情報処理装置20の構成について、図5を参照して説明する。図5は、情報処理装置20の構成を示すブロック図である。情報処理装置20は、複数の用例から、機械学習モデル群COM0を用いて人工用例を生成する装置である。

20

【0041】

(機械学習モデル群)

機械学習モデル群COM0は、例示的实施形態1における機械学習モデル群とほぼ同様に構成される。ただし、本例示的实施形態では、機械学習モデル群COM0に含まれる複数の機械学習モデルは、そのうち少なくとも2つが、互いに異なる機械学習アルゴリズムを用いて生成されたものである。

【0042】

また、機械学習モデル群COM0は、人工用例を用いて訓練する訓練対象である機械学習モデルを含む。また、機械学習モデル群COM0に含まれる機械学習モデルは、少なくとも1つが決定木である。ここでは、訓練対象である機械学習モデルが、決定木である。

30

【0043】

情報処理装置20は、図5に示すように、取得部21と、訓練部22と、選択部23と、生成部24と、ラベル付与部25と、出力部26と、制御部27とを含む。取得部21は、請求の範囲に記載した取得手段を実現する構成の一例である。訓練部22は、請求の範囲に記載した訓練手段を実現する構成の一例である。選択部23は、請求の範囲に記載した選択手段を実現する構成の一例である。生成部24は、請求の範囲に記載した生成手段を実現する構成の一例である。ラベル付与部25は、請求の範囲に記載したラベル付与手段を実現する構成の一例である。出力部26は、請求の範囲に記載した出力手段を実現する構成の一例である。

【0044】

取得部21は、例示的实施形態1における取得部11と同様に構成される。

40

【0045】

訓練部22は、例示的实施形態1における訓練部12とほぼ同様に構成されるが、機械学習モデル群COM0を複数のグループに分けて訓練する点が異なる。訓練部22による訓練処理の詳細については後述する。

【0046】

選択部23は、例示的实施形態1における選択部23とほぼ同様に構成されるが、予測結果のばらつきの評価対象とする訓練用例の詳細が異なる。評価対象とする訓練用例の詳細については後述する。

【0047】

50

生成部 2 4 は、選択部 2 3 が選択した訓練用例と、特徴量空間において選択した訓練用例の近傍に存在する用例とを合成して人工用例を生成する。生成部 2 4 は、一例として、上記式 (1) により 2 つの訓練用例を合成して人工用例を生成する。

【 0 0 4 8 】

ラベル付与部 2 5 は、複数の訓練用例及び人工用例の一部又は全部にラベルを付与する。ラベル付与部 2 5 は、一例として、ユーザ操作を受け付ける入力装置から出力される情報に基づきラベルを付与してもよい。また、一例として、ラベル付与部 2 5 は、用例を入力としてラベルを出力するよう訓練された機械学習モデルに、訓練用例および人工用例を入力することにより得られるラベルを付与してもよい。この場合、ラベルを出力する機械学習モデルは、機械学習モデル群 C O M 0 に含まれる各機械学習モデルとは異なる機械学習モデルである。また、ラベルを出力する機械学習モデルは、機械学習モデル群に含まれる少なくとも 1 つの機械学習モデルより予測精度の高いモデルであることが望ましい。例えば、機械学習モデル群に含まれる、訓練対象の機械学習モデルが決定木であれば、ラベルを出力する機械学習モデルは、ランダムフォレストであってもよい。

10

【 0 0 4 9 】

出力部 2 6 は、生成部 2 4 が生成した人工用例を出力する。出力部 2 6 は一例として、生成部 2 4 が生成した人工用例を、外部記憶装置等の記録媒体に格納してもよい。また、出力部 2 6 は一例として、表示装置等の出力装置に人工用例を出力してもよい。

【 0 0 5 0 】

制御部 2 7 は、情報処理装置 2 0 の各部を制御する。本例示的实施形態において制御部 2 7 は特に、生成部 2 4 が生成した人工用例を複数の訓練用例に追加して、取得部 2 1 、訓練部 2 2 、選択部 2 3 、および生成部 2 4 を再度機能させる。

20

【 0 0 5 1 】

< 情報処理方法の流れ >

本例示的实施形態に係る情報処理方法 S 2 0 の流れについて、図 6 を参照して説明する。図 6 は、情報処理方法 S 2 0 の流れを示すフロー図である。

【 0 0 5 2 】

(ステップ S 2 0 1)

ステップ S 2 0 1 (取得処理) において、取得部 2 1 は、複数の訓練用例を取得する。取得する複数の訓練用例は、観測により得られた用例を含んでいてもよいし、人工用例を含んでいてもよい。

30

【 0 0 5 3 】

(ステップ S 2 0 2)

ステップ S 2 0 2 において、ラベル付与部 2 5 は、取得部 2 1 が取得した複数の訓練用例の各々にラベルを付与する。

【 0 0 5 4 】

(ステップ S 2 0 3)

ステップ S 2 0 3 (訓練処理) において、訓練部 2 2 は、取得部 2 1 が取得した複数の訓練用例の一部または全部を用いて、複数の機械学習モデル群の各々を訓練する。機械学習モデル群の各々を訓練する訓練処理の詳細については後述する。

40

【 0 0 5 5 】

(ステップ S 2 0 4)

ステップ S 2 0 4 (選択処理) において、選択部 2 3 は、取得部 2 1 が取得した複数の訓練用例のうち、予測結果にばらつきがある訓練用例を 1 つ以上選択する。選択部 2 3 が行う選択処理については後述する。予測結果にばらつきがある訓練用例を選択する処理の詳細については後述する。

【 0 0 5 6 】

(ステップ S 2 0 5)

ステップ S 2 0 5 (生成処理) において、生成部 2 4 は、選択部 2 3 が選択した訓練用例を含む複数の訓練用例を、合成対象として特定する。また、生成部 2 4 は、合成対象と

50

して特定した複数の訓練用例を合成して人工用例を生成する。生成部 24 が行う生成処理の詳細については後述する。

【0057】

(ステップ S206)

ステップ S206 において、ラベル付与部 25 は、生成部 24 が生成した人工用例の各々にラベルを付与する。ステップ S207 において、制御部 27 は、訓練処理を終了するかを判定する。制御部 27 は、一例として、ステップ S203 ~ S206 の処理を実行した回数が所定の閾値以上である場合、訓練処理を終了すると判定する。一方、ステップ S203 ~ S206 の処理を実行した回数が所定の閾値未満である場合、訓練処理を終了しないと判定する。訓練処理を終了しない場合(ステップ S207 にて NO)、制御部 27 はステップ S208 の処理に進む。一方、訓練処理を終了する場合(ステップ S207 にて YES)、制御部 27 はステップ S209 の処理に進む。

10

【0058】

(ステップ S208)

ステップ S208 において、制御部 27 は、これまでに実行したステップ S206 で生成された 1 以上の人工用例を複数の訓練用例に追加する。ステップ S208 の処理を終えると、制御部 27 は、ステップ S203 の処理に戻る。換言すると、制御部 27 は、人工用例を複数の訓練用例に追加して、取得部 21、訓練部 22、選択部 23、および生成部 24 を再度機能させる。

【0059】

(ステップ S209)

ステップ S209 において、出力部 26 は、これまでに実行したステップ S206 で生成された 1 以上の人工用例を出力する。

20

【0060】

< 訓練対象の機械学習モデルの訓練 >

このようにして情報処理方法 S20 を用いて生成された 1 つ以上の人工用例は、訓練対象の機械学習モデルを訓練するために用いられる。訓練対象の機械学習モデルを訓練する処理は、例えば、訓練部 22 が実行してもよい。

【0061】

(訓練処理、選択処理の具体例)

ステップ S203 ~ S204 における訓練処理及び選択処理の具体例について、図 7 を参照して説明する。図 7 は、情報処理方法 S20 の具体例を模式的に示す図である。

30

【0062】

図 7 に示すように、ステップ S203 において、訓練部 22 は、機械学習モデル群 COM0 を、複数のグループ COM_i ($i = 1, 2, \dots, M$ 、 M は 2 以上の整数) に分割して訓練を行う。以降、分割した各グループを、機械学習モデル群 COM_i と記載する。また、機械学習モデル群 COM_i には、複数の機械学習モデル m_{i-j} ($j = 1, 2, \dots$) が含まれる。以降、機械学習モデル群 COM_i に含まれる複数の機械学習モデルを、機械学習モデル m_{i-j} と記載する。機械学習モデル群 COM_i に含まれる複数の機械学習モデル m_{i-j} は、全てが同一の機械学習アルゴリズムにより生成されたモデルであってもよいし、そのうち少なくとも 2 つが互いに異なる機械学習アルゴリズムにより生成されたモデルであってもよい。また、機械学習モデル群 COM_{i1} に含まれる機械学習モデル m_{i1-j} の個数は、機械学習モデル群 COM_{i2} に含まれる機械学習モデル m_{i2-j} の個数と同一であってもよいし、異なってもよい ($i1 = 1, 2, \dots, M$ 、 $i2 = 1, 2, \dots, M$ 、 $i1 \neq i2$)。

40

【0063】

ステップ S203 において、訓練部 22 は、ステップ S201 において取得部 21 が取得した訓練用例群 T から、訓練用例群 D_i を抽出する。訓練用例群 D_i は、訓練用例群 T の一部である。例えば、訓練部 22 は、ランダムサンプリングにより訓練用例群 D_i を抽出してもよい。訓練用例群 D_{i1} と D_{i2} とは、含まれる訓練用例が全て同一であっても

50

よいし、一部または全部が異なってもよい。訓練部 2 2 は、訓練用例群 D_i を用いて、機械学習モデル群 COM_i に含まれる各機械学習モデル m_{i-j} を訓練することを、 $i = 1, 2, \dots, M$ について繰り返す。

【0064】

ステップ S 2 0 4 において、選択部 2 3 は、機械学習モデル群 COM_i を用いて、予測結果にばらつきがある訓練用例を選択することを、 $i = 1, 2, \dots, M$ について繰り返す。具体的には、選択部 2 3 は、機械学習モデル群 COM_i の訓練に用いなかった各訓練用例（すなわち、訓練用例群 T のうち訓練用例群 D_i 以外）について、予測結果のばらつきを評価する。これにより、選択部 2 3 は、このような評価対象の訓練用例のうち、予測結果にばらつきがある訓練用例を選択する。図 7 の例では、選択部 2 3 は、機械学習モデル群 COM_1 を用いて、予測結果にばらつきがある訓練用例 t_1, t_2, \dots を選択している。また、選択部 2 3 は、機械学習モデル群 COM_2 を用いて、予測結果にばらつきがある訓練用例 t_{11}, t_{12}, \dots を選択している。

10

【0065】

選択部 2 3 は例えば、評価対象の各訓練用例について、QBC (query by committee) の手法における投票結果のエントロピー (vote entropy) の指標を用いて、ばらつきの評価を行う。例えば、以下の式 (2) は、投票結果のエントロピーが最大である訓練用例 \hat{x} を示す式である。

【数 2】

$$\hat{x} = \operatorname{argmax}_x \left(- \sum_y \frac{V(y)}{C} \log \frac{V(y)}{C} \right) \quad \dots (2)$$

20

式 (2) において、 C は、機械学習モデル群 COM_i における機械学習モデル m_{i-j} の総数を示す。 $V(y)$ は、機械学習モデル群 COM_i においてラベル y を予測した機械学習モデル m_{i-j} の数を示す。選択部 2 3 は、式 (2) が示す訓練用例 \hat{x} を、予測結果にばらつきがある訓練用例として選択してもよい。この場合、各機械学習モデル群 COM_i について選択部 2 3 が選択する、予測結果にばらつきがある訓練用例の数は 1 つである。換言すると、この場合、選択部 2 3 は、 M 個の機械学習モデル群 COM_i を用いて、予測結果にばらつきがある訓練用例を M 個選択する。また、選択部 2 3 は、各機械学習モデル群 COM_i について、投票結果のエントロピーが大きい順に所定数の訓練用例を選択してもよいし、投票結果のエントロピーが閾値以上の訓練用例を選択してもよい。この場合、各機械学習モデル群 COM_i について選択部 2 3 が選択する、予測結果にばらつきがある訓練用例の数は複数でありうる。換言すると、この場合、選択部 2 3 は、 M 個の機械学習モデル群 COM_i を用いて、予測結果にばらつきがある訓練用例を M 個以上選択する。さらに、選択部 2 3 は、このようにして選択した、予測結果にばらつきがある M 個以上の訓練用例の中から、ランダムに 1 つ又は所定数を選択してもよいし、投票結果のエントロピーが大きい順に 1 つ又は所定数を選択してもよい。

30

【0066】

(生成処理の具体例)

ステップ S 2 0 5 における生成処理の具体例について説明する。ステップ S 2 0 5 において、生成部 2 4 は、予測結果にばらつきがある訓練用例を用いて、第 1 生成処理 S 3 0、第 2 生成処理 S 4 0、及び第 3 生成処理 S 5 0 の何れかを実行することにより、人工用例を生成する。第 1 生成処理 S 3 0 は、予測結果にばらつきがある訓練用例とその近傍の訓練用例とを合成して人工用例を生成する処理である。第 2 生成処理は、予測結果にばらつきがある 2 つ以上の訓練用例を合成して人工用例を生成する処理である。第 3 生成処理は、第 1 生成処理及び第 2 生成処理の何れかを選択的に実行する処理である。

40

【0067】

(第 1 生成処理)

50

第1生成処理S30について、図8を参照して説明する。図8は、第1生成処理S30の流れを示すフロー図である。図8において、第1生成処理S30は、ステップS301～S302を含む。ここで、先行して実施されたステップS204では、予測結果にばらつきがある1または複数の訓練用例が選択されている。生成部24は、予測結果にばらつきがある1または複数の訓練用例のそれぞれ（以下では、当該訓練用例と記載）について、以下のステップS301～S302を実行する。

【0068】

（ステップS301）

ステップS301において、生成部24は、当該訓練用例の近傍の訓練用例を選択する。近傍の訓練用例は、予測結果にばらつきがある訓練用例であってもよいし、予測結果のばらつきが小さい訓練用例であってもよい。例えば、近傍の訓練用例は、訓練用例群Tのうち、当該訓練用例との特徴量空間における距離が最も近い訓練用例であってもよい。また、例えば、近傍の訓練用例は、訓練用例群Tのうち、当該訓練用例との特徴量空間における距離が閾値以下の訓練用例であってもよい。

10

【0069】

（ステップS302）

ステップS302において、生成部24は、当該訓練用例と、ステップS301で選択した近傍の訓練用例とを合成して人工用例を生成する。例えば、図7の例では、予測結果にばらつきがある訓練用例t1とその近傍の訓練用例とを合成して人工用例tv1-1が生成される。また、予測結果にばらつきがある訓練用例t2とその近傍の訓練用例とを合成して人工用例tv1-2が生成される。

20

【0070】

ここで、生成部24は、合成処理の一例として、上記式(1)を用いてもよい。また、生成部24は、合成処理の他の例として、MUNGE（参考文献1参照）、SMOTE（参考文献2参照）等の公知の技術を用いてもよい。

【0071】

[参考文献1] Bucilua, C., Caruana, R. and Niculescu-Mizil, A., "Model Compression", Proc. ACM SIGKDD, pp. 535-541 (2006)

[参考文献2] Chawla, N. V., Bowyer, K. W., Hall, L. O., and Kegelmeyer, W. P., "SMOTE: Synthetic minority over-sampling technique", Journal of Artificial Intelligent Research, 16, 321-357 (2002).

30

（第2生成処理）

第2生成処理S40について、図9を参照して説明する。図9は、第2生成処理S40の流れを説明するフロー図である。図9に示すように、第2生成処理S40は、ステップS401～S402を含む。なお、第2生成処理は、ステップS204において、予測結果にばらつきがある訓練用例が複数選択されている場合に実行可能である。生成部24は、生成部24は、予測結果にばらつきがある複数の訓練用例のそれぞれ（以下、当該訓練用例と記載）について、以下のステップS401～S402を実行する。

【0072】

（ステップS401）

ステップS401において、選択部23は、予測結果にばらつきがある複数の訓練用例のうち、当該訓練用例とは異なる他の訓練用例を選択する。換言すると、選択部23は、複数の訓練用例のうち、複数の予測結果にばらつきがある2つ以上の訓練用例を選択する。例えば、選択部23は、そのような他の訓練用例を、予測結果にばらつきがある複数の訓練用例からランダムに選択してもよい。また、例えば、選択部23は、そのような他の訓練用例として、予測結果にばらつきがある複数の訓練用例のうち当該訓練用例との特徴量空間における距離が最も小さいもの、又は距離が閾値以下のものを選択してもよい。なお、当該訓練用例が、既に合成に用いられている場合、当該訓練用例に関するステップS401～S402の処理は、実行されなくてもよい。

40

【0073】

50

(ステップ S 4 0 2)

ステップ S 4 0 2 において、生成部 2 4 は、当該訓練用例と、ステップ S 4 0 1 で選択した他の訓練用例とを合成して人工用例を生成する。例えば、図 7 の例では、予測結果にばらつきがある訓練用例 t_{11} 及び t_{12} を合成して人工用例 t_{v2-1} が生成される。ここで、生成部 2 4 が合成する 2 以上の訓練用例は、この例のように、同一の機械学習モデル群 COM i を用いて選択したものであってもよいし、そのうち少なくとも 1 つが他とは異なる機械学習モデル群 COM i を用いて選択したものであってもよい。例えば、図 7 の例では、生成部 2 4 は、予測結果にばらつきのある訓練用例 t_1 、 t_2 、...、 t_{11} 、 t_{12} 、... の中から 2 つ以上の訓練用例を選択し、選択した訓練用例を合成して人工用例 t_{v1-1} 、 t_{v1-2} 、 t_{v2-1} 、又は t_{v2-2} を生成してもよい。なお、ステップ S 4 0 2 における合成処理に用いる手法については、ステップ S 3 0 2 で説明した通りであるため、詳細な説明を繰り返さない。

10

【0074】

(第3生成処理)

第3生成処理について、図 10 を参照して説明する。図 10 は、第3生成処理 S 5 0 の流れを説明するフロー図である。図 10 において、第3生成処理 S 5 0 は、ステップ S 5 0 1 ~ S 5 0 3 を含む。ここで、ステップ S 2 0 4 では、予測結果にばらつきがある 1 または複数の訓練用例が選択されている。生成部 2 4 は、予測結果にばらつきがある 1 または複数の訓練用例のそれぞれ（以下では、当該訓練用例と記載）について、以下のステップ S 5 0 1 ~ S 5 0 3 を実行する。

20

【0075】

(ステップ S 5 0 1)

ステップ S 5 0 1 において、生成部 2 4 は、第1生成処理及び第2生成処理の何れかを選択する。例えば、生成部 2 4 は、ランダム関数により決定した確率 p を用いて第1生成処理を選択し、第1生成処理を選択しなかった場合に第2生成処理を選択してもよい。なお、第1生成処理及び第2生成処理の何れかを選択する手法は、確率 p を用いる手法に限らず、他の手法であってもよい。

【0076】

(ステップ S 5 0 2 ~ S 5 0 4)

ステップ S 5 0 2 において、生成部 2 4 は、いずれを選択したかを判定する。第1生成処理を選択した場合、生成部 2 4 はステップ S 5 0 3 の処理に進み、第1生成処理を実行する。一方、第2生成処理を選択した場合、生成部 2 4 はステップ S 5 0 4 の処理に進み、第2生成処理を実行する。第1生成処理及び第2生成処理の詳細については、上述した通りである。

30

【0077】

<本例示的实施形態の効果>

本例示的实施形態は、予測結果にばらつきがある訓練用例と、その近傍にある訓練用例とを合成して人工用例を生成する第1生成処理を実行する、との構成を有する。

【0078】

ここで、第1生成処理により生成される人工用例は、予測結果にばらつきがある訓練用例の近傍に生成される。予測結果にばらつきがある訓練用例は、特徴量空間において訓練用例が不足している領域にあると考えられる。したがって、このような人工用例は、訓練用例が不足している領域に生成される。

40

【0079】

また、本例示的实施形態は、予測結果にばらつきがある 2 つ以上の訓練用例を合成して人工用例を生成する第2生成処理を実行する、との構成を有する。

【0080】

ここで、第2生成処理により生成される人工用例は、訓練用例が不足している領域にある訓練用例同士が合成されたものである。したがって、このような人工用例が存在する領域も、訓練用例が不足している可能性が高い。

50

【 0 0 8 1 】

また、本例示的实施形態は、第 1 生成処理と、第 2 生成処理との何れかを選択して実行することにより人工用例を生成する第 3 生成処理を実行する、との構成を有する。

【 0 0 8 2 】

ここで、例えば、第 3 生成処理により生成される人工用例は、第 1 生成処理又は第 2 生成処理により生成される。第 1 生成処理により人工用例が生成される領域と、第 2 生成処理により人工用例が生成される領域とは、異なり得る。したがって、第 3 生成処理により複数の人工用例を生成する場合、これらの人工用例は、訓練用例が不足しているより多様な領域に分散して生成される可能性が高くなる。

【 0 0 8 3 】

その結果、本例示的实施形態は、第 1 生成処理、第 2 生成処理、及び第 3 生成処理の何れかを実行することにより、訓練用例が十分である領域に偏って過剰に人工用例を生成することがなく、訓練用例がより不足している領域に人工用例を生成することができる。

【 0 0 8 4 】

また、本例示的实施形態は、機械学習モデル群が、訓練対象の機械学習モデルを含んでいる。これにより、本例示的实施形態は、訓練対象の機械学習モデルの精度向上により効果的な人工用例を生成することができる。

【 0 0 8 5 】

また、本例示的实施形態は、機械学習モデル群のうち少なくとも 2 つが、互いに異なる機械学習アルゴリズムにより生成されるモデルである、との構成を採用している。

【 0 0 8 6 】

これにより、予測結果にばらつきがある訓練用例をより精度よく選択することができる。

【 0 0 8 7 】

また、本例示的实施形態は、訓練対象の機械学習モデルが決定木であり、サポートベクタマシンではない。また、本例示的实施形態は、そのような訓練対象の機械学習モデルが機械学習モデル群 C O M 0 に含まれている。したがって、サポートベクタマシンの決定境界近傍に人工用例を生成する非特許文献 1 に記載の技術と比較して、本例示的实施形態は、訓練対象の機械学習モデルの精度向上により効果的な人工用例を生成することができる。

【 0 0 8 8 】

〔 例示的实施形態 3 〕

本発明の例示的实施形態 3 について、図面を参照して詳細に説明する。なお、例示的实施形態 2 にて説明した構成要素と同じ機能を有する構成要素については、同じ符号を付記してその説明を繰り返さない。本例示的实施形態は、例示的实施形態 2 における生成部 2 4 を次のように変形した形態である。

【 0 0 8 9 】

＜ 生成部の構成 ＞

本例示的实施形態において、生成部 2 4 は、複数の人工用例を生成する。また、生成部 2 4 は、生成した複数の人工用例のうち、類似条件を満たす 2 つの人工用例を 1 つの人工用例に統合する。ここで、類似条件とは、用例が類似することを示す条件である。類似条件は、例えばコサイン類似度が閾値以上であることであってもよいし、特徴量空間における距離が閾値以下であることであってもよい。ただし、類似条件はこれらに限られない。統合する処理の詳細については後述する。

【 0 0 9 0 】

＜ 情報処理方法の流れ ＞

本例示的实施形態における情報処理方法 S 2 0 A について、図 1 1 を参照して説明する。図 1 1 は、例示的实施形態 3 に係る情報処理方法 S 2 0 A の流れを説明するフロー図である。図 1 1 に示す情報処理方法 S 2 0 A は、例示的实施形態 2 に係る情報処理方法 S 2 0 とほぼ同様に構成されるが、ステップ S 2 0 5 A をさらに含む点が異なる。

10

20

30

40

50

【 0 0 9 1 】

(ステップ S 2 0 5 A)

ステップ S 2 0 5 A において、生成部 2 4 は、ステップ S 2 0 5 において生成した人工用例のうち、類似する 2 つの人工用例を統合する。具体的には、生成部 2 4 は、今回のステップ S 2 0 5 において生成した人工用例と、前回までのステップ S 2 0 5 において生成した人工用例の何れかが類似条件を満たすか否かを判定する。類似条件を満たすと判定した場合、生成部 2 4 は、類似条件を満たす 2 つの人工用例を統合する。

【 0 0 9 2 】

(統合処理の具体例)

統合処理の一例として、2 つの人工用例を合成する処理が挙げられる。この場合、生成部 2 4 は、2 つの人工用例を合成して 1 つの人工用例を生成し、類似条件を満たした元の 2 つの人工用例を削除する。また、統合処理の他の例として、2 つの人工用例のうち一方を削除する処理が挙げられる。なお、統合処理は、類似条件を満たす 2 つの人工用例の代わりに、当該 2 つの人工用例を参照して生成した 1 つの人工用例を採用する処理であればよく、上述した処理に限られない。なお、人工用例を削除するとは、ステップ S 2 0 6 でラベルを付与する対象、及びステップ S 2 0 8 で訓練用例に追加する対象から削除することである。これにより、統合された人工用例に対して、ラベルが付与されるとともに訓練用例に追加される。

10

【 0 0 9 3 】

< 本例示的实施形態の効果 >

本例示的实施形態においては、生成部が、複数の人工用例を生成し、生成した複数の人工用例のうち、類似条件を満たす 2 つの人工用例を 1 つの人工用例に統合する、との構成が採用されている。

20

【 0 0 9 4 】

ここで、訓練用例が不足している領域に存在する複数の用例が類似する場合、それらの用例を用いて機械学習モデルを訓練することは、機械学習モデルの精度向上において効率的ではない。したがって、本例示的实施形態は、類似条件を満たす人工用例を統合することにより、訓練用例が不足している領域に、より効率的に機械学習モデルの精度を向上させることができる人工用例を生成することができる。

【 0 0 9 5 】

〔 例示的实施形態 4 〕

本発明の例示的实施形態 4 について、図面を参照して詳細に説明する。なお、例示的实施形態 2 にて説明した構成要素と同じ機能を有する構成要素については、同じ符号を付記してその説明を繰り返さない。本例示的实施形態は、例示的实施形態 2 における生成部 2 4 を次のように変形した形態である。

30

【 0 0 9 6 】

< 生成部の構成 >

本例示的实施形態において、生成部 2 4 は、生成した人工用例のうち、訓練後の機械学習モデル群 C O M 0 を用いて得られる複数の予測結果にばらつきがある人工用例を出力する。ここで、ばらつきがある人工用例は、ばらつきの評価結果が「ばらつきが大きい」ことを示す人工用例である。ばらつきの評価の詳細については、上述した通りであるため、詳細を繰り返さない。換言すると、生成部 2 4 は、生成した人工用例のばらつきを、訓練後の機械学習モデル群 C O M 0 を用いて事後評価し、事後評価により予測結果にばらつきがある人工用例を採用する。

40

【 0 0 9 7 】

< 情報処理方法の流れ >

本例示的实施形態における情報処理方法 S 2 0 B について、図 1 2 を参照して説明する。図 1 2 は、例示的实施形態 4 に係る情報処理方法 S 2 0 B の流れを説明するフロー図である。図 1 2 に示す情報処理方法 S 2 0 B は、例示的实施形態 2 に係る情報処理方法 S 2 0 とほぼ同様に構成されるが、ステップ S 2 0 5 B をさらに含む点が異なる。

50

【0098】

(ステップS205B)

ステップS205Bにおいて、生成部24は、ステップS205において生成した人工用例を事後評価する。

【0099】

具体的には、生成部24は、当該人工用例について、機械学習モデル群COM0を用いて予測結果のばらつきを評価する。例えば、図7に示した例では、生成部24は、人工用例tv1-1について、機械学習モデル群COM1を用いて予測結果のばらつきを評価する。このように、ばらつきの評価に用いる機械学習モデル群COM1は、当該人工用例tv1-1を生成するために参照した訓練用例t1の評価に用いたものであることが望ましい。機械学習モデル群COM0を用いて予測結果のばらつきを評価する処理の詳細については、例示的实施形態2において説明した通りである。

10

【0100】

生成部24は、ステップS205において生成した人工用例について、予測結果のばらつきが大きくないと評価した場合には、当該人工用例を削除する。ここで、人工用例を削除するとは、ステップS206でラベルを付与する対象、及びステップS208で訓練用例に追加する対象から削除することである。これにより、予測結果にばらつきがある人工用例に対してラベルが付与されるとともに、訓練用例に追加される。

【0101】

<本例示的实施形態の効果>

本例示的实施形態においては、生成部が、生成した人工用例のうち、訓練後の機械学習モデル群を用いて得られる複数の予測結果にばらつきがある人工用例を出力する、との構成が採用されている。

20

【0102】

ここで、予測結果にばらつきがある訓練用例を含む複数の訓練用例を合成して得られた人工用例は、必ずしも予測結果にばらつきがあるとは限らない。換言すると、このようにして生成した人工用例は、予測結果のばらつきが小さい可能性がある。予測結果のばらつきが小さい訓練用例を用いて機械学習モデルを訓練することは、機械学習モデルの精度向上において効率的ではない。したがって、本例示的实施形態は、生成した人工用例を事後評価することにより、訓練用例が不足している領域に、より効率的に機械学習モデルの精度を向上させることができる人工用例を生成することができる。

30

【0103】

[例示的实施形態5]

本発明の例示的实施形態5について、図面を参照して詳細に説明する。なお、例示的实施形態2にて説明した構成要素と同じ機能を有する構成要素については、同じ符号を付記してその説明を繰り返さない。

【0104】

本例示的实施形態は、例示的实施形態2における機械学習モデル群COM0の構成、及び情報処理方法S20におけるステップS203~S204を次のように変形した形態である。本例示的实施形態について、図13を参照して説明する。図13は、本例示的实施形態に係る情報処理方法を模式的に説明する図である。

40

【0105】

(機械学習モデル群)

図13に示すように、本例示的实施形態では、機械学習モデル群COM0は、機械学習モデルmj(j=1,2,...,M)を含む。各機械学習モデルmjは、同一の機械学習アルゴリズムによって生成されるモデルである。例えば、各機械学習モデルmjは、決定木であってもよい。

【0106】

(ステップS203)

本例示的实施形態のステップS203において、訓練部22は、ステップS201で取

50

得部 2 1 が取得した訓練用例群 T から、訓練用例群 D j を抽出する。訓練用例群 D j は、訓練用例群 T の一部である。例えば、訓練部 2 2 は、ランダムサンプリングにより訓練用例群 D j を抽出してもよい。訓練部 2 2 は、訓練用例群 D j を用いて、機械学習モデル m j を訓練することを、j = 1, 2, ..., M について繰り返す。

【0107】

ここで、訓練用例群 D j 1 と D j 2 とは、含まれる訓練用例が全て同一であってもよいが、一部または全部が異なることが望ましい (j 1 = 1, 2, ..., M、j 2 = 1, 2, ..., M、j 1 ≠ j 2)。少なくとも一部が異なる訓練用例群 D j 1 と D j 2 とをそれぞれ用いることにより、機械学習モデル群 m j 1 及び m j 2 は、それぞれを構成するパラメータが互いに異なるように訓練される。

10

【0108】

(ステップ S 2 0 4)

本例示的实施形態のステップ S 2 0 4 において、選択部 2 3 は、機械学習モデル群 C O M 0 を用いて、訓練用例群 T に含まれる各訓練用例について、予測結果のばらつきを評価する。また、選択部 2 3 は、予測結果にばらつきがある訓練用例を選択する。図 1 3 の例では、選択部 2 3 は、機械学習モデル群 C O M 0 を用いて、予測結果にばらつきがある訓練用例 t 1, t 3, ... を選択している。

【0109】

ステップ S 2 0 5 の処理は、例示的实施形態 2 で説明した通りである。すなわち、図 1 3 の例では、予測結果にばらつきがある訓練用例 t 1、t 2、... の各々について、第 1 生成処理、第 2 生成処理、及び第 3 生成処理の何れかが実行される。これにより、人工用例 t v 1、t v 2、... が生成される。

20

【0110】

< 本例示的实施形態の効果 >

本例示的实施形態は、機械学習モデル群を構成する機械学習モデルとして、全て同一の機械学習アルゴリズムによって生成されたモデルを用い、取得した訓練用例群の中から予測結果にばらつきがある訓練用例を選択する、との構成を採用している。

【0111】

これにより、本例示的实施形態は、取得した訓練用例群の全てに亘って訓練用例が不足している領域に、人工用例を生成することができる。

30

【0112】

また、本例示定期実施形態は、機械学習モデル群に含まれる機械学習モデルが全て決定木である場合、当該機械学習モデルの精度向上により効果的な人工用例を生成することができる。その理由について説明する。決定木は、訓練用例の小さな変更に対して木の構造が大きく変化する。そのため、複数の決定木を含む機械学習モデル群を用いることにより、予測結果にばらつきがある訓練用例をより精度よく選択することができるためである。

【0113】

[ソフトウェアによる実現例]

情報処理装置 1 0, 2 0 の一部又は全部の機能は、集積回路 (IC チップ) 等のハードウェアによって実現してもよいし、ソフトウェアによって実現してもよい。

40

【0114】

後者の場合、情報処理装置 1 0, 2 0 は、例えば、各機能を実現するソフトウェアであるプログラムの命令を実行するコンピュータによって実現される。このようなコンピュータの一例 (以下、コンピュータ C と記載する) を図 1 5 に示す。コンピュータ C は、少なくとも 1 つのプロセッサ C 1 と、少なくとも 1 つのメモリ C 2 と、を備えている。メモリ C 2 には、コンピュータ C を情報処理装置 1 0, 2 0 として動作させるためのプログラム P が記録されている。コンピュータ C において、プロセッサ C 1 は、プログラム P をメモリ C 2 から読み取って実行することにより、情報処理装置 1 0, 2 0 の各機能が実現される。

50

【 0 1 1 5 】

プロセッサ C 1 としては、例えば、C P U (Central Processing Unit)、G P U (Graphic Processing Unit)、D S P (Digital Signal Processor)、M P U (Micro Processing Unit)、F P U (Floating point number Processing Unit)、P P U (Physics Processing Unit)、マイクロコントローラ、又は、これらの組み合わせなどを用いることができる。メモリ C 2 としては、例えば、フラッシュメモリ、H D D (Hard Disk Drive)、S S D (Solid State Drive)、又は、これらの組み合わせなどを用いることができる。

【 0 1 1 6 】

なお、コンピュータ C は、プログラム P を実行時に展開したり、各種データを一時的に記憶したりするための R A M (Random Access Memory) を更に備えていてもよい。また、コンピュータ C は、他の装置との間でデータを送受信するための通信インタフェースを更に備えていてもよい。また、コンピュータ C は、キーボードやマウス、ディスプレイやプリンタなどの入出力機器を接続するための入出力インタフェースを更に備えていてもよい。

【 0 1 1 7 】

また、プログラム P は、コンピュータ C が読み取り可能な、一時的でない有形の記録媒体 M に記録することができる。このような記録媒体 M としては、例えば、テープ、ディスク、カード、半導体メモリ、又はプログラマブルな論理回路などを用いることができる。コンピュータ C は、このような記録媒体 M を介してプログラム P を取得することができる。また、プログラム P は、伝送媒体を介して伝送することができる。このような伝送媒体としては、例えば、通信ネットワーク、又は放送波などを用いることができる。コンピュータ C は、このような伝送媒体を介してプログラム P を取得することもできる。

【 0 1 1 8 】

〔 付 記 事 項 1 〕

本発明は、上述した実施形態に限定されるものでなく、請求項に示した範囲で種々の変更が可能である。例えば、上述した実施形態に開示された技術的手段を適宜組み合わせ得られる実施形態についても、本発明の技術的範囲に含まれる。

【 0 1 1 9 】

〔 付 記 事 項 2 〕

上述した実施形態の一部又は全部は、以下のようにも記載され得る。ただし、本発明は、以下の記載する態様に限定されるものではない。

【 0 1 2 0 】

(付 記 1)

複数の訓練用例を取得する取得手段と、
 用例を入力として予測結果を出力する機械学習モデルを複数含む機械学習モデル群を、前記複数の訓練用例を用いて訓練する訓練手段と、
 前記複数の訓練用例のうち、訓練後の前記機械学習モデル群を用いて得られる複数の予測結果にばらつきがある訓練用例を選択する選択手段と、
 前記複数の訓練用例のうち、前記選択した訓練用例を含む 2 つ以上の訓練用例を合成して人工用例を生成する生成手段と、
 を備えた情報処理装置。

【 0 1 2 1 】

上記構成により、予測結果にばらつきがある訓練用例を用いて人工用例を生成するので、決定境界の近傍のように偏った領域に人工用例を生成することがなく、訓練用例が不足する領域に精度よく人工用例を生成することができる。

【 0 1 2 2 】

(付 記 2)

前記生成手段は、前記選択した訓練用例と、特徴量空間において前記選択した訓練用例の近傍に存在する用例とを合成して前記人工用例を生成する、付記 1 に記載の情報処理装

置。

【0123】

上記構成により、予測結果にばらつきがある訓練用例の近傍に人工用例を生成するので、訓練用例が不足する領域に精度よく人工用例を生成することができる。

【0124】

(付記3)

前記選択手段は、前記複数の訓練用例のうち、2つ以上の前記複数の予測結果にばらつきがある訓練用例を選択し、

前記生成手段は、2つ以上の前記選択した訓練用例を合成して前記人工用例を生成する、付記1に記載の情報処理装置。

【0125】

上記構成により、予測結果にばらつきがある訓練用例同士を合成して人工用例を生成するので、訓練用例が不足する領域に精度よく人工用例を生成することができる。

【0126】

(付記4)

前記生成手段は、

前記選択した訓練用例と、特徴量空間において前記選択した訓練用例の近傍に存在する用例とを合成する第1生成処理と、

2つ以上の前記選択した訓練用例を合成して前記人工用例を生成する第2生成処理と

、
の何れかを実行することにより前記人工用例を生成する、付記1に記載の情報処理装置。

【0127】

上記構成により、第1生成処理及び第2生成処理の何れかを選択的に用いて人工用例を生成するので、複数の人工用例を生成する場合には、訓練用例が不足するより多様な領域に人工用例を生成することができる。

【0128】

(付記5)

前記人工用例を前記複数の訓練用例に追加して、前記取得手段、前記訓練手段、前記選択手段、及び前記生成手段を再度機能させる、付記1から4の何れか1つに記載の情報処理装置。

【0129】

上記構成により、生成した人工用例を加えた訓練用例を用いて機械学習モデル群を訓練することを繰り返すので、予測結果にばらつきがある訓練用例をより精度よく選択することができる。その結果、人工用例がより不足している領域に人工用例を生成することができる。

【0130】

(付記6)

前記生成手段は、

複数の前記人工用例を生成し、

複数の前記人工用例のうち類似条件を満たす2つの人工用例を1つの人工用例に統合する、付記1から4の何れか1つに記載の情報処理装置。

【0131】

上記構成により統合した人工用例を用いて機械学習モデルを訓練する場合、既に用いた人工用例に類似する人工用例を用いて訓練することを避けられる。したがって、より効率的に機械学習モデルの精度を向上させることができる人工用例を生成することができる。

【0132】

(付記7)

前記生成手段は、前記人工用例のうち、訓練後の前記機械学習モデル群を用いて得られる複数の予測結果にばらつきがある人工用例を出力する、付記1から6の何れか1つに記

10

20

30

40

50

載の情報処理装置。

【0133】

上記構成により出力した人工用例を用いて機械学習モデルを訓練する場合、予測結果のばらつきが小さい人工用例を用いて訓練することを避けられる。したがって、より効率的に機械学習モデルの精度を向上させることができる人工用例を生成することができる。

【0134】

(付記8)

前記機械学習モデル群は、前記人工用例を用いて訓練する訓練対象の機械学習モデルを含む、付記1から7の何れか1つに記載の情報処理装置。

【0135】

上記構成により、生成した人工用例を用いて訓練対象の機械学習モデルを訓練すれば、より効果的に訓練対象の機械学習モデルの精度を向上させることができる。

【0136】

(付記9)

前記機械学習モデル群のうち少なくとも2つは、互いに異なる機械学習アルゴリズムを用いる、付記1から8の何れか1つに記載の情報処理装置。

【0137】

上記構成により、予測結果にばらつきがあるより多様な訓練用例を選択することができる。

【0138】

(付記10)

前記機械学習モデル群のそれぞれは、同一の機械学習アルゴリズムを用いる、付記1から8の何れか1つに記載の情報処理装置。

【0139】

上記構成により、予測結果にばらつきがある訓練用例をより精度よく選択することができる。

【0140】

(付記11)

前記機械学習モデル群のうち少なくとも1つは決定木である、付記1から10の何れか1つに記載の情報処理装置。

【0141】

上記構成により、決定木の精度をより効果的に向上させることができる人工用例を生成することができる。

【0142】

(付記12)

前記複数の訓練用例及び前記人工用例の一部又は全部にラベルを付与するラベル付与手段をさらに備える、付記1から11の何れか1つに記載の情報処理装置。

【0143】

上記構成により、用例にラベルが付与されていることを前提する訓練手法を用いて、機械学習モデル群又は訓練対象の機械学習モデルを訓練することができる。

【0144】

(付記13)

複数の訓練用例を取得すること、
用例を入力として予測結果を出力する機械学習モデルを複数含む機械学習モデル群を、前記複数の訓練用例を用いて訓練すること、
前記複数の訓練用例のうち、訓練後の前記機械学習モデル群を用いて得られる複数の予測結果にばらつきがある訓練用例を選択すること、及び、
前記複数の訓練用例のうち、前記選択した訓練用例を含む2つ以上の訓練用例を合成して人工用例を生成すること、
を含む情報処理方法。

10

20

30

40

50

【 0 1 4 5 】

上記構成により、付記 1 と同様の効果を奏する。

【 0 1 4 6 】

(付記 1 4)

コンピュータを情報処理装置として機能させるためのプログラムであって、前記コンピュータを、

複数の訓練用例を取得する取得手段と、

用例を入力として予測結果を出力する機械学習モデルを複数含む機械学習モデル群を、前記複数の訓練用例を用いて訓練する訓練手段と、

前記複数の訓練用例のうち、訓練後の前記機械学習モデル群を用いて得られる複数の予測結果にばらつきがある訓練用例を選択する選択手段と、 10

前記複数の訓練用例のうち、前記選択した訓練用例を含む 2 つ以上の訓練用例を合成して人工用例を生成する生成手段と、

として機能させるプログラム。

【 0 1 4 7 】

上記構成により、付記 1 と同様の効果を奏する。

【 0 1 4 8 】

(付記 1 5)

付記 1 4 に記載のプログラムが記録された、コンピュータ読み取り可能な記録媒体。

【 0 1 4 9 】

上記の構成によれば、付記 1 と同様の効果を奏する。 20

【 0 1 5 0 】

〔 付記事項 3 〕

上述した実施形態の一部又は全部は、更に、以下のように表現することもできる。

【 0 1 5 1 】

少なくとも 1 つのプロセッサを備え、前記プロセッサは、

複数の訓練用例を取得する取得処理と、

用例を入力として予測結果を出力する機械学習モデルを複数含む機械学習モデル群を、前記複数の訓練用例を用いて訓練する訓練処理と、

前記複数の訓練用例のうち、訓練後の前記機械学習モデル群を用いて得られる複数の予測結果にばらつきがある訓練用例を選択する選択処理と、 30

前記複数の訓練用例のうち、前記選択した訓練用例を含む 2 つ以上の訓練用例を合成して人工用例を生成する生成処理と、を実行する情報処理装置。

【 0 1 5 2 】

なお、この情報処理装置は、更にメモリを備えていてもよく、このメモリには、前記取得処理と、前記訓練処理と、前記選択処理と、前記生成処理とを前記プロセッサに実行させるためのプログラムが記憶されていてもよい。また、このプログラムは、コンピュータ読み取り可能な一時的でない有形の記録媒体に記録されていてもよい。

【 符号の説明 】

【 0 1 5 3 】

- 1 0、2 0 情報処理装置
- 1 1、2 1 取得部
- 1 2、2 2 訓練部
- 1 3、2 3 選択部
- 1 4、2 4 生成部
- 2 5 ラベル付与部
- 2 6 出力部
- 2 7 制御部

40

50