

(19)日本国特許庁(JP)

## (12)特許公報(B2)

(11)特許番号  
特許第7424507号  
(P7424507)

(45)発行日 令和6年1月30日(2024.1.30)

(24)登録日 令和6年1月22日(2024.1.22)

(51)国際特許分類	F I
G 0 6 N 20/00 (2019.01)	G 0 6 N 20/00
G 0 6 N 3/096(2023.01)	G 0 6 N 3/096

請求項の数 8 (全24頁)

(21)出願番号	特願2022-556825(P2022-556825)	(73)特許権者	000005223 富士通株式会社 神奈川県川崎市中原区上小田中4丁目1番1号
(86)(22)出願日	令和2年10月16日(2020.10.16)	(74)代理人	110002147 弁理士法人酒井国際特許事務所
(86)国際出願番号	PCT/JP2020/039191	(72)発明者	金月 寛彰 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
(87)国際公開番号	WO2022/079919	審査官	多賀 実
(87)国際公開日	令和4年4月21日(2022.4.21)		
審査請求日	令和4年12月21日(2022.12.21)		

最終頁に続く

(54)【発明の名称】 検知プログラム、検知方法および検知装置

## (57)【特許請求の範囲】

## 【請求項1】

第1の複数のデータの入力に応じて第1の機械学習モデルから出力された第1の結果と前記第1の複数のデータとに基づいた機械学習によって生成された第2の機械学習モデルに第2の複数のデータを入力し、

前記第2の複数のデータの入力に応じて前記第2の機械学習モデルから出力された第2の結果を取得し、

前記第2の結果と前記第2の機械学習モデルの損失関数の勾配とに基づいて算出される値と閾値との比較に基づいて、前記第1の複数のデータの分布と前記第2の複数のデータの分布との差を検知する、

処理をコンピュータに実行させることを特徴とする検知プログラム。

## 【請求項2】

前記第1の複数のデータからデータのペアとなる第1ペアデータを選択し、前記第1ペアデータのデータをそれぞれ前記第2の機械学習モデルに入力することで、第1ペアデータのスコア差を算出し、前記スコア差が所定のスコア値未満となるように前記第2の機械学習モデルの出力を調整するハイパーパラメータを算出する処理を更にコンピュータに実行させることを特徴とする請求項1に記載の検知プログラム。

## 【請求項3】

前記第2の結果を取得する処理は、前記第2の複数のデータからデータのペアとなる第2ペアデータを選択し、前記第2ペアデータのデータをそれぞれ前記第2の機械学習モデ

ルに入力することで、前記第 2 ペアデータのデータのスコア差を算出し、前記検知する処理は、前記第 2 ペアデータに関する損失関数の勾配によって、前記第 2 ペアデータのスコア差を除算した値と、閾値とを基にして、前記第 1 の複数のデータの分布と前記第 2 の複数のデータの分布との差を検知することを特徴とする請求項 2 に記載の検知プログラム。

【請求項 4】

前記第 1 の複数のデータからデータのペアとなる第 1 ペアデータを選択し、前記第 1 ペアデータのデータをそれぞれ前記第 2 の機械学習モデルに入力することで、前記第 2 ペアデータのデータのスコア差を算出し、前記第 1 ペアデータに関する損失関数の勾配によって、前記第 1 ペアデータのスコア差を除算した値を算出する処理を繰り返し実行し、算出した複数の値を基にして前記閾値を算出する処理を更にコンピュータに実行させることを特徴とする請求項 3 に記載の検知プログラム。

10

【請求項 5】

前記検知する処理によって差が検知されない場合に、前記第 2 の複数のデータを前記第 1 の機械学習モデルに入力して、前記第 2 の複数のデータを分類する処理を更にコンピュータに実行させることを特徴とする請求項 1 に記載の検知プログラム。

【請求項 6】

前記検知する処理によって差が検知された場合に、前記第 1 の機械学習モデルに対して再度機械学習を実行する処理を更にコンピュータに実行させることを特徴とする請求項 1 に記載の検知プログラム。

【請求項 7】

第 1 の複数のデータの入力に応じて第 1 の機械学習モデルから出力された第 1 の結果と前記第 1 の複数のデータとに基づいた機械学習によって生成された第 2 の機械学習モデルに第 2 の複数のデータを入力し、

20

前記第 2 の複数のデータの入力に応じて前記第 2 の機械学習モデルから出力された第 2 の結果を取得し、

前記第 2 の結果と前記第 2 の機械学習モデルの損失関数の勾配とに基づいて算出される値と閾値との比較に基づいて、前記第 1 の複数のデータの分布と前記第 2 の複数のデータの分布との差を検知する、

処理をコンピュータが実行することを特徴とする検知方法。

【請求項 8】

第 1 の複数のデータの入力に応じて第 1 の機械学習モデルから出力された第 1 の結果と前記第 1 の複数のデータとに基づいた機械学習によって生成された第 2 の機械学習モデルに第 2 の複数のデータを入力し、前記第 2 の複数のデータの入力に応じて前記第 2 の機械学習モデルから出力された第 2 の結果を取得する取得部と、

30

前記第 2 の結果と前記第 2 の機械学習モデルの損失関数の勾配とに基づいて算出される値と閾値との比較に基づいて、前記第 1 の複数のデータの分布と前記第 2 の複数のデータの分布との差を検知する検知部と

を有することを特徴とする検知装置。

【発明の詳細な説明】

【技術分野】

40

【0001】

本発明は、運用中の機械学習モデルの精度劣化を検知する検知プログラム等に関する。

【背景技術】

【0002】

近年、企業等で利用されている情報システムに対して、データの判定機能、分類機能等を有する機械学習モデルの導入が進んでいる。以下、情報システムを「システム」と表記する。機械学習モデルは、システム開発時に学習させた訓練データの通りに判定、分類を行うため、システム運用中に入力データの傾向が変化すると、機械学習モデルの精度が劣化する。

【0003】

50

図 2 1 は、入力データの傾向の変化による機械学習モデルの劣化を説明するための図である。ここで説明する機械学習モデルは、入力データを第 1 クラス、第 2 クラス、第 3 クラスのいずれかに分類するモデルであり、システム運用前に、訓練データに基づき、予め訓練されているものとする。

【 0 0 0 4 】

図 2 1 において、分布 1 A は、システム運用初期の入力データの分布を示す。分布 1 B は、システム運用初期から T 1 時間経過した時点の入力データの分布を示す。分布 1 C は、システム運用初期から更に T 2 時間経過した時点の入力データの分布を示す。時間経過に伴って、入力データの傾向（特徴量等）が変化するものとする。たとえば、入力データが画像であれば、同一の被写体を撮影した画像であっても、季節や時間帯に応じて、入力データの傾向が変化する。

10

【 0 0 0 5 】

決定境界 3 は、モデル適用領域 3 a ~ 3 c の境界を示すものである。たとえば、モデル適用領域 3 a は、第 1 クラスに属する訓練データが分布する領域である。モデル適用領域 3 b は、第 2 クラスに属する訓練データが分布する領域である。モデル適用領域 3 c は、第 3 クラスに属する訓練データが分布する領域である。

【 0 0 0 6 】

星印は、第 1 クラスに属する入力データであり、機械学習モデルに入力した際に、モデル適用領域 3 a に分類されることが正しい。三角印は、第 2 クラスに属する入力データであり、機械学習モデルに入力した際に、モデル適用領域 3 b に分類されることが正しい。丸印は、第 3 クラスに属する入力データであり、機械学習モデルに入力した際に、モデル適用領域 3 a に分類されることが正しい。

20

【 0 0 0 7 】

分布 1 A では、全ての入力データが正常なモデル適用領域に分布している。すなわち、星印の入力データがモデル適用領域 3 a に位置し、三角印の入力データがモデル適用領域 3 b に位置し、丸印の入力データがモデル適用領域 3 c に位置している。

【 0 0 0 8 】

分布 1 B では、入力データの傾向が変化したため、全ての入力データが、正常なモデル適用領域に分布しているものの、星印の入力データの分布がモデル適用領域 3 b の方向に変化している。

30

【 0 0 0 9 】

分布 1 C では、入力データの傾向が更に変化し、星印の一部の入力データが、決定境界 3 を跨いで、モデル適用領域 3 b に移動しており、適切に分類されておらず、正解率が低下している（機械学習モデルの精度が劣化している）。

【 0 0 1 0 】

ここで、運用中の機械学習モデルの精度劣化を検知する技術として、 $T^2$  統計量（Hotelling's T-square）を用いる従来技術がある。この従来技術では、入力データおよび正常データ（訓練データ）のデータ群を主成分分析し、入力データの  $T^2$  統計量を算出する。 $T^2$  統計量は、標準化した各主成分の原点からデータまでの距離の二乗を合計したものである。従来技術は、入力データ群の  $T^2$  統計量の分布の変化を基にして、機械学習モデルの精度劣化を検知する。たとえば、入力データ群の  $T^2$  統計量は、異常値データの割合に対応する。

40

【先行技術文献】

【非特許文献】

【 0 0 1 1 】

【文献】A.Shabbak and H. Midi,"An Improvement of the Hotelling Statistic in Monitoring Multivariate Quality Characteristics",Mathematical Problems in Engineering (2012) 1-15.

【発明の概要】

【発明が解決しようとする課題】

50

## 【 0 0 1 2 】

しかしながら、上述した従来技術では、機械学習モデルの精度劣化を引き起こす可能性のあるデータの分布の変化を検知することができないという問題がある。

## 【 0 0 1 3 】

たとえば、元々の情報量が非常に大きい高次元（数千～数万次元）データでは、従来技術のように、主成分分析により次元を削減すると、ほとんどの情報が失われてしまう。そのため、分類や判定を行うための重要な情報（特徴量）まで落ちてしまい、機械学習モデルの精度劣化を引き起こす可能性のあるデータの分布の変化を検知することができない。

## 【 0 0 1 4 】

1つの側面では、本発明は、機械学習モデルの精度劣化を引き起こす可能性のあるデータの分布の変化を検知することができる検知プログラム、検知方法および検知装置を提供することを目的とする。

10

## 【課題を解決するための手段】

## 【 0 0 1 5 】

1つの態様において、検知プログラムは、コンピュータに次の処理を実行させる。コンピュータは、第1の複数のデータの入力に応じて第1の機械学習モデルから出力された第1の結果と第1の複数のデータとに基づいた機械学習によって生成された第2の機械学習モデルに第2の複数のデータを入力する。コンピュータは、第2の複数のデータの入力に応じて第2の機械学習モデルから出力された第2の結果を取得する。コンピュータは、第2の結果と第2の機械学習モデルの損失関数の勾配とに基づいて算出される値と閾値との比較に基づいて、第1の複数のデータの分布と第2の複数のデータの分布との差を検知する。

20

## 【発明の効果】

## 【 0 0 1 6 】

機械学習モデルの精度劣化を引き起こす可能性のあるデータの分布の変化を検知することができる。

## 【図面の簡単な説明】

## 【 0 0 1 7 】

【図1】図1は、参考技術を説明するための図である。

【図2】図2は、精度劣化予測の一例を示す図である。

30

【図3】図3は、コンセプトドリフトの一例を示す図である。

【図4】図4は、インスペクターモデルの基本的な仕組みを説明するための図である。

【図5】図5は、参考技術の問題を説明するための図である。

【図6】図6は、統計的検定の問題を説明するための図である。

【図7】図7は、本実施形態に係る検知装置の処理を説明するための図である。

【図8】図8は、知識蒸留を説明するための図である。

【図9】図9は、本実施形態に係る検知装置の構成を示す機能ブロック図である。

【図10】図10は、訓練データセットのデータ構造の一例を示す図である。

【図11】図11は、運用モデルの一例を説明するための図である。

【図12】図12は、疑似サンプルテーブルのデータ構造の一例を示す図である。

40

【図13】図13は、蒸留データテーブルのデータ構造の一例を示す図である。

【図14】図14は、運用データセットテーブルのデータ構造の一例を示す図である。

【図15】図15は、本実施形態に係る特徴空間の決定境界を説明するための図である。

【図16】図16は、ハイパーパラメータに応じたスコア差の分布を示す図である。

【図17】図17は、本実施形態に係る検知装置の処理手順を示すフローチャートである。

【図18】図18は、各機械学習モデルの決定境界の性質を示す図(1)である。

【図19】図19は、各機械学習モデルの決定境界の性質を示す図(2)である。

【図20】図20は、本実施形態に係る検知装置と同様の機能を実現するコンピュータのハードウェア構成の一例を示す図である。

【図21】図21は、入力データの傾向の変化による機械学習モデルの劣化を説明するた

50

めの図である。

【発明を実施するための形態】

【0018】

以下に、本願の開示する検知プログラム、検知方法および検知装置の実施例を図面に基  
づいて詳細に説明する。なお、この実施例によりこの発明が限定されるものではない。

【0019】

本実施形態の説明を行う前に、機械学習モデルの精度劣化を検知する参考技術について  
説明する。参考技術では、異なる条件でモデル適用領域を狭めた複数の監視器を用いて、  
機械学習モデルの精度劣化を検知する。以下の説明では、監視器を「インスペクターモデ  
ル」と表記する。

10

【0020】

図1は、参考技術を説明するための図である。機械学習モデル10は、訓練データを用  
いて機械学習を実行することで生成される機械学習モデルである。参考技術では、機械学  
習モデル10の精度劣化を検知する。たとえば、訓練データは、機械学習モデル10のパ  
ラメータを訓練する場合に用いられるものであり、正解ラベルが対応付けられる。

【0021】

インスペクターモデル11A, 11B, 11Cは、それぞれ異なる条件でモデル適用領  
域が狭められ、異なる決定境界を有する。参考技術では、訓練データに何らかの改変を加  
え、改変を加えた訓練データを用いて、インスペクターモデル11A~11Cを訓練して  
いる。

20

【0022】

インスペクターモデル11A~11Cは、それぞれ決定境界が異なるため、同一の入力  
データを入力しても、出力結果が異なる場合がある。参考技術では、インスペクターモデ  
ル11A~11Cの出力結果の違いを基にして、機械学習モデル10の精度劣化を検知す  
る。図1に示す例では、インスペクターモデル11A~11Cを示すが、他のインスペク  
ターモデルを用いて、精度劣化を検知してもよい。インスペクターモデル11A~11C  
にはDNN(Deep Neural Network)を利用する。

【0023】

参考技術では、インスペクターモデル11A~11Cの出力結果が全て同じである場合  
に、機械学習モデル10の精度が劣化していないと判定する。一方、参考技術では、イン  
スペクターモデル11A~11Cの出力結果が異なる場合に、機械学習モデル10の精度  
劣化を検知する。

30

【0024】

図2は、精度劣化予測の一例を示す図である。図2のグラフの縦軸は、精度に対応する  
軸であり、横軸は時刻に対応する軸である。図2に示すように、時間経過に伴って、精度  
が低下しており、時刻t1において、精度の許容限界を下回る。たとえば、参考技術では  
、時刻t1において、精度劣化(許容限界を下回ったこと)を検知する。

【0025】

時間経過に伴う入力データの分布(特徴量)の変化をコンセプトドリフトと呼ぶ。図3  
は、コンセプトドリフトの一例を示す図である。図3の縦軸は、第1の特徴量に対応する  
軸であり、横軸は、第2の特徴量に対応する軸である。たとえば、機械学習モデル10の  
運用開始時において、第1クラスに対応する第1データの分布を分布A<sub>1</sub>とし、第2クラ  
スに対応する第2データの分布を分布Bとする。

40

【0026】

時間経過に伴って、第1データの分布A<sub>1</sub>が、分布A<sub>2</sub>に変化する場合がある。オリジ  
ナルの機械学習モデル10は、第1データの分布を、分布A<sub>1</sub>として学習を行っているた  
め、時間経過に伴って精度が下がり、再訓練が必要となる。

【0027】

コンセプトドリフトが発生するデータには、スパムメール、電気需要予測、株価予測、  
ポーカーハンドの戦略手順、画像等が含まれる。たとえば、画像は、季節や時間帯によっ

50

て、同一の被写体であっても、画像の特徴量が異なる。

【0028】

ここで、上述した参考技術では、機械学習モデル10の精度劣化を検知するために、複数のインスペクターモデル11A～11Cを訓練している。そして、複数のインスペクターモデル11A～11Cを訓練するためには、機械学習モデル10や、機械学習モデル10の訓練時に用いた、訓練データに何らかの改変を加えることができるという条件が必須である。たとえば、機械学習モデル10が確信度を算出するモデルであること等、機械学習モデル10が特定の機械学習モデルであることが求められる。

【0029】

図4は、インスペクターモデルの基本的な仕組みを説明するための図である。たとえば、インスペクターモデルは、第1クラスに属する訓練データの分布 $A_1$ と、第2クラスに属する訓練データの分布Bとの境界となる決定境界5の機械学習を実行することで、作成される。時間経過に伴う、運用データに対する機械学習モデル10の精度劣化を検出するためには、決定境界5の危険領域5aを監視し、危険領域5aに含まれる運用データの数が増加したか（または減少したか）否かを特定し、運用データの数が増加した（または減少した）場合に、精度劣化を検知する。

10

【0030】

ここで、上述した参考技術の問題について説明する。図5は、参考技術の問題を説明するための図である。たとえば、インスペクターモデル等を用いて、危険領域5aを監視する手法を用いるとケース1またはケース2のようになる。

20

【0031】

ケース1のように、時間経過に伴って、データの分布Aと、データの分布Bとが近づく方向に変化した場合に、危険領域5aのデータの数が変化するため、参考技術によって、精度劣化を検知することができる。

【0032】

一方、ケース2のように、時間経過に伴って、データの分布Aと、データの分布Bとが離れていく方向に変化した場合、危険領域5aに含まれるデータの数が変化しないため、参考技術によって、精度劣化を検知することができない。

【0033】

ケース2のように、データの分布Aと、データの分布Bとが離れていく場合、データの答えそのものが変化している可能性があるが、参考技術で説明したインスペクターモデルを用いても、各インスペクターモデルから出力される結果に変化はない。たとえば、インスペクターモデルは、データが、第1クラスの領域に含まれていれば、決定境界からいくら離れても、データを第1クラスに分類する。参考技術では、データの答えが変化していないことを暗黙的に仮定している。

30

【0034】

ところで、決定境界によらないデータのシフトは統計的検定で検出するが可能である。たとえば、統計的検定には、Student検定、Kolmogorov Smirnov検定、L2距離を用いる方法、コサイン距離を用いる方法、KL距離を用いる方法、ワーサースタイン距離を用いる方法等がある。

40

【0035】

しかし、統計的検定を用いると、分類に関係のない特徴量（ノイズ）の移動もすべて検知してしまうため、誤検知が多いという問題がある。図6は、統計的検定の問題を説明するための図である。図6では、x軸、y軸、z軸からなる特徴空間に位置するデータ群6aと、6a-1を用いて説明する。ここで、決定境界7が、xy平面に位置していると、z軸方向の変化は、分類結果と無関係となる。しかし、統計的検定を用いると、時間経過に伴って、データ群6b-1がz軸方向に移動し、データ群6b-2に変化すると、かかる変化を検知し、誤検知となる。

【0036】

次に、本実施形態に係る検知装置の処理の一例について説明する。図7は、本実施形態

50

に係る検知装置の処理を説明するための図である。図7では、学習フェーズおよび運用フェーズにおける検知装置の処理について説明する。

【0037】

検知装置の学習フェーズの処理について説明する。検知装置は、訓練データセット141を入力として、運用モデル50の機械学習を実行する。訓練データセット141は、複数の訓練データを含み、複数の訓練データには、正解ラベルが付与される。検知装置は、機械学習を実行済みの運用モデル50に訓練データセット141を入力した場合の出力と、訓練データセット141とを入力として、インスペクターモデル55の機械学習を実行する。たとえば、検知装置は、知識蒸留(KD: Knowledge Distiller)によって、インスペクターモデル55の機械学習を実行する。

10

【0038】

図8は、知識蒸留を説明するための図である。知識蒸留では、Teacherモデル7Aの出力値を模倣するような、Studentモデル7Bを構築する。Teacherモデル7Aは、図7の運用モデル50に対応する。Studentモデル7Bは、図7のインスペクターモデル55に対応する。たとえば、訓練データ6が与えられ、訓練データ6には正解ラベル「犬」が付与されているものとする。説明の便宜上、Teacherモデル7AおよびStudentモデル7BをNN(Neural Network)とするが、これに限定されるものではない。

【0039】

検知装置は、訓練データ6を入力した際のTeacherモデル7Aの出力結果が、正解ラベル「犬」に近づくように、Teacherモデル7Aのパラメータを訓練(誤差逆伝播法による機械学習を実行)する。また、検知装置は、訓練データ6を入力した際のStudentモデル7Bの出力結果が、訓練データ6を入力した際のTeacherモデル7Aの出力結果に近づくように、Studentモデル7Bのパラメータを訓練する。Teacherモデル7Aの出力を「ソフトターゲット(Soft Target)」と呼ぶ。訓練データの正解ラベルを「ハードターゲット(Hard Target)」と呼ぶ。

20

【0040】

上記のように、Teacherモデル7Aを、訓練データ6とハードターゲットとを用いて訓練し、Studentモデル7Bを、訓練データ6とソフトターゲットとを用いて訓練する手法を、知識蒸留と呼ぶ。検知装置は、他の訓練データについても同様にして、Teacherモデル7AおよびStudentモデル7Bを訓練する。たとえば、検知装置は、訓練データセット141と、運用モデル50から出力されるソフトターゲットとを用いて、インスペクターモデル55の機械学習を実行する。

30

【0041】

図7の説明に戻る。検知装置の運用フェーズの処理について説明する。検知装置は、運用データセットC0に含まれる複数の運用データをインスペクターモデル55に入力し、インスペクターモデル55から出力される結果を取得する。検知装置は、インスペクターモデル55から出力される結果およびインスペクターモデル55の損失係数の勾配を基に算出される値と、閾値とを比較し、コンセプトドリフトを検知する。

【0042】

たとえば、インスペクターモデル55から出力される結果およびインスペクターモデル55の損失係数の勾配を基に算出される値は、モデル適用領域の境界を示す決定境界からの距離を示す。以下の説明では、インスペクターモデル55から出力される結果およびインスペクターモデル55の損失係数の勾配を基に算出される値を「評価値」と表記する。評価値が、閾値以上である場合には、インスペクターモデル55に入力した運用データセットが、決定境界から大きく離れており、コンセプトドリフトが発生していることを意味する。検知装置は、評価値が閾値以上となった場合に、コンセプトドリフトを検知する。

40

【0043】

検知装置は、ドリフトを検知していない場合には、運用データセットC0を、運用モデル50に入力して、運用データセットC0のデータの属するクラスを予測する。一方、検知装置は、ドリフトを検知した場合には、新たな訓練データセットによって、運用モデル

50

50に対して機械学習を再度実行する。

【0044】

上記のように、本実施形態に係る検知装置は、知識蒸留を用いて、運用モデル50の監視器となるインスペクターモデル55の機械学習を実行する。検知装置は、インスペクターモデル55に、運用データセットC0を入力し、インスペクターモデル55から出力される結果およびインスペクターモデル55の損失係数の勾配を基に算出される評価値と、閾値とを比較し、コンセプトドリフトを検知する。これによって、時間経過に伴い、データの分布が、決定境界から離れる方向に変化しても、かかるデータの分布の変化を検知することができ、運用モデル50の精度劣化を検知することができる。

【0045】

次に、本実施形態に係る検知装置の構成の一例について説明する。図9は、本実施形態に係る検知装置の構成を示す機能ブロック図である。図9に示すように、この検知装置100は、通信部110、入力部120、出力部130、記憶部140、制御部150を有する。

【0046】

通信部110は、ネットワークを介して、外部装置(図示略)とデータ通信を実行する。たとえば、通信部110は、外部装置から後述する訓練データセット141等を受信する。

【0047】

入力部120は、データを入力するための装置またはインタフェースである。たとえば、入力部120は、マウス、およびキーボードなどである。

【0048】

出力部130は、画面を表示するディスプレイなどである。

【0049】

記憶部140は、データや制御部150が実行するプログラムなどを記憶する記憶装置の一例であり、たとえば、ハードディスクやメモリなどである。記憶部140は、訓練データセット141と、運用モデルデータ142、疑似サンプルテーブル143、蒸留データテーブル144、インスペクターモデルデータ145、運用データセットテーブル146とを有する。

【0050】

訓練データセット141は、複数の訓練データを含む。図10は、訓練データセットのデータ構造の一例を示す図である。図10に示すように、この訓練データセットは、レコード番号と、訓練データと、正解ラベルとを対応付ける。レコード番号は、訓練データと、正解ラベルとの組を識別する番号である。訓練データは、メールスパムのデータ、電気需要予測、株価予測、ポーカーハンドのデータ、画像データ等に対応し、複数の次元の特徴量を含む。正解ラベルは、第1クラスまたは第2クラスを一意に識別する情報である。

【0051】

運用モデルデータ142は、運用モデル50(機械学習モデル)のデータである。本実施形態の運用モデル50は、所定の分類アルゴリズムによって、入力データを、複数のクラスに分類する。本実施形態では、運用モデル50を、NNとして説明を行う。

【0052】

図11は、運用モデルの一例を説明するための図である。図11に示すように、運用モデル50は、ニューラルネットワークの構造を有し、入力層50a、隠れ層50b、出力層50cを持つ。入力層50a、隠れ層50b、出力層50cは、複数のノードがエッジで結ばれる構造となっている。隠れ層50b、出力層50cは、活性化関数と呼ばれる関数とバイアス値とを持ち、エッジには、重みが設定される。以下の説明では、バイアス値および重みを「パラメータ」と表記する。

【0053】

入力層50aに含まれる各ノードに、データ(データの特徴量)を入力すると、隠れ層20bを通過して、出力層20cから、各クラスの確率が出力される。

10

20

30

40

50

## 【 0 0 5 4 】

疑似サンプルテーブル 1 4 3 は、訓練データセット 1 4 1 を基に生成される複数の疑似サンプルを保持する。図 1 2 は、疑似サンプルテーブルのデータ構造の一例を示す図である。図 1 2 に示すように、この疑似サンプルテーブル 1 4 3 は、サンプル番号と、疑似サンプルとを対応付ける。サンプル番号は、疑似サンプルを識別する情報である。疑似サンプルは、訓練データの特徴量をスケージングしたデータである。

## 【 0 0 5 5 】

蒸留データテーブル 1 4 4 は、疑似サンプルテーブル 1 4 3 の各疑似サンプルを、運用モデル 5 0 に入力した場合の出力結果（ソフトターゲット）を格納する。図 1 3 は、蒸留データテーブルのデータ構造の一例を示す図である。図 1 3 に示すように、この蒸留データテーブル 1 4 4 は、サンプル番号と、疑似サンプルと、ソフトターゲットとを対応付ける。サンプル番号、疑似サンプルの説明は、図 1 2 で行ったサンプル番号、疑似サンプルの説明と同様である。ソフトターゲットは、疑似サンプルを運用モデル 5 0 に入力した場合の出力結果である。たとえば、ソフトターゲットは、複数のクラスのうち、いずれかのクラスとなる。

10

## 【 0 0 5 6 】

インスペクターモデルデータ 1 4 5 は、インスペクターモデル 5 5 のデータである。インスペクターモデル 5 5 は、図 1 1 で説明した運用モデル 5 0 と同様にして、ニューラルネットワークの構造を有し、入力層、隠れ層、出力層を持つ。インスペクターモデル 5 5 には、パラメータが設定される。インスペクターモデル 5 5 のパラメータは、知識蒸留によって訓練される。

20

## 【 0 0 5 7 】

運用データセットテーブル 1 4 6 は、時間経過に伴って、追加される運用データセットを有する。図 1 4 は、運用データセットテーブルのデータ構造の一例を示す図である。図 1 4 に示すように、運用データセットテーブル 1 4 6 は、データ識別情報と、運用データセットとを有する。データ識別情報は、運用データセットを識別する情報である。運用データセットは、複数の運用データを含む。運用データは、メールスパムのデータ、電気需要予測、株価予測、ポーカーハンドのデータ、画像データ等に対応する。

## 【 0 0 5 8 】

図 9 の説明に戻る。制御部 1 5 0 は、検知装置 1 0 0 全体を司る処理部であり、生成部 1 5 1、算出部 1 5 2、取得部 1 5 3、検知部 1 5 4 を有する。制御部 1 5 0 は、たとえば、プロセッサ等である。

30

## 【 0 0 5 9 】

生成部 1 5 1 は、運用モデルデータ 1 4 2 を生成する処理、疑似サンプルテーブル 1 4 3 を生成する処理、蒸留データテーブル 1 4 4 を生成する処理、インスペクターモデルデータ 1 4 5 を生成する処理を実行する。

## 【 0 0 6 0 】

生成部 1 5 1 が、運用モデルデータ 1 4 2 を生成する処理について説明する。生成部 1 5 1 は、訓練データセット 1 4 1 を入力として、運用モデル 5 0 の機械学習を実行する。たとえば、生成部 1 5 1 は、訓練データセットの訓練データを、運用モデル 5 0 の入力層に入力した場合、出力層の出力結果が、入力した訓練データの正解ラベルに近づくように、運用モデル 5 0 のパラメータを訓練する。たとえば、生成部 1 5 1 は、誤差逆伝播法による機械学習を実行する。生成部 1 5 1 は、機械学習を実行した運用モデル 5 0 のデータ（運用モデルデータ 1 4 2）を、記憶部 1 4 0 に登録する。

40

## 【 0 0 6 1 】

図 1 5 は、本実施形態に係る特徴空間の決定境界を説明するための図である。特徴空間 3 0 は、訓練データセット 1 4 1 の各訓練データを可視化したものある。特徴空間 3 0 の横軸は、第 1 特徴量の軸に対応し、縦軸は、第 2 特徴量の軸に対応する。ここでは説明の便宜上、2 軸で各訓練データを示すが、訓練データは、多次元のデータであるものとする。たとえば、丸印の訓練データに対応する正解ラベルを「第 1 クラス」とし、三角印の訓

50

練データに対応する正解ラベルを「第2クラス」とする。

【0062】

たとえば、訓練データセット141によって、運用モデル50を学習すると、特徴空間30は、決定境界31によって、モデル適用領域31Aと、モデル適用領域31Bとに分類される。たとえば、運用モデル50が、NNである場合、運用モデル50に運用データを入力すると、第1クラスの確率と、第2クラスの確率とが出力される。第1クラスの確率が、第2クラスよりも大きい場合には、データは、第1クラスに分類される。第2クラスの確率が、第1クラスよりも大きい場合には、データは、第2クラスに分類される。

【0063】

図15では、訓練データの正解ラベルが「第1クラス」または「第2クラス」となる場合について説明するが、他のクラスの正解ラベルが付与されてもよい。n種類のクラスが存在する場合には、特徴空間30には、n個のモデル適用領域が設定される。運用モデル50が、NNである場合、運用モデル50に運用データを入力すると、各クラスの確率が出力される。

10

【0064】

生成部151が、疑似サンプルテーブル143を生成する処理について説明する。生成部151は、訓練データセット141に含まれる各訓練データに対してデータ変換を実行する。たとえば、生成部151は、訓練データの各次元の特徴量の値が、0以上、1未満の値に含まれるようなデータ変換(Min-Max Scaling)を実行する。

【0065】

生成部151は、データ変換後の各訓練データのうち、各次元の特徴量の値が、-m以上、1+m未満となる訓練データをランダムに選択する。「m」は、marginであり、あらかじめ任意の実数が設定される。上記処理によって生成部151がランダムに選択したデータ変換後の訓練データを「疑似サンプル」と表記する。たとえば、疑似サンプルの特徴量の値の範囲は、式(1)によって定義される。特徴量の次元数をnとする。

20

【0066】

【数1】

$$[0,1] \in \mathbb{R}^n \quad \dots(1)$$

30

【0067】

生成部151は、サンプル番号と、疑似サンプルと、正解ラベルとを対応付けて、疑似サンプルテーブル143に登録する。疑似サンプルの正解ラベルは、疑似サンプルに対応するデータ変換前の訓練データの正解ラベルとする。

【0068】

生成部151が、蒸留データテーブル144を生成する処理について説明する。生成部151は、疑似サンプルテーブル143の疑似サンプルを、運用モデル50に入力し、運用モデル50の出力結果(ソフトターゲット)を取得する。生成部151は、サンプル番号、疑似サンプル、ソフトターゲットを、蒸留データテーブル144に登録する。

【0069】

生成部151は、疑似サンプルテーブル143の各疑似サンプルについて、上記処理を繰り返し実行することで、ソフトターゲットを取得し、蒸留データテーブル144に登録する。

40

【0070】

たとえば、蒸留データテーブル144に登録される疑似サンプルと、ソフトターゲットとの組の集合を、疑似データ集合とすると、疑似データ集合は、式(2)によって定義される。式(2)の位置a1に示す記号を「Dハット」と表記する。式(2)のa2に示す記号を「xハット」と表記する。Dハットは、疑似データ集合を示す。xハットは、疑似サンプルを示す。f(xハット)は、運用モデル50から出力されるソフトターゲットである。式(2)の位置a3に示す記号を「花文字X」と表記する。花文字Xは、入力空間

50

を示す。

【 0 0 7 1 】

【数 2】

$$\hat{D} = \{\hat{x}, f(\hat{x})\} \quad (\forall \hat{x} \in \mathcal{X}; \mathcal{X} \in \mathbb{R}^d) \quad \dots (2)$$

$\begin{array}{ccc} \uparrow & \uparrow & \uparrow \\ a1 & a2 & a3 \end{array}$

【 0 0 7 2 】

生成部 1 5 1 が、インスペクターモデルデータ 1 4 5 を生成する処理について説明する。生成部 1 5 1 は、蒸留データテーブル 1 4 4 を取得し、蒸留データテーブル 1 4 4 を基  
 10  
 にして、インスペクターモデル 5 5 のパラメータを訓練する。たとえば、生成部 1 5 1 は、蒸留データテーブル 1 4 4 の疑似サンプルを、インスペクターモデル 5 5 の入力層に入力した場合に、出力層の出力結果が、ソフトターゲットに近づくように、インスペクターモデル 5 5 のパラメータを訓練する。たとえば、生成部 1 5 1 は、誤差逆伝播法による機械学習を実行する。生成部 1 5 1 は、機械学習を実行したインスペクターモデル 5 5 のデータ（インスペクターモデルデータ 1 4 5）を、記憶部 1 4 0 に登録する。

【 0 0 7 3 】

上記の生成部 1 5 1 が、インスペクターモデル 5 5 を訓練する処理は、式 ( 3 ) に示す損失関数  $\mathcal{L}^*$  を最小化させるように、インスペクターモデル 5 5 を訓練することに対応する。式 ( 3 ) において、 $f(X; \theta_1)$  は、運用モデル 5 0 の出力に対応し、 $X$  は、D  
 20  
 ハットに対応する。 $\theta_1$  は、運用モデル 5 0 のパラメータを示し、訓練済みのパラメータとなる。 $g(X; \theta_2)$  は、インスペクターモデル 5 5 の出力に対応し、 $X$  は、D ハットに対応する。 $\theta_2$  は、インスペクターモデル 5 5 のパラメータを示し、訓練対象のパラメータとなる。

【 0 0 7 4 】

【数 3】

$$\theta_2^* = \underset{\theta_2}{\operatorname{argmin}} \mathcal{L}(f(X; \theta_1), g(X; \theta_2)) \quad \dots (3)$$

30

【 0 0 7 5 】

生成部 1 5 1 は、後述する検知部 1 5 4 から、ドリフトを検知した旨を受け付けた場合には、運用モデル 5 0 およびインスペクターモデル 5 5 の機械学習を再度実行する。たとえば、生成部 1 5 1 は、最新の訓練データセット 1 4 1 を、外部装置から取得し、最新の訓練データセット 1 4 1 を用いて、運用モデル 5 0 およびインスペクターモデル 5 5 を再度訓練する。

【 0 0 7 6 】

算出部 1 5 2 は、温度付きソフトマックス ( softmax ) を用いて、インスペクターモデル 5 5 の出力をスケールするためのハイパーパラメータを算出する。データ  $i$  を入力した場合の、温度付きソフトマックスを用いたインスペクターモデル 5 5 の出力  $g_i$  は、式  
 40  
 ( 4 ) によって定義される。式 ( 4 ) において、「 $z_i$ 」は、データ  $i$  を入力した際の、インスペクターモデル 5 5 の出力であって、通常ソフトマックスを用いたインスペクターモデル 5 5 の出力を示す。「 $T$ 」は、ハイパーパラメータを示す。以下の算出部 1 5 2 の説明では、温度付きソフトマックスを用いたインスペクターモデル 5 5 の出力を、スコアと表記する。

【 0 0 7 7 】

【数 4】

$$g_i = \frac{\exp(z_i/T)}{\sum_j \exp(z_j/T)} \quad \dots (4)$$

50

## 【 0 0 7 8 】

算出部 1 5 2 は、疑似サンプルテーブル 1 4 3 から、疑似サンプルのペアを選択する。ここでは、疑似サンプルのペアを、第 1 疑似サンプル、第 2 疑似サンプルとする。算出部 1 5 2 は、第 1 疑似サンプルを、温度付きソフトマックスを用いたインスペクターモデル 5 5 に入力して、第 1 スコアを算出する。算出部 1 5 2 は、第 2 疑似サンプルを、温度付きソフトマックスを用いたインスペクターモデル 5 5 に入力して、第 2 スコアを算出する。算出部 1 5 2 は、第 1 スコアと第 2 スコアとの差分の絶対値を、スコア差として算出する。

## 【 0 0 7 9 】

算出部 1 5 2 は、疑似サンプルテーブル 1 4 3 から、異なる疑似サンプルのペアを選択し、スコア差を算出する処理を繰り返し実行する。算出部 1 5 2 は、複数のスコア差のうち、最大のスコア差が閾値  $T_{hs}$  未満となるような、ハイパーパラメータを探索する。閾値  $T_{hs}$  は、あらかじめ設定される。

10

## 【 0 0 8 0 】

図 1 6 は、ハイパーパラメータに応じたスコア差の分布を示す図である。図 1 6 において、グラフ G 1 - 1 は、ハイパーパラメータが  $T = 1$  の場合の、スコア差の頻度を示すヒストグラムである。グラフ G 1 - 1 の横軸はスコア差に対応し、グラフ G 1 - 1 の縦軸は頻度に対応する。グラフ G 1 - 1 に示す例では、スコア差  $0 \sim 0.1$  において、わずかに頻度があり、スコア差  $0.9 \sim 1.0$  において、頻度が集中している。

## 【 0 0 8 1 】

グラフ G 1 - 2 は、ハイパーパラメータが  $T = 1$  の場合の、決定境界からの距離を色によって表現している。グラフ G 1 - 2 上の色が、バー b a の下方向の色に近づくほど、決定境界からの距離が近いことを示す。グラフ G 1 - 2 上の色が、バー b a の上方向の色に近づくほど、決定境界からの距離が遠いことを示す。グラフ G 1 - 2 では、色がほぼ二極化しており、決定境界からの距離を細かく分類できていない。

20

## 【 0 0 8 2 】

グラフ G 2 - 1 は、ハイパーパラメータが  $T = 1.0000$  の場合の、スコア差の頻度を示すヒストグラムである。グラフ G 2 - 1 の横軸はスコア差に対応し、グラフ G 2 - 1 の縦軸は頻度に対応する。グラフ G 2 - 1 に示す例では、グラフ G 1 - 1 と比較して、スコア差の頻度がまんべんなく分布している。

30

## 【 0 0 8 3 】

グラフ G 2 - 2 は、ハイパーパラメータが  $T = 1.0000$  の場合の、決定境界からの距離を色によって表現している。グラフ G 2 - 2 上の色が、バー b a の下方向の色に近づくほど、決定境界からの距離が近いことを示す。グラフ G 2 - 2 上の色が、バー b a の上方向の色に近づくほど、決定境界からの距離が遠いことを示す。グラフ G 2 - 2 では、グラフ G 2 - 1 と比較して、決定境界からの距離を細かく分類できている。

## 【 0 0 8 4 】

算出部 1 5 2 が、最大のスコア差が閾値  $T_{hs}$  未満となるようなハイパーパラメータを探索すると、スコア差と頻度との関係は、グラフ G 2 - 1 の関係に近づき、グラフ G 2 - 2 で示すように、決定境界からの距離を細かく分類できるようになる。算出部 1 5 2 は、算出（探索）したハイパーパラメータの情報を、検知部 1 5 4 に出力する。

40

## 【 0 0 8 5 】

検知部 1 5 4 は、インスペクターモデル 5 5 に運用データを入力して、決定境界からの距離を算出し、決定境界からの距離を基にして、訓練データセット 1 4 1 の分布と、運用データセットの分布との差を検知する。検知部 1 5 4 は、訓練データセット 1 4 1 の分布と、運用データセットの分布との差を、ドリフトとして検知し、ドリフトを検知した旨を、生成部 1 5 1 に出力する。

## 【 0 0 8 6 】

決定境界からの距離は、式 ( 5 ) によって近似することができる。式 ( 5 ) の位置 b 1 に示す記号を単に「 $d \sim$ 」と表記する。

50

【 0 0 8 7 】

【 数 5 】

$$\tilde{d}_{g,x,t,\{i,j\}}^{\uparrow b1} = \frac{|g_i(x_t) - g_j(x_t)|}{\|\nabla_x g_i(x_t) - \nabla_x g_j(x_t)\|_q} \quad \dots (5)$$

【 0 0 8 8 】

検知部 154 は、運用データセットから、運用データのペアとなる運用データ i と、運用データ j とを選択するものとする。式 (5) において、 $g_i(x_t)$  は、運用データ i を入力することで「温度付きソフトマックスを用いたインスペクターモデル 55」から出力される出力結果である。 $g_j(x_t)$  は、運用データ j を入力することで「温度付きソフトマックスを用いたインスペクターモデル 55」から出力される出力結果である。

10

【 0 0 8 9 】

$x g_i(x_t)$  は、偏微分した「温度付きソフトマックスを用いたインスペクターモデル 55」に運用データ i を入力した際の出力結果であり、運用データ i における損失関数の勾配に対応する。 $x g_j(x_t)$  は、偏微分した「温度付きソフトマックスを用いたインスペクターモデル 55」に運用データ j を入力した際の出力結果であり、運用データ j における損失関数の勾配に対応する。

20

【 0 0 9 0 】

式 (5) の分母は、運用データ i における損失関数の勾配と、運用データ j における損失関数の勾配との差分の  $q$ -norm を示す。 $q$ -norm は、 $p$ -norm の dual-norm であり、 $p$  と  $q$  とは式 (6) の関係が成り立つ。たとえば、式 (6) により、 $q = 1$  なら  $p = \infty$ 、 $q = 2$  なら  $p = 2$ 、 $q = \infty$  なら  $p = 1$  となる。

【 0 0 9 1 】

$$1/p + 1/q = 1 \quad \dots (6)$$

【 0 0 9 2 】

$p$ -norm は、式 (7) によって示される。

【 0 0 9 3 】

【 数 6 】

$$\|x\|_p = \sqrt[p]{\sum |x_i|^p} \quad \dots (7)$$

30

【 0 0 9 4 】

たとえば、 $q = 1$  を計算する場合には、 $p = \infty$  を計算すればよく、無限大ノルムは、式 (8) によって示される。なお、 $1$ -norm は、式 (9) によって示され、 $2$ -norm は、式 (10) によって示される。

【 0 0 9 5 】

【 数 7 】

$$\|x\|_\infty = \max_i |x_i| \quad \dots (8)$$

40

【 数 8 】

$$\|x\|_1 = \sum |x_i| \quad \dots (9)$$

【 数 9 】

$$\|x\|_2 = \sqrt{\sum |x_i|^2} \quad \dots (10)$$

50

## 【 0 0 9 6 】

検知部 1 5 4 は、運用データセットから、運用データのペアを再度選択し、再度選択した運用データのペアと、式 ( 5 ) とを基にして、 $d \sim$  を算出する処理を繰り返し実行する。検知部 1 5 4 は、算出した複数の  $d \sim$  の平均値を算出する。複数の  $d \sim$  の平均値は、上述した評価値に対応する。検知部 1 5 4 は、評価値が、閾値  $t_h$  以上となった場合に、訓練データセット 1 4 1 の分布と、運用データセットの分布との差を検知し、ドリフトとして検知する。検知部 1 5 4 は、ドリフトを検知した旨を、外部装置に通知してもよい。

## 【 0 0 9 7 】

検知部 1 5 4 は、運用データセットテーブル 1 4 6 に複数の運用データセットが登録されている場合には、各運用データセットについて、上記処理を繰り返し実行する。

10

## 【 0 0 9 8 】

ところで、検知部 1 5 4 は、次の処理を実行して、閾値  $t_h$  を算出しておく。検知部 1 5 4 は、疑似サンプルテーブル 1 4 3 から、疑似サンプルのペアを選択し、選択した疑似サンプルのペアと、式 ( 5 ) とを基にして、 $d \sim$  を算出する処理を繰り返し実行する。検知部 1 5 4 は、算出した複数の  $d \sim$  を基にして、 $d \sim$  の標準偏差を算出し、算出した標準偏差を、上記の閾値  $T_h$  として設定する。

## 【 0 0 9 9 】

分類部 1 5 5 は、運用データセットの運用データを、運用モデル 5 0 に入力することで、運用データの属するクラスを特定する。分類部 1 5 5 は、運用データセットの他の運用データに対して、上記処理を繰り返し実行することで、複数の運用データを複数のクラスに分類する。

20

## 【 0 1 0 0 】

次に、本実施形態に係る検知装置 1 0 0 の処理手順の一例について説明する。図 1 7 は、本実施形態に係る検知装置の処理手順を示すフローチャートである。たとえば、検知装置 1 0 0 は、運用データセットテーブル 1 4 6 に新たな運用データセットが登録されるたびに、図 1 7 の処理を繰り返し実行する。

## 【 0 1 0 1 】

図 1 7 に示すように、検知装置 1 0 0 の生成部 1 5 1 は、訓練データセット 1 4 1 を入力として、運用モデルの機械学習を実行する ( ステップ S 1 0 1 ) 。

## 【 0 1 0 2 】

生成部 1 5 1 は、訓練データセット 1 4 1 を基にして、疑似サンプルテーブル 1 4 3 を生成する ( ステップ S 1 0 2 ) 。

生成部 1 5 1 は、疑似サンプルテーブル 1 4 3 の疑似サンプルを運用モデルに入力することで、蒸留データテーブル 1 4 4 を生成する ( ステップ S 1 0 3 ) 。

30

## 【 0 1 0 3 】

生成部 1 5 1 は、知識蒸留によって、運用モデル 5 0 を模倣するインスペクターモデル 5 5 の機械学習を実行する ( ステップ S 1 0 4 ) 。

検知装置 1 0 0 の算出部 1 5 2 は、インスペクターモデル 5 5 のハイパーパラメータを算出する ( ステップ S 1 0 5 ) 。

## 【 0 1 0 4 】

検知装置 1 0 0 の取得部 1 5 3 は、運用データセットの運用データのペアをインスペクターモデル 5 5 に入力し、インスペクターモデルの出力結果を取得する ( ステップ S 1 0 6 ) 。

検知装置 1 0 0 の検知部 1 5 4 は、式 ( 5 ) を基にして、ドリフトを検知したか否かを判定する ( ステップ S 1 0 7 ) 。

40

## 【 0 1 0 5 】

検知部 1 5 4 は、ドリフトを検知した場合には ( ステップ S 1 0 8 , Y e s ) 、ステップ S 1 0 1 に移行する。検知部 1 5 4 は、ドリフトを検知していない場合には ( ステップ S 1 0 8 , N o ) 、ステップ S 1 0 9 に移行する。検知装置 1 0 0 の分類部 1 5 5 は、運用データを運用モデル 5 0 に入力し、運用データをクラスに分類する ( ステップ S 1 0 9 ) 。

## 【 0 1 0 6 】

50

次に、本実施形態に係る検知装置 100 の効果について説明する。検知装置 100 は、知識蒸留を用いて、運用モデル 50 の監視器となるインスペクターモデル 55 の機械学習を実行する。検知装置 100 は、インスペクターモデル 55 に、運用データセットを入力し、インスペクターモデル 55 から出力される結果およびインスペクターモデル 55 の損失係数の勾配を基に算出される評価値と、閾値とを比較し、コンセプトドリフトを検知する。これによって、時間経過に伴い、データの分布が、決定境界から離れる方向に変化しても、かかるデータの分布の変化を検知することができ、運用モデル 50 の精度劣化を検知することができる。

#### 【0107】

検知装置 100 は、異なる疑似サンプルのペアを選択し、スコア差を算出する処理を繰り返し実行する。検知装置 100 は、複数のスコア差のうち、最大のスコア差が閾値  $T_h$  未満となるような、ハイパーパラメータを算出する。これによって、決定境界からの距離を、段階的に数値化することができる。

10

#### 【0108】

検知装置 100 は、式 (5) を基にして、評価値を算出して、閾値  $T_h$  との比較を行い、データの分布の変化を検知する。また、検知装置 100 は、訓練データセットを用いて、閾値  $T_h$  を算出しておく。これによって、ドリフトを精度よく検知することができる。

#### 【0109】

検知装置 100 は、ドリフトを検知しない場合に、運用データセットの運用データを運用モデル 50 に入力して、運用データを複数のクラスに分類する。このため、ドリフト発生前の運用モデル 50 によって、運用データを適切に、複数のクラスに分類することができる。

20

#### 【0110】

検知装置 100 は、ドリフトを検知した場合に、新たな訓練データセットによって、運用モデルの機械学習を再度実行する。これによって、ドリフトに対応した運用モデル 50 を再度生成することができる。

#### 【0111】

続いて、同一の訓練データセットを複数種類の機械学習モデルにそれぞれ入力した場合の決定境界の性質について説明する。図 18 は、各機械学習モデルの決定境界の性質を示す図 (1) である。図 18 に示す例では、ある訓練データセット 15 を用いて、サポートベクターマシン (Soft-Margin SVM)、ランダムフォレスト (Random Forest)、NN に対してそれぞれ機械学習を実行した例を示す。

30

#### 【0112】

そうすると、訓練したサポートベクターマシンにデータセットを入力した場合の分布は、分布 20A となり、各データは、決定境界 21A で第 1 クラス、第 2 クラスに分類される。訓練したランダムフォレストにデータセットを入力した場合の分布は、分布 20B となり、各データは、決定境界 21B で第 1 クラス、第 2 クラスに分類される。訓練した NN にデータセットを入力した場合の分布は、分布 20C となり、各データは、決定境界 21C で第 1 クラス、第 2 クラスに分類される。

40

#### 【0113】

図 18 に示すように、同一の訓練データセットで訓練を行った場合でも、機械学習モデルの種類によっては、決定境界の性質が違ってくる。

#### 【0114】

図 19 は、各機械学習モデルの決定境界の性質を示す図 (2) である。図 19 では、ある訓練データセット 35 を用いて、複数種類の機械学習モデルを訓練した例を示す。ここでは、機械学習モデルとして、Nearest Neighbors、RBF SVM、Gaussian Process、Random Forest、Neural Net、Gradient Boosting Tree、Naive Bayes を示す。

#### 【0115】

訓練した Nearest Neighbors にデータセットを入力した場合の分布は、分布 40A と

50

なる。各データは、決定境界 4 1 A で第 1 クラス、第 2 クラスに分類される。訓練した Nearest Neighbors のインスペクターモデルの分布は、分布 4 2 A となり、各データは、決定境界 4 3 A で第 1 クラス、第 2 クラスに分類される。G 4 2 A は、訓練した Nearest Neighbors のインスペクターモデルを基に計算した決定境界からの距離を示す。G 4 2 A では、同じ色の等高線が、同じ距離を示す。インスペクターモデルを NN とする。

【 0 1 1 6 】

訓練した RBF SVM にデータセットを入力した場合の分布は、分布 4 0 B となる。各データは、決定境界 4 1 B で第 1 クラス、第 2 クラスに分類される。訓練した RBF SVM のインスペクターモデルの分布は、分布 4 2 B となり、各データは、決定境界 4 3 B で第 1 クラス、第 2 クラスに分類される。G 4 2 B は、訓練した RBF SVM のインスペクターモデルを基に計算した決定境界からの距離を示す。インスペクターモデルを NN とする。

10

【 0 1 1 7 】

訓練した Gaussian Process にデータセットを入力した場合の分布は、分布 4 0 C となる。各データは、決定境界 4 1 C で第 1 クラス、第 2 クラスに分類される。訓練した Gaussian Process のインスペクターモデルの分布は、分布 4 2 C となり、各データは、決定境界 4 3 C で第 1 クラス、第 2 クラスに分類される。G 4 2 C は、訓練した Gaussian Process のインスペクターモデルを基に計算した決定境界からの距離を示す。インスペクターモデルを NN とする。

【 0 1 1 8 】

訓練した Random Forest にデータセットを入力した場合の分布は、分布 4 0 D となる。各データは、決定境界 4 1 D で第 1 クラス、第 2 クラスに分類される。訓練した Random Forest のインスペクターモデルの分布は、分布 4 2 D となり、各データは、決定境界 4 3 D で第 1 クラス、第 2 クラスに分類される。G 4 2 D は、訓練した Random Forest のインスペクターモデルを基に計算した決定境界からの距離を示す。インスペクターモデルを NN とする。

20

【 0 1 1 9 】

訓練した Neural Net にデータセットを入力した場合の分布は、分布 4 0 E となる。各データは、決定境界 4 1 E で第 1 クラス、第 2 クラスに分類される。訓練した Neural Net のインスペクターモデルの分布は、分布 4 2 E となり、各データは、決定境界 4 3 E で第 1 クラス、第 2 クラスに分類される。G 4 2 E は、訓練した Neural Net のインスペクターモデルを基に計算した決定境界からの距離を示す。インスペクターモデルを NN とする。

30

【 0 1 2 0 】

訓練した Gradient Boosting Tree にデータセットを入力した場合の分布は、分布 4 0 F となる。各データは、決定境界 4 1 F で第 1 クラス、第 2 クラスに分類される。訓練した Gradient Boosting Tree のインスペクターモデルの分布は、分布 4 2 F となり、各データは、決定境界 4 3 F で第 1 クラス、第 2 クラスに分類される。G 4 2 F は、訓練した Gradient Boosting Tree のインスペクターモデルを基に計算した決定境界からの距離を示す。インスペクターモデルを NN とする。

【 0 1 2 1 】

訓練した Naive Bayes にデータセットを入力した場合の分布は、分布 4 0 G となる。各データは、決定境界 4 1 G で第 1 クラス、第 2 クラスに分類される。訓練した Naive Bayes のインスペクターモデルの分布は、分布 4 2 G となり、各データは、決定境界 4 3 G で第 1 クラス、第 2 クラスに分類される。G 4 2 G は、訓練した Naive Bayes のインスペクターモデルを基に計算した決定境界からの距離を示す。インスペクターモデルを NN とする。

40

【 0 1 2 2 】

本実施形態で説明した検知装置 1 0 0 では、機械学習モデルのアーキテクチャにかかわらず、インスペクターモデルを用いて、決定境界からの距離を近似計算することができる。

【 0 1 2 3 】

次に、本実施形態に示した検知装置 1 0 0 と同様の機能を実現するコンピュータのハー

50

ドウェア構成の一例について説明する。図 20 は、本実施形態に係る検知装置と同様の機能を実現するコンピュータのハードウェア構成の一例を示す図である。

【0124】

図 20 に示すように、コンピュータ 200 は、各種演算処理を実行する CPU 201 と、ユーザからのデータの入力を受け付ける入力装置 202 と、ディスプレイ 203 とを有する。また、コンピュータ 200 は、記憶媒体からプログラム等を読み取る読み取り装置 204 と、有線または無線ネットワークを介して、外部装置等との間でデータの授受を行うインタフェース装置 205 とを有する。コンピュータ 200 は、各種情報を一時記憶する RAM 206 と、ハードディスク装置 207 とを有する。そして、各装置 201 ~ 207 は、バス 208 に接続される。

10

【0125】

ハードディスク装置 207 は、生成プログラム 207 a、算出プログラム 207 b、取得プログラム 207 c、検知プログラム 207 d、分類プログラム 207 e を有する。CPU 201 は、生成プログラム 207 a、算出プログラム 207 b、取得プログラム 207 c、検知プログラム 207 d、分類プログラム 207 e を読み出して RAM 206 に展開する。

【0126】

生成プログラム 207 a は、生成プロセス 206 a として機能する。算出プログラム 207 b は、算出プロセス 206 b として機能する。取得プログラム 207 c は、取得プロセス 206 c として機能する。検知プログラム 207 d は、検知プロセス 206 d として機能する。分類プログラム 207 e は、分類プロセス 206 e として機能する。

20

【0127】

生成プロセス 206 a の処理は、生成部 151 の処理に対応する。算出プロセス 206 b の処理は、算出部 152 の処理に対応する。取得プロセス 206 c の処理は、取得部 153 の処理に対応する。検知プロセス 206 d の処理は、検知部 154 の処理に対応する。分類プロセス 206 e の処理は、分類部 155 の処理に対応する。

【0128】

なお、各プログラム 207 a ~ 207 e については、必ずしも最初からハードディスク装置 507 に記憶させておかなくてもよい。例えば、コンピュータ 200 に挿入されるフレキシブルディスク (FD)、CD-ROM、DVD ディスク、光磁気ディスク、IC カードなどの「可搬用の物理媒体」に各プログラムを記憶させておく。そして、コンピュータ 200 が各プログラム 207 a ~ 207 e を読み出して実行するようにしてもよい。

30

【符号の説明】

【0129】

- 100 検知装置
- 110 通信部
- 120 入力部
- 130 出力部
- 140 記憶部
- 141 訓練データセット
- 142 運用モデルデータ
- 143 疑似サンプルテーブル
- 144 蒸留データテーブル
- 145 インспекターモデルデータ
- 146 運用データセットテーブル
- 150 制御部
- 151 生成部
- 152 算出部
- 153 取得部
- 154 検知部

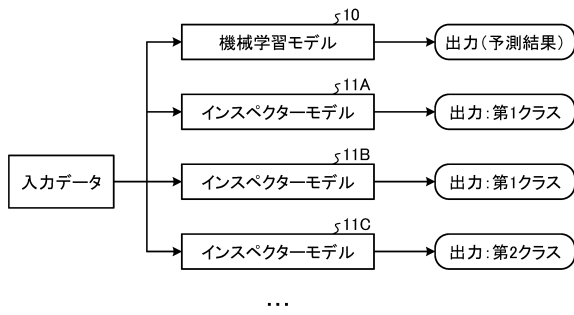
40

50

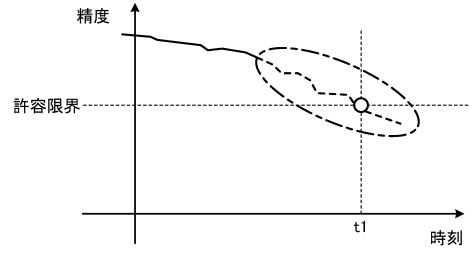
1 5 5 分類部

【図面】

【図 1】

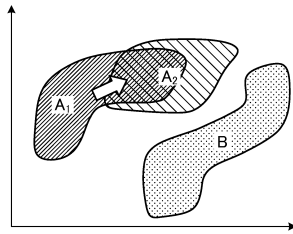


【図 2】

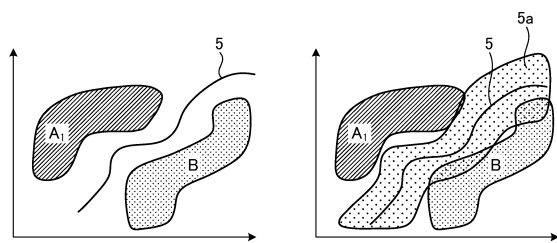


10

【図 3】



【図 4】



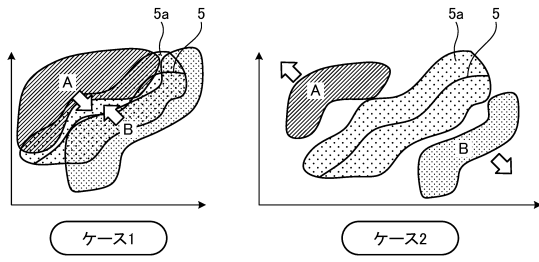
20

30

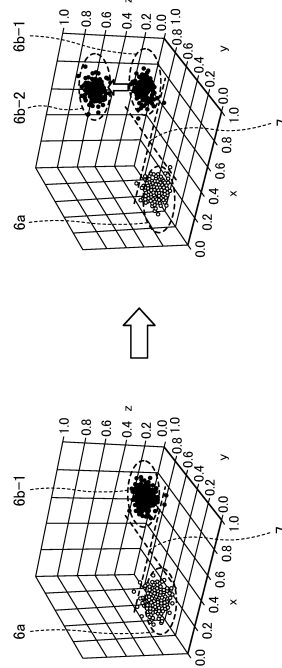
40

50

【図5】



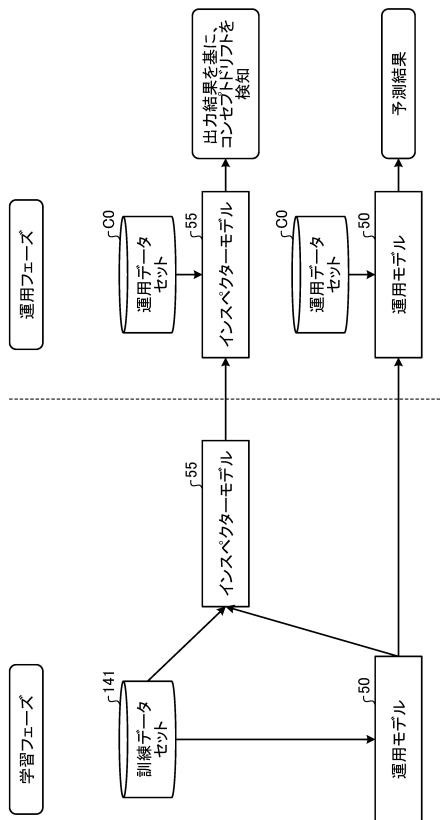
【図6】



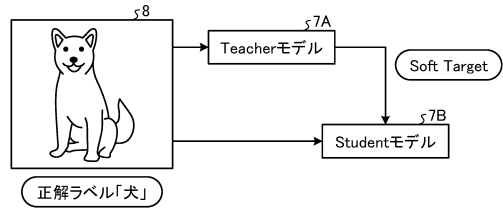
10

20

【図7】



【図8】

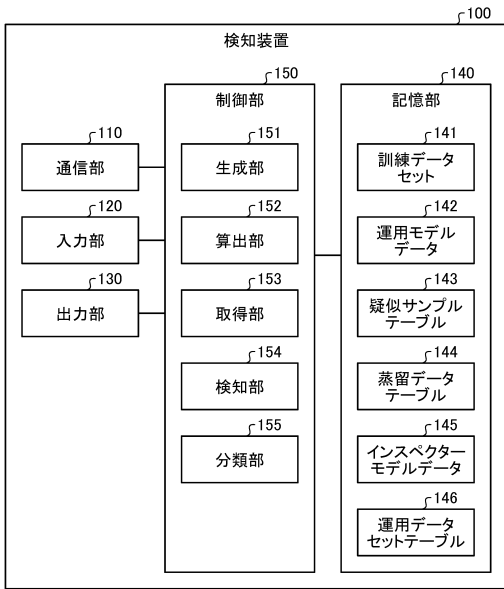


30

40

50

【図 9】



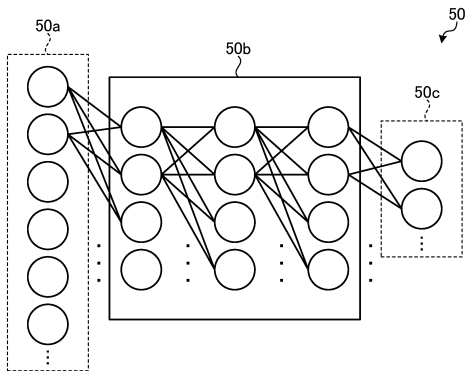
【図 10】

レコード番号	訓練データ	正解ラベル
1001	レコード番号「1001」 訓練データ	第1クラス
1002	レコード番号「1002」 訓練データ	第1クラス
1003	レコード番号「1003」 訓練データ	第1クラス
...	...	...
1050	レコード番号「1050」 訓練データ	第2クラス
1051	レコード番号「1051」 訓練データ	第2クラス
...	...	...

10

20

【図 11】



【図 12】

サンプル番号	疑似サンプル
pe1001	サンプル番号「pe1001」の 疑似サンプル
pe1002	サンプル番号「pe1002」の 疑似サンプル
pe1003	サンプル番号「pe1003」の 疑似サンプル
...	...
pe1050	サンプル番号「pe1050」の 疑似サンプル
pe1051	サンプル番号「pe1051」の 疑似サンプル
...	...

30

40

50

【 図 1 3 】

r144

サンプル番号	疑似サンプル	ソフトターゲット
pe1001	サンプル番号「pe1001」疑似サンプル	サンプル番号「pe1001」疑似サンプルを運用モデルに入力した際の出力結果
pe1002	サンプル番号「pe1002」疑似サンプル	サンプル番号「pe1002」疑似サンプルを運用モデルに入力した際の出力結果
pe1003	サンプル番号「pe1003」疑似サンプル	サンプル番号「pe1003」疑似サンプルを運用モデルに入力した際の出力結果
...	...	...
pe1050	サンプル番号「pe1050」疑似サンプル	サンプル番号「pe1050」疑似サンプルを運用モデルに入力した際の出力結果
pe1051	サンプル番号「pe1051」疑似サンプル	サンプル番号「pe1051」疑似サンプルを運用モデルに入力した際の出力結果
...	...	...

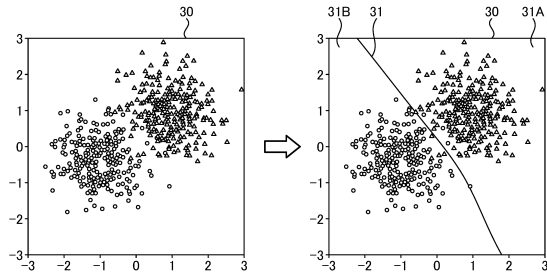
【 図 1 4 】

r146

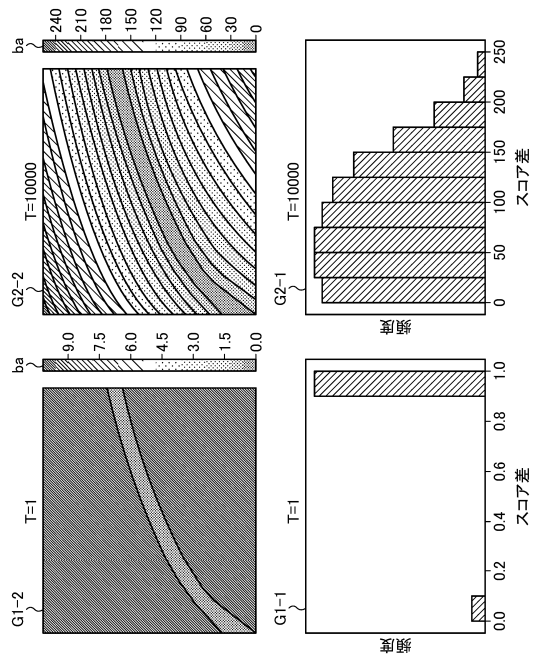
データ識別情報	運用データセット
C0	データ識別情報「C0」の運用データセット
C1	データ識別情報「C1」の運用データセット
C2	データ識別情報「C2」の運用データセット
C3	データ識別情報「C3」の運用データセット

10

【 図 1 5 】



【 図 1 6 】



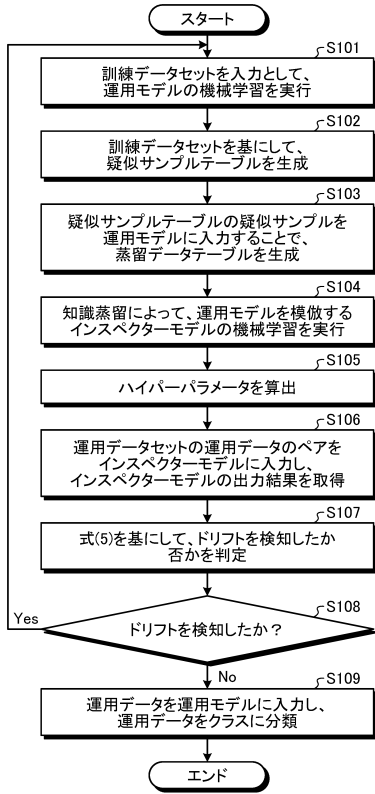
20

30

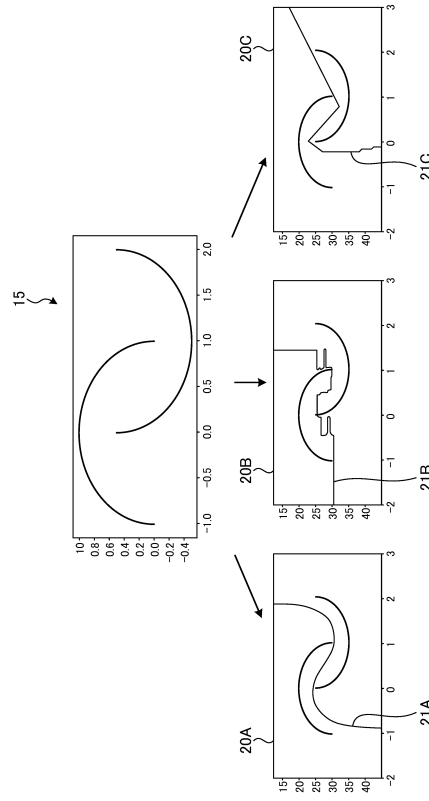
40

50

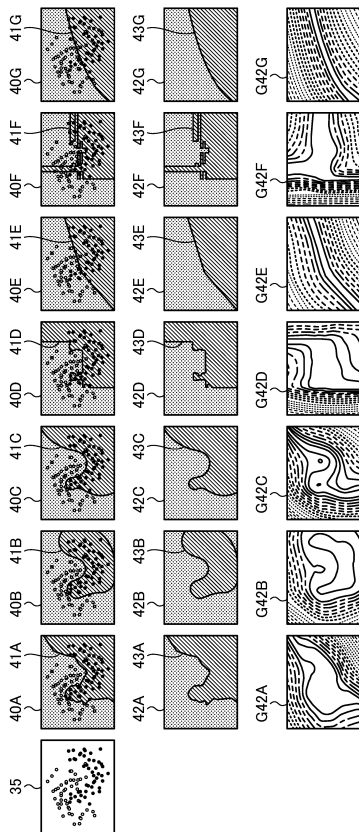
【 図 1 7 】



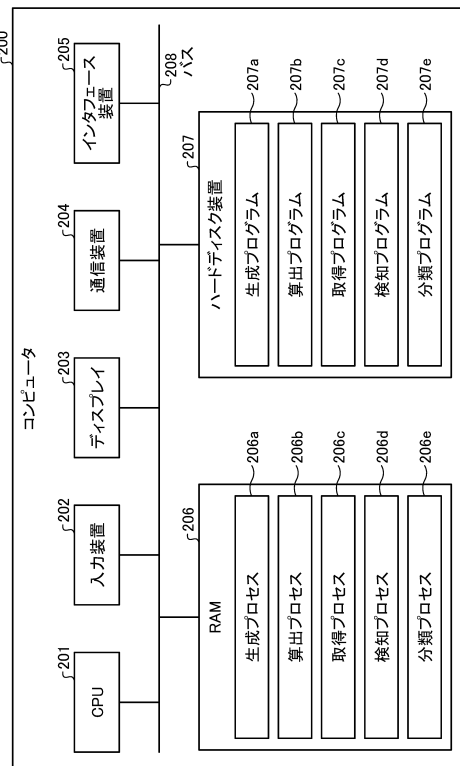
【 図 1 8 】



【 図 1 9 】



【 図 2 0 】



10

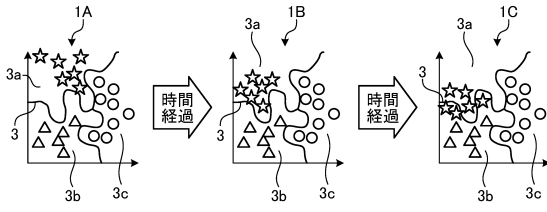
20

30

40

50

【 図 2 1 】



10

20

30

40

50

## フロントページの続き

- (56)参考文献 米国特許出願公開第 2 0 1 7 / 0 3 3 0 1 0 9 ( U S , A 1 )  
 CERQUEIRA, Vitor et al. , "Unsupervised Concept Drift Detection Using a Student-Teacher Approach" , Lecture Notes in Computer Science [online] , スイス , Springer Nature , 2020年10月15日 , volume 12323 , pp. 190-204 , [検索日 2020.11.27], インターネット: URL: [https://link.springer.com/chapter/10.1007/978-3-030-61527-7\\_13](https://link.springer.com/chapter/10.1007/978-3-030-61527-7_13)  
 LU, Jie et al. , "Learning under Concept Drift: A Review" , IEEE Transactions on Knowledge and Data Engineering [online] , 米国 , IEEE , 2019年12月 , vol. 31, no. 12 , pp. 2346-2363 , [検索日 2020.11.27], インターネット: URL: <https://ieeexplore.ieee.org/stamp/stamp.jsp?tp=&number=8496795>  
 石田 勉 ほか , 「ラベルなし運用データに対するコンセプトドリフト検知手法の評価」 , 人工知能学会 第34回全国大会(2020) 論文集 [online] , 一般社団法人 人工知能学会 , 2020年06月 , セッションID: 4Rin1-05, pp. 1-4 , [検索日 2020.11.27], インターネット: URL: [https://www.jstage.jst.go.jp/article/pjsai/JSAI2020/0/JSAI2020\\_4Rin105/\\_pdf/-char/ja](https://www.jstage.jst.go.jp/article/pjsai/JSAI2020/0/JSAI2020_4Rin105/_pdf/-char/ja)  
 石井 遊哉 ほか , 「Knowledge Distillationにおける温度パラメータの適正化に関する検討」 , 情報処理学会研究報告 コンピュータビジョンとイメージメディア ( C V I M ) [online] , 情報処理学会 , 2019年02月28日 , vol. 2019-CVIM-216, no. 10 , pp. 1-5 , [検索日 2019.03.01], インターネット: URL: [https://ipsj.ixsq.nii.ac.jp/ej/?action=repository\\_uri&item\\_id=194906&file\\_id=1&file\\_no=1](https://ipsj.ixsq.nii.ac.jp/ej/?action=repository_uri&item_id=194906&file_id=1&file_no=1)
- (58)調査した分野 (Int.Cl. , D B 名)  
 G 0 6 N 3 / 0 0 - 9 9 / 0 0  
 G 0 6 F 1 6 / 0 0 - 1 6 / 9 5 8  
 G 0 6 F 1 8 / 0 0 - 1 8 / 4 0  
 G 0 6 F 2 1 / 0 0 - 2 1 / 8 8  
 G 0 6 T 7 / 0 0  
 G 0 6 V 1 0 / 7 0 - 1 0 / 8 6