

(12) 特許協力条約に基づいて公開された国際出願

(19) 世界知的所有権機関  
国際事務局

(43) 国際公開日  
2024年2月29日(29.02.2024)



(10) 国際公開番号

WO 2024/042736 A1

(51) 国際特許分類:

G06N 3/0455 (2023.01)

(21) 国際出願番号:

PCT/JP2023/005451

(22) 国際出願日:

2023年2月16日(16.02.2023)

(25) 国際出願の言語:

日本語

(26) 国際公開の言語:

日本語

(30) 優先権データ:

特願 2022-133545 2022年8月24日(24.08.2022) JP

(71) 出願人:株式会社日立製作所(HITACHI, LTD.)

[JP/JP]; 〒1008280 東京都千代田区丸の内一丁目6番6号 Tokyo (JP).

(72) 発明者: 関根 理 敏 (SEKINE, Masatoshi);

〒1008280 東京都千代田区丸の内一丁目6番6号 株式会社日立製作所内 Tokyo (JP).

(74) 代理人: 弁理士法人サンネクスト国際特許事務所(SUNNEXT INTERNATIONAL PATENT OFFICE);

〒1400002 東京都品川区東品川二丁目3番12号 シーフォートスクエア センタービルディング16階 Tokyo (JP).

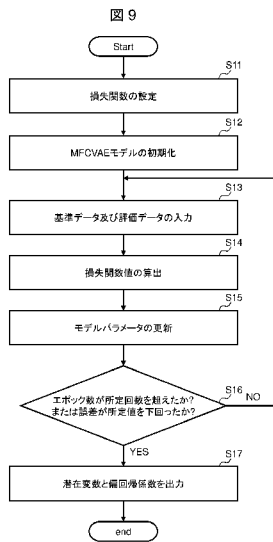
目3番12号 シーフォートスクエア センタービルディング16階 Tokyo (JP).

(81) 指定国(表示のない限り、全ての種類の国内保護が可能): AE, AG, AL, AM, AO, AT, AU, AZ, BA, BB, BG, BH, BN, BR, BW, BY, BZ, CA, CH, CL, CN, CO, CR, CU, CV, CZ, DE, DJ, DK, DM, DO, DZ, EC, EE, EG, ES, FI, GB, GD, GE, GH, GM, GT, HN, HR, HU, ID, IL, IN, IQ, IR, IS, IT, JM, JO, KE, KG, KH, KN, KP, KR, KW, KZ, LA, LC, LK, LR, LS, LU, LY, MA, MD, MG, MK, MN, MW, MX, MY, MZ, NA, NG, NI, NO, NZ, OM, PA, PE, PG, PH, PL, PT, QA, RO, RS, RU, RW, SA, SC, SD, SE, SG, SK, SL, ST, SV, SY, TH, TJ, TM, TN, TR, TT, TZ, UA, UG, US, UZ, VC, VN, WS, ZA, ZM, ZW.

(84) 指定国(表示のない限り、全ての種類の広域保護が可能): ARIPO (BW, CV, GH, GM, KE, LR, LS, MW, MZ, NA, RW, SC, SD, SL, ST, SZ, TZ, UG, ZM, ZW), ユーラシア (AM, AZ, BY, KG, KZ, RU, TJ, TM), ヨーロッパ (AL, AT, BE, BG, CH, CY, CZ, DE, DK, EE, ES, FI, FR, GB, GR, HR, HU, IE, IS,

(54) Title: INFORMATION PROCESSING METHOD, INFORMATION PROCESSING SYSTEM, AND INFORMATION PROCESSING PROGRAM

(54) 発明の名称: 情報処理方法、情報処理システム、及び情報処理プログラム



S11 Set loss function  
S12 Initialize MFCVAE model  
S13 input reference data and assessment data  
S14 Calculate loss function value  
S15 Update model parameter  
S16 Has number of epochs exceeded prescribed number of times or has error become less than prescribed value?  
S17 Output latent variable and partial regression coefficient

(57) Abstract: The present invention increases interpretability of attribute values and attributes of training data and test data for a learning model. Accordingly, regression models, in which latent variables outputted from MFCVAE in response to input of input data including reference data having a plurality of attributes assigned with attribute values are defined as explanatory variables and the attribute values are defined as objective variables, are set for the respective attributes. From the latent variables and the attribute values, prediction values of the attribute values at which prediction errors with respect to the attribute values become minimum and regression coefficients of the regression models are calculated for the respective attributes. On the basis of the prediction values and the regression coefficients for the respective attributes, function values of loss functions, of the MFCVAE, obtained by adding additional terms based on indices which are for the respective attributes and of which the values become smaller as the adaptation of the latent variables and the attribute values to the regression models becomes better are calculated. By error backpropagation based on the function values, a model parameter of the MFCVAE is updated. As described above, model training of MFCVAE is executed.

IT, LT, LU, LV, MC, ME, MK, MT, NL, NO, PL, PT,  
RO, RS, SE, SI, SK, SM, TR), OAPI (BF, BJ, CF,  
CG, CI, CM, GA, GN, GQ, GW, KM, ML, MR, NE,  
SN, TD, TG).

添付公開書類：

一 国際調査報告（条約第21条(3)）

---

(57) 要約：学習モデルの訓練データ及びテストデータの属性と属性値の解釈性を高める。このため、データの複数の属性に属性値が付与されている基準データを含んだ入力データの入力に対してMFCVAEから出力された潜在変数を説明変数とし、属性値を目的変数とする回帰モデルを属性毎に設定する。潜在変数と属性値とから、属性値に対する予測誤差が最小となる属性値の予測値及び回帰モデルの回帰係数を属性毎に算出する。属性毎の予測値及び回帰係数に基づいて、潜在変数及び属性値の回帰モデルへの適合が良いほど小さい値を取る属性毎の指標に基づく追加項を追加したMFCVAEの損失関数の関数値を算出する。関数値に基づく誤差逆伝搬によってMFCVAEのモデルパラメータを更新する。以上のようにMFCVAEのモデル学習を実行する。

## 明 細 書

発明の名称：

情報処理方法、情報処理システム、及び情報処理プログラム

### 技術分野

[0001] 本発明は、情報処理方法、情報処理システム、及び情報処理プログラムに関する。

### 背景技術

[0002] AI (Artificial Intelligence) モデルの品質は、利用するデータの品質に依存する。AI モデルの品質を保証するためには、AI モデルを構築する際の訓練データや推論を行う際のテストデータが有する属性に関する情報を評価することが有用である。そのために、変分オートエンコーダ技術を用いて、学習モデルのエンコーダから潜在変数（特徴量）を抽出し、データの属性情報の内容を明らかにしたり、お互い類似する属性情報をもつデータを抽出したりすること等が行われている。

[0003] 例えば特許文献1では、訓練データから同一のセマンティック特徴に対応する3つの画像を抽出し、3つの画像の各画像について、セマンティック特徴に対応する潜在変数の損失関数を最小化するように変分オートエンコーダのパラメータを更新する。これにより、同一のセマンティック特徴を有する異なる画像の識別性を高めている。

[0004] また例えば非特許文献1では、各潜在変数が入力値に対して与える情報が一意となるように潜在変数の独立性を高めることで、潜在変数の変化に対応する属性の内容やその大きさの変化の解釈性を高めている。その結果、例えば手書き文字データにおいて、ある潜在変数の変化に対して、文字の角度が左斜めから右斜めへ連続変化することが分かる。

[0005] また例えば非特許文献2では、直行する基底の線形結合で潜在変数を表現し、学習によって得られた基底の係数とデータの属性の変化とを対応付けることで、基底の係数の変化に対応する属性の内容やその大きさの変化の解釈

性を高めている。その結果、例えば顔画像データにおいて、ある基底の係数の変化に対して、髪の毛の色が金色から黒色に連続変化することが分かる。

## 先行技術文献

### 特許文献

[0006] 特許文献1：特開2019-75108号公報

### 非特許文献

[0007] 非特許文献1：Shuyang Gao, Rob Brekelmans, Greg Ver Steeg, Aram Galstyan, “Auto-Encoding Total Correlation Explanation, ” Proceedings of the 22nd International Conference on Artificial Intelligence and Statistics pages 1157-1166., [online], Proceedings of Machine Learning Research (PMLR) 2019., [令和4年8月1日検索], インターネット<URL : <https://arxiv.org/abs/1802.05822>>

非特許文献2：Jin-Young Kim, Sung-Bae Cho, “BasisVAE: Orthogonal Latent Space for Deep Disentangled Representation, ” [online], International Conference on Learning Representations (ICLR) 2020., [令和4年8月1日検索], インターネット<URL : <https://arxiv.org/abs/1802.05822>>

## 発明の概要

### 発明が解決しようとする課題

[0008] しかしながら上述の従来技術では、潜在変数と、それに対応する属性や属性値がユーザの解釈に依存し、定性的にしか評価できないため、データの属性と属性値の解釈性が依然として低いという問題があった。

[0009] 本願の開示の一側面では、訓練データやテストデータの潜在変数に対応する属性と属性値の解釈性を高めることを目的とする。

### 課題を解決するための手段

[0010] 本願の開示の一側面では、処理部と記憶部とを有する情報処理システムが実行する情報処理方法であって、前記処理部が、データの複数の属性に属性

値が付与されている基準データを含んだ入力データを、該データの前記複数の属性のそれぞれに関する潜在変数を入力するMF C V A E (Multi-Facet Clustering Variational Auto-Encoder) に入力する第1ステップと、前記入力データの入力に対して前記MF C V A Eから出力された前記潜在変数を説明変数とし、前記属性値を目的変数とする回帰モデルを前記属性毎に設定する第2ステップと、前記潜在変数と前記属性値とから、該属性値に対する予測誤差が最小となる前記属性値の予測値及び前記回帰モデルの回帰係数を前記属性毎に算出する第3ステップと、前記第3ステップによって算出された前記属性毎の前記予測値及び前記回帰係数に基づいて、前記潜在変数及び前記属性値の前記回帰モデルへの適合が良いほど小さい値を取る指標を前記属性毎に算出する第4ステップと、前記MF C V A Eによるデータ再構成の誤差を表す再構成誤差項と、前記潜在変数の分布に制約を与える正則化項と、を有する前記MF C V A Eの損失関数に、前記属性毎の前記指標に基づく追加項を追加した損失関数の関数値を算出する第5ステップと、前記第5ステップによって算出された前記関数値に基づく誤差逆伝搬によって前記MF C V A Eのモデルパラメータを更新する第6ステップと、を実行し、前記第1ステップから前記第6ステップまでを、前記予測誤差又はエポック回数が所定条件を充足するまでこの順序で繰り返すことで前記MF C V A Eのモデル学習を実行する、ことを特徴とする。

### 発明の効果

[0011] 本願の開示の一側面によれば、学習モデルの訓練データやテストデータの属性と属性値の解釈性を高めることができる。前述した以外の課題、構成及び効果は、以下の実施形態の説明により明らかにされる。

### 図面の簡単な説明

[0012] [図1]従来技術 (MF C V A E) の問題点を説明するための図。  
[図2]基準データと評価データ (文字データの場合) を示す図。  
[図3]基準データと評価データ (一般データの場合) を示す図。  
[図4]実施形態1に係る情報処理システムのモデル学習時の動作を説明するた

めの図。

[図5]実施形態1に係る情報処理システムの評価データに対する属性値付与時の動作を説明するための図。

[図6]実施形態1に係る情報処理システムの属性値を指定したデータ生成時の動作を説明するための図。

[図7]実施形態1に係る情報処理システムのモデル学習時の決定係数の推移を説明するための図。

[図8]実施形態1に係る情報処理システムの構成を示すブロック図。

[図9]実施形態1に係る特徴量抽出処理を示すフローチャート。

[図10]実施形態1に係る属性値付与処理を示すフローチャート。

[図11]実施形態1に係るデータ生成処理を示すフローチャート。

[図12]実施形態1に係るデータ品質評価処理を示すフローチャート。

[図13]実施形態1に係る属性及び属性値の出力例1（データに対する属性及び属性値）を示す図。

[図14]実施形態1に係るデータ、属性、及び属性値の出力例2（各属性及び属性値に対するデータ数）を示す図。

[図15]実施形態2に係る情報処理システムの構成を示すブロック図。

[図16]コンピュータのハードウェア構成を示す図。

### 発明を実施するための形態

[0013] 以下、図面を参照して本願の開示に係る実施形態を説明する。実施形態は、図面も含めて本願を説明するための例示である。実施形態では、説明の明確化のため、適宜、省略及び簡略化がされている。特に限定しない限り、実施形態の構成要素は単数でも複数でもよい。また、ある実施形態と他の実施形態を組み合わせた形態も、本願に係る実施形態に含まれる。

[0014] 同一又は類似の構成要素には、同一の符号を付与し、既出に対する後出の実施形態での説明を省略する、又は差分を中心とした説明のみを行う場合がある。また、同一又は類似の構成要素が複数ある場合には、同一の符号に異なる添字を付して説明する場合がある。また、これらの複数の構成要素を区

別する必要がない場合には、添字を省略して説明する場合がある。

[0015] 以下の実施形態では、各種情報をテーブル形式で説明するが、各種情報はテーブル形式以外のデータ形式であってもよい。また、例えば、「XX情報」「XXテーブル」「XXリスト」「XXキュー」等の各種呼称は、これらは互換可能である。例えば「XXテーブル」は、「XXリスト」と呼んでもよい。また、識別情報について説明する際に、「識別情報」「識別子」「名」「ID」「番号」等の表現を用いるが、これらは互換可能である。

[0016] (従来技術の問題点)

実施形態の説明に先立ち、実施形態が前提とする従来技術(MFCVAE: Multi-Facet Clustering Variational Autoencoders)の問題点について説明する。図1は、従来技術の問題点を説明するための図である。MFCVAEは、複数の観点での潜在変数を出力可能な、拡張された変分オートエンコーダ(VAE: Variational Auto-Encoder)である。変分オートエンコーダとは、ニューラルネットワークを使い、潜在変数の空間として確率分布を仮定した生成モデルである。MFCVAEにおける観点とは、MFCVAEが出力する潜在変数(ベクトル)の種類であり、文字データの例では「文字の種類」「字形(太さ、角度等)」等が該当する。

[0017] なお、変分オートエンコーダは、文献1「Diederik P Kingma, Max Welling, “Auto-Encoding Variational Bayes,” May 2014., [令和4年8月1日検索], インターネット<URL: <https://arxiv.org/abs/1312.6114>>」に開示されている。MFCVAEは、文献2「Fabian Falck et.al, “Multi-Facet Clustering Variational Autoencoders, Oct. 2021. [令和4年8月1日検索], インターネット<URL: <https://arxiv.org/abs/2106.05241>>」に開示されている。

[0018] 以下、特徴とは、データを特徴づける情報(属性、属性値、潜在変数等)である。特徴量とは、定量的に表現可能な特徴の値である。属性とは、データを特徴づける性質(文字データの例では「太さ」「傾き」「ノイズ量」「文字の砕け度合い」等)である。属性値とは、属性の度合いを示す値(文字

データの例では属性「太さ」に対する「1 mm」、属性「傾き」に対する「10度」、属性「ノイズ量」に対する「10%」、属性「文字の砕け度合い」に対する「レベル2」等)である。属性値は、連続値でも離散値でもよい。潜在変数とは、変分オートエンコーダ関連技術において、エンコーダから出力される特徴量である。変分オートエンコーダ関連技術とは、VAEやMFCVAEを含む変分ベイズアルゴリズムを有する変分オートエンコーダ技術全般を指す。

[0019] MFCVAEは、一つの潜在変数の変化に対して、複数の属性の属性値が変化する。このため、ある潜在変数の変化に対して変化する属性の対応付けが困難であった。図1を参照して、手書き文字の場合について属性及び属性値を例に説明する。

[0020] 図1に示すように、潜在変数1を横軸、潜在変数2を縦軸に取った座標系において、潜在変数1及び潜在変数2のグループ101は、属性「太さ」が細い文字に対応した潜在変数のグループである。グループ102は、属性「太さ」が中程度の文字に対応した潜在変数のグループである。グループ103は、属性「太さ」が太い文字に対応した潜在変数のグループである。グループ104は、属性「ノイズ量」が少ない文字に対応した潜在変数のグループである。グループ105は、属性「ノイズ量」が中程度の文字に対応した潜在変数のグループである。グループ106は、属性「ノイズ量」が多い文字に対応した潜在変数のグループである。ここで属性「太さ」の「太い」「中程度」「細い」、及び属性「ノイズ量」の「多い」「中程度」「少ない」は例示的表現に過ぎず、定量的な表現又はこれに付与したラベルの一例である。

[0021] 図1に示す例では、潜在変数の変化に対して、全ての属性値が一様に増加あるいは減少するものではない。よって、ある潜在変数の変化に対する属性値の変化の対応付けが困難である。例えば図1では、潜在変数1の値が小さいグループ101からグループ103、グループ102へと二段階にわたり増加すると、属性「太さ」の属性値が「細い」、「太い」、「中程度」と変

化する。しかし、潜在変数 1 の値の様な増加に対して、属性値の変化が一樣でない。潜在変数 1 の変化に関する属性「ノイズ量」の属性値の変化も同様である。その結果、ある潜在変数の変化に対する属性値の変化の解釈がしにくい。

[0022] 一方、潜在変数 2 の値が小さいグループ 103 からグループ 102、グループ 101 へと二段階にわたり増加すると、属性「太さ」の属性値が「太い」、「中程度」、「細い」と変化するように、潜在変数 1 の値の様な増加に対して、属性値が一樣に減少する。潜在変数 2 に関する属性「ノイズ量」の属性値の変化も同様である。その結果、潜在変数の変化に対して、属性値が一樣に減少するため、ある潜在変数の変化に対する属性値の変化の解釈がしやすい。

[0023] (基準データと評価データ)

先ず、文字データの基準データと評価データを説明する。図 2 は、基準データと評価データ（文字データの場合）を示す図である。図 2 において、各行をデータという。それぞれのデータに対して、「文字の種類」「太さ」「傾き」といった各属性について、属性値が格納されているデータと格納されていないデータがある。「文字の種類」は該当の文字のイメージデータである。

[0024] ある属性において、属性値が格納されているデータが、その属性に関する基準データであり、属性値が格納されていないデータが、その属性に関する評価データである。基準データは、評価データの属性値を求めるために用いられる属性値が既知のデータである。評価データは、属性値が未知の属性に関して属性値を求めて付与される対象のデータである。

[0025] 文字データの場合、手書き文字データに対して、一般に「A」や「B」といった文字の種類属性値を付与するのは容易であるが、「太さ」や「傾き」などの属性値を付与するのは容易ではない。そこでゴシック体といった「太さ」や「傾き」といった属性の属性値が変更可能である活字を用いて、変分オートエンコーダのモデルを学習する。

[0026] 図2の例では「太さ」や「傾き」等の属性値を持つ「データ属性」“活字”の「データ番号」“1”“2”“3”のデータが基準データ、属性値を持たない「データ属性」“手書き文字”「データ番号」“4”“5”のデータが評価データとなる。より一般には、基準データは活字及び手書き文字を含み、評価データは手書き文字を含む。

[0027] 次に、一般データの基準データと評価データを説明する。図3は、基準データと評価データ（一般データの場合）を示す図である。図2では、データ属性毎に基準データと評価データが分かれていた。図3では基準データと評価データを一般化し、それぞれのデータに対して、各属性に関して属性値が格納されているデータが該当の属性に関する基準データであり、各属性に関して属性値が格納されていないデータが該当の属性に関する評価データである。

[0028] 「属性1」に関して、基準データは「データ番号」“1”“2”“3”“4”のデータであり、評価データは「データ番号」“5”のデータである。同様に、「属性2」に関して、基準データは「データ番号」“1”“2”“5”のデータであり、評価データは「データ番号」“3”“4”のデータである。

[0029] 以下の実施形態の目的は、各属性において、基準データの属性値を利用して評価データの属性値を推定することである。

[0030] [実施形態1]

(モデル学習時及びデータ再構成時の動作)

図4は、実施形態1に係る情報処理システム1のモデル学習時の動作を説明するための図である。本実施形態では、属性に対して属性値が予め付与された基準データ201を用い、基準データ201の潜在変数を基に属性値を重回帰モデルで予測できるように、変分オートエンコーダ関連技術において新規の損失関数を用いてモデル学習を実行する。

[0031] 情報処理システム1は、MFCVAE2を有する。MFCVAE2は、エンコーダ203とデコーダ205とを含んで構成される。

- [0032] MFCVAE2のモデル学習時には、エンコーダ203に入力されるデータ（訓練データセット）は、基準データ201のみを含むか、あるいは基準データ201及び評価データ202の両方を含む。エンコーダ203は、MFCVAE2の中間出力である潜在変数204を出力する。MFCVAE2は、潜在変数204を説明変数とし、属性値の正解値（正解ラベル）208を目的変数とする重回帰モデルを設定する。なお重回帰モデルに限らず、線形回帰モデル及び非線形回帰モデルの何れでもよい。重回帰モデルは、モデル計算の負荷が少ないという利点があることから、本実施形態では重回帰モデルを採用する。
- [0033] 情報処理システム1は、潜在変数204と属性値の正解値208とから、属性値の予測値と正解値との平均二乗誤差を最小化する偏回帰係数209及びその時の属性値の予測値210を求める。属性値の予測値210は、偏回帰係数209を各係数とする潜在変数204の一次結合として算出される。属性値の正解値208と属性値の予測値210とに基づいて、重回帰モデルへの適合（重回帰モデルの当てはまり）が良いほど小さい値を取る適合度を表す指標となり得る決定係数211や予測誤差212が求められる。決定係数211や予測誤差212は、重回帰モデルへの適合度を表す指標を含む損失関数213で学習される。このようにしてエンコーダ203が学習される。
- [0034] 潜在変数204は、デコーダ205に入力される。デコーダ205は、デコーダ205によって再構成された基準データ201である再構成基準データ206と、デコーダ205によって再構成された評価データ202である再構成評価データ207とを出力する。再構成評価データ207は、属性値が付与されたデータとなっている。
- [0035] ここで、偏回帰係数209、損失関数213、決定係数211、予測誤差212を説明する。
- [0036] 従来技術のMFCVAEでは、目的関数である変分下限（Evidence Lower Bound（ELBO））は、式（1）のように表される。従来手法のMFCV

A Eでは、式（1）の変分下限の符号をマイナスにした負の損失関数が最小化されるようにMFCVAEモデルのパラメータが学習される（上述の文献2参照）。式（1）において“D”は訓練データセット、“x”が訓練データセットに含まれる訓練データ、“z→”は潜在変数、“θ”はエンコーダのパラメータ、“φ”はデコーダのパラメータ、“KL（A|B）”は分布Aと分布BのKLダイバージェンスを表す。

[0037] [数1]

$$\begin{aligned} \mathcal{L}^{MFCVAE}(D; \theta, \phi) \\ = \mathbb{E}_{x \sim D} \left[ \mathbb{E}_{q_\phi(\bar{z}|x)} \log p_\theta(x|\bar{z}) - \sum_{j=1}^J \left[ \mathbb{E}_{q_\phi(c_j|x)} KL(q_\phi(z_j|x)|p_\theta(z_j|c_j)) + KL(q_\phi(c_j|x)|p_\theta(c_j)) \right] \right] \cdots(1) \end{aligned}$$

[0038] これに対して本実施形態では、観点jでの自由度調整済み決定係数を $R_{f,j}^2$ 、各自由度調整済み決定係数 $R_{f,j}^2$ の重み係数を $\gamma_j$  ( $>0$ )として、式（2）のように目的関数を設定する。式（2）の右辺は、式（1）の右辺の期待値E[\*]のカッコ内の式に第3項 $\gamma_j R_{f,j}^2$ が追加されたものである。

[0039] なお、式（2）の第3項では、自由度調整済み決定係数 $R_{f,j}^2$ ではなく、後述の決定係数 $R_j^2$ が採用されてもよい。決定係数 $R_j^2$ 、自由度調整済み決定係数 $R_{f,j}^2$ は、決定係数211の一例である。

[0040] [数2]

$$\begin{aligned} \mathcal{L}^{Proposed}(D; \theta, \phi) \\ = \mathbb{E}_{x \sim D} \left[ \mathbb{E}_{q_\phi(\bar{z}|x)} \log p_\theta(x|\bar{z}) - \sum_{j=1}^J \left[ \mathbb{E}_{q_\phi(c_j|x)} KL(q_\phi(z_j|x)|p_\theta(z_j|c_j)) + KL(q_\phi(c_j|x)|p_\theta(c_j)) \right] \right. \\ \left. + \gamma_j R_{f,j}^2 \right] \cdots(2) \end{aligned}$$

[0041] ここで、式（2）の期待値E[\*]のカッコ内の第1項は、MFCVAEによるデータ再構成の誤差を表す再構成誤差項である。式（2）の期待値E[\*]のカッコ内の第2項は、MFCVAEの潜在変数のバラつきを抑制する等の潜在変数の分布に制約を与える正則化項である。式（2）の期待値E[\*]のカッコ内の第3項は、属性毎の潜在変数及び属性値に対する重回帰モデルへの適合が良いほど小さい値を取る指標に基づく追加項である。

[0042] 損失関数  $L_{loss}$  (損失関数  $L_{loss}$  と表す) は負の目的関数であるから、式 (2) の目的関数を用いて式 (3) のように表される。

[0043] [数3]

$$Loss = -\mathcal{L}^{Proposed}(D; \theta, \phi) \cdots (3)$$

[0044] 本実施形態では、損失関数  $L_{loss}$  が最小化、すなわち決定係数  $R_{f,i}^2$  が最大化されるように、MFCVAEモデルのパラメータが学習される。

[0045] 次に、観点  $j$  の決定係数  $R_j^2$ 、及び観点  $j$  の自由度調整済み決定係数  $R_{f,j}^2$  の算出方法を説明する。

[0046] データ数が  $N$  個、データの属性の種類が  $j = 1, 2, \dots, J$  の  $J$  個、ある属性  $j$  における潜在変数の次元数を  $K_j$  とする。またデータの各属性がMFCVAEの  $J$  個の各観点と一対一に対応しているものとする。インデックス番号  $n$  であるデータ  $n$  に属性  $j$  の属性値が付与されていれば、データ  $n$  は属性  $j$  において基準データである。一方データ  $n$  に属性  $j$  の属性値が付与されていなければ、データ  $n$  は属性  $j$  において評価データである。

[0047] ある属性  $j$  に関する基準データのインデックスの集合を  $B_j$ 、集合  $B_j$  の要素数を  $M_j$  とする。集合  $B_j = \{b_{j,1}, b_{j,2}, \dots, b_{j,M_j}\}$  とする。ある属性  $j$  に関する基準データ  $n$  における潜在変数を  $z_{n,j} = \{z_{n,j,1}, z_{n,j,2}, \dots, z_{n,j,K_j}\}$ 、属性値の正解値を  $y_{n,j}$ 、潜在変数  $z_{n,j}$  を説明変数とする。また、属性値を目的変数とした重回帰モデルの偏回帰係数を  $w_j = \{w_{j,0}, w_{j,1}, w_{j,2}, \dots, w_{j,K_j}\}^T$ 、属性値の予測値を  $\hat{y}_{n,j}$  とする。属性値の予測値  $\hat{y}_{n,j}$  は、式 (4) のように表される。

[0048] [数4]

$$\begin{aligned} \hat{y}_{n,j} &= w_{j,0} + w_{j,1}z_{n,j,1} + w_{j,2}z_{n,j,2} \cdots + w_{j,K_j}z_{n,j,K_j} \\ &= (1, z_{n,j,1}, z_{n,j,2}, \dots, z_{n,j,K_j}) (w_{j,0}, w_{j,1}, \dots, w_{j,K_j})^T = Z_{n,j}w_j \cdots (4) \end{aligned}$$

[0049] ただし、式 (5) のように潜在変数ベクトル  $Z_{n,j}$  を定義した。

[数5]

$$Z_{n,j} = (1, z_{n,j,1}, z_{n,j,2}, \dots, z_{n,j,K_j}) \cdots (5)$$

[0050] ある属性  $j$  に関する属性値  $y_{n,j}$  (ただし  $n \in B_j$ ) と重回帰モデルによる属性値の予測値  $\hat{y}_{n,j}$  との予測誤差である平均二乗誤差  $MSE_j$  は、式 (6) のように表される。平均二乗誤差  $MSE_j$  は、予測誤差 212 の一例である。

[数6]

$$MSE_j = \frac{1}{M_j} \sum_{n \in B_j} (\hat{y}_{n,j} - y_{n,j})^2 = \frac{1}{M_j} \sum_{n \in B_j} (z_{n,j} w_j - y_{n,j})^2 \cdots (6)$$

[0051] ここで平均二乗誤差  $MSE_j$  を最小化する  $w_j$  は、式 (6) の右辺を  $w_j$  で偏微分してゼロとおく ( $\nabla w_j = 0$ ) ことで、式 (7) のように、偏回帰係数  $w_j$  は、潜在変数  $Z_j$  と属性値  $y_j$  の関数となる。偏回帰係数  $w_j$  は、偏回帰係数 209 の一例である。

[数7]

$$w_j = (Z_j^T Z_j)^T Z_j^T y_j \cdots (7)$$

[0052] ただし式 (7) において、潜在変数  $Z_j$  と属性値  $y_j$  を、式 (8) と式 (9) に示すようにおいた。

[数8]

$$Z_j = \{z_{b_{j,1},j}, z_{b_{j,2},j}, \dots, z_{b_{j,M_j},j}\}^T \cdots (8)$$

[数9]

$$y_j = \{y_{b_{j,1},j}, y_{b_{j,2},j}, \dots, y_{b_{j,M_j},j}\}^T \cdots (9)$$

[0053] なお、ある属性  $j$  に関する属性値  $y_{n,j}$  (ただし  $n \in B_j$ ) と重回帰モデルによる属性値の予測値  $\hat{y}_{n,j}$  との予測誤差は、平均二乗誤差に限らず、平均誤差、平均絶対誤差、平均平方二乗誤差、平均誤差率、平均絶対誤差率等を採用することもできる。

[0054] ある属性  $j$  に関する決定係数  $R_j^2$  は、説明変数が目的変数をどれくらい説明しているかを表す。決定係数  $R_j^2$  は、属性値の平均値  $\bar{y}_{n,j}$  を用いて式 (

10) のように表される。

[数10]

$$R_j^2 = 1 - \frac{\sum_{n \in B_j} (y_{n,j} - \hat{y}_{n,j})^2}{\sum_{n \in B_j} (y_{n,j} - \bar{y}_{n,j})^2} \dots (10)$$

[0055] また決定係数は、説明変数の数が増えるほど1に近づくという性質を持っているため、説明変数の数が多い場合には、この性質を補正した自由度調整済み決定係数  $R_{f,j}^2$  が採用されてもよい。自由度調整済み決定係数  $R_{f,j}^2$  は、説明変数の数を  $p$  とし、基準データのサンプル数は  $M_j$  であるので、式(11)のように表される。

[数11]

$$R_{f,j}^2 = 1 - \frac{\frac{\sum_{n \in B_j} (y_{n,j} - \hat{y}_{n,j})^2}{M_j - p - 1}}{\frac{\sum_{n \in B_j} (y_{n,j} - \bar{y}_{n,j})^2}{M_j - 1}} \dots (11)$$

[0056] なお、重み係数  $\gamma_j$  は、再構成誤差項、正則化項、及び決定係数  $R_j^2$  ( $R_{f,j}^2$ ) の絶対値の比較から求めることができる。具体的には、決定係数  $R_j^2$  に対する重み係数  $\gamma_j$  は、 $|\gamma_j R_j^2|$  のオーダーが式(2)の右辺の期待値  $E[*]$  のカッコ内の式の再構成誤差項と正則化項の各絶対値のオーダーと同じになるように定められる。同様に、自由度調整済み決定係数  $R_{f,j}^2$  に対する重み係数  $\gamma_j$  も、 $|\gamma_j R_{f,j}^2|$  のオーダーが、式(2)の右辺の期待値  $E[*]$  のカッコ内の式の再構成誤差項と正則化項の各絶対値のオーダーと同じになるように定められる。

[0057] 基準データ201及び評価データ202の再構成時には、デコーダ205は、基準データ201及び評価データ202の属性及び属性値(評価データ202の場合は付与された属性値)と、エンコーダ203の学習の最後のエポックで得た偏回帰係数209を用いる。そして、基準データ201及び評価データ202の属性及び属性値と、偏回帰係数209とを用いて、式(4)から、潜在変数(潜在変数ベクトル  $Z_{n,j}$ )を算出する。そして、デコーダ205は、算出した潜在変数を入力として、入力された基準データ201及

び評価データ202をそれぞれ再構成した再構成基準データ206及び再構成評価データ207を出力する。

[0058] (属性値付与時の動作)

図5は、実施形態1に係る情報処理システム1の評価データ202に対する属性値付与時の動作を説明するための図である。情報処理システム1は、評価データ202への属性値付与時には、先ず評価データ202を学習済みのエンコーダ203に入力し、潜在変数204を得る。情報処理システム1は、MFCVAE2(図4)の学習の最終エポックで得た偏回帰係数209を用い、属性値の予測値210を、偏回帰係数209を各係数とする潜在変数204の一次結合式で算出する。情報処理システム1は、属性値の予測値210を評価データ202に付与する。

[0059] (データ生成時の動作)

図6は、実施形態1に係る情報処理システム1の属性値を指定したデータ生成時の動作を説明するための図である。データ生成とは、変分オートエンコーダ関連技術において、潜在変数を入力として、デコーダからデータを出力することをいう。情報処理システム1は、生成させたい属性値401を持つデータ405の生成時には、ユーザが生成させたい属性及び属性値401と、MFCVAE2(図4)の学習の最後で得た偏回帰係数209とから潜在変数204を算出する。そして情報処理システム1は、算出した潜在変数204をデコーダ205に入力することで、生成させたい属性値401を持つデータ405を生成する。

[0060] なお、データ生成の際に、指定された属性及び属性値401に該当する基準データ201が存在する場合には、この基準データ201に対応するデータを再構成したデータ405として採用する。指定された属性及び属性値401に該当する基準データ201が存在しない場合に、生成させたい属性及び属性値401と、偏回帰係数209とから潜在変数204を算出する。そして、算出した潜在変数204をデコーダ205に入力することで、生成させたい属性値401を持つデータ405を生成する。

[0061] (モデル学習時の決定係数の推移)

図7は、実施形態1に係る情報処理システム1のモデル学習時の決定係数の推移を説明するための図である。図7のグラフでは、潜在変数を横軸、属性値を縦軸に取り、属性値の実際の値を点で表し、属性値の予測値を直線で表している。情報処理システム1は、MFCVAE2の損失関数に、潜在変数を説明変数、属性値を目的関数とする重回帰モデルの決定係数を含む追加項を追加し、決定係数が高くなるようにMFCVAE2を学習させる。その結果、学習の初期では決定係数は低い(図7(a))が、学習のエポック数が進行して学習の中期(図7(b))、後期(図7(c))と推移するに従って、決定係数は高くなり、潜在変数に基づく属性値の予測精度が高くなる。

[0062] (実施形態1に係る情報処理システム1の構成)

図8は、実施形態1に係る情報処理システム1の構成を示すブロック図である。情報処理システム1は、データ記憶部602、特徴量抽出部603、属性値付与部608、データ生成部614、及びデータ品質評価部612を有する。

[0063] データ記憶部602は、メモリ又はストレージであり、基準データ201と評価データ202の入力を受け付け、蓄積する。データ記憶部602は、情報処理システム1に含まれる装置であっても、情報処理システム1の外部装置であっても何れでもよい。

[0064] 特徴量抽出部603は、MFCVAE2のデータ記憶部602に格納されている基準データ201を元に、MFCVAE2のモデル学習を実行する。また特徴量抽出部603は、MFCVAE2のデータ記憶部602に格納されている評価データ202の属性推定を行う。また特徴量抽出部603は、属性値を指定したデータ生成を行う。特徴量抽出部603は、回帰モデル適合度評価部604、損失算出部605、モデル更新部606、及びエンコーダ部607を有する。特徴量抽出部603の処理機能は、図9を参照して後述する。

- [0065] 属性値付与部608は、評価データ202の属性推定を行い、評価データ202の属性及び属性値611を出力する。属性値付与部608は、属性値推定部609と、属性及び属性値出力部610とを有する。属性値付与部608の処理機能は、図10を参照して後述する。
- [0066] データ生成部614は、属性値付与部608によって出力された対象データの属性及び属性値611を入力として、属性値を指定したデータ生成を行い、生成したデータ405を出力する。データ生成部614は、潜在変数算出部615、デコード部616、及びデータ出力部617を有する。データ生成部614の処理機能は、図11を参照して後述する。
- [0067] データ品質評価部612は、属性値付与部608によって出力された対象データ（基準データ201、評価データ202）の属性及び属性値611に基づいて対象データの品質評価を行い、データ品質評価結果613を出力する。
- [0068] データ品質評価部612は、評価の対象データの属性及び属性値611を用いて、一例として、下記のような観点で対象データの品質評価を行う。品質評価については“機械学習品質マネジメントガイドライン”、国立研究開発法人産業技術総合研究所、[令和4年8月1日検索]、インターネット<URL : [https://www.aist.go.jp/aist\\_j/press\\_release/pr2020/pr20200630\\_2/pr20200630\\_2.html](https://www.aist.go.jp/aist_j/press_release/pr2020/pr20200630_2/pr20200630_2.html)>を参照すればよい。
- (1) データ設計の十分性：データを用いる対象のシステムが対応すべき様々な状況に対して十分な訓練データやテストデータを確保していること。
  - (2) データセットの被覆性：基準を定めて網羅したそれぞれのケースに対してそれぞれのケースに対応する入力の可能性に対して抜け漏れなく、レアケース及び通常ケースそれぞれに正しく推論できる学習に必要な十分な量のデータが与えられていること。
  - (3) データの均一性：全体として推論性能の期待値を最大化するように、訓練データを偏り無く用意すること。
- [0069] なお、データ品質評価部612の処理機能は、図12を参照して後述する

。

[0070] また、特徴量抽出部603、属性値付与部608、データ生成部614、及びデータ品質評価部612は、1つのコンピュータ上に実現されていてもよいし、異なるコンピュータ上に実現されてもよく、これらの統合分散の形態は適宜変更可能である。

[0071] (実施形態1に係る特徴量抽出処理)

図9は、実施形態1に係る特徴量抽出処理を示すフローチャートである。特徴量抽出処理は、特徴量抽出部603(図8)によって、ユーザ指示を契機として実行される。

[0072] 先ずステップS11では、回帰モデル適合度評価部604は、基準データ201の潜在変数を説明変数とし、属性値を目的変数としたMFCVAEモデル(本実施形態では重回帰モデル)の当てはまりの良さを示す指標を損失関数 $L_{oss}$ に設定する。損失関数 $L_{oss}$ に設定される指標は、本実施形態では、観点 $j$ での自由度調整済み決定係数 $R_{f, j}^2$ である。

[0073] 次にステップS12では、回帰モデル適合度評価部604は、MFCVAEモデルの初期化を行う。次にステップS13では、回帰モデル適合度評価部604は、基準データ201及び評価データ202をMFCVAEモデルに入力する。ステップS13では、回帰モデル適合度評価部604は、少なくとも基準データ201をMFCVAEモデルに入力すればよい。

[0074] 次にステップS14では、損失算出部605は、式(3)に基づいて損失関数 $L_{oss}$ の関数値を算出する。回帰モデル適合度評価部604は、損失算出部605による損失関数 $L_{oss}$ の関数値の算出の前段階として、次の処理を行う。すなわち、回帰モデル適合度評価部604は、ステップS13の入力データの入力に対してMFCVAEモデルから出力された潜在変数を説明変数とし、属性値を目的変数とする重回帰モデルを属性毎に設定する。次に、回帰モデル適合度評価部604は、潜在変数と属性値とから、属性値に対する予測誤差が最小となる属性値の予測値及び重回帰モデルの回帰係数を属性毎に算出する。次に、回帰モデル適合度評価部604は、算出された

属性毎の予測値及び回帰係数に基づいて、潜在変数及び属性値に対する重回帰モデルへの適合が良いほど小さい値を取る指標を属性毎に算出する。その後、損失算出部605は、ステップS14で、損失関数Lossの関数値を算出する。

[0075] なお、入力データが基準データ201及び評価データ202を含む場合、ステップS14では、回帰モデル適合度評価部604は、損失関数Lossの追加項を、基準データ201を用いて計算する。一方、損失算出部605は、再構成誤差項及び正則化項を、基準データ201及び評価データ202の何れか一方又は両方を用いて計算する。これは、損失関数Lossの追加項は、潜在変数と属性値との重回帰モデルへの適合度に基づくことから、属性値を含む基準データのみ損失関数Lossの追加項を計算可能なためである。

[0076] 次にステップS15では、モデル更新部606は、誤差逆伝搬によりMFCVAEモデルのパラメータを更新する。次にステップS16では、モデル更新部606は、所定条件（エポック数が所定回数を超えたか、又はMFCVAEモデルによる推定値と実際の値の誤差が所定値を下回った）が充足されたかを判定する。モデル更新部606は、所定条件が充足された場合（ステップS16YES）にステップS17に処理を移し、所定条件が充足されていない場合（ステップS16NO）にステップS13に処理を戻す。

[0077] ステップS17では、エンコーダ部607は、学習済みのMFCVAEモデルのエンコーダ203に評価データ202を入力し、潜在変数204と偏回帰係数209を出力する。

[0078] （実施形態1に係る属性値付与処理）

図10は、実施形態1に係る属性値付与処理を示すフローチャートである。属性値付与処理は、属性値付与部608（図8）によって、ユーザ指示を契機として実行される。

[0079] 先ずステップS21では、属性値推定部609は、特徴量抽出部603（エンコーダ203）から得られた潜在変数204と偏回帰係数209から、

評価データ 202 の属性値の予測値 210 を算出する。次にステップ S 22 では、属性及び属性値出力部 610 は、基準データ 201 と、ステップ S 21 で属性値の予測値 210 が算出された評価データ 202 との属性及び属性値から、属性毎の各属性値の出現頻度のヒストグラムを求める（後述の図 14 参照）。そして属性及び属性値出力部 610 は、このヒストグラムをもとに各属性値が基準データ 201 及び評価データ 202 の各属性において出現する確率をデータ含有率として求め、結果を出力する。

[0080]（実施形態 1 に係るデータ生成処理）

図 11 は、実施形態 1 に係るデータ生成処理を示すフローチャートである。データ生成処理は、データ生成部 614（図 8）によって、ユーザ指示を契機として実行される。

[0081] 先ずステップ S 31 では、データ生成部 614 は、ユーザによる生成させたい属性及び属性値 401 の入力を受け付ける。次にステップ S 32 では、潜在変数算出部 615 は、ステップ S 31 で入力を受け付けた属性及び属性値 401 と偏回帰係数 209 から、潜在変数 204 を計算して出力する。次にステップ S 33 では、デコード部 616（デコーダ 205）は、ステップ S 32 で計算された潜在変数 204 を入力として生成させたい属性及び属性値 401 を持つデータ 405（例えば文字データ）を再構成する。次にステップ S 34 では、データ出力部 617 は、デコード部 616（デコーダ 205）によって再構成されたデータ 405 を出力する。

[0082] なお、潜在変数算出部 615 は、ステップ S 31 で入力を受け付けた属性及び属性値 401 に該当する基準データ 201 が存在する場合には、ステップ S 32 をスキップし、ステップ S 33 でこの基準データ 201 に対応するデータを再構成したデータ 405 とする。

[0083]（実施形態 1 に係るデータ品質評価処理）

図 12 は、実施形態 1 に係るデータ品質評価処理を示すフローチャートである。データ品質評価処理は、データ品質評価部 612（図 8）によって、ユーザ指示を契機として実行される。

[0084] ステップS 4 1では、データ品質評価部6 1 2は、属性値付与部6 0 8によって出力された属性及び属性値6 1 1に関して、例えば上述の(1) データ設計の充分性、(2) データセットの被覆性、(3) データセットの均一性の少なくとも一つの観点で評価する。次にステップS 4 2では、データ品質評価部6 1 2は、ステップS 4 1のデータ品質評価結果6 1 3を出力する。

[0085] (属性及び属性値の出力例1)

図1 3は、属性及び属性値の出力例1 (データに対する属性及び属性値)を示す図である。図1 3は、属性値付与部6 0 8の属性値推定部6 0 9 (図8)によって、例えば図2又は図3に示す属性値が付与されていなかったデータに属性値が付与され、属性及び属性値出力部6 1 0によって出力されたものである。

[0086] (属性及び属性値の出力例2)

図1 4は、データ、属性、及び属性値の出力例2 (各属性及び属性値に対するデータ数)を示す図である。図1 4は、図1 3の表示方法を変えた出力例である。図1 4は、属性値付与部6 0 8の属性及び属性値出力部6 1 0 (図8)によって出力される属性毎の属性値のヒストグラムである。この表示によって、属性毎に例えば上述の(2) データセットの被覆性や(3) データの均一性を確認できる。(2) データセットの被覆性は、図1 4のヒストグラムの各属性の属性値が所定の広い範囲に分布しかつ各度数が何れも所定数以上であることで充足されると考えられる。(3) データの均一性は、図1 4のヒストグラムの各属性の属性値が所定の広い範囲に均等に分布していることで充足されると考えられる。このような分析によって、属性値に対して不足しているデータを確認することが可能となる。

[0087] 例えば図1 4のヒストグラム1 1 0 1は、属性1の度数分布を示す。ヒストグラム1 1 0 1は、ヒストグラム1 1 0 2、1 1 0 3と比較してデータの分布範囲が広い又は同等であるが、この分布範囲に存在しない属性値がある。この点でヒストグラム1 1 0 1は、(2) データセットの被覆性が充足さ

れていないと言える。またヒストグラム1101は、属性値の分布が均一でない。属性値の分布の均一性は、属性値の分散や標準偏差といったバラつきを表す統計値に基づいて判断できる。この点でヒストグラム1101は、(3)データの均一性が充足されていないと言える。

[0088] また図14のヒストグラム1102は、属性2の度数分布を示す。ヒストグラム1102は、ヒストグラム1101、1103と比較してデータの分布範囲が狭く、この分布範囲に存在しない属性値がある。この点でヒストグラム1102は、(2)データセットの被覆性が充足されていないと言える。またヒストグラム1102は、属性値の分布が均一でない。この点でヒストグラム1102は、(3)データの均一性が充足されていないと言える。

[0089] また図14のヒストグラム1103は、属性Jの度数分布を示す。ヒストグラム1103は、ヒストグラム1101、1102と比較してデータの分布範囲が広い又は同等であるが、この分布範囲に存在しない属性値がある。この点でヒストグラム1103は、(2)データセットの被覆性が充足されていないと言える。またヒストグラム1103は、ヒストグラム1101、1102と比較して属性値の分布が均一でない。この点でヒストグラム1102は、(3)データの均一性が充足されていないと言える。

[0090] なお、図14の各グラフを、「データ数」に代えて「各属性値が基準データ201及び評価データ202の各属性において出現するデータ含有率」を縦軸とするグラフとしてもよい。

[0091] [実施形態1の効果]

本実施形態では、訓練データやテストデータの属性をユーザが明示的に指定し、定量的な属性値で表すことにより、ユーザにとって解釈性の高い属性分析が可能となる。このため、訓練データやテストデータの中に不足しているデータや、誤判別の多いデータの特徴を発見しやすい。

[0092] また本実施形態では、従来技術のように、得られた潜在変数の持つ属性をユーザが解釈する（偏在変数は太さ又は角度に依存する等）ではなく、ユーザが潜在変数に持たせるべき属性を明示的に指定できるため、ユーザの意図

に従った属性分析が可能となる。

[0093] また、従来技術では属性が定性的にしか分からないため、異なるデータセットや異なるモデルで学習したデータ間の属性は比較できなかった。しかし、本実施形態では属性値が定量的に求まるため、異なるデータセットや異なるモデルで学習したデータ間の属性の比較が可能となる。

[0094] また、本実施形態では、属性値の予測に限定した場合であっても、教師あり学習による回帰モデルを使って属性値を推定する場合よりも少ないデータ量又は学習量で属性値を付与できる。

[0095] また、本実施形態では、データ生成では、複数の属性に対する属性値を指定してデータ生成できるため、必要とされるデータを容易に生成できる。また、データ生成の際に、指定された属性及び属性値に該当する基準データが存在する場合にはこの基準データに対応するデータを再構成したデータとして採用する。これにより、属性及び属性値と偏回帰係数から潜在変数を算出しデコーダでデコードしてデータを再構成する場合と比較して、速やかにデータを再構成できる。

[0096] [実施形態2]

実施形態1では、1つの情報処理システム1を用いてモデル学習(図4、図9)、属性値付与及び属性と属性値の関係の出力(図5、図10)、及びデータ品質評価処理(図12)を実行する例を示した。しかし、モデル学習、属性値付与、及びデータ品質評価処理は、図15に示す複数の情報処理システム1(1-1, 1-2, ..., 1-n)で並列に実行されてもよい。

[0097] 例えばモデル学習を複数の情報処理システム1で実行する場合、ステップS13~S16(図9)を、複数の情報処理システム1毎にそれぞれ異なる基準データ201を含む入力データを用いて実行してもよい。そして、複数の情報処理システム1の少なくとも1つが、各情報処理システム1によって得られたMFCVAEモデルの学習結果をマージして出力する。

[0098] また複数の情報処理システム1毎に得られたMFCVAEのモデルのそれぞれの学習結果に基づいて、情報処理システム1毎に各入力データへの属性

値付与及び属性と属性値の関係の出力（ステップS 2 1～S 2 3（図1 0））を実行してもよい。そして、複数の情報処理システム1の少なくとも1つが、各情報処理システム1によって得られた属性と属性値の関係（図1 4）をマージして出力する。

[0099] 本実施形態では、従来技術と比較して、潜在変数が定量的に求まるため、モデル毎に別システムで並列処理しても計算結果をマージできることから、複数システムでモデル学習、属性値付与、及び属性と属性値の関係の出力の各処理の負荷分散が可能となる。よって、従来と比較して短い時間で、これらの処理を完了させ、必要とされるデータを生成することができる。

[0100] （実施形態の適用例）

実施形態は、上述のように手書き文字の文字認識に適用できる。その他、実施形態は、全てのデータにラベル付与するのが困難であり、一部のデータのみ正確な属性値（ラベル）が付与されており、残りのデータに属性値（ラベル）を付与したいといったケースであれば適用できる。

[0101] 例えば、工場設備の振動データに対する回転数のラベルの付与がある。前提として過去に取得した工場設備の振動データには回転数のラベルが付与されておらず、新たに回転数が計測できる装置を導入し、過去に取得した工場設備のデータに対して回転数のラベルを付与するような場合である。

[0102] また、画像における被写体の角度予測がある。角度のラベル付与された少量のデータから、未知の画像の被写体の角度を予測する場合である。この適用例は、ロボットが物をつかむときの把持の方向制御等に利用できる。

[0103] また、楽曲の印象評価を行う場合がある。予めユーザが評価した楽曲の印象（楽しい、悲しい、うれしい、寂しいなど）から、未知の楽曲の印象のラベル付与を行うことができる。

[0104] また、学会論文の研究分野の可視化を行う場合がある。予め各分野との関係度（画像認識分野との関連度が3 0、強化学習分野との関連度が5 0、・・・）が分かっている論文を基に、未知の論文の各分野との関係度を推定する場合である。

[0105] (コンピュータ1000のハードウェア)

図16は、コンピュータ1000の構成を示すハードウェア図である。例えば、情報処理システム1、あるいは特徴量抽出部603、属性値付与部608、データ生成部614、及びデータ品質評価部612等の情報処理システム1を適宜分散した各システムは、コンピュータ1000によって実現される。

[0106] コンピュータ1000は、バス等の内部通信線1009を介して相互に接続されたCPUをはじめとするプロセッサ1001、主記憶装置1002、補助記憶装置1003、ネットワークインタフェース1004、入力装置1005、及び出力装置1006を備える。

[0107] プロセッサ1001は、コンピュータ1000全体の動作制御を司る。また主記憶装置1002は、例えば揮発性の半導体メモリから構成され、プロセッサ1001のワークメモリとして利用される。補助記憶装置1003は、非一時的記憶媒体の一例であり、ハードディスク装置、SSD (Solid State Drive)、又はフラッシュメモリ等の大容量の不揮発性の記憶装置から構成され、各種プログラムやデータを長期間保持するために利用される。

[0108] 補助記憶装置1003に格納された実行可能プログラム1100がコンピュータ1000の起動時や必要時に主記憶装置1002にロードされ、主記憶装置1002にロードされた実行可能プログラム1100をプロセッサ1001が実行する。これにより、各種処理を実行するシステムが実現される。

[0109] なお、実行可能プログラム1100は、非一時的記録媒体に記録され、媒体読み取り装置によって非一時的記録媒体から読み出されて、主記憶装置1002にロードされてもよい。または、実行可能プログラム1100は、ネットワークを介して外部のコンピュータから取得されて、主記憶装置1002にロードされてもよい。

[0110] ネットワークインタフェース1004は、コンピュータ1000をシステム内の各ネットワークに接続する、あるいは他のコンピュータと通信するた

めのインタフェース装置である。ネットワークインタフェース1004は、例えば、有線LAN (Local Area Network) や無線LAN等のNIC (Network Interface Card) から構成される。

[0111] 入力装置1005は、キーボードや、マウス等のポインティングデバイス等から構成され、ユーザがコンピュータ1000に各種指示や情報を入力するために利用される。出力装置1006は、例えば、液晶ディスプレイ又は有機EL (Electro Luminescence) ディスプレイ等の表示装置や、スピーカ等の音声出力装置から構成され、必要時に必要な情報をユーザに提示するために利用される。

[0112] なお、本発明は前述した実施形態に限定されるものではなく、添付した特許請求の範囲の趣旨内における様々な変形例及び同等の構成が含まれる。例えば、前述した実施形態は本発明を分かりやすく説明するために詳細に説明したものであり、必ずしも説明した全ての構成を備えるものに本発明は限定されない。また、ある実施形態の構成の一部を他の実施形態の構成に置き換えてもよい。また、ある実施形態の構成に他の実施形態の構成を加えてもよい。また、各実施形態の構成の一部について、他の構成の追加、削除、又は置換をしてもよい。

[0113] また、前述した各構成、機能、処理部、処理手段等は、それらの一部又は全部を、例えば集積回路で設計する等により、ハードウェアで実現してもよい。あるいは、プロセッサがそれぞれの機能を実現するプログラムを解釈し実行することにより、ソフトウェアで実現してもよい。

[0114] 各機能を実現するプログラム、テーブル、ファイル等の情報は、メモリ、ハードディスク、SSD (Solid State Drive) 等の記憶装置、又は、IC (Integrated Circuit) カード、SDカード、DVD (Digital Versatile Disc) の非一時的記録媒体に格納することができる。

[0115] また、制御線や情報線は説明上必要と考えられるものを示しており、実装上必要な全ての制御線や情報線を示しているとは限らない。実際には、ほとんど全ての構成が相互に接続されていると考えてよい。

## 符号の説明

[0116] 1 : 情報処理システム、201 : 基準データ、202 : 評価データ、204 : 潜在変数、603 : 特徴量抽出部、608 : 属性値付与部、612 : データ品質評価部、614 : データ生成部、1000 : コンピュータ。

## 請求の範囲

- [請求項1] 処理部と記憶部とを有する情報処理システムが実行する情報処理方法であって、
- 前記処理部が、
- データの複数の属性に属性値が付与されている基準データを含んだ入力データを、該データの前記複数の属性のそれぞれに関する潜在変数を入力するMFCVAE (Multi-Facet Clustering Variational Auto-Encoder) に入力する第1ステップと、
- 前記入力データの入力に対して前記MFCVAEから出力された前記潜在変数を説明変数とし、前記属性値を目的変数とする回帰モデルを前記属性毎に設定する第2ステップと、
- 前記潜在変数と前記属性値とから、該属性値に対する予測誤差が最小となる前記属性値の予測値及び前記回帰モデルの回帰係数を前記属性毎に算出する第3ステップと、
- 前記第3ステップによって算出された前記属性毎の前記予測値及び前記回帰係数に基づいて、前記潜在変数及び前記属性値の前記回帰モデルへの適合が良いほど小さい値を取る指標を前記属性毎に算出する第4ステップと、
- 前記MFCVAEによるデータ再構成の誤差を表す再構成誤差項と、前記潜在変数の分布に制約を与える正則化項と、を有する前記MFCVAEの損失関数に、前記属性毎の前記指標に基づく追加項を追加した損失関数の関数値を算出する第5ステップと、
- 前記第5ステップによって算出された前記関数値に基づく誤差逆伝搬によって前記MFCVAEのモデルパラメータを更新する第6ステップと、を実行し、
- 前記第1ステップから前記第6ステップまでを、前記予測誤差又はエポック回数が所定条件を充足するまでこの順序で繰り返すことで前記MFCVAEのモデル学習を実行する、ことを特徴とする情報処理

方法。

- [請求項2] 請求項1に記載の情報処理方法であって、  
前記回帰モデルは、重回帰モデルである、ことを特徴とする情報処理方法。
- [請求項3] 請求項1に記載の情報処理方法であって、  
前記属性毎の前記指標は、前記回帰モデルの決定係数である、ことを特徴とする情報処理方法。
- [請求項4] 請求項1に記載の情報処理方法であって、  
前記属性毎の前記指標は、前記予測誤差である、ことを特徴とする情報処理方法。
- [請求項5] 請求項4に記載の情報処理方法であって、  
前記予測誤差は、平均二乗誤差である、ことを特徴とする情報処理方法。
- [請求項6] 請求項1に記載の情報処理方法であって、  
前追加項は、前記属性毎の前記指標に、前記属性毎の重み係数を乗算した項であり、  
前記処理部が、  
前記属性毎の前記重み係数を、前記属性毎に、前記指標と前記再構成誤差項及び前記正則化項との各絶対値のオーダーが等しくなるように決定する、ことを特徴とする情報処理方法。
- [請求項7] 請求項1に記載の情報処理方法であって、  
前記入力データは、前記基準データ及び前記属性に前記属性値が付与されていない評価データを含み、  
前記処理部が、前記第5ステップにおいて、  
前記追加項を、前記基準データを用いて計算し、  
前記再構成誤差項及び前記正則化項を、前記基準データ及び前記評価データの何れか一方又は両方を用いて計算する、ことを特徴とする情報処理方法。

- [請求項8] 請求項1に記載の情報処理方法であって、  
前記入力データは、前記基準データ及び前記属性に前記属性値が付与されていない評価データを含み、  
前記処理部が、  
前記第1ステップから前記第6ステップまでを繰り返すことでモデル学習済みの前記MFCVAEに前記評価データを入力し、前記評価データに関する前記潜在変数を取得する第7ステップと、  
前記第7ステップによって取得された前記潜在変数と、前記MFCVAEのモデル学習の最終エポック時における前記回帰係数とに基づいて、前記評価データの前記属性値が付与されていない前記属性の前記属性値の予測値を算出して該評価データに該属性値として付与する第8ステップと、を実行することを特徴とする情報処理方法。
- [請求項9] 請求項8に記載の情報処理方法であって、  
前記処理部が、  
前記基準データと前記第8ステップによって前記予測値が付与された前記評価データとの前記属性及び前記属性値に関する情報を出力する第9ステップ、を実行することを特徴とする情報処理方法。
- [請求項10] 請求項8に記載の情報処理方法であって、  
前記処理部が、  
前記基準データと前記第8ステップによって前記予測値が付与された前記評価データとを用いて、データの設計の十分性、データの被覆性、又はデータの均一性を含む観点に従って前記入力データを評価する第10ステップ、を実行することを特徴とする情報処理方法。
- [請求項11] 請求項1に記載の情報処理方法であって、  
複数の前記情報処理システムの各前記処理部が、  
前記第1ステップから前記第6ステップまでを、前記予測誤差又はエポック回数が所定条件を充足するまで繰り返すことで前記MFCVAEをモデル学習することを、それぞれ異なる前記入力データを用い

て実行し、

各前記処理部によって得られた前記M F C V A Eのモデルの学習結果をマージして出力する第11ステップ、を実行することを特徴とする情報処理方法。

[請求項12]

請求項8に記載の情報処理方法であって、

複数の前記情報処理システムの各前記処理部が、

前記第1ステップから前記第8ステップまでを、それぞれ異なる前記入力データを用いて実行し、

各前記処理部によって得られた、前記基準データと前記第8ステップによって前記予測値が付与された前記評価データの前記属性及び前記属性値に関する情報をマージして出力する第12ステップ、を実行することを特徴とする情報処理方法。

[請求項13]

請求項1に記載の情報処理方法であって、

前記処理部が、

指定された前記属性及び前記属性値を、前記第1ステップから前記第6ステップまでを繰り返すことでモデル学習済みの前記M F C V A Eに入力し、入力された前記属性及び前記属性値と前記回帰係数とから前記潜在変数を算出し、該潜在変数を基に、入力された前記属性及び前記属性値に対応する前記データを再構成する第13ステップ、を実行することを特徴とする情報処理方法。

[請求項14]

請求項13に記載の情報処理方法であって、

前記処理部が、

前記指定された前記属性及び前記属性値に該当する前記基準データが存在する場合には、該基準データに対応する前記データを再構成したデータとして採用し、

前記指定された前記属性及び前記属性値に該当する前記基準データが存在しない場合に、前記第13ステップを実行する、ことを特徴とする情報処理方法。

- [請求項15] 請求項1に記載の情報処理方法であって、  
前記基準データは活字及び手書き文字を含み、前記評価データは手書き文字を含む、ことを特徴とする情報処理方法。
- [請求項16] データの複数の属性に属性値が付与されている基準データを含んだ入力データを、該データの前記複数の属性のそれぞれに関する潜在変数を入力するMF C V A E (Multi-Facet Clustering Variational Auto-Encoder) に入力し、  
前記入力データの入力に対して前記MF C V A E から出力された前記潜在変数を説明変数とし、前記属性値を目的変数とする回帰モデルを前記属性毎に設定し、  
前記潜在変数と前記属性値とから、該属性値に対する予測誤差が最小となる前記属性値の予測値及び前記回帰モデルの回帰係数を前記属性毎に算出し、  
算出された前記属性毎の前記予測値及び前記回帰係数に基づいて、前記潜在変数及び前記属性値の前記回帰モデルへの適合が良いほど小さい値を取る指標を前記属性毎に算出する回帰モデル適合度評価部と、  
前記MF C V A E によるデータ再構成の誤差を表す再構成誤差項と、前記潜在変数の分布に制約を与える正則化項と、を有する前記MF C V A E の損失関数に、前記属性毎の前記指標に基づく追加項を追加した損失関数の関数値を算出する損失算出部と、  
前記損失算出部によって算出された前記関数値に基づく誤差逆伝搬によって前記MF C V A E のモデルパラメータを更新するモデル更新部と、を有し、  
前記回帰モデル適合度評価部、前記損失算出部、及び前記モデル更新部は、前記予測誤差又はエポック回数が所定条件を充足するまでこの順序で処理を順次繰り返すことで前記MF C V A E のモデル学習を実行する、ことを特徴とする情報処理システム。

[請求項17] コンピュータを情報処理システムとして機能させるための情報処理プログラムであって、

前記コンピュータを、

データの複数の属性に属性値が付与されている基準データを含んだ入力データを、該データの前記複数の属性のそれぞれに関する潜在変数を入力するMF C V A E (Multi-Facet Clustering Variational Auto-Encoder) に入力し、

前記入力データの入力に対して前記MF C V A E から出力された前記潜在変数を説明変数とし、前記属性値を目的変数とする回帰モデルを前記属性毎に設定し、

前記潜在変数と前記属性値とから、該属性値に対する予測誤差が最小となる前記属性値の予測値及び前記回帰モデルの回帰係数を前記属性毎に算出し、

算出された前記属性毎の前記予測値及び前記回帰係数に基づいて、前記潜在変数及び前記属性値の前記回帰モデルへの適合が良いほど小さい値を取る指標を前記属性毎に算出する回帰モデル適合度評価部と、

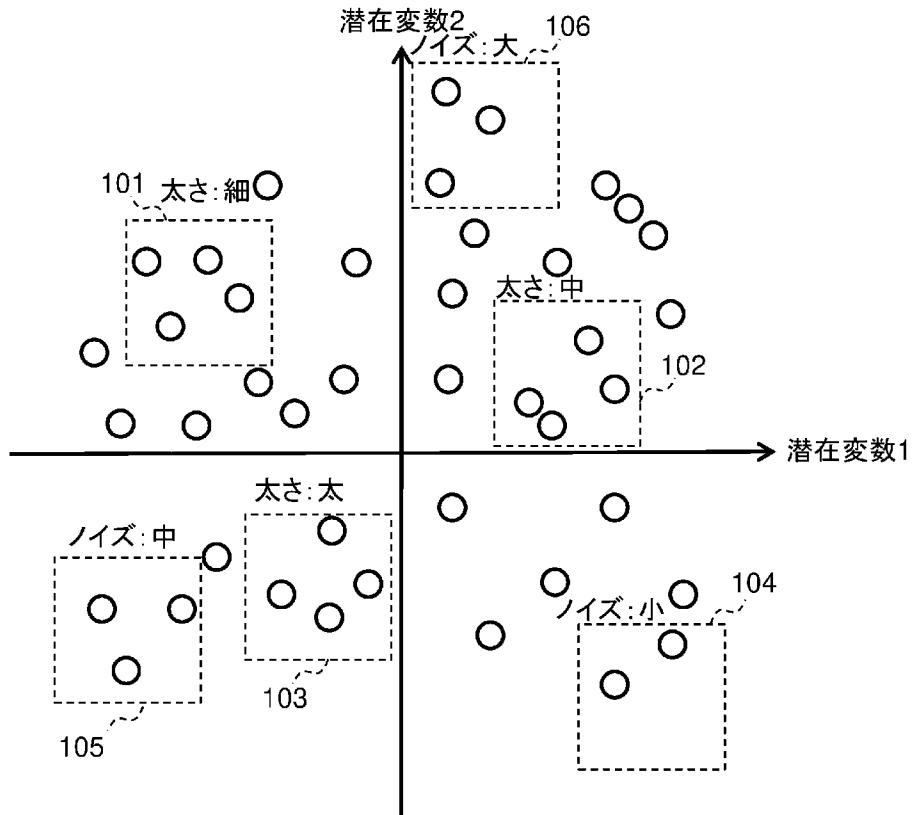
前記MF C V A E によるデータ再構成の誤差を表す再構成誤差項と、前記潜在変数の分布に制約を与える正則化項と、を有する前記MF C V A E の損失関数に、前記属性毎の前記指標に基づく追加項を追加した損失関数の関数値を算出する損失算出部と、

前記損失算出部によって算出された前記関数値に基づく誤差逆伝搬によって前記MF C V A E のモデルパラメータを更新するモデル更新部と、して機能させ、

前記回帰モデル適合度評価部、前記損失算出部、及び前記モデル更新部は、前記予測誤差又はエポック回数が所定条件を充足するまでこの順序で処理を順次繰り返すことで前記MF C V A E のモデル学習を実行する、ことを特徴とする情報処理プログラム。

[図1]

図 1



[図2]

図 2

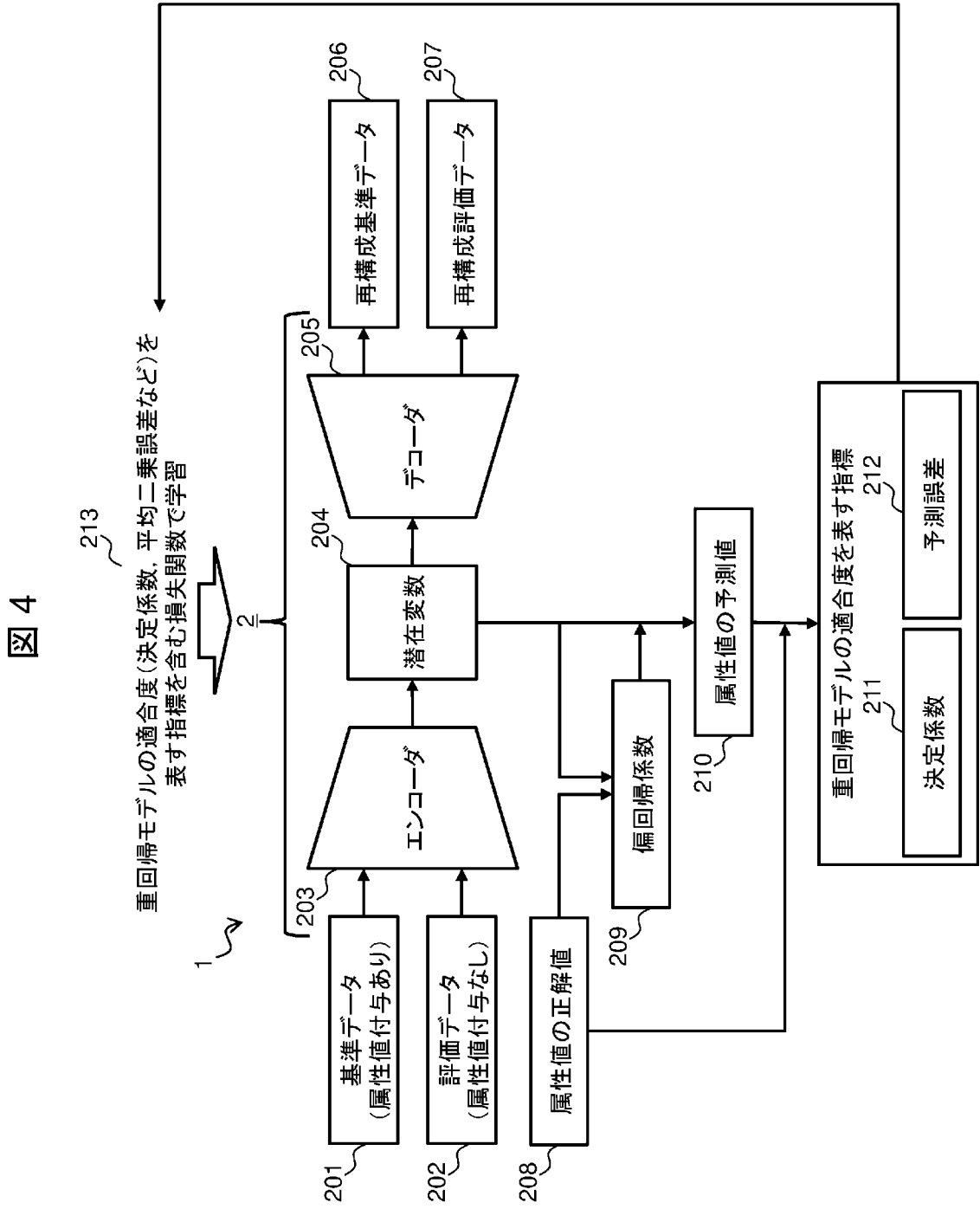
データ番号 (データ属性)	属性				
	文字の種類	太さ(mm)	傾き(度)	...	砕け度合い (レベル)
1 (活字)	あ	0.3	3.0	...	3
2 (活字)	い	0.2	1.0	...	2
3 (活字)	う	0.5	5.0	...	4
4 (手書き文字)	あ			...	
5 (手書き文字)	か			...	
...					

[図3]

図 3

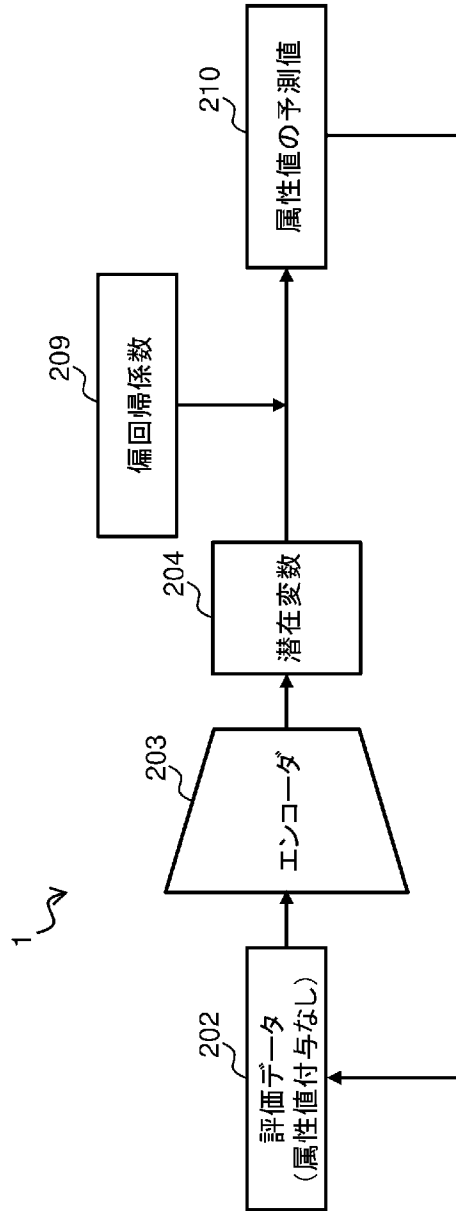
データ番号	属性				
	属性1	属性2	属性3	...	属性J
1	0.1	0.3		...	
2	0.5	0.2		...	
3	0.6		3.0	...	
4	0.7		1.0		2
5		0.5	5.0	...	3
...					

[図4]



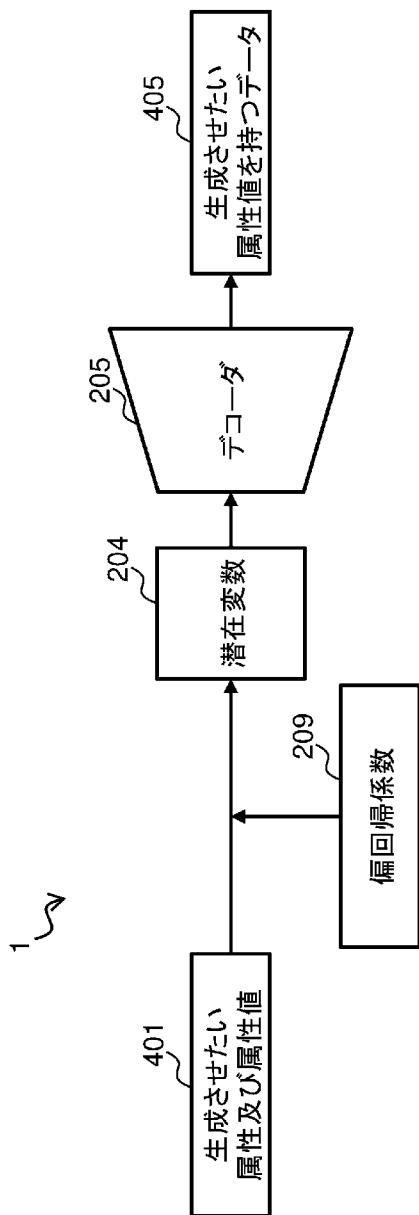
[図5]

図 5



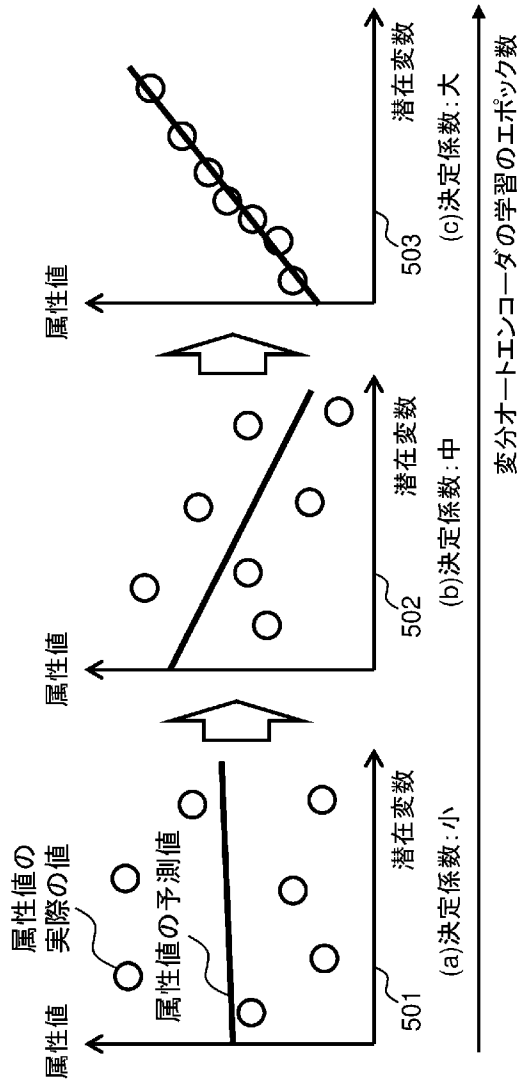
[図6]

図 6

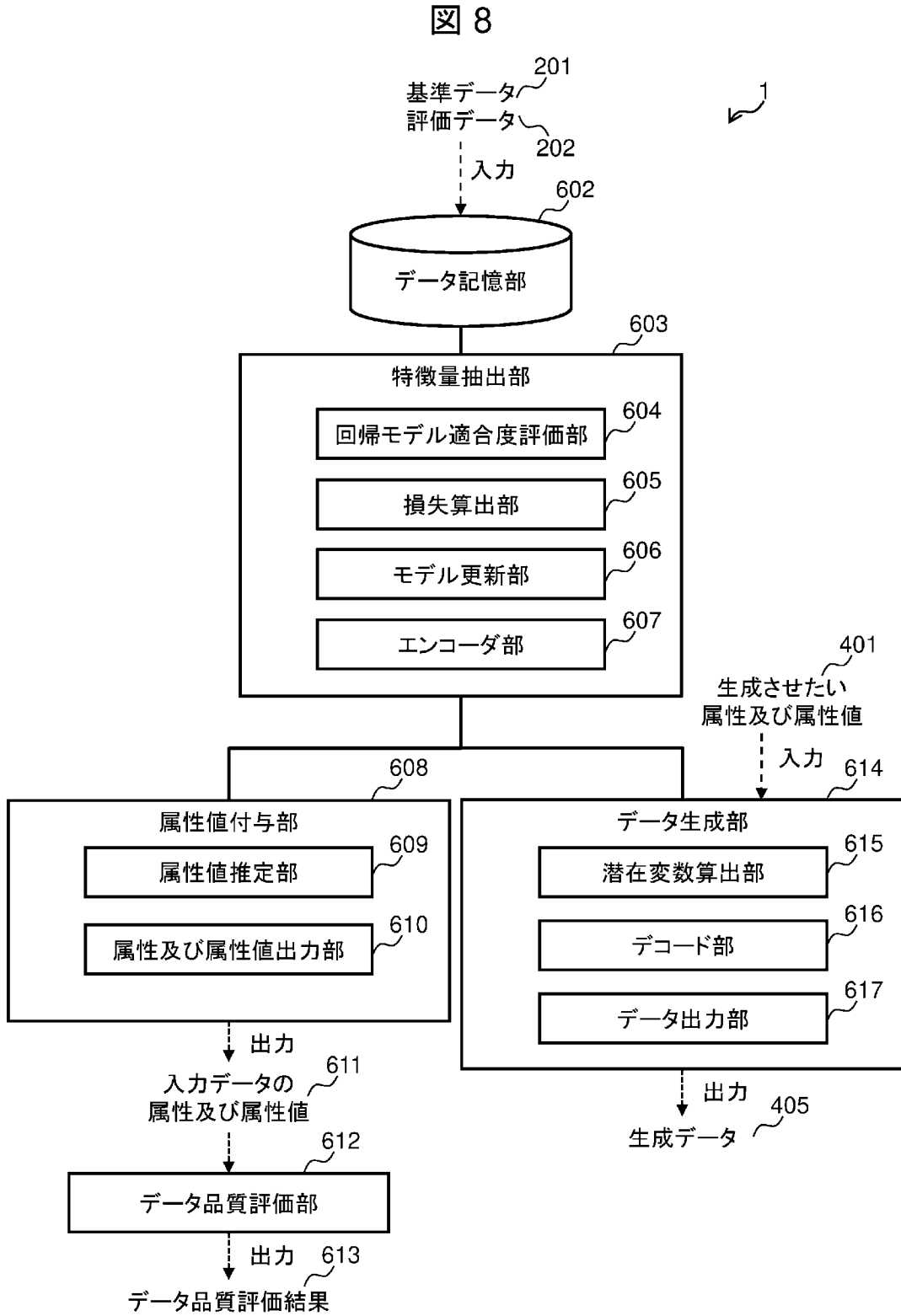


[図7]

図7

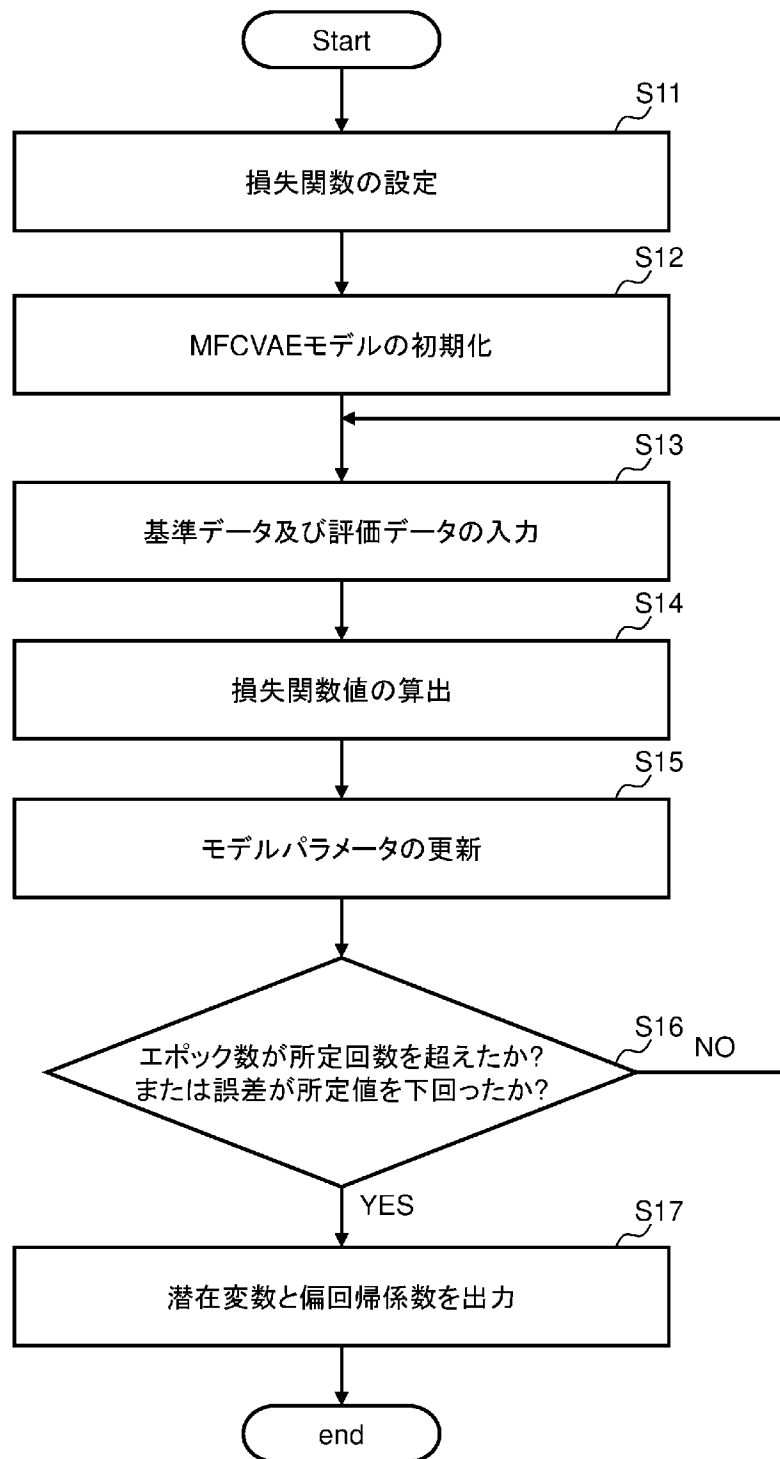


[図8]



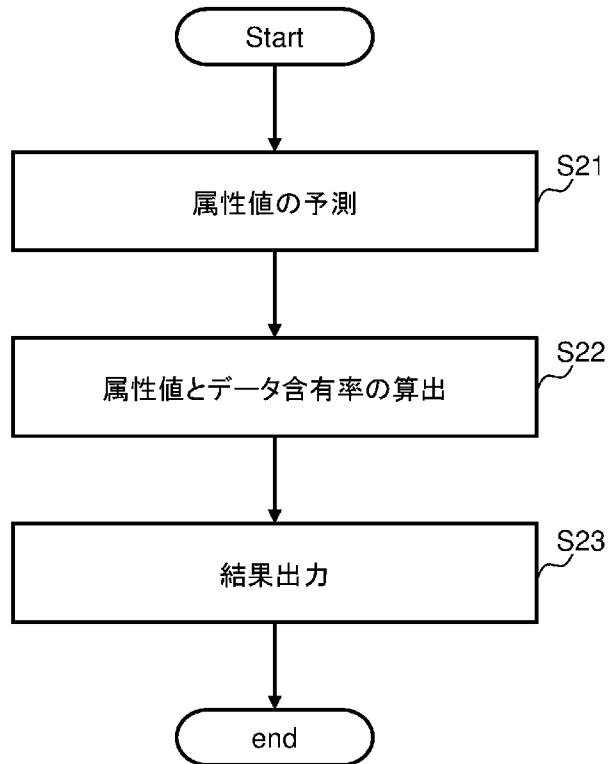
[図9]

図 9



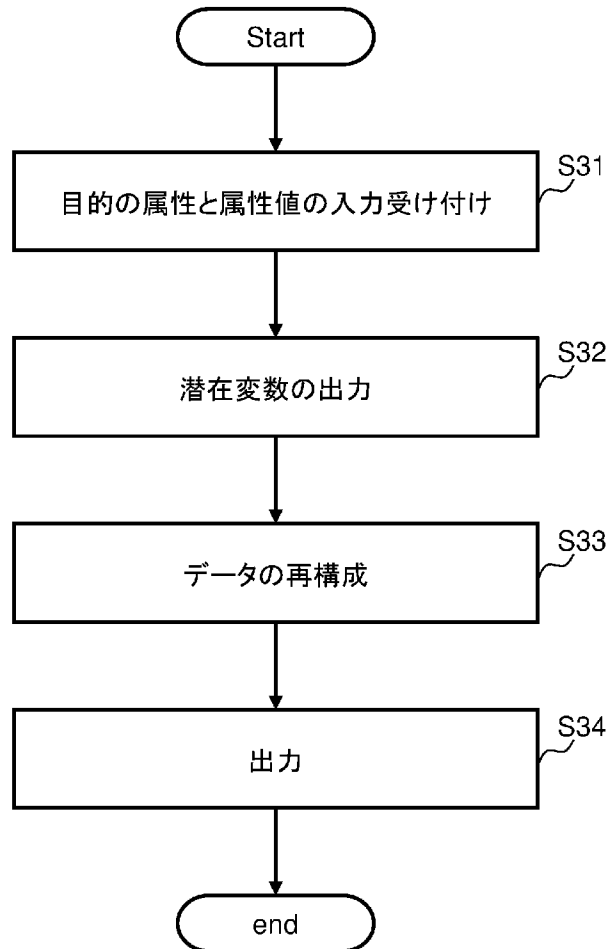
[図10]

図 10



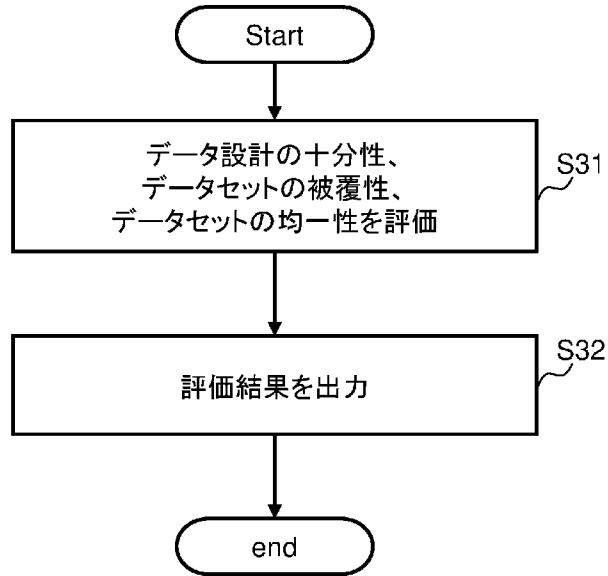
[図11]

図 11



[図12]

図 12

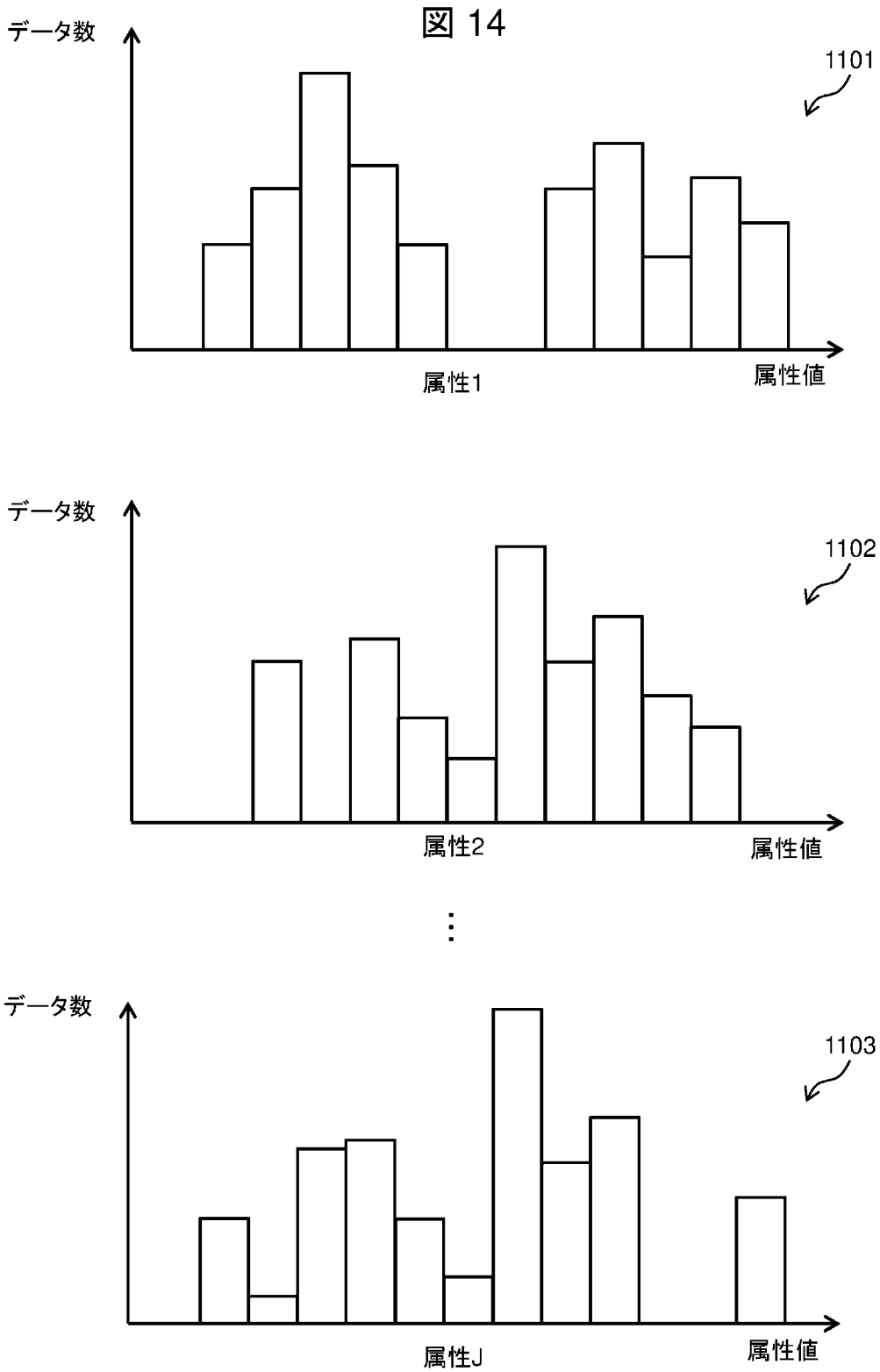


[図13]

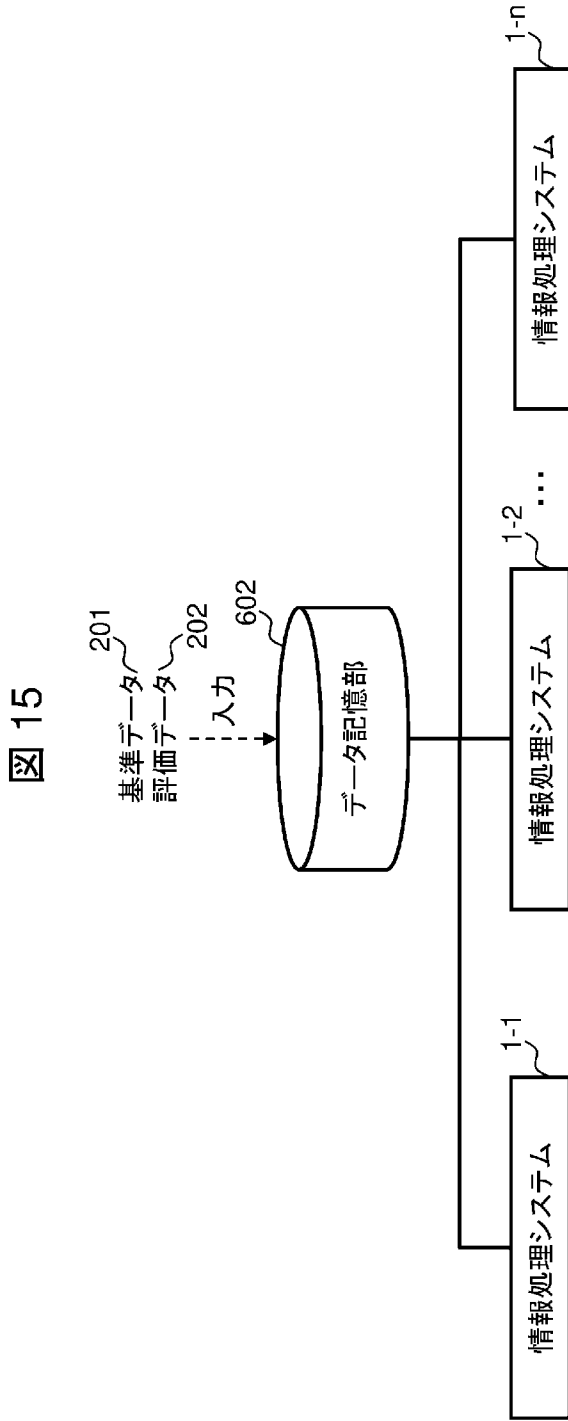
図 13

データ番号/属性	属性1	属性2	属性3	...	属性J
1	0.1	0.3	2.0	...	4
2	0.5	0.2	1.0	...	2
3	0.6	0.5	3.0	...	3
4	0.7	0.3	1.0		2
5	0.2	0.5	5.0	...	3

[図14]

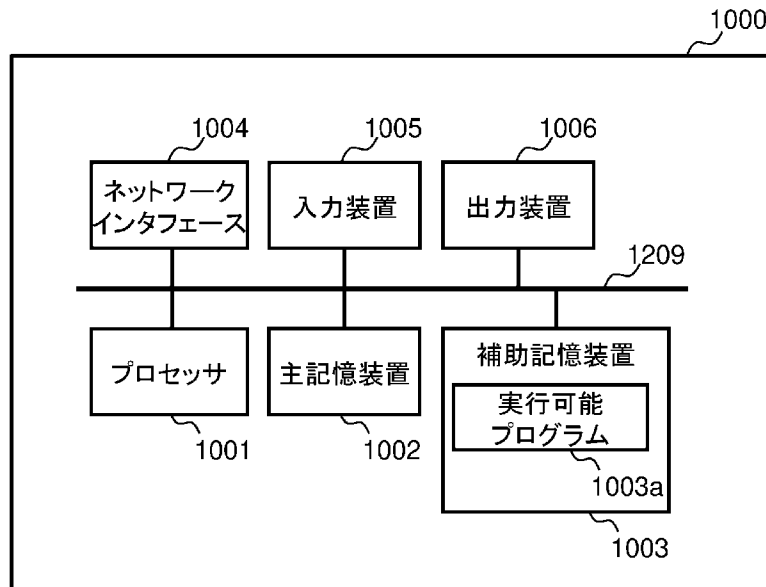


[図15]



[図16]

図 16



## INTERNATIONAL SEARCH REPORT

International application No.

PCT/JP2023/005451

<b>A. CLASSIFICATION OF SUBJECT MATTER</b>		
G06N 3/0455(2023.01)i FI: G06N3/0455		
According to International Patent Classification (IPC) or to both national classification and IPC		
<b>B. FIELDS SEARCHED</b>		
Minimum documentation searched (classification system followed by classification symbols) G06N3/00-99/00		
Documentation searched other than minimum documentation to the extent that such documents are included in the fields searched Published examined utility model applications of Japan 1922-1996 Published unexamined utility model applications of Japan 1971-2023 Registered utility model specifications of Japan 1996-2023 Published registered utility model applications of Japan 1994-2023		
Electronic data base consulted during the international search (name of data base and, where practicable, search terms used)		
<b>C. DOCUMENTS CONSIDERED TO BE RELEVANT</b>		
Category*	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
A	FALCK, Fabian et al. Multi-Facet Clustering Variational Autoencoders. arXiv [online]. 29 October 2021, [retrieval date 25 April 2023], Internet: <URL: https://arxiv.org/abs/2106.05241v2> abstract, section 3	1-17
A	JP 2020-144799 A (FUJITSU LTD) 10 September 2020 (2020-09-10) paragraphs [0026]-[0032], fig. 3-4	1-17
P, A	関根 理敏ほか, データセットの多種多様な属性情報抽出に向けた多面的クラスタリング変分オートエンコーダの手書き文字データへの適用, 情報処理学会 シンポジウム ソフトウェアエンジニアリングシンポジウム 2022, 29 August 2022, pp. 145-146 entire text, all drawings, (SEKINE, Masatoshi et al.), non-official translation (Application of multi-facet clustering variational autoencoders to handwritten character data, so as to extract a wide variety of attribute information from datasets. IPSJ symposium: Software Engineering Symposium 2022.)	1-17
<input type="checkbox"/> Further documents are listed in the continuation of Box C. <input checked="" type="checkbox"/> See patent family annex.		
* Special categories of cited documents: "A" document defining the general state of the art which is not considered to be of particular relevance "E" earlier application or patent but published on or after the international filing date "L" document which may throw doubts on priority claim(s) or which is cited to establish the publication date of another citation or other special reason (as specified) "O" document referring to an oral disclosure, use, exhibition or other means "P" document published prior to the international filing date but later than the priority date claimed "T" later document published after the international filing date or priority date and not in conflict with the application but cited to understand the principle or theory underlying the invention "X" document of particular relevance; the claimed invention cannot be considered novel or cannot be considered to involve an inventive step when the document is taken alone "Y" document of particular relevance; the claimed invention cannot be considered to involve an inventive step when the document is combined with one or more other such documents, such combination being obvious to a person skilled in the art "&" document member of the same patent family		
Date of the actual completion of the international search <b>25 April 2023</b>		Date of mailing of the international search report <b>09 May 2023</b>
Name and mailing address of the ISA/JP <b>Japan Patent Office (ISA/JP) 3-4-3 Kasumigaseki, Chiyoda-ku, Tokyo 100-8915 Japan</b>		Authorized officer  Telephone No.

**INTERNATIONAL SEARCH REPORT**  
**Information on patent family members**

International application No.

**PCT/JP2023/005451**

Patent document cited in search report	Publication date (day/month/year)	Patent family member(s)	Publication date (day/month/year)
JP 2020-144799 A	10 September 2020	US 2020/0285690 A1 paragraphs [0039]-[0045], fig. 3-4 EP 3706029 A1 CN 111666708 A	

---

A. 発明の属する分野の分類（国際特許分類（IPC）） G06N 3/0455(2023.01)i FI: G06N3/0455		
B. 調査を行った分野 調査を行った最小限資料（国際特許分類（IPC）） G06N3/00-99/00 最小限資料以外の資料で調査を行った分野に含まれるもの 日本国実用新案公報 1922-1996年 日本国公開実用新案公報 1971-2023年 日本国実用新案登録公報 1996-2023年 日本国登録実用新案公報 1994-2023年		
国際調査で使用した電子データベース（データベースの名称、調査に使用した用語）		
C. 関連すると認められる文献		
引用文献の カテゴリー*	引用文献名 及び一部の箇所が関連するときは、その関連する箇所の表示	関連する 請求項の番号
A	FALCK, Fabian et al., Multi-Facet Clustering Variational Autoencoders, arXiv [online], 2021.10.29, [検索日 2023.04.25], インターネット: <URL: https://arxiv.org/abs/2106.05241v2> Abstract, 第3節	1-17
A	JP 2020-144799 A (富士通株式会社) 10.09.2020 (2020-09-10) [0026]-[0032], 図3-4	1-17
P, A	関根 理敏ほか, データセットの多種多様な属性情報抽出に向けた多面的クラスタリング変分オートエンコーダの手書き文字データへの適用, 情報処理学会 シンポジウム ソフトウェアエンジニアリングシンポジウム 2022, 2022.08.29, pp.145-146 全文、全図	1-17
<input type="checkbox"/> C欄の続きにも文献が列挙されている。 <input checked="" type="checkbox"/> パテントファミリーに関する別紙を参照。		
* 引用文献のカテゴリー “A” 特に関連のある文献ではなく、一般的技術水準を示すもの “E” 国際出願日前の出願または特許であるが、国際出願日以後に公表されたもの “L” 優先権主張に疑義を提起する文献又は他の文献の発行日若しくは他の特別な理由を確立するために引用する文献（理由を付す） “O” 口頭による開示、使用、展示等に言及する文献 “P” 国際出願日前で、かつ優先権の主張の基礎となる出願の日の後に公表された文献 “T” 国際出願日又は優先日後に公表された文献であって出願と抵触するものではなく、発明の原理又は理論の理解のために引用するもの “X” 特に関連のある文献であって、当該文献のみで発明の新規性又は進歩性がないと考えられるもの “Y” 特に関連のある文献であって、当該文献と他の1以上の文献との、当業者にとって自明である組合せによって進歩性がないと考えられるもの “&” 同一パテントファミリー文献		
国際調査を完了した日	25.04.2023	国際調査報告の発送日
		09.05.2023
名称及びあて先	日本国特許庁(ISA/JP) 〒100-8915 日本国 東京都千代田区霞が関三丁目4番3号	権限のある職員（特許庁審査官）  山本 俊介 5B 5087  電話番号 03-3581-1101 内線 3545

国際調査報告  
パテントファミリーに関する情報

国際出願番号

PCT/JP2023/005451

引用文献	公表日	パテントファミリー文献	公表日
JP 2020-144799 A	10.09.2020	US 2020/0285690 A1 [0039]-[0045], FIGs.3-4 EP 3706029 A1 CN 111666708 A	