

(19)日本国特許庁(JP)

(12)公開特許公報(A)

(11)公開番号

特開2022-144738  
(P2022-144738A)

(43)公開日 令和4年10月3日(2022.10.3)

|                         |                     |            |
|-------------------------|---------------------|------------|
| (51)国際特許分類              | F I                 | テーマコード(参考) |
| G 0 6 F 16/35 (2019.01) | G 0 6 F 16/35       | 5 B 1 7 5  |
| G 0 6 N 20/00 (2019.01) | G 0 6 N 20/00 1 6 0 |            |

審査請求 未請求 請求項の数 7 O L (全17頁)

|          |                           |          |   |
|----------|---------------------------|----------|---|
| (21)出願番号 | 特願2021-45884(P2021-45884) | (71)出願人  | 000006150<br>京セラドキュメントソリューションズ株式会社<br>大阪府大阪市中央区玉造1丁目2番28号 |
| (22)出願日  | 令和3年3月19日(2021.3.19)      | (74)代理人  | 100140796<br>弁理士 原口 貴志                                    |
|          |                           | (72)発明者  | 庄司 秀典<br>大阪府大阪市中央区玉造1丁目2番28号 京セラドキュメントソリューションズ株式会社内       |
|          |                           | Fターム(参考) | 5B175 DA01 FA03   |

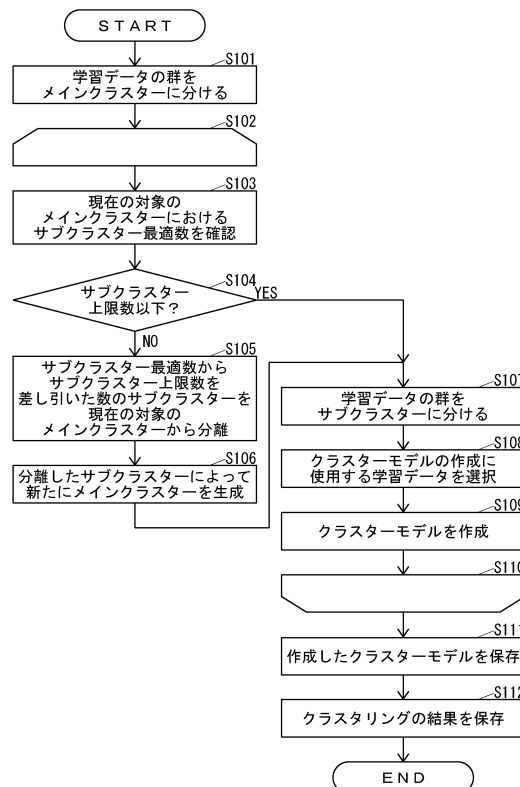
(54)【発明の名称】 情報抽出システムおよび情報抽出プログラム

(57)【要約】

【課題】 情報抽出モデルの作成のための計算量を低減することができる情報抽出システムおよび情報抽出プログラムを提供する。

【解決手段】 情報抽出システムは、請求書データから情報を抽出するための情報抽出モデルであるクラスターモデルの作成のための学習データの群をクラスタリングすることによって、学習データのそれぞれをいずれかのメインクラスターに分け(S101)、メインクラスター毎に学習データを使用して学習を実行することによって、メインクラスター毎のクラスターモデルを作成する(S109)ことを特徴とする。

【選択図】 図3



## 【特許請求の範囲】

## 【請求項 1】

文書のデータから情報を抽出するための情報抽出モデルの作成のための学習データの群をクラスタリングすることによって、前記学習データのそれぞれをいずれかのメインクラスターに分ける文書クラスタリング部と、

前記メインクラスター毎に前記学習データを使用して学習を実行することによって、前記メインクラスター毎の前記情報抽出モデルを作成するモデル学習部と  
を備えることを特徴とする情報抽出システム。

## 【請求項 2】

前記文書クラスタリング部は、前記メインクラスター内の前記学習データの群をクラスタリングすることによって、前記メインクラスター内の前記学習データのそれぞれをいずれかのサブクラスターに分け、

前記モデル学習部は、前記情報抽出モデルの作成に使用する前記学習データを前記サブクラスター毎に選択し、選択した前記学習データを使用して学習を実行することによって、前記メインクラスター毎の前記情報抽出モデルを作成することを特徴とする請求項 1 に記載の情報抽出システム。

## 【請求項 3】

前記モデル学習部は、重心が前記メインクラスターの重心に最も近い前記サブクラスターにおいて、重心が前記メインクラスターの重心に最も近い前記学習データを、前記情報抽出モデルの作成に使用する前記学習データとして選択することを特徴とする請求項 2 に記載の情報抽出システム。

## 【請求項 4】

前記モデル学習部は、重心が前記メインクラスターの重心に最も近い前記サブクラスター以外の前記サブクラスターのそれぞれにおいて、重心が前記メインクラスターの重心から最も遠い前記学習データを、前記情報抽出モデルの作成に使用する前記学習データとして選択することを特徴とする請求項 3 に記載の情報抽出システム。

## 【請求項 5】

前記文書クラスタリング部は、前記メインクラスターにおける前記サブクラスターの最適数をクラスター数自動推定法によって確認し、確認した前記最適数が特定の上限数を超える場合に、前記最適数から前記上限数を差し引いた数の前記サブクラスターを、このメインクラスターから分離することを特徴とする請求項 2 から請求項 4 までのいずれかに記載の情報抽出システム。

## 【請求項 6】

前記文書クラスタリング部は、前記最適数から前記上限数を差し引いた数の前記サブクラスターを前記メインクラスターから分離する場合に、重心がこのメインクラスターの重心から遠い前記サブクラスターを優先して、このメインクラスターから分離することを特徴とする請求項 5 に記載の情報抽出システム。

## 【請求項 7】

文書のデータから情報を抽出するための情報抽出モデルの作成のための学習データの群をクラスタリングすることによって、前記学習データのそれぞれをいずれかのメインクラスターに分ける文書クラスタリング部と、

前記メインクラスター毎に前記学習データを使用して学習を実行することによって、前記メインクラスター毎の前記情報抽出モデルを作成するモデル学習部と  
をコンピューターに実現させることを特徴とする情報抽出プログラム。

## 【発明の詳細な説明】

## 【技術分野】

## 【0001】

本発明は、文書のデータから特定の項目に対する値を抽出する情報抽出システムおよび情報抽出プログラムに関する。

## 【背景技術】

10

20

30

40

50

## 【 0 0 0 2 】

従来、文書のデータから情報を抽出するための情報抽出モデルを使用して文書のデータから情報を抽出する情報抽出システムが知られている（例えば、特許文献 1、2 参照。）。

## 【 先行技術文献 】

## 【 特許文献 】

## 【 0 0 0 3 】

【 特許文献 1 】 米国特許出願公開第 2 0 1 4 / 0 1 7 7 9 5 1 号明細書

【 特許文献 2 】 特許第 6 6 2 9 9 4 2 号公報

## 【 発明の概要 】

## 【 発明が解決しようとする課題 】

## 【 0 0 0 4 】

しかしながら、従来の情報抽出システムにおいては、情報抽出モデルの作成のための計算量が多いという問題がある。

## 【 0 0 0 5 】

そこで、本発明は、情報抽出モデルの作成のための計算量を低減することができる情報抽出システムおよび情報抽出プログラムを提供することを目的とする。

## 【 課題を解決するための手段 】

## 【 0 0 0 6 】

本発明の情報抽出システムは、文書のデータから情報を抽出するための情報抽出モデルの作成のための学習データの群をクラスタリングすることによって、前記学習データのそれぞれをいずれかのメインクラスターに分ける文書クラスタリング部と、前記メインクラスター毎に前記学習データを使用して学習を実行することによって、前記メインクラスター毎の前記情報抽出モデルを作成するモデル学習部とを備えることを特徴とする。

## 【 0 0 0 7 】

この構成により、本発明の情報抽出システムは、メインクラスター毎に情報抽出モデルを作成するので、情報抽出モデル毎の特徴を単純化することができ、その結果、情報抽出モデル毎に必要な学習データの数を低減することができる。したがって、本発明の情報抽出システムは、情報抽出モデルの作成のための計算量を低減することができる。

## 【 0 0 0 8 】

本発明の情報抽出システムにおいて、前記文書クラスタリング部は、前記メインクラスター内の前記学習データの群をクラスタリングすることによって、前記メインクラスター内の前記学習データのそれぞれをいずれかのサブクラスターに分け、前記モデル学習部は、前記情報抽出モデルの作成に使用する前記学習データを前記サブクラスター毎に選択し、選択した前記学習データを使用して学習を実行することによって、前記メインクラスター毎の前記情報抽出モデルを作成しても良い。

## 【 0 0 0 9 】

この構成により、本発明の情報抽出システムは、情報抽出モデルの作成に使用する学習データをサブクラスター毎に選択し、選択した学習データを使用して学習を実行することによって、メインクラスター毎の情報抽出モデルを作成するので、情報抽出モデル毎に必要な学習データの数を低減することができ、その結果、情報抽出モデルの作成のための計算量を低減することができる。

## 【 0 0 1 0 】

本発明の情報抽出システムにおいて、前記モデル学習部は、重心が前記メインクラスターの重心に最も近い前記サブクラスターにおいて、重心が前記メインクラスターの重心に最も近い前記学習データを、前記情報抽出モデルの作成に使用する前記学習データとして選択しても良い。

## 【 0 0 1 1 】

この構成により、本発明の情報抽出システムは、重心がメインクラスターの重心に最も近いサブクラスターにおいて、重心がメインクラスターの重心に最も近い学習データを、

10

20

30

40

50

情報抽出モデルの作成に使用する学習データとして選択するので、メインクラスターの特徴を最も強く表す学習データを使用して情報抽出モデルを作成することができ、その結果、メインクラスターの特徴が適切に反映された情報抽出モデルを作成することができる。

【0012】

本発明の情報抽出システムにおいて、前記モデル学習部は、重心が前記メインクラスターの重心に最も近い前記サブクラスター以外の前記サブクラスターのそれぞれにおいて、重心が前記メインクラスターの重心から最も遠い前記学習データを、前記情報抽出モデルの作成に使用する前記学習データとして選択しても良い。

【0013】

この構成により、本発明の情報抽出システムは、重心がメインクラスターの重心に最も近いサブクラスター以外のサブクラスターのそれぞれにおいて、重心がメインクラスターの重心から最も遠い学習データを、情報抽出モデルの作成に使用する学習データとして選択するので、メインクラスターにおいて広範囲に散らばった学習データを使用して情報抽出モデルを作成することができ、その結果、メインクラスターの特徴が適切に反映された情報抽出モデルを作成することができる。

10

【0014】

本発明の情報抽出システムにおいて、前記文書クラスタリング部は、前記メインクラスターにおける前記サブクラスターの最適数をクラスター数自動推定法によって確認し、確認した前記最適数が特定の上限数を超える場合に、前記最適数から前記上限数を差し引いた数の前記サブクラスターを、このメインクラスターから分離しても良い。

20

【0015】

この構成により、本発明の情報抽出システムは、メインクラスターにおけるサブクラスターの最適数が特定の上限数を超える場合に、最適数から上限数を差し引いた数のサブクラスターを、このメインクラスターから分離するので、情報抽出モデル毎に必要な学習データの数を低減することができ、その結果、情報抽出モデルの作成のための計算量を低減することができる。

【0016】

本発明の情報抽出システムにおいて、前記文書クラスタリング部は、前記最適数から前記上限数を差し引いた数の前記サブクラスターを前記メインクラスターから分離する場合に、重心がこのメインクラスターの重心から遠い前記サブクラスターを優先して、このメインクラスターから分離しても良い。

30

【0017】

この構成により、本発明の情報抽出システムは、最適数から上限数を差し引いた数のサブクラスターをメインクラスターから分離する場合に、重心がこのメインクラスターの重心から遠いサブクラスターを優先して、このメインクラスターから分離するので、メインクラスターの特徴を強く表す学習データを使用して情報抽出モデルを作成することができ、その結果、メインクラスターの特徴が適切に反映された情報抽出モデルを作成することができる。

【0018】

本発明の情報抽出プログラムは、文書のデータから情報を抽出するための情報抽出モデルの作成のための学習データの群をクラスタリングすることによって、前記学習データのそれぞれをいずれかのメインクラスターに分ける文書クラスタリング部と、前記メインクラスター毎に前記学習データを使用して学習を実行することによって、前記メインクラスター毎の前記情報抽出モデルを作成するモデル学習部とをコンピューターに実現させることを特徴とする。

40

【0019】

この構成により、本発明の情報抽出プログラムを実行するコンピューターは、メインクラスター毎に情報抽出モデルを作成するので、情報抽出モデル毎の特徴を単純化することができ、その結果、情報抽出モデル毎に必要な学習データの数を低減することができる。したがって、本発明の情報抽出プログラムを実行するコンピューターは、情報抽出モデル

50

の作成のための計算量を低減することができる。

【発明の効果】

【0020】

本発明の情報抽出システムおよび情報抽出プログラムは、情報抽出モデルの作成のための計算量を低減することができる。

【図面の簡単な説明】

【0021】

【図1】本発明の一実施の形態に係る情報抽出システムのブロック図である。

【図2】図1に示す記憶部に記憶される情報抽出モデルの一例を示す図である。

【図3】クラスターモデルを作成する場合の図1に示す情報抽出システムの動作のフローチャートである。 10

【図4】図3に示す動作において学習データの群をメインクラスターに分ける処理のイメージを示す図である。

【図5】図3に示す動作においてメインクラスターからサブクラスターを分離する処理のイメージを示す図である。

【図6】図3に示す動作においてクラスターモデルの作成に使用する学習データを選択する処理のイメージを示す図である。

【図7】請求書データから特定の項目に対する値を抽出する場合の図1に示す情報抽出システムの動作のフローチャートである。

【図8】クラスターモデルを更新する場合の図1に示す情報抽出システムの動作の一部のフローチャートである。 20

【図9】図8に示す動作の続きの動作のフローチャートである。

【発明を実施するための形態】

【0022】

以下、本発明の実施の形態について、図面を用いて説明する。

【0023】

まず、本発明の一実施の形態に係る情報抽出システムの構成について説明する。

【0024】

図1は、本実施の形態に係る情報抽出システム10のブロック図である。

【0025】

図1に示すように、情報抽出システム10は、種々の操作が入力される例えばキーボード、マウスなどの操作デバイスである操作部11と、種々の情報を表示する例えばLCD(Liquid Crystal Display)などの表示デバイスである表示部12と、LAN、インターネットなどのネットワーク経由で、または、ネットワークを介さずに有線または無線によって直接に、外部の装置と通信を行う通信デバイスである通信部13と、各種の情報を記憶する例えば半導体メモリー、HDD(Hard Disk Drive)などの不揮発性の記憶デバイスである記憶部14と、情報抽出システム10全体を制御する制御部15とを備えている。情報抽出システム10は、例えば、PC(Personal Computer)またはサーバーによって構成されても良いし、プリンター専用機などの画像形成装置によって構成されても良い。 30 40

【0026】

記憶部14は、文書としての請求書のデータ(以下「請求書データ」という。)から情報を抽出するための情報抽出モデルを使用して請求書データから情報を抽出するための情報抽出プログラム14aを記憶している。情報抽出プログラム14aは、例えば、情報抽出システム10の製造段階で情報抽出システム10にインストールされていても良いし、USB(Universal Serial Bus)メモリーなどの外部の記憶媒体から情報抽出システム10に追加でインストールされても良いし、ネットワーク上から情報抽出システム10に追加でインストールされても良い。

【0027】

記憶部14は、複数のフォーマットの請求書を学習済みの情報抽出モデル(以下「ベー 50

スモデル」という。) 14bを記憶している。ベースモデル14bは、情報抽出システム10の利用者に情報抽出システム10を提供する者が用意しても良い。

【0028】

記憶部14は、後述のメインクラスター毎の情報抽出モデル(以下「クラスターモデル」という。)14cを記憶可能である。クラスターモデルによる値の抽出の対象の請求書データ(以下「抽出対象データ」という。)は、請求書内の文字と、請求書内の文字以外の素性とを含む請求書データである。請求書内の文字以外の素性には、請求書における各文字の座標が含まれる。また、請求書内の文字以外の素性には、例えば、請求書内の画像と、請求書における各画像の座標とが含まれても良い。請求書内の文字と、請求書における各文字の座標とは、例えば、請求書の画像に対するOCR(Optical Character Recognition)処理によって取得されることが可能である。請求書内の画像と、請求書における各画像の座標とは、これらを請求書の画像から取得することが可能なシステムによって取得されることが可能である。

10

【0029】

記憶部14は、メインクラスターのクラスタリングの結果(以下「クラスタリング結果」という。)14dを記憶可能である。

【0030】

制御部15は、例えば、CPU(Central Processing Unit)と、プログラムおよび各種のデータを記憶しているROM(Read Only Memory)と、制御部15のCPUの作業領域として用いられるメモリーとしてのRAM(Random Access Memory)とを備えている。制御部15のCPUは、記憶部14または制御部15のROMに記憶されているプログラムを実行する。

20

【0031】

制御部15は、情報抽出プログラム14aを実行することによって、請求書データをクラスタリングする文書クラスタリング部15aと、クラスターモデルを作成するモデル学習部15bと、クラスターモデルを使用して請求書データから特定の項目に対する値を抽出するデータ抽出実行部15cとを実現する。

【0032】

文書クラスタリング部15aにおいてクラスタリングに使用されるアルゴリズムとしては、例えば、DBSCAN、g-means、エルボー法など、クラスターの数を自動で決定することが可能なアルゴリズムが採用される。文書クラスタリング部15aにおいてクラスタリングに使用される素性としては、例えば、単語ベクトル、単語の座標が採用される。単語ベクトルとしては、例えば、one-hotベクトル、tf-idf、word2vecなどのベクトル表現が採用される。

30

【0033】

モデル学習部15bにおいてクラスターモデルの作成に使用されるアルゴリズムとしては、例えば、LSTM、Transformerなどの自然言語処理を使用したアルゴリズムをベースにしたものが採用される。モデル学習部15bにおいてクラスターモデルの作成に使用される素性としては、例えば、テキスト情報、文字の座標が採用される。

【0034】

データ抽出実行部15cによって値を抽出する対象の文書には、文書毎に値の記載の位置が異なる場合がない定型文書と、文書毎に値の記載の位置が異なる場合がある準定型文書とが含まれるが、非定型文書は含まれない。

40

【0035】

文書クラスタリング部15a、モデル学習部15bおよびデータ抽出実行部15cにおいてデータの距離の計算に使用されるアルゴリズムとしては、例えば、コサイン距離、マンハッタン距離、ユークリッド距離が採用される。

【0036】

図2は、記憶部14に記憶される情報抽出モデル20の一例を示す図である。

【0037】

50

図 2 に示す情報抽出モデル 20 は、抽出対象データ 40 における「請求書内の文字」に基づいて各文字を取得し (S 21)、S 21 において取得した各文字に対して、各文字に基づいたベクトル情報を付与し (S 22)、S 22 の出力を Bi-LSTM に入力する (S 23)。

【0038】

また、情報抽出モデル 20 は、抽出対象データ 40 における「請求書内の文字」に基づいて各単語を取得し (S 24)、S 24 において取得した各単語に対して、各単語に基づいたベクトル情報を付与する (S 25)。

【0039】

また、情報抽出モデル 20 は、抽出対象データ 40 における「請求書における各文字の座標」に基づいて各単語の座標を取得し (S 26)、S 26 において取得した各単語の座標を全結合層に入力する (S 27)。

10

【0040】

そして、情報抽出モデル 20 は、S 23 の出力と、S 25 の出力と、S 27 の出力とを連結する (S 28)。

【0041】

次いで、情報抽出モデル 20 は、S 28 の出力を Bi-LSTM に入力し (S 29)、S 29 の出力を全結合層に入力し (S 30)、S 30 の出力を全結合層に入力し (S 31)、S 31 の出力を CRF に入力する (S 32)。

【0042】

20

次に、情報抽出システム 10 の動作について説明する。

【0043】

まず、クラスターモデルを作成する場合の情報抽出システム 10 の動作について説明する。

【0044】

図 3 は、クラスターモデルを作成する場合の情報抽出システム 10 の動作のフローチャートである。

【0045】

利用者は、クラスターモデルの作成のための学習データの群を用意し、用意した学習データの群を使用した学習を、操作部 11 から、または、図示していないコンピューターから通信部 13 を介して、情報抽出システム 10 に指示することができる。ここで、学習データは、請求書内の文字と、請求書内の文字以外の素性と、請求書から抽出されることを利用者が希望する項目に対する正解ラベルとを含む、請求書毎の請求書データである。請求書内の文字以外の素性には、請求書における各文字の座標が含まれる。また、請求書内の文字以外の素性には、例えば、請求書内の画像と、請求書における各画像の座標とが含まれても良い。請求書から抽出されることを利用者が希望する項目とは、例えば、文書が請求書である場合には、請求先、請求日、締切日、請求金額などである。文書から抽出されることを利用者が希望する項目に対する正解ラベルは、請求書内の文字と、請求書内の文字以外の素性から、利用者によって選択された値である。請求書内の文字と、請求書における各文字の座標とは、例えば、請求書の画像に対する OCR 処理によって取得されることが可能である。請求書内の画像と、請求書における各画像の座標とは、これらを請求書の画像から取得することが可能なシステムによって取得されることが可能である。

30

40

【0046】

情報抽出システム 10 の制御部 15 は、学習データの群を使用した学習が指示されると、図 3 に示す動作を実行する。

【0047】

図 3 に示すように、文書クラスタリング部 15 a は、学習データの群をクラスタリングすることによって、学習データのそれぞれをいずれかのメインクラスターに分ける (S 101)。

【0048】

50

図4は、図3に示す動作において学習データの群をメインクラスターに分ける処理のイメージを示す図である。なお、図4(b)において、学習データは、学習データ自身が所属するメインクラスター毎のマークで表示されている。

【0049】

図4に示すように、文書クラスタリング部15aは、学習データの群をクラスタリングするために、学習データの対象の請求書内の文字を学習データ同士で比較することができるよう学習データを図4(a)に示すようにベクトル化する。

【0050】

次いで、文書クラスタリング部15aは、ベクトル化した学習データの群をクラスタリングすることによって、学習データのそれぞれを図4(b)に示すようにメインクラスターA～Eのいずれかに分ける(S101)。

【0051】

図3に示すように、制御部15は、S101の処理の後、図3に示す動作の今回の実行において未だS103の処理の対象にしていないメインクラスターの1つを対象にする(S102)。

【0052】

次いで、文書クラスタリング部15aは、現在の対象のメインクラスターにおけるサブクラスターの最適数(以下「サブクラスター最適数」という。)をクラスター数自動推定法によって確認する(S103)。

【0053】

次いで、文書クラスタリング部15aは、S103において確認したサブクラスター最適数が、サブクラスターの上限数(以下「サブクラスター上限数」という。)以下であるか否かを判断する(S104)。ここで、サブクラスター上限数は、本実施の形態において例えば5である。

【0054】

文書クラスタリング部15aは、S103において確認したサブクラスター最適数がサブクラスター上限数以下ではないとS104において判断すると、S103において確認したサブクラスター最適数からサブクラスター上限数を差し引いた数のサブクラスターを現在の対象のメインクラスターから分離する(S105)。ここで、文書クラスタリング部15aは、重心が現在の対象のメインクラスターの重心から遠いサブクラスターを優先して現在の対象のメインクラスターから分離する。なお、メインクラスターの重心は、例えば、このメインクラスターに所属する学習データの文書ベクトルの平均値である。同様に、サブクラスターの重心は、例えば、このサブクラスターに所属する学習データの文書ベクトルの平均値である。

【0055】

文書クラスタリング部15aは、S105の処理の後、S105において現在の対象のメインクラスターから分離したサブクラスターによって新たにメインクラスターを生成する(S106)。すなわち、文書クラスタリング部15aは、S105において現在の対象のメインクラスターから分離したサブクラスターを新たなメインクラスターにする。

【0056】

図5は、図3に示す動作においてメインクラスターからサブクラスターを分離する処理のイメージを示す図である。なお、図5は、図4(b)に示すメインクラスターBの例である。図5(a)、(b)において、学習データは、学習データ自身が所属するサブクラスター毎のマークで表示されている。図5(c)において、学習データは、学習データ自身が所属するメインクラスター毎のマークで表示されている。

【0057】

図5(a)に示すように、文書クラスタリング部15aは、メインクラスターBにおけるサブクラスター最適数を確認する(S103)。図5(a)に示す例では、文書クラスタリング部15aは、メインクラスターBにおけるサブクラスター最適数をクラスター数自動推定法によって7と確認している。

10

20

30

40

50



## 【 0 0 5 8 】

次いで、文書クラスタリング部 1 5 a は、S 1 0 3 において確認したサブクラスター最適数がサブクラスター上限数以下ではない場合に ( S 1 0 4 で N O )、S 1 0 3 において確認したサブクラスター最適数からサブクラスター上限数を差し引いた数のサブクラスターを図 5 ( b ) に示すようにメインクラスター B から分離する ( S 1 0 5 )。すなわち、文書クラスタリング部 1 5 a は、サブクラスター F、G をメインクラスター B から分離する。図 5 ( b ) に示す例は、サブクラスター上限数が 5 の場合の例である。

## 【 0 0 5 9 】

文書クラスタリング部 1 5 a は、S 1 0 5 の処理の後、S 1 0 5 においてメインクラスター B から分離したサブクラスター F、G を図 5 ( c ) に示すようにそれぞれ新たなメインクラスター F、G にする ( S 1 0 6 )。

10

## 【 0 0 6 0 】

図 3 に示すように、文書クラスタリング部 1 5 a は、S 1 0 3 において確認した最適数がサブクラスター上限数以下であると S 1 0 4 において判断するか、S 1 0 6 の処理が終了すると、現在の対象のメインクラスター内の学習データの群をサブクラスター最適数でクラスタリングすることによって、現在の対象のメインクラスター内の学習データのそれぞれをいずれかのサブクラスターに分ける ( S 1 0 7 )。

## 【 0 0 6 1 】

次いで、モデル学習部 1 5 b は、現在の対象のメインクラスター内のサブクラスターから、クラスターモデルの作成に使用する学習データを選択する ( S 1 0 8 )。ここで、モデル学習部 1 5 b は、現在の対象のメインクラスター内のサブクラスターのうち、重心が現在の対象のメインクラスターの重心に最も近いサブクラスターにおいて、重心が現在の対象のメインクラスターの重心に最も近い学習データを、クラスターモデルの作成に使用する学習データとして選択する。また、モデル学習部 1 5 b は、現在の対象のメインクラスター内のサブクラスターのうち、重心が現在の対象のメインクラスターの重心に最も近いサブクラスター以外のサブクラスターのそれぞれにおいて、重心が現在の対象のメインクラスターの重心から最も遠い学習データを、クラスターモデルの作成に使用する学習データとして選択する。なお、学習データの重心は、例えば、この学習データの文書ベクトルである。

20

## 【 0 0 6 2 】

図 6 は、図 3 に示す動作においてクラスターモデルの作成に使用する学習データを選択する処理のイメージを示す図である。なお、図 6 は、図 5 ( c ) に示すメインクラスター B の例である。なお、図 6 において、学習データは、学習データ自身が所属するサブクラスター毎のマークで表示されている。

30

## 【 0 0 6 3 】

図 6 に示すように、モデル学習部 1 5 b は、メインクラスター B 内のサブクラスターのうち、重心がメインクラスター B の重心に最も近いサブクラスター D において、重心がメインクラスター B の重心に最も近い学習データをクラスターモデルの作成に使用する学習データとして選択するとともに、メインクラスター B 内のサブクラスターのうち、サブクラスター D 以外のサブクラスターのそれぞれにおいて、重心がメインクラスター B の重心から最も遠い学習データをクラスターモデルの作成に使用する学習データとして選択する ( S 1 0 8 )。なお、図 6 において、右上にチェックマークが付されている学習データが、クラスターモデルの作成に使用する学習データとして選択されたものである。

40

## 【 0 0 6 4 】

図 3 に示すように、モデル学習部 1 5 b は、S 1 0 8 の処理の後、S 1 0 8 において選択した学習データを使用して学習を実行することによって、現在の対象のメインクラスター用のクラスターモデルを作成する ( S 1 0 9 )。ここで、モデル学習部 1 5 b は、ベースモデル 1 4 b を基にしてクラスターモデルを作成する。

## 【 0 0 6 5 】

文書クラスタリング部 1 5 a は、S 1 0 9 の処理の後、図 3 に示す動作の今回の実行に

50

において未だ S 1 0 3 の処理の対象にしていないメインクラスターが存在する場合には、図 3 に示す動作の今回の実行において未だ S 1 0 3 の処理の対象にしていないメインクラスターの 1 つを対象にして ( S 1 1 0 )、S 1 0 3 の処理を実行する。

【 0 0 6 6 】

モデル学習部 1 5 b は、S 1 0 9 の処理の後、図 3 に示す動作の今回の実行において未だ S 1 0 3 の処理の対象にしていないメインクラスターが存在しない場合には、図 3 に示す動作の今回の実行において新たに作成した全てのクラスターモデルを記憶部 1 4 に保存する ( S 1 1 1 )。

【 0 0 6 7 】

次いで、文書クラスタリング部 1 5 a は、図 3 に示す動作におけるメインクラスターのクラスタリングの結果をクラスタリング結果 1 4 d に保存して ( S 1 1 2 )、図 3 に示す動作を終了する。

【 0 0 6 8 】

次に、請求書データから特定の項目に対する値を抽出する場合の情報抽出システム 1 0 の動作について説明する。

【 0 0 6 9 】

図 7 は、請求書データから特定の項目に対する値を抽出する場合の情報抽出システム 1 0 の動作のフローチャートである。

【 0 0 7 0 】

利用者は、抽出対象データを用意し、用意した抽出対象データからの特定の項目に対する値の抽出を、操作部 1 1 から、または、図示していないコンピューターから通信部 1 3 を介して、情報抽出システム 1 0 に指示することができる。ここで、特定の項目とは、クラスターモデルの作成時に使用された学習データにおいて正解ラベルに対する項目、すなわち、請求書から抽出されることを利用者が希望する項目である。

【 0 0 7 1 】

情報抽出システム 1 0 の制御部 1 5 は、抽出対象データからの特定の項目に対する値の抽出が指示されると、図 7 に示す動作を実行する。

【 0 0 7 2 】

図 7 に示すように、文書クラスタリング部 1 5 a は、クラスタリング結果 1 4 d を使用して、抽出対象データが所属するメインクラスターを判定する ( S 1 2 1 )。

【 0 0 7 3 】

データ抽出実行部 1 5 c は、S 1 2 1 の処理の後、抽出対象データが所属するメインクラスターが S 1 2 1 において特定されたか否かを判断する ( S 1 2 2 )。

【 0 0 7 4 】

データ抽出実行部 1 5 c は、抽出対象データが所属するメインクラスターが S 1 2 1 において特定されたと S 1 2 2 において判断すると、抽出対象データが所属すると S 1 2 1 において特定されたメインクラスター用のクラスターモデルを使用して請求書データから特定の項目に対する値を抽出して ( S 1 2 3 )、図 7 に示す動作を終了する。

【 0 0 7 5 】

データ抽出実行部 1 5 c は、抽出対象データが所属するメインクラスターが S 1 2 1 において特定されなかった、すなわち、抽出対象データがいずれのメインクラスターにも所属しない外れ値であると S 1 2 2 において判断すると、抽出対象データに適合するクラスターモデルが存在しないことを利用者に通知する ( S 1 2 4 )。ここで、利用者への通知の方法としては、例えば、抽出対象データからの特定の項目に対する値の抽出が操作部 1 1 から指示された場合には、表示部 1 2 における表示でも良いし、抽出対象データからの特定の項目に対する値の抽出が図示していないコンピューターから通信部 1 3 を介して指示された場合には、通信部 1 3 を介した、このコンピューターへの出力でも良い。

【 0 0 7 6 】

データ抽出実行部 1 5 c は、S 1 2 4 の処理の後、抽出対象データに最も近いメインクラスター用のクラスターモデルを使用して抽出対象データから特定の項目に対する値を抽

10

20

30

40

50

出して ( S 1 2 5 )、図 7 に示す動作を終了する。

【 0 0 7 7 】

なお、 S 1 2 3 または S 1 2 5 において抽出された値は、様々な用途に活用されることが可能である。例えば、 S 1 2 3 または S 1 2 5 において抽出された値は、抽出対象データの基になった請求書の画像ファイルのファイル名に使用されても良い。

【 0 0 7 8 】

次に、クラスターモデルを更新する場合の情報抽出システム 1 0 の動作について説明する。

【 0 0 7 9 】

図 8 は、クラスターモデルを更新する場合の情報抽出システム 1 0 の動作の一部のフローチャートである。図 9 は、図 8 に示す動作の続きの動作のフローチャートである。 10

【 0 0 8 0 】

利用者は、クラスターモデルの更新のための学習データ ( 以下「追加データ」という。 ) を用意し、用意した追加データを使用した学習を、操作部 1 1 から、または、図示していないコンピューターから通信部 1 3 を介して、情報抽出システム 1 0 に指示することができる。ここで、利用者は、例えば、クラスターモデルを使用して抽出された値が適切ではなかった請求書データに、正解ラベルを付与することによって、追加データとしても良い。

【 0 0 8 1 】

情報抽出システム 1 0 の制御部 1 5 は、追加データを使用した学習が指示されると、図 8 および図 9 に示す動作を実行する。 20

【 0 0 8 2 】

図 8 および図 9 に示すように、文書クラスタリング部 1 5 a は、クラスタリング結果 1 4 d を使用して、追加データが所属するメインクラスターを判定する ( S 1 4 1 ) 。

【 0 0 8 3 】

文書クラスタリング部 1 5 a は、 S 1 4 1 の処理の後、追加データが所属するメインクラスターが S 1 4 1 において特定されたか否かを判断する ( S 1 4 2 ) 。

【 0 0 8 4 】

文書クラスタリング部 1 5 a は、追加データが所属するメインクラスターが S 1 4 1 において特定されたと S 1 4 2 において判断すると、追加データが所属すると S 1 4 1 において特定されたメインクラスターに追加データを追加する ( S 1 4 3 ) 。

【 0 0 8 5 】

次いで、文書クラスタリング部 1 5 a は、追加データが所属すると S 1 4 1 において特定されたメインクラスターを対象にする ( S 1 4 4 ) 。

【 0 0 8 6 】

次いで、文書クラスタリング部 1 5 a は、現在の対象のメインクラスターにおけるサブクラスター最適数をクラスター数自動推定法によって確認する ( S 1 4 5 ) 。

【 0 0 8 7 】

次いで、文書クラスタリング部 1 5 a は、 S 1 4 5 において確認したサブクラスター最適数がサブクラスター上限数以下であるか否かを判断する ( S 1 4 6 ) 。

【 0 0 8 8 】

文書クラスタリング部 1 5 a は、 S 1 4 6 の処理の後、 S 1 4 5 において確認したサブクラスター最適数がサブクラスター上限数以下ではないと S 1 4 6 において判断すると、 S 1 4 5 において確認したサブクラスター最適数からサブクラスター上限数を差し引いた数のサブクラスターを現在の対象のメインクラスターから分離する ( S 1 4 7 ) 。

ここで、文書クラスタリング部 1 5 a は、重心が現在の対象のメインクラスターの重心から遠いサブクラスターを優先して現在の対象のメインクラスターから分離する。

【 0 0 8 9 】

文書クラスタリング部 1 5 a は、 S 1 4 7 の処理の後、 S 1 4 7 において現在の対象のメインクラスターから分離したサブクラスターによって新たにメインクラスターを生成す 50

る（S 1 4 8）。すなわち、文書クラスタリング部 1 5 a は、S 1 4 7 において現在の対象のメインクラスターから分離したサブクラスターを新たなメインクラスターにする。

【0 0 9 0】

文書クラスタリング部 1 5 a は、S 1 4 5 において確認した最適数がサブクラスター上限数以下であると S 1 4 6 において判断するか、S 1 4 8 の処理が終了すると、現在の対象のメインクラスター内の学習データの群をサブクラスター最適数でクラスタリングすることによって、現在の対象のメインクラスター内の学習データのそれぞれをいずれかのサブクラスターに分ける（S 1 4 9）。

【0 0 9 1】

次いで、モデル学習部 1 5 b は、現在の対象のメインクラスター内のサブクラスターから、クラスターモデルの作成に使用する学習データを選択する（S 1 5 0）。ここで、モデル学習部 1 5 b は、現在の対象のメインクラスター内のサブクラスターのうち、重心が現在の対象のメインクラスターの重心に最も近いサブクラスターにおいて、重心が現在の対象のメインクラスターの重心に最も近い学習データを、クラスターモデルの作成に使用する学習データとして選択する。また、モデル学習部 1 5 b は、現在の対象のメインクラスター内のサブクラスターのうち、重心が現在の対象のメインクラスターの重心に最も近いサブクラスター以外のサブクラスターのそれぞれにおいて、重心が現在の対象のメインクラスターの重心から最も遠い学習データを、クラスターモデルの作成に使用する学習データとして選択する。

10

【0 0 9 2】

モデル学習部 1 5 b は、S 1 5 0 の処理の後、S 1 5 0 において選択された学習データを使用して学習を実行することによって、現在の対象のメインクラスター用のクラスターモデルを作成する（S 1 5 1）。ここで、モデル学習部 1 5 b は、ベースモデル 1 4 b を基にしてクラスターモデルを作成する。

20

【0 0 9 3】

文書クラスタリング部 1 5 a は、S 1 5 1 の処理の後、図 8 および図 9 に示す動作の今回の実行において新たに生成したメインクラスターに、図 8 および図 9 に示す動作の今回の実行において未だ S 1 4 5 の処理の対象にしていないメインクラスターが存在する場合には、図 8 および図 9 に示す動作の今回の実行において新たに生成したメインクラスターのうち、図 8 および図 9 に示す動作の今回の実行において未だ S 1 4 5 の処理の対象に

30

【0 0 9 4】

データ抽出実行部 1 5 c は、S 1 5 1 の処理の後、図 8 および図 9 に示す動作の今回の実行において新たに生成したメインクラスターに、図 8 および図 9 に示す動作の今回の実行において未だ S 1 4 5 の処理の対象にしていないメインクラスターが存在しない場合には、図 8 および図 9 に示す動作の今回の実行において新たに作成した全てのクラスターモデルが、クラスターモデル自身の対象のメインクラスターに含まれる全ての学習データに対して特定の程度以上に高い精度で特定の項目に対する値を抽出することができるか否かを判断する（S 1 5 3）。ここで、データ抽出実行部 1 5 c は、高い精度で特定の項目に

40

【0 0 9 5】

モデル学習部 1 5 b は、図 8 および図 9 に示す動作の今回の実行において新たに作成した全てのクラスターモデルが、クラスターモデル自身の対象のメインクラスターに含まれる全ての学習データに対して特定の程度以上に高い精度で特定の項目に対する値を抽出できると S 1 5 3 において判断すると、追加データが所属すると S 1 4 1 において特定されたメインクラスター用のクラスターモデルを記憶部 1 4 から削除し（S 1 5 4）、図 8 および図 9 に示す動作の今回の実行において新たに作成した全てのクラスターモデルを記憶部 1 4 に保存する（S 1 5 5）。

50

## 【 0 0 9 6 】

文書クラスタリング部 1 5 a は、図 8 および図 9 に示す動作の今回の実行において新たに作成したいずれかのクラスターモデルが、クラスターモデル自身の対象のメインクラスターに含まれるいずれかの学習データに対して特定の程度以上に高い精度で特定の項目に対する値を抽出することができないと S 1 5 3 において判断すると、図 8 および図 9 に示す動作の今回の実行におけるこれまでのクラスタリングの結果を全て廃棄する ( S 1 5 6 )。したがって、文書クラスタリング部 1 5 a は、追加データが現在所属するメインクラスターから追加データを分離する。

## 【 0 0 9 7 】

文書クラスタリング部 1 5 a は、追加データが所属するメインクラスターが S 1 4 1 において特定されなかった、すなわち、追加データがいずれのメインクラスターにも所属しない外れ値であると S 1 4 2 において判断するか、 S 1 5 6 の処理が終了すると、追加データによって新たにメインクラスターを生成する ( S 1 5 7 )。

10

## 【 0 0 9 8 】

モデル学習部 1 5 b は、 S 1 5 7 の処理の後、追加データを使用して学習を実行することによって、追加データが所属するメインクラスター用のクラスターモデルを作成する ( S 1 5 8 )。ここで、モデル学習部 1 5 b は、ベースモデル 1 4 b を基にしてクラスターモデルを作成する。

## 【 0 0 9 9 】

モデル学習部 1 5 b は、 S 1 5 8 の処理の後、 S 1 5 8 において新たに作成したクラスターモデルを記憶部 1 4 に保存する ( S 1 5 9 )。

20

## 【 0 1 0 0 】

文書クラスタリング部 1 5 a は、 S 1 5 5 または S 1 5 9 の処理の後、図 8 および図 9 に示す動作におけるメインクラスターのクラスタリングの結果をクラスタリング結果 1 4 d に保存して ( S 1 6 0 )、図 8 および図 9 に示す動作を終了する。

## 【 0 1 0 1 】

以上に説明したように、情報抽出システム 1 0 は、メインクラスター毎に情報抽出モデルとしてのクラスターモデルを作成する ( S 1 0 9、 S 1 5 1 および S 1 5 8 ) ので、クラスターモデル毎の特徴を単純化することができ、その結果、クラスターモデル毎に必要な学習データの数を低減することができる。したがって、情報抽出システム 1 0 は、クラスターモデルの作成のための計算量を低減することができる。

30

## 【 0 1 0 2 】

情報抽出システム 1 0 は、クラスターモデルの作成に使用する学習データをサブクラスター毎に選択し ( S 1 0 8 および S 1 5 0 )、選択した学習データを使用して学習を実行することによって、メインクラスター毎のクラスターモデルを作成する ( S 1 0 9 および S 1 5 1 ) ので、クラスターモデル毎に必要な学習データの数を低減することができ、その結果、クラスターモデルの作成のための計算量を低減することができる。

## 【 0 1 0 3 】

情報抽出システム 1 0 は、重心がメインクラスターの重心に最も近いサブクラスターにおいて、重心がメインクラスターの重心に最も近い学習データを、クラスターモデルの作成に使用する学習データとして選択する ( S 1 0 8 および S 1 5 0 ) ので、メインクラスターの特徴を最も強く表す学習データを使用してクラスターモデルを作成することができ、その結果、メインクラスターの特徴が適切に反映されたクラスターモデルを作成することができる。

40

## 【 0 1 0 4 】

情報抽出システム 1 0 は、重心がメインクラスターの重心に最も近いサブクラスター以外のサブクラスターのそれぞれにおいて、重心がメインクラスターの重心から最も遠い学習データを、クラスターモデルの作成に使用する学習データとして選択する ( S 1 0 8 および S 1 5 0 ) ので、メインクラスターにおいて広範囲に散らばった学習データを使用してクラスターモデルを作成することができ、その結果、メインクラスターの特徴が適切に

50

反映されたクラスターモデルを作成することができる。

【0105】

情報抽出システム10は、メインクラスターにおけるサブクラスター最適数がサブクラスター上限数を超える場合に、サブクラスター最適数からサブクラスター上限数を差し引いた数のサブクラスターを、このメインクラスターから分離する（S105およびS147）ので、クラスターモデル毎に必要な学習データの数を実減することができ、その結果、クラスターモデルの作成のための計算量を低減することができる。

【0106】

情報抽出システム10は、クラスター最適数からクラスター上限数を差し引いた数のサブクラスターをメインクラスターから分離する場合に、重心がこのメインクラスターの重心から遠いサブクラスターを優先して、このメインクラスターから分離する（S105およびS147）ので、メインクラスターの特徴を強く表す学習データを使用して情報抽出モデルを作成することができ、その結果、メインクラスターの特徴が適切に反映された情報抽出モデルを作成することができる。

10

【0107】

情報抽出システム10は、クラスターモデルの作成のための計算量を低減することができるので、例えば、一般的なPCの計算リソースでも深層学習の学習処理を実行することができる。したがって、情報抽出システム10は、情報を抽出する対象の文書が、例えば個人情報や取引情報など、保護すべき情報が含まれる、例えば請求書などの文書である場合に、文書のデータをローカル環境外にアップロードすることなく、ローカル環境内の一般的なPCでクラスターモデルを作成することができる。

20

【0108】

以上においては、モデル学習部15bは、クラスターモデルを更新する場合に、ベースモデル14bを基にしてクラスターモデルを作成する。しかしながら、モデル学習部15bは、クラスターモデルを更新する場合に、更新の対象のクラスターモデルが記憶部14に既に記憶されている場合には、更新の対象のクラスターモデルを基にして新たなクラスターモデルを作成しても良い。

【0109】

以上においては、情報抽出システム10は、請求書のデータから情報を抽出する。しかしながら、情報抽出システム10は、請求書の場合と同様にして、例えば答案用紙など、請求書以外の種類の文書のデータから情報を抽出することが可能である。なお、情報抽出システム10は、文書の種類毎のベースモデルを使用しても良いし、複数の種類の文書に共通のベースモデルを使用しても良い。ここで、情報抽出システム10は、文書の種類毎のベースモデルを使用する方が、複数の種類の文書に共通のベースモデルを使用するよりも、情報の抽出の精度を向上することができる。しかしながら、情報抽出システム10は、複数の種類の文書に共通のベースモデルを使用する方が、文書の種類毎のベースモデルを使用するよりも、ベースモデルの用意の労力を低減することができる。

30

【符号の説明】

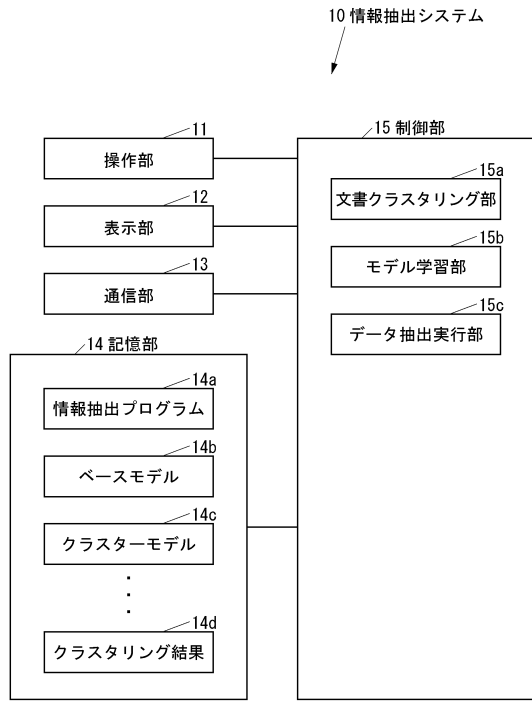
【0110】

- 10 情報抽出システム
- 14 a 情報抽出プログラム
- 14 c クラスターモデル（情報抽出モデル）
- 15 a 文書クラスタリング部
- 15 b モデル学習部
- 20 情報抽出モデル

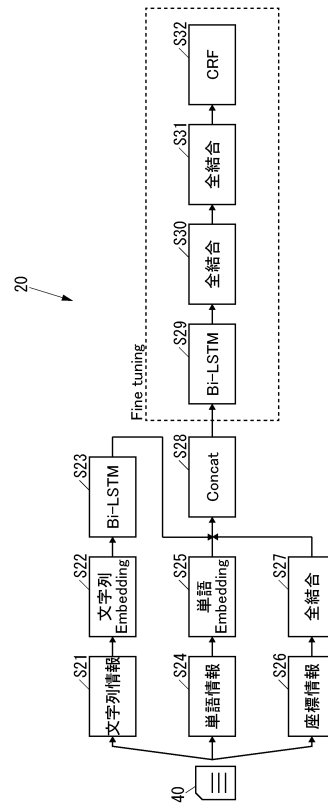
40

【 図 面 】

【 図 1 】



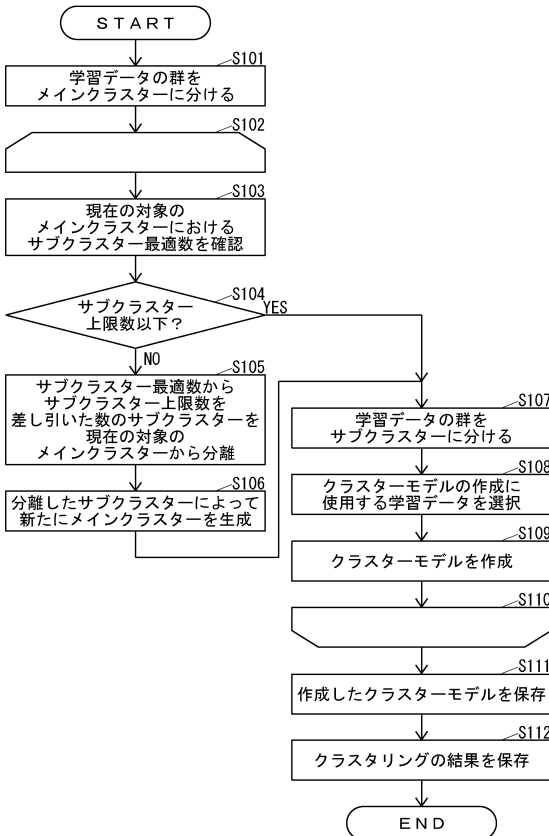
【 図 2 】



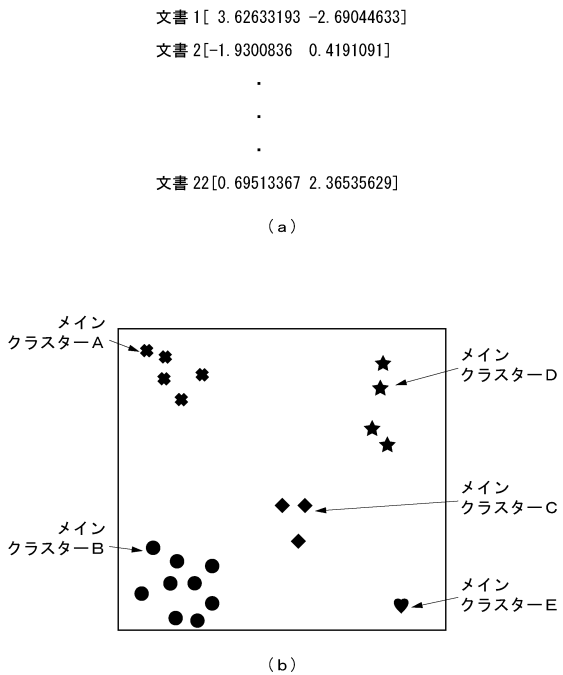
10

20

【 図 3 】



【 図 4 】

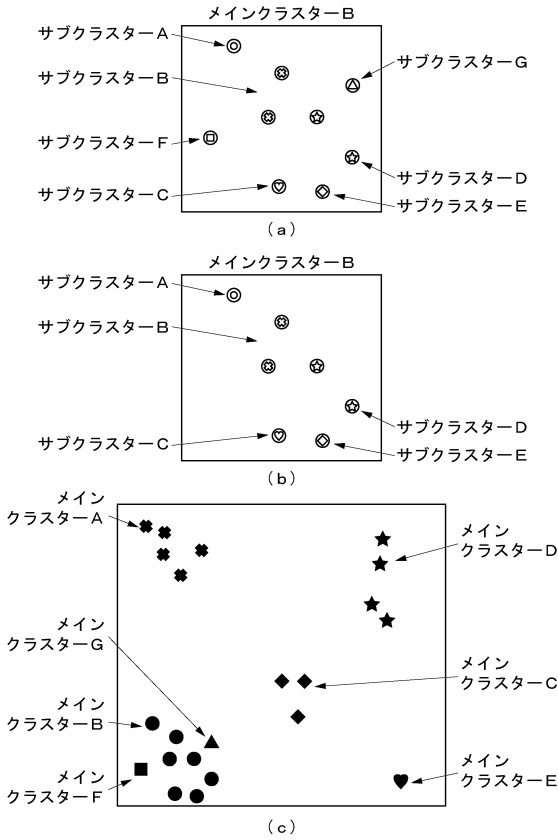


30

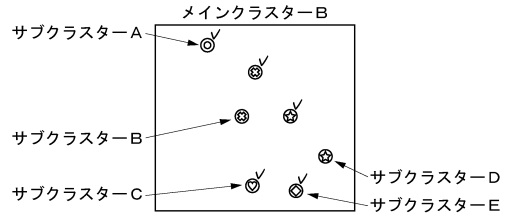
40

50

【 図 5 】



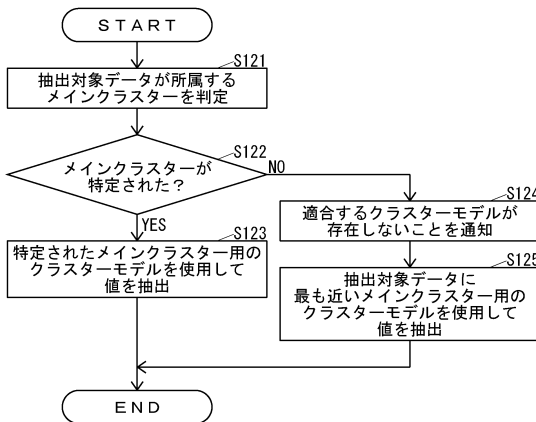
【 図 6 】



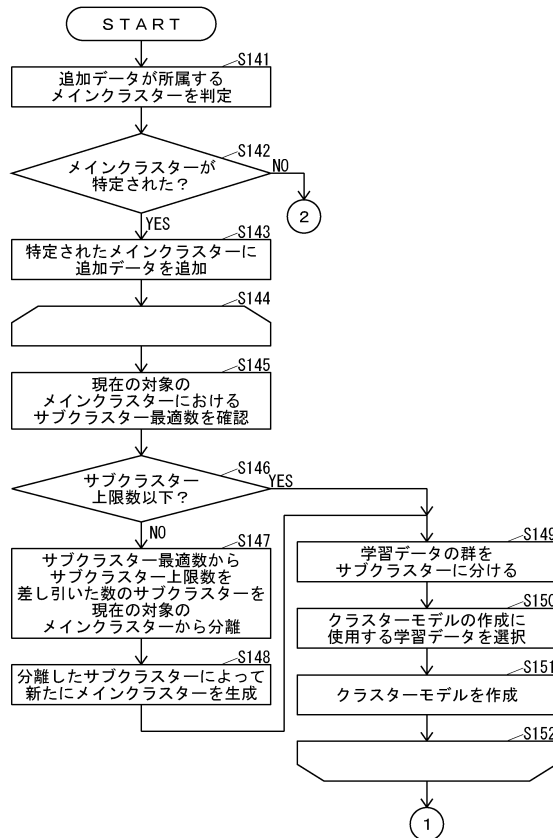
10

20

【 図 7 】



【 図 8 】



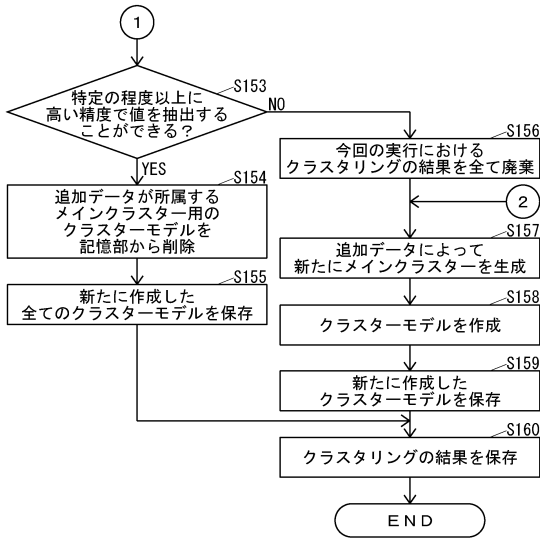
30

40

50



【 図 9 】



10

20

30

40

50