

(19) 日本国特許庁(JP)

(12) 公開特許公報(A)

(11) 特許出願公開番号

特開2017-146745
(P2017-146745A)

(43) 公開日 平成29年8月24日(2017.8.24)

(51) Int.Cl.	F I	テーマコード (参考)
G06F 17/30 (2006.01)	G06F 17/30	210D
G06N 99/00 (2010.01)	G06F 17/30	170B
	G06F 17/30	220C
	G06N 99/00	153

審査請求 未請求 請求項の数 10 O L (全 18 頁)

(21) 出願番号 特願2016-27352 (P2016-27352)
(22) 出願日 平成28年2月16日 (2016.2.16)

(71) 出願人 000001007
キヤノン株式会社
東京都大田区下丸子3丁目30番2号
(74) 代理人 100114775
弁理士 高岡 亮一
(74) 代理人 100121511
弁理士 小田 直
(72) 発明者 宮内 崇
東京都大田区下丸子3丁目30番2号 キ
ヤノン株式会社内

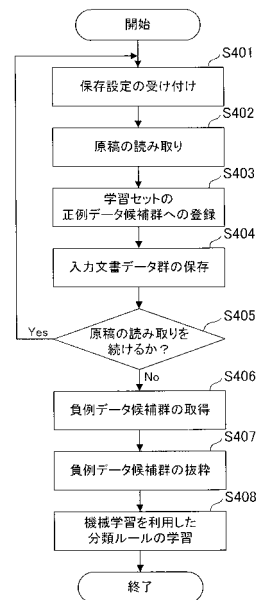
(54) 【発明の名称】 情報処理装置、制御方法、情報処理システム、およびプログラム

(57) 【要約】

【課題】高精度な分類ルールの構築を可能とする高品質な教師データを効率的に生成する情報処理装置を提供する。

【解決手段】MFP101は、分類するクラスごとの画像データを正例データとして受け付け、受け付けた前記画像データに付与された文書ファイル情報に含まれる情報のうち少なくとも1つが一致する文書ファイル情報が付与された画像データを負例データとして取得し、前記正例データ及び負例データを用いて、画像データを種別ごとに分類するために用いる分類ルールを生成する。

【選択図】図4



【特許請求の範囲】**【請求項 1】**

分類する種別ごとの画像データを正例データとして受け付ける受付手段と、
受け付けた前記画像データに付与されたファイル情報に含まれる情報のうち少なくとも
1つが一致するファイル情報が付与された画像データを負例データとして取得する取得手
段と、

前記正例データ及び負例データを用いて、画像データを種別ごとに分類するために用い
る分類ルールを生成する生成手段と、を備える

ことを特徴とする情報処理装置。

【請求項 2】

前記ファイル情報は、少なくとも画像データのメタ情報を含む

ことを特徴とする請求項 1 に記載の情報処理装置。

【請求項 3】

前記メタ情報は、画像データのタイトル、作成者名、ファイル形式、作成デバイス、生
成日時、または当該画像データが含むキーワードのうち少なくとも 1つを含む

ことを特徴とする請求項 2 に記載の情報処理装置。

【請求項 4】

前記取得手段は、前記正例データに付与されたファイル情報に含まれる情報のうち、当
該正例データにおいて共起性が高い情報を含むファイル情報が付与された画像データを負
例データとして取得する

ことを特徴とする請求項 1 乃至 3 のいずれか 1 項に記載の情報処理装置。

【請求項 5】

前記取得手段は、前記負例データとして取得した画像データのうち、前記正例データに
付与されたファイル情報に含まれる情報との一致率が高いファイル情報が付与された画像
データを前記負例データとして使用しない

ことを特徴とする請求項 1 乃至 4 のいずれか 1 項に記載の情報処理装置。

【請求項 6】

前記取得手段は、前記負例データとして取得した画像データを、当該画像データのキー
ワードに基づき種別ごとに分類し、それぞれの種別において分類された画像データの個数
が上限の値よりも多い場合は、当該分類された画像データの個数が当該上限の値以下とな
るように当該画像データを削除する

ことを特徴とする請求項 1 乃至 5 のいずれか 1 項に記載の情報処理装置。

【請求項 7】

前記取得手段は、

受け付けた前記画像データに付与されたファイル情報に含まれる情報を項目ごとに表
示する画面を有し、

前記画面において指定された項目ごとの値が当該画面において指定された条件を満た
す画像データを前記負例データとして取得する

ことを特徴とする請求項 1 乃至 6 のいずれか 1 項に記載の情報処理装置。

【請求項 8】

情報処理装置とサーバとを備えるシステムであって、

前記情報処理装置は、

分類する種別ごとの画像データを正例データとして受け付ける受付手段と、

受け付けた前記画像データに付与されたファイル情報に含まれる情報のうち少なくと
も 1つが一致するファイル情報が付与された画像データを、前記サーバから負例データと
して取得する取得手段と、

前記正例データ及び負例データを用いて、画像データを種別ごとに分類するために用
いる分類ルールを生成する生成手段と、を備え、

前記サーバは、

前記情報処理装置の要求に応じて、画像データを前記情報処理装置に送信する送信手

10

20

30

40

50

段を備える

ことを特徴とする情報処理システム。

【請求項 9】

分類する種別ごとの画像データを正例データとして受け付ける受付工程と、

受け付けた前記画像データに付与されたファイル情報に含まれる情報のうち少なくとも一つが一致するファイル情報が付与された画像データを、負例データとして取得する取得工程と、

前記正例データ及び負例データを用いて、画像データを種別ごとに分類するために用いる分類ルールを生成する生成工程と、を備える

ことを特徴とする情報処理装置の制御方法。

10

【請求項 10】

請求項 1 乃至 7 のいずれか 1 項に記載の情報処理装置が備える各手段としてコンピュータを機能させるためのプログラム。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、情報処理装置、制御方法、情報処理システム、およびプログラムに関する。

【背景技術】

【0002】

文書を扱うワークフローの効率化を実現する技術の 1 つとして、機械学習を利用した画像分類が提案されている。機械学習を利用した画像分類は、一般的に学習と分類（運用）の 2 つのプロセスを有し、画像データ群（教師データ、学習セット）を与えることで分類ルールを学習によって構築し、構築した分類ルールに基づいて入力画像を分類する。

20

【0003】

データを複数の種別に分類するには、データが学習セットとして与えた種別のいずれであるかを分類できればよい場合が多いが、文書を扱う場合には、学習したい種別でもない文書を「該当なし」と分類したいというニーズがある。例えば、MFP のスキャナによって大量の文書が読み込まれた際に、特定の種別の文書のみをあらかじめ指定されたフォルダに格納し、その他の種別の文書は「該当なし」に分類し、まとめて一か所のフォルダに格納するようなケースが考えられる。

30

【0004】

機械学習では、学習セットとして与えられたデータに基づいて分類ルールを構築するため、学習セット内のデータは、運用時に入力されるデータと特徴量が近い方がよい。また、「該当なし」の分類を実現するには、本来分類したい種別のデータ（正例データ）に加えて、「その他」の種別であるデータ（負例データ）を用意した方がよく、負例データとしては、実際に分類時に入力される可能性の高い文書を用意することが望ましい。

【0005】

しかし、ユーザが多くの種別の文書を扱っている場合に、本来分類したい種別のデータ（正例データ）以外の大量な文書データを負例データとして用意するのは、ユーザにとって大きな負担となってしまう。また、機械学習では正例データと負例データに同じ種別のデータが混在していると正しく分類ルールを構築することができない。そのため、初めて学習セットを用意する際だけでなく、正例データの種別を追加する度に、負例データの中に新しく追加した種別の正例データが混在していないかを確認する必要がある。

40

【0006】

特許文献 1 は、正例の文書（正例データ）から特徴語を抽出し、ファイルサーバから取り出した負例候補文書から、当該正例の特徴語をなるべく含まず、かつ当該正例の特徴語以外の特徴語を多く含む文書を負例として選択する文書分類システムを開示している。

【先行技術文献】

【特許文献】

【0007】

50

【特許文献1】特開2014-96086号公報

【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、特許文献1のように正例データと同じ種別である可能性の低い文書データを除くだけでは、効率よく高精度な分類器を構築することは困難である。一般に、学習セットのデータ量に応じて学習時間が増加する。このため、例えば、ユーザが用意したデータからその場で分類ルールを構築するシステムの場合には、学習セットを絞り込む必要がある。しかし、ファイルサーバからランダムに一定数のファイルを選ぶ等、学習に利用するデータを一律に削減してしまうと、実際に分類時に入力される可能性の高いデータも減り、分類精度が低下してしまう。

10

【0009】

本発明は、高精度な分類ルールの構築を可能とする高品質な教師データを効率的に生成する情報処理装置の提供を目的とする。

【課題を解決するための手段】

【0010】

本発明の一実施形態の情報処理装置は、分類する種別ごとの画像データを正例データとして受け付ける受付手段と、受け付けた前記画像データに付与されたファイル情報に含まれる情報のうち少なくとも1つが一致するファイル情報が付与された画像データを負例データとして取得する取得手段と、前記正例データ及び負例データを用いて、画像データを種別ごとに分類するために用いる分類ルールを生成する生成手段と、を備える。

20

【発明の効果】

【0011】

本発明の情報処理装置によれば、高精度な分類ルールの構築を可能とする高品質な教師データを生成することができる。

【図面の簡単な説明】

【0012】

【図1】第1実施形態における情報処理システム構成を示す図である。

【図2】MFPの構成例を示す図である。

【図3】サーバのハードウェア構成の一例を示す図である。

30

【図4】MFPが分類ルールを学習する処理を説明するためのフローチャートである。

【図5】負例データを構築する処理を示すフローチャートである。

【図6】文書ファイル情報の一例を示す図である。

【図7】負例データ候補群の絞り込み条件をユーザが編集する画面を示す図である。

【図8】文書ファイル情報による絞り込み結果の例を示す図である。

【図9】学習セットを用いた機械学習の一例を示す図である。

【図10】特徴量の算出方法について説明する図である。

【図11】画像データからパッチ画像を切り出す方法について説明する図である。

【図12】学習セットを生成し、分類ルールを学習する処理を説明する図である。

【発明を実施するための形態】

40

【0013】

以下、本発明を実施するための形態について図面などを参照して説明する。

(第1実施形態)

図1は、本実施形態における情報処理システム構成を示す図である。

第1実施形態における情報処理システムは、情報処理装置であるMFP101及びサーバ102を備える。

【0014】

LAN103には、MFP101が接続されている。また、LAN103は、インターネット104を経由してサービスを提供するサーバ102と接続されている。MFP101及びサーバ102は、LAN103を介して互いに接続されており、画像データや各種

50

情報の送受信を行う。なお、MFP101とサーバ102とは、互いに接続され、画像データや各種情報を送受信できればよく、有線により直接接続されていてもよく、また、無線通信により接続されていてもよい。

【0015】

サーバ102は、MFP101から入力された画像データを格納し、MFP101から指定された条件を満たす画像データをMFP101に送信するファイルサーバとして機能する。なお、本実施形態では、分類ルールを学習する際に使用する学習セットの生成や、当該学習セットを用いた分類ルールの構築はMFP101が実行するが、同様の処理をサーバ102が実行してもよい。

【0016】

図2は、MFP101の構成例を示す図である。

図2(A)は、MFP101のハードウェア構成の一例を示す図である。図2(A)に示すように、MFP101は、コントローラ20、画像読取部201、画像出力部205、及び操作部207を備える。コントローラ20は、装置制御部200、画像処理部202、記憶部203、CPU204、及びネットワークI/F部206を備える。

【0017】

装置制御部200は、MFP101内およびネットワークI/F部206を経由した外部とのデータの受け渡しや、操作部207からの操作の受け付けを行う。画像読取部201は、原稿の画像を読み取り、画像データをコントローラ20に出力する。画像処理部202は、画像読取部201や外部から入力される画像データを含む印刷情報を中間情報(以下「オブジェクト」と呼ぶ)に変換し、記憶部203のオブジェクトバッファに格納する。

【0018】

オブジェクトは、テキスト、グラフィック、イメージの属性を持つ。さらに、オブジェクトバッファに格納したオブジェクトに基づきビットマップデータを生成し、記憶部203のバッファに格納する。その際、色変換処理、濃度調整処理、トナー総量制御処理、ビデオカウント処理、プリンタガンマ補正処理、ディザなどの疑似中間調処理を行う。記憶部203は、ROM(Read Only Memory)、RAM(Random Access Memory)、HDD(Hard Disk Drive)などから構成される。

【0019】

ROMは、CPU204が実行する各種の制御プログラムや画像処理プログラムを格納する。RAMは、CPU204がデータや各種情報を格納する参照領域や作業領域として用いられる。また、RAMおよびHDDは、上述したオブジェクトバッファなどに用いられる。コントローラ20は、RAMおよびHDD上で画像データを蓄積し、ページのソートや、ソートされた複数ページにわたる原稿を蓄積し、複数部プリント出力を行う。

【0020】

なお、記憶部203を構成するHDDは、ファイルサーバとして機能し、画像読取部201やネットワークI/F部206経由で入力された画像データが蓄積されているものとする。画像出力部205は、記録紙などの記録媒体にカラー画像を形成して出力する。ネットワークI/F部206は、MFP101をLAN103に接続し、インターネット104や他の装置との間で各種情報を送受信する。操作部207は、タッチパネルや操作ボタンを備え、ユーザからの操作を受け付けて装置制御部200へ該操作の情報を送信する。

【0021】

図2(B)は、MFP101の外観の一例を示す図である。画像読取部201は、複数の受光画素を有している。各受光画素の感度が夫々異なっていると、たとえ原稿上の各画素の濃度が同じであったとしても、各画素が夫々違う濃度であると認識されてしまう。そのため、画像読取部201では、最初に白板(一様に白い板)を露光走査し、露光走査して得られた反射光の量を電気信号に変換してコントローラに出力している。

10

20

30

40

50

【 0 0 2 2 】

なお、画像処理部 2 0 2 内には、各受光画素から得られた電気信号を元に、各受光画素の感度の違いを認識し、その違いを利用して、原稿上の画像をスキャンして得られた電気信号の値を補正する、公知のシェーディング補正処理部を有する。さらに、シェーディング補正部は、コントローラ内の CPU 2 0 4 からゲイン調整の情報を受取ると、当該情報に応じたゲイン調整を行う。

【 0 0 2 3 】

ゲイン調整は、原稿を露光走査して得られた電気信号の値を、どのように 0 ~ 2 5 5 の輝度信号値に割り付けるかを調整するために用いられる。このゲイン調整により、原稿を露光走査して得られた電気信号の値を高い輝度信号値に変換したり、低い輝度信号値に変換したりすることができるようになっている。すなわち、ゲイン調整により、読み取り信号のダイナミックレンジの調整が可能である。

10

【 0 0 2 4 】

続いて、この原稿上の画像をスキャンする構成について説明する。

画像読取部 2 0 1 は、原稿上の画像を露光走査して得られた反射光を受光画素に入力することで画像の情報を電気信号に変換する。さらに電気信号をレッド R , グリーン G , およびブルー B の各色からなる輝度信号に変換し、当該輝度信号を画像としてコントローラ 2 0 に対して出力する。

【 0 0 2 5 】

なお、原稿は原稿フィーダ 2 1 1 のトレイ 2 1 2 にセットされる。ユーザが操作部 2 0 7 から読み取り開始を指示すると、コントローラ 2 0 から画像読取部 2 0 1 に原稿読み取り指示が与えられる。画像読取部 2 0 1 は、この指示を受けると原稿フィーダ 2 1 1 のトレイ 2 1 2 から原稿を 1 枚ずつフィードして、原稿の読み取り動作を行う。なお、原稿の読み取り方法は、原稿フィーダ 2 1 1 による自動送り方式に限られるものではなく、原稿を不図示のガラス面上に載置し露光部を移動させることで原稿の走査を行う方法であってもよい。

20

【 0 0 2 6 】

画像出力部 2 0 5 は、コントローラ 2 0 から受取った画像を用紙上に形成する画像形成デバイスである。なお、本実施形態では、画像形成方式は、感光体ドラムや感光体ベルトを用いた電子写真方式であるが、これに限られるものではない。例えば、微少ノズルアレイからインクを吐出して用紙上に印字するインクジェット方式などであっても本発明は適用可能である。また、画像出力部 2 0 5 には、異なる用紙サイズまたは異なる用紙向きを選択可能とする複数の用紙カセット 2 1 3 乃至 2 1 5 が設けられている。排紙トレイ 2 1 6 には印字後の用紙が排出される。

30

【 0 0 2 7 】

図 3 は、サーバのハードウェア構成の一例を示す図である。

サーバ 1 0 2 は、CPU 3 0 1、RAM 3 0 2、ROM 3 0 3、ネットワーク I / F 部 3 0 4、HDD 3 0 5、及びデータバス 3 0 6 を備える。CPU 3 0 1 は、ROM 3 0 3 に記憶された制御プログラムを読み出して RAM 3 0 2 にロードし、各種制御処理を実行する。RAM 3 0 2 は、CPU 3 0 1 の実行するプログラムや、ワークメモリ等の一時記憶領域として用いられる。

40

【 0 0 2 8 】

ネットワーク I / F 部 3 0 4 は、サーバ 1 0 2 をインターネット 1 0 4 に接続し、他の装置との間で各種情報を送受信する。HDD 3 0 5 は、画像データや特徴量データ、各種プログラム等を格納する。ネットワーク I / F 部 3 0 4 を介して MFP 1 0 1 から受信した画像データは、データバス 3 0 6 を介して CPU 3 0 1、RAM 3 0 2、及び ROM 3 0 3 に送受信される。

【 0 0 2 9 】

CPU 3 0 1 が ROM 3 0 3 や HDD 3 0 5 に格納された画像処理プログラムを実行することによって、画像データに対する画像処理が実現される。また、HDD 3 0 5 は、ネ

50

ットワーク I / F 部 3 0 4 を介して M F P 1 0 1 以外の外部装置からもデータの入力が可能であり、すでに文書の画像データを含む大量のファイルが格納されているものとする。

【 0 0 3 0 】

< 第 1 実施形態の詳細説明 >

図 4 は、学習セットを生成し、分類ルールを学習する処理を説明するフローチャートである。

図 4 に示す処理は、M F P 1 0 1 およびサーバ 1 0 2 にて実行される。M F P 1 0 1 において実行される処理は、C P U 2 0 4 が記憶部 2 0 3 に格納されている処理プログラムをロードして実行することにより実現される。また、サーバ 1 0 2 において実行される処理は、C P U 3 0 1 が H D D 3 0 5 に格納されている処理プログラムを R A M 3 0 2 にロ

10

【 0 0 3 1 】

なお、本実施形態では、ユーザが M F P 1 0 1 を用いて文書（原稿）をスキャンし、その種別毎に文書の画像データをサーバ 1 0 2 に格納するというワークフローの中で、同時に M F P 1 0 1 内で画像データの分類ルールを学習するシステムを想定している。このように、文書を扱うワークフローに機械学習を利用した分類ルールを応用すると、スキャナを備えた M F P などの入力機器から入力された文書の格納先や配布先の自動決定、ファイル名の自動付与などが可能になる。また、ユーザ毎に用意した文書から学習することで、個別にカスタマイズされた分類ルールを構築することも可能になる。

【 0 0 3 2 】

20

なお、文書のスキャン及びサーバへの格納と、分類ルールの学習を行うタイミングは上記のワークフローに限られるものではなく、文書のスキャン及びサーバへの格納と、分類ルールの学習が別々に実行されてもよい。第 2 実施形態では、すでにサーバ 1 0 2 に格納されたデータを分類ルールの学習時に取得する場合について説明する。また、分類ルールの学習は、データを読み込んだ M F P 1 0 1 で必ずしも行う必要はなく、例えば画像データを格納したサーバで本実施形態の分類に係る処理を行ってもよい。

【 0 0 3 3 】

ステップ S 4 0 1 において、M F P 1 0 1 は、ユーザから操作部 2 0 7 経由で画像データの保存設定を受付ける。なお、画像データの保存設定は、M F P 1 0 1 において読み込んだ画像データの保存先を示すフォルダのパスや、保存時のファイル名、ファイル形式などのことである。

30

【 0 0 3 4 】

ステップ S 4 0 2 において、M F P 1 0 1 は、操作部 2 0 7 からユーザの指示を受け付けると、原稿フィード 2 1 1 のトレイ 2 1 2 から原稿を 1 枚ずつフィードして、画像読取部 2 0 1 で原稿を読み取る。なお、本実施形態では、トレイ 2 1 2 にセットされる原稿は、同一種別の文書とする。また、同一種別の文書は、分類ルールにおいて同一のクラスに分類される文書とする。

【 0 0 3 5 】

ステップ S 4 0 3 において、M F P 1 0 1 は、ステップ S 4 0 2 で画像読取部 2 0 1 が読み込んだ画像データ群を、記憶部 2 0 3 に学習セットの正例データ候補群として格納する。画像データ群を格納する際には、各画像データに文書ファイル情報を付与する。文書ファイル情報は、後述する負例データ候補群の絞り込みに利用する。文書ファイル情報としては、タイトルや作成者名、ファイル形式、作成ツール、作成デバイス、変換ツール、キーワード、生成日時、更新日時など、アプリケーションで電子ファイルを作成する際に付与される一般的なメタ情報を利用する。

40

【 0 0 3 6 】

キーワードとは、文書ファイルの特徴を表す文字列群であり、本実施形態では、原稿を読み込む際に文字認識を行い、その結果を利用する。例えば、タイトルとなる最初のページの上部中央や、ヘッダーやフッター、表内の項目など文書の特徴的な位置にある文字列、他の文字と比べてフォントの異なる文字列など、特徴的な文字列をキーワードとして利

50

用する。

【 0 0 3 7 】

また、文書ファイル情報用のキーワード群と対応する項目とを辞書として保持しておき、文字認識を行った結果、辞書内のキーワードに当てはまる文字列が含まれる場合に、当該キーワードに対応する項目を文書ファイル情報のキーワードとして付与してもよい。文書ファイル情報用のキーワード群としては、「決裁書」や「申請書」、「注文書」といった一般的に利用される文書のタイトルや、企業名リストを利用する。

【 0 0 3 8 】

なお、文書ファイル情報は、上記のようなメタ情報に限定されるものではなく、文字認識の過程等で得られる文字列の位置情報やフォントサイズなどの文書構造情報を用いてもよい。また、MFP101での読み取り時に付与された読取解像度や色、割り当てなどのスキャン設定を用いてもよい。また、本実施形態では、文書ファイル情報と共に画像データ群が格納されるが、これに限定されるものではなく、例えば、読み込まれた画像データ群から算出される特徴量のデータを格納してもよい。

10

【 0 0 3 9 】

ステップS404において、MFP101は、ネットワークI/F部206を通じて画像読取部201で読み込まれた画像データ群をサーバ102に送信する。サーバ102は、LAN103およびインターネット104を経由してMFP101から画像データ群を受信する。サーバ102のCPU301は、ステップS401において設定された画像データの保存設定に基づき、受け付けた画像データをHDD305に記録する。

20

【 0 0 4 0 】

ステップS405において、MFP101は、原稿の読み取りを続けるか否かの指示を、操作部207を介してユーザから受け付ける。原稿の読み取りを続ける場合には、処理はステップS401に戻る。原稿の読み取りを続けない場合には、処理はステップS406に進む。なお、原稿の読み取りを続けるか否かの判断は、上記の方法に限るものではない。例えば、ステップS401での原稿の読み取り回数をカウントし、あらかじめ操作部207を介してユーザによって設定された原稿の読み取り回数に達するまで原稿の読み取りを続けてもよい。

【 0 0 4 1 】

ステップS406において、MFP101は、記憶部203に格納されている文書の画像データおよび、インターネット104およびLAN103を経由してサーバ102から取得した文書の画像データを、負例データ候補群として記憶部203に格納する。ステップS407において、MFP101は、ステップS406にて取得した負例データ候補群のファイルを抜粋する。負例データ候補群の抜粋処理の詳細については、図5を用いて後述する。

30

【 0 0 4 2 】

ステップS408において、MFP101は、ステップS403にて格納した正例データ候補群、およびステップS407にて格納した負例データ候補群を学習セットとして機械学習を利用した分類ルールの学習に用いる。本実施形態において、分類ルールの学習については、図9～11を用いて後述する。

40

【 0 0 4 3 】

< 負例データ候補群の抜粋処理に係る詳細説明 (ステップS407) >

運用時に入力される可能性の低い文書データは、運用時の分類精度に寄与しない無駄なデータとなってしまう。例えば、サーバからランダムに選ばれた50個の文書データの中に使われていないデータが5個、別の業務で利用するデータが10個含まれていた場合、分類ルールの構築に有効なデータが35個となってしまう。このように、ランダムにデータを取得するだけでは、実際に分類時に入力される可能性の高い文書を減らしてしまう要因となる。

【 0 0 4 4 】

また、データの冗長性を考慮していない場合も、実際に分類時に入力される可能性の高

50

い文書を減らしてしまう要因となる。例えば、負例データとして利用する文書データ50個が、5種類各10個の文書である場合と、50種類各1個の文書である場合には、前者の方が分類時に入力される可能性の高い文書を減らしてしまう。本実施形態では、負例データ候補群の抜粋処理により、高精度な分類を可能とする負例データを取得することが可能となる。

【0045】

図5は、負例データ候補群から負例データを構築する処理を示すフローチャートである。

詳細には、図5に示す処理は、分類ルールの構築に使用する学習セットの一部である負例データを、ステップS406にて取得した負例データ候補群から抜粋する処理である。図5に示す処理は、MFP101のCPU204が、記憶部203に格納されている処理プログラムをロードして実行することで実現される。

【0046】

ステップS501において、MFP101は、ステップS403で記憶部203に記録された正例データ候補群から、ステップS403で付与された文書ファイル情報およびユーザの指示に基づき、負例データ候補群の絞り込み条件を取得する。ステップS502において、MFP101は、ステップS501で取得した絞り込み条件に基づき、ステップS406で取得した負例データ候補群を絞り込む(抜粋する)。ステップS501およびステップS502の詳細については、図5～図8を用いて後述する。

【0047】

ステップS503において、MFP101は、ステップS502で抜粋した負例データ候補群から、冗長なデータを削減する。冗長なデータの特定には、例えば、文書ファイル情報の1つであるキーワードを特徴量としたクラスタリングを利用する。これは、同じキーワードで構成される文書は、同じ種別の文書である可能性が高いため、同じ種別の文書であると判定するためである。

【0048】

なお、冗長なデータの特定は、上記の方法に限るものではない。例えば、キーワード以外の特徴量として文書構造情報に基づきタイトル文字列やタイトル文字列のフォントサイズ、タイトル文字列の位置等を特徴量としたクラスタリングを利用してもよく、また、それ以外の方法を用いてもよい。そして、同じ種別であると判定された文書が大量にある場合は、それらの中から一部を抜粋して、残りの文書は削除することにより冗長なデータを削減することができる。このとき、例えば、あらかじめ文書のデータ容量や個数の上限を決めておき、当該データ容量や個数が上限を超えた場合に、それらが上限の値以下となるように文書を削除すればよい。

【0049】

ステップS504において、MFP101は、ステップS503で冗長なデータが削減された負例データ候補群から正例データ候補群に含まれる種別の可能性がある文書を削除する。正例データ候補群に含まれる種別であるか否かの判定は、正例データの文書とキーワードが一致する確率(一致率)に基づいて行う。なお、正例データ候補群に含まれる種別であるか否かの判定は、上記の方法に限るものではない。

【0050】

ここでの判定は、分類ルールを用いて「その他」に分類するか否かを判定する際の精度は必要としていない。文書構造情報の一致率や、画像特徴量の一致率を利用して、正例の種別であると疑わしい文書を削除できればよい。また、すでに分類ルールを一度構築しており、正例データの種別を追加する場合であれば、構築済みの分類ルールを適用して正例データの種別であるか否かを判定してもよい。

【0051】

なお、本実施形態では、サーバ102から取得した画像データ群をMFP101が絞り込む処理を実行することにより負例データを作成したが、これに限られるものではない。例えば、図6を用いて説明する絞り込み条件に従って、サーバ102がデータの絞り込み

10

20

30

40

50

を行い、作成した負例データをMFP101に送信してもよい。

【0052】

< 絞り込み条件の取得および絞り込み処理の詳細説明 (ステップS501、S502) >
絞り込み条件の取得および絞り込みの処理は、MFP101のCPU204が実行する処理である。絞り込み条件の取得について、図6および図7を用いて説明する。

図6は、文書ファイル情報の一例を示す図である。正例データ候補群として与えられた3種類の文書に関して、文書ファイル情報を示している。図7は、正例データ候補群の文書ファイル情報による絞り込み条件をユーザが確認および編集するための画面の一例を示す図である。

【0053】

図7の画面は、ラジオボタン701および702を有する。ラジオボタン701および702により、絞り込み条件を設定するか否かを切り替えることができる。ボタン703は、負例データ候補群の取得を指示(要求)するためのボタンであり、ラジオボタン701および702の状態に応じて取得する処理を切り替える。

【0054】

具体的には、ラジオボタン701が選択されている場合には、条件式フィールド705および条件フィールド706において設定された内容に基づいて、記憶部203およびHDD305内の文書を絞り込んで取得する。ラジオボタン702が選択されている場合には、記憶部203およびHDD305内の文書を絞り込まずにそのまま取得する。ボタン704は、絞り込み条件の自動取得を指示するためのボタンである。

【0055】

ボタン704によって絞り込み条件の自動取得が指示されると、MFP101は、正例データ候補群の文書ファイル情報から絞り込み条件を取得して、条件式フィールド705および条件フィールド706に表示する。具体的には、条件式フィールド705および条件フィールド706には、図6に示した正例データ候補群の文書ファイル情報に基づいて、正例データ候補群の文書ファイル情報と1つでも共通の項目を含む文書が取得できる条件式が示される。条件フィールド706には、正例データ候補群の文書ファイル情報の各項目が条件として表示される。

【0056】

また、条件式フィールド705には、条件フィールド706の各条件が、和集合を表す「+」で結合された条件式が入力されている。すなわち、正例データ候補群の文書ファイル情報の各項目と1つでも共通の項目を含む文書が抽出される。なお、条件式の自動取得では、上記のように文書ファイル情報の各項目の和集合を抽出する方法に限られるものではない。例えば、正例データ候補群の間で、文書ファイル情報の共起性を計算し、共起性の高い文書ファイル情報の項目の組み合わせを含む文書を絞り込むように条件を表示してもよい。すなわち、正例データ候補群において付与されている頻度が高い文書ファイル情報の項目の組み合わせを使用して、文書を絞り込むようにしてもよい。

【0057】

条件式フィールド705および条件フィールド706は、条件を表示するだけでなく、ユーザによる編集も受け付ける。ユーザは、ボタン704を用いて自動取得した条件を修正したい場合には編集すればよく、また、ユーザ所望の文書を絞り込むための条件を任意に設定することも可能である。ボタン709によって条件式フィールド705および条件フィールド706表示された条件をクリアすることも可能である。

【0058】

また、図7に示す例では、条件フィールド706に条件番号7までの条件が一覧されているが、これらの数は可変であり、上限も現在表示されている10個に限られるものではない。ボタン710によって、条件の追加が指示されると、条件の数(行数)を増やすことが可能である。また、条件式フィールド705において、条件フィールド706に表示されている条件番号と括弧や演算子を用いて多項演算のように条件式を入力することも可能である。例えば、和集合であれば「+」の演算子で表記し、積集合であれば「*」の演

10

20

30

40

50

算子で表記する。

【0059】

また、条件式フィールド705および条件フィールド706で表現される絞り込み条件は、ファイルに保存または読み込みが可能である。ボタン707は、絞り込み条件をファイルに保存するためのボタンであり、ボタン707が押下されると条件式フィールド705および条件フィールド706に表示されている絞り込み条件がテキストファイル形式にて保存される。

【0060】

また、ボタン708は、絞り込み条件をファイルから読み込むためのボタンであり、ボタン708が押下されるとファイルから読み込んだ絞り込み条件が、条件式フィールド705および条件フィールド706に表示される。なお、絞り込み条件を保存するファイルの形式は、テキストファイル形式に限られるものではなく、条件を表現することができれば特に限定されない。例えば、XML形式に保存してもよい。

10

【0061】

図8は、記憶部203およびHDD305内の文書ファイルを、上記の絞り込み条件によって絞り込んだ結果の一例を示す図である。

図7に示した条件によって絞り込んだ場合に、負例データとして採用されるデータの1つがデータ801である。

【0062】

文書ファイル情報の項目802（作成者名）、項目803（形式）、及び項目804（作成デバイス）が、それぞれ条件711、712、713と一致するため、採用される。一方、負例データとして採用されないデータの1つがデータ805である。文書ファイル情報の項目が、条件フィールド706に示す条件のいずれにも一致しないため、負例データとして採用されず、負例データ候補群から削除される。

20

【0063】

<機械学習を利用した分類ルールの構築の詳細説明（ステップS408）>

次に、本実施形態で分類ルールの構築に利用する機械学習の手法について説明する。本実施形態では、機械学習の手法としてReal AdaBoostと呼ばれる公知の手法を利用する。Real AdaBoostは、大量の特徴量から、与えられた学習セットの分類に適した特徴量を選択して、その特徴量を組み合わせて分類器（分類ルール）を構成することが可能な手法である。

30

【0064】

画像の分類時に大量の特徴量を利用すると、特徴量の計算負荷のためにパフォーマンスが低下する可能性がある。Real AdaBoostのように、分類に適した特徴量を選択して、一部の特徴量だけを利用し、分類器を構成できることは、大きな利点である。ただし、Real AdaBoostは、2クラス分類器であり、2種類のラベルがついたデータを分類するものである。つまり、このままでは、3種類以上の種別の画像データの分類には利用できない。

【0065】

そこで、本実施形態では、2クラス分類器を多クラス分類器に拡張するOVA（One - Versus - All）と呼ばれる公知の方法を利用する。OVAは、1つのクラス（対象クラス）とそれ以外のクラスを分類する分類器をクラスの数だけ作成し、それぞれの分類器の出力を、対象クラスの信頼度とする。すなわち、1つの分類器では、その分類器が分類するクラスに属するデータを正例データとし、それ以外のクラスに属するデータを負例データとして分類ルールを学習する。

40

【0066】

各分類器は、その分類器が対象とする1つのクラスのデータが入力された場合に、出力する信頼度が高くなるように学習を行う。分類の際には、分類したいデータをすべての分類器に入力し、信頼度が最大であったクラスを分類先とする。また、すべての分類器の出力する信頼度が小さい場合や、複数の分類器が出力する信頼度が高くなった場合には、「

50

該当なし」や「不明」といった判定を行う。

【0067】

図9は、学習セットを用いた機械学習の一例を示す図である。

この例では、学習セットとして、正例データ候補群の3つのクラス(種別)の文書(文書A、文書B、文書C)および負例データ候補群の「その他」の文書(文書A、文書B、文書Cではない文書)のそれぞれに対応する特徴量が用意されているものとする。この文書A、文書B、文書Cの3種類のクラスを分類するために、OVAでは3種類の分類器を用意する。3種類の分類器はそれぞれ、文書Aとそれ以外のクラスに文書を分類するための文書A分類器、文書Bとそれ以外のクラスに文書を分類するための文書B分類器、文書Cとそれ以外のクラスに文書を分類するための文書C分類器である。

10

【0068】

ここで、文書A分類器を構築する方法について説明する。まず、MFP101のCPU204は、分類ルールを学習するにあたって必要となる正例データおよび負例データを、学習セットの中から取得する。文書A用分類器では、正例データは文書Aのデータであり、負例データはそれ以外のクラスのデータである。したがって、CPU204は、正例データ候補群の中から、文書Aのラベルが付与された画像データを取得し、正例データとする。

【0069】

また、CPU204は、正例データ候補群の中から、文書A以外(文書B、文書C)のラベルの付与された画像データを、負例データとして取得する。さらに、CPU204は、負例データ候補群の中から、画像データを負例データとして取得する。このとき、負例データ候補群の中に、正例データである文書Aのデータが混ざっている場合には、正しく分類ルールを学習することができない。このため、上記のステップS504の処理により文書Aである可能性の高いものは取り除かれているものとする。

20

【0070】

CPU204は、取得した正例データおよび負例データの特徴量に基づき、Real AdaBoostを利用して文書A分類器を構築する。文書A分類器では、文書Aの特徴量が入力された場合に、大きい出力値(信頼度)が出力され、それ以外のクラスの文書の特徴量が入力された場合に、小さい出力値(信頼度)が出力される。文書B分類器、文書C分類器についても同様である。

30

【0071】

なお、本実施形態で利用可能な機械学習の手法は、上記の手法に限定されるものではない。Support Vector MachineやRandom Forest等の公知の手法を利用してもよい。また、特徴量選択の枠組みが機械学習の手法に含まれていない場合に、分類時の分類速度を向上させたい場合には、主成分分析や判別分析を利用した特徴量選択等の公知の特徴量選択を行ってもよい。機器学習の手法が2クラス分類器である場合は、OVA以外の、All-Versus-All(AVA)やError-Correcting Output-Coding(ECOC)等の公知の手法を用いてもよい。

【0072】

<分類ルールの構築に利用する特徴量の詳細>

本実施形態において分類ルールの構築に利用する特徴量について、図10および図11を用いて説明する。

40

【0073】

図10は、特徴量の算出方法について説明する図である。

本実施形態において特徴量は、入力画像1001内から切り出されたパッチ画像1002に対して勾配情報に基づき算出される9次元の特徴量である。MFP101のCPU204は、パッチ画像1002内の各画素について注目し、注目画素に隣接する画素の階調値から、勾配強度および勾配方向を算出する。

【0074】

50

そして、CPU 204は、勾配強度に基づいてエッジ判定を行うことで、勾配強度が一定値以上の画素をエッジ画素、一定値より小さい画素を非エッジ画素と判定する。エッジ判定の結果、画素1003は、非エッジ画素と判定された画素の一例であり、画素1004は、エッジ画素と判定された画素の一例である。エッジ画素である画素1004内の矢印は、勾配方向を表す。

【0075】

勾配方向は、文字や罫線の線の方向を表現するため、180度回転した角度は同一方向とみなして、0～180度に正規化される。CPU 204は、エッジ画素群から勾配方向を22.5度毎の8方向に量子化し、方向ごとの勾配強度積算値/パッチ画素数を計算して8ピンのヒストグラムを作成する。また、CPU 204は、非エッジ画素群から、非エッジ画素数/パッチ画素数を計算し、エッジ画素群から作成したヒストグラムと合わせて、1つのパッチ画像から9次元の特徴量を算出する。

10

【0076】

エッジ画素と非エッジ画素を利用することで、罫線や文字の情報だけでなく、文書画像の大きな特徴である余白部分を表現することが可能になる。これまでの説明は、1つのパッチ画像1002における特徴量の説明であるが、実際には、1つの入力画像から複数のパッチ画像を切り出して利用することにより、多数の特徴量を利用する。

【0077】

図11は、読み取った画像データからパッチ画像を切り出す方法について説明する図である。

20

CPU 204は、入力画像1101から余白をカットし、ノイズが表れやすい画像端1102を削除する。CPU 204は、余白カット後の画像1103を縮小することで、マルチスケール(複数の解像度の)画像を作成する。マルチスケールの画像を用意するのは、解像度ごとにエッジの構造が変わるためであり、画像読取部201の読取解像度や文書の解像度が多少異なっても対応できるようにするためである。

【0078】

画像1104は、余白カット後の画像1103を1/4に縮小した画像である。余白カット後の画像1103および縮小した画像1104から、パッチサイズと切り出し位置を変えながら、パッチ画像を切り出す。具体的には、まず、縮小した画像1104から、均等に16分割して得られる1/16サイズのパッチ画像16枚と、均等に64分割して得られる1/64サイズのパッチ画像64枚から、合計80枚のパッチ画像を作成する。

30

【0079】

また、余白カット後の画像1103から、同様に分割して80枚のパッチ画像を作成することで、1枚の入力画像1101から、合計160枚のパッチ画像が得られる。各パッチ画像から9次元の特徴量を算出するため、1枚の入力画像1101から $9 \times 160 = 1440$ 次元の特徴量を算出することが可能となる。

【0080】

なお、画像解像度、パッチサイズ、パッチ切り出し位置に関するパラメータは、上記の数字に限定されるものではない。また、算出する特徴量として、原稿の色の情報を利用するために、色ヒストグラムや色分散等を特徴量としてもよい。また、分類ルールの構築に利用する特徴量は、上記のような画像データに関する特徴量に限定されるものではない。例えば、負例データ候補群の絞り込みに利用するメタ情報や文書構造情報などの文書ファイル情報を利用してもよい。

40

【0081】

また、本実施形態では、文書をMFP 101により画像データとして読み込み、当該画像データを分類する場合について説明したが、これに限られるものではない。例えば、テキスト形式のデータに対しても、本発明の正例データを用いた負例データの絞り込みは適用可能である。

【0082】

以上のように、本実施形態によれば、高精度な分類ルールの構築を可能とする高品質な

50

負例データを効率的に生成することができる。

【0083】

(第2実施形態)

第1実施形態では、トレイ212にセットされ画像読取部201により一度に読み取られる原稿を正例データとして利用することを想定していた。これに対して、本実施形態では、トレイ212にセットされ画像読取部201により一度に読み取られる原稿に加え、すでにサーバ102上に格納された文書を正例データとして利用する場合を想定する。以下、第1実施形態との差分についてのみ説明する。

【0084】

<第2実施形態の詳細説明>

図12は、学習セットを生成し、分類ルールを学習する処理を説明するフローチャートである。

図12に示す処理は、MFP101およびサーバ102にて実行される。MFP101において実行される処理は、CPU204が記憶部203に格納されている処理プログラムをロードして実行することにより実現される。また、サーバ102において実行される処理は、CPU301がHDD305に格納されている処理プログラムをRAM302にロードして実行することにより実現される。

【0085】

なお、本実施形態では、ユーザがMFP101を用いて文書(原稿)をスキャンし、その種別毎に文書の画像データをサーバ102に保存するという業務フローの中で、同時にMFP101内で画像データの分類ルールを学習するシステムを想定している。さらに、本実施形態では、分類ルールの学習に利用する文書をサーバ102から取得することを想定している。

【0086】

ステップS1201において、MFP101は、正例データとして利用する文書を、原稿フィーダ211から読み込むか、サーバ102から選択するかを受け付ける。原稿フィーダ211から読み込む場合には、処理はステップS1202に進み、サーバ102から選択する場合には、処理はステップS1205に進む。ステップS1202およびステップS1203は、図4のステップS401およびステップS402と同様である。また、ステップS1204は、図4のステップS404と同様である。

【0087】

ステップS1205において、MFP101は、ユーザから操作部207経由でサーバ102のHDD305内のどの文書を利用するかを指示を受け付ける。サーバ102のCPU301は、ユーザの指示に基づきHDD305内の画像データ群を、インターネット104およびLAN103を経由してMFP101に送信する。ステップS1206において、MFP101は、ステップS1203にて画像読取部201で読み込まれた画像データ群、または、ステップS1205にてサーバ102から受信した画像データ群を、記憶部203に学習セットの正例データ候補群として格納する。

【0088】

格納する際には、各画像データに負例データ候補群の絞り込みにて利用する文書ファイル情報を付与する。画像読取部201で読み込まれた画像データ群に付与する文書ファイル情報は、図4のステップS403で付与する文書ファイル情報と同様である。一方、サーバ102から受信した画像データ群には、すでに文書ファイル情報が付与されている場合にはその文書ファイル情報を利用する。また、文書ファイル情報が不足している場合には不足している項目について、図4のステップS403で付与する文書ファイル情報と同様の文書ファイル類情報を付与する。

【0089】

ステップS1207において、MFP101は、正例データの登録を続けるか否かの指示を、操作部207を介してユーザから受け付ける。正例データの登録を続ける場合には、処理はステップS1201に戻る。正例データの登録を続けない場合には、処理はステ

10

20

30

40

50

ステップ S 1 2 0 8 に進む。なお、正例データの登録を続けるか否かの判断は、上記の方法に限られるものではない。例えば、ステップ S 1 2 0 6 における正例データの登録数をカウントし、あらかじめ操作部 2 0 7 を介してユーザによって設定された正例データの登録数に達するまで正例データの登録を続けてもよい。ステップ S 1 2 0 8 ~ ステップ S 1 2 1 0 は、図 4 のステップ S 4 0 6 ~ ステップ S 4 0 8 と同様である。

【 0 0 9 0 】

このように、本実施形態によれば、負例データを作成する際に、MFPから入力された画像データ（正例データ）から得られる文書ファイル情報に加えて、サーバから取得された画像データにすでに付与されている文書ファイル情報を利用することができる。これにより、大量の文書の中からデータの容量を抑えつつ、分類時に入力される可能性の高いデータを負例データとして収集することができ、高精度の分類ルールを効率よく生成することが可能となる。

10

【 0 0 9 1 】

（その他の実施形態）

本発明は、上述の実施形態の 1 以上の機能を実現するプログラムを、ネットワーク又は記憶媒体を介してシステム又は装置に供給し、そのシステム又は装置のコンピュータにおける 1 つ以上のプロセッサがプログラムを読み出し実行する処理でも実現可能である。また、1 以上の機能を実現する回路（例えば、ASIC）によっても実現可能である。

【 0 0 9 2 】

以上、本発明の好ましい実施形態について説明したが、本発明は、これらの実施形態に限定されず、その要旨の範囲内で種々の変形および変更が可能である。

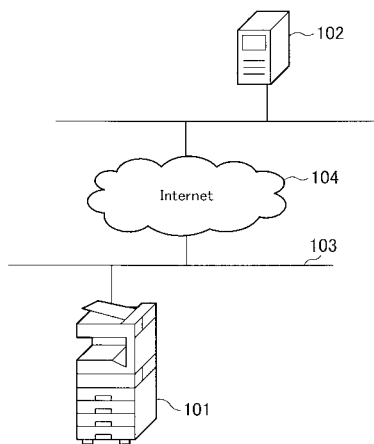
20

【 符号の説明 】

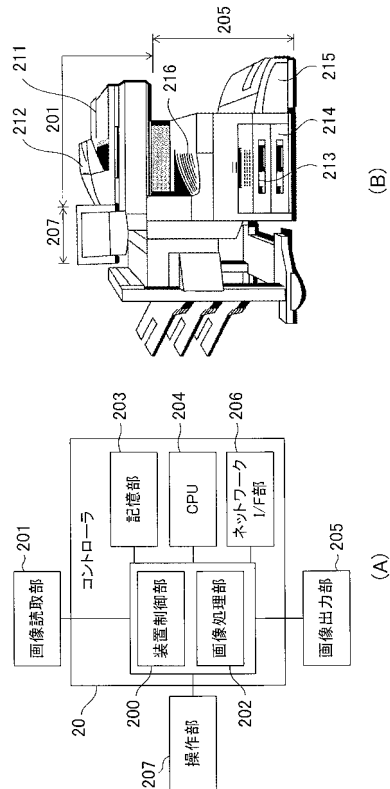
【 0 0 9 3 】

- 1 0 1 M F P
- 1 0 2 サーバ

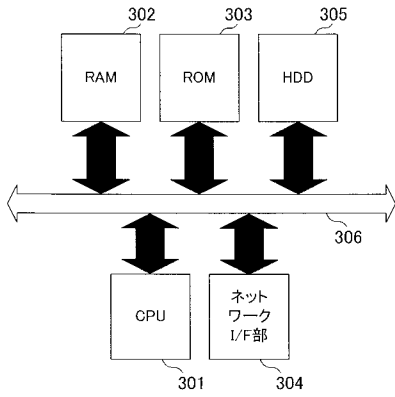
【 図 1 】



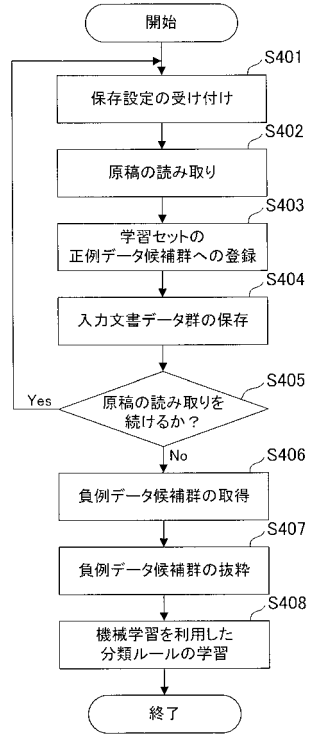
【 図 2 】



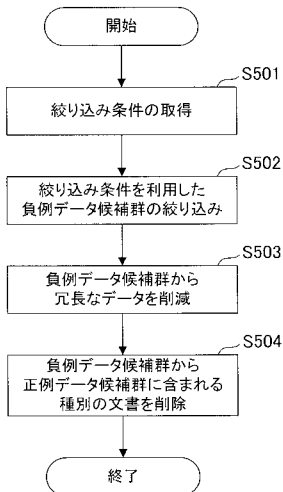
【 図 3 】



【 図 4 】



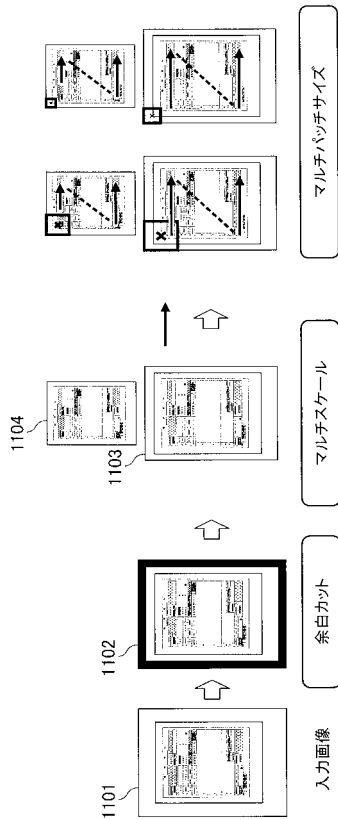
【 図 5 】



【 図 6 】

クラス	ファイル名	タイトル	作成者名	形式	キーワード	作成デバイス	生成日時
文書A	文書A_001.pdf	発注書	ユーザA	PDF	企業A, 商品A, 商品コード	デバイスA	2015/10/25
文書A	文書A_002.pdf	発注書	ユーザA	PDF	企業A, 商品A, 商品コード	デバイスA	2015/10/25
文書A	文書A_003.pdf	発注書	ユーザA	PDF	企業A, 商品A, 商品コード	デバイスA	2015/10/25
文書B	文書B_001.pdf	注文書	ユーザA	PDF	企業B, 商品B, 商品コード	デバイスA	2015/10/25
文書B	文書B_002.pdf	注文書	ユーザA	PDF	企業B, 商品B, 商品コード	デバイスA	2015/10/25
文書C	文書C_001.pdf	発注書	ユーザA	PDF	企業C, 商品番号, 商品名	デバイスA	2015/10/25
文書C	文書C_002.pdf	発注書	ユーザA	PDF	企業C, 商品番号, 商品名	デバイスA	2015/10/25
文書C	文書C_003.pdf	発注書	ユーザA	PDF	企業C, 商品番号, 商品名	デバイスA	2015/10/25
文書C	文書C_004.pdf	発注書	ユーザA	PDF	企業C, 商品番号, 商品名	デバイスA	2015/10/25

【 図 1 1 】



【 図 1 2 】

