

【特許請求の範囲】**【請求項 1】**

画像データが示す画像において対象が含まれる第 1 の領域と、該対象が分類されるカテゴリとを認識する画像処理装置であって、

畳み込みニューラルネットワークを用いて、入力された前記画像データが分類されるカテゴリを認識する認識手段と、

前記認識手段における前記畳み込みニューラルネットワークの所定の層の出力結果を示す第 1 の出力データに基づいて、前記画像データが示す画像に含まれる 1 以上の候補領域を示す 1 以上の候補領域画像データを作成する候補領域作成手段と

を有し、

10

前記認識手段は、

前記候補領域作成手段により作成された前記 1 以上の候補領域画像データがそれぞれ分類されるカテゴリを認識する、画像処理装置。

【請求項 2】

前記第 1 の出力データは、前記畳み込みニューラルネットワークの前記所定の層のネットワークパラメータから特定されるフィルタ毎の第 2 の出力データを含み、

前記第 1 の出力データから所定の個数の前記第 2 の出力データを決定する決定手段を有し、

前記候補領域作成手段は、

前記決定手段で決定された前記第 2 の出力データに基づいて、前記 1 以上の候補領域データを作成する、請求項 1 記載の画像処理装置。

20

【請求項 3】

前記決定手段は、

前記第 2 の出力データの代表データ値の昇順に、前記所定の個数の前記第 2 の出力データを決定する、請求項 2 記載の画像処理装置。」

【請求項 4】

前記第 2 の出力データが示す画像を 1 以上の第 2 の領域に分割する分割手段を有し、

前記候補領域作成手段は、

前記分割手段により分割された前記 1 以上の第 2 の領域のそれぞれについて、該第 2 の領域を囲む最小の矩形領域を前記候補領域とする、請求項 2 又は 3 に記載の画像処理装置

30

【請求項 5】

前記分割手段は、

微分処理により前記 1 以上の第 2 の領域の境界を検出し、該検出された境界に基づいて分割する、請求項 4 記載の画像処理装置。

【請求項 6】

前記分割手段は、

前記微分処理に Sobel フィルタを用いる、請求項 5 記載の画像処理装置。

【請求項 7】

所定の閾値以下のデータ値を削除する閾値手段を有し、

40

前記分割手段は、

前記閾値手段により所定の閾値以下のデータ値を削除した前記第 2 の出力データが示す画像を 1 以上の領域に分割する、請求項 4 ないし 6 のいずれか 1 項に記載の画像処理装置

【請求項 8】

画像データが示す画像において対象が含まれる第 1 の領域と、該対象が分類されるカテゴリとを認識する画像処理装置による画像処理方法であって、

畳み込みニューラルネットワークを用いて、入力された前記画像データが分類されるカテゴリを認識する認識手順と、

前記認識手順における前記畳み込みニューラルネットワークの所定の層の出力結果を示

50

す第1の出力データに基づいて、前記画像データが示す画像に含まれる1以上の候補領域を示す1以上の候補領域画像データを作成する候補領域作成手順と

を有し、

前記認識手順は、

前記候補領域作成手順により作成された前記1以上の候補領域画像データがそれぞれ分類されるカテゴリを認識する、画像処理方法。

【請求項9】

画像データが示す画像において対象が含まれる第1の領域と、該対象が分類されるカテゴリとを認識する画像処理装置を、

畳み込みニューラルネットワークを用いて、入力された前記画像データが分類されるカテゴリを認識する認識手段、

前記認識手段における前記畳み込みニューラルネットワークの所定の層の出力結果を示す第1の出力データに基づいて、前記画像データが示す画像に含まれる1以上の候補領域を示す1以上の候補領域画像データを作成する候補領域作成手段

として機能させ、

前記認識手段は、

前記候補領域作成手段により作成された前記1以上の候補領域画像データがそれぞれ分類されるカテゴリを認識する、プログラム。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、画像処理装置、画像処理方法、及びプログラムに関する。

【背景技術】

【0002】

デジタルカメラや携帯情報端末等の機器において、撮影された画像中の被写体が属するカテゴリ（例えば、「人」、「動物」、「車」等）を分類する技術が知られている。

【0003】

また、画像中において、被写体が占める領域と、当該被写体が分類されるカテゴリとを認識する技術が知られている（例えば特許文献1及び非特許文献1参照）。このような技術では、被写体が占める領域の候補である候補領域に対して、カテゴリを分類するための処理を行うことで、被写体が占める領域と、当該被写体が分類されるカテゴリとを認識する。

【発明の概要】

【発明が解決しようとする課題】

【0004】

しかしながら上記の従来技術では、被写体が占める領域と、当該被写体が分類されるカテゴリとの認識処理に多くの時間を要する場合があった。例えば、候補領域の数が多い場合には、それぞれの候補領域に対してカテゴリを分類するための処理を行うため、認識処理に多くの時間を要することがある。

【0005】

本発明の実施形態は、認識処理の処理時間の削減を支援することを目的とする。

【課題を解決するための手段】

【0006】

上記目的を達成するため、本発明の実施の形態では、画像データが示す画像において対象が含まれる第1の領域と、該対象が分類されるカテゴリとを認識する画像処理装置であって、畳み込みニューラルネットワークを用いて、入力された前記画像データが分類されるカテゴリを認識する認識手段と、前記認識手段における前記畳み込みニューラルネットワークの所定の層の出力結果を示す第1の出力データに基づいて、前記画像データが示す画像に含まれる1以上の候補領域を示す1以上の候補領域画像データを作成する候補領域作成手段とを有し、前記認識手段は、前記候補領域作成手段により作成された前記1以上

10

20

30

40

50

の候補領域画像データがそれぞれ分類されるカテゴリを認識する。

【発明の効果】

【0007】

本発明の実施形態によれば、認識処理の処理時間の削減を支援することができる。

【図面の簡単な説明】

【0008】

【図1】本実施形態の画像処理装置のハードウェア構成の一例を示す図である。

【図2】本実施形態の画像処理装置の機能構成の一例を示す図である。

【図3】本実施形態の画像処理装置の認識処理のフローチャートの一例を示す図である。

【図4】本実施形態の畳み込みニューラルネットワーク処理のフローチャートの一例を示す図である。 10

【図5】本実施形態の入力画像データの加工処理の一例を示す図である。

【図6】本実施形態の第1層の畳み込み処理の一例を示す図である。

【図7】本実施形態の第1層のネットワークパラメータの一例を示す図である。

【図8】本実施形態の第1層のフィルタの一例を示す図である。

【図9】本実施形態の第1層のプーリング処理の一例を示す図である。

【図10】本実施形態の第2層の畳み込み処理の一例を示す図である。

【図11】本実施形態の第2層のネットワークパラメータの一例を示す図である。

【図12】本実施形態の第2層のフィルタの一例を示す図である。

【図13】本実施形態の候補領域の作成処理のフローチャートの一例を示す図である。 20

【図14】本実施形態の微分処理の一例を示す図である。

【図15】本実施形態の閾値処理の一例を示す図である。

【図16】本実施形態の領域分割の一例を示す図である。

【図17】本実施形態の最小矩形の一例を示す図である。

【図18】本実施形態のカテゴリ分類処理のフローチャートの一例を示す図である。

【図19】本実施形態の第3層の全結合処理の一例を示す図である。

【図20】本実施形態の第3層のネットワークパラメータの一例を示す図である。

【図21】本実施形態の正規化処理の一例を示す図である。

【発明を実施するための形態】

【0009】

30

本実施形態は、画像データが示す画像において、当該画像の被写体を示す対象（例えば、人や物体等）を含む領域と、当該対象が分類されるカテゴリとを認識するものである。ここで、カテゴリとは、例えば、「人」、「動物」、「車」、「花」、「料理」等の対象が分類される種別のことである。

【0010】

以降では、画像データに対して、上述した認識を行う処理（認識処理）を実行する画像処理装置10について説明する。なお、本実施形態の画像処理装置10は、例えば、デジタルカメラ、スマートフォン、タブレット端末、ゲーム機器、ノート型PC、デスクトップ型PC等である。

【0011】

40

<ハードウェア構成>

まず、本実施形態の画像処理装置10のハードウェア構成について、図1を参照しながら説明する。図1は、本実施形態の画像処理装置のハードウェア構成の一例を示す図である。

【0012】

本実施形態の画像処理装置10は、入力装置11、表示装置12、CPU（Central Processing Unit）13、及びROM（Read Only Memory）14を有する。また、本実施形態の画像処理装置10は、RAM（Random Access Memory）15、インタフェース装置16、記憶装置17、及び撮像装置18を有する。これら各ハードウェアは、バスBにより相互に接続されている。

50

【 0 0 1 3 】

入力装置 1 1 は、キーボードやマウス、タッチパネル、各種ボタン等を含み、画像処理装置 1 0 に各種信号を入力するのに用いられる。表示装置 1 2 は、ディスプレイ等を含み、各種の処理結果を表示する。特に、表示装置 1 2 には、本実施形態の認識処理の処理結果が表示される。すなわち、表示装置 1 2 には、入力された画像データが示す画像において、被写体等の対象が含まれる領域と、当該対象が分類されるカテゴリと示す処理結果が表示される。

【 0 0 1 4 】

CPU 1 3 は、例えば記憶装置 1 7 や ROM 1 4 等からプログラムやデータを RAM 1 5 上に読み出して、各種処理を実行する演算装置である。ROM 1 4 は、電源を切ってもデータを保持することができる不揮発性の半導体メモリである。RAM 1 5 は、プログラムやデータを一時保存することができる揮発性の半導体メモリである。

10

【 0 0 1 5 】

インタフェース装置 1 6 は、外部装置とのインタフェースである。外部装置には、例えば、CD (Compact Disk) や DVD (Digital Versatile Disk)、SD メモリカード (SD memory card)、USB メモリ (Universal Serial Bus memory) 等の記録媒体がある。画像処理装置 1 0 は、インタフェース装置 1 6 を介して、本実施形態の認識処理の処理対象となる画像データを記録媒体から読み取ることができる。

【 0 0 1 6 】

記憶装置 1 7 は、プログラムやデータを格納している HDD (Hard Disk Drive) や SSD (Solid State Drive) 等の不揮発性のメモリである。記憶装置 1 7 に格納されるプログラムやデータには、本実施形態の認識処理を実行する画像処理プログラム 2 0 がある。また、本実施形態の認識処理の処理対象となる画像データが格納されても良い。

20

【 0 0 1 7 】

撮像装置 1 8 は、カメラ等であり、本実施形態の認識処理の処理対象となる画像データを作成する。

【 0 0 1 8 】

本実施形態の画像処理装置 1 0 は、上記ハードウェア構成により後述する各種処理を実現することができる。

【 0 0 1 9 】

< 機能構成 >

次に、本実施形態の画像処理装置 1 0 の機能構成について、図 2 を参照しながら説明する。図 2 は、本実施形態の画像処理装置の機能構成の一例を示す図である。

30

【 0 0 2 0 】

本実施形態の画像処理装置 1 0 は、CNN 処理部 1 1 0、候補領域作成処理部 1 2 0、正規化処理部 1 3 0、及び出力部 1 4 0 を有する。これら各部は、画像処理装置 1 0 にインストールされた画像処理プログラム 2 0 が、CPU 1 3 に実行させる処理により実現される。

【 0 0 2 1 】

CNN 処理部 1 1 0 は、ネットワークパラメータ 1 0 0 0 に基づいて、畳み込みニューラルネットワーク (CNN: Convolutional Neural Network) 処理を行う。畳み込みニューラルネットワークは、一般に、 n を 3 以上の任意の自然数として、畳み込み処理及びプーリング処理を行う第 1 層 ~ 第 $n - 2$ 層と、畳み込み処理を行う第 $n - 1$ 層と、全結合処理を行う第 n 層とを含む。

40

【 0 0 2 2 】

ここで、ネットワークパラメータ 1 0 0 0 は、教師あり学習の手法により、学習データに基づいて畳み込みニューラルネットワークの各層毎に予め学習されたデータである。教師あり学習の手法には、例えば誤差逆伝播法 (Backpropagation) を用いれば良い。

【 0 0 2 3 】

このようなネットワークパラメータ 1 0 0 0 は、例えば記憶装置 1 7 等に格納され、パ

50

イアスデータ 1100 及び重みデータ 1200 が含まれる。なお、以降では、第 n 層のネットワークパラメータ 1000 を「ネットワークパラメータ 1000 - n」と表す。したがって、第 n 層のバイアスデータ 1100 及び重みデータ 1200 はそれぞれ「バイアスデータ 1100 - n」及び「重みデータ 1200 - n」と表される。ネットワークパラメータ 1000 の詳細については後述する。

【0024】

CNN 処理部 110 は、入力画像を示す画像データ 510 に対して、畳み込みニューラルネットワーク処理を行い、予め設定された第 N 層における畳み込み処理の処理結果を示す出力データ 520 を出力する。

【0025】

ここで、本実施形態では、N = 2 であるものとして説明する。N = 2 の場合、出力データ 520 は、例えば、28 × 28 × 64 チャンネルの画像データとして表すことができる。換言すれば、出力データ 520 は、64 個の 28 × 28 チャンネルの画像データの集合として表すことができる。なお、N の値は、画像処理プログラム 20 の設計者等により予め設定される。N の値は、例えば、2 ~ 20 程度が好ましい。

【0026】

また、CNN 処理部 110 は、後述する候補領域作成処理部 120 により作成された候補領域画像データ 530 に対して、畳み込みニューラルネットワーク処理を行い、出力結果を正規化処理部 130 に出力する。

【0027】

さらに、CNN 処理部 110 は、加工部 111、畳み込み処理部 112、プーリング処理部 113、及び全結合処理部 114 を有する。加工部 111 は、CNN 処理部 110 に入力された画像データの加工処理を行う。畳み込み処理部 112 は、畳み込みニューラルネットワークの各層において畳み込み処理を行う。プーリング処理部 113 は、畳み込みニューラルネットワークの各層においてプーリング処理を行う。全結合処理部 114 は、全結合処理を行う。

【0028】

ここで、CNN 処理部 110 は、全結合処理部 114 をカテゴリの組毎に有しているものとする。カテゴリの組とは、カテゴリと、当該カテゴリ以外を示すカテゴリとのペアである。具体的には、カテゴリの組は、「人」「人以外」、「車」「車以外」、「動物」「動物以外」等の、あるカテゴリと、当該カテゴリ以外を示すカテゴリとのペアである。なお、以降では、複数の全結合処理部 114 を区別して表す場合は、「全結合処理部 114 - 1」、「全結合処理部 114 - 2」等と表す。

【0029】

候補領域作成処理部 120 は、出力データ 520 に基づいて、1 以上の候補領域画像データ 530 を作成する。候補領域画像データ 120 とは、画像データ 510 が示す画像において、対象が含まれる領域の候補を示すデータである。なお、以降では、複数の候補領域画像データ 530 を区別して表す場合は、「候補領域画像データ 530 - 1」、「候補領域画像データ 530 - 2」等と表す。

【0030】

ここで、候補領域作成処理部 120 は、データ決定部 121、境界決定部 122、閾値処理部 123、領域分割部 124、及び候補領域作成部 125 を有する。

【0031】

データ決定部 121 は、例えば 64 個の 28 × 28 チャンネルのデータとして表される出力データ 520 から所定の M 個の 28 × 28 チャンネルのデータを決定する。ここで、M の値は、画像処理プログラム 20 の設計者等により予め設定される。M の値は、例えば、3 ~ 20 程度が好ましい。

【0032】

境界決定部 122 は、データ決定部 121 により決定されたそれぞれのデータに対して、所定の微分処理を行い、領域分割部 124 が分割する領域の境界を決定する。

10

20

30

40

50

【 0 0 3 3 】

閾値処理部 1 2 3 は、閾値処理を行う。閾値処理とは、予め設定された閾値以下のデータを削除（すなわち、「0」とする）する処理である。なお、このような閾値は、画像処理プログラム 2 0 の設計者等により予め設定される。閾値の値は、例えば、1 0 ~ 5 0 程度が好ましい。

【 0 0 3 4 】

領域分割部 1 2 4 は、境界決定部 1 2 2 により決定された境界に基づいて、データ決定部 1 2 1 により決定されたデータが示す画像を、複数の領域に分割する。

【 0 0 3 5 】

候補領域作成部 1 2 5 は、領域分割部 1 2 4 により分割された複数の領域に基づいて、候補領域を作成し、作成した候補領域を示す候補領域画像データ 5 3 0 を出力する。

10

【 0 0 3 6 】

例えば、候補領域作成部 1 2 5 は、領域分割部 1 2 4 により分割された複数の領域のうちの一の領域に基づいて、候補領域画像データ 5 3 0 - 1 を出力する。同様に、候補領域作成部 1 2 5 は、領域分割部 1 2 4 により分割された複数の領域のうち他の領域に基づいて、候補領域画像データ 5 3 0 - 2 を出力する。

【 0 0 3 7 】

このように、本実施形態の候補領域作成処理部 1 2 0 は、出力データ 5 2 0 に基づいて候補領域画像データ 5 3 0 を作成する。これにより、本実施形態では、認識処理の精度の低下を防ぎつつ、候補領域を削減させることができる。したがって、本実施形態では、認識処理の処理時間を削減させることができる。

20

【 0 0 3 8 】

正規化処理部 1 3 0 は、CNN 処理部 1 1 0 による処理結果を正規化する。CNN 処理部 1 1 0 の各全結合処理部 1 1 4 による処理結果を比較することができる。以降では、正規化処理部 1 3 0 により正規化された、全結合処理部 1 1 4 の処理結果を「確信度」と表す。

【 0 0 3 9 】

例えば、カテゴリの組「人」「人以外」に対応する全結合処理部 1 1 4 の確信度は、CNN 処理部 1 1 0 に入力された画像データが示す画像が、カテゴリ「人」に分類される度合いを示す第 1 の値と、カテゴリ「人以外」に分類される度合いを示す第 2 の値との組で表される。

30

【 0 0 4 0 】

同様に、カテゴリの組「車」「車以外」に対応する全結合処理部 1 1 4 の確信度は、CNN 処理部 1 1 0 に入力された画像データが示す画像が、カテゴリ「車」に分類される度合いを示す第 1 の値と、カテゴリ「人以外」に分類される度合いを示す第 2 の値との組で表される。

【 0 0 4 1 】

出力部 1 4 0 は、認識結果 5 4 0 を出力する。ここで、認識結果 5 4 0 には、候補領域画像データ 5 3 0 から選択された結果画像データ 5 4 1 と、当該結果画像データ 5 4 1 のカテゴリを示すカテゴリ情報 5 4 2 とが含まれる。なお、出力部 1 4 0 は、候補領域画像データ 5 3 0 の確信度に基づいて、当該候補領域画像データ 5 3 0 から結果画像データ 5 4 1 を選択するとともに、当該結果画像データ 5 4 1 のカテゴリを決定してカテゴリ情報 5 4 2 を作成する。

40

【 0 0 4 2 】

これにより、画像データ 5 1 0 が示す画像において、対象が含まれる領域の画像と、当該対象が分類されるカテゴリとが出力される。

【 0 0 4 3 】

< 処理の詳細 >

次に、本実施形態の画像処理装置 1 0 の認識処理の詳細について、図 3 を参照しながら説明する。図 3 は、本実施形態の画像処理装置の認識処理のフローチャートの一例を示す

50

図である。

【 0 0 4 4 】

画像処理装置 1 0 は、画像データ 5 1 0 を入力する（ステップ S 3 1）。画像処理装置 1 0 は、例えば、記憶装置 1 7 に格納されている画像データ 5 1 0 を入力しても良いし、撮像装置 1 8 により生成された画像データ 5 1 0 を入力しても良い。また、画像処理装置 1 0 は、例えば、ネットワーク経由でダウンロードした画像データ 5 1 0 を入力しても良い。

【 0 0 4 5 】

画像処理装置 1 0 は、CNN 処理部 1 1 0 により、入力された画像データ 5 1 0 に対して、予め設定された第 N 層の畳み込み処理までの畳み込みニューラルネットワーク処理を行う（ステップ S 3 2）。この畳み込みニューラルネットワーク処理についての詳細については、後述する。ここでは、本ステップの畳み込みニューラルネットワーク処理において、第 N 層の畳み込み処理の処理結果を示す出力データ 5 2 0 が得られたものとして説明を続ける。

10

【 0 0 4 6 】

なお、上述したように、 $N = 3$ である場合、出力データ 5 2 0 は、例えば 6 4 個の 28×28 チャンネルのデータとして表される。

【 0 0 4 7 】

画像処理装置 1 0 は、候補領域作成処理部 1 2 0 により、出力データ 5 2 0 を入力して候補領域の作成処理を行う（ステップ S 3 3）。この候補領域の作成処理において、候補領域作成処理部 1 2 0 は、出力データ 5 2 0 に基づいて、1 以上の候補領域画像データ 5 3 0 を作成する。この候補領域の作成処理の詳細については、後述する。ここでは、本ステップの候補領域の作成処理において、1 以上の候補領域画像データ 5 3 0 が得られたものとして説明を続ける。

20

【 0 0 4 8 】

画像処理装置 1 0 は、CNN 処理部 1 1 0 及び正規化処理部 1 3 0 により、一の候補領域画像データ 5 3 0 を入力し、当該一の候補領域画像データ 5 3 0 のカテゴリを分類するカテゴリ分類処理を行う（ステップ S 3 4）。このカテゴリ分類処理により、入力された一の候補領域画像データ 5 3 0 の確信度が得られる。このカテゴリ分類処理の詳細については、後述する。ここでは、本ステップのカテゴリ分類処理において、一の候補領域画像データ 5 3 0 の確信度が得られたものとして説明を続ける。

30

【 0 0 4 9 】

画像処理装置 1 0 は、CNN 処理部 1 1 0 及び正規化処理部 1 3 0 により、すべての候補領域画像データ 5 3 0 の確信度が得られたか否かを判定する（ステップ S 3 5）。確信度が得られていない（すなわち、カテゴリ分類処理を行っていない）候補領域画像データ 5 3 0 が存在する場合には、ステップ S 3 4 に戻る。すなわち、画像処理装置 1 0 は、候補領域画像データ 5 3 0 - 1、候補領域画像データ 5 3 0 - 2、・・・等に対して、それぞれの確信度を順に取得する。

【 0 0 5 0 】

一方、すべての候補領域画像データ 5 3 0 の確信度が得られた場合には、ステップ S 3 6 に進む。

40

【 0 0 5 1 】

画像処理装置 1 0 は、出力部 1 4 0 により、得られた確信度に基づいて候補領域画像データ 5 3 0 から結果画像データ 5 4 1 を選択するとともに、当該結果画像データ 5 4 1 のカテゴリを決定してカテゴリ情報 5 4 2 を作成する。（ステップ S 3 6）。すなわち、出力部 1 4 0 は、認識結果 5 4 0 を決定する。

【 0 0 5 2 】

出力部 1 4 0 は、すべての候補領域画像データ 5 3 0 を結果画像データ 5 4 1 と選択しても良いし、候補領域画像データ 5 3 0 のうちの一部を結果画像データ 5 4 1 と選択しても良い。

50

【 0 0 5 3 】

また、出力部 1 4 0 は、例えば、候補領域画像データ 5 3 0 が示す画像のうち、一部が重畳している画像が存在する場合に、当該重畳している画像が示す候補領域画像データ 5 3 0 のうち、最も確信度が高い候補領域画像データ 5 3 0 を結果画像データ 5 4 1 と選択しても良い。より具体的には、例えば、候補領域画像データ 5 3 0 - 1 が示す第 1 の画像と、候補領域画像データ 5 3 0 - 2 が示す第 2 の画像と、候補領域画像データ 5 3 0 - 3 が示す第 3 の画像とが、少なくとも一部の領域において重畳しているものとする。この場合、第 1 の画像の確信度の第 1 の値と、第 2 の画像の確信度の第 1 の値と、第 3 の画像の確信度の第 1 の値とを比較し、最も値が高い画像を示す候補領域画像データ 5 3 0 を結果画像データ 5 4 1 と選択すれば良い。

10

【 0 0 5 4 】

なお、ステップ S 3 6 において、出力部 1 4 0 は、2 以上の認識結果 5 4 0 を決定しても良い。すなわち、出力部 1 4 0 は、候補領域画像データ 5 3 0 から 2 以上の結果画像データ 5 4 1 を選択するとともに、当該 2 以上の結果画像データ 5 4 1 のそれぞれのカテゴリ情報 5 4 2 を作成しても良い。これにより、例えば、画像データ 5 1 0 が示す画像において、複数の対象（例えば、「人」と「車」等）が写っている場合にも、それぞれの対象が含まれる領域の画像と、それぞれの対象が分類されるカテゴリとを決定することができる。

【 0 0 5 5 】

画像処理装置 1 0 は、出力部 1 4 0 により、決定された認識結果 5 4 0 を出力する（ステップ S 3 7）。このとき、出力部 1 4 0 は、例えば表示装置 1 2 に認識結果 5 4 0 を出力すれば良い。これにより、画像データ 5 1 0 が示す画像において、対象が含まれる領域の画像と、当該対象が分類されるカテゴリとが表示装置 1 2 に表示される。

20

【 0 0 5 6 】

次に、図 3 のステップ S 3 2 の畳み込みニューラルネットワーク処理について、図 4 を参照しながら説明する。図 4 は、本実施形態の畳み込みニューラルネットワーク処理のフローチャートの一例を示す図である。

【 0 0 5 7 】

加工部 1 1 1 は、入力された画像データ 5 1 0 の加工処理を行う（ステップ S 4 1）。この加工処理は、入力された画像データ 5 1 0 を、畳み込み処理部 1 1 2 が処理可能な形式とするための処理である。

30

【 0 0 5 8 】

ここで、加工処理について、図 5 を参照しながら説明する。図 5 は、本実施形態の入力画像データの加工処理の一例を示す図である。なお、入力された画像データ 5 1 0 の色空間が RGB 色空間である（すなわち、画像データ 5 1 0 の色チャンネルが 3 チャンネルである）ものとして説明する。ただし、画像データ 5 1 0 の色空間は、RGB 色空間に限られず、例えば、CMK 色空間、HSV 色空間、HLS 色空間等であっても良い。

【 0 0 5 9 】

Step 4 1 1) 加工部 1 1 1 は、入力された画像データ 5 1 0 を 64×64 (ピクセル) となるように縮小する、このとき、加工部 1 1 1 は、画像データ 5 1 0 の長辺が 64 (ピクセル) となるように縮小を行う。また、加工部 1 1 1 は、短辺が縮小された結果 64 (ピクセル) に満たない部分については値 0 (すなわち、RGB の各色成分の値が 0) でパディングして 64 (ピクセル) とする。なお、画像データ 5 1 0 を縮小するためのアルゴリズムには、例えば、バイリニア法を用いれば良い。

40

【 0 0 6 0 】

Step 4 1 2) 加工部 1 1 1 は、Step S 4 1 1 で得られた 64×64 の画像データの各画素値から、所定の値を減算した画像データを生成する。

【 0 0 6 1 】

ここで、所定の値は、各学習データに含まれる画像データ（以降、「学習画像データ」という）の各画素値の平均値である。すなわち、学習画像データの画素位置 (i, j) に

50

おける各学習画像データの画素値の平均値を $M(i, j)$ とした場合、上記の Step 4 1 1 において得られた 64×64 の画像データの各画素位置 (i, j) の画素値から $M(i, j)$ を減算する。ここで、 $i, j = 1, \dots, 64$ である。

【0062】

Step 4 1 3) 加工部 1 1 1 は、Step 4 1 2 で得られた画像データの中心の 56×56 (ピクセル) の画像データ以外を 0 クリアする。換言すれば、Step 4 1 2 において得られた画像データの周辺 4 ピクセル分を 0 クリアする。なお、図 5 において、網掛け部分が 0 クリアした部分である。

【0063】

そして、加工部 1 1 1 は、図 5 の Step 4 1 3 で得られた 64×64 (ピクセル) の画像データ (この画像データを「画像データ 5 1 1」とする。) を畳み込み処理部 1 1 2 に出力する。

【0064】

CNN 処理部 1 1 0 は、畳み込みニューラルネットワークの層を示す変数 n を 1 とする (ステップ S 4 2)。

【0065】

畳み込み処理部 1 1 2 は、画像データ 5 1 1 を入力して、第 1 層の畳み込み処理を行う (ステップ S 4 3)。

【0066】

ここで、第 1 層の畳み込み処理について、図 6 を参照しながら説明する。図 6 は、本実施形態の第 1 層の畳み込み処理の一例を示す図である。

【0067】

Step 4 3 1) 畳み込み処理部 1 1 2 は、画像データ 5 1 1 を入力する。ここで、入力した画像データ 5 1 1 の色空間は RGB 色空間であるため、色チャンネルは $64 \times 64 \times 3$ チャンネルである。

【0068】

Step 4 3 2) 畳み込み処理部 1 1 2 は、重みデータ 1 2 0 0 - 1 からフィルタを生成し、画像データ 5 1 1 の中心の 56×56 (ピクセル) の部分に対して、生成したフィルタを用いてフィルタ処理を行う。ここで、重みデータ 1 2 0 0 - 1 のデータ構成及び当該重みデータ 1 2 0 0 - 1 から生成されるフィルタ $1300f_j - 1$ ($j = 1, \dots, 64$) のデータ構成について説明する。

【0069】

図 7 (b) は、第 1 層の重みデータ 1 2 0 0 - 1 の一例を示す図である。図 7 (b) に示すように、第 1 層の重みデータ 1 2 0 0 - 1 は、 75×64 の行列で表される。なお、重みデータ 1 2 0 0 - 1 の各値 $w_1(i, j)$ は、上述したように、学習データに基づいて予め学習された値である。

【0070】

次に、重みデータ 1 2 0 0 - 1 から生成されるフィルタ $1300f_j - 1$ ($j = 1, \dots, 64$) について説明する。図 8 は、本実施形態の第 1 層のフィルタの一例を示す図である。

【0071】

図 8 に示すように、各フィルタ $1300f_j - 1$ ($j = 1, \dots, 64$) は、 5×5 の行列の 3 つの組で表される。換言すれば、各フィルタ $1300f_j - 1$ ($j = 1, \dots, 64$) は、 $5 \times 5 \times 3$ で表される。

【0072】

ここで、重みデータ 1 2 0 0 - 1 の $w_1(1, 1) \sim w_1(25, 1)$ 、 $w_1(26, 1) \sim w_1(50, 1)$ 、及び $w_1(51, 1) \sim w_1(75, 1)$ からフィルタ $1300f_1 - 1$ が生成される。同様に、重みデータ 1 2 0 0 - 1 の $w_1(1, 2) \sim w_1(25, 2)$ 、 $w_1(26, 2) \sim w_1(50, 2)$ 、及び $w_1(51, 2) \sim w_1(75, 2)$ からフィルタ $1300f_2 - 1$ が生成される。 $j = 3, \dots, 64$ の場合も同様で

10

20

30

40

50

ある。

【0073】

以上のように生成された各フィルタ $1300f_j - 1$ ($j = 1, \dots, 64$) を用いて、畳み込み処理部 112 は、画像データ 511 に対してフィルタ処理を行う。畳み込み処理部 112 は、例えば以下のようにしてフィルタ処理を行う。

【0074】

(1) 画像データ 511 の中心 $56 \times 56 \times 3$ の部分に対してフィルタ $1300f_1 - 1$ をかける (すなわち、画像データ 511 とフィルタ $1300f_1 - 1$ の対応する値の乗算を行う)。

【0075】

これは、例えば、Rチャンネルを固定し、フィルタ $1300f_1 - 1$ のRチャンネル用フィルタの中心を、画像データ 511 のRチャンネルの 56×56 の部分に対して、左上から5ずつ右にずらしながら行う。そして、フィルタ $1300f_1 - 1$ のRチャンネル用フィルタの中心が画像データ 511 のRチャンネルの 56×56 の部分の右端まで辿り着いたら、当該Rチャンネル用フィルタの中心を下に5ずらして、再度、左端から行えば良い。

【0076】

(2) 次に、画像データ 511 のGチャンネルに対しても、上記(1)と同様の方法でフィルタ $1300f_1 - 1$ のGチャンネル用フィルタをかける。画像データ 511 のBチャンネルに対しても同様である。

【0077】

(3) フィルタ $1300f_2 - 1 \sim$ フィルタ $1300f_{64} - 1$ についても、上記と同様に、画像データ 511 のRGBの各チャンネルに対してフィルタ処理を順に行う。

【0078】

以上のフィルタ処理により、画像データ 511 から $64 \times 64 \times 3 \times 64$ チャンネルの画像データが生成される。

【0079】

Step 433) 畳み込み処理部 112 は、Step 432 で得られた $64 \times 64 \times 3 \times 64$ チャンネルの画像データの各RGB成分を加算する。この結果、 $64 \times 64 \times 64$ チャンネルの画像データが得られる。

【0080】

Step 434) 畳み込み処理部 112 は、Step 433 で得られた $64 \times 64 \times 64$ チャンネルの画像データの各画素値に対して、バイアスデータ $1100 - 1$ を加算する。

【0081】

ここで、図7(a)は、第1層のバイアスデータ $1100 - 1$ の一例を示す図である。図7(a)に示すように、バイアスデータ $1100 - 1$ は、 1×64 の行列により表される。そこで、畳み込み処理部 112 は、1つめの 64×64 チャンネルの画像データの各画素値に対してバイアスデータ $1100 - 1$ のデータ値 $b_1(1)$ を加算する。同様に、2つ目の 64×64 チャンネルの画像データの各画素値に対してバイアスデータ $1100 - 1$ のデータ値 $b_1(2)$ を加算する。以降、同様に、64個すべての 64×64 チャンネルの画像データの各画素値に対して、それぞれ、バイアスデータ $1100 - 1$ のデータ値を加算する。

【0082】

Step 435) 畳み込み処理部 112 は、Step 434 で得られた $64 \times 64 \times 64$ チャンネルの画像データに対して、所定の活性化関数を適用して出力画像データを得る。所定の活性化関数としては、例えば、任意の画素値 x に対して、 $f(x) = \max(0, x)$ で定義される関数が挙げられる。

【0083】

そして、 $64 \times 64 \times 64$ チャンネルの画像データに対して、活性化関数を適用した後

10

20

30

40

50

、ステップ S 4 1 の加工処理において 0 クリアした部分は取り除き、画像データの中心の 56×56 部分をプーリング処理部 1 1 3 に出力する。したがって、第 1 層において、畳み込み処理部 1 1 2 がプーリング処理部 1 1 3 に出力する画像データの色チャンネルは、 $56 \times 56 \times 64$ である。このようにして得られた $56 \times 56 \times 64$ チャンネルの画像データを「画像データ 5 1 2」と表す。なお、ステップ S 4 1 の加工処理において 0 クリアした部分は、Step 4 3 3 又は Step 4 3 4 で取り除いても良い。

【0084】

プーリング処理部 1 1 3 は、画像データ 5 1 2 を入力して、第 1 層のプーリング処理を行う（ステップ S 4 4）。

【0085】

ここで、第 1 層のプーリング処理について、図 9 を参照しながら説明する。図 9 は、本実施形態の第 1 層のプーリング処理の一例を示す図である。

【0086】

Step 4 4 1) プーリング処理部 1 1 3 は、 $56 \times 56 \times 64$ チャンネルの画像データ 5 1 2 を入力する。

【0087】

Step 4 4 2) プーリング処理部 1 1 3 は、画像データ 5 1 2 の 3×3 の領域内の最大値を出力する処理を繰り返し行い、 $28 \times 28 \times 64$ の画像データ（この画像データを以降「画像データ 5 1 3」とする）を生成する。これは、例えば、以下のようにして行う。

【0088】

(1) 画像データ 5 1 3 の 1 つの 56×56 の画像データ（1 つのチャンネルを固定した 56×56 の画像データ）について、左上を中心とした 3×3 の領域における画素値の最大値を得る。そして、この最大値を、画像データ 5 1 3 の画素位置 (1 , 1) の画素値とする。

【0089】

(2) 次に、 3×3 の領域を右に 2 ずつ移動させながら、それぞれの領域内における画素値の最大値を得て、それぞれ、画像データ 5 1 3 の画素位置 (1 , 2) ~ (1 , 28) の画素値とする。

【0090】

(3) 続いて、 3×3 の領域の中心を下に 2 移動させ、左端から同様に 2 ずつ領域の中心を移動させながら、それぞれの領域内における画素値の最大値を得て、それぞれ、画像データ 5 1 3 の画素位置 (2 , 1) ~ (2 , 28) の画素値とする。以降、同様に、(3 , 1) ~ (28 , 28) の画素値を得る。

【0091】

(4) 上記の (1) ~ (3) を、すべての 56×56 の画像データについて行う。すなわち、上記の (1) ~ (3) を、64 個の 56×56 の画像データについて行う。

【0092】

Step 4 4 3) プーリング処理部 1 1 3 は、画像データ 5 1 3 を第 2 層の畳み込み処理部 1 1 2 に出力する。

【0093】

次に、CNN 処理部 1 1 0 は、畳み込みニューラルネットワークの層を示す変数 n に 1 を加算する（ステップ S 4 5）。

【0094】

次に、CNN 処理部 1 1 0 は、変数 n が、予め設定された N と等しいか否かを判定する（ステップ S 4 6）。変数 n が N と等しい場合、CNN 処理部 1 1 0 は、ステップ S 4 7 に進む。

【0095】

一方、変数 n が N と等しくない場合（すなわち、変数 n が N より小さい場合）、CNN 処理部 1 1 0 は、ステップ S 4 3 に戻る。すなわち、この場合、CNN 処理部 1 1 0 は、

10

20

30

40

50

畳み込みニューラルネットワークの次の層の畳み込み処理及びプーリング処理を行う。

【0096】

本実施形態では、 $N = 2$ であるため、CNN処理部110は、ステップS47に進むものとする。

【0097】

畳み込み処理部112は、画像データ513を入力して、第2層の畳み込み処理を行う(ステップS47)。

【0098】

ここで、第2層の畳み込み処理について、図10を参照しながら説明する。図10は、本実施形態の第2層の畳み込み処理の一例を示す図である。なお、第2層の畳み込み処理は、第1層の畳み込み処理と各データのチャンネル数が異なること以外は同様である。より一般には、第 n 層の畳み込み処理は、他の層の畳み込み処理と各データのチャンネル数が異なること以外は同様である。

10

【0099】

Step471)畳み込み処理部112は、画像データ513を入力する。ここで、入力した画像データ513の色チャンネルは、上述した通り、 $28 \times 28 \times 64$ チャンネルである。

【0100】

Step472)畳み込み処理部112は、重みデータ1200-2からフィルタを生成し、画像データ513に対して、生成したフィルタを用いてフィルタ処理を行う。ここで、重みデータ1200-2のデータ構成及び当該重みデータ1200-2から生成されるフィルタ $1300f_j - 2$ ($j = 1, \dots, 64$)のデータ構成について説明する。

20

【0101】

図11(b)は、第2層の重みデータ1200-2の一例を示す図である。図11(b)に示すように、第2層の重みデータ1200-2は、 1600×64 の行列で表される。なお、重みデータ1200-2の各値 $w_2(i, j)$ は、上述したように、学習データに基づいて予め学習された値である。

【0102】

次に、重みデータ1200-2から生成されるフィルタ $1300f_j - 2$ ($j = 1, \dots, 64$)について説明する。図12は、本実施形態の第2層のフィルタの一例を示す図である。

30

【0103】

図12に示すように、各フィルタ $1300f_j - 2$ ($j = 1, \dots, 64$)は、 5×5 の行列の64個の組で表される。換言すれば、各フィルタ $1300f_j - 2$ ($j = 1, \dots, 64$)は、 $5 \times 5 \times 64$ で表される。

【0104】

ここで、重みデータ1200-2の $w_2(1, 1) \sim w_2(25, 1)$ 、 \dots 、 $w_2(1576, 1) \sim w_2(1600, 1)$ からフィルタ $1300f_1 - 2$ が生成される。同様に、重みデータ1200-2の $w_2(1, 2) \sim w_2(25, 2)$ 、 \dots 、 $w_2(1576, 2) \sim w_2(1600, 2)$ からフィルタ $1300f_2 - 2$ が生成される。 $j = 3, \dots, 64$ の場合も同様である。

40

【0105】

以上のように生成された各フィルタ $1300f_j - 2$ ($j = 1, \dots, 64$)を用いて、畳み込み処理部112は、画像データ513に対してフィルタ処理を行う。畳み込み処理部112は、例えば以下のようにしてフィルタ処理を行う。

【0106】

(1)画像データ513に対してフィルタ $1300f_1 - 2$ をかける(すなわち、画像データ513とフィルタ $1300f_1 - 2$ の対応する値の乗算を行う)。

【0107】

これは、例えば、1つのチャンネルを固定し、フィルタ $1300f_1 - 2$ の中心を、画

50

像データ513の 28×28 の部分の左上から5ずつ右にずらしながら行う。そして、フィルタ $1300f_1 - 2$ の中心が画像データ513の 28×28 の部分の右端まで辿り着いたら、フィルタ $1300f_1 - 2$ の中心を下に5ずらして、再度、左端から行えば良い。

【0108】

(2)次に、画像データ513の他のチャンネルに対しても、上記(1)と同様の方法でフィルタ $1300f_1 - 2$ をかける。この処理をすべてのチャンネル1~64に対して繰り返す。

【0109】

(3)フィルタ $1300f_2 - 2$ ~フィルタ $1300f_{64} - 2$ についても、上記と同様に、1~64のチャンネル毎に、画像データ513の 28×28 の部分に対して、フィルタ処理を順に行う。

10

【0110】

以上のフィルタ処理により、画像データ513から $28 \times 28 \times 64 \times 64$ チャンネルの画像データが生成される。

【0111】

Step473)畳み込み処理部112は、Step472で得られた画像データの 28×28 の部分について、各画素値を1~64チャンネルのそれぞれについて加算する。この結果、 $28 \times 28 \times 64$ チャンネルの画像データが得られる。

【0112】

Step474)畳み込み処理部112は、Step473で得られた $28 \times 28 \times 64$ チャンネルの画像データの各画素値に対して、バイアスデータ1100-2を加算する。

20

【0113】

ここで、図11(a)は、第2層のバイアスデータ1100-2の一例を示す図である。図11(a)に示すように、バイアスデータ1100-2は、 1×64 の行列により表される。そこで、畳み込み処理部112は、1つめの 28×28 チャンネルの画像データの各画素値に対してバイアスデータ1100-2のデータ値 $b_2(1)$ を加算する。同様に、2つ目の 28×28 チャンネルの画像データの各画素値に対してバイアスデータ1100-2のデータ値 $b_2(2)$ を加算する。以降、同様に、64個すべての 28×28 チャンネルの画像データの各画素値に対して、それぞれ、バイアスデータ1100-2のデータ値を加算する。

30

【0114】

Step475)畳み込み処理部112は、Step474で得られた $28 \times 28 \times 64$ チャンネルの画像データに対して、所定の活性化関数を適用して出力画像データを得る。所定の活性化関数としては、例えば、任意の画素値 x に対して、 $f(x) = \max(0, x)$ で定義される関数が挙げられる。このようにして得られた出力画像データが、出力データ520である。このように本実施形態の出力データ520は、 $28 \times 28 \times 64$ チャンネルの画像データである。

【0115】

なお、上記の説明で示されるように、出力データ520は、フィルタ $1300f_j - 2$ の各 j ($j = 1, \dots, 64$)に対応する 28×28 の画像データ(出力データ)の集合とすることができる。すなわち、出力データ520には、フィルタ $1300f_1 - 2$ に対応する 28×28 の出力データ520-1、 \dots 、フィルタ $1300f_{64} - 2$ に対応する 28×28 の出力データ520-64が含まれる。

40

【0116】

次に、図3のステップS33の候補領域の作成処理について、図13を参照しながら説明する。図13は、本実施形態の候補領域の作成処理のフローチャートの一例を示す図である。

【0117】

50

候補領域作成処理部 120 のデータ決定部 121 は、出力データ 520 に含まれる出力データ 520 - 1, …, 出力データ 520 - 64 のそれぞれについて代表値 a_1, \dots, a_{64} を決定する (ステップ S131)。

【0118】

ここで、代表値 a_1, \dots, a_{64} としては、出力データ 520 - 1, …, 出力データ 520 - 64 それぞれのデータ値の最大値とすれば良い。例えば、出力データ 520 - 1 に含まれるデータ値の最大値を代表値 a_1 とすれば良い。他の出力データ 520 - 2, …, 出力データ 520 - 64 についても同様である。ただし、代表値 a_1, \dots, a_{64} は、最大値に限られず、例えば、平均値等を用いても良い。

【0119】

候補領域作成処理部 120 のデータ決定部 121 は、代表値 a_1, \dots, a_{64} に基づいて、出力データ 520 - 1, …, 出力データ 520 - 64 から所定の M 個のデータを決定する (ステップ S132)。ここで、データ決定部 121 は、代表値 a_1, \dots, a_{64} の値が大きい順に (昇順に)、上位 M 個の代表値に対応する出力データを決定すれば良い。

【0120】

以降では、 $M = 3$ として、データ決定部 121 により、出力データ 520 - 2、出力データ 520 - 43、及び出力データ 520 - 47 が決定されたものとする。

【0121】

なお、M の値を大きくすることで、認識処理の精度を向上させることができるが、処理速度は低下する。一方で、M の値を小さくすることで、認識処理の精度は低下するものの処理速度が向上する。したがって、M は、画像処理プログラム 20 の設計者等により、認識対象の画像データ 510 の性質や、認識処理に求められる精度等に応じて適切な値が予め設定される。

【0122】

候補領域作成処理部 120 は、データ決定部 121 により決定された M 個の出力データ 520 のうちの出力データを取得する (ステップ S133)。すなわち、本実施形態では、データ決定部 121 は、出力データ 520 - 2、出力データ 520 - 43、及び出力データ 520 - 47 から一の出力データを取得する。以降では、候補領域作成処理部 120 は、出力データ 520 - 2 を取得したものとして説明する。

【0123】

候補領域作成処理部 120 の境界決定部 122 は、取得された出力データ 520 - 2 について、微分処理を行って、領域分割部 124 により分割される領域の境界を決定する (ステップ S134)。

【0124】

ここで、境界決定部 122 により決定される領域の境界について、図 14 を参照しながら説明する。図 14 は、本実施形態の微分処理の一例を示す図である。

【0125】

図 14 では、一例として、出力データ 520 - 2 について、微分処理を行った場合を示している。図 14 に示すように、境界決定部 122 により微分処理を行い、微分値が負から正に変わる部分を、出力データ 520 - 1 の出力値の谷間として検出する。そして、境界決定部 122 は、検出された出力値の谷間を、境界 D1 及び境界 D2 として決定する。ここで、微分処理には、例えば Sobel フィルタを用いれば良い。

【0126】

候補領域作成処理部 120 の閾値処理部 123 は、閾値処理を行う (ステップ S135)。すなわち、閾値処理部 123 は、予め設定された閾値 (例えば、閾値 = 30) 以下のデータを削除する。

【0127】

ここで、閾値処理部 123 による閾値処理について、図 15 を参照しながら説明する。図 15 は、本実施形態の閾値処理の一例を示す図である。図 15 では、一例として、出力

10

20

30

40

50

データ520-2に対して閾値処理を行った場合を示している。図15に示すように、閾値処理部123は、閾値処理を行って所定の閾値以下のデータ値を削除することにより、出力データ520-2から出力データ521-2を作成する。なお、図15に示す出力データ521において、網掛けで示した部分がデータ値を削除した部分である。

【0128】

候補領域作成処理部120の領域分割部124は、境界決定部122により決定された境界に基づいて、ステップS133で取得された一の出力データが示す画像を複数の領域に分割する(ステップS136)。

【0129】

ここで、領域分割部124により分割される領域について、図16を参照しながら説明する。図16は、本実施形態の領域分割の一例を示す図である。図16では、出力データ521-2が示す画像を境界D1及び境界D2に基づいて分割した例を示している。図16に示すように、出力データ521-2が示す画像は、境界D1及び境界D2に基づいて、領域S1、領域S2、領域S3、及び領域S4に分割される。

10

【0130】

候補領域作成処理部120の候補領域作成部125は、領域分割部124により分割された領域S1~S4について、各領域を含む最小矩形を特定し、当該特定された最小矩形に基づいて候補領域を示す候補領域画像データ530を作成する(ステップS137)。

【0131】

ここで、一例として、領域S1を囲む最小矩形B1を図17に示す。このように最小矩形とは、領域分割部124により分割された領域された領域に外接する矩形のことである。したがって、候補領域作成部125は、各領域S1~S4について、それぞれ最小矩形を特定する。

20

【0132】

そして、候補領域作成部125は、画像データ510が示す画像において、当該特定された最小矩形によって囲まれる領域と対応する領域を候補領域として候補領域画像データ530を作成する。このとき、候補領域作成部125は、画像データ510が示す画像において、最小矩形によって囲まれる領域と対応する領域を、当該画像データ510の解像度を考慮した上で候補領域として候補領域画像データ530を作成する。

【0133】

候補領域作成処理部120は、ステップS132で決定されたすべての出力データに対して、候補領域画像データ530を作成したか否かを判定する(ステップS138)。すなわち、候補領域作成処理部120は、出力データ520-2、出力データ520-43、及び出力データ520-47に対して、ステップS133~ステップS138の処理が実行されたか否かを判定する。

30

【0134】

ステップS132で決定されたすべての出力データに対して、候補領域画像データ530が作成された場合、候補領域作成処理部120は、処理を終了させる。一方、ステップS132で決定された出力データのうち、候補領域画像データ530が作成されていない出力データがある場合、候補領域作成処理部120は、ステップS133に戻る。

40

【0135】

これにより、本実施形態の画像処理装置10では、入力された画像データ510が示す画像において、対象が含まれる領域の候補である候補領域を示す候補領域画像データ530が作成される。しかも、本実施形態の画像処理装置10では、畳み込みニューラルネットワークの第N層における出力データ520を用いて、候補領域画像データ530が作成される。このため、本実施形態の画像処理装置10では、認識処理の精度の低下を防ぎつつ、候補領域を削減させることができる。

【0136】

次に、図3のステップS34のカテゴリ分類処理について、図18を参照しながら説明する。図18は、本実施形態のカテゴリ分類処理のフローチャートの一例を示す図である

50

。

【0137】

CNN処理部110は、1以上の候補領域画像データ530から一の候補領域画像データ530を入力し、入力された候補領域画像データ530に対して、畳み込みニューラルネットワーク処理を行う(ステップS181)。すなわち、CNN処理部110は、入力された候補領域画像データ530に対して、図4で示した畳み込みニューラルネットワーク処理を行う。

【0138】

なお、ステップS181において、CNN処理部110は、予め設定された第N層までの畳み込みニューラルネットワーク処理を行っても良いし、Nより大きい任意の自然数をLとして、第L層までの畳み込みニューラルネットワーク処理を行っても良い。

10

【0139】

ここでは、ステップS181において、CNN処理部110は、第N層までの畳み込みニューラルネットワーク処理を行ったものとして説明する。したがって、ステップS181の処理結果として、CNN処理部110の畳み込み処理部112は、出力データ520と同じデータ構成である $28 \times 28 \times 64$ チャンネルの出力データ531を全結合処理部114に出力する。

【0140】

次に、CNN処理部110の全結合処理部114は、出力データ531を入力して、全結合処理を行う。なお、全結合処理部114は、上述したように、カテゴリの組毎に存在する。したがって、各全結合処理部114は、それぞれ、出力データ531を入力する。

20

【0141】

例えば、カテゴリ数が「人」、「動物」、「車」の3つである場合、全結合処理部114は、カテゴリの組「人」「人以外」に対応する全結合処理部114-1、カテゴリの組「動物」「動物以外」に対応する全結合処理部114-2、及びカテゴリの組「車」「車以外」に対応する全結合処理部114-3の3つが存在する。

【0142】

ここで、全結合処理について、図19を参照しながら説明する。図19は、本実施形態の第3層の全結合処理の一例を示す図である。

【0143】

Step1821)全結合処理部114は、出力データ531を入力する。ここで、入力した出力データ531の色チャンネルは、上述したように、 $28 \times 28 \times 64$ である。

30

【0144】

Step1822)全結合処理部114は、出力データ531の各データ値をベクトル値に変換する。すなわち、 $28 \times 28 \times 64$ チャンネルの出力データ531の各データ値を50176行1列のベクトル値に変換する。ここで、ベクトル値の各成分の値を x_1, \dots, x_{50176} とする。

【0145】

Step1823)全結合処理部114は、それぞれ、バイアスデータ1100-3及び重みデータ1200-3を用いて、積和演算を行う。

40

【0146】

ここで、バイアスデータ1100-3及び重みデータ1200-3について、図20を参照しながら説明する。図20は、本実施形態の第3層のネットワークパラメータの一例を示す図である。

【0147】

図20(a)は、第3層のバイアスデータ1100-3の一例を示す図である。図20(a)に示すように、第3層のバイアスデータ1100-3は、カテゴリ毎のバイアスデータ1100-3₁, バイアスデータ1100-3₂, ...を含む。また、カテゴリ毎のバイアスデータ1100-3_kは、1行2列のベクトル値である。なお、ベクトルの各成分の値 $b_3(k, j)$ は、上述したように、学習データに基づいて予め学習された値で

50

ある。

【 0 1 4 8 】

ここで、 k は、カテゴリを示す数値であるとする。例えば、 $k = 1$ のときカテゴリ「人」を示し、 $k = 2$ のときカテゴリ「動物」を示し、 $k = 3$ のときカテゴリ「車」を示す等である。また、 j は、カテゴリに分類されるか否かを示す数値である。例えば、 $j = 1$ のときは該当のカテゴリに分類される場合を示し、 $j = 2$ のときは該当のカテゴリに分類されない場合（すなわち、該当のカテゴリ以外のカテゴリに分類される場合）を示す。

【 0 1 4 9 】

図 2 0 (b) は、第 3 層の重みデータ 1 2 0 0 - 3 の一例を示す図である。図 2 0 (b) に示すように、第 3 層の重みデータ 1 2 0 0 - 3 は、カテゴリ毎の重みデータ 1 2 0 0 - 3₁ , 重みデータ 1 2 0 0 - 3₂ , … を含む。また、カテゴリ毎の重みデータ 1 2 0 0 - 3_k は、5 0 1 7 6 行 2 列の行列である。なお、この行列の各成分の値 $w_3(i, j, k)$ は、上述したように、学習データに基づいて予め学習された値である。

10

【 0 1 5 0 】

図 1 9 の説明に戻り、全結合処理部 1 1 4 は、それぞれ以下の積和演算を行う。すなわち、カテゴリ k に対して、全結合処理部 1 1 4 - k は、以下の積和演算を行う。

【 0 1 5 1 】

【 数 1 】

20

$$y_j(k) = b_3(k, j) + \sum_{i=1}^{50176} w_3(i, j, k) x_i$$

ここで、 j 及び k の意味は上述した通りである。

【 0 1 5 2 】

Step 1 8 2 4) 全結合処理部 1 1 4 は、Step 1 8 2 3 で得られた $2 \times 1 \times |k|$ のデータを正規化処理部 1 3 0 に出力する。なお、 $|k|$ は、カテゴリ数である。

【 0 1 5 3 】

なお、上記の積和演算の結果が、入力された候補領域画像データ 5 3 0 がカテゴリ k に分類される場合 ($j = 1$ の場合) の算出結果と、当該候補領域画像データ 5 3 0 がカテゴリ k 以外のカテゴリに分類される場合 ($j = 2$ の場合) の算出結果である。

30

【 0 1 5 4 】

これにより、候補領域画像データ 5 3 0 が、あるカテゴリ k に分類されるか否かを数値として判定することができる。例えば、あるカテゴリ k について、 $y_1(k)$ の値が 0 . 7、 $y_2(k)$ の値が 0 . 3 である場合、当該候補領域画像データ 5 3 0 は、カテゴリ k に分類される場合が高いと判定することができる。換言すれば、あるカテゴリ k について、 $y_1(k)$ の値が $y_2(k)$ の値より高い場合、入力された候補領域画像データ 5 3 0 はカテゴリ k に分類される可能性が高いといえる。

40

【 0 1 5 5 】

ただし、上記の算出結果では、各全結合処理部 1 1 4 の出力結果同士の比較ができない場合があるため、次のステップ S 1 8 3 において正規化処理を行う。

【 0 1 5 6 】

正規化処理部 1 3 0 は、全結合処理部 1 1 4 により出力された $2 \times 1 \times |k|$ のデータを入力して、正規化処理を行う (ステップ S 1 8 3) 。

【 0 1 5 7 】

ここで、正規化処理について、図 2 1 を参照しながら説明する。図 2 1 は、本実施形態の正規化処理の一例を示す図である。

【 0 1 5 8 】

50

Step 1831) 正規化処理部 130 は、全結合処理部 114 により出力された $2 \times 1 \times |k|$ のデータを入力する。

【0159】

Step 1832) 正規化処理部 130 は、 $(y_1(k), y_2(k))$ について、カテゴリ毎に以下の式により正規化を行う。

【0160】

【数 2】

$$z_j(k) = \frac{\exp(y_j(k))}{\exp(y_1(k)) + \exp(y_2(k))}$$

10

このようにして得られた $2 \times 1 \times |k|$ が確信度である、このように正規化処理を行うことにより、すべてのカテゴリにおける確信度は 0 以上 1 以下の値に正規化される。このため、異なるカテゴリ同士の確信度を比較することが可能となる。例えば、 $k = 1$ をカテゴリ「人」、 $k = 2$ をカテゴリ「動物」とした場合において、 $z_1(1) = 0.8$ 、 $z_2(1) = 0.2$ 、 $z_1(2) = 0.6$ 、 $z_2(2) = 0.4$ であるとき、入力された候補領域画像データ 530 は、カテゴリ「人」に分類される可能性が高いと言える。

20

【0161】

Step 1833) 正規化処理部 130 は、各カテゴリの確信度を出力部 140 に出力する。

【0162】

以上により、本実施形態の画像処理装置 10 では、入力された画像データが示す画像において、被写体等を示す対象が含まれる領域の候補となる候補領域画像データを作成する。しかも、本実施形態の画像処理装置 10 では、畳み込みニューラルネットワークの予め設定された層の出力結果に基づいて、候補領域画像データを作成することにより、認識処理の精度の低下を防ぎつつ、候補領域画像データの数の削減を図ることができる。

30

【0163】

したがって、本実施形態の画像処理装置 10 は、入力された画像データが示す画像において、対象が含まれる領域と、当該対象が分類されるカテゴリとを識別する識別処理の処理時間を削減することができる。

【0164】

本発明は、具体的に開示された上記の実施形態に限定されるものではなく、特許請求の範囲から逸脱することなく、種々の変形や変更が可能である。

【符号の説明】

【0165】

10	画像処理装置
20	画像処理プログラム
110	CNN 処理部
111	加工部
112	畳み込み処理部
113	プーリング処理部
114	全結合処理部
120	候補領域作成処理部
121	データ決定部
122	境界決定部
123	閾値処理部
124	領域分割部

40

50

- 1 2 5 候補領域作成部
- 1 3 0 正規化処理部
- 1 4 0 出力部

【先行技術文献】

【特許文献】

【0166】

【特許文献1】特許第4322913号公報

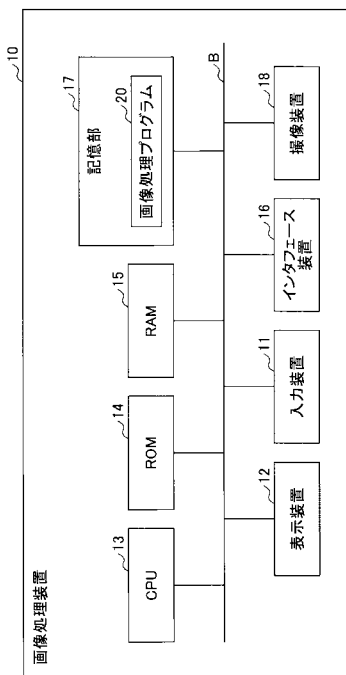
【非特許文献】

【0167】

【非特許文献1】Rich feature hierarchies for accurate object detection and semantic segmentation. Ross Girshick, Jeff Donahue, Trevor Darrell, Jitendra Malik. CVPR 2014.

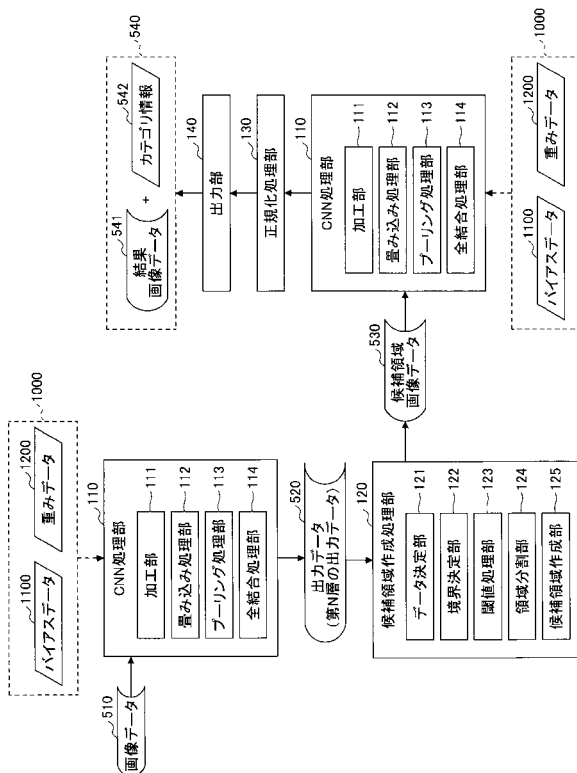
10

【図1】



本実施形態の画像処理装置のハードウェア構成の一例を示す図

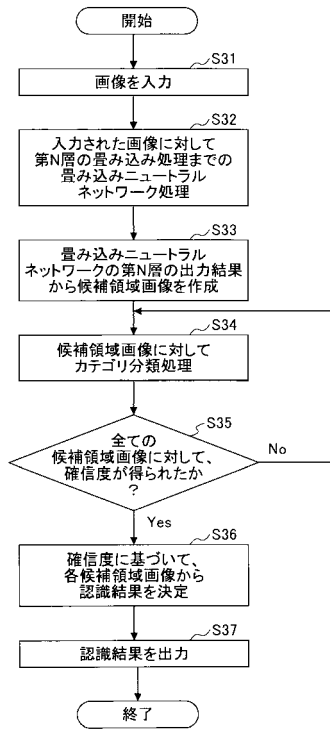
【図2】



本実施形態の画像処理装置の機能構成の一例を示す図

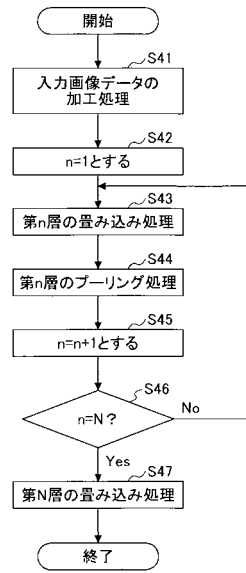
【 図 3 】

本実施形態の画像処理装置の認識処理のフローチャートの一例を示す図



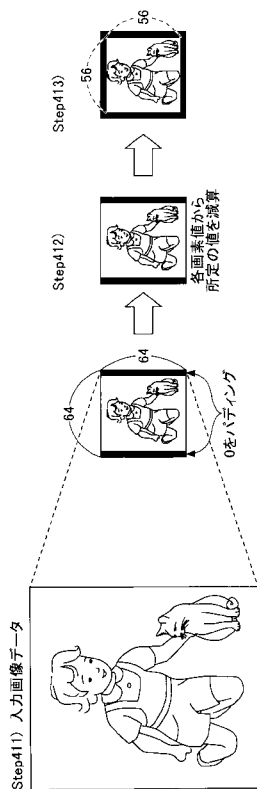
【 図 4 】

本実施形態の畳み込みニューラルネットワーク処理のフローチャートの一例を示す図



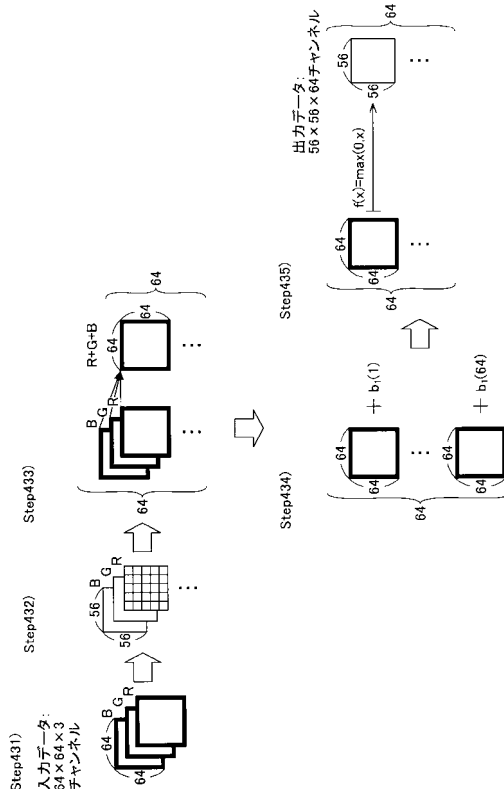
【 図 5 】

本実施形態の入力画像データの加工処理の一例を示す図



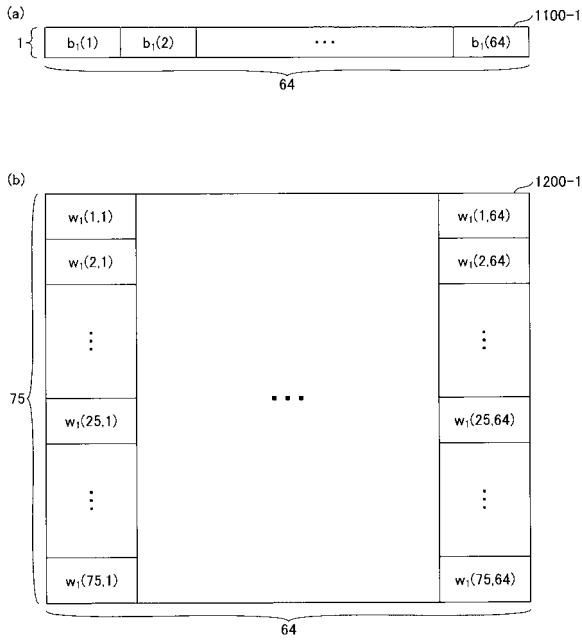
【 図 6 】

本実施形態の第1層の畳み込み処理の一例を示す図



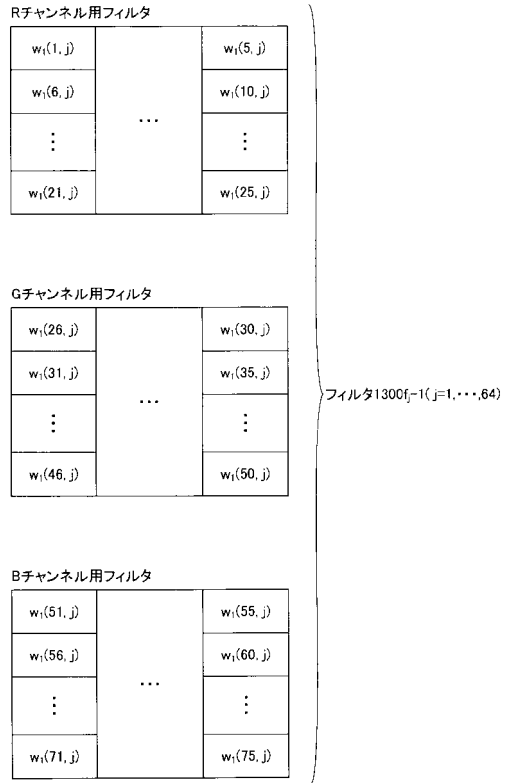
【 図 7 】

本実施形態の第1層のネットワークパラメータの一例を示す図



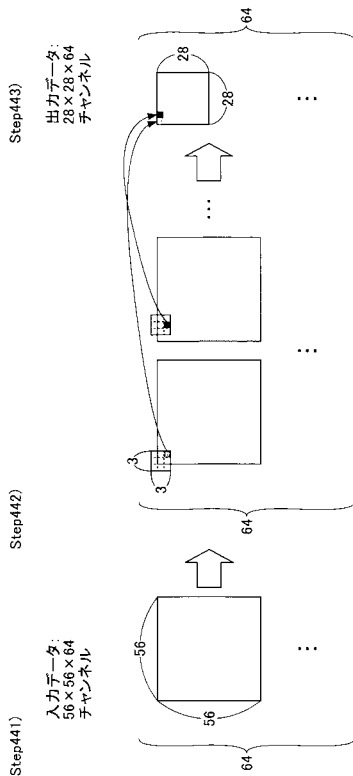
【 図 8 】

本実施形態の第1層のフィルタの一例を示す図



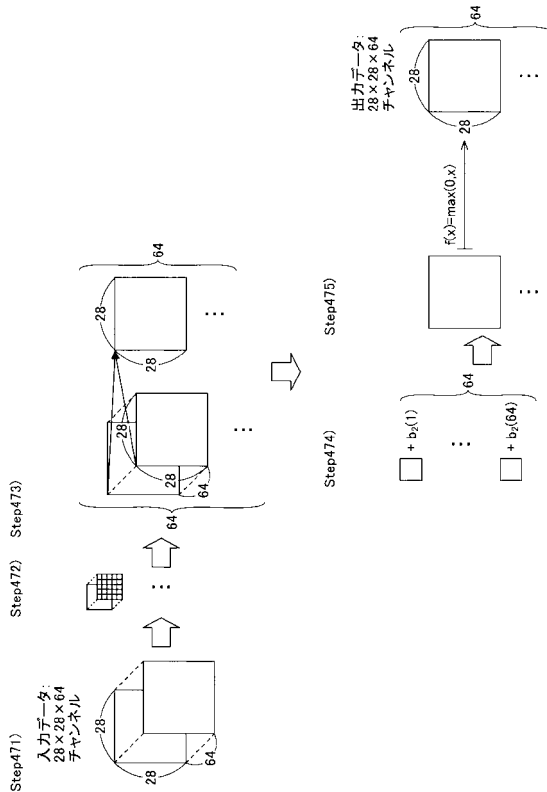
【 図 9 】

本実施形態の第1層のプーリング処理の一例を示す図



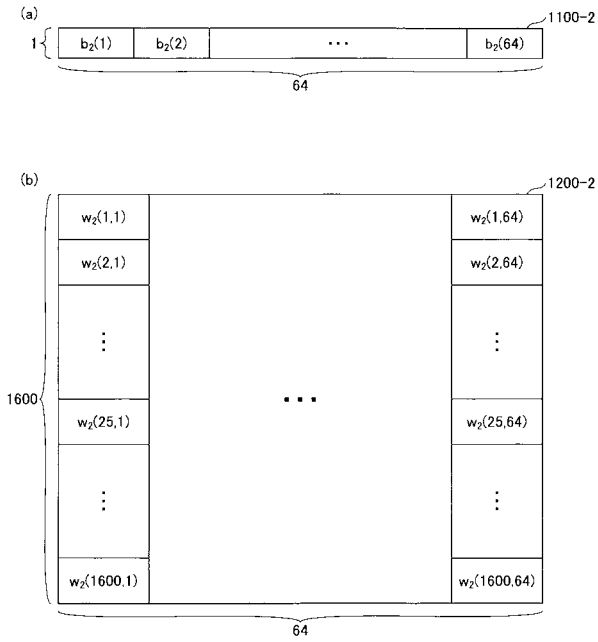
【 図 10 】

本実施形態の第2層の畳み込み処理の一例を示す図



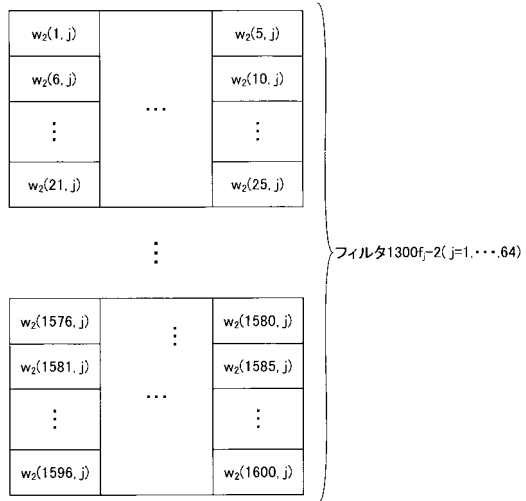
【 図 1 1 】

本実施形態の第2層のネットワークパラメータの一例を示す図



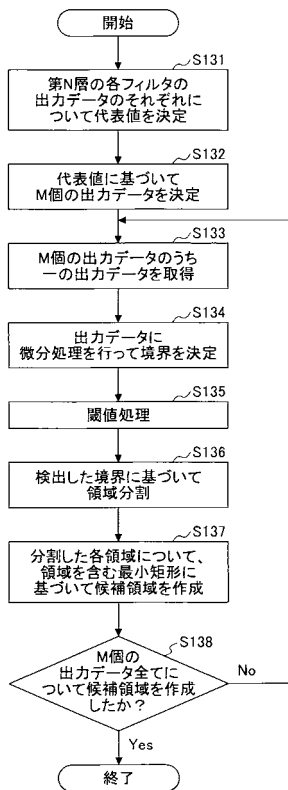
【 図 1 2 】

本実施形態の第2層のフィルタの一例を示す図



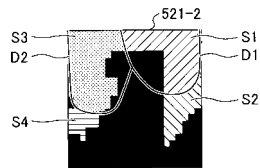
【 図 1 3 】

本実施形態の候補領域の作成処理のフローチャートの一例を示す図



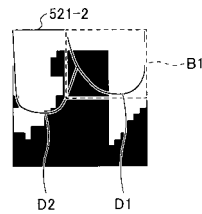
【 図 1 6 】

本実施形態の領域分割の一例を示す図



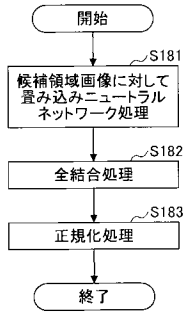
【 図 1 7 】

本実施形態の最小矩形の一例を示す図



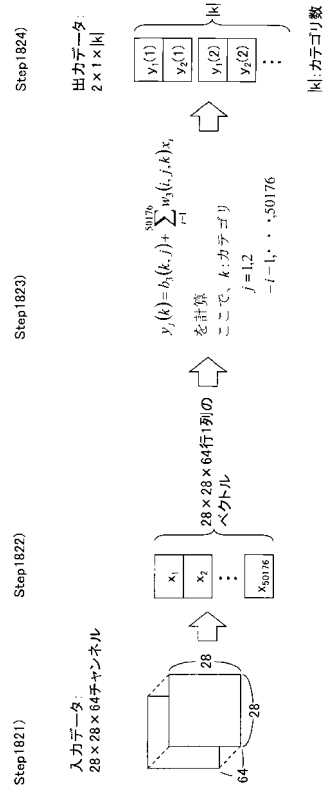
【 図 1 8 】

本実施形態のカテゴリ分類処理のフローチャートの一例を示す図



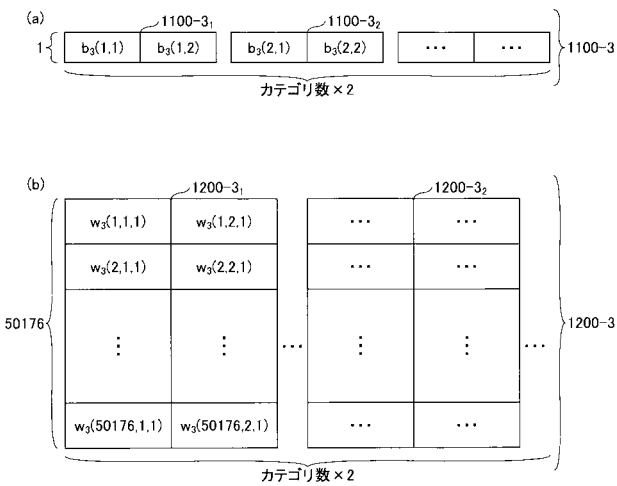
【 図 1 9 】

本実施形態の第3層の全結合処理の一例を示す図



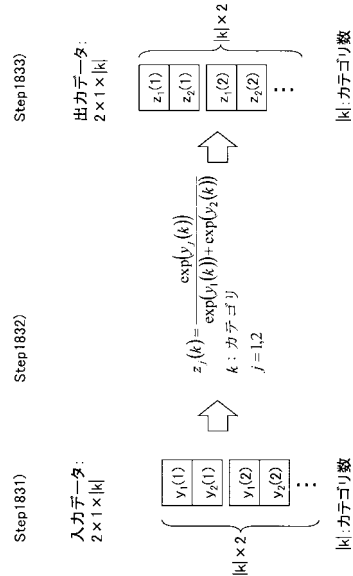
【 図 2 0 】

本実施形態の第3層のネットワークパラメータの一例を示す図

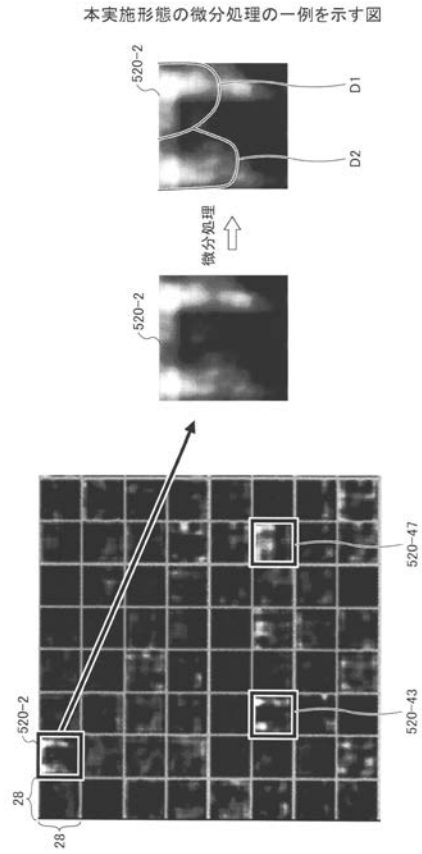


【 図 2 1 】

本実施形態の正規化処理の一例を示す図



【 図 1 4 】



【 図 1 5 】

