

(19) 日本国特許庁 (JP)

(12) 特 許 公 報 (B2)

(11) 特許番号

特許第6859332号
(P6859332)

(45) 発行日 令和3年4月14日 (2021.4.14)

(24) 登録日 令和3年3月29日 (2021.3.29)

(51) Int. Cl.

F I

G 0 6 N 3 / 0 8 (2006.01)

G 0 6 N 3 / 0 8 1 4 0

請求項の数 10 (全 23 頁)

| | | | |
|--------------------|-------------------------------|-----------|---|
| (21) 出願番号 | 特願2018-515936 (P2018-515936) | (73) 特許権者 | 595020643 |
| (86) (22) 出願日 | 平成28年9月7日 (2016.9.7) | | クゥアルコム・インコーポレイテッド |
| (65) 公表番号 | 特表2018-533138 (P2018-533138A) | | Q U A L C O M M I N C O R P O R A T E D |
| (43) 公表日 | 平成30年11月8日 (2018.11.8) | | アメリカ合衆国、カリフォルニア州 9 2 |
| (86) 国際出願番号 | PCT/US2016/050539 | | 1 2 1 - 1 7 1 4、サン・ディエゴ、モア |
| (87) 国際公開番号 | W02017/058479 | | ハウス・ドライブ 5 7 7 5 |
| (87) 国際公開日 | 平成29年4月6日 (2017.4.6) | (74) 代理人 | 100108855 |
| 審査請求日 | 令和1年8月14日 (2019.8.14) | | 弁理士 蔵田 昌俊 |
| (31) 優先権主張番号 | 62/234,559 | (74) 代理人 | 100109830 |
| (32) 優先日 | 平成27年9月29日 (2015.9.29) | | 弁理士 福原 淑弘 |
| (33) 優先権主張国・地域又は機関 | 米国 (US) | (74) 代理人 | 100158805 |
| (31) 優先権主張番号 | 15/081,780 | | 弁理士 井関 守三 |
| (32) 優先日 | 平成28年3月25日 (2016.3.25) | (74) 代理人 | 100112807 |
| (33) 優先権主張国・地域又は機関 | 米国 (US) | | 弁理士 岡田 貴志 |

最終頁に続く

(54) 【発明の名称】 選択的バックプロパゲーション

(57) 【特許請求の範囲】

【請求項 1】

機械学習モデルのためのクラス間のトレーニングデータの平衡を変更する方法であって、

最も少数のメンバーをもつクラスの例の数と現在のクラスの例の数との比からファクタを決定することと、

前記現在のクラスで前記機械学習モデルをトレーニングする間、前記決定されたファクタに基づいて、前記現在のクラスに関連する、バックプロパゲーションプロセスの勾配を変更することと、ここにおいて、前記変更することが、前記最も少数のメンバーをもつ前記クラスの前記例のサンプリングに少なくとも部分的に基づいて前記勾配を選択的に適用することを備え、サンプリング確率は、前記決定されたファクタに基づいて決定される、を備える、方法。

【請求項 2】

前記決定されたファクタに基づいて前記変更することが、前記決定されたファクタで前記勾配をスケールリングすることを備える、請求項 1 に記載の方法。

【請求項 3】

前記最も少数のメンバーをもつ前記クラスの前記サンプリングが、各トレーニングエポックから固定数の例を選択することによって行われる、請求項 1 に記載の方法。

【請求項 4】

前記サンプリングが、トレーニングエポック中の例の交換なしに行われる、請求項 1 に

記載の方法。

【請求項 5】

機械学習モデルのためのクラス間のトレーニングデータの平衡を変更するための装置であって、

最も少数のメンバーをもつクラスの例の数と現在のクラスの例の数との比からファクタを決定するための手段と、

前記現在のクラスで前記機械学習モデルをトレーニングする間、前記決定されたファクタに基づいて、前記現在のクラスに関連する、バックプロパゲーションプロセスの勾配を変更するための手段と、ここにおいて、前記変更するための手段が、前記最も少数のメンバーをもつ前記クラスの前記例のサンプリングに少なくとも部分的に基づいて前記勾配を選択的に適用するための手段を備え、サンプリング確率は、前記決定されたファクタに基づいて決定される、を備える、装置。

10

【請求項 6】

前記決定されたファクタに基づいて前記変更するための手段が、前記決定されたファクタで前記勾配をスケールリングするための手段を備える、請求項 5 に記載の装置。

【請求項 7】

前記クラスの前記サンプリングが、各トレーニングエポックから固定数の例を選択することによって行われる、請求項 5 に記載の装置。

【請求項 8】

前記サンプリングが、トレーニングエポック中の例の交換なしに行われる、請求項 5 に記載の装置。

20

【請求項 9】

メモリをさらに備え、

前記決定するための手段および前記変更するための手段が、前記メモリに結合された少なくとも 1 つのプロセッサを備える、請求項 5 に記載の装置。

【請求項 10】

機械学習モデルのためのクラス間のトレーニングデータの平衡を変更するための非一時的コンピュータ可読媒体であって、前記非一時的コンピュータ可読媒体がそれに記録されたプログラムコードを有し、前記プログラムコードが、実行されると請求項 1 ~ 4 のいずれか一項に記載の方法を実施する、非一時的コンピュータ可読媒体。

30

【発明の詳細な説明】

【技術分野】

【0001】

関連出願の相互参照

[0001]本出願は、その開示全体が参照により本明細書に明確に組み込まれる、2015 年 9 月 29 日に提出された、「SELECTIVE BACKPROPAGATION」と題する米国仮特許出願第 62/234,559 号の利益を主張する。

【背景技術】

【0002】

[0002]本開示のいくつかの態様は、一般に機械学習に関し、より詳細には、機械学習モデルのためのクラス間のトレーニングデータの平衡を変更することに関する。

40

【0003】

[0003]人工ニューロン（たとえば、ニューロンモデル）の相互結合されたグループを備え得る人工ニューラルネットワークは、計算デバイスであるか、または計算デバイスによって実施されるべき方法を表す。

【0004】

[0004]畳み込みニューラルネットワークは、フィードフォワード人工ニューラルネットワークのタイプである。畳み込みニューラルネットワークは、各々が受容野を有し、入力空間を集散的にタイリングするニューロンの集合を含み得る。畳み込みニューラルネットワーク（CNN：convolutional neural network）は多数の適用例を有する。特に、C N

50

Nは、パターン認識および分類の領域内で広く使用されている。

【0005】

[0005]深層信念ネットワークおよび深層畳み込みネットワーク (deep convolutional network) など、深層学習アーキテクチャは、層状 (layered) ニューラルネットワークアーキテクチャであり、ニューロンの第1の層の出力はニューロンの第2の層への入力になり、ニューロンの第2の層の出力はニューロンの第3の層になり、入力し、以下同様である。深層ニューラルネットワークは、特徴の階層 (hierarchy) を認識するようにトレーニングされ得、したがって、それらはオブジェクト認識適用例においてますます使用されている。畳み込みニューラルネットワークのように、これらの深層学習アーキテクチャにおける計算は、1つまたは複数の計算チェーンにおいて構成され得る処理ノードの集団にわたって分散され得る。これらの多層アーキテクチャは、一度に1つの層をトレーニングされ得、バックプロパゲーション (back propagation) を使用して微調整され得る。

10

【0006】

[0006]他のモデルも、オブジェクト認識のために利用可能である。たとえば、サポートベクターマシン (SVM) は、分類のために適用され得る学習ツールである。サポートベクターマシンは、データをカテゴリー分類する分離超平面 (separating hyperplane) (たとえば、決定境界 (decision boundary)) を含む。超平面は、教師あり学習によって定義される。所望の超平面は、トレーニングデータのマージンを増加させる。言い換えれば、超平面は、トレーニング例との最大の最小距離を有するべきである。

20

【0007】

[0007]これらのソリューションは、いくつかの分類ベンチマーク上で優れた結果を達成するが、それらの計算複雑さは極めて高いことがある。さらに、モデルのトレーニングが難しいことがある。

【発明の概要】

【0008】

[0008]一態様では、機械学習モデルのためのクラス間のトレーニングデータの平衡を変更する方法が開示される。本方法は、最も少数のメンバーをもつクラスの例の数と現在のクラスの例の数との比に基づいて、モデルをトレーニングする間、バックプロパゲーションプロセスの勾配を変更することを含む。

30

【0009】

[0009]別の態様は、機械学習モデルのためのクラス間のトレーニングデータの平衡を変更するための装置を開示する。本装置は、最も少数のメンバーをもつクラスの例の数と現在のクラスの例の数との比に基づいて、勾配を変更するためのファクタを決定するための手段を含む。本装置はまた、決定されたファクタに基づいて、現在のクラスに関連する勾配を変更するための手段を含む。

【0010】

[0010]別の態様は、メモリと、メモリに結合された少なくとも1つのプロセッサとを有するワイヤレス通信を開示する。(1つまたは複数の) プロセッサは、最も少数のメンバーをもつクラスの例の数と現在のクラスの例の数との比に基づいて、モデルをトレーニングする間、バックプロパゲーションプロセスの勾配を変更するように構成される。

40

【0011】

[0011]別の態様は、それに記録された非一時的プログラムコードを有する非一時的コンピュータ可読媒体であって、(1つまたは複数の) プロセッサによって実行されたとき、(1つまたは複数の) プロセッサに、最も少数のメンバーをもつクラスの例の数と現在のクラスの例の数との比に少なくとも部分的に基づいて、モデルをトレーニングする間、バックプロパゲーションプロセスの勾配を変更する動作を実行させる、非一時的コンピュータ可読媒体を開示する。

【0012】

[0012]本開示の追加の特徴および利点が、以下で説明される。本開示は、本開示の同じ目的を実行するための他の構造を変更または設計するための基礎として容易に利用され得

50

ることを、当業者は諒解されたい。また、そのような等価な構成が、添付の特許請求の範囲に記載の本開示の教示から逸脱しないことを、当業者は了解されたい。さらなる目的および利点とともに、本開示の編成と動作の方法の両方に関して、本開示を特徴づけると考えられる新規の特徴は、添付の図に関連して以下の説明を検討するとより良く理解されよう。ただし、図の各々は、例示および説明のみの目的で与えられたものであり、本開示の限界を定めるものではないことを明確に理解されたい。

【図面の簡単な説明】

【0013】

[0013]本開示の特徴、特性、および利点は、全体を通じて同様の参照符号が同様のものを指す図面とともに、以下に記載される発明を実施するための形態を読めばより明らかになる。

10

【図1】[0014]本開示のいくつかの態様による、汎用プロセッサを含むシステムオンチップ(SOC)を使用してニューラルネットワークを設計する例示的なインプリメンテーションを示す図。

【図2】[0015]本開示の態様による、システムの例示的なインプリメンテーションを示す図。

【図3A】[0016]本開示の態様による、ニューラルネットワークを示す図。

【図3B】[0017]本開示の態様による、例示的な深層畳み込みネットワーク(DCN)を示すブロック図。

【図4】[0018]本開示の態様による、人工知能(AI)機能をモジュール化し得る例示的なソフトウェアアーキテクチャを示すブロック図。

20

【図5】[0019]本開示の態様による、スマートフォン上のAIアプリケーションのランタイム動作を示すブロック図。

【図6】[0020]本開示の態様による、トレーニングデータを平衡させるための方法を示す図。

【図7】[0021]本開示の態様による、トレーニングデータを平衡させるための全体的例を示す図。

【図8】[0022]本開示の態様による、トレーニングデータを平衡させるための方法を示す図。

【発明を実施するための形態】

30

【0014】

[0023]添付の図面に関して以下に記載される発明を実施するための形態は、様々な構成を説明するものであり、本明細書で説明される概念が実施され得る構成のみを表すものではない。発明を実施するための形態は、様々な概念の完全な理解を与えるための具体的な詳細を含む。ただし、これらの概念はこれらの具体的な詳細なしに実施され得ることが当業者には明らかであろう。いくつかの事例では、そのような概念を不明瞭にしないように、よく知られている構造および構成要素がブロック図の形式で示される。

【0015】

[0024]これらの教示に基づいて、本開示の範囲は、本開示の他の態様とは無関係にインプリメントされるにせよ、本開示の他の態様と組み合わせてインプリメントされるにせよ、本開示のいかなる態様をもカバーするものであることを、当業者なら諒解されたい。たとえば、記載された態様をいくつか使用しても、装置はインプリメントされ得るか、または方法は実施され得る。さらに、本開示の範囲は、記載された本開示の様々な態様に加えてまたはそれらの態様以外に、他の構造、機能、または構造および機能を使用して実施されるそのような装置または方法をカバーするものとする。開示される本開示のいずれの態様も、請求項の1つまたは複数の要素によって実施され得ることを理解されたい。

40

【0016】

[0025]「例示的」という単語は、本明細書では「例、事例、または例示の働きをすること」を意味するために使用される。「例示的」として本明細書で説明されるいかなる態様も、必ずしも他の態様よりも好適または有利であると解釈されるべきであるとは限らない

50

。

【 0 0 1 7 】

[0026] 本明細書では特定の態様が説明されるが、これらの態様の多くの変形および置換は本開示の範囲内に入る。好適な態様のいくつかの利益および利点が説明されるが、本開示の範囲は特定の利益、使用、または目的に限定されるものではない。むしろ、本開示の態様は、様々な技術、システム構成、ネットワーク、およびプロトコルに広く適用可能であるものとし、それらのいくつかは、例として、図および好適な態様についての以下の説明において示される。発明を実施するための形態および図面は、本開示を限定するものではなく説明するものにすぎず、本開示の範囲は添付の特許請求の範囲およびその均等物によって定義される。

10

選択的バックプロパゲーション

[0027] 本開示の態様は、機械学習モデルにおいてクラス間のトレーニングデータの平衡を変更することを対象とする。特に、入力段においてトレーニングデータを操作し、各クラスについての例の数を調節するのではなく、本開示の態様は、勾配段における調節を対象とする。本開示の様々な態様では、データセット中のクラス例頻度に基づいて勾配を調節するかまたは選択的に適用するために、コスト関数を変更するために、選択的バックプロパゲーションが利用される。特に、勾配は、各クラスについての例の実際のまたは予想される頻度に基づいて調節され得る。

【 0 0 1 8 】

[0028] 図 1 に、本開示のいくつかの態様による、汎用プロセッサ (CPU) またはマルチコア汎用プロセッサ (CPU) 102 など、少なくとも 1 つのプロセッサを含み得るシステムオンチップ (SOC) 100 を使用する、上述の選択的バックプロパゲーションの例示的なインプリメンテーションを示す。変数 (たとえば、ニューラル信号およびシナプス荷重)、計算デバイスに関連するシステムパラメータ (たとえば、重みをもつニューラルネットワーク)、遅延、周波数ピン情報、およびタスク情報が、ニューラル処理ユニット (NPU) 108 に関連するメモリブロックに記憶されるか、CPU 102 に関連するメモリブロックに記憶されるか、グラフィックス処理ユニット (GPU) 104 に関連するメモリブロックに記憶されるか、デジタル信号プロセッサ (DSP) 106 に関連するメモリブロックに記憶されるか、専用メモリブロック 118 に記憶され得るか、または複数のブロックにわたって分散され得る。汎用プロセッサ 102 において実行される命令が、CPU 102 に関連するプログラムメモリからロードされ得るか、または専用メモリブロック 118 からロードされ得る。

20

30

【 0 0 1 9 】

[0029] SOC 100 はまた、GPU 104、DSP 106 など、特定の機能に適合された追加の処理ブロックと、第 4 世代ロングタームエボリューション (4G LTE (登録商標)) 接続性、無認可 Wi-Fi (登録商標) 接続性、USB 接続性、Bluetooth (登録商標) 接続性などを含み得る接続性ブロック 110 と、たとえば、ジェスチャーを検出および認識し得るマルチメディアプロセッサ 112 とを含み得る。一つのインプリメンテーションでは、NPU は、CPU、DSP、および / または GPU においてインプリメントされる。SOC 100 はまた、センサプロセッサ 114、画像信号プロセッサ (ISP)、および / または全地球測位システムを含み得るナビゲーション 120 を含み得る。

40

【 0 0 2 0 】

[0030] SOC 100 は ARM 命令セットに基づき得る。本開示の別の態様では、汎用プロセッサ 102 にロードされる命令は、機械学習モデルをトレーニングする間、バックプロパゲーションプロセスの勾配を変更するためのコードを備え得る。変更することは、最も少数のメンバーをもつクラスの例の数と現在のクラスの例の数との比に基づく。変更することは、現在のクラスに関連する勾配に適用される。

【 0 0 2 1 】

[0031] 図 2 に、本開示のいくつかの態様による、システム 200 の例示的なインプリメ

50

ンテーションを示す。図2に示されているように、システム200は、本明細書で説明される方法の様々な動作を実施し得る複数のローカル処理ユニット202を有し得る。各ローカル処理ユニット202は、ローカル状態メモリ204と、ニューラルネットワークのパラメータを記憶し得るローカルパラメータメモリ206とを備え得る。さらに、ローカル処理ユニット202は、ローカルモデルプログラムを記憶するためのローカル(ニューロン)モデルプログラム(LMP)メモリ208と、ローカル学習プログラムを記憶するためのローカル学習プログラム(LLP)メモリ210と、ローカル接続メモリ212とを有し得る。さらに、図2に示されているように、各ローカル処理ユニット202は、ローカル処理ユニットのローカルメモリのための構成を与えるための構成プロセッサユニット214、およびローカル処理ユニット202間のルーティングを与えるルーティング接続処理ユニット216とインターフェースし得る。

10

【0022】

[0032]深層学習アーキテクチャは、各層において連続的により高い抽象レベルで入力を表現するように学習し、それにより、入力データの有用な特徴表現を蓄積することによって、オブジェクト認識タスクを実施し得る。このようにして、深層学習は、旧来の機械学習の主要なボトルネックに対処する。深層学習の出現より前に、オブジェクト認識問題に対する機械学習手法は、場合によっては浅い分類器(shallow classifier)と組み合わせて、人的に設計された特徴に大きく依拠していることがある。浅い分類器は、たとえば、入力がどのクラスに属するかを予測するために、特徴ベクトル成分の重み付き和がしきい値と比較され得る2クラス線形分類器であり得る。人的に設計された特徴は、領域の専門知識をもつ技術者によって特定の領域に適合されたテンプレートまたはカーネルであり得る。対照的に、深層学習アーキテクチャは、人間の技術者が設計し得るものと同様である特徴を表現するように学習するが、トレーニングを通してそれを行い得る。さらに、深層ネットワークは、人間が考慮していないことがある新しいタイプの特徴を表現し、認識するように学習し得る。

20

【0023】

[0033]深層学習アーキテクチャは特徴の階層を学習し得る。たとえば、視覚データが提示された場合、第1の層は、エッジなど、入力ストリーム中の比較的単純な特徴を認識するように学習し得る。別の例では、聴覚データが提示された場合、第1の層は、特定の周波数におけるスペクトル電力を認識するように学習し得る。第1の層の出力を入力として取る第2の層は、視覚データの場合の単純な形状、または聴覚データの場合の音の組合せなど、特徴の組合せを認識するように学習し得る。たとえば、上位層は、視覚データ中の複雑な形状、または聴覚データ中の単語を表現するように学習し得る。さらに上位の層は、共通の視覚オブジェクトまたは発話フレーズを認識するように学習し得る。

30

【0024】

[0034]深層学習アーキテクチャは、自然階層構造を有する問題に適用されたとき、特にうまく機能し得る。たとえば、原動機付き車両の分類は、ホイール、フロントガラス、および他の特徴を認識するための第1の学習から恩恵を受け得る。これらの特徴は、車、トラック、および飛行機を認識するために、異なる方法で、上位層において組み合わせられ得る。

40

【0025】

[0035]ニューラルネットワークは、様々な結合性パターンを用いて設計され得る。フィードフォワードネットワークでは、情報が下位層から上位層に受け渡され、所与の層における各ニューロンは、上位層におけるニューロンに通信する。上記で説明されたように、フィードフォワードネットワークの連続する層において、階層表現が蓄積され得る。ニューラルネットワークはまた、リカレントまたは(トップダウンとも呼ばれる)フィードバック結合を有し得る。リカレント結合では、所与の層におけるニューロンからの出力は、同じ層における別のニューロンに通信され得る。カレントアーキテクチャは、ニューラルネットワークに順次配信される入力データチャンクのうちの2つ以上にわたるパターンを認識するのに役立ち得る。所与の層におけるニューロンから下位層におけるニューロンへ

50

の結合は、フィードバック（またはトップダウン）結合と呼ばれる。高レベルの概念の認識が、入力の特定の低レベルの特徴を弁別することを助け得るとき、多くのフィードバック結合をもつネットワークが役立ち得る。

【 0 0 2 6 】

[0036]図 3 A を参照すると、ニューラルネットワークの層間の結合は全結合 3 0 2 または局所結合 3 0 4 であり得る。全結合ネットワーク 3 0 2 では、第 1 の層におけるニューロンは、第 2 の層における各ニューロンが第 1 の層におけるあらゆるニューロンから入力を受信するように、その出力を第 2 の層におけるあらゆるニューロンに通信し得る。代替的に、局所結合ネットワーク 3 0 4 では、第 1 の層におけるニューロンは、第 2 の層における限られた数のニューロンに結合され得る。畳み込みネットワーク 3 0 6 は、局所結合であり得、第 2 の層における各ニューロンのための入力に関連する結合強度が共有されるようにさらに構成される（たとえば、3 0 8）。より一般的には、ネットワークの局所結合層は、層における各ニューロンが同じまたは同様の結合性パターンを有するように構成されるが、異なる値を有し得る結合強度で構成され得る（たとえば、3 1 0、3 1 2、3 1 4、および 3 1 6）。局所結合の結合性パターンは、所与の領域中の上位層ニューロンが、ネットワークへの総入力のうちの制限された部分のプロパティにトレーニングを通して調整された入力を受信し得るので、上位層において空間的に別個の受容野を生じ得る。

10

【 0 0 2 7 】

[0037]局所結合ニューラルネットワークは、入力の空間ロケーションが有意味である問題に好適であり得る。たとえば、車載カメラからの視覚特徴を認識するように設計されたネットワーク 3 0 0 は、画像の下側部分対上側部分とのそれらの関連付けに依存して、異なるプロパティをもつ上位層ニューロンを発達させ得る。画像の下側部分に関連するニューロンは、たとえば、車線区分線を認識するように学習し得るが、画像の上側部分に関連するニューロンは、交通信号、交通標識などを認識するように学習し得る。

20

【 0 0 2 8 】

[0038]深層畳み込みネットワーク（DCN）が、教師あり学習を用いてトレーニングされ得る。トレーニング中に、DCNは、速度制限標識のクロップされた画像 3 2 6 など、画像を提示され得、次いで、出力 3 2 2 を生成するために、「フォワードパス」が計算され得る。出力 3 2 2 は、「標識」、「6 0」、および「1 0 0」など、特徴に対応する値のベクトルであり得る。ネットワーク設計者は、DCNが、出力特徴ベクトルにおけるニューロンのうちのいくつか、たとえば、トレーニングされたネットワーク 3 0 0 のための出力 3 2 2 に示されているように「標識」および「6 0」に対応するニューロンについて、高いスコアを出力することを希望し得る。トレーニングの前に、DCNによって生成された出力は不正確である可能性があり、したがって、実際の出力とターゲット出力との間で誤差が計算され得る。次いで、DCNの重みは、DCNの出力スコアがターゲットとより密接に整合されるように調節され得る。

30

【 0 0 2 9 】

[0039]重みを調節するために、学習アルゴリズムは、重みのための勾配ベクトルを計算し得る。勾配は、重みがわずかに調節された場合に、誤差が増加または減少する量を示し得る。最上層において、勾配は、最後から 2 番目の層における活性化されたニューロンと出力層におけるニューロンとを結合する重みの値に直接対応し得る。下位層では、勾配は、重みの値と、上位層の計算された誤差勾配とに依存し得る。次いで、重みは、誤差を低減するように調節され得る。重みを調節するこの様式は、それがニューラルネットワークを通して「バックワードパス」を伴うので、「バックプロパゲーション」と呼ばれることがある。

40

【 0 0 3 0 】

[0040]実際には、重みの誤差勾配は、計算された勾配が真の誤差勾配を近似するように、少数の例にわたって計算され得る。この近似方法は、確率的勾配降下（stochastic gradient descent）と呼ばれることがある。システム全体の達成可能な誤差レートが減少し

50

なくなるまで、または誤差レートがターゲットレベルに達するまで、確率的勾配降下が繰り返され得る。

【 0 0 3 1 】

[0041]学習の後に、D C Nは新しい画像 3 2 6 を提示され得、ネットワークを通したフォワードパスは、D C Nの推論または予測と見なされ得る出力 3 2 2 をもたらし得る。

【 0 0 3 2 】

[0042]深層信念ネットワーク (D B N : deep belief network) は、隠れノードの複数の層を備える確率モデルである。D B Nは、トレーニングデータセットの階層表現を抽出するために使用され得る。D B Nは、制限ボルツマンマシン (R B M : Restricted Boltzmann Machine) の層を積層することによって取得され得る。R B Mは、入力セットにわたる確率分布を学習することができる人工ニューラルネットワークのタイプである。R B Mは、各入力にそれぞれにカテゴリ分類されるべきクラスに関する情報の不在下で確率分布を学習することができるので、R B Mは、教師なし学習においてしばしば使用される。ハイブリッド教師なしおよび教師ありパラダイムを使用して、D B Nの下部R B Mは、教師なし様式でトレーニングされ得、特徴抽出器として働き得、上部R B Mは、(前の層からの入力とターゲットクラスとの同時分布上で) 教師あり様式でトレーニングされ得、分類器として働き得る。

【 0 0 3 3 】

[0043]深層畳み込みネットワーク (D C N) は、追加のプーリング層および正規化層で構成された、畳み込みネットワークのネットワークである。D C Nは、多くのタスクに関して最先端の性能を達成している。D C Nは、入力と出力ターゲットの両方が、多くの標本について知られており、勾配降下方法の使用によってネットワークの重みを変更するために使用される、教師あり学習を使用してトレーニングされ得る。

【 0 0 3 4 】

[0044]D C Nは、フィードフォワードネットワークであり得る。さらに、上記で説明されたように、D C Nの第 1 の層におけるニューロンから次の上位層におけるニューロンのグループへの結合は、第 1 の層におけるニューロンにわたって共有される。D C Nのフィードフォワードおよび共有結合は、高速処理のために活用され得る。D C Nの計算負担は、たとえば、リカレントまたはフィードバック結合を備える同様のサイズのニューラルネットワークのそれよりもはるかに少ないことがある。

【 0 0 3 5 】

[0045]畳み込みネットワークの各層の処理は、空間的に不変のテンプレートまたは基底投射と見なされ得る。入力が、カラー画像の赤色、緑色、および青色チャネルなど、複数のチャネルに最初に分解された場合、その入力に関してトレーニングされた畳み込みネットワークは、画像の軸に沿った 2 つの空間次元と、色情報をキャプチャする第 3 の次元とをもつ、3 次元であると思われ得る。畳み込み結合の出力は、後続の層 3 1 8 および 3 2 0 において特徴マップを形成すると考えられ、特徴マップ (たとえば、3 2 0) の各要素が、前の層 (たとえば、3 1 8) における様々なニューロンから、および複数のチャネルの各々から入力を受信し得る。特徴マップにおける値は、整流 (rectification)、 $\max(0, x)$ など、非線形性を用いてさらに処理され得る。隣接するニューロンからの値は、さらにプーリングされ得、これは、ダウンサンプリングに対応し、さらなる局所不変性と次元削減とを与え得る。白色化に対応する正規化はまた、特徴マップにおけるニューロン間のラテラル抑制によって適用され得る。

【 0 0 3 6 】

[0046]深層学習アーキテクチャの性能は、より多くのラベリングされたデータポイントが利用可能となるにつれて、または計算能力が増加するにつれて、向上し得る。現代の深層ニューラルネットワークは、ほんの 1 5 年前に一般的な研究者にとって利用可能であったものより数千倍も大きいコンピューティングリソースを用いて、ルーチン的にトレーニングされる。新しいアーキテクチャおよびトレーニングパラダイムが、深層学習の性能をさらに高め得る。整流された線形ユニット (rectified linear unit) は、勾配消失 (van

10

20

30

40

50

ishing gradients)として知られるトレーニング問題を低減し得る。新しいトレーニング技法は、過学習(over-fitting)を低減し、したがって、より大きいモデルがより良い汎化を達成することを可能にし得る。カプセル化技法は、所与の受容野においてデータを抽出し、全体的性能をさらに高め得る。

【0037】

[0047]図3Bは、例示的な深層畳み込みネットワーク350を示すブロック図である。深層畳み込みネットワーク350は、結合性および重み共有に基づく、複数の異なるタイプの層を含み得る。図3Bに示されているように、例示的な深層畳み込みネットワーク350は、複数の畳み込みブロック(たとえば、C1およびC2)を含む。畳み込みブロックの各々は、畳み込み層と、正規化層(LNorm)と、プーリング層とで構成され得る。畳み込み層は、1つまたは複数の畳み込みフィルタを含み得、これは、特徴マップを生成するために入力データに適用され得る。2つの畳み込みブロックのみが示されているが、本開示はそのように限定しておらず、代わりに、設計選好に従って、任意の数の畳み込みブロックが深層畳み込みネットワーク350中に含まれ得る。正規化層は、畳み込みフィルタの出力を正規化するために使用され得る。たとえば、正規化層は、白色化またはラテラル抑制を行い得る。プーリング層は、局所不変性および次元削減のために、空間にわたってダウンサンプリングアグリゲーションを行い得る。

【0038】

[0048]たとえば、深層畳み込みネットワークの並列フィルタバンクは、高性能および低電力消費を達成するために、随意にARM命令セットに基づいて、SOC100のCPU102またはGPU104にロードされ得る。代替実施形態では、並列フィルタバンクは、SOC100のDSP106またはISP116にロードされ得る。さらに、DCNは、センサー114およびナビゲーション120に専用の処理ブロックなど、SOC上に存在し得る他の処理ブロックにアクセスし得る。

【0039】

[0049]深層畳み込みネットワーク350はまた、1つまたは複数の全結合層(たとえば、FC1およびFC2)を含み得る。深層畳み込みネットワーク350は、ロジスティック回帰(LR)層をさらに含み得る。深層畳み込みネットワーク350の各層の間には、更新されるべき重み(図示せず)がある。各層の出力は、第1の畳み込みブロックC1において供給された入力データ(たとえば、画像、オーディオ、ビデオ、センサーデータおよび/または他の入力データ)から階層特徴表現を学習するために、深層畳み込みネットワーク350中の後続の層の入力として働き得る。

【0040】

[0050]図4は、人工知能(AI)機能をモジュール化し得る例示的なソフトウェアアーキテクチャ400を示すブロック図である。アーキテクチャを使用して、SOC420の様々な処理ブロック(たとえば、CPU422、DSP424、GPU426および/またはNPU428)に、アプリケーション402のランタイム動作中に計算をサポートすることを実施させ得るアプリケーション402が設計され得る。

【0041】

[0051]AIアプリケーション402は、たとえば、デバイスが現在動作するロケーションを示すシーンの検出および認識を与え得る、ユーザ空間404において定義されている機能と呼び出すように構成され得る。AIアプリケーション402は、たとえば、認識されたシーンがオフィス、講堂、レストラン、または湖などの屋外環境であるかどうかに応じて別様に、マイクロフォンおよびカメラを構成し得る。AIアプリケーション402は、現在のシーンの推定を与えるために、SceneDetectアプリケーションプログラミングインターフェース(API)406において定義されているライブラリに関連するコンパイルされたプログラムコードへの要求を行い得る。この要求は、たとえば、ビデオおよび測位データに基づくシーン推定を与えるように構成された深層ニューラルネットワークの出力に最終的に依拠し得る。

【0042】

10

20

30

40

50

[0052]さらに、ランタイムフレームワークのコンパイルされたコードであり得るランタイムエンジン408が、AIアプリケーション402にとってアクセス可能であり得る。AIアプリケーション402は、たとえば、ランタイムエンジンに、特定の時間間隔における、またはアプリケーションのユーザインターフェースによって検出されたイベントによってトリガされた、シーン推定を要求させ得る。シーンを推定させられたとき、ランタイムエンジンは、SOC420上で実行している、Linux（登録商標）カーネル412など、オペレーティングシステム410に信号を送り得る。オペレーティングシステム410は、CPU422、DSP424、GPU426、NPU428、またはそれらの何らかの組合せ上で、計算を実施させ得る。CPU422は、オペレーティングシステムによって直接アクセスされ得、他の処理ブロックは、DSP424のための、GPU426のための、またはNPU428のためのドライバ414~418など、ドライバを通してアクセスされ得る。例示的な例では、深層ニューラルネットワークは、CPU422およびGPU426など、処理ブロックの組合せ上で動作するように構成され得るか、または存在する場合、NPU428上で動作させられ得る。

【0043】

[0053]図5は、スマートフォン502上のAIアプリケーションのランタイム動作500を示すブロック図である。AIアプリケーションは、画像506のフォーマットを変換し、次いで画像508をクロップおよび/またはリサイズするように（たとえば、JAV A（登録商標）プログラミング言語を使用して）構成され得る前処理モジュール504を含み得る。次いで、前処理された画像は、視覚入力に基づいてシーンを検出および分類するように（たとえば、Cプログラミング言語を使用して）構成され得るSceneDetectバックエンドエンジン512を含んでいる分類アプリケーション510に通信され得る。SceneDetectバックエンドエンジン512は、スケーリング516およびクロッピング518によって、画像をさらに前処理514するように構成され得る。たとえば、画像は、得られた画像が224ピクセル×224ピクセルであるように、スケーリングされ、クロップされ得る。これらの次元は、ニューラルネットワークの入力次元にマッピングし得る。ニューラルネットワークは、SOC100の様々な処理ブロックに、深層ニューラルネットワークを用いて画像ピクセルをさらに処理させるように、深層ニューラルネットワークブロック520によって構成され得る。次いで、深層ニューラルネットワークの結果は、しきい値処理522され、分類アプリケーション510中の指数平滑化ブロック524を通され得る。次いで、平滑化された結果は、スマートフォン502の設定および/またはディスプレイの変更を生じ得る。

【0044】

[0054]一構成では、機械学習モデルは、機械学習モデルをトレーニングする間、バックプロパゲーションプロセスの勾配を変更するために構成される。モデルは、変更手段のための手段、および/または決定するための手段を含む。一態様では、変更手段、および/または決定手段は、具陳された機能を実行するように構成された、汎用プロセッサ102、汎用プロセッサ102に関連するプログラムメモリ、メモリブロック118、ローカル処理ユニット202、およびまたはルーティング接続処理ユニット216であり得る。別の構成では、上述の手段は、上述の手段によって具陳された機能を実行するように構成された任意のモジュールまたは任意の装置であり得る。

【0045】

[0055]別の態様では、変更手段は、勾配をスケーリングするための手段を含み得る。随意に、変更手段は、勾配を選択的に適用するための手段を含み得る。

【0046】

[0056]本開示のいくつかの態様によれば、各ローカル処理ユニット202は、モデルの所望の1つまたは複数の機能的特徴に基づいてモデルのパラメータを決定し、決定されたパラメータがさらに適合、調整および更新されるように、1つまたは複数の機能的特徴を所望の機能的特徴のほうへ発達させるように構成され得る。

【0047】

10

20

30

40

50

[0057]多くの機械学習プロセスでは、学習された分類関数の出力と所望の出力との間の誤差を定量化するために、コスト関数を使用される。機械学習プロセスの目的は、このコスト関数を最小限に抑えるように、学習された分類関数のパラメータを変えることである。分類問題では、コスト関数は、しばしば、何らかの入力に関連する実際のクラスラベルと、その入力に関数を適用することによって達成される予測されたクラスラベルとのログ確率ペナルティ関数である。トレーニングは、学習された分類関数のパラメータを変更するプロセスである。トレーニング中に、例示的な入力とそれらの関連するラベルとが、機械学習プロセスに提示される。プロセスは、現在の学習された分類関数パラメータが与えられれば予測されたラベルを見つけ、コスト関数を評価し、学習された分類関数のパラメータを何らかの更新学習則に従って変える。

10

【0048】

[0058]トレーニングプロセス中に、不平衡トレーニングデータの使用が、(1つまたは複数の)分類器をバイアスし得る。各クラスラベルのほぼ等しい数の例があるように、「学習則」など、ルールが、トレーニングデータを平衡させるための試みとして利用され得る。トレーニングデータが、あるクラスの多数の例と別のクラスの少数の例とを含んでいる場合、分類関数のパラメータは、より多数の例をもつクラスのほうへバイアスされる方法で、よりしばしば更新される。極端に言うと、第1のクラスの100万個の例と第2のクラスの1つのみの例とを用いてバイナリ分類器をトレーニングしている場合、分類器は、単に常に第1のクラスを予測することによって極めてうまく機能する。別の例では、犬認識器がトレーニングされている。この例では、トレーニングデータは、合計1000個の例を含み、ここで、例のうちの990個は犬であり、例のうちの10個は猫である。分類器は、画像を犬として分類するために学習し得、これは、トレーニングセットに対して高精度で高い再現率を生じることになる。しかしながら、分類器が何も学ばなかった可能性が高い。

20

【0049】

[0059]一般に、クラス間のトレーニングデータの「平衡化」は、各クラスについてのトレーニング例の相対頻度(relative frequency)が、トレーニング中に使用されない新しい例に分類器を適用するときに遭遇すると予想される相対頻度に一致することを保証することによって対処される。しかしながら、この手法は、いくつかの欠点を有する。第1に、それは、将来のデータセット中のクラス例の相対頻度が知られていると仮定する。しかしながら、これは、決定することが常に容易であるとは限らない。第2に、トレーニングデータは、各クラスのあまりに多くの例またはあまりに少数の例を含んでいることがある。トレーニング例を平衡させるために、データは、捨てられるかまたは繰り返される。データを捨てることによって、いくつかのクラスについて有益なトレーニングデータが除外され得、これは、分類器がそのクラスに関連する入力変動を十分に表すのを妨げ得る。簡単な方法でデータを繰り返すことによって、データを段階に分けるためにより多くのディスクスペースが使用される。特に、目的が、データのすべてを使用することである場合、あらゆるクラスが、完全な平衡のために最小公倍数まで繰り返されるであろう。さらに、各例が2つまたはそれ以上のラベルについて正としてラベリングされ得るマルチラベルデータの場合、すべてのラベルにわたる平衡は、複雑なスケジューリング訓練になり、単純な繰り返しは十分でないことがある。

30

40

【0050】

[0060]本開示の態様は、機械学習モデルにおいてクラス間のトレーニングデータを平衡させることを対象とする。特に、入力段においてトレーニングデータを操作し、各クラスについての例の数を調節するのではなく、本開示の態様は、勾配段における調節を対象とする。

【0051】

[0061]後方への誤差伝播とも呼ばれるバックプロパゲーションが、コスト関数の勾配を計算するために利用され得る。特に、バックプロパゲーションは、誤差を0のより近くに低減するために重み値をどのように調節するかを決定することを含む。本開示の様々な態

50

様では、選択的バックプロパゲーションは、データセット中のクラス例頻度に基づいて勾配を調節するかまたは選択的に適用するための、何らかの所与のコスト関数への変更である。画像が入力され、勾配が、バックプロパゲーションを実行するために適用されようとしている後に、勾配は、各クラスについての例の頻度に基づいて調節され得る。

【 0 0 5 2 】

[0062]一態様では、調節は、トレーニングデータセット中の例の最小数 (

【 0 0 5 3 】

【 数 1 】

$$\min_{\mathbb{C}} N_c$$

【 0 0 5 4 】

)とトレーニングデータセット中のすべての例の数と (N_c 、たとえば、最も少数のメンバーをもつクラスの例の数と現在のクラスの例の数と)の比である、相対クラス頻度 f_c に関係する。(頻度ファクタ (frequency factor) ととも呼ばれる) 相対クラス頻度は、次のように表され得る。

【 0 0 5 5 】

【 数 2 】

$$f_c = \frac{\min_{\mathbb{C}} N_c}{N_c} \quad \mathbb{C} \ni \text{すべての概念} \quad (1)$$

【 0 0 5 6 】

[0063]例の最小数は、実際のまたは予想される数に基づき得る。さらに、トレーニングデータセット中のすべての例の数は、予想される数の例の実際の数に基づき得る。再び、犬認識器がトレーニングされている猫/犬の例を参照すると、犬の990個の例と猫の10個の例とがある。犬のための各クラスについての頻度ファクタは、 $10/990$ であり、ここで、10は例の最小数であり、990は対象のクラスについての例の数である。猫のための各クラスについてのファクタは、 $10/10$ である。調節ファクタ (たとえば、相対クラス頻度) は、例の最小数を有するクラスについて値「1」であり、すべての他のクラスについて1よりも小さくなり得る。

【 0 0 5 7 】

[0064]頻度ファクタが決定されると、バックプロパゲーション勾配が変更される。変更は、各クラスについて勾配をスケールリングすることを含み得る。スケールリングは、次のように表され得る。

【 0 0 5 8 】

【 数 3 】

$$\text{スケールリング: } \frac{dE_{\text{applied}}}{dx} = f_c \frac{dE}{dx} \quad (2)$$

【 0 0 5 9 】

[0065]スケールリングインプリメンテーションでは、勾配は、頻度ファクタ (たとえば、相対クラス頻度) によって乗算され得る。勾配は、特定のパラメータに関する誤差の導関数である。あるクラスの多くの例がある一例では、そのクラスの過剰学習を防ぐために、勾配の分数のみが毎回適用される。連続する、犬の10個の例がある犬/猫の例では、勾配の10分の1のみが適用される。目的は、モデルが猫よりも犬のはるかに多い例を参照したので、モデルがすべての画像を過剰学習し、犬としてラベリングするのを防ぐことである。スケールリングは、特定のクラスのすべての重み中のすべての勾配に等しく適用される。

【 0 0 6 0 】

[0066]変更は、画像からサンプリングするためのファクタを使用することをも含み得る。サンプリングは、次のように表され得る。

【 0 0 6 1 】

10

20

30

40

【数 4】

$$\text{サンプリング: } \frac{dE_{\text{applied}}}{dx} = \begin{cases} 0, & s=0 \text{ の場合} \\ \frac{dE}{dx}, & s=1 \text{ の場合} \end{cases} \quad (3)$$

【 0 0 6 2 】

[0067]ここで、勾配は、クラス例のサンプリングに基づいて選択的に適用される。一例では、サンプリングはランダムに適用される。スケーリングファクタの値は、サンプルがそれから引き出されるベルヌーイ分布の確率パラメータとして使用され得る。この分布からのサンプリングは0または1を生成し、1をサンプリングすることの確率が、第1の方法において説明されたスケーリングファクタに等しい。例の最小数をもつクラスの場合、サンプリングは、1を生成する。コインフリップ (coin flip) が1を生成するとき、そのクラスについての誤差勾配がバックプロパゲートされる。コインフリップが0を生成するとき、そのクラスについての勾配がバックプロパゲートされない場合が、事実上、0に設定される。言い換えれば、画像は、多くの例があるときに単に時々勾配を返送するために、勾配段においてサンプリングされる。最小数の例があるとき、勾配は毎回返送される。これは、入力を調節するのではなく勾配を調節することによって分類器がそれから学習している例の等化を与える。一態様では、画像を前方伝搬する前に、それは、クラスが、現在のエポックのためにその画像を使用するように設定されるかどうかを検査される。各エポックについて、セットは再シャッフルされ得る。

10

【 0 0 6 3 】

[0068]サンプリングは、個々ベース、エポックベース、またはトレーニングコーパスベースで適用され得る。上記で提示されたように、個々ベースでは、トレーニングエポック中に提示される他の画像に依存しない各画像のためのランダム結果がベルヌーイ分布から生成される。いくつかのエポックは、サンプリングのランダム性質により、各クラスについて、所望の数よりも多いまたは少ない例を参照し得る。

20

【 0 0 6 4 】

[0069]エポックベースの場合、スケールファクタは、すべてのクラス例から各クラスについてランダムに選択される。各エポック中に各クラスについて固定数の例が使用される。たとえば、10個の例が各クラスから選択され得る。それらの例のみが、特定のエポック中にバックプロパゲートされる。

30

【 0 0 6 5 】

[0070]トレーニングコーパスベースの場合、頻度ファクタは、分類器にまだ提示されていないものから各クラスについて各エポックのためにランダムに選択される。例は、交換なしにサンプリングされる。以下の例示的な例では、1000個の犬の例があり、各エポックでは、10個のサンプルがランダムに選択される。第1のエポックでは、10個の例が、合計1000個の例から選択される。次のエポックでは、前の10個の選択された例が除去され、10個の例が、残りの990個の例から選択される。これは、例のすべてが使い尽くされるまで続き、各エポック中に各クラスについて同数の例が使用され、トレーニングの過程にわたってすべての利用可能な例が使用されることを保証する。次回、データを巡回するとき、同じ順序が維持され得るか、または代替的に、異なる順序が使用され得る。別の構成では、例は、交換を用いてサンプリングされる。

40

【 0 0 6 6 】

[0071]多くの場合、トレーニングの開始の前にトレーニングコーパス全体が利用可能であり、fcファクタは、トレーニングセッションにわたって静的であり、トレーニングが始まる前に各クラスについて計算され得る。しかしながら、トレーニングが始まった後にクラスが追加されるか、またはトレーニング中にトレーニング例がアドホックに供給される場合、fcファクタは、時間とともに変化しているかまたはトレーニングの開始時に未知であることがある。この状況では、各例が提示された後に、各クラスについての例の数(Nc)のランニングカウントが、保たれ、更新され得る。fcファクタは、次いで、特定のクラス(c)についてNcの各更新の後にオンザフライで計算される。

50

【 0 0 6 7 】

[0072]別の態様では、各クラスについてネットワークの変化の量を等化するために、および各クラスが分類器によって比較的同様に推測される可能性があることを保証するために、クラスの相対頻度（たとえば、頻度ファクタ）が利用される。相対頻度クラスは、データセット中のクラスの一様分布を促進する。他のクラスよりもいくつかのクラスのより多くがあるという知られている予想がある場合、頻度ファクタは調節され得る。たとえば、実世界において犬よりも多くの猫がいることが知られているが、トレーニングデータが犬の1000個の例と猫の10個の例を含む場合、頻度ファクタは、実世界予想を考慮するように調節され得る。実世界において犬よりも猫を見る可能性が10倍高いと知られている場合、頻度ファクタは、猫についてファクタ10で乗算され、犬についてファクタ1で乗算され得る。本質的に、頻度ファクタ（ F_c ）は、実世界において存在するものの均一な予想をターゲットにするように学習段において操作され得る。頻度ファクタは、次のように調節され得る。

10

【 0 0 6 8 】

【 数 5 】

$$f_c = \frac{\min_c p(c)}{p(c)} \frac{\min_c N_c}{N_c}, \quad (4)$$

【 0 0 6 9 】

ここで、 $p(c)$ は、実世界（または「野生」）における特定のクラスを観測する予想される確率である。

20

【 0 0 7 0 】

[0073]図6は、機械学習モデルのためのクラス間のトレーニングデータを平衡させるための方法600を示す。ブロック602において、プロセスは、最も少数のメンバーをもつクラスの例の数と現在のクラスの例の数との比に基づいて、勾配を変更するためのファクタを決定する。最も少数のメンバーは、実際のまたは予想されるメンバーの数に基づき得る。同様に、現在のクラスの例の数は、例の実際のまたは予想される数に基づき得る。ブロック604において、プロセスは、決定されたファクタに基づいて、現在のクラスに関連する勾配を変更する。

【 0 0 7 1 】

30

[0074]図7は、機械学習モデルのためのクラス間のトレーニングデータを平衡させるための全体的方法700を示す。ブロック702において、トレーニングデータを評価する。ブロック704において、クラス中の例の頻度を決定する。ブロック706において、決定された頻度に基づいて勾配を更新する。更新は、ブロック710において、各クラスについて勾配にスケーリングファクタを適用することによって実行され得る。代替的に、更新は、ブロック708において、クラス例のサンプルに基づいて勾配を選択的に適用することによって実行され得る。選択的サンプリング更新は、ブロック712において個々ベースで、ブロック714においてエポックベースで、またはブロック716においてトレーニングコーパスベースで実行され得る。

【 0 0 7 2 】

40

[0075]図8は、本開示の態様による、トレーニングデータを平衡させるための方法800を示す。ブロック802において、プロセスは、モデルをトレーニングする間、バックプロパゲーションプロセスの勾配を変更する。変更は、最も少数のメンバーをもつクラスの例の数と現在のクラスの例の数との比に基づく。

【 0 0 7 3 】

[0076]いくつかの態様では、方法600、700、および800は、SOC100（図1）またはシステム200（図2）によって実行され得る。すなわち、方法1100および1200の要素の各々は、たとえば、限定はしないが、SOC100またはシステム200または1つまたは複数のプロセッサ（たとえば、CPU102およびローカル処理ユニット202）および/あるいは本明細書中に含まれる他の構成要素によって実行され得

50

る。いくつかの態様では、方法600および700は、SOC420(図4)または1つまたは複数のプロセッサ(たとえば、CPU422)および/あるいは本明細書に含まれる他の構成要素によって実行され得る。

【0074】

[0077]上記で説明された方法の様々な動作は、対応する機能を実施することが可能な任意の好適な手段によって実施され得る。それらの手段は、限定はしないが、回路、特定用途向け集積回路(ASIC)、またはプロセッサを含む、様々な(1つまたは複数の)ハードウェアおよび/またはソフトウェア構成要素および/またはモジュールを含み得る。概して、図に示されている動作がある場合、それらの動作は、同様の番号をもつ対応するカウンターパートのミーンズプラスファンクション構成要素を有し得る。

10

【0075】

[0078]本明細書で使用する「決定すること」という用語は、多種多様なアクションを包含する。たとえば、「決定すること」は、計算すること(calculating)、計算すること(computing)、処理すること、導出すること、調査すること、ルックアップすること(たとえば、テーブル、データベースまたは別のデータ構造においてルックアップすること)、確認することなどを含み得る。さらに、「決定すること」は、受信すること(たとえば、情報を受信すること)、アクセスすること(たとえば、メモリ中のデータにアクセスすること)などを含み得る。さらに、「決定すること」は、解決すること、選択すること、選定すること、確立することなどを含み得る。

【0076】

20

[0079]本明細書で使用する、項目のリスト「のうちの少なくとも1つ」を指す句は、単一のメンバーを含む、それらの項目の任意の組合せを指す。一例として、「a、b、またはcのうちの少なくとも1つ」は、a、b、c、a-b、a-c、b-c、およびa-b-cを包含するものとする。

【0077】

[0080]本開示に関連して説明された様々な例示的な論理ブロック、モジュールおよび回路は、汎用プロセッサ、デジタル信号プロセッサ(DSP)、特定用途向け集積回路(ASIC)、フィールドプログラマブルゲートアレイ信号(FPGA)または他のプログラマブル論理デバイス(PLD)、個別ゲートまたはトランジスタ論理、個別ハードウェア構成要素、あるいは本明細書で説明された機能を実施するように設計されたそれらの任意の組合せを用いてインプリメントまたは実施され得る。汎用プロセッサはマイクロプロセッサであり得るが、代替として、プロセッサは、任意の市販のプロセッサ、コントローラ、マイクロコントローラ、または状態機械であり得る。プロセッサはまた、コンピューティングデバイスの組合せ、たとえば、DSPとマイクロプロセッサとの組合せ、複数のマイクロプロセッサ、DSPコアと連携する1つまたは複数のマイクロプロセッサ、あるいは任意の他のそのような構成としてインプリメントされ得る。

30

【0078】

[0081]本開示に関連して説明された方法またはアルゴリズムのステップは、ハードウェアで直接実施されるか、プロセッサによって実行されるソフトウェアモジュールで実施されるか、またはその2つの組合せで実施され得る。ソフトウェアモジュールは、当技術分野で知られている任意の形態の記憶媒体中に常駐し得る。使用され得る記憶媒体のいくつかの例としては、ランダムアクセスメモリ(RAM)、読取り専用メモリ(ROM)、フラッシュメモリ、消去可能プログラマブル読取り専用メモリ(EPROM)、電気消去可能プログラマブル読取り専用メモリ(EEPROM(登録商標))、レジスタ、ハードディスク、リムーバブルディスク、CD-ROMなどがある。ソフトウェアモジュールは、単一の命令、または多数の命令を備え得、いくつかの異なるコードセグメント上で、異なるプログラム間で、および複数の記憶媒体にわたって分散され得る。記憶媒体は、プロセッサがその記憶媒体から情報を読み取ることができ、その記憶媒体に情報を書き込むことができるように、プロセッサに結合され得る。代替として、記憶媒体はプロセッサと一体であり得る。

40

50

【 0 0 7 9 】

[0082]本明細書で開示された方法は、説明された方法を達成するための1つまたは複数のステップまたはアクションを備える。本方法のステップおよび/またはアクションは、特許請求の範囲から逸脱することなく、互いに交換され得る。言い換えれば、ステップまたはアクションの特定の順序が指定されない限り、特定のステップおよび/またはアクションの順序および/または使用は特許請求の範囲から逸脱することなく変更され得る。

【 0 0 8 0 】

[0083]説明された機能は、ハードウェア、ソフトウェア、ファームウェア、またはそれらの任意の組合せでインプリメントされ得る。ハードウェアでインプリメントされる場合、例示的なハードウェア構成はデバイス中に処理システムを備え得る。処理システムは、バスアーキテクチャを用いてインプリメントされ得る。バスは、処理システムの特定の適用例および全体的な設計制約に応じて、任意の数の相互接続バスおよびブリッジを含み得る。バスは、プロセッサと、機械可読媒体と、バスインターフェースとを含む様々な回路を互いにリンクし得る。バスインターフェースは、ネットワークアダプタを、特に、バスを介して処理システムに接続するために使用され得る。ネットワークアダプタは、信号処理機能をインプリメントするために使用され得る。いくつかの態様では、ユーザインターフェース（たとえば、キーパッド、ディスプレイ、マウス、ジョイスティックなど）もバスに接続され得る。バスはまた、タイミングソース、周辺機器、電圧調整器、電力管理回路など、様々な他の回路をリンクし得るが、それらは当技術分野でよく知られており、したがってこれ以上説明されない。

【 0 0 8 1 】

[0084]プロセッサは、機械可読媒体に記憶されたソフトウェアの実行を含む、バスおよび一般的な処理を管理することを担当し得る。プロセッサは、1つまたは複数の汎用および/または専用プロセッサを用いてインプリメントされ得る。例としては、マイクロプロセッサ、マイクロコントローラ、DSPプロセッサ、およびソフトウェアを実行することができる他の回路がある。ソフトウェアは、ソフトウェア、ファームウェア、ミドルウェア、マイクロコード、ハードウェア記述言語などの名称にかかわらず、命令、データ、またはそれらの任意の組合せを意味すると広く解釈されたい。機械可読媒体は、例として、ランダムアクセスメモリ(RAM)、フラッシュメモリ、読取り専用メモリ(ROM)、プログラマブル読取り専用メモリ(PROM)、消去可能プログラマブル読取り専用メモリ(EPROM)、電気消去可能プログラマブル読取り専用メモリ(EEPROM)、レジスタ、磁気ディスク、光ディスク、ハードドライブ、または他の好適な記憶媒体、あるいはそれらの任意の組合せを含み得る。機械可読媒体はコンピュータプログラム製品において実施され得る。コンピュータプログラム製品はパッケージング材料を備え得る。

【 0 0 8 2 】

[0085]ハードウェアインプリメンテーションでは、機械可読媒体は、プロセッサとは別個の処理システムの一部であり得る。しかしながら、当業者なら容易に理解するように、機械可読媒体またはその任意の部分は処理システムの外部にあり得る。例として、機械可読媒体は、すべてバスインターフェースを介してプロセッサによってアクセスされ得る、伝送線路、データによって変調された搬送波、および/またはデバイスとは別個のコンピュータ製品を含み得る。代替的に、または追加として、機械可読媒体またはその任意の部分は、キャッシュおよび/または汎用レジスタファイルがそうであり得るように、プロセッサに統合され得る。局所構成要素など、説明された様々な構成要素は、特定のロケーションを有するものとして説明され得るが、それらはまた、分散コンピューティングシステムの一部として構成されているいくつかの構成要素など、様々な方法で構成され得る。

【 0 0 8 3 】

[0086]処理システムは、すべて外部バスアーキテクチャを介して他のサポート回路と互いにリンクされる、プロセッサ機能を提供する1つまたは複数のマイクロプロセッサと、機械可読媒体の少なくとも一部を提供する外部メモリとをもつ汎用処理システムとして構成され得る。代替的に、処理システムは、本明細書で説明されたニューロンモデルとニュー

10

20

30

40

50

ーラルシステムのモデルとをインプリメントするための1つまたは複数の神経形態学的プロセッサを備え得る。別の代替として、処理システムは、プロセッサをもつ特定用途向け集積回路(A S I C)と、バスインターフェースと、ユーザインターフェースと、サポート回路と、単一のチップに統合された機械可読媒体の少なくとも一部分とを用いて、あるいは1つまたは複数のフィールドプログラマブルゲートアレイ(F P G A)、プログラマブル論理デバイス(P L D)、コントローラ、状態機械、ゲート論理、個別ハードウェア構成要素、もしくは他の好適な回路、または本開示全体にわたって説明された様々な機能を実施することができる回路の任意の組合せを用いて、インプリメントされ得る。当業者は、特定の適用例と、全体的なシステムに課される全体的な設計制約とに応じて、どのようにしたら処理システムについて説明された機能を最も良くインプリメントし得るかを理解されよう。

10

【0084】

[0087]機械可読媒体はいくつかのソフトウェアモジュールを備え得る。ソフトウェアモジュールは、プロセッサによって実行されたときに、処理システムに様々な機能を実施させる命令を含む。ソフトウェアモジュールは、送信モジュールと受信モジュールとを含み得る。各ソフトウェアモジュールは、単一の記憶デバイス中に常駐するか、または複数の記憶デバイスにわたって分散され得る。例として、トリガイイベントが発生したとき、ソフトウェアモジュールがハードドライブからR A Mにロードされ得る。ソフトウェアモジュールの実行中、プロセッサは、アクセス速度を高めるために、命令のいくつかをキャッシュにロードし得る。次いで、1つまたは複数のキャッシュラインが、プロセッサによる実行のために汎用レジスタファイルにロードされ得る。以下でソフトウェアモジュールの機能に言及する場合、そのような機能は、そのソフトウェアモジュールからの命令を実行したときにプロセッサによってインプリメントされることが理解されよう。さらに、本開示の態様が、そのような態様をインプリメントするプロセッサ、コンピュータ、機械、または他のシステムの機能に改善を生じることを諒解されたい。

20

【0085】

[0088]ソフトウェアでインプリメントされる場合、機能は、1つまたは複数の命令またはコードとしてコンピュータ可読媒体上に記憶されるか、あるいはコンピュータ可読媒体を介して送信され得る。コンピュータ可読媒体は、ある場所から別の場所へのコンピュータプログラムの転送を可能にする任意の媒体を含む、コンピュータ記憶媒体と通信媒体の両方を含む。記憶媒体は、コンピュータによってアクセスされ得る任意の利用可能な媒体であり得る。限定ではなく例として、そのようなコンピュータ可読媒体は、R A M、R O M、E E P R O M、C D - R O Mまたは他の光ディスクストレージ、磁気ディスクストレージまたは他の磁気ストレージデバイス、あるいは命令またはデータ構造の形態の所望のプログラムコードを搬送または記憶するために使用され得、コンピュータによってアクセスされ得る、任意の他の媒体を備えることができる。さらに、いかなる接続もコンピュータ可読媒体と適切に呼ばれる。たとえば、ソフトウェアが、同軸ケーブル、光ファイバーケーブル、ツイストペア、デジタル加入者回線(D S L)、または赤外線(I R)、無線、およびマイクロ波などのワイヤレス技術を使用して、ウェブサイト、サーバ、または他のリモートソースから送信される場合、同軸ケーブル、光ファイバーケーブル、ツイストペア、D S L、または赤外線、無線、およびマイクロ波などのワイヤレス技術は、媒体の定義に含まれる。本明細書で使用されるディスク(disk)およびディスク(disc)は、コンパクトディスク(disc)(C D)、レーザーディスク(登録商標)(disc)、光ディスク(disc)、デジタル多用途ディスク(disc)(D V D)、フロッピー(登録商標)ディスク(disk)、およびB l u - r a y(登録商標)ディスク(disc)を含み、ディスク(disk)は、通常、データを磁氣的に再生し、ディスク(disc)は、データをレーザーで光学的に再生する。したがって、いくつかの態様では、コンピュータ可読媒体は非一時的コンピュータ可読媒体(たとえば、有形媒体)を備え得る。さらに、他の態様では、コンピュータ可読媒体は一時的コンピュータ可読媒体(たとえば、信号)を備え得る。上記の組合せもコンピュータ可読媒体の範囲内に含まれるべきである。

30

40

50

【 0 0 8 6 】

[0089]したがって、いくつかの態様は、本明細書で提示された動作を実施するためのコンピュータプログラム製品を備え得る。たとえば、そのようなコンピュータプログラム製品は、本明細書で説明された動作を実行するために1つまたは複数のプロセッサによって実行可能である命令をその上に記憶した（および／または符号化した）コンピュータ可読媒体を備え得る。いくつかの態様では、コンピュータプログラム製品はパッケージング材料を含み得る。

【 0 0 8 7 】

[0090]さらに、本明細書で説明された方法および技法を実行するためのモジュールおよび／または他の適切な手段は、適用可能な場合にユーザ端末および／または基地局によってダウンロードされ、および／または他の方法で取得され得ることを諒解されたい。たとえば、そのようなデバイスは、本明細書で説明された方法を実行するための手段の転送を可能にするためにサーバに結合され得る。代替的に、本明細書で説明された様々な方法は、ユーザ端末および／または基地局が記憶手段（たとえば、RAM、ROM、コンパクトディスク（CD）またはフロッピーディスクなどの物理記憶媒体など）をデバイスに結合するかまたは与えると様々な方法を得ることができるよう、記憶手段によって提供され得る。その上、本明細書で説明された方法および技法をデバイスに提供するための任意の他の好適な技法が利用され得る。

【 0 0 8 8 】

[0091]特許請求の範囲は、上記で示された厳密な構成および構成要素に限定されないことを理解されたい。上記で説明された方法および装置の構成、動作および詳細において、特許請求の範囲から逸脱することなく、様々な改変、変更および変形が行われ得る。

以下に本願の出願当初の特許請求の範囲に記載された発明を付記する。

【 C 1 】

機械学習モデルのためのクラス間のトレーニングデータの平衡を変更する方法であって

、

最も少数のメンバーをもつクラスの例の数と現在のクラスの例の数との比に少なくとも部分的に基づいて、前記モデルをトレーニングする間、バックプロパゲーションプロセスの勾配を変更することを備える、方法。

【 C 2 】

前記変更することが、前記勾配をスケールリングすることを備える、C 1 に記載の方法。

【 C 3 】

前記変更することが、前記クラス例のサンプリングに少なくとも部分的に基づいて前記勾配を選択的に適用することを備える、C 1 に記載の方法。

【 C 4 】

前記クラスの前記サンプリングが、各トレーニングエポックから固定数の例を選択することによって行われる、C 3 に記載の方法。

【 C 5 】

前記サンプリングが、トレーニングエポック中の例の交換なしに行われる、C 1 に記載の方法。

【 C 6 】

機械学習モデルのためのクラス間のトレーニングデータの平衡を変更するための装置であって、

最も少数のメンバーをもつクラスの例の数と現在のクラスの例の数との比に少なくとも部分的に基づいて、勾配を変更するためのファクタを決定するための手段と、

前記決定されたファクタに基づいて、前記現在のクラスに関連する前記勾配を変更するための手段とを備える、装置。

【 C 7 】

前記変更手段が、前記勾配をスケールリングするための手段を備える、C 6 に記載の装置

。

10

20

30

40

50

[C 8]

前記変更手段が、前記クラス例のサンプリングに少なくとも部分的に基づいて前記勾配を選択的に適用するための手段を備える、C 6 に記載の装置。

[C 9]

前記クラスの前記サンプリングが、各トレーニングエポックから固定数の例を選択することによって行われる、C 8 に記載の装置。

[C 10]

前記サンプリングが、トレーニングエポック中の例の交換なしに行われる、C 6 に記載の装置。

[C 11]

機械学習モデルのためのクラス間のトレーニングデータの平衡を変更するための装置であって、

メモリと、

前記メモリに結合された少なくとも1つのプロセッサと、前記少なくとも1つのプロセッサが、最も少数のメンバーをもつクラスの例の数と現在のクラスの例の数との比に少なくとも部分的に基づいて、前記モデルをトレーニングする間、バックプロパゲーションプロセスの勾配を変更するように構成された、を備える、装置。

[C 12]

前記少なくとも1つのプロセッサが、前記勾配をスケールリングすることによって変更するように構成された、C 11 に記載の装置。

[C 13]

前記少なくとも1つのプロセッサが、前記クラス例のサンプリングに少なくとも部分的に基づいて、前記勾配を選択的に適用することによって変更するように構成された、C 11 に記載の装置。

[C 14]

前記クラスの前記サンプリングが、各トレーニングエポックから固定数の例を選択することによって行われる、C 13 に記載の装置。

[C 15]

前記サンプリングが、トレーニングエポック中の例の交換なしに行われる、C 11 に記載の装置。

[C 16]

機械学習モデルのためのクラス間のトレーニングデータの平衡を変更するための非一時的コンピュータ可読媒体であって、前記非一時的コンピュータ可読媒体がそれに記録されたプログラムコードを有し、前記プログラムコードが、

最も少数のメンバーをもつクラスの例の数と現在のクラスの例の数との比に少なくとも部分的に基づいて、前記モデルをトレーニングする間、バックプロパゲーションプロセスの勾配を変更するためのプログラムコードを備える、非一時的コンピュータ可読媒体。

[C 17]

変更するための前記プログラムコードが、前記勾配をスケールリングするためのプログラムコードを備える、C 16 に記載の非一時的コンピュータ可読媒体。

[C 18]

変更するための前記プログラムコードが、前記クラス例のサンプリングに少なくとも部分的に基づいて前記勾配を選択的に適用するためのプログラムコードを備える、C 16 に記載の非一時的コンピュータ可読媒体。

[C 19]

前記クラスの前記サンプリングが、各トレーニングエポックから固定数の例を選択することによって行われる、C 18 に記載の非一時的コンピュータ可読媒体。

[C 20]

前記サンプリングが、トレーニングエポック中の例の交換なしに行われる、C 16 に記載の非一時的コンピュータ可読媒体。

10

20

30

40

50

【図 1】

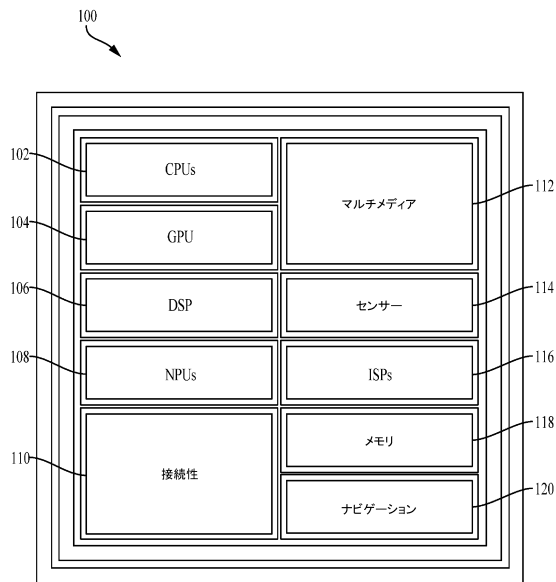


FIG. 1

【図 2】

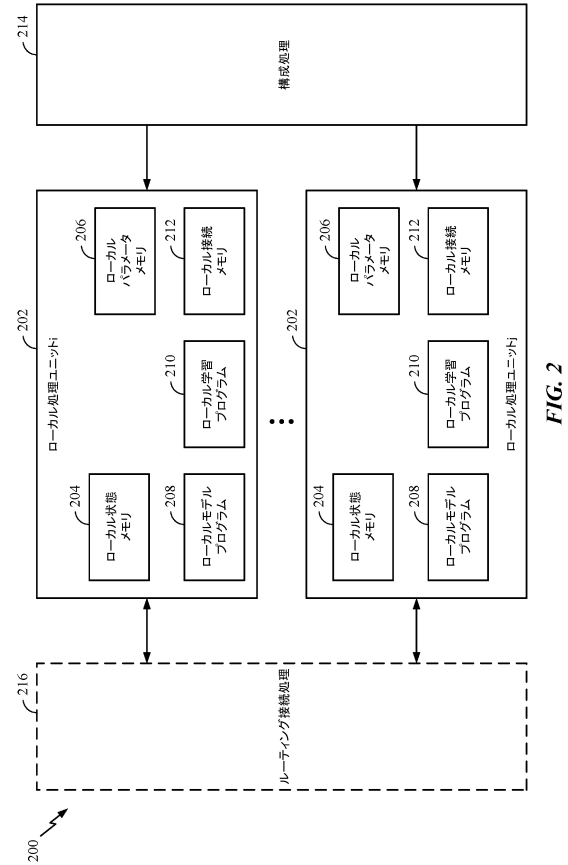


FIG. 2

【図 3 A】

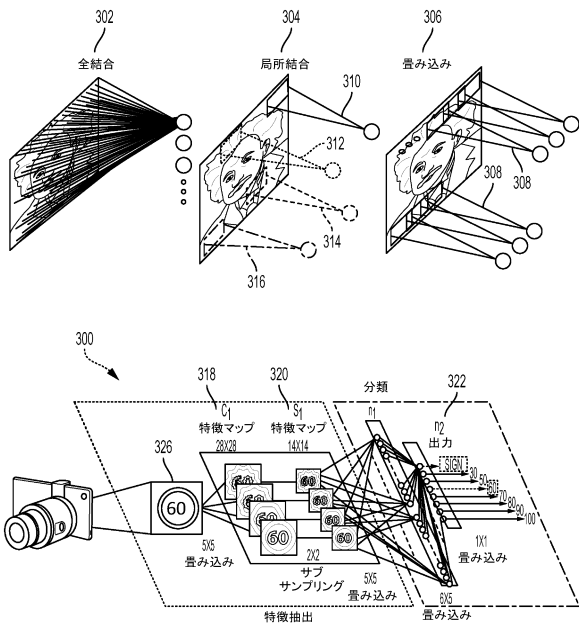


FIG. 3A

【図 3 B】

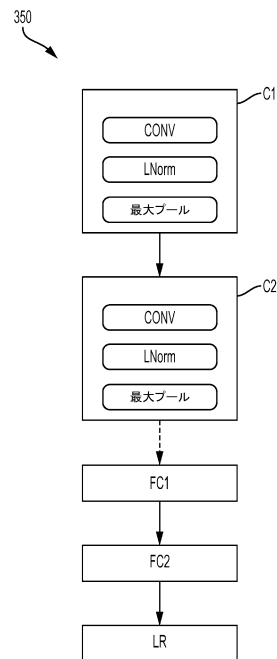


FIG. 3B

【図4】

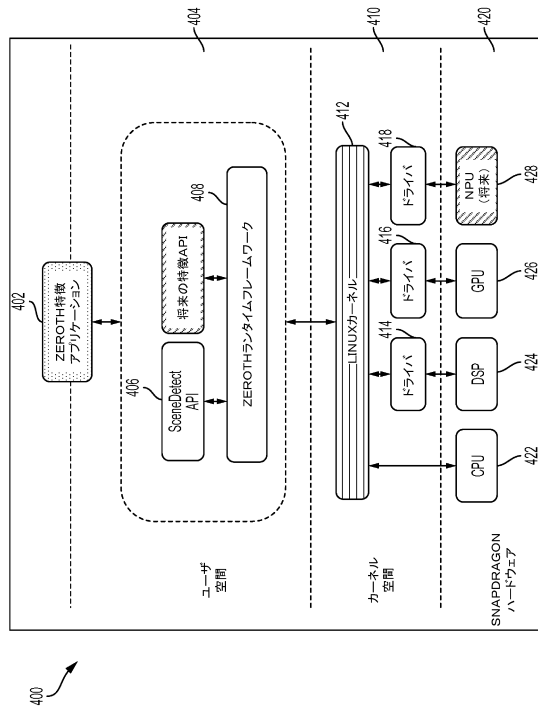


FIG. 4

【図5】

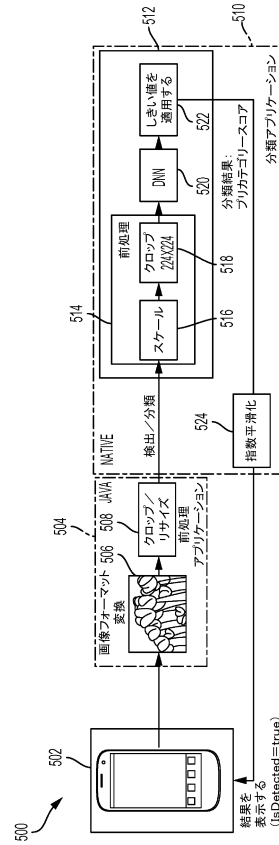


FIG. 5

【図6】

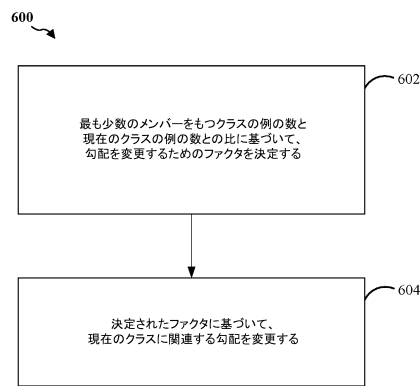


FIG. 6

【図7】

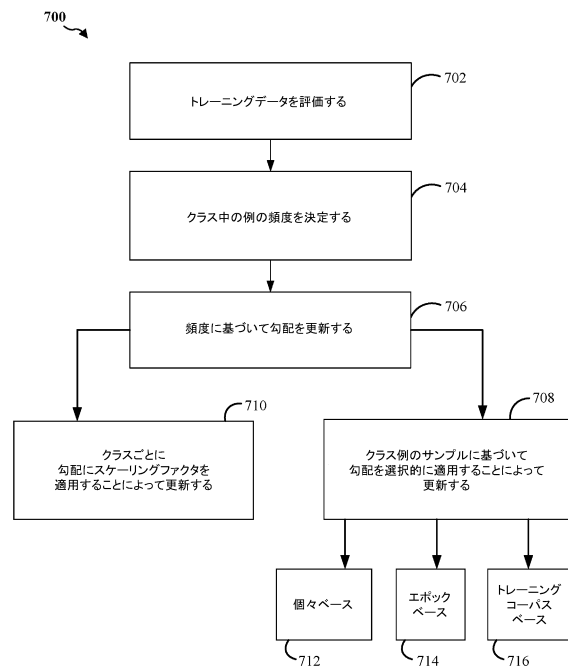


FIG. 7

【図 8】

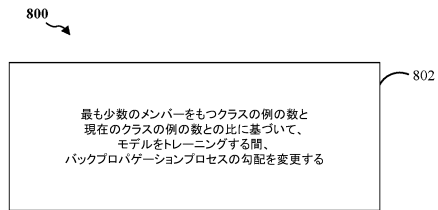


FIG. 8

フロントページの続き

(74)代理人 100184332

弁理士 中丸 慶洋

(72)発明者 トワル、レーガン・ブライス

アメリカ合衆国、カリフォルニア州 9 2 1 2 1 - 1 7 1 4、サン・ディエゴ、モアハウス・ドライブ 5 7 7 5

(72)発明者 ジュリアン、デイビッド・ジョナサン

アメリカ合衆国、カリフォルニア州 9 2 1 0 1、サン・ディエゴ、ファースト・アベニュー 5 1 0、ユニット 5 0 4

審査官 杉浦 孝光

(56)参考文献 特開 2 0 0 9 - 1 2 2 8 5 1 (J P , A)

OH, Sang-Hoon , Error back-propagation algorithm for classification of imbalanced data , Neurocomputing , 2 0 1 1 年 , vol.74, Issue.6, pp.1058-1061 , U R L , <https://www.sciencedirect.com/science/article/abs/pii/S0925231210005084>

(58)調査した分野(Int.Cl. , D B 名)

G 0 6 N 3 / 0 0 - 9 9 / 0 0