

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第6862426号  
(P6862426)

(45) 発行日 令和3年4月21日(2021.4.21)

(24) 登録日 令和3年4月2日(2021.4.2)

(51) Int. Cl. F I  
**GO 6 N 20/00 (2019.01)** GO 6 N 20/00  
**GO 6 N 3/08 (2006.01)** GO 6 N 3/08

請求項の数 15 (全 31 頁)

(21) 出願番号	特願2018-510504 (P2018-510504)	(73) 特許権者	595020643
(86) (22) 出願日	平成28年8月11日 (2016. 8. 11)		クアアルコム・インコーポレイテッド
(65) 公表番号	特表2018-529159 (P2018-529159A)		QUALCOMM INCORPORATED
(43) 公表日	平成30年10月4日 (2018. 10. 4)		ED
(86) 国際出願番号	PCT/US2016/046576		アメリカ合衆国、カリフォルニア州 92
(87) 国際公開番号	W02017/034820		121-1714、サン・ディエゴ、モア
(87) 国際公開日	平成29年3月2日 (2017. 3. 2)		ハウス・ドライブ 5775
審査請求日	令和1年7月19日 (2019. 7. 19)	(74) 代理人	100108855
(31) 優先権主張番号	62/209, 859		弁理士 蔵田 昌俊
(32) 優先日	平成27年8月25日 (2015. 8. 25)	(74) 代理人	100109830
(33) 優先権主張国・地域又は機関	米国 (US)		弁理士 福原 淑弘
(31) 優先権主張番号	14/863, 410	(74) 代理人	100158805
(32) 優先日	平成27年9月23日 (2015. 9. 23)		弁理士 井関 守三
(33) 優先権主張国・地域又は機関	米国 (US)	(74) 代理人	100112807
			弁理士 岡田 貴志

最終頁に続く

(54) 【発明の名称】 トレーニングされた機械学習モデルのパフォーマンスを改善するための方法

(57) 【特許請求の範囲】

【請求項 1】

トレーニングされた機械学習モデルのパフォーマンスを改善するためのコンピュータによりインプリメントされる方法であって、

第1の目的関数を有する第1の分類器に、第2の目的関数を有する第2の分類器を追加すること、前記第1の目的関数は、前記第2の目的関数とは異なり、前記第1の目的関数は、微分可能であり、前記第2の目的関数は、微分不可能であり、前記第2の目的関数は、前記第1の分類器の誤差を直接的に低減させるために使用される、

を備える、コンピュータによりインプリメントされる方法。

【請求項 2】

前記第2の目的関数は、前記第1の分類器と前記第2の分類器との誤差間の差の関数である、請求項1に記載のコンピュータによりインプリメントされる方法。

【請求項 3】

より高い複雑度のモデルからの確率の混合に少なくとも部分的に基づいて、前記第2の目的関数を決定することをさらに備える、請求項1に記載のコンピュータによりインプリメントされる方法。

【請求項 4】

前記第1の分類器を再トレーニングすることなく、前記第2の分類器を追加することをさらに備える、請求項1に記載のコンピュータによりインプリメントされる方法。

【請求項 5】

前記第 1 の分類器の外部に前記第 2 の分類器を追加することをさらに備える、請求項 1 に記載のコンピュータによりインプリメントされる方法。

【請求項 6】

アイデンティティ値に、前記第 1 の分類器によってトレーニングされたモデルによって生成される特徴に対する重みを割り当てることをさらに備える、請求項 1 に記載のコンピュータによりインプリメントされる方法。

【請求項 7】

ゼロに、高い複雑度のモデルの確率ベクトルによって生成される特徴に対する重みを割り当てることをさらに備える、請求項 6 に記載のコンピュータによりインプリメントされる方法。

10

【請求項 8】

前記第 2 の分類器の確率ベクトルによって生成される特徴に対する重みを割り当てることをさらに備える、請求項 1 に記載のコンピュータによりインプリメントされる方法。

【請求項 9】

ゼロに、高い複雑度のモデルの確率ベクトルによって生成される特徴に対する重みを割り当てることをさらに備える、請求項 1 に記載のコンピュータによりインプリメントされる方法。

【請求項 10】

前記第 2 の分類器の確率ベクトルによって生成される特徴に対する重みを割り当てることをさらに備える、請求項 9 に記載のコンピュータによりインプリメントされる方法。

20

【請求項 11】

固定された温度 T によって、より高い複雑度のモデルによって生成される確率ベクトルをスケールリングすることをさらに備える、請求項 1 に記載のコンピュータによりインプリメントされる方法。

【請求項 12】

トレーニングされた機械学習モデルのパフォーマンスを改善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

を備え、前記少なくとも 1 つのプロセッサは、

第 1 の目的関数を有する第 1 の分類器に、第 2 の目的関数を有する第 2 の分類器を追加すること、前記第 1 の目的関数は、微分可能であり、前記第 2 の目的関数は、微分不可能であり、前記第 1 の目的関数は、前記第 2 の目的関数とは異なり、前記第 2 の目的関数は、前記第 1 の分類器の誤差を直接的に低減させるために使用される、

30

を行うように構成される、装置。

【請求項 13】

i) 前記第 2 の目的関数は、前記第 1 の分類器と前記第 2 の分類器との誤差間の差の関数である、

ii) 前記少なくとも 1 つのプロセッサは、より高い複雑度のモデルからの確率の混合に少なくとも部分的に基づいて、前記第 2 の目的関数を決定するようにさらに構成される、

40

iii) 前記少なくとも 1 つのプロセッサは、前記第 1 の分類器を再トレーニングすることなく、前記第 2 の分類器を追加するようにさらに構成される、

iv) 前記少なくとも 1 つのプロセッサは、前記第 1 の分類器の外部に前記第 2 の分類器を追加するようにさらに構成される、

v) 前記少なくとも 1 つのプロセッサは、アイデンティティ値に、前記第 1 の分類器によってトレーニングされたモデルによって生成される特徴に対する重みを割り当てるようにさらに構成される、

vi) 前記少なくとも 1 つのプロセッサは、前記第 2 の分類器の確率ベクトルによって生成される特徴に対する重みを割り当てるようにさらに構成される、

vii) 前記少なくとも 1 つのプロセッサは、ゼロに、高い複雑度のモデルの確率ベクトル

50

ルによって生成される特徴に対する重みを割り当てるようにさらに構成される、

v i i i ) 前記少なくとも1つのプロセッサは、固定された温度Tによって、より高い複雑度のモデルによって生成される確率ベクトルをスケーリングするようにさらに構成される、

のうちの1つを備える、請求項12に記載の装置。

【請求項14】

トレーニングされた機械学習モデルのパフォーマンスを改善するための装置であって、第1の目的関数を有する第1の分類器に、第2の目的関数を有する第2の分類器を追加するための手段と、前記第1の目的関数は、前記第2の目的関数とは異なり、前記第1の目的関数は、微分可能であり、前記第2の目的関数は、微分不可能であり、前記第2の目的関数は、前記第1の分類器の誤差を直接的に低減させるために使用される、

10

前記トレーニングされた機械学習モデルを介して受信される入力に少なくとも部分的に基づいて、前記第2の分類器から特徴ベクトルを出力するための手段と

を備える装置。

【請求項15】

トレーニングされた機械学習モデルのパフォーマンスを改善するためのプログラムコードをその上に符号化された非一時的なコンピュータ可読媒体であって、前記プログラムコードは、プロセッサによって実行され、第1の目的関数を有する第1の分類器に、第2の目的関数を有する第2の分類器を追加するためのプログラムコードを備え、前記第1の目的関数は、前記第2の目的関数とは異なり、前記第1の目的関数は、微分可能であり、前記第2の目的関数は、微分不可能であり、前記第2の目的関数は、前記第1の分類器の誤差を直接的に低減させるために使用される、非一時的なコンピュータ可読媒体。

20

【発明の詳細な説明】

【関連出願への相互参照】

【0001】

[0001]本願は、2015年8月25日に出願された「トレーニングされた機械学習モデルのパフォーマンスを改善するための方法(METHOD FOR IMPROVING PERFORMANCE OF A TRAINED MACHINE LEARNING MODEL)」と題する米国仮特許出願第62/209,859号の利益を主張し、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。

【技術分野】

30

【0002】

[0002]本開示のある特定の態様は、一般に機械学習に関し、さらに詳細には、トレーニングされた機械学習モデルのパフォーマンスを改善するシステムおよび方法に関する。

【背景技術】

【0003】

[0003]相互結合された人工ニューロン(例えば、ニューロンモデル)のグループを備え得る、人工ニューラルネットワークのような機械学習モデルは、計算デバイスであるか、または計算デバイスによって実行される方法を表す。

【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(feed-forward)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かつ入力空間を集合的にタイリングする(collectively tile)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 networks)は、多数のアプリケーションを有する。特に、CNNは、パターン認識および分類(classification)の分野において幅広く使用されてきた。

40

【0005】

[0005]ディープピリフネットワーク(deep belief networks)およびディープ畳み込みネットワークのような、ディープラーニングアーキテクチャ(deep learning architectures)は、第1の層のニューロンの出力が第2の層のニューロンへの入力となり、第2の層のニューロンの出力が第3の層のニューロンへの入力となるというような、層状の二

50

ューラルネットワークアーキテクチャ (layered neural networks architectures) である。ディープニューラルネットワークは、特徴の階層 (hierarchy of features) 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、オブジェクト認識アプリケーション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と同様に、これらのディープラーニングアーキテクチャにおける計算は、処理ノードの集団 (population) にわたって分散され得、これは、1つまたは複数の計算チェーンにおいて構成され得る。これらの多層アーキテクチャは、一度に1層ずつ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 (fine-tuned) され得る。

【0006】

[0006]他のモデルもまた、オブジェクト認識のために利用可能である。例えば、サポートベクターマシン (SVM: support vector machines) は、分類に適用されることができ、学習ツールである。サポートベクターマシンは、データをカテゴリ化する (categorizes) 分離超平面 (例えば、決定境界) を含む。超平面は、教師あり学習 (supervised learning) 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大させる。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 (greatest minimum distance) を有すべきである。

10

【0007】

[0007]これらの解決策は、多数の (a number of) 分類ベンチマークに対して優れた結果を達成するが、それらの計算複雑度は、極めて高くなり得る。加えて、モデルのトレーニングは、困難であり得る。

20

【発明の概要】

【0008】

[0008]本開示のある態様では、トレーニングされた機械学習モデルのパフォーマンスを改善するための方法が提示される。方法は、第1の目的関数を有する第1の分類器に、第2の目的関数を有する第2の分類器を追加することを備える。第2の目的関数は、第1の分類器の誤差 (errors) を直接的に低減させるために使用される。

【0009】

[0009]別の態様では、トレーニングされた機械学習モデルのパフォーマンスを改善するための装置が提示される。装置は、メモリと、メモリに結合された少なくとも1つのプロセッサを含む。(1つまたは複数の) プロセッサは、第1の目的関数を有する第1の分類器に、第2の目的関数を有する第2の分類器を追加するように構成される。第2の目的関数は、第1の分類器の誤差を直接的に低減させるために使用される。

30

【0010】

[0010]さらに別の態様では、トレーニングされた機械学習モデルのパフォーマンスを改善するための装置が提示される。装置は、第1の目的関数を有する第1の分類器に、第2の目的関数を有する第2の分類器を追加するための手段を含む。第2の目的関数は、第1の分類器の誤差を直接的に低減させるために使用される。装置は、トレーニングされた機械学習モデルを介して受信される入力に基づいて、第2の分類器から特徴ベクトルを出力するための手段をさらに含む。

【0011】

40

[0011]さらになお別の態様では、非一時的なコンピュータ可読媒体が提示される。非一時的なコンピュータ可読媒体は、トレーニングされた学習機械モデルのパフォーマンスを改善するためのプログラムコードをその上に符号化されている。プログラムコードは、プロセッサによって実行され、第1の目的関数を有する第1の分類器に、第2の目的関数を有する第2の分類器を追加するためのプログラムコードを含む。第2の目的関数は、第1の分類器の誤差を直接的に低減させるために使用される。

【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。本開示が、本開示と同じ目的を実行するための他の構造を修正 (modifying) または設計するための基礎として容易に利用され得ることが、当業者によって理解されるべきである。また、そのような同等

50

の構造が、添付された特許請求の範囲に示される本開示の教示から逸脱しないことも、当業者によって認識されるべきである。さらなる目的および利点と共に、その構成および動作の方法の両方について、本開示の特徴であると考えられる新規の特徴は、添付の図面に関連して考慮されるとき、以下の説明からより良く理解されるであろう。しかしながら、図面の各々は、例示および説明のみの目的で提供されており、本開示の限定の定義として意図されるものではないことが、明確に理解されるべきである。

【図面の簡単な説明】

【0013】

[0013]本開示の特徴、性質、および利点は、同様の参照符号が全体を通して同様のものを指す図面と共に考慮されるとき、以下に示される詳細な説明からより明らかになるであろう。

10

【図1】[0014]図1は、本開示のある特定の態様による、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となるインプリメンテーションを例示する。

【図2】[0015]図2は、本開示の態様による、システムの例となるインプリメンテーションを例示する。

【図3A】[0016]図3Aは、本開示の態様による、ニューラルネットワークを例示する図である。

【図3B】[0017]図3Bは、本開示の態様による、例示的なディープ畳み込みネットワーク(DCN)を例示するブロック図である。

20

【図4】[0018]図4は、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例示するブロック図である。

【図5】[0019]図5は、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(run-time operation)を例示するブロック図である。

【図6A】[0020]図6Aは、本開示の態様による、機械学習モデルのパフォーマンスを改善するために、第1の分類器に第2の分類器を追加するためのバリエーションを例示するブロック図である。

【図6B】[0020]図6Bは、本開示の態様による、機械学習モデルのパフォーマンスを改善するために、第1の分類器に第2の分類器を追加するためのバリエーションを例示するブロック図である。

30

【図7】[0021]図7は、本開示の態様による、トレーニングされた機械学習モデルのパフォーマンスを改善するための例示的な分類器の概略図である。

【図8】[0022]図8は、本開示の態様による、トレーニングされた機械学習モデルのパフォーマンスを改善するための方法を例示する。

【図9】[0023]図9は、本開示の態様による、トレーニングされた機械学習モデルのパフォーマンスを改善するための方法を例示するブロック図である。

【発明の詳細な説明】

【0014】

[0024]添付された図面に関連して以下に示される詳細な説明は、様々な構成の説明として意図され、ここで説明される概念が実施され得る唯一の構成を表すようには意図されない。詳細な説明は、様々な概念の完全な理解を提供することを目的とした特定の詳細を含む。しかしながら、これらの概念が、これらの特定の詳細なしで実施され得ることは、当業者にとって明らかであろう。いくつかの事例では、周知の構造およびコンポーネントが、このような概念を曖昧にすることを避けるために、ブロック図形式で示される。

40

【0015】

[0025]本教示に基づき、当業者は、本開示の範囲が、本開示のその他任意の態様と組み合わせられてインプリメントされようと、あるいは独立してインプリメントされようと、本開示の任意の態様をカバーするように意図されていることを理解すべきである。例えば、示される任意の数の態様を使用して、装置がインプリメントされ得、または方法が実施され得る。加えて、本発明の範囲は、示される本開示の様々な態様に加えて、またはそれ以

50

外の、他の構造、機能、または構造と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ように意図される。開示される本開示の任意の態様が、請求項の1つまたは複数の要素によって具現化され得ることが理解されるべきである。

【0016】

[0026]「例示的(exemplary)」という用語は、ここで、「例、事例、または例示を提供する」という意味で使用される。「例示的」であるとしてここで説明される任意の態様は、必ずしも他の態様よりも好ましいまたは有利であるようには解釈されるべきでない。

【0017】

[0027]特定の態様がここで説明されるが、これらの態様の多くの変形および置換が、本開示の範囲内に含まれる。好ましい態様のいくつかの利益および利点が述べられるが、本開示の範囲は、特定の利益、用途または目的に限定されるようには意図されない。むしろ、本開示の態様は、異なる技術、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ように意図されており、そのうちのいくつかは、図面および好ましい態様の以下の説明において、例として例示される。詳細な説明および図面は、限定ではなく、本開示の単なる例示であり、本開示の範囲は、添付された特許請求の範囲およびそれらの同等物によって定義されている。

【0018】

[0028]本開示の態様が、トレーニングされたより低い複雑度の機械学習モデルのパフォーマンスを改善することに向けられる。本開示の態様によると、モデルのパフォーマンスは、低い複雑度の分類器の分類誤差の数を直接的に最小化するまたは低減させるように構成された第2の分類器を追加することによって改善され得る。すなわち、標準的な技法(例えば、勾配降下法(gradient descent))を使用して、典型的なコスト関数(例えば、平方和(SSE: sum of squares)、または負の対数尤度(negative log likelihood))によって与えられるような誤差の関数(function of errors)を最小化するのではなく、追加された分類器のための新しい目的関数が、誤差の数を直接的に最小化するまたは低減させるように定義される。例えば、7つの正確な分類と3つの不正確な分類と共に、分類動作(classification operations)が実行された場合、目的関数は、3つの誤差をゼロに低減させるように設計され得る。

【0019】

[0029]加えて、本開示の態様によると、トレーニングされたより低い複雑度の機械学習モデルのパフォーマンスは、より高い複雑度のモデルの軟確率(soft probabilities)を使用してさらに改善され得る。

【0020】

軟確率

[0030]軟確率は、確率ベクトルの未知の値(dark values)または非最大確率値(non-maximum probability values)である。多くの従来の分類システムでは、確率ベクトルは、クラスラベルを予測するために使用される。このような従来のシステムでは、クラスラベルは、確率ベクトルにおける最高または最大確率値を使用して予測される。非最大確率値または軟確率は、無視される。

【0021】

[0031]例えば、機械学習モデル $M(W)$ が、 $N$ 個のサンプルの入力データ $X^{tr} = [x_0, x_1, x_2, \dots, x_{N-1}]$ 、ここで、

【0022】

【数1】

$$x_i \in \mathbb{R}^D$$

【0023】

である、と、対応する $N$ 個のトレーニングサンプルの $C$ 個のラベル付けされた出力データ( $C$ -labeled output data) $y^{tr} = [y_0, y_1, y_2, \dots, y_{N-1}]$ 、ここで、 $y_i$

$[0, C-1]$ である、とから成るトレーニングデータを使用してトレーニングされる、教師あり機械学習の分類問題を考慮する。典型的に、機械学習モデル(例えば、ニュー

10

20

30

40

50

ラルネットワーク)のアーキテクチャを定義するパラメータ およびこのモデルをトレーニングするための学習プロセスのパラメータは、予め決定されている。その後、トレーニングデータ  $\{x^{tr}, y^{tr}\}$  は、モデル  $M$  の重み  $W$  を学習するために使用される。このトレーニングは、

【0024】

【数2】

$$p_j \in \mathbb{Z}_2^C$$

【0025】

となるように、 $P = [p_0, p_1, \dots, p_{N-1}]$  に、1-K符号化(1-K encoding)を使用して、ラベル付けされたデータ  $y = [y_0, y_1, \dots, y_{N-1}]$  を符号化することを含み得、ここで、 $y_j = k$  の場合、 $p_{jk} = 1$  であり、

【0026】

【数3】

$$\sum_{k=0}^{C-1} p_{jk} = 1$$

【0027】

である。

【0028】

[0032] 入力  $x$  を与えられると、機械学習モデル  $M$  は、出力確率に対する推定値 (estimate) を生成し、それは、

【0029】

【数4】

$$\hat{p} = M_\lambda(x, W) \quad (1)$$

【0030】

として表され得、結果として、

【0031】

【数5】

$$C = \sum_{i=0}^{N-1} \sum_{j=0}^{C-1} p_{ij} \log(\hat{p}_{ij}). \quad (2)$$

【0032】

によって与えられるマルチクラス交差エントロピー関数を最小化するようにする。

【0033】

[0033] 出力クラスラベルは、次のように得られる：

【0034】

【数6】

$$\hat{y} = \underset{j}{\operatorname{argmax}} [\hat{p}] \quad (3)$$

【0035】

[0034] したがって、硬確率 (hard-probability) と呼ばれる、ベクトル

【0036】

【数7】

$$\hat{p}$$

【0037】

の最大値の指標のみが、推論のために利用され、非最大値は、無視される。

【0038】

[0035] 本開示の態様は、分類パフォーマンスを改善するために軟確率を利用する。いくつかの態様では、軟確率は、温度スケールリング (temperature scaling) を使用して抽出され得る。例えば、ニューラルネットワークモデルによって生成される確率

10

20

30

40

50

【 0 0 3 9 】

【 数 8 】

 $\hat{p}$ 

【 0 0 4 0 】

は、次のようなソフトマックス関数 (softmax function) を介した推定である：

【 0 0 4 1 】

【 数 9 】

$$\hat{p}_k = \frac{\exp(a_{out,k})}{\sum_{j=0}^{C-1} \exp(a_{out,j})} \quad (4)$$

10

【 0 0 4 2 】

ここで、 $a_{out} = [a_{out,0}, a_{out,1}, \dots, a_{out,C-1}]$  は、ニューラルネットワークの出力ノードから出てきた活性化値 (activation values) である。

【 0 0 4 3 】

[0036] トレーニングされた機械学習モデル (例えば、ニューラルネットワーク) によって生成される出力確率は、次のように、軟確率において隠されている情報を抽出するために、温度  $T$  によってスケールリングされ得る：

【 0 0 4 4 】

【 数 1 0 】

$$\hat{p}_k^{Te} = \frac{\exp\left(\frac{a_{out,k}}{T}\right)}{\sum_{j=0}^{C-1} \exp\left(\frac{a_{out,j}}{T}\right)} \quad (5)$$

20

【 0 0 4 5 】

[0037] 1つの目的が、トレーニングされたモデルによって生成される確率ベクトル

【 0 0 4 6 】

【 数 1 1 】

 $\hat{p}$ 

【 0 0 4 7 】

の分布を緩和する (soften) ことである。温度  $T$  を介してスケールリングすることは、確率の分布を平坦化し (flattens)、それによって、軟確率における情報が活用されることを可能にする。

30

【 0 0 4 8 】

[0038] いったん抽出されると、軟確率は、分類パフォーマンスを改善するために使用され得る。例えば、 $b_m$  および  $W_m$  が軟確率における情報を共にプーリング (pooling together) するために使用されるバイアスおよび重みのセットを表す一例では、混合確率 (mixture probability) が、

【 0 0 4 9 】

【 数 1 2 】

$$\tilde{p}^{Te} = \frac{1}{1 + \exp(-(W_m \hat{p}^{Te} + b_m))} \quad (6)$$

40

【 0 0 5 0 】

によって与えられ得る。

【 0 0 5 1 】

[0039] 混合確率は、次のように、トレーニングされた機械学習モデルによって出力クラスラベルを予測するために使用され得る：

【 0 0 5 2 】

【 数 1 3 】

$$\tilde{y} = \underset{j}{\operatorname{argmax}} [\tilde{p}^{Te}] \quad (7)$$

【 0 0 5 3 】

50



[0040] トレーニングデータ  $\{X^{tr}, y^{tr}\}$  は、軟確率の混合を生成するために使用される重みおよびバイアスについての値を推定するために使用され得る。出力ラベルが硬確率のみを使用して予測されるとき、トレーニングされた機械学習モデルによって生成されるフラクショナルトレーニング誤差 (fractional training error)  $e_d$  (式3)、および出力ラベルが軟確率を使用して予測されるとき、フラクショナルトレーニング誤差 ( $e$ ) (式7) が、以下によって与えられる:

【0054】

【数14】

$$e_d = \frac{1}{N} \sum_{j=0}^{N-1} \mathbb{I}_{\hat{y}_j \neq y_j} \quad (8) \quad 10$$

$$e = \frac{1}{N} \sum_{j=0}^{N-1} \mathbb{I}_{\hat{y}_j \neq \hat{y}_j} \quad (9)$$

【0055】

[0041] コスト関数  $C$  が、分類誤差を低減させるために使用され得る。すなわち、コスト関数  $C$  は、軟確率の混合によって生成された出力ラベルに関する予測された値を使用したときのトレーニングデータに対する誤差 (the error on the training data) が、コスト関数が正の非ゼロ値をとる (takes on) 確率を使用することによって得られる誤差よりも低くなるように設計され得る。コスト関数は、次のように表され得る:

【0056】

【数15】

$$C = \max(0, (e_d - e) / e_d) \quad (10) \quad 20$$

【0057】

[0042] 軟確率の混合についての改善されたまたは最適な重みおよびバイアスが、次の最適化問題を解くことによって得られ得る:

【0058】

【数16】

$$\{W_m^*, b_m^*\} = \underset{\{W_m, b_m\}}{\operatorname{argmin}} [1 - C] \quad (11) \quad 30$$

【0059】

[0043] 式11の最適化問題は、初期状態  $\{W_m(0), b_m(0)\} = \{I, 0\}$  で (with) 勾配値を使用しない標準の制約なし最適化プロセス (standard unconstrained optimization processes) のうちの任意のものを使用して解かれ得る。いくつかの態様では、最適化技法はまた、軟確率を生成するための改善されたまたは最適な温度を決定するために用いられ得る。例えば、式11の最適化問題は、次のように修正され得る:

【0060】

【数17】

$$\{T^*, W_m^*, b_m^*\} = \underset{\{T, W_m, b_m\}}{\operatorname{argmin}} [1 - C] \quad (12) \quad 40$$

【0061】

[0044] 標準の制約なし最小化プロセスを使用することは、ある解をもたらし、それは、温度の初期選択の周辺の  $C$  についての局所極小値である (a local minima for  $C$  around the initial choice of the temperature)。収束ストラテジが、温度の初期選択の周辺の局所極小値から脱出するために使用され得る。例えば、いくつかの態様では、このストラテジは、パラメータの初期セット:  $\{T(0), W_m(0), b_m(0)\}$  で始まり、式11を使用して、重みおよびバイアスについての最適値

【0062】

【数 1 8】

$$\{W_m^{*,T(0)}, b_m^{*,T(0)}\}$$

【0063】

を求め得る。初期状態  $T(0)$  から開始し、コスト関数： $C = \max(0, (e - e) / e)$  を最適化し、ここで、 $e$  は、

【0064】

【数 1 9】

$$\{T'(0), W_m^{*,T(0)}, b_m^{*,T(0)}\}$$

10

【0065】

を用いて式 11 を使用して計算され、 $e$  は、

【0066】

【数 2 0】

$$\{T'(0), W_m^{*,T(0)}, b_m^{*,T(0)}\}$$

【0067】

を用いて式 11 を使用して計算される。このシーケンスは、収束するまで繰り返され得る。

【0068】

20

[0045]いくつかの態様では、アンサンブル平均化 (ensemble averaging) が、複数の機械学習モデルにわたって、および/または単一の機械学習モデルにおける複数のロジスティック回帰層 (logistic regression layers) にわたって、インプリメントされ得る。一例では、複数の機械学習モデル ( $M > 1$ ) が、 $M$  個のトレーニングされたモデルによって生成される出力確率

【0069】

【数 2 1】

$$\{\hat{p}_0, \hat{p}_1 \dots \hat{p}_{M-1}\}$$

【0070】

を有する (with) トレーニングデータを使用してトレーニングされる。これらのモデルの各々について、軟確率の最適な混合が、上記のプロシージャ最適化技法および/または収束ストラテジを使用して生成され得る。結果として得られる混合確率

30

【0071】

【数 2 2】

$$\{\tilde{p}_0^{Te_0}, \tilde{p}_1^{Te_1} \dots \tilde{p}_{M-1}^{Te_{M-1}}\}$$

【0072】

は、

【0073】

【数 2 3】

40

$$y^{\text{pred}} = \underset{j}{\operatorname{argmax}} \left[ \sum_k w_k \tilde{p}_k^{Te_k} \right] \quad (13)$$

【0074】

として出力ラベルを予測するために使用され得る。

【0075】

[0046]  $\{w_k\}$  についての 1 つの選択が、 $k = (1, 2, \dots, M - 1)$  について、 $w_k = 1 / M$  である。代替として、上記の最適化技法および収束ストラテジ、または他の同様の技法が、マルチモデル確率混合重み  $\{w_k\}$  の最適なセットを推定するために使用され得る。

50

## 【 0 0 7 6 】

[0047]別の例では、複数のロジスティック回帰出力層を有してではあるが、単一の機械学習モデルにおいて、最適化技法、収束ストラテジ、および同様のものは、このモデルの異なるロジスティック回帰層から結果として得られる軟確率を改善するまたは最適化するために使用され得る。

## 【 0 0 7 7 】

[0048]いくつかの態様では、推論は、クラスの数が多い(例えば、 $C = 1$ )ときに、軟確率を使用して改善され得る。軟確率の最適な混合を生成するためのパラメータの数は、 $C^2$ としてスケーリングされ、推論のための軟確率の混合を推定するとき問題になることがある。このケースでは、有用な情報を含むと考えられる各クラスについての最も高い軟確率のサブセット  $P \ll C$  が、分類パフォーマンスを改善するために活用され得る。順に、式 11 は、推定されるべきパラメータの総数が  $P(P + 1)$  となるように、重みおよびバイアスを得るために解かれ得る。推論時間またはおよそ推論時間において、上位  $P$  個の軟確率 (the top  $P$  soft probabilities) の指標が、追跡され、最適な重みおよびバイアスを使用して推定された混合確率を介して付加され得る。

10

## 【 0 0 7 8 】

[0049]図 1 は、本開示のある特定の態様による、汎用プロセッサ (CPU) またはマルチコア汎用プロセッサ (CPUs) 102 を含み得る、システムオンチップ (SOC) 100 を使用した、前述のトレーニングされた機械学習モデルのパフォーマンスを改善する方法の例となるインプリメンテーションを例示する。変数 (例えば、モデルの重み)、計算デバイスに関連付けられたシステムパラメータ (例えば、重みを有する機械学習モデル)、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 に関連付けられたメモリブロック中、CPU 102 に関連付けられたメモリブロック中、グラフィックス処理ユニット (GPU) 104 に関連付けられたメモリブロック中、デジタルシグナルプロセッサ (DSP) 106 に関連付けられたメモリブロック中、専用メモリブロック 118 中に記憶され得、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は、CPU 102 に関連付けられた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

## 【 0 0 7 9 】

[0050]SOC 100 はまた、GPU 104、DSP 106、接続性ブロック 110、これは、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、免許不要の Wi-Fi 接続性、USB 接続性、Bluetooth (登録商標) 接続性、および同様のものを含み得る、および、例えば、ジェスチャを検出および認識し得るマルチメディアプロセッサ 112 のような、特定の機能に合わせられた (tailored to) 追加の処理ブロックを含み得る。1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。

30

## 【 0 0 8 0 】

[0051]SOC 100 は、ARM 命令セットに基づき得る。本開示のある態様では、汎用プロセッサ 102 にロードされる命令は、第 1 の目的関数 (例えば、コスト) を有する第 1 の分類器に、第 2 の目的関数 (例えば、コスト) を有する第 2 の分類器を追加するためのコードを備え得る。第 2 の目的関数は、第 1 の分類器の誤差を直接的に低減させるために使用される。

40

## 【 0 0 8 1 】

[0052]図 2 は、本開示のある特定の態様による、システム 200 の例となるインプリメンテーションを例示する。図 2 に例示されるように、システム 200 は、ここで説明される方法の様々な動作を実行し得る複数のローカル処理ユニット (local processing units) 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、

50

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。加えて、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に例示されるように、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提供するための構成プロセッサユニット214と、およびローカル処理ユニット202間のルーティングを提供するルーティング接続処理ユニット216とインタフェースし得る。

#### 【0082】

[0053]ディープラーニングアーキテクチャは、各層において連続的により高度な抽象化レベル(successively higher levels of abstraction)で入力表現する(represent)ことを学習することによってオブジェクト認識タスクを実行し得、それにより、入力データの有用な特徴表現を構築(building up)する。このようにして、ディープラーニングは、従来の機械学習の主要なボトルネック(major bottleneck)に対処する。ディープラーニングが出現する前は、オブジェクト認識問題に対する機械学習アプローチは、ことによるとシャロー分類器(shallow classifier)との組合せにおいて、人間によって設計された特徴(human engineered features)に依存するところが大きかった。シャロー分類器は、2クラス線形分類器であり得、例えば、そこで、特徴ベクトル成分の加重和(weighted sum)が、どのクラスに入力が属するか予測するためにしきい値と比較され得る。人間によって設計された特徴は、領域の専門知識を有するエンジニアによって、特定の問題領域に合わせられたテンプレートまたはカーネルであり得る。ディープラーニングアーキテクチャは、対照的に、トレーニングを通じてであるが、人間のエンジニアが設計し得るものと同様の特徴を表現することを学習し得る。さらに、ディープネットワークは、人間が考慮することがなかったであろう新しいタイプの特徴を表現および認識することを学習し得る。

#### 【0083】

[0054]ディープラーニングアーキテクチャは、特徴の階層を学習し得る。例えば、視覚データが提示された場合、第1の層は、入力ストリームにおける、エッジのような、比較的単純な特徴を認識することを学習し得る。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パワーを認識することを学習し得る。第1の層の出力を入力として受ける第2の層は、視覚データについては単純な形状または聴覚データについては音の組合せのような、特徴の組合せを認識することを学習し得る。例えば、上位層は、視覚データにおける複雑な形状または聴覚データにおける単語を表現することを学習し得る。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ことを学習し得る。

#### 【0084】

[0055]ディープラーニングアーキテクチャは、自然階層構造を有する問題に適用される時、特によく機能し得る。例えば、モーターの付いた乗り物の分類は、車輪、風防ガラス、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。

#### 【0085】

[0056]ニューラルネットワークのような機械学習モデルは、様々な結合パターン(connectivity patterns)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、情報は、下位層から上位層へ渡されるとともに、所与の層における各ニューロンが、上位層におけるニューロンに伝達する。階層的な表現が、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した層において構築され得る。ニューラルネットワークはまた、再帰型(recurrent)結合またはフィードバック(トップダウンとも呼ばれる)結合を有し得る。再帰型結合では、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伝達され得る。再帰型アーキテクチャは、シーケンスにおいてニュー

10

20

30

40

50

ラルネットワークに伝達される入力データチャンクのうちの1つよりも多くにまたがる (span) パターンを認識するのに役立つ。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。多くのフィードバック結合を有するネットワークは、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区別することを支援し得るときに役立つ。

【0086】

[0057] 図3A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302または局所結合304であり得る。全結合ネットワーク302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の層における全てのニューロンから入力を受信するように、その出力を第2の層における全てのニューロンに伝達し得る。代替として、局所結合されたネットワーク304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306は、局所結合され得、第2の層における各ニューロンのための入力に関連付けられた結合強度が共有されるようにさらに構成される (例えば、308)。より一般的には、ネットワークの局所結合された層は、層における各ニューロンが、異なる値を有し得る結合強度を持ってではあるが、同じまたは同様の結合パターンを有するように構成され得る (例えば、310、312、314、および316)。局所結合された結合パターンは、所与の領域における上位層のニューロンが、ネットワークへの総入力の制限された部分の特性 (properties) にトレーニングを通じて調整される入力を受信し得るので、上位層における空間的に別個の受容野 (spatially distinct receptive fields) を生じさせ得る。

10

20

【0087】

[0058] 局所結合されたニューラルネットワークは、入力の空間的ロケーションが意味をもつ問題によく適し得る。例えば、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300は、異なる特性を有する上位層のニューロンを、画像の下部対上部 (the lower versus the upper portion) とのそれらの関連付けに依存して発達 (develop) させ得る。例えば、画像の下部に関連付けられたニューロンは、車線区分線を認識することを学習し得、一方、画像の上部に関連付けられたニューロンは、交通信号、交通標識、および同様のものを認識することを学習し得る。

【0088】

[0059] DCNは、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、DCNは、速度制限標識のクロップされた画像 (cropped image) 326のような、画像を提示され得、その後、「フォワードパス (forward pass)」が、出力322を生成するために計算され得る。出力322は、「標識」、「60」、および「100」のような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、例えば、トレーニングされたネットワーク300についての出力322において示される「標識」および「60」に対応するもののような、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について、高いスコアを出力することを望み得る。トレーニング前は、DCNによって生成される出力は、不正確である可能性が高く、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。その後、DCNの重みは、DCNの出力スコアがターゲットにより密接に合わせられる (aligned) ように調整され得る。

30

40

【0089】

[0060] 重みを調整するために、学習アルゴリズムが、重みについての勾配ベクトルを計算し得る。勾配は、重みがわずかに調整された場合に、誤差が増大または低減するであろう量を示し得る。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。下位層において、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。その後、重みは、誤差を低減させるように調整され得る。重みを調整するこの方法は、それがニューラルネットワークを通じた「バックワードパス (backward pass)」を伴うので、「バックプロパゲーション」と呼ばれ得る。

【0090】

50

[0061]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 (true error gradient) 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、確率的勾配降下法 (stochastic gradient descent) と呼ばれ得る。確率的勾配降下法は、システム全体の達成可能な誤差率 (error rate) の低減が止まるまで、または誤差率がターゲットレベルに達するまで繰り返され得る。

【 0 0 9 1 】

[0062]学習後、DCNは、新しい画像326を提示され得、ネットワークを通じたフォワードパスが、DCNの推論または予測と見なされ得る出力322をもたらし得る。

【 0 0 9 2 】

[0063]ディープピラミッドネットワーク (DBN) は、隠れノードの複数の層 (multiple layers of hidden nodes) を備える確率モデルである。DBNは、トレーニングデータセットの階層的な表現 (hierarchical representation) を抽出するために使用され得る。DBNは、制限付きボルツマンマシン (RBM: Restricted Boltzmann Machines) の層を積み上げること (stacking up) によって取得され得る。RBMは、入力の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。RBMは、各入力のカテゴリ化されるべきクラスについての情報がない状態で確率分布を学習することができるので、RBM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、DBNの下方のRBM (bottom RBM) は、教師なしの方法でトレーニングされ得、かつ特徴抽出器として機能し得、また、上方のRBM (top RBM) は、(ターゲットクラスおよび前の層からの入力の同時分布 (joint distribution) で) 教師ありの方法でトレーニングされ得、かつ分類器として機能し得る。

【 0 0 9 3 】

[0064]ディープ畳み込みネットワーク (DCN) は、畳み込みネットワークのネットワークであり、追加のプーリング層および正規化層を用いて構成される。DCNは、多くのタスクについて最先端のパフォーマンスを達成している。DCNは、入力ターゲットおよび出力ターゲットの両方が、多くのエグゼンプラー (exemplars) について知られており、かつ勾配降下法の使用によってネットワークの重みを修正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ることができ得る。

【 0 0 9 4 】

[0065]DCNは、フィードフォワードネットワークであり得る。加えて、上記で説明されたように、DCN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。DCNのフィードフォワード結合および共有結合は、高速処理に活用され得る。DCNの計算負担は、例えば、再帰型結合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、はるかに少なくなり得る。

【 0 0 9 5 】

[0066]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射影 (basis projection) であると見なされ得る。入力が、最初にカラー画像の赤、緑、および青チャネルのような複数のチャネルに分解される場合には、その入力で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の次元とを有する、3次元であると見なされ得る。畳み込み結合の出力は、後続の層318および320において特徴マップを形成すると考えられ得るとともに、特徴マップ (例えば、320) の各要素が、複数のチャネルの各々から、および前の層 (例えば、318) におけるある範囲のニューロン (a range of neurons) から入力を受信する。特徴マップにおける値は、正規化 (rectification)、すなわち  $\max(0, x)$  のような、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、追加の局所不変性および次元削減 (dimensionality reduction) を提供し得る。白色化に対応する正規化がまた、特徴マップにおけるニューロン間の側抑制 (lateral inhibition) を通じて適用され得る。

## 【 0 0 9 6 】

[0067]ディープラーニングアーキテクチャのパフォーマンスは、より多くのラベル付けされたデータポイントが利用可能になるにつれて、または計算能力が増大するにつれて増大し得る。現代のディープニューラルネットワークは、わずか15年前に一般的な研究者に利用可能であったものよりも何千倍も大きいコンピューティングリソースを用いて、日常的に(routinely)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は、ディープラーニングのパフォーマンスをさらに高め得る。整流された線形ユニットは、勾配消失(vanishing gradients)として知られるトレーニング課題を低減し得る。新しいトレーニング技法は、過剰適合(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象化(abstarct)し、全体的なパフォーマンスをさらに高め得る。

10

## 【 0 0 9 7 】

[0068]図3Bは、例示的なディープ畳み込みネットワーク350を例示するブロック図である。ディープ畳み込みネットワーク350は、結合性および重みの共有に基づいて、複数の異なるタイプの層を含み得る。図3Bに示されるように、例示的なディープ畳み込みネットワーク350は、複数の畳み込みブロック(例えば、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、正規化層(LN or m)、およびプーリング層で構成され得る。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するものではなく、代わりに、設計の選好に従って、任意の数の畳み込みブロックがディープ畳み込みネットワーク350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。例えば、正規化層は、白色化または側抑制を提供し得る。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提供し得る。

20

## 【 0 0 9 8 】

[0069]ディープ畳み込みネットワークの、例えば、並列フィルタバンク(parallel filter banks)は、高いパフォーマンスおよび低い電力消費を達成するために、オプションとしてARM命令セットに基づいて、SOC100のCPU102またはGPU104上に搭載(loaded on)され得る。代替的な実施形態では、並列フィルタバンクは、SOC100のDSP106またはISP116上に搭載され得る。加えて、DCNは、センサ114およびナビゲーション120に専用の処理ブロックのような、SOC上に存在し得る他の処理ブロックにアクセスし得る。

30

## 【 0 0 9 9 】

[0070]ディープ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(例えば、FC1およびFC2)を含み得る。ディープ畳み込みネットワーク350は、ロジスティック回帰(LR: logistic regression)層をさらに含み得る。ディープ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(例えば、画像、音声、ビデオ、センサデータ、および/または他の入力データ)から階層的な特徴表現を学習するために、ディープ畳み込みネットワーク350における後続の層の入力として機能し得る。

40

## 【 0 1 0 0 】

[0071]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を例示するブロック図である。このアーキテクチャを使用して、SOC420の様々な処理ブロック(例えば、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サポート計算(supporting computations)を実行させ得るアプリケーション402が設計され得る。

## 【 0 1 0 1 】

50

[0072] A I アプリケーション 4 0 2 は、例えば、そこで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もたらす (provide for) 得る、ユーザ空間 4 0 4 において定義された機能呼び出すように構成され得る。A I アプリケーション 4 0 2 は、例えば、認識されたシーンがオフィス、講堂、レストラン、または湖のような屋外環境であるかどうかによって異なるように、マイクロフォンおよびカメラを構成し得る。A I アプリケーション 4 0 2 は、現在のシーンの推定を提供するために、S c e n e D e t e c t アプリケーションプログラミングインタフェース (A P I) 4 0 6 において定義されたライブラリに関連付けられた、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、例えば、ビデオおよび測位データに基づいてシーン推定を提供するように構成されたディープニューラルネットワークの出力に最終的に依拠し得る。

10

#### 【 0 1 0 2 】

[0073] ランタイムフレームワーク (Runtime Framework) 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 4 0 8 が、A I アプリケーション 4 0 2 にとってさらにアクセス可能であり得る。A I アプリケーション 4 0 2 は、例えば、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る、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、順に、S O C 4 2 0 上で実行中の L i n u x (登録商標) カーネル 4 1 2 のような、オペレーティングシステム 4 1 0 に信号を送り得る。オペレーティングシステム 4 1 0 は、順に、C P U 4 2 2、D S P 4 2 4、G P U 4 2 6、N P U 4 2 8、またはこれらの何らかの組合せ上で、計算を実行させ得る。C P U 4 2 2 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、D S P 4 2 4 のための、G P U 4 2 6 のための、または N P U 4 2 8 のためのドライバ 4 1 4 ~ 4 1 8 のような、ドライバを通じてアクセスされ得る。例示的な例では、ディープニューラルネットワークは、C P U 4 2 2 および G P U 4 2 6 のような、処理ブロックの組合せ上で実行するように構成され得るか、または、存在する場合、N P U 4 2 8 上で実行され得る。

20

#### 【 0 1 0 3 】

[0074] 図 5 は、スマートフォン 5 0 2 上の A I アプリケーションのランタイム動作 5 0 0 を例示するブロック図である。A I アプリケーションは、画像 5 0 6 のフォーマットを変換し、その後、画像 5 0 8 をクロップおよび/またはリサイズするように (例えば、J A V 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 5 0 4 を含み得る。その後、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 (例えば、C プログラミング言語を使用して) 構成され得る S c e n e D e t e c t バックエンドエンジン 5 1 2 を含む分類アプリケーション 5 1 0 に通信され得る。S c e n e D e t e c t バックエンドエンジン 5 1 2 は、スケールアップ 5 1 6 およびクロッピング 5 1 8 によって、画像をさらに前処理 5 1 4 するように構成され得る。例えば、画像は、結果として得られる画像が 2 2 4 ピクセル x 2 2 4 ピクセルとなるように、スケールアップされ、クロップされ得る。これらの次元 (dimensions) 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、S O C 1 0 0 の様々な処理ブロックに、ディープ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、ディープニューラルネットワークブロック 5 2 0 によって構成され得る。その後、ディープニューラルネットワークの結果は、しきい値処理 (thresholded) 5 2 2 され、分類アプリケーション 5 1 0 内の指数平滑化 (exponential smoothing) ブロック 5 2 4 を通され得る。その後、平滑化された結果は、スマートフォン 5 0 2 の設定および/または表示の変更を生じ得る。

30

40

#### 【 0 1 0 4 】

[0075] 一構成では、機械学習モデルが、第 1 の目的 (例えば、コスト) 関数を有する第 1 の分類器に、第 2 の目的 (例えば、コスト) 関数を有する第 2 の分類器を追加するために構成され、第 2 の目的関数は、第 1 の分類器の誤差を直接的に低減させるために使用される。機械学習モデルはまた、トレーニングされた機械学習モデルを介して受信される入

50



力に基づいて、第2の分類器から特徴ベクトルを出力するために構成される。機械学習モデルは、追加する手段および/または出力する手段を含む。一態様では、追加する手段および/または出力する手段は、記載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に関連付けられたプログラムメモリ、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、上述された手段は、これら上述された手段によって、記載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。

【0105】

[0076]本開示のある特定の態様によると、各ローカル処理ユニット202は、ネットワークの所望の1つまたは複数の機能的特徴に基づいて、ネットワークのパラメータを決定することと、決定されたパラメータがさらに適合、調整、および更新されるにつれて、所望の機能的特徴に向けて1つまたは複数の機能的特徴を発展させることと、を行うように構成され得る。

10

【0106】

[0077]図6Aおよび図6Bは、ニューラルネットワークモデルのような機械学習モデルのパフォーマンスを改善するために、第1の分類器に第2の分類器を追加するためのバリエーションを例示するブロック図である。図6Aおよび図6Bを参照すると、第2の分類器602が、トレーニングされた機械学習モデル606の第1の分類器604に追加され得る。いくつかの態様では、機械学習モデル606は、局所結合された(L-C)層を含むディープ畳み込みネットワーク(DCN)または別の機械学習モデルを備え得る。この機械学習モデルは、複雑度が低くなり得る。いくつかの例示的な態様では、10億回未満の積和演算(MAC)を有する機械学習モデルが、低い複雑度のモデルであると考えられ得る。他方では、10億回を上回る積和演算を有する機械学習モデルが、高い複雑度のモデルであると考えられ得る。むしろ、他の測定基準もまた、モデルの相対的な複雑度を決定するために使用され得る(例えば、パラメータの数、ステージ(層)の数、および/またはステージのタイプ)。

20

【0107】

[0078]トレーニングされた機械学習モデル606は、入力(例えば、画像)(図示せず)を受信するように構成され得る。機械学習モデル606は、入力から特徴のセットを抽出するために、画像を処理し得る。入力に対応する特徴ベクトルが、第1の分類器604に供給され得る。第1の分類器604は、微分可能な(例えば、勾配が決定できる)目的関数を用いて構成され得、これは、分類精度(classification accuracy)を改善するために使用され得る。順に、第1の分類器604は、出力クラスラベルを決定するために使用され得る確率ベクトル $P_c$ を生成し得る。

30

【0108】

[0079]第1の分類器604のパフォーマンスおよび精度を改善するために、第2の分類器602が追加され得る。第2の分類器602は、微分不可能な(non-differentiable)(例えば、勾配がない)目的関数を用いて構成され得る。目的関数は、第1の分類器604によって生成される誤差の数を直接的に低減させるように構成され得る。すなわち、第1の分類器604のためのコスト関数または誤差の関数を最小化するように試みるのではなく、第2の分類器602は、誤差の総数を低減させる。例えば、いくつかの態様では、第2の分類器602のための目的関数は、

40

【0109】

【数24】

$$\text{Objective function: } \operatorname{argmax} [\max(0, (e_d - e))] \quad (14)$$

【0110】

として表され得る。

【0111】

[0080]目的関数は、上記で説明されたような制約なし最小化技法を使用して、第2の分類器602のための重みおよびバイアスの項を決定するために使用され得る。したがって

50

、第2の分類器602からの出力クラスラベルは、第2の分類器602のみを介して生成されるよりも少ない誤差を含み得る。

【0112】

[0081]この構成は、分類パフォーマンスにおける改善が、予めトレーニングされた機械学習モデルを再トレーニングすることなく達成され得るので、特に有益であり得る。代わりに、パフォーマンスは、第2の分類器602を再トレーニングすることのみによって改善され得る。

【0113】

[0082]いくつかの態様では、図6Bに示されるように、第2の分類器602は、代替として、(例えば、トレーニングされた機械学習モデルからのモデルの層として)トレーニングされた機械学習モデル606内に提供され得る。さらに、いくつかの態様では、(図6Aおよび図6Bに示される)機械学習モデル606のパフォーマンスは、高い複雑度のモデル608を介して供給される軟確率を使用してさらに改善され得る。

10

【0114】

[0083]図7は、本開示の態様による、トレーニングされた機械学習モデル(例えば、ニューラルネットワーク)のパフォーマンスを改善するための例示的な分類器700の概略図を提示する。図7を参照すると、微分不可能な目的関数 $\mathcal{O}$ が、ニューラルネットワークの分類器(回帰)層の出力において追加される。目的関数は、所与のトレーニング(またはテスト)データセットに対する目的関数についての最大の非ゼロ値が、トレーニング(テスト)誤差の数が、元のトレーニングされたニューラルネットワークにつ

20

。

【0115】

[0084]入力

【0116】

【数25】

$$X \in \mathbb{R}^D$$

【0117】

を与えられると、機械学習モデル702は、入力をC個のクラスのうちの1つに分類するように構成され得る。ワンホット符号化(one-hot encoding)のような、符号化スキームを使用して、クラスラベルは、所与のクラスラベル $1 < C$ について、 $P = [p_1 p_2 \dots p_C]^T$ であるように、確率ベクトル

30

【0118】

【数26】

$$P \in \mathbb{Z}_2^C$$

【0119】

によって表され得、ここで、 $i = 1$ の場合、 $p_i = 1$ であり、

【0120】

【数27】

$$\sum_{i=1}^C p_i = 1$$

40

【0121】

である。トレーニングされた機械学習モデル(例えば、ニューラルネットワーク)

【0122】

【数28】

$$M: X \in \mathbb{R}^D \rightarrow Z \in \mathbb{R}^C$$

【0123】

を与えられると、推定された確率ベクトル

【0124】

【数29】

$$\hat{P}$$

50

【 0 1 2 5 】

が、

【 0 1 2 6 】

【 数 3 0 】

$$\hat{P} = \sigma(Z)$$

【 0 1 2 7 】

として、Zから得られ得、ここで、 $\sigma$  は、ソフトマックス非線形性 (soft-max nonlinearity) である。

【 0 1 2 8 】

[0085] 上記で説明されたように、従来のアプローチは、

10

【 0 1 2 9 】

【 数 3 1 】

$$\hat{I} = \operatorname{argmax} [\hat{P}]$$

【 0 1 3 0 】

としてクラスラベルを予測するために、

【 0 1 3 1 】

【 数 3 2 】

$$\hat{P}$$

【 0 1 3 2 】

20

を使用する。U個のトレーニングサンプルを有する所与のデータセットについて、その後、トレーニング誤差は、

【 0 1 3 3 】

【 数 3 3 】

$$e_d^{tr} = \frac{1}{U} \sum_{i=1}^U \mathbb{1}_{i \neq \hat{I}}$$

【 0 1 3 4 】

として得られ、V個のテストサンプルに対するテスト誤差は、同様に、

【 0 1 3 5 】

【 数 3 4 】

30

$$e_d^{ts} = \frac{1}{V} \sum_{i=1}^V \mathbb{1}_{i \neq \hat{I}}$$

【 0 1 3 6 】

として得られる。

【 0 1 3 7 】

【 数 3 5 】

$$e_d^{tr}$$

【 0 1 3 8 】

および

40

【 0 1 3 9 】

【 数 3 6 】

$$e_d^{ts}$$

【 0 1 4 0 】

についての値は、モデルMの優良性 (goodness) または精度を決定する。トレーニングモデルMについての1つの優良性または精度の測定基準が、

【 0 1 4 1 】

【数 3 7】

$$e_d^{tr} = 0$$

【 0 1 4 2】

および

【 0 1 4 3】

【数 3 8】

$$e_d^{ts} \ll 1$$

【 0 1 4 4】

である。本開示の態様は、それに対して

【 0 1 4 5】

【数 3 9】

$$e_d^{tr} \neq 0$$

【 0 1 4 6】

である、トレーニングされたモデルMのパフォーマンスを改善することを目的とする。

【 0 1 4 7】

[0086]本開示の態様によると、トレーニングされたモデル702を介して生成される特徴表現が、分類器700に供給され得る。分類器700は、特徴ベクトルZを受信し、それは、モデルの重み $W_z$ と混合(mixed)されて、新しい特徴ベクトル

【 0 1 4 8】

【数 4 0】

$$Z_s = W_z^T Z$$

【 0 1 4 9】

を生成し得る。その後、特徴ベクトル $Z_s$ は、確率ベクトル $P_s = (Z_s)$ を推定するために使用され得る。その後、確率特徴ベクトル

【 0 1 5 0】

【数 4 1】

$$P_f = W_p^T P_s$$

【 0 1 5 1】

が、

【 0 1 5 2】

【数 4 2】

$$e^{tr} = \frac{1}{U} = \sum_{i=1}^U \mathbb{1}_{\hat{l}_f \neq l_i}$$

【 0 1 5 3】

ここで、

【 0 1 5 4】

【数 4 3】

$$\hat{l}_f = \operatorname{argmax} [P_f]$$

【 0 1 5 5】

である、としてトレーニングセットに対する推定予測誤差を計算するために使用され得る。パラメータ  $W = [W_z, W_p]$  は、次の目的関数上で最適化することによって推定される：

【 0 1 5 6】

10

20

30

40

【数 4 4】

$$O = \text{MAX} \left( 0, (e_d^{tr} - e^{tr}) \right) \quad (14)$$

【 0 1 5 7】

[0087]いくつかの態様では、高い複雑度のモデル 7 0 4 が、機械学習モデル 7 0 2 に軟確率ベクトル  $P_H$  を提供し得る。軟確率ベクトルは、モデルの重み  $W_h$  と混合され得る。順に、確率特徴ベクトル

【 0 1 5 8】

【数 4 5】

$$P_f = W_p^T P_s + W_h^T P_H \quad 10$$

【 0 1 5 9】

は、

【 0 1 6 0】

【数 4 6】

$$e^{tr} = \frac{1}{U} = \sum_{i=1}^U \mathbb{1}_{i_f \neq i}$$

【 0 1 6 1】

ここで、

【 0 1 6 2】

【数 4 7】

$$\hat{i}_f = \text{argmax} [P_f] \quad 20$$

【 0 1 6 3】

である、としてトレーニングセットに対する推定予測誤差を計算するために使用され得る。パラメータ  $\theta = [W_z, W_p, W_h, T]$  は、式 1 4 の目的関数上で最適化することによって推定され得る。

【 0 1 6 4】

[0088]  $O$  が微分不可能な関数であると仮定すると、制約なし最小化プロセスは、 $\theta^* = \text{argmax} [O]$  として、最適な  $\theta^*$  を求めるために使用され得る。 $O$  に関する非ゼロ収束値は、

【 0 1 6 5】

【数 4 8】

$$e_d^{tr} < e^{tr}$$

【 0 1 6 6】

ということ暗示し、したがって、パラメータの追加のセットを推定するという代償を払って、元のモデルよりも良いパフォーマンスを有する、結果として得られるモデルを生成する。

【 0 1 6 7】

[0089]いくつかの態様では、 $\theta$  (例えば、 $W_z$ 、 $W_p$ 、 $W_h$ 、または  $T$ ) におけるパラメータのうちいくつかは、アプリアリに設定され得る。したがって、新しいパラメータのうちいくつかの追加による過剰適合 (overfitting) の問題が、緩和または低減され得る。

【 0 1 6 8】

[0090]いくつかの態様では、様々な単純化が、設計の選好に従ってパフォーマンスを改善しながら用いられ得る。例えば、トレーニングされた学習モデルによって生成される特徴に対応する重みが、アイデンティティ値 (identity value) に設定され得る。したがって、トレーニングされた機械学習モデルによって生成される特徴ベクトルの混合は、考慮されないであろう。一方、第 2 の例では、トレーニングされた機械学習モデルを介して生成される特徴ベクトルの混合のみが考慮され得る。

## 【 0 1 6 9 】

[0091]第3の例では、トレーニングされた学習モデルによって生成される特徴に対応する重みが、アイデンティティ値に設定され得、高い複雑度のモデル704からの利用可能な軟確率情報は、無視され得る。

## 【 0 1 7 0 】

[0092]第4の例では、高い複雑度のモデル704からの軟確率 $P_H$ が、所与の温度値（例えば、 $T = \quad, \quad > 1$ ）によって再スケーリング（rescaled）され得る。

## 【 0 1 7 1 】

[0093]図8は、トレーニングされた機械学習モデルのパフォーマンスを改善するための方法800を例示する。ブロック802において、プロセスは、第1の目的関数（例えば、コスト）を有する第1の分類器に、第2の目的関数（例えば、コスト）を有する第2の分類器を追加する。第2の目的関数は、第1の分類器の誤差を直接的に低減させるために使用される。

10

## 【 0 1 7 2 】

[0094]第1の目的関数は、微分可能であり、第2の目的関数は、微分不可能である。いくつかの態様では、第2の目的関数は、第1の分類器と第2の分類器との誤差間の差の関数であり得る。他の態様では、第2の目的関数は、より高い複雑度のモデルからの確率の混合に基づいて決定され得る。

## 【 0 1 7 3 】

[0095]いくつかの態様では、第2の分類器は、第1の分類器の外部に（externally）追加され得る。代替として、第2の分類器は、第1の分類器（例えば、第1の分類器の層）の内部に組み込まれ得る。さらに、第2の分類器は、第1の分類器を再トレーニングすることなく追加され得る。

20

## 【 0 1 7 4 】

[0096]ブロック804において、プロセスは、トレーニングされた機械学習モデルを介して受信される入力に基づいて、第2の分類器から特徴ベクトルを出力する。

## 【 0 1 7 5 】

[0097]いくつかの態様では、プロセスは、過剰適合の問題を低減または緩和するために、様々な単純化をインプリメントし得る。例えば、プロセスは、アイデンティティ値に、第1の分類器によってトレーニングされたモデルによって生成される特徴に対する重みを割り当て（assign）得る。プロセスはまた、ゼロに、より高い複雑度のモデルの確率ベクトルによって生成される特徴に対する重みを割り当て得る。プロセスは、第2の分類器の確率ベクトルによって生成される特徴に対する重みをさらに割り当て得る。プロセスは、ゼロに、より高い複雑度のモデルの確率ベクトルによって生成される特徴に対する重みを割り当て得る。プロセスは、第2の分類器の確率ベクトルによって生成される特徴に対する重みをさらに割り当て得る。プロセスはまた、固定された温度 $T$ によって、より高い複雑度のモデルによって生成される確率ベクトルをスケーリングし得る。

30

## 【 0 1 7 6 】

[0098]図9は、本開示の態様による、トレーニングされた機械学習モデルのパフォーマンスを改善するための方法900を例示するブロック図である。ブロック902において、プロセスは、トレーニングされた機械学習モデルを介して、機械確率ベクトルを機械学習モデル（例えば、分類器）において受信する。確率ベクトルは、トレーニングされた機械学習モデルにおいて受信される入力に対応する。ブロック904において、モデルの重みおよびバイアスのような、機械学習モデルのパラメータは、トレーニングされた機械学習モデルの誤差を直接的に低減させる目的関数に基づいて計算され得る。すなわち、目的関数は、トレーニングされた機械学習モデルについての誤差の関数ではなく、誤差の数を直接的に低減させるように設計される。したがって、機械学習モデルの目的関数は、微分不可能である。

40

## 【 0 1 7 7 】

[0099]いくつかの態様では、高い複雑度のモデルおよび/またはトレーニングされた機

50

械学習モデルからの軟確率は、パラメータを計算するために使用され得る。

【0178】

[00100]ブロック906において、プロセスは、機械学習モデルのパラメータを更新し得る。その後、機械学習モデルは、ブロック908において、受信された確率ベクトルに対応する入力についての出力クラスラベルを生成し得る。したがって、更新の後に続く分類誤差は、同じ入力についてトレーニングされた機械学習モデルによって生成されるそれらよりも少なくなり得る。したがって、トレーニングされた機械学習モデルのパフォーマンスは、改善され得る。

【0179】

[00101]上記で説明された方法の様々な動作は、対応する機能を実行することが可能な任意の適切な手段によって実行され得る。これら手段は、それに限定されるものではないが、回路、特定用途向け集積回路(AASIC)、またはプロセッサを含む、様々なハードウェアおよび/またはソフトウェアの(1つまたは複数の)コンポーネントおよび/または(1つまたは複数の)モジュールを含み得る。一般に、図中に例示された動作がある場合、これらの動作は、同様の番号付けを有する、対応する対をなすミーンズプラスファンクションのコンポーネントを有し得る。

【0180】

[00102]ここで使用される場合、「決定すること」という用語は、幅広いアクションを包含する。例えば、「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(例えば、表、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)、確定すること、および同様のことを含み得る。加えて、「決定すること」は、受信すること(例えば、情報を受信すること)、アクセスすること(例えば、メモリにおけるデータにアクセスすること)、および同様のことを含み得る。さらに、「決定すること」は、解決すること、選択すること、選ぶこと、確立すること、および同様のことを含み得る。

【0181】

[00103]ここで使用される場合、アイテムのリスト「のうちの少なくとも1つ」を指す表現は、単一のメンバ(members)を含む、それらのアイテムの任意の組合せを指す。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をカバーするように意図される。

【0182】

[00104]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、デジタルシグナルプロセッサ(DSP)、特定用途向け集積回路(AASIC)、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ロジック、個別ハードウェアコンポーネント、あるいはここで説明された機能を実行するように設計されたこれらの任意の組合せを用いてインプリメントまたは実行され得る。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、このプロセッサは、任意の商業的に利用可能なプロセッサ、コントローラ、マイクロコントローラまたはステートマシンであり得る。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、例えば、DSPとマイクロプロセッサの組合せ、複数のマイクロプロセッサ、DSPコアと連携した1つまたは複数のマイクロプロセッサ、あるいはその他任意のこのような構成としてインプリメントされ得る。

【0183】

[00105]本開示に関連して説明されたアルゴリズムまたは方法のステップは、直接ハードウェアにおいて、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において、またはこれら2つの組合せにおいて、具現化され得る。ソフトウェアモジュールは、当該技術分野で知られている任意の形態の記憶媒体内に存在し得る。使用され得る記憶媒体のいくつかの例は、ランダムアクセスメモリ(RAM)、読取専用メモリ(ROM)、フラッシュメモリ、消去可能なプログラマブル読取専用メモリ(EPROM)、電氣的に消去可

10

20

30

40

50

能なプログラマブル読取専用メモリ（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、等を含む。ソフトウェアモジュールは、単一の命令、または多くの命令を備え得、いくつかの異なるコードセグメントにわたって、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、プロセッサが記憶媒体から情報を読み取り、および/または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、プロセッサと一体化され得る。

**【0184】**

[00106]ここ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。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置き換えられ得る。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明記されない限り、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修正され得る。

**【0185】**

[00107]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、またはこ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、例となる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アプリケーションおよび全体的な設計制約に依存して、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、機械可読媒体、およびバスインタフェースを含む様々な回路を共にリンクし得る。バスインタフェースは、特に、バスを介してネットワークアダプタを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。ある特定の態様では、ユーザインタフェース（例えば、キーボード、ディスプレイ、マウス、ジョイスティック、等）がまた、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、電圧レギュレータ、電力管理回路、および同様のもののような、様々な他の回路をリンクし得、これらは、当該技術分野において周知であり、したがって、これ以上は説明されない。

**【0186】**

[00108]プロセッサは、バスの管理と、機械可読媒体上に記憶されたソフトウェアの実行を含む汎用処理とを担い得る。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その他の回路を含む。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、またはその他の方法で呼ばれるかにかかわらず、命令、データ、またはこれらの任意の組合せを意味するように広く解釈されるべきである。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、読取専用メモリ（ROM）、プログラマブル読取専用メモリ（PROM）、消去可能なプログラマブル読取専用メモリ（EPROM）、電氣的に消去可能なプログラマブル読取専用メモリ（EEPROM）、レジスタ、磁気ディスク、光ディスク、ハードドライブ、またはその他任意の適切な記憶媒体、あるいはこ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具現化され得る。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。

**【0187】**

[00109]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、当業者が容易に理解するであろうように、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、伝送路（transmission line）、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、これら全ては、バスインタフェースを通じてプロセッサによってアクセスされ得る。代替として、またはこれに加えて、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

10

20

30

40

50



スタファイルでのケースでそうであり得るように、プロセッサに組み込まれ得る。ローカルコンポーネントのような、説明された様々なコンポーネントは、特定のロケーションを有するものとして説明されているが、それらはまた、分散型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ある特定のコンポーネントのように、様々な方法で構成され得る。

【0188】

[00110]処理システムは、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分を提供する外部メモリとを有し、全てが外部バスアーキテクチャを通じて他のサポート回路と共にリンクされている、汎用処理システムとして構成され得る。代替として、処理システムは、ここで説明されたニューロンモデルおよびニューラルシステムのモデル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、プロセッサと、バスインタフェースと、ユーザインタフェースと、サポート回路と、単一のチップに組み込まれた機械可読媒体の少なくとも一部分とを有する特定用途向け集積回路(AASIC)を用いて、または、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(FPGA)、プログラマブル論理デバイス(PLD)、コントローラ、ステートマシン、ゲート論理、個別ハードウェアコンポーネント、またはその他任意の適切な回路、あるいは本開示全体にわたって説明された様々な機能を実行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。当業者であれば、特定のアプリケーションおよびシステム全体に課せられる全体的な設計制約に依存して、処理システムに関して説明された機能をインプリメントするのに最良の方法を認識するであろう。

【0189】

[00111]機械可読媒体は、多数のソフトウェアモジュールを備え得る。これら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ると、様々な機能を処理システムに実行させる命令を含む。これらソフトウェアモジュールは、送信モジュールおよび受信モジュールを含み得る。各ソフトウェアモジュールは、単一の記憶デバイス内に存在し得るか、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、ソフトウェアモジュールは、トリガリングイベントが生じたときに、ハードドライブからRAM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、アクセス速度を増大させるために、命令のうちいくつかをキャッシュにロードし得る。その後、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するとき、プロセッサによってインプリメントされるということが理解されるであろう。さらに、本開示の態様が、プロセッサ、コンピュータ、機械、またはこのような態様をインプリメントする他のシステムの機能に改善をもたらすことが理解されるべきである。

【0190】

[00112]ソフトウェアでインプリメントされる場合、これら機能は、コンピュータ可読媒体上で、1つまたは複数の命令またはコードとして送信または記憶され得る。コンピュータ可読媒体は、1つの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容易にする任意の媒体を含む通信媒体とコンピュータ記憶媒体との両方を含む。記憶媒体は、コンピュータによってアクセスされることができ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、このようなコンピュータ可読媒体は、RAM、ROM、EEPROM、CD-ROMまたは他の光ディスク記憶装置、磁気ディスク記憶装置またはその他の磁気記憶デバイス、あるいは、データ構造または命令の形式で所望のプログラムコードを記憶または搬送するために使用されることができ、かつ、コンピュータによってアクセスされることができるその他任意の媒体を備えることができる。また、任意の接続は、厳密にはコンピュータ可読媒体と称される。例えば、ソフトウェアが、同軸ケーブル、光ファイバケーブル、ツイストペア、デジタル加入者回線(DSL)、または赤外線(IR)、無線、およびマイクロ波のようなワイヤレス技術を使用して、ウェブサイト、サーバ、また

は他の遠隔ソースから送信される場合には、同軸ケーブル、光ファイバーケーブル、ツイストペア、DSL、または赤外線、無線、およびマイクロ波のようなワイヤレス技術は、媒体の定義に含まれる。ここで使用される場合、ディスク(disk)およびディスク(disc)は、コンパクトディスク(CD)、レーザーディスク(登録商標)、光ディスク、デジタル多目的ディスク(DVD)、フロッピー(登録商標)ディスク、およびブルーレイ(登録商標)ディスクを含み、ここでディスク(disks)は、通常磁気的にデータを再生し、一方ディスク(disks)は、レーザーを用いて光学的にデータを再生する。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、非一時的なコンピュータ可読媒体(例えば、有形媒体)を備え得る。加えて、他の態様では、コンピュータ可読媒体は、一時的なコンピュータ可読媒体(例えば、信号)を備え得る。上記の組合せもまた、コンピュータ可読媒体の範囲内に含まれるべきである。

10

【0191】

[00113]したがって、ある特定の態様は、ここで提示された動作を実行するためのコンピュータプログラム製品を備え得る。例えば、このようなコンピュータプログラム製品は、その上に命令が記憶された(および/または符号化された)コンピュータ可読媒体を備え得、これら命令は、ここ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。ある特定の態様では、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。

【0192】

[00114]さらに、ここ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、適宜、ユーザ端末および/または基地局によって、ダウンロードされ得ること、および/または、別の方法で取得され得ることが理解されるべきである。例えば、このようなデバイスは、ここで説明された方法を実行するための手段の転送を容易にするためにサーバに結合されることができる。代替として、ここ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、デバイスに記憶手段を結合または提供する際に、様々な方法を得ることができるよう、記憶手段(例えば、RAM、ROM、コンパクトディスク(CD)またはフロッピーディスクのような物理記憶媒体、等)を介して提供されることができる。さらに、ここ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その他任意の適切な技法が、利用されることができる。

20

【0193】

[00115]特許請求の範囲は、上記に例示された厳密な構成およびコンポーネントに限定されないことが理解されるべきである。様々な修正、変更、および変形が、特許請求の範囲から逸脱することなく、上記で説明された方法および装置の配置、動作および詳細において行われ得る。

30

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。

【C1】 トレーニングされた機械学習モデルのパフォーマンスを改善するための方法であって、

第1の目的関数を有する第1の分類器に、第2の目的関数を有する第2の分類器を追加すること、前記第2の目的関数は、前記第1の分類器の誤差を直接的に低減させるために使用される、

40

を備える方法。

【C2】 前記第1の目的関数は、微分可能である、C1に記載の方法。

【C3】 前記第2の目的関数は、微分不可能である、C1に記載の方法。

【C4】 前記第2の目的関数は、前記第1の分類器と前記第2の分類器との誤差間の差の関数である、C1に記載の方法。

【C5】 より高い複雑度のモデルからの確率の混合に少なくとも部分的に基づいて、前記第2の目的関数を決定することをさらに備える、C1に記載の方法。

【C6】 前記第1の分類器を再トレーニングすることなく、前記第2の分類器を追加することをさらに備える、C1に記載の方法。

【C7】 前記第1の分類器の外部に前記第2の分類器を追加することをさらに備える、

50

C 1 に記載の方法。

[ C 8 ] アイデンティティ値に、前記第 1 の分類器によってトレーニングされたモデルによって生成される特徴に対する重みを割り当てることをさらに備える、C 1 に記載の方法。

[ C 9 ] ゼロに、より高い複雑度のモデルの確率ベクトルによって生成される特徴に対する重みを割り当てることをさらに備える、C 8 に記載の方法。

[ C 10 ] 前記第 2 の分類器の確率ベクトルによって生成される特徴に対する重みを割り当てることをさらに備える、C 1 に記載の方法。

[ C 11 ] ゼロに、より高い複雑度のモデルの確率ベクトルによって生成される特徴に対する重みを割り当てることをさらに備える、C 1 に記載の方法。

[ C 12 ] 前記第 2 の分類器の確率ベクトルによって生成される特徴に対する重みを割り当てることをさらに備える、C 11 に記載の方法。

[ C 13 ] 固定された温度 T によって、より高い複雑度のモデルによって生成される確率ベクトルをスケールリングすることをさらに備える、C 1 に記載の方法。

[ C 14 ] トレーニングされた機械学習モデルのパフォーマンスを改善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

を備え、前記少なくとも 1 つのプロセッサは、

第 1 の目的関数を有する第 1 の分類器に、第 2 の目的関数を有する第 2 の分類器を追加すること、前記第 2 の目的関数は、前記第 1 の分類器の誤差を直接的に低減させるために使用される、

を行うように構成される、装置。

[ C 15 ] 前記第 1 の目的関数は、微分可能である、C 14 に記載の装置。

[ C 16 ] 前記第 2 の目的関数は、微分不可能である、C 14 に記載の装置。

[ C 17 ] 前記第 2 の目的関数は、前記第 1 の分類器と前記第 2 の分類器との誤差間の差の関数である、C 14 に記載の装置。

[ C 18 ] 前記少なくとも 1 つのプロセッサは、より高い複雑度のモデルからの確率の混合に少なくとも部分的に基づいて、前記第 2 の目的関数を決定するようにさらに構成される、C 14 に記載の装置。

[ C 19 ] 前記少なくとも 1 つのプロセッサは、前記第 1 の分類器を再トレーニングすることなく、前記第 2 の分類器を追加するようにさらに構成される、C 14 に記載の装置

。

[ C 20 ] 前記少なくとも 1 つのプロセッサは、前記第 1 の分類器の外部に前記第 2 の分類器を追加するようにさらに構成される、C 14 に記載の装置。

[ C 21 ] 前記少なくとも 1 つのプロセッサは、アイデンティティ値に、前記第 1 の分類器によってトレーニングされたモデルによって生成される特徴に対する重みを割り当てるようにさらに構成される、C 14 に記載の装置。

[ C 22 ] 前記少なくとも 1 つのプロセッサは、ゼロに、より高い複雑度のモデルの確率ベクトルによって生成される特徴に対する重みを割り当てるようにさらに構成される、C 21 に記載の装置。

[ C 23 ] 前記少なくとも 1 つのプロセッサは、前記第 2 の分類器の確率ベクトルによって生成される特徴に対する重みを割り当てるようにさらに構成される、C 14 に記載の装置。

[ C 24 ] 前記少なくとも 1 つのプロセッサは、ゼロに、より高い複雑度のモデルの確率ベクトルによって生成される特徴に対する重みを割り当てるようにさらに構成される、C 14 に記載の装置。

[ C 25 ] 前記少なくとも 1 つのプロセッサは、前記第 2 の分類器の確率ベクトルによって生成される特徴に対する重みを割り当てるようにさらに構成される、C 24 に記載の装置。

10

20

30

40

50

[ C 2 6 ] 前記少なくとも1つのプロセッサは、固定された温度Tによって、より高い複雑度のモデルによって生成される確率ベクトルをスケーリングするようにさらに構成される、C 1 4に記載の装置。

[ C 2 7 ] トレーニングされた機械学習モデルのパフォーマンスを改善するための装置であって、

第1の目的関数を有する第1の分類器に、第2の目的関数を有する第2の分類器を追加するための手段と、前記第2の目的関数は、前記第1の分類器の誤差を直接的に低減させるために使用される、

前記トレーニングされた機械学習モデルを介して受信される入力に少なくとも部分的に基づいて、前記第2の分類器から特徴ベクトルを出力するための手段と

を備える装置。

[ C 2 8 ] トレーニングされた機械学習モデルのパフォーマンスを改善するためのプログラムコードをその上に符号化された非一時的なコンピュータ可読媒体であって、前記プログラムコードは、プロセッサによって実行され、第1の目的関数を有する第1の分類器に、第2の目的関数を有する第2の分類器を追加するためのプログラムコードを備え、前記第2の目的関数は、前記第1の分類器の誤差を直接的に低減させるために使用される、非一時的なコンピュータ可読媒体。

【 図 1 】

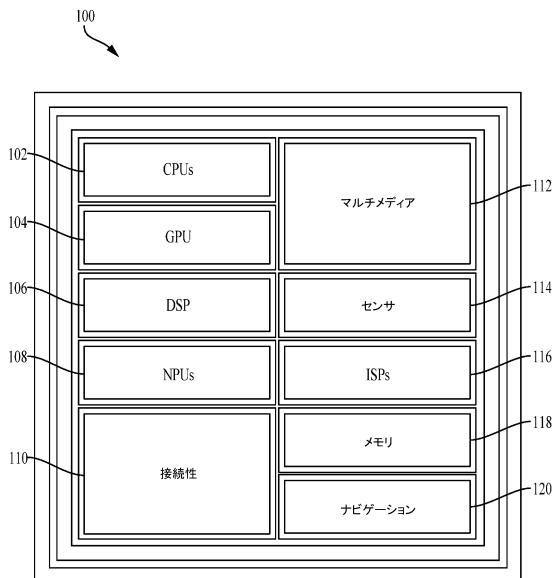


FIG. 1

【 図 2 】

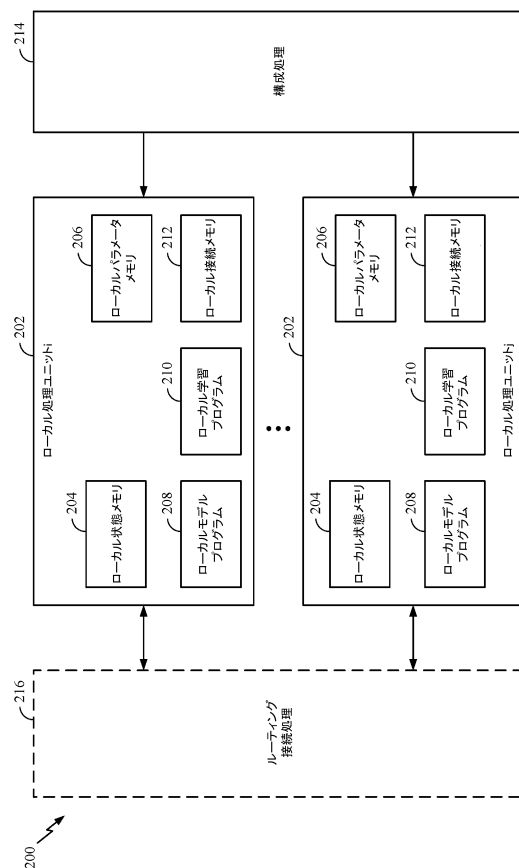


FIG. 2

【図3A】

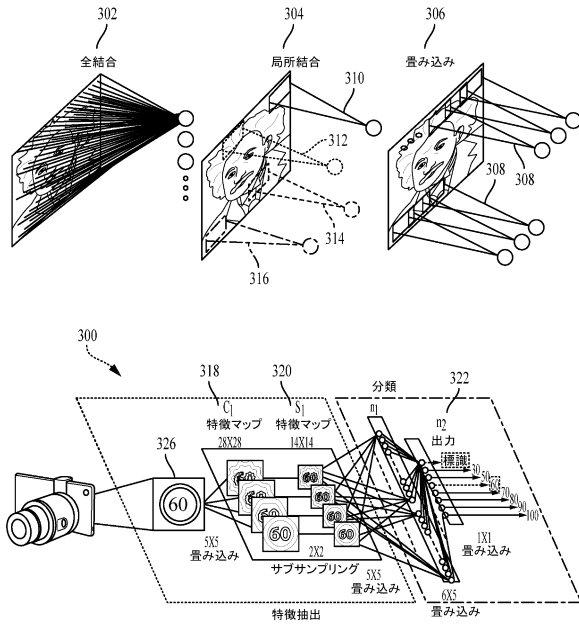


FIG. 3A

【図3B】

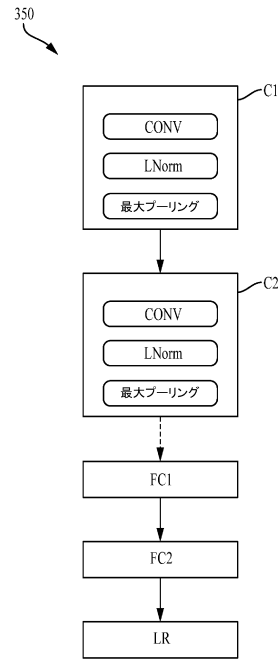


FIG. 3B

【図4】

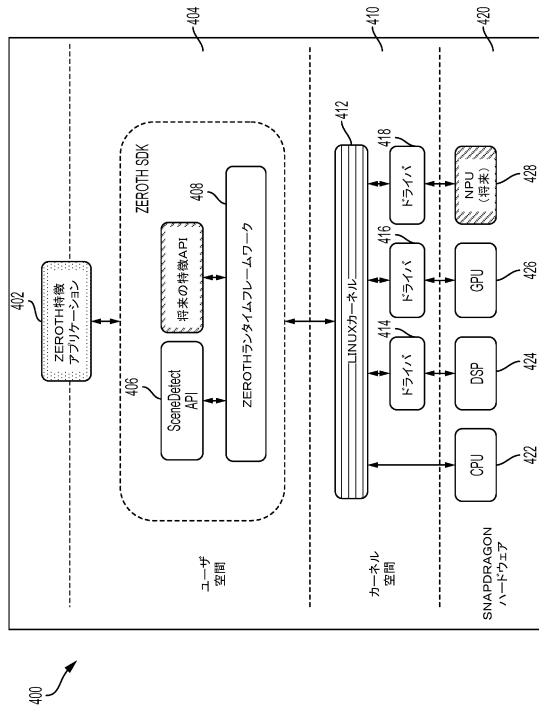


FIG. 4

【図5】

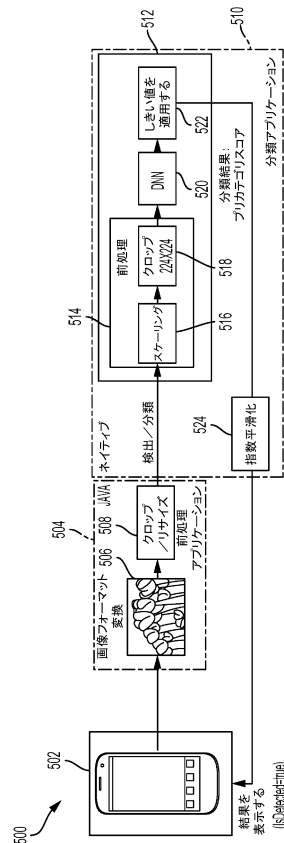


FIG. 5

【 図 6 A 】

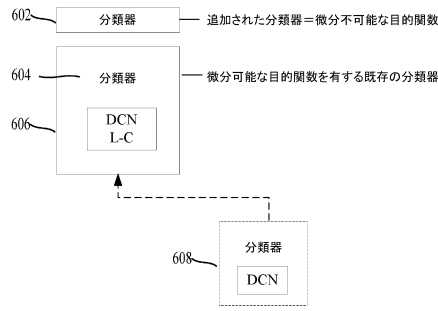


FIG. 6A

【 図 6 B 】

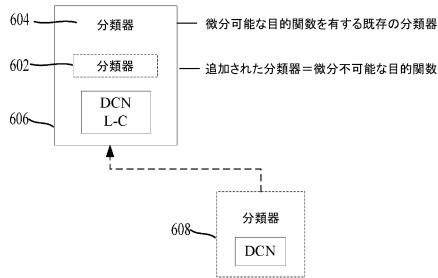


FIG. 6B

【 図 7 】

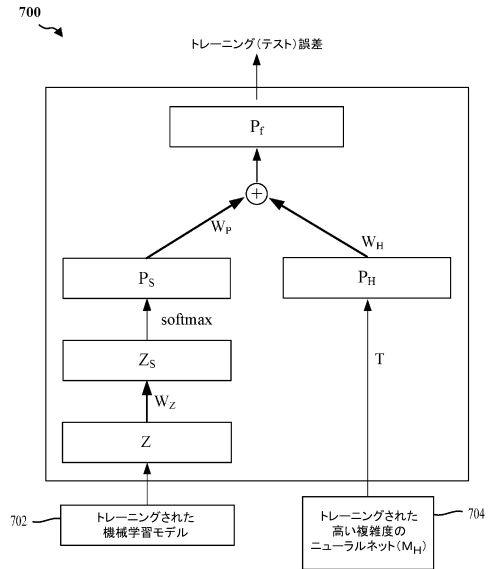


FIG. 7

【 図 8 】

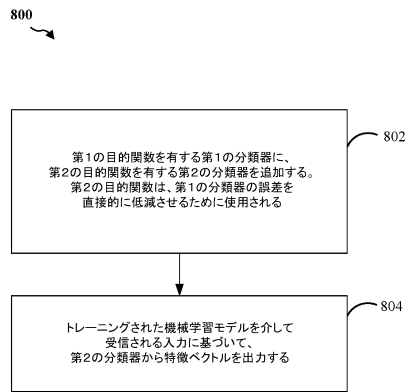


FIG. 8

【 図 9 】

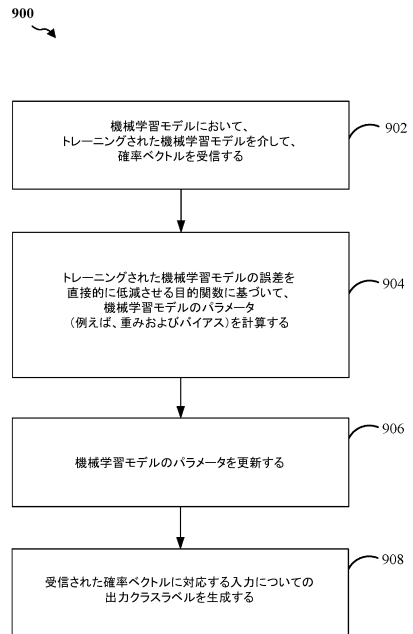


FIG. 9

---

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 タラティ、サチン・スバシュ

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 バルタク、アニケット

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

審査官 多賀 実

(56)参考文献 特開平05 - 290013 (JP, A)

特開平09 - 185394 (JP, A)

特開2011 - 048262 (JP, A)

Geoffrey Hinton et al., "Distilling the Knowledge in a Neural Network", arXiv.org [online], Cornell University, 2015年 3月, arXiv:1503.02531, pp.1-9, [令和2年8月28日検索], インターネット<URL:https://arxiv.org/pdf/1503.02531>

(58)調査した分野(Int.Cl., DB名)

G06N 3/00 - 99/00