

(19)日本国特許庁(JP)

(12)特許公報(B2)

(11)特許番号
特許第7581370号
(P7581370)

(45)発行日 令和6年11月12日(2024.11.12)

(24)登録日 令和6年11月1日(2024.11.1)

(51)国際特許分類 F I
G 0 6 N 3/0495(2023.01) G 0 6 N 3/0495

請求項の数 62 (全59頁)

(21)出願番号	特願2022-562943(P2022-562943)	(73)特許権者	504315750 フラウンホーファー ゲゼルシャフト ツァ フェルデルング デア アンゲヴァ ンテン フォルシュング アインゲトラ ゲナー フェライン
(86)(22)出願日	令和3年4月13日(2021.4.13)	(74)代理人	100133411 弁理士 山本 龍郎
(65)公表番号	特表2023-522886(P2023-522886 A)	(74)代理人	100067677 弁理士 山本 彰司
(43)公表日	令和5年6月1日(2023.6.1)	(72)発明者	ヴィーデマン, ジモン ドイツ連邦共和国 ベルリン 1 0 5 8 7 アインシュタインウーファー 3 7 ハー ハーイー, ハイブリッヒ - ヘルツ - イン 最終頁に続く
(86)国際出願番号	PCT/EP2021/059592		
(87)国際公開番号	WO2021/209469		
(87)国際公開日	令和3年10月21日(2021.10.21)		
審査請求日	令和4年12月14日(2022.12.14)		
(31)優先権主張番号	20169502.0		
(32)優先日	令和2年4月14日(2020.4.14)		
(33)優先権主張国・地域又は機関	欧州特許庁(EP)		

(54)【発明の名称】 ニューラルネットワークパラメーターの表現の改良された概念

(57)【特許請求の範囲】

【請求項1】

NN表現(110)を生成する装置(100)であって、前記NN表現(110)は、NNパラメーター(130)を表すために、量子化パラメーター(142)及び量子化値(152)を含み、前記装置(100)は、

量子化パラメーター(142)から、

前記量子化パラメーター(142)によって導出された被除数と精度パラメーター(145)によって導出された除数との間の除算の剰余に基づく乗数(144)と、

前記除算の商の丸めに基づくビットシフト数(146)と、

が導出可能であるように、NNパラメーター(120)の前記量子化パラメーター(142)及び量子化値(152)を決定(140)することによって、前記NNパラメーター(120)を量子化された値(130)に量子化するように構成され、

それにより、前記NNパラメーター(120)の前記量子化された値(130)は、前記量子化値(152)と前記乗数(144)に依存する因数(148)との間の積であって、前記ビットシフト数(146)に依存するビット数だけビットシフトされた積に対応する、装置。

【請求項2】

NN表現(110)からNNパラメーターを導出する装置であって、

前記NN表現(110)から量子化パラメーター(142)を導出することと、

前記NN表現(110)から量子化値(152)を導出することと、

10

20

前記量子化パラメーター(142)から、

前記量子化パラメーター(142)によって導出された被除数と精度パラメーター(145)によって導出された除数との間の除算の剰余に基づく乗数(144)と、

前記除算の商の丸めに基づくビットシフト数(146)と、
を導出することと、

を行うように構成され、

前記NNパラメーター(130)は、前記量子化値(152)と前記乗数(144)に依存する因数(148)との間の積であって、前記ビットシフト数(146)に依存するビット数だけビットシフトされた積に対応する、装置。

【請求項3】

前記NN表現(110)から前記精度パラメーター(145)を導出するように更に構成されている、請求項2に記載の装置。

【請求項4】

前記NNパラメーター(130)は、

一对のニューロン(10)間のニューロン間活性化フィードフォワード(12)に重み付けする重みパラメーター、

ニューラルネットワーク層のアフィン変換をパラメーター化するバッチノルムパラメーター、及び

所定のニューラルネットワークニューロン(10)のインバウンドニューロン間活性化フィードフォワード(12)の和にバイアスをかけるバイアス、

のうちの1つである、請求項2又は3に記載の装置。

【請求項5】

前記NNパラメーター(130)は、NNの複数のニューロン間活性化フィードフォワード(122)のうちの単一のニューロン間活性化フィードフォワード(12)に関して前記NNをパラメーター化し、前記装置は、前記複数のニューロン間活性化フィードフォワード(122)の各々について、前記NN表現(110)から、対応するNNパラメーター(130)を、

前記複数のニューロン間活性化フィードフォワード(122)の各々(12)について、

前記NN表現(110)から前記それぞれのニューロン間活性化フィードフォワード(12)に関連付けられた関連量子化パラメーター(142)を導出することと、

前記NN表現(110)から前記それぞれのニューロン間活性化フィードフォワード(12)に関連付けられた関連量子化値(152)を導出することと、

前記関連量子化パラメーター(142)から、

前記関連量子化パラメーター(142)によって導出された被除数と、前記それぞれのニューロン間活性化フィードフォワード(12)に関連付けられた関連精度パラメーター(145)によって導出された除数との間の除算の剰余に基づいて、前記それぞれのニューロン間活性化フィードフォワード(12)に関連付けられた関連乗数(144)と、

前記除算の商の丸めに基づいて、前記それぞれのニューロン間活性化フィードフォワード(12)に関連付けられた関連ビットシフト数(146)と、を導出することと、
によって導出するように構成され、

前記それぞれのニューロン間活性化フィードフォワード(12)の前記対応するNNパラメーター(130)は、前記関連量子化値(152)と前記関連乗数(144)に依存する因数(148)との間の積であって、前記関連ビットシフト数(146)に依存するビット数だけビットシフトされた積に対応する、請求項2～4のいずれか一項に記載の装置。

【請求項6】

前記装置は、NNの複数のニューロン間活性化フィードフォワード(122)をニューロン間活性化フィードフォワードのサブグループ(122a、122b)に細分するように構成され、それにより、各サブグループ(122a、122b)が、前記NNの関連する一对のNN層に関連付けられ、前記関連する一对のNN層の間のニューロン間活性化フィ

10

20

30

40

50

ードフォワードを含み、前記関連する一対の層以外の更なる一対のNN層の間のニューロン間活性化フィードフォワードを除外し、2つ以上のサブグループ(122a、122b)が、所定のNN層に関連付けられ、

前記NNパラメーター(130)は、前記NNの前記複数のニューロン間活性化フィードフォワード(122)のうちの単一のニューロン間活性化フィードフォワード(12)に関して前記NNをパラメーター化し、前記装置は、前記複数のニューロン間活性化フィードフォワード(122)の各々について、前記NN表現(110)から、対応するNNパラメーター(130)を、

ニューロン間活性化フィードフォワードの各サブグループ(122a、122b)について、

前記NN表現(110)から前記それぞれのサブグループ(122a、122b)に関連付けられた関連量子化パラメーター(142)を導出することと、

前記関連量子化パラメーター(142)から、

前記関連量子化パラメーター(142)によって導出された被除数と、前記それぞれのサブグループ(122a、122b)に関連付けられた関連精度パラメーター(145)によって導出された除数との間の除算の剰余に基づいて、前記それぞれのサブグループ(122a、122b)に関連付けられた関連乗数(144)と、

前記除算の商の丸めに基づいて、前記それぞれのサブグループ(122a、122b)に関連付けられた関連ビットシフト数(146)と、を導出することと、

前記複数のニューロン間活性化フィードフォワード(122)の各々について、

前記NN表現(110)から前記それぞれのニューロン間活性化フィードフォワード(12)に関連付けられた関連量子化値(152)を導出することと、
によって導出するように構成され、

前記それぞれのニューロン間活性化フィードフォワード(12)の前記対応するNNパラメーター(130)は、前記それぞれのニューロン間活性化フィードフォワード(12)が含まれる前記サブグループ(122a、122b)に関連付けられた、前記関連量子化値(152)と前記関連乗数(144)に依存する因数(148)との間の積であって、前記それぞれのニューロン間活性化フィードフォワード(12)が含まれる前記サブグループ(122a、122b)の前記関連ビットシフト数(146)に依存するビット数だけビットシフトされた積に対応する、請求項2~4のいずれか一項に記載の装置。

【請求項7】

前記関連精度パラメーター(145)は、前記NNにわたって又は各NN層内でグローバルに等しく評価される、請求項5又は6に記載の装置。

【請求項8】

前記NN表現(110)から前記関連精度パラメーター(145)を導出するように構成されている、請求項5~7のいずれか一項に記載の装置。

【請求項9】

前記NN表現(110)から前記関連量子化パラメーター(142)を基準量子化パラメーターとの差分の形態で導出するように構成されている、請求項5~8のいずれか一項に記載の装置。

【請求項10】

前記量子化パラメーター(142)から、前記乗数(144)及び前記ビットシフト数(146)を、

$$mul = k + QP \% k$$

【数1】

$$shift = \lfloor QP/k \rfloor$$

に従って導出するように構成され、式中、mulは、前記乗数(144)であり、shiftは、前記ビットシフト数(146)であり、QPは、前記量子化パラメーター(14

10

20

30

40

50

2) であり、 k は、前記精度パラメータ (145) であり、
【数 2】

[]

は、そのオペランド以下の最大の整数を生成するフロア演算子であり、 $\%$ は、 $x \% y$ に対して

【数 3】

$$x - y \cdot \lfloor x/y \rfloor$$

10

を生成するモジュロ演算子であり、それにより、前記 NN パラメータ (130) は、
 $(mul/k) \cdot 2^{shift \cdot p}$

であり、式中、 P は、前記量子化値 (152) である、請求項 2 ~ 9 のいずれか一項に記載の装置。

【請求項 11】

前記精度パラメータ (145) は、2 の累乗である、請求項 2 ~ 10 のいずれか一項に記載の装置。

【請求項 12】

コンテキスト適応型二値算術復号化の使用によって、又は

前記 NN 表現 (110) から前記量子化パラメータ (142) を表すビットを直接読み出すことによって、又は

前記装置のコンテキスト適応型二値デコーダの等確率バイパスモードを介して、前記 NN 表現 (110) から前記量子化パラメータ (142) を表すビットを導出することによって、

前記 NN 表現 (110) から前記量子化パラメータ (142) を導出するように構成されている、請求項 2 ~ 11 のいずれか一項に記載の装置。

【請求項 13】

二値化方式を使用してピンストリングを非二値化することによって、前記 NN 表現 (110) から前記量子化パラメータ (142) を導出するように構成されている、請求項 2 ~ 12 のいずれか一項に記載の装置。

【請求項 14】

前記二値化方式は、指数ゴロム符号である、請求項 13 に記載の装置。

【請求項 15】

固定小数点表現の形態で前記 NN 表現 (110) から前記量子化パラメータ (142) を導出するように構成されている、請求項 2 ~ 14 のいずれか一項に記載の装置。

【請求項 16】

前記精度パラメータ (145) は、 2^t であり、前記固定小数点表現のビット長は、前記 NN に対して一定になるように設定されるか、又は前記 NN に対して一定である基底ビット長と t との和になるように設定される、請求項 15 に記載の装置。

【請求項 17】

整数値シンタックス要素として前記 NN 表現 (110) から前記量子化パラメータ (142) を導出するように構成されている、請求項 2 ~ 16 のいずれか一項に記載の装置。

【請求項 18】

前記 NN 表現 (110) から前記精度パラメータ (145) を表すビットを直接読み出すことによって、又は前記装置のコンテキスト適応型二値デコーダの等確率バイパスモードを介して前記 NN 表現 (110) から前記精度パラメータ (145) を表すビットを導出することによって、前記 NN 表現 (110) から前記精度パラメータ (145) を導出するように構成されている、請求項 2 ~ 17 のいずれか一項に記載の装置。

50

【請求項 19】

固定小数点表現の形態で前記 NN 表現 (110) から前記量子化値 (152) を導出するように構成されている、請求項 2 ~ 18 のいずれか一項に記載の装置。

【請求項 20】

二値化方式に従ってピンストリングから前記量子化値 (152) を非二値化し、コンテキスト適応算術復号化を使用して前記 NN 表現 (110) から前記ピンストリングのビットを復号化することによって、前記 NN 表現 (110) から前記量子化値 (152) を導出するように構成されている、請求項 2 ~ 19 のいずれか一項に記載の装置。

【請求項 21】

二値化方式に従ってピンストリングから前記量子化値 (152) を非二値化し、コンテキスト適応算術復号化を使用して前記 NN 表現 (110) から前記ピンストリングの第 1 のビットを復号化し、等確率バイパスモードを使用して前記ピンストリングの第 2 のビットを復号化することによって、前記 NN 表現 (110) から前記量子化値 (152) を導出するように構成されている、請求項 2 ~ 20 のいずれか一項に記載の装置。

10

【請求項 22】

NN (20) を使用して推論を実行するデバイス (400) であって、前記デバイス (400) は、

前記 NN (20) をパラメータ化するように構成された NN パラメトライザー (410) であって、前記 NN パラメトライザー (410) は、請求項 2 ~ 21 のいずれか一項に記載の NN 表現 (110) から NN パラメータ (130) を導出する装置 (300) を備える、NN パラメトライザー (410) と、

20

前記 NN (20) を使用して NN 入力 (440) に基づいて推論出力 (430) を計算するように構成された計算ユニット (420) と、
を備える、デバイス。

【請求項 23】

請求項 22 に記載のデバイス (400) であって、

前記 NN パラメトライザー (410) は、
前記装置 (300) を介して、第 1 の NN パラメータ (130) 及び第 2 の NN パラメータ (130) のうちの少なくとも 1 つを導出することであって、それにより、前記第 1 の NN パラメータ (130) は、第 1 の量子化値 (152) と第 1 の因数 (148) との積であって、第 1 のビット数 (146) だけビットシフトされた積に対応し、前記第 2 の NN パラメータ (130) は、第 2 の量子化値 (152) と第 2 の因数 (148) との積であって、第 2 のビット数 (146) だけビットシフトされた積に対応することと、

30

前記第 1 の NN パラメータ (130) の第 1 の量子化値 (152) によって形成され、前記第 1 の乗数 (144) で重み付けされた第 1 の加数と、前記第 1 の NN パラメータ (130) の第 2 の量子化値 (152) によって形成され、前記第 2 の乗数 (144) で重み付けされ、前記第 1 のビット数 (146) 及び前記第 2 のビット数 (146) の差だけビットシフトされた第 2 の加数との間の和を形成することと、前記第 1 の加数及び前記第 2 の加数の前記和を、前記第 1 のビット数 (146) 及び前記第 2 のビット数 (146) のうちの 1 つに依存するビット数だけビットシフトすることと、によって、前記第 1 の NN パラメータ (130) 及び前記第 2 の NN パラメータ (130) に加算を施して、前記 NN (20) の最終 NN パラメータ (130) を生成することと、
を行うように構成されている、デバイス。

40

【請求項 24】

請求項 23 に記載のデバイス (400) であって、

前記第 1 の NN パラメータ (130) は、前記 NN (20) のベース層表現を表し、前記第 2 の NN パラメータ (130) は、前記 NN のエンハンスメント層表現を表すか、又は

前記第 1 の NN パラメータ (130) は、前記 NN (20) の現在の表現を表し、前

50

前記第2のNNパラメーター(130)は、前記現在のNN表現(110)の更新を表すか、又は

前記第1のNNパラメーター(130)は、所定のニューラルネットワークニューロン(10)のインバウンドニューロン間活性化フィードフォワード(12)の和にバイアスをかけるバイアスを表し、前記第2のNNパラメーター(130)は、ニューラルネットワーク層のアフィン変換をパラメーター化するバッチノルムパラメーターを表す、デバイス。

【請求項25】

請求項22～24のいずれか一項に記載のデバイス(400)であって、

前記NNパラメトライザー(410)は、

前記装置(300)を介して、第3のNNパラメーター(130)及び第4のNNパラメーター(130)のうちの少なくとも1つを導出するように構成され、それにより、前記第3のNNパラメーター(130)は、第3の量子化値(152)と第3の因数(148)との積であって、第3のビット数(146)だけビットシフトされた積に対応し、前記第4のNNパラメーター(130)は、第4の量子化値(152)と第4の因数(148)との積であって、第4のビット数(146)だけビットシフトされた積に対応し、

前記計算ユニット(420)は、前記計算を実行する際に、

前記第3のNNパラメーター(130)及び前記第4のNNパラメーター(130)に乗算を施して、前記第3のNNパラメーター(130)の第3の量子化値(152)によって形成される第1の因数と、前記第3の乗数(144)によって形成される第2の因数と、前記第4のNNパラメーター(130)の第4の量子化値(152)によって形成される第3の因数と、前記第4の乗数(144)によって形成される第4の因数との積であって、前記第3のビット数によって形成される第1の加数と前記第4のビット数によって形成される第2の加数とを含む和に対応するビット数だけビットシフトされた積を形成することによって積を生成するように構成される、デバイス。

【請求項26】

請求項25に記載のデバイス(400)であって、

前記第3のNNパラメーター(130)は、第1のNN層(114)の第1のニューロン(101)から第2のNN層(116)の第2のニューロン(102)へのニューロン間活性化フィードフォワード(12)を重み付けする重みパラメーターを表し、前記第4のNNパラメーター(130)は、バッチノルムパラメーターを表す、デバイス。

【請求項27】

請求項26に記載のデバイス(400)であって、前記バッチノルムパラメーターは、前記第2のNN層に対して前記第1のニューロン(101)の活性化フィードフォワード増幅を調整するものである、デバイス。

【請求項28】

請求項22～27のいずれか一項に記載のデバイス(400)であって、

活性化について第5の量子化パラメーター(142)及び第5の量子化値(152)を決定することによって、前記活性化を量子化された値(130)に量子化することによって前記NN入力(440)を量子化し、前記第5の量子化パラメーター(142)から、前記第5の量子化パラメーター(142)によって導出された被除数と前記活性化に関連付けられた精度パラメーター(145)によって導出された除数との間の除算の剰余に基づく第5の乗数(144)と、前記除算の商の丸めに基づく第5のビットシフト数(146)とを導出することにより、前記第5の量子化値(152)と前記第5の乗数(144)に依存する因数(148)との間の積であって、前記第5のビットシフト数(146)に依存する第5のビット数だけビットシフトされた積に対応する前記量子化された値(130)が得られるように更に構成されている、デバイス。

【請求項29】

請求項28に記載のデバイス(400)であって、

前記NNパラメトライザー(410)は、

10

20

30

40

50

前記装置(300)を介して、第6のNNパラメーター(130)を導出することであって、それにより、前記第6のNNパラメーター(130)は、第6の量子化値(152)と第6の因数(148)との積であって、第6のビット数(146)だけビットシフトされた積に対応することと、

前記第6のNNパラメーター(130)及び前記活性化に乗算を施して、前記第6のNNパラメーター(130)の第6の量子化値(152)によって形成される第1の因数と、前記第6の乗数(144)によって形成される第2の因数と、前記第5の量子化値(152)によって形成される第3の因数と、前記第5の乗数(144)によって形成される第4の因数との積であって、前記第6のビット数によって形成される第1の加数と前記第4のビット数(146)によって形成される第2の加数と、を含む和に対応するビット数だけビットシフトされた積を形成することによって積を生成することと、
を行うように構成されている、デバイス。

10

【請求項30】

NN(20)を使用して推論を実行するデバイス(500)であって、前記デバイス(500)は、前記NN(20)を使用してNN入力(440)に基づいて推論出力(430)を計算するように構成され、前記NN(20)は、一対のNN層と、前記一対のNN層のうちの第1のNN層から前記一対のNN層のうちの第2のNN層へのニューロン間活性化フィードフォワード(12)とを含み、前記デバイス(500)は、

前記第1のNN層の前記ニューラルネットワークニューロン(10)の前記活性化から行列 X (532)を形成(530)することと、

20

$s \cdot W' * X$ を計算(540)することであって、式中、 $*$ は、行列乗算を表し、 W' は、 n 及び m N である次元 $n \times m$ の重み行列(544)であり、 s は、長さ n の転置ベクトル(546)であり、 \cdot は、 \cdot の一方の側の行列と \cdot の他方の側の転置ベクトルとの間の列に関するアダマール乗算を示すことと、

によって、前記第1のNN層の前記ニューラルネットワークニューロン(10)の前記活性化に基づいて、前記第2のNN層の前記ニューラルネットワークニューロン(10)の活性化を計算するように構成され、

前記デバイス(500)が、NN表現(110)から W' (544)を導出するように構成されたNNパラメトライザー(410)を備え、前記NNパラメトライザー(410)は、請求項2~21のいずれか一項に記載のNN表現(110)からNNパラメーター(130)を導出する装置(300)を備えている、デバイス。

30

【請求項31】

請求項30に記載のデバイス(500)であって、 n ビット固定小数点演算を使用して前記行列乗算を計算(540)して、内積を生成し、 $m > n$ である m ビット固定小数点演算を使用して前記内積を s (546)と乗算するように構成されている、デバイス。

【請求項32】

請求項30又は31に記載のデバイス(500)であって、 s (546)は、 W' (544)を符号化するためのより高い圧縮及び/又はより高い推論忠実度に関する W' (544)の最適化の結果である、デバイス。

【請求項33】

請求項30~32のいずれか一項に記載のデバイス(500)であって、前記NNパラメトライザー(410)は、 W' (544)に関連するNNパラメーター(130)と比較して異なる量子化パラメーター(142)を使用して、前記NN表現(110)から s (546)を導出するように更に構成されている、デバイス。

40

【請求項34】

NN(20)のバッチノルム演算子(710)のNNパラメーターをNN表現(110)に符号化する装置(600)であって、前記バッチノルム演算子(710)は、

【数4】

50

$$\frac{W * X + b - \mu}{\sqrt{\sigma^2 + \epsilon}} \cdot \gamma + \beta$$

として定義され、式中、

μ 、 σ^2 、 ϵ 、及び β は、バッチノルムパラメータであり、

W は、重み行列であり、

X は、NN層の活性化から導出される入力行列であり、

b は、バイアスを形成する転置ベクトルであり、

ϵ は、ゼロ除算回避のための定数であり、

\cdot は、 \cdot の一方の側の行列と他方の側の転置ベクトルとの間の列に関するアダマール乗算を示し、

$*$ は、行列乗算を示し、

前記装置(600)は、

b 、 μ 、 σ^2 又は ϵ 、及び (610)を受信することと、

【数5】

$$1) \beta' := \beta + \frac{(b - \mu) \cdot \gamma}{\sqrt{\sigma^2 + \epsilon}}$$

$$2) \gamma' := \gamma \cdot \frac{\sqrt{\theta + \epsilon}}{\sqrt{\sigma^2 + \epsilon}}$$

を計算(620)することと、

β' 及び γ' を前記バッチノルム演算子(710)のNNパラメータとして前記NN表現(110)に符号化することであって、前記バッチノルム演算子(710)を

【数6】

$$\frac{W * X + b' - \mu'}{\sqrt{\sigma'^2 + \epsilon'}} \cdot \gamma' + \beta'$$

として定義し、

3) $\sigma'^2 := \sigma^2$

4) $\mu' := 0$

5) $b' := 0$ であり、ここで、 ϵ' は所定のパラメータであることと、

を行うように構成されており、

装置(600)が、 β' 及び γ' を量子化して前記NN表現(110)に符号化する、請求項1に記載の装置(100)を更に備える、装置。

【請求項35】

請求項34に記載の装置(600)であって、前記所定のパラメータは、1又は1 - ϵ' である、装置。

【請求項36】

NNのバッチノルム演算子(710)のNNパラメータをNN表現(110)に符号化する装置(600)であって、前記バッチノルム演算子(710)は、

【数7】

$$\frac{W * X + b - \mu}{\sqrt{\sigma^2}} \cdot \gamma + \beta$$

として定義され、式中、

μ 、 σ^2 、 β 、及び γ は、バッチノルムパラメータであり、

W は、重み行列であり

X は、NN層の活性化から導出される入力行列であり、

b は、バイアスを形成する転置ベクトルであり、

\cdot は、 \cdot の一方の側の行列と他方の側の転置ベクトルとの間の列に関するアダマール乗算を示し、

$*$ は、行列乗算を示し、

前記装置 (600) は、

b 、 μ 、 σ^2 又は β 、 γ 、及び (610) を受信することと、

【数8】

$$1) \beta' := \beta + \frac{(b-\mu) \cdot \gamma}{\sqrt{\sigma^2}}$$

$$2) \gamma' := \gamma \cdot \frac{1}{\sqrt{\sigma^2}}$$

を計算 (620) することと、

β' 及び γ' を前記バッチノルム演算子 (710) のNNパラメータとして前記NN表現 (110) に符号化することであって、前記バッチノルム演算子 (710) を

【数9】

$$\frac{W \cdot X + b' - \mu'}{\sqrt{\sigma'^2}} \cdot \gamma' + \beta'$$

として定義し、

3) $\sigma'^2 := 1$

4) $\mu' := 0$

5) $b' := 0$ であることと、

を行うように構成され、

装置 (600) が、 β' 及び γ' を量子化して前記NN表現 (110) に符号化する、請求項1に記載の装置 (100) を更に備える、装置。

【請求項37】

NNのバッチノルム演算子 (710) のNNパラメータをNN表現 (110) に符号化する装置 (600) であって、前記バッチノルム演算子 (710) は、

【数10】

$$\frac{W \cdot X - \mu}{\sqrt{\sigma^2 + \epsilon}} \cdot \gamma + \beta$$

として定義され、式中、

μ 、 σ^2 、 β 、及び γ は、バッチノルムパラメータであり、

W は、重み行列であり、

X は、NN層の活性化から導出される入力行列であり、

ϵ は、ゼロ除算回避のための定数であり、

\cdot は、 \cdot の一方の側の行列と他方の側の転置ベクトルとの間の列に関するアダマール乗算を示し、

$*$ は、行列乗算を示し、

10

20

30

40

50

前記装置(600)は、
μ、²又は、及び(610)を受信することと、

【数11】

$$1) \beta' := \beta - \frac{\mu * \gamma}{\sqrt{\sigma^2 + \epsilon}}$$

$$2) \gamma' := \gamma \cdot \frac{\sqrt{\theta + \epsilon}}{\sqrt{\sigma^2 + \epsilon}}$$

を計算(620)することと、

10

及びを前記バッチノルム演算子(710)のNNパラメータとして前記NN表現(110)に符号化することであって、前記バッチノルム演算子(710)を

【数12】

$$\frac{W * X - \mu'}{\sqrt{\sigma'^2 + \epsilon}} \cdot \gamma' + \beta'$$

として定義し、

3) ² :=

20

4) μ' := 0であり、ここで、は所定のパラメータであることと、

を行うように構成され、

装置(600)が、及びを量子化して前記NN表現(110)に符号化する、請求項1に記載の装置(100)を更に備える、装置。

【請求項38】

請求項37に記載の装置(600)であって、前記所定のパラメータは、1又は1 - である、装置。

【請求項39】

NNのバッチノルム演算子(710)のNNパラメータをNN表現(110)に符号化する装置(600)であって、前記バッチノルム演算子(710)は、

30

【数13】

$$\frac{W * X - \mu}{\sqrt{\sigma^2}} \cdot \gamma + \beta$$

として定義され、式中、

μ、²、及びは、バッチノルムパラメータであり、

Wは、重み行列であり、

Xは、NN層の活性化から導出される入力行列であり、

40

・は、・の一方の側の行列と他方の側の転置ベクトルとの間の列に関するアダマール乗算を示し、

*は、行列乗算を示し、

前記装置(600)は、

μ、²又は、及び(610)を受信することと、

【数14】

50

$$1) \beta' := \beta - \frac{\mu * \gamma}{\sqrt{\sigma^2}}$$

$$2) \gamma' := \gamma \cdot \frac{1}{\sqrt{\sigma^2}}$$

を計算 (620) することと、

' 及び ' を前記バッチノルム演算子 (710) の NN パラメータとして前記 NN 表現 (110) に符号化することであって、前記バッチノルム演算子 (710) を

【数 15】

$$\frac{W * X - \mu'}{\sqrt{\sigma'^2}} \cdot \gamma' + \beta'$$

として定義し、

$$3) \sigma'^2 := 1$$

$$4) \mu' := 0 \text{ であることと、}$$

を行うように構成され、

装置 (600) が、' 及び ' を量子化して前記 NN 表現 (110) に符号化する、請求項 1 に記載の装置 (100) を更に備える、装置。

【請求項 40】

請求項 34 ~ 39 のいずれか一項に記載の装置 (600) であって、

前記表現において、 σ'^2 の全ての成分が互いに等しいこと、及びその値を示し、及び / 又は

前記表現において、 μ' の全ての成分が互いに等しいこと、及びその値を示し、及び / 又は

前記表現において、存在する場合、 b' の全ての成分が互いに等しいこと、及びその値を示すように更に構成されている、装置。

【請求項 41】

請求項 34 ~ 39 のいずれか一項に記載の装置 (600) であって、2つのバッチノルム符号化モードの間で切り替え可能であるように更に構成され、第1のバッチノルム符号化モードでは、前記装置 (600) は、' 及び ' の前記計算及び前記符号化を実行するように構成され、第2のバッチノルム符号化モードでは、前記装置 (600) は、前記受信された μ 、 σ'^2 又は、 μ 、 σ'^2 、及び、並びに存在する場合、 b を符号化するように構成されている、装置。

【請求項 42】

NN 表現 (110) から NN のバッチノルム演算子 (710) の NN パラメータを復号化する装置 (700) であって、前記バッチノルム演算子 (710) は、

【数 16】

$$\frac{W * X + b - \mu}{\sqrt{\sigma^2 + \epsilon}} \cdot \gamma + \beta$$

として定義され、式中、

μ 、 σ^2 、 ϵ 、及び β は、バッチノルムパラメータであり、

W は、重み行列であり、

X は、NN 層の活性化から導出される入力行列であり、

b は、バイアスを形成する転置ベクトルであり、

ϵ は、ゼロ除算回避のための定数であり、

\cdot は、 \cdot の一方の側の行列と他方の側の転置ベクトルとの間の列に関するアダマール乗

10

20

30

40

50

算を示し、

* は、行列乗算を示し、

前記装置 (7 0 0) は、

前記 NN 表現 (1 1 0) から (7 2 2) 及び (7 2 4) を導出 (7 2 0) することと、

その全ての成分に適用される 1 つのシグナリング (7 3 4) によって、

1) $\sigma^2 :=$

2) $\mu' := 0$

3) $b' := 0$ であり、ここで、 は所定のパラメータであることを推論又は導出 (7 3 0) することと、

を行うように構成され、

前記 NN 表現 (1 1 0) から 及び を導出する、請求項 2 ~ 2 1 のいずれか一項に記載の装置 (3 0 0) を更に備える、装置。

【請求項 4 3】

請求項 4 2 に記載の装置 (7 0 0) であって、前記所定のパラメータは、 1 又は 1 - である、装置。

【請求項 4 4】

NN 表現 (1 1 0) から NN のバッチノルム演算子 (7 1 0) の NN パラメータを復号化する装置 (7 0 0) であって、前記バッチノルム演算子 (7 1 0) は、

【数 1 7】

$$\frac{W * X + b - \mu}{\sqrt{\sigma^2}} \cdot \gamma + \beta$$

として定義され、式中、

μ 、 σ^2 、 及び は、バッチノルムパラメータであり、

W は、重み行列であり、

X は、NN 層の活性化から導出される入力行列であり、

b は、バイアスを形成する転置ベクトルであり、

\cdot は、 \cdot の一方の側の行列と他方の側の転置ベクトルとの間の列に関するアダマール乗算を示し、

* は、行列乗算を示し、

前記装置 (7 0 0) は、

前記 NN 表現 (1 1 0) から (7 2 2) 及び (7 2 4) を導出 (7 2 0) することと、

その全ての成分に適用される 1 つのシグナリング (7 3 4) によって、

1) $\sigma^2 := 1$

2) $\mu := 0$

3) $b := 0$ であることを推論又は導出 (7 3 0) することと、

を行うように構成され、

前記 NN 表現 (1 1 0) から 及び を導出する、請求項 2 ~ 2 1 のいずれか一項に記載の装置 (3 0 0) を更に備える、装置。

【請求項 4 5】

NN 表現 (1 1 0) から NN のバッチノルム演算子 (7 1 0) の NN パラメータを復号化する装置 (7 0 0) であって、前記バッチノルム演算子 (7 1 0) は、

【数 1 8】

10

20

30

40

50

$$\frac{W * X - \mu}{\sqrt{\sigma^2 + \epsilon}} \cdot \gamma + \beta$$

として定義され、式中、

μ 、 σ^2 、 ϵ 、及び β は、バッチノルムパラメーターであり、

W は、重み行列であり、

X は、NN層の活性化から導出される入力行列であり、

\cdot は、 \cdot の一方の側の行列と他方の側の転置ベクトルとの間の列に関するアダマール乗

算を示し、

$*$ は、行列乗算を示し、

前記装置(700)は、

前記NN表現(110)から(722)及び(724)を導出(720)することと、

その全ての成分に適用される1つのシグナリング(734)によって、

1) $\sigma^2 :=$

2) $\mu := 0$ であり、ここで、 ϵ は所定のパラメーターであることを推論又は導出(730)することと、

を行うように構成され、

前記NN表現(110)から及びを導出する、請求項2~21のいずれか一項に記載の装置(300)を更に備える、装置。

【請求項46】

請求項45に記載の装置(700)であって、前記所定のパラメーターは、1又は1-である、装置。

【請求項47】

NN表現(110)からNNのバッチノルム演算子(710)のNNパラメーターを復号化する装置(700)であって、前記バッチノルム演算子(710)は、

【数19】

$$\frac{W * X - \mu}{\sqrt{\sigma^2}} \cdot \gamma + \beta$$

として定義され、式中、

μ 、 σ^2 、 ϵ 、及び β は、バッチノルムパラメーターであり、

W は、重み行列であり、

X は、NN層の活性化から導出される入力行列であり、

\cdot は、 \cdot の一方の側の行列と他方の側の転置ベクトルとの間の列に関するアダマール乗算を示し、

$*$ は、行列乗算を示し、

前記装置(700)は、

前記NN表現(110)から(722)及び(724)を導出(720)することと、

その全ての成分に適用される1つのシグナリング(734)によって、

1) $\sigma^2 := 1$

2) $\mu := 0$ であることを推論又は導出(730)することと、

を行うように構成され、

前記NN表現(110)から及びを導出する、請求項2~21のいずれか一項に記載

10

20

30

40

50

載の装置(300)を更に備える、装置。

【請求項48】

請求項42～47のいずれか一項に記載の装置(700)であって、

前記表現から、 μ^2 の全ての成分が互いに等しいこと、及びその値を導出し、及び/又は前記表現から、 μ' の全ての成分が互いに等しいこと、及びその値を導出し、及び/又は存在する場合、前記表現から、 b' の全ての成分が互いに等しいこと、及びその値を導出するように更に構成されている、装置。

【請求項49】

請求項42～48のいずれか一項に記載の装置(700)であって、2つのバッチノルム符号化モード間で切り替え可能であるように更に構成され、第1のバッチノルム符号化モードでは、前記装置(700)は、前記導出することと、前記推論又は導出することと、を実行するように構成され、第2の第1のバッチノルム符号化モードでは、前記装置(700)は、 μ 、 μ^2 又は、 μ' 、及び、 b 、並びに存在する場合、 b を前記表現から復号化するように構成されている、装置。

10

【請求項50】

NN表現を生成する方法であって、前記NN表現は、NNパラメータを表すために、量子化パラメータ及び量子化値を含み、

量子化パラメータから、

前記量子化パラメータによって導出された被除数と精度パラメータによって導出された除数との間の除算の剰余に基づく乗数と、

前記除算の商の丸めに基づくビットシフト数と、
が導出可能であるように、NNパラメータの前記量子化パラメータ及び量子化値を決定することによって、前記NNパラメータを量子化された値に量子化することによって、それにより、前記NNパラメータの前記量子化された値は、前記量子化値と前記乗数に依存する因数との積であって、前記ビットシフト数に依存するビット数だけビットシフトされた積に対応することを含み、方法。

20

【請求項51】

NN表現からNNパラメータを導出する方法であって、

前記NN表現から量子化パラメータを導出することと、

前記NN表現から量子化値を導出することと、

前記量子化パラメータから、

前記量子化パラメータによって導出された被除数と精度パラメータによって導出された除数との間の除算の剰余に基づく乗数と、

前記除算の商の丸めに基づくビットシフト数と、

を導出することと、

を含み、

前記NNパラメータは、前記量子化値と前記乗数に依存する因数との積であって、前記ビットシフト数に依存するビット数だけビットシフトされた積に対応する、方法。

30

【請求項52】

NNを使用して推論を実行する方法であって、

NN表現からNNパラメータを導出するために請求項51に記載の方法を使用して、前記NNをパラメータ化することと、

前記NNを使用して、NN入力に基づいて推論出力を計算することと、

を含む、方法。

40

【請求項53】

NNを使用して推論を実行する方法であって、前記NNを使用してNN入力に基づいて推論出力を計算することを含み、前記NNは、一対のNN層と、前記一対のNN層のうちの第1のNN層から前記一対のNN層のうちの第2のNN層へのニューロン間活性化フィードフォワードとを含み、前記方法は、

50

前記第 1 の NN 層の前記ニューラルネットワークニューロンの前記活性化から行列 X を形成することと、

$s \cdot W' * X$ を計算することであって、式中、* は、行列乗算を表し、W' は、n 及び m N である次元 n x m の重み行列であり、s は、長さ n の転置ベクトルであり、・ は、・ の一方の側の行列と・ の他方の側の転置ベクトルとの間の列に関するアダマール乗算を示すことと、

よって、前記第 1 の NN 層の前記ニューラルネットワークニューロンの活性化に基づいて、前記第 2 の NN 層の前記ニューラルネットワークニューロンの活性化を計算することを含み、

前記方法が、請求項 5 1 に記載の NN 表現 (1 1 0) から NN パラメーターを導出する方法を用いて、NN 表現から W' を導出する、方法。

10

【請求項 5 4】

NN のバッチノルム演算子の NN パラメーターを NN 表現に符号化する方法であって、前記バッチノルム演算子は、

【数 2 0】

$$\frac{W * X + b - \mu}{\sqrt{\sigma^2 + \epsilon}} \cdot \gamma + \beta$$

20

として定義され、式中、

μ 、 σ^2 、及び ϵ は、バッチノルムパラメーターであり、

W は、重み行列であり、

X は、NN 層の活性化から導出される入力行列であり、

b は、バイアスを形成する転置ベクトルであり、

ϵ は、ゼロ除算回避のための定数であり、

・ は、・ の一方の側の行列と他方の側の転置ベクトルとの間の列に関するアダマール乗算を示し、

* は、行列乗算を示し、

前記方法は、

30

b 、 μ 、 σ^2 又は ϵ 、及び β を受信することと、

【数 2 1】

$$1) \beta' := \beta + \frac{(b - \mu) \cdot \gamma}{\sqrt{\sigma^2 + \epsilon}}$$

$$2) \gamma' := \gamma \cdot \frac{\sqrt{\sigma^2 + \epsilon}}{\sqrt{\sigma^2 + \epsilon}} d$$

を計算することと、

β' 及び γ' を前記バッチノルム演算子の NN パラメーターとして前記 NN 表現に符号化することであって、前記バッチノルム演算子を

40

【数 2 2】

$$\frac{W * X + b' - \mu'}{\sqrt{\sigma'^2 + \epsilon}} \cdot \gamma' + \beta'$$

として定義し、

3) $\sigma'^2 :=$

50

4) $\mu' := 0$

5) $b' := 0$ であり、ここで、 σ は所定のパラメータであることと、を含み、前記方法が、 μ' 及び γ' を量子化して前記 NN 表現に符号化する、請求項 50 に記載の方法を更に備える、方法。

【請求項 55】

NN のバッチノルム演算子の NN パラメータを NN 表現に符号化する方法であって、前記バッチノルム演算子は、

【数 23】

$$\frac{W * X + b - \mu}{\sqrt{\sigma^2}} \cdot \gamma + \beta$$

10

として定義され、式中、

μ 、 σ^2 、 γ 、及び β は、バッチノルムパラメータであり、

W は、重み行列であり、

X は、NN 層の活性化から導出される入力行列であり、

b は、バイアスを形成する転置ベクトルであり、

\cdot は、 \cdot の一方の側の行列と他方の側の転置ベクトルとの間の列に関するアダマール乗算を示し、

20

$*$ は、行列乗算を示し、

前記方法は、

b 、 μ 、 σ^2 又は γ 、 β 、及び σ を受信することと、

【数 24】

$$1) \beta' := \beta + \frac{(b - \mu) * \gamma}{\sqrt{\sigma^2}}$$

$$2) \gamma' := \gamma \cdot \frac{1}{\sqrt{\sigma^2}}$$

30

を計算することと、

μ' 及び γ' を前記バッチノルム演算子の NN パラメータとして前記 NN 表現に符号化することであって、前記バッチノルム演算子を

【数 25】

$$\frac{W * X + b' - \mu'}{\sqrt{\sigma'^2}} \cdot \gamma' + \beta'$$

として定義し、

40

3) $\sigma'^2 := 1$

4) $\mu' := 0$

5) $b' := 0$ であることと、を含み、

前記方法が、 μ' 及び γ' を量子化して前記 NN 表現に符号化する、請求項 50 に記載の方法を更に備える、方法。

【請求項 56】

NN のバッチノルム演算子の NN パラメータを NN 表現に符号化する方法であって、前記バッチノルム演算子は、

【数 26】

50

$$\frac{W * X - \mu}{\sqrt{\sigma^2 + \epsilon}} \cdot \gamma + \beta$$

として定義され、式中、

μ 、 σ^2 、 ϵ 、及び β は、バッチノルムパラメーターであり、

W は、重み行列であり、

X は、NN層の活性化から導出される入力行列であり、

ϵ は、ゼロ除算回避のための定数であり、

\cdot は、 \cdot の一方の側の行列と他方の側の転置ベクトルとの間の列に関するアダマール乗算を示し、

$*$ は、行列乗算を示し、

前記方法は、

μ 、 σ^2 又は ϵ 、 β 、及び γ を受信することと、

【数 2 7】

$$1) \beta' := \beta - \frac{\mu * \gamma}{\sqrt{\sigma^2 + \epsilon}}$$

$$2) \gamma' := \gamma \cdot \frac{\sqrt{\sigma^2 + \epsilon}}{\sqrt{\sigma^2 + \epsilon}}$$

を計算することと、

β' 及び γ' を前記バッチノルム演算子のNNパラメーターとして前記NN表現に符号化することであって、前記バッチノルム演算子を

【数 2 8】

$$\frac{W * X - \mu'}{\sqrt{\sigma'^2 + \epsilon}} \cdot \gamma' + \beta'$$

として定義し、

3) $\sigma'^2 := \sigma^2$

4) $\mu' := 0$ であり、ここで、 ϵ は所定のパラメーターであることと、を含み、

前記方法が、 β' 及び γ' を量子化して前記NN表現に符号化する、請求項 5 0 に記載の方法を更に備える、方法。

【請求項 5 7】

NNのバッチノルム演算子のNNパラメーターをNN表現に符号化する方法であって、前記バッチノルム演算子は、

【数 2 9】

$$\frac{W * X - \mu}{\sqrt{\sigma^2}} \cdot \gamma + \beta$$

として定義され、式中、

μ 、 σ^2 、 ϵ 、及び β は、バッチノルムパラメーターであり、

W は、重み行列であり、

X は、NN層の活性化から導出される入力行列であり、

\cdot は、 \cdot の一方の側の行列と他方の側の転置ベクトルとの間の列に関するアダマール乗

10

20

30

40

50

算を示し、

* は、行列乗算を示し、

前記方法は、

μ 、 σ^2 又は σ 、及び γ を受信することと、

【数 3 0】

$$1) \beta' := \beta - \frac{\mu * \gamma}{\sqrt{\sigma^2}}$$

$$2) \gamma' := \gamma \cdot \frac{1}{\sqrt{\sigma^2}}$$

10

s を計算することと、

' μ ' 及び ' γ ' を前記バッチノルム演算子の NN パラメータとして前記 NN 表現に符号化することと、前記バッチノルム演算子を

【数 3 1】

$$\frac{W * X - \mu'}{\sqrt{\sigma'^2}} \cdot \gamma' + \beta'$$

20

として定義し、

3) $\sigma'^2 := 1$

4) $\mu' := 0$ であることと、を含み、

前記方法が、' μ ' 及び ' γ ' を量子化して前記 NN 表現に符号化する、請求項 5 0 に記載の方法を更に備える、方法。

【請求項 5 8】

NN 表現から NN のバッチノルム演算子の NN パラメータを復号化する方法と、前記バッチノルム演算子は、

【数 3 2】

$$\frac{W * X + b - \mu}{\sqrt{\sigma^2 + \epsilon}} \cdot \gamma + \beta$$

30

として定義され、式中、

μ 、 σ^2 、 ϵ 、及び γ は、バッチノルムパラメータであり、

W は、重み行列であり、

X は、NN 層の活性化から導出される入力行列であり、

b は、バイアスを形成する転置ベクトルであり、

ϵ は、ゼロ除算回避のための定数であり、

\cdot は、 \cdot の一方の側の行列と他方の側の転置ベクトルとの間の列に関するアダマール乗算を示し、

40

* は、行列乗算を示し、

前記方法は、

前記 NN 表現から ' μ ' 及び ' γ ' を導出することと、

その全ての成分に適用される 1 つのシグナリングによって、

1) $\sigma'^2 := 1$

2) $\mu' := 0$

3) $b' := 0$ であり、ここで、' ϵ ' は所定のパラメータであることを推論又は導出することと、を含み、

50

前記NN表現から σ^2 及び μ を導出する、請求項51に記載の方法を更に備える、方法。

【請求項59】

NN表現からNNのバッチノルム演算子のNNパラメーターを復号化する方法であって、前記バッチノルム演算子は、

【数33】

$$\frac{W * X + b - \mu}{\sqrt{\sigma^2}} \cdot \gamma + \beta$$

10

として定義され、式中、

μ 、 σ^2 、 γ 、及び β は、バッチノルムパラメーターであり、

Wは、重み行列であり、

Xは、NN層の活性化から導出される入力行列であり、

bは、バイアスを形成する転置ベクトルであり、

\cdot は、 \cdot の一方の側の行列と他方の側の転置ベクトルとの間の列に関するアダマール乗算を示し、

* は、行列乗算を示し、

前記方法は、

前記NN表現から σ^2 及び μ を導出することと、

20

その全ての成分に適用される1つのシグナリングによって、

1) $\sigma^2 := 1$

2) $\mu := 0$

3) $b := 0$ であることを推論又は導出することと、を含み、

前記NN表現から σ^2 及び μ を導出する、請求項51に記載の方法を更に備える、方法。

【請求項60】

NN表現からNNのバッチノルム演算子のNNパラメーターを復号化する方法であって、前記バッチノルム演算子は、

【数34】

30

$$\frac{W * X - \mu}{\sqrt{\sigma^2 + \epsilon}} \cdot \gamma + \beta$$

として定義され、式中、

μ 、 σ^2 、 γ 、及び β は、バッチノルムパラメーターであり、

Wは、重み行列であり、

Xは、NN層の活性化から導出される入力行列であり、

ϵ は、ゼロ除算回避のための定数であり、

\cdot は、 \cdot の一方の側の行列と他方の側の転置ベクトルとの間の列に関するアダマール乗算を示し、

40

* は、行列乗算を示し、

前記方法は、

前記NN表現から σ^2 及び μ を導出することと、

その全ての成分に適用される1つのシグナリングによって、

1) $\sigma^2 :=$

2) $\mu := 0$ であり、ここで、 ϵ は所定のパラメーターであることを推論又は導出することと、を含み、

前記NN表現から σ^2 及び μ を導出する、請求項51に記載の方法を更に備える、方法。

【請求項61】

50

NN表現からNNのバッチノルム演算子のNNパラメーターを復号化する方法であって、前記バッチノルム演算子は、

【数 3 5】

$$\frac{W * X - \mu}{\sqrt{\sigma^2}} \cdot \gamma + \beta$$

として定義され、式中、

μ 、 σ^2 、 γ 、及び β は、バッチノルムパラメーターであり、

10

Wは、重み行列であり、

Xは、NN層の活性化から導出される入力行列であり、

\cdot は、 \cdot の一方の側の行列と他方の側の転置ベクトルとの間の列に関するアダマール乗算を示し、

* は、行列乗算を示し、

前記方法は、

前記NN表現から μ 及び σ^2 を導出することと、

その全ての成分に適用される1つのシグナリングによって、

1) $\sigma^2 := 1$

2) $\mu := 0$ であることを推論又は導出することと、を含み、

20

前記NN表現から μ 及び σ^2 を導出する、請求項51に記載の方法を更に備える、方法。

【請求項62】

コンピュータプログラムがコンピュータで実行される時、請求項50～61のいずれか一項に記載の方法を実行するコンピュータプログラム。

【発明の詳細な説明】

【技術分野】

【0001】

本発明による実施形態は、ニューラルネットワークパラメーターの表現の改善された概念を使用して、ニューラルネットワークパラメーターを符号化又は復号化する装置及び方法に関する。推論及び/又は記憶ビットレート最適化に関する改善を達成することができる。

30

【背景技術】

【0002】

ニューラルネットワークは、その最も基本的な形態において、一連のアフィン変換とそれ続く要素ごとの非線形関数とを構成する。図1に示すように、それらは有向非巡回グラフとして表すことができる。各ノードは、エッジのそれぞれの重み値との乗算によって次のノードに順方向伝搬される特定の値を伴う。次に、全ての入力値が単純に集約される。

【0003】

図1は、フィードフォワードニューラルネットワークのグラフ表現の例を示している。具体的には、この2層ニューラルネットワークは、4次元入力ベクトルを実線に写像する非線形関数である。

40

【0004】

数学的には、図1のニューラルネットワークは次の方法で出力を算出する。

$$\text{output} = L_2(L_1(\text{input}))$$

ここで、

$$L_i(X) = N_i(B_i(X))$$

であり、式中、 B_i は層iのアフィン変換であり、 N_i は層iの何らかの非線形関数である。

【0005】

バイアス層

50

いわゆる「バイアス層」の場合、 B_i は、層*i*に関連する重みパラメーター（エッジ重み） W_i と層*i*の入力 X_i との行列乗算と、それに続くバイアス b_i との和である。

$$B_i(X) = W_i * X_i + b_i$$

W_i は、次元 $n_i \times k_i$ を有する重み行列であり、 X_i は、次元 $k_i \times m_i$ を有する入力行列である。バイアス b_i は、長さ n_i の転置ベクトルである。演算子 $*$ は、行列乗算を示すものとする。バイアス b_i との和は、行列の列に対する要素ごとの演算である。より正確には、 $W_i * X_i + b_i$ は、 b_i が $W_i * X_i$ の各列に追加されることを意味する。

【0006】

いわゆる畳み込み層は、非特許文献1に記載されているように、それらを行列 - 行列積としてキャストすることによって使用することもできる。

【0007】

以降、所与の入力から出力を算出する手順を推論と称する。また、中間結果を、隠れ層又は隠れ活性化値と称し、これは、例えば上記の第1の内積 + 非線形性の算出等、線形変換 + 要素ごとの非線形性を構成する。

【0008】

通常、ニューラルネットワークは、数百万のパラメーターを含むため、その表現のために数百メガバイトを必要とし得る。したがって、その推論手順には、大きな行列間の多くの内積演算の計算が含まれるため、その実行には高い計算リソースが必要となる。したがって、これらの内積を実行する複雑性を低減することが非常に重要である。

【0009】

バッチノルム層

ニューラルネットワーク層のアフィン変換のより洗練された変形例として、以下のような、いわゆるバイアス及びバッチノルム演算が挙げられる。

式1：

【数1】

$$BN(X) = \frac{B(X) - \mu}{\sqrt{\sigma^2 + \epsilon}} \cdot \gamma + \beta = \frac{W * X + b - \mu}{\sqrt{\sigma^2 + \epsilon}} \cdot \gamma + \beta$$

ここで、 μ 、 σ^2 、 ϵ 、及び γ は、バッチノルムパラメーターを示す。なお、層インデックス*i*はここでは無視する。 W は、次元 $n \times k$ を有する重み行列であり、 X は、次元 $k \times m$ を有する入力行列である。バイアス b 及びバッチノルムパラメーター μ 、 σ^2 、 ϵ 、及び γ は、長さ n の転置ベクトルである。演算子 $*$ は、行列乗算を示す。ベクトルを有する行列に対する他の全ての演算（加算、乗算、除算）は、行列の列に対する要素ごとの演算であることに留意されたい。例えば、 $X \cdot$ は、 X の各列が γ と要素ごとに乗算されることを意味する。 ϵ は、0による除算を避けるために必要な小さなスカラー数（0.001等）である。ただし、0であってもよい。

【0010】

b の全てのベクトル要素が0に等しい場合、式1はバッチノルム層を指す。

対照的に、 μ 及び σ^2 の全てのベクトル要素が0に設定され、 ϵ 及び γ の全ての要素が1に設定される場合、バッチノルムのない層（バイアスのみ）が処理される。

【0011】

パラメーターの効率的な表現

パラメーター W 、 b 、 μ 、 σ^2 、 ϵ 、及び γ は、集合的に層のパラメーターを示すものとする。それらは通常、ビットストリーム中でシグナリングされる必要がある。例えば、それらは32ビットの浮動小数点数として表すことができ、又は整数表現に量子化することができる。 ϵ は通常、ビットストリーム中でシグナリングされないことに留意されたい。

【0012】

10

20

30

40

50

かかるパラメータを符号化するための特に効率的な手法では、各値がいわゆる量子化ステップサイズ値の整数倍として表される均一再構成量子化器が用いられる。対応する浮動小数点数は、整数に、通常は単一の浮動小数点数である量子化ステップサイズを乗算することによって再構築することができる。しかしながら、ニューラルネットワーク推論のための効率的な実施態様（すなわち、入力に対するニューラルネットワークの出力の算出）では、可能な限り整数演算が用いられる。したがって、パラメータを浮動小数点表現に再構成する必要があることは望ましくない場合がある。

【先行技術文献】

【非特許文献】

【0013】

【文献】“cuDNN: Efficient Primitives for Deep Learning” (Sharan Chetlur, et al.; arXiv: 1410.0759, 2014)

【発明の概要】

【発明が解決しようとする課題】

【0014】

したがって、ニューラルネットワークパラメータの効率的な符号化及び/又は復号化をサポートするために、かかるパラメータの表現の概念を改善することが望まれている。ニューラルネットワークパラメータが符号化されるビットストリームを低減して、シグナル化コストを低減することが望ましい場合がある。加えて、又は代替として、ニューラルネットワーク推論を改善するために計算リソースの複雑性を低減することが望ましい場合があり、例えば、ニューラルネットワーク推論のための効率的な実施態様を達成することが望ましい場合がある。

【課題を解決するための手段】

【0015】

これは、本出願の独立請求項の主題によって達成される。

【0016】

本発明による更なる実施の形態は、本出願の従属請求項の主題によって定義される。

本発明の第1の態様によれば、本出願の発明者らは、ニューラルネットワーク(NN)表現が直面する1つの問題が、ニューラルネットワークが数百万のパラメータを含み、したがってその表現のために数百メガバイトを必要とし得るという事実から生じることを認識した。したがって、その推論手順には、大きな行列間の多くの内積演算の計算が含まれるため、その実行には高い計算リソースが必要となる。本出願の第1の態様によれば、この困難は、NNパラメータの量子化を使用することによって克服され、浮動小数点演算をほとんど又は全く用いずに推論が可能となる。本発明者らは、乗数及びビットシフト数を導出することができる量子化パラメータを決定することが有利であることを見出した。これは、ビットレートに関して、32ビット浮動小数点値の代わりに量子化パラメータ及び量子化値(quantization value)のみをシグナリングすることが効率的であるという着想に基づいている。NNパラメータの量子化された値(quantized value)は、乗数、ビットシフト数、及び量子化値を使用して算出することができるため、浮動小数点領域の代わりに整数領域において、計算、例えばNNパラメータの和及び/又はNNパラメータとベクトルとの乗算を実行することが可能である。したがって、提示するNN表現を用いて、推論の効率的な計算を達成することができる。

【0017】

したがって、本出願の第1の態様によれば、NN表現、例えばデータストリームを生成する装置は、量子化パラメータから乗数及びビットシフト数が導出可能であるように、NNパラメータの量子化パラメータ及び量子化値を決定することによって、NNパラメータを量子化された値に量子化するように構成される。生成されたNN表現は、NN表現、例えばデータストリームからNNパラメータ、例えばNNパラメータの量子化された値を導出する装置によって読み出され/復号化され得る。NNパラメータを導出する装置は、NN表現から量子化パラメータ及び量子化値を導出し、量子化パラメータ

10

20

30

40

50

ーから乗数及びビットシフト数を導出するように構成される。乗数は、量子化パラメータによって導出された被除数と精度パラメータによって導出された除数との間の除算の剰余に基づいて量子化パラメータから導出可能であり、例えば、精度パラメータは、デフォルト値に設定されてもよく、又は、自然数若しくは2の累乗等の精度パラメータの複数の異なる整数値が、NN全体に対して、又は各層等のNNの各部分に対して装置によってテストされてもよく、量子化誤差及びビットレートの観点から、そのラグランジュ和に関して最良のものがテストされ、精度パラメータとして最良の値を取得し、NN表現でこの選択がシグナリングされる。ビットシフト数は、除算の商の丸めに基づいて量子化パラメータから導出可能である。NNパラメータを導出する装置の場合のNNパラメータ、又はNN表現を生成する装置の場合のNNパラメータの量子化された値は、
(例えば、少なくとも、シフトの場合に符号の別個の処理を伴う量子化された値の絶対値に関して、又はさらに、積、その因数、及びシフトに関してそれぞれ2の補数表現及び2の補数演算を使用する場合等、絶対値と符号の両方に関して)量子化値と乗数に依存する因数との間の積であって、ビットシフト数に依存するビット数だけビットシフトされた積に対応する。デジタルデータは、上述したように、NNパラメータを表すために、量子化パラメータ及び量子化値を含むNN表現を定義することができる。

10

【0018】

NNパラメータを導出する装置によって導出されるNNパラメータは、NN表現を生成する装置によって生成されるNNパラメータの量子化された値に対応することに留意されたい。これは、NNパラメータを導出する装置には、元のNNパラメータが見えないため、NNパラメータを導出する装置から見て、NNパラメータの量子化された値をNNパラメータとみなすためである。

20

【0019】

一実施の形態は、NNを使用して推論を実行するデバイスに関し、該デバイスは、NNをパラメータ化するように構成されたNNパラメライザーを備える。NNパラメライザーは、上述したように、NN表現からNNパラメータを導出する装置を備える。加えて、デバイスは、NNを使用してNN入力に基づいて推論出力を計算するように構成された計算ユニットを備える。上述したように、NNパラメータは、乗数、ビットシフト数及び量子化値に基づいて導出することができ、そのため、浮動小数点領域の代わりに整数領域において、計算、例えばNNパラメータの和及び/又はNNパラメータとベクトルとの乗算を実行することが可能である。したがって、推論の効率的な計算が、デバイスによって達成され得る。

30

【0020】

本発明の第2の態様によれば、本出願の発明者らは、ニューラルネットワーク(NN)を使用して推論を実行するときに遭遇する1つの問題が、推論に使用される重み行列が量子化誤差を有する可能性があり、そのため、低いレベルの精度しか達成されないという事実起因することを認識した。本出願の第1の態様によれば、この困難は、重み行列 W' の各列と要素ごとに乗算される転置ベクトル s 、例えばスケールリングファクターを使用することによって克服される。本発明者らは、算術符号化方法が、重み行列のスケールリングを使用することによってより高い符号化利得をもたらすこと、及び/又は重み行列のスケールリングがニューラルネットワーク性能結果を増加させる、例えば、より高い精度を達成することを見出した。これは、量子化誤差を低減し、それにより量子化されたニューラルネットワークの予測性能を向上させるために、例えば重み行列、例えば量子化された重み行列に依存して、転置ベクトル s を効率的に適応させることができるという着想に基づいている。さらに、本発明者らは、重みパラメータを転置ベクトル s と重み行列 W' との合成として因数分解することで、両者を独立に量子化でき、例えば転置ベクトル s と重み行列 W' の量子化には、様々な量子化パラメータを使用できるため、表現の効率を高めることができることを見出した。これは、性能の観点から有益であるが、ハードウェア効率の観点からも有益である。

40

【0021】

50

したがって、本出願の第2の態様によれば、NNを使用して推論を実行するデバイスは、NNを使用してNN入力に基づいて推論出力を計算するように構成される。NNは、一对のNN層と、一对のNN層のうちの第1のNN層から一对のNN層のうちの第2のNN層へのニューロン間活性化フィードフォワードとを含む。デバイスは、第1のNN層のニューラルネットワークニューロンの活性化から行列Xを形成し、 $s \cdot W' * X$ を計算することによって、第1のNN層のニューラルネットワークニューロンの活性化に基づいて第2のNN層のニューラルネットワークニューロンの活性化を計算するように構成される。演算子*は、行列乗算を示し、W'は、n及びm Nである次元n x mの重み行列であり、sは、長さnの転置ベクトルであり、演算子・は、・の一方の側の行列と・の他方の側の転置ベクトルとの間の列に関するアダマール乗算を示す。

10

【0022】

本発明の第3の態様によれば、本出願の発明者らは、バッチノルム層を使用するとき遭遇する1つの問題が、バッチノルム演算子のバッチノルムパラメーター/要素が通常は浮動小数点表現であるという事実から生じることを認識した。しかしながら、ニューラルネットワーク推論のための効率的な実施態様(すなわち、入力に対するニューラルネットワークの出力の算出)では、可能な限り整数演算が用いられる。この困難は、所定の定数値をバッチノルムパラメーター/要素に、例えばb及びμ及び σ^2 又は ϵ に割り当てることによって克服される。本発明者らは、バッチノルムパラメーター/要素が所定の定数値を有する場合、それらを極めて効率的に圧縮することができることを見出した。これは、全ての要素/パラメーターが所定の定数値を有するかを示す単一のフラグの使用を可能にし、それにより、それらを所定の定数値に設定することができるという着想に基づいている。加えて、バッチノルム演算子の結果は、所定の定数値を使用することによって変更されないことを見出した。

20

【0023】

したがって、本出願の第3の態様によれば、第1の実施の形態は、NNのバッチノルム演算子のNNパラメーターをNN表現に符号化する装置に関する。バッチノルム演算子は、

$$\frac{W * X + b - \mu}{\sqrt{\sigma^2 + \epsilon}} \cdot \gamma + \beta$$

30

として定義され、式中、

μ、 σ^2 、 ϵ 、及びβは、バッチノルムパラメーター、例えば、各出力ノードについて1つの成分を含む転置ベクトルであり、

Wは、重み行列であり、例えば、その各行は1つの出力ノードに対するものであり、それぞれの行の各成分はXの1つの行に関連付けられており、

Xは、NN層の活性化から導出される入力行列であり、

bは、バイアスを形成する転置ベクトル、例えば、各出力ノードに対して1つの成分を含む転置ベクトルであり、

40

σは、ゼロ除算回避のための定数であり、

・は、・の一方の側の行列と他方の側の転置ベクトルとの間の列に関するアダマール乗算を示し、

*は、行列乗算を示す。

装置は、b及びμ及びσ²及びε及びβを受信し、

【数3】

50

$$\beta' := \beta + \frac{(b - \mu) * \gamma}{\sqrt{\sigma^2 + \epsilon}}$$

及び

【数 4】

$$\gamma' := \gamma \cdot \frac{\sqrt{\theta + \epsilon}}{\sqrt{\sigma^2 + \epsilon}}$$

10

を計算するように構成される。

加えて、装置は、 β' 及び γ' を NN 表現に符号化するように構成され、例えば、出力ノードごとに 1 つの成分を含む転置ベクトルもバッチノルム演算子の NN パラメータとして符号化するように構成され、バッチノルム演算子を

【数 5】

$$\frac{W * X + b' - \mu'}{\sqrt{\sigma'^2 + \epsilon}} \cdot \gamma' + \beta'$$

20

として定義し、

$\sigma'^2 := \sigma^2$ 、 $\mu' := 0$ 、 $b' := 0$ であり、ここで、 σ^2 は所定のパラメータである。

【0024】

NN 表現から NN のバッチノルム演算子の NN パラメータを復号化する並列装置は、NN 表現から β' 及び γ' を導出し、その全ての成分に適用される 1 つのシグナリングによって、 $\sigma'^2 := \sigma^2$ 及び $\mu' := 0$ 及び $b' := 0$ を推論又は導出するように構成され、ここで、 σ^2 は所定のパラメータである。装置は、例えば、1 つのシグナリング、例えばフラグを読み出し、そこから $\sigma'^2 := \sigma^2$ 及び $\mu' := 0$ 及び $b' := 0$ を推論又は導出するように構成される。バッチノルム演算子は、第 3 の態様の第 1 の実施の形態に関して上記で説明したように定義される。

30

【0025】

したがって、本出願の第 3 の態様によれば、第 2 の実施の形態は、NN のバッチノルム演算子の NN パラメータを NN 表現に符号化する装置に関する。バッチノルム演算子は、

【数 6】

$$\frac{W * X + b - \mu}{\sqrt{\sigma^2}} \cdot \gamma + \beta$$

40

として定義され、式中、

μ 、 σ^2 、 b 、及び γ は、バッチノルムパラメータ、例えば、各出力ノードについて 1 つの成分を含む転置ベクトルであり、

W は、重み行列であり、例えば、その各行は 1 つの出力ノードに対するものであり、それぞれの行の各成分は X の 1 つの行に関連付けられており、

50

X は、NN 層の活性化から導出される入力行列であり、
 b は、バイアスを形成する転置ベクトル、例えば、各出力ノードに対して 1 つの成分を含む転置ベクトルであり、
 ・ は、・ の一方の側の行列と他方の側の転置ベクトルとの間の列に関するアダマール乗算を示し、

* は、行列乗算を示す。

装置は、b 及び μ 及び γ 及び σ^2 又は β を受信し、
 【数 7】

$$\beta' := \beta + \frac{(b - \mu) * \gamma}{\sqrt{\sigma^2}} \quad 10$$

及び

【数 8】

$$\gamma' := \gamma \cdot \frac{1}{\sqrt{\sigma^2}} \quad 20$$

を計算するように構成される。

加えて、装置は、バッチノルム演算子の NN パラメータとして γ' 及び β' を NN 表現に符号化するように構成され、バッチノルム演算子を

【数 9】

$$\frac{W * X + b' - \mu'}{\sqrt{\sigma'^2}} \cdot \gamma' + \beta'$$

として定義し、

$\sigma'^2 := 1$ 、 $\mu' := 0$ 、 $b' := 0$ である。

【0026】

NN 表現から NN のバッチノルム演算子の NN パラメータを復号化する並列装置は、NN 表現から γ' 及び β' を導出し、その全ての成分に適用される 1 つのシグナリングによって、 $\sigma'^2 := 1$ 及び $\mu' := 0$ 及び $b' := 0$ を推論又は導出するように構成される。装置は、例えば、1 つのシグナリング、例えばフラグを読み出し、そこから $\sigma'^2 := 1$ 及び $\mu' := 0$ 及び $b' := 0$ を推論又は導出するように構成される。バッチノルム演算子は、第 3 の態様の第 2 の実施の形態に関して上記で説明したように定義される。

【0027】

したがって、本出願の第 3 の態様によれば、第 3 の実施の形態は、NN のバッチノルム演算子の NN パラメータを NN 表現に符号化する装置に関する。バッチノルム演算子は、

【数 10】

$$\frac{W * X - \mu}{\sqrt{\sigma^2 + \epsilon}} \cdot \gamma + \beta$$

として定義され、式中、

μ 、 σ^2 、 γ 、及び β は、バッチノルムパラメータ、例えば、各出力ノードについて 1 つの成分を含む転置ベクトルであり、

10

20

30

40

50

Wは、重み行列であり、例えば、その各行は1つの出力ノードに対するものであり、それぞれの行の各成分はXの1つの行に関連付けられており、

Xは、NN層の活性化から導出される入力行列であり、

は、ゼロ除算回避のための定数であり、

・は、・の一方の側の行列と他方の側の転置ベクトルとの間の列に関するアダマール乗算を示し、

*は、行列乗算を示す。

装置は、 μ 及び γ 及び σ^2 又は ϵ を受信し、

【数 1 1】

$$\beta' := \beta - \frac{\mu * \gamma}{\sqrt{\sigma^2 + \epsilon}}$$

10

及び

【数 1 2】

$$\gamma' := \gamma \cdot \frac{\sqrt{\theta + \epsilon}}{\sqrt{\sigma^2 + \epsilon}}$$

20

を計算するように構成される。

加えて、装置は、バッチノルム演算子のNNパラメータとして γ' 及び β' をNN表現に符号化するように構成され、バッチノルム演算子を

【数 1 3】

$$\frac{W * X - \mu'}{\sqrt{\sigma'^2 + \epsilon}} \cdot \gamma' + \beta'$$

30

として定義し、

$\sigma'^2 := \sigma^2$ 、 $\mu' := 0$ であり、 ϵ は所定のパラメータである。

【0 0 2 8】

NN表現からNNのバッチノルム演算子のNNパラメータを復号化する並列装置は、NN表現から γ' 及び β' を導出し、その全ての成分に適用される1つのシグナリングによって、 $\sigma'^2 := \sigma^2$ 及び $\mu' := 0$ を推論又は導出するように構成され、 ϵ は所定のパラメータである。装置は、例えば、1つのシグナリング、例えばフラグを読み出し、そこから $\sigma'^2 := \sigma^2$ 及び $\mu' := 0$ を推論又は導出するように構成される。バッチノルム演算子は、第3の態様の第3の実施の形態に関して上記で説明したように定義される。

40

【0 0 2 9】

したがって、本出願の第3の態様によれば、第4の実施の形態は、NNのバッチノルム演算子のNNパラメータをNN表現に符号化する装置に関する。バッチノルム演算子は、

【数 1 4】

$$\frac{W * X - \mu}{\sqrt{\sigma^2}} \cdot \gamma + \beta$$

として定義され、式中、

50

μ 、 σ^2 、 γ 、及び β は、バッチノルムパラメーター、例えば、各出力ノードについて 1 つの成分を含む転置ベクトルであり、

W は、重み行列であり、例えば、その各行は 1 つの出力ノードに対するものであり、それぞれの行の各成分は X の 1 つの行に関連付けられており、

X は、 NN 層の活性化から導出される入力行列であり、

\cdot は、 \cdot の一方の側の行列と他方の側の転置ベクトルとの間の列に関するアダマール乗算を示し、

$*$ は、行列乗算を示す。

装置は、 μ 及び γ 及び σ^2 又は β を受信し、
【数 15】

10

$$\beta' := \beta - \frac{\mu * \gamma}{\sqrt{\sigma^2}}$$

及び

【数 16】

$$\gamma' := \gamma \cdot \frac{1}{\sqrt{\sigma^2}}$$

20

を計算するように構成される。

加えて、装置は、バッチノルム演算子の NN パラメーターとして γ' 及び β' を NN 表現に符号化するように構成され、バッチノルム演算子を

【数 17】

$$\frac{W * X - \mu'}{\sqrt{\sigma'^2}} \cdot \gamma' + \beta'$$

30

として定義し、

$\sigma'^2 := 1$ 、 $\mu' := 0$ である。

【0030】

NN 表現から NN のバッチノルム演算子の NN パラメーターを復号化する並列装置は、 NN 表現から γ 及び β を導出し、その全ての成分に適用される 1 つのシグナリングによって、 $\sigma^2 := 1$ 及び $\mu := 0$ を推論又は導出するように構成される。装置は、例えば、1 つのシグナリング、例えばフラグを読み出し、そこから $\sigma^2 := 1$ 及び $\mu := 0$ を推論又は導出するように構成される。バッチノルム演算子は、第 3 の態様の第 4 の実施の形態に関して上記で説明したように定義される。

40

【0031】

以下の方法は、上述の原理に従って動作する。

【0032】

一実施の形態は、 NN 表現を生成する方法であって、量子化パラメーターから、量子化パラメーターによって導出された被除数と精度パラメーターによって導出された除数との間の除算の剰余に基づいて乗数が導出可能であり、量子化パラメーターから、除算の商の丸めに基づいてビットシフト数が導出可能であるように、 NN パラメーターの量子化パラメーター及び量子化値を決定することによって、 NN パラメーターを量子化された値に量子化することを含む方法に関する。量子化パラメーターは、 NN パラメーターの量子化された値が、量子化値と乗数に依存する因数との積であって、ビットシフト数に依存するピ

50

ット数だけビットシフトされた積に対応するように決定される。

【 0 0 3 3 】

一実施の形態は、NN表現からNNパラメータを導出する方法であって、NN表現から量子化パラメータ及び量子化値を導出することを含む方法に関する。加えて、方法は、量子化パラメータから、量子化パラメータによって導出された被除数と精度パラメータによって導出された除数との間の除算の剰余に基づいて乗数を導出することと、量子化パラメータから、除算の商の丸めに基づいてビットシフト数を導出することを含む。NNパラメータは、量子化値と乗数に依存する因数との積であって、ビットシフト数に依存するビット数だけビットシフトされた積に対応する。

【 0 0 3 4 】

一実施の形態は、NNを使用して推論を実行する方法であって、NN表現からNNパラメータを導出するために、NNパラメータを導出する上述の方法を使用してNNをパラメータ化することを含む方法に関する。加えて、推論を実行する方法は、NNを使用してNN入力に基づいて推論出力を計算することを含む。

【 0 0 3 5 】

一実施の形態は、NNを使用して推論を実行する方法であって、NNを使用してNN入力に基づいて推論出力を計算することを含む方法に関する。NNは、一对のNN層と、一对のNN層のうちの第1のNN層から一对のNN層のうちの第2のNN層へのニューロン間活性化フィードフォワードとを含む。方法は、第1のNN層のニューラルネットワークニューロンの活性化から行列Xを形成することと、 $s \cdot W' * X$ を計算することであって、式中、*は、行列乗算を表し、W'は、n及びm NNである次元n x mの重み行列であり、sは、長さnの転置ベクトルであり、 \cdot は、 \cdot の一方の側の行列と \cdot の他方の側の転置ベクトルとの間の列に関するアダマール乗算を示すこととによって、第1のNN層のニューラルネットワークニューロンの活性化に基づいて、第2のNN層のニューラルネットワークニューロンの活性化を計算することを含む。

【 0 0 3 6 】

一実施の形態は、NNのバッチノルム演算子のNNパラメータをNN表現に符号化する方法に関し、バッチノルム演算子は、

【数 1 8】

$$\frac{W * X + b - \mu}{\sqrt{\sigma^2 + \epsilon}} \cdot \gamma + \beta$$

として定義され、式中、 μ 、 σ^2 、 ϵ 、及び β は、バッチノルムパラメータであり、Wは、重み行列であり、Xは、NN層の活性化から導出される入力行列であり、bは、バイアスを形成する転置ベクトルであり、 ϵ は、ゼロ除算回避のための定数であり、 \cdot は、 \cdot の一方の側の行列と他方の側の転置ベクトルとの間の列に関するアダマール乗算を示し、*は、行列乗算を示す。方法は、b、 μ 、 σ^2 、及び β を受信することと、

【数 1 9】

$$\beta' := \beta + \frac{(b - \mu) * \gamma}{\sqrt{\sigma^2 + \epsilon}}$$

及び

【数 2 0】

10

20

30

40

50

$$\gamma' := \gamma \cdot \frac{\sqrt{\theta + \epsilon}}{\sqrt{\sigma^2 + \epsilon}}$$

を計算することと、を含む。

加えて、方法は、 γ' 及び β' をバッチノルム演算子の NN パラメータとして NN 表現に符号化することであって、バッチノルム演算子を

【数 2 1】

$$\frac{W * X + b' - \mu'}{\sqrt{\sigma'^2 + \epsilon}} \cdot \gamma' + \beta'$$

10

として定義し、 $\sigma'^2 := \sigma^2$ 、 $\mu' := 0$ 、及び $b' := 0$ であり、ここで、 γ' は所定のパラメータであることを含む。

【0 0 3 7】

一実施の形態は、NN のバッチノルム演算子の NN パラメータを NN 表現に符号化する方法に関し、バッチノルム演算子は、

20

【数 2 2】

$$\frac{W * X + b - \mu}{\sqrt{\sigma^2}} \cdot \gamma + \beta$$

として定義され、式中、 μ 、 σ^2 、 γ 、及び β は、バッチノルムパラメータであり、 W は、重み行列であり、 X は、NN 層の活性化から導出される入力行列であり、 b は、バイアスを形成する転置ベクトルであり、 \cdot は、 \cdot の一方の側の行列と他方の側の転置ベクトルとの間の列に関するアダマール乗算を示し、 $*$ は、行列乗算を示す。方法は、 b 、 μ 、 σ^2 、及び γ を受信することと、

30

【数 2 3】

$$\beta' := \beta + \frac{(b - \mu) * \gamma}{\sqrt{\sigma^2}}$$

及び

【数 2 4】

$$\gamma' := \gamma \cdot \frac{1}{\sqrt{\sigma^2}}$$

40

を計算することと、を含む。加えて、方法は、 γ' 及び β' をバッチノルム演算子の NN パラメータとして NN 表現に符号化することであって、バッチノルム演算子を

【数 2 5】

$$\frac{W * X + b' - \mu'}{\sqrt{\sigma'^2}} \cdot \gamma' + \beta'$$

として定義し、 $\sigma'^2 := 1$ 、 $\mu' := 0$ 、及び $b' := 0$ であることを含む。

50

【 0 0 3 8 】

一実施の形態は、NNのバッチノルム演算子のNNパラメータをNN表現に符号化する方法に関し、バッチノルム演算子は、

【数 2 6】

$$\frac{W \cdot X - \mu}{\sqrt{\sigma^2 + \epsilon}} \cdot \gamma + \beta$$

として定義され、式中、 μ 、 σ^2 、 ϵ 、及び β は、バッチノルムパラメータであり、 W は、重み行列であり、 X は、NN層の活性化から導出される入力行列であり、 \cdot は、ゼロ除算回避のための定数であり、 \cdot は、 \cdot の一方の側の行列と他方の側の転置ベクトルとの間の列に関するアダマール乗算を示し、 $*$ は、行列乗算を示す。方法は、 μ 、 σ^2 、 ϵ 、及び β を受信することと、

【数 2 7】

$$\beta' := \beta - \frac{\mu \cdot \gamma}{\sqrt{\sigma^2 + \epsilon}}$$

及び

【数 2 8】

$$\gamma' := \gamma \cdot \frac{\sqrt{\theta + \epsilon}}{\sqrt{\sigma^2 + \epsilon}}$$

を計算することと、 θ を含む。加えて、方法は、 β' 及び γ' をバッチノルム演算子のNNパラメータとしてNN表現に符号化することであって、バッチノルム演算子を

【数 2 9】

$$\frac{W \cdot X - \mu'}{\sqrt{\sigma'^2 + \epsilon}} \cdot \gamma' + \beta'$$

として定義し、 $\sigma'^2 := \sigma^2 + \theta$ 及び $\mu' := 0$ であり、ここで、 θ は所定のパラメータであることを含む。

【 0 0 3 9 】

一実施の形態は、NNのバッチノルム演算子のNNパラメータをNN表現に符号化する方法に関し、バッチノルム演算子は、

【数 3 0】

$$\frac{W \cdot X - \mu}{\sqrt{\sigma^2}} \cdot \gamma + \beta$$

として定義され、式中、 μ 、 σ^2 、 ϵ 、及び β は、バッチノルムパラメータであり、 W は、重み行列であり、 X は、NN層の活性化から導出される入力行列であり、 \cdot は、 \cdot の一方の側の行列と他方の側の転置ベクトルとの間の列に関するアダマール乗算を示し、 $*$ は、行列乗算を示す。方法は、 μ 、 σ^2 、 ϵ 、及び β を受信することと、

【数 3 1】

$$\beta' := \beta - \frac{\mu \cdot \gamma}{\sqrt{\sigma^2}}$$

10

20

30

40

50

及び

【数 3 2】

$$\gamma' := \gamma \cdot \frac{1}{\sqrt{\sigma^2}}$$

を計算することと、を含む。加えて、方法は、 σ^2 及び μ' をバッチノルム演算子の NN パラメータとして NN 表現に符号化することであって、バッチノルム演算子を

【数 3 3】

$$\frac{W * X - \mu'}{\sqrt{\sigma'^2}} \cdot \gamma' + \beta'$$

10

として定義し、 $\sigma'^2 := 1$ 、及び $\mu' := 0$ であることを含む。

【0040】

一実施の形態は、NN 表現から NN のバッチノルム演算子の NN パラメータを復号化する方法に関し、バッチノルム演算子は、

【数 3 4】

$$\frac{W * X + b - \mu}{\sqrt{\sigma^2 + \epsilon}} \cdot \gamma + \beta$$

20

として定義され、式中、 μ 、 σ^2 、 ϵ 、及び b は、バッチノルムパラメータであり、 W は、重み行列であり、 X は、NN 層の活性化から導出される入力行列であり、 b は、バイアスを形成する転置ベクトルであり、 ϵ は、ゼロ除算回避のための定数であり、 \cdot は、 \cdot の一方の側の行列と他方の側の転置ベクトルとの間の列に関するアダマール乗算を示し、 $*$ は、行列乗算を示す。方法は、NN 表現から σ^2 及び μ' を導出することと、その全ての成分に適用される 1 つのシグナリングによって、 $\sigma'^2 := 1$ 、 $\mu' := 0$ 、及び $b' := 0$ であり、ここで、 σ' は所定のパラメータであることを推論又は導出することを含む。

30

【0041】

一実施の形態は、NN 表現から NN のバッチノルム演算子の NN パラメータを復号化する方法に関し、バッチノルム演算子は、

【数 3 5】

$$\frac{W * X + b - \mu}{\sqrt{\sigma^2}} \cdot \gamma + \beta$$

として定義され、式中、 μ 、 σ^2 、 ϵ 、及び b は、バッチノルムパラメータであり、 W は、重み行列であり、 X は、NN 層の活性化から導出される入力行列であり、 b は、バイアスを形成する転置ベクトルであり、 \cdot は、 \cdot の一方の側の行列と他方の側の転置ベクトルとの間の列に関するアダマール乗算を示し、 $*$ は、行列乗算を示す。方法は、NN 表現から σ^2 及び μ' を導出することと、その全ての成分に適用される 1 つのシグナリングによって、 $\sigma'^2 := 1$ 、 $\mu' := 0$ 、及び $b' := 0$ であることを推論又は導出することを含む。

40

【0042】

一実施の形態は、NN 表現から NN のバッチノルム演算子の NN パラメータを復号化する方法に関し、バッチノルム演算子は、

【数 3 6】

50

$$\frac{W \cdot X - \mu}{\sqrt{\sigma^2 + \epsilon}} \cdot \gamma + \beta$$

として定義され、式中、 μ 、 σ^2 、 ϵ 、及び β は、バッチノルムパラメータであり、 W は、重み行列であり、 X は、 $N \times N$ 層の活性化から導出される入力行列であり、 ϵ は、ゼロ除算回避のための定数であり、 \cdot は、 \cdot の一方の側の行列と他方の側の転置ベクトルとの間の列に関するアダマール乗算を示し、 $*$ は、行列乗算を示す。方法は、 $N \times N$ 表現から及び β を導出することと、その全ての成分に適用される1つのシグナリングによって、 $\sigma^2 := \sigma^2$ 及び $\mu := 0$ であり、ここで、 ϵ は所定のパラメータであることを推論又は導出することを含む。

10

【0043】

一実施の形態は、 $N \times N$ 表現から $N \times N$ のバッチノルム演算子の $N \times N$ パラメータを復号化する方法に関し、バッチノルム演算子は、

【数37】

$$\frac{W \cdot X - \mu}{\sqrt{\sigma^2}} \cdot \gamma + \beta$$

20

として定義され、式中、 μ 、 σ^2 、 ϵ 、及び β は、バッチノルムパラメータであり、 W は、重み行列であり、 X は、 $N \times N$ 層の活性化から導出される入力行列であり、 \cdot は、 \cdot の一方の側の行列と他方の側の転置ベクトルとの間の列に関するアダマール乗算を示し、 $*$ は、行列乗算を示す。方法は、 $N \times N$ 表現から及び β を導出することと、その全ての成分に適用される1つのシグナリングによって、 $\sigma^2 := 1$ 及び $\mu := 0$ であることを推論又は導出することを含む。

【0044】

上述したように、これらの方法は、上述した装置又はデバイスと同じ考察に基づいている。方法は、装置又はデバイスに関しても説明される全ての特徴及び機能で完成され得る。

30

【0045】

一実施の形態は、上述したように、 $N \times N$ 表現を生成する方法又は装置によって生成された $N \times N$ 表現を定義するデジタルデータを含むデジタル記憶媒体に関する。

【0046】

一実施の形態は、上述の方法のうちの1つを実行するコンピュータプログラムに関する。

【0047】

一実施の形態は、上述したように、 $N \times N$ 表現を生成する方法又は装置によって生成されるデータストリームに関する。

【0048】

図面は、必ずしも縮尺通りではなく、代わりに、概して、本発明の原理を図示することに重点が置かれている。以下の説明では、本発明の種々の実施形態が、以下の図面を参照して説明される。

40

【図面の簡単な説明】

【0049】

【図1】ニューラルネットワークを示す図である。

【図2】本発明の一実施形態による、 $N \times N$ 表現を生成する装置、 $N \times N$ 表現を定義するデジタルデータ、及び $N \times N$ 表現から $N \times N$ パラメータを導出する装置を概略的に示す図である。

【図3】フィードフォワードニューラルネットワークを概略的に示す図である。

【図4】本発明の一実施形態による、 $N \times N$ パラメトライザーを使用して推論を実行するデ

50

バイスを概略的に示す図である。

【図 5】本発明の一実施形態による、ベクトル及び行列の合成として重みパラメータを因数分解することによって推論を実行するデバイスを概略的に示す図である。

【図 6】本発明の一実施形態による、NNパラメータをNN表現に符号化する装置及びNN表現からNNパラメータを復号化する装置を概略的に示す図である。

【図 7】行列XとWとの間の可能な関係を概略的に示す図である。

【発明を実施するための形態】

【0050】

同一若しくは同等の要素、又は同一若しくは同等の機能を有する要素は、異なる図に現れる場合であっても、以下の説明において同一又は同等の参照番号によって示される。

10

【0051】

以下の説明では、本発明の実施形態のより完全な説明を提供するために、複数の詳細が記載される。しかしながら、本発明の実施形態がこれらの具体的な詳細なしに実施され得ることは、当業者には明らかであろう。他の例では、本発明の実施形態を不明瞭なことを回避するために、周知の構造及びデバイスは、詳細にではなくブロック図の形態で示される。加えて、本明細書で後に説明される異なる実施形態の特徴は、特に別段の記載がない限り、互いに組み合わせることができる。

【0052】

以下では、少数の浮動小数点演算のみを用いた、又は更には浮動小数点演算を全く用いない推論を可能にする層のパラメータの量子化ステップサイズを表し、シグナリングする効率的な手法を提示する。つまり、ビットレートの点で効率的な表現であり、更に推論の効率的な計算に利用することができる。

20

【0053】

図 2 は、NN表現 110 を生成する装置 100 を示している。装置 100 は、量子化パラメータ 142 を決定 (140) することと、NNパラメータ 120 の量子化値 152 を決定 (150) することとによって、NNパラメータ 120 を量子化された値 130 に量子化するように構成される。量子化値 152 は、量子化パラメータ 142 に基づいて決定 (150) することができる。量子化パラメータ 142 の決定 (140) は、量子化パラメータ決定器によって実行することができる。量子化値 152 の決定 (150) は、量子化値決定器によって実行することができる。

30

【0054】

量子化パラメータ 142 が決定 (140) され、それにより、量子化パラメータ 142 から、乗数 144 及びビットシフト数 146 が導出可能である。量子化パラメータ 142 の決定 (140) において、装置 100 は、例えば、乗数 144 及びビットシフト数 146 が、決定された量子化パラメータ 142 から導出可能であるかどうかを既にチェックしている場合がある。

【0055】

任意選択で、装置 100 は、量子化パラメータ 142 から乗数 144 を導出し、量子化パラメータ 142 からビットシフト数 146 を導出して、例えば装置 100 による量子化された値 130 の決定を可能にするように構成することができる。しかし、量子化された値 130 は量子化パラメータ 142 及び量子化値 152 によって表すことができるので、これは必要ではない。装置 100 は、量子化された値 130 を明示的に決定する必要はない。

40

【0056】

一実施形態によれば、生成されたNN表現 110 は、決定された量子化パラメータ 142 及び決定された量子化値 152 を含むことができ、それにより、NNパラメータ 120、すなわちNNパラメータ 120 の量子化された値 130 は、NN表現 110 から導出可能である。例えば、装置 100 は、量子化パラメータ 142 及び量子化値 152 をNN表現 110 に符号化するように構成されてもよい。

【0057】

50

乗数 1 4 4 は、量子化パラメーター 1 4 2 によって導出される被除数と、精度パラメーター k 1 4 5 によって導出される除数との間の除算の剰余に基づいて、量子化パラメーター 1 4 2 から導出可能である。

【 0 0 5 8 】

ビットシフト数 1 4 6 は、除算の商の丸めに基づいて、すなわち、量子化パラメーター 1 4 2 によって導出された被除数と精度パラメーター k 1 4 5 によって導出された除数との間の除算の商の丸めに基づいて、量子化パラメーター 1 4 2 から導出可能である。

【 0 0 5 9 】

量子化パラメーター 1 4 2 の決定 (1 4 0) が実行され、それにより、NNパラメーター 1 2 0 の量子化された値 1 3 0 は、量子化値 1 5 2 と乗数 1 4 4 に依存する因数 1 4 8 との間の積であって、ビットシフト数 1 4 6 に依存するビット数だけビットシフトされた積に対応する。NNパラメーター 1 2 0 の量子化された値 1 3 0 は、例えば、少なくとも、シフトの場合には符号の別個の処理を伴う量子化された値の絶対値に関して、又は更には、積、その因数、及びシフトに対してそれぞれ 2 の補数表現及び 2 の補数演算を使用する場合等、絶対値と符号との両方に関して、積に対応する。これは、ユニット 1 5 0 に例示的かつ概略的に示されている。

10

【 0 0 6 0 】

一実施形態によれば、装置 1 0 0 は、NNパラメーターのための浮動小数点表現を使用して NN 2 0 を訓練することによって、かつ量子化誤差を低減することを目的とする反復最適化方式によって NNパラメーターのための量子化パラメーター 1 4 2 及び量子化値 1 5 2 を決定することによって、NNパラメーター、例えば NNパラメーター 1 2 0 の量子化された値 1 3 0 を提供するように構成される。

20

【 0 0 6 1 】

装置 1 0 0 とは別に、図 1 は、NN表現 1 1 0 を定義するデジタルデータ 2 0 0 と、NN表現 1 1 0 から NNパラメーター、すなわち NNパラメーター 1 2 0 の量子化された値 1 3 0 を導出する装置 3 0 0 とを示している。デジタルデータ 2 0 0 及び装置 3 0 0 が NNパラメーターの元の値を見ることがないという事実から、量子化された値 1 3 0 は、この文脈では NNパラメーターの値として理解される。このため、NNパラメーターは、デジタルデータ 2 0 0 及び装置 3 0 0 の以下の説明のために、1 3 0 として示される。本明細書で説明される NNパラメーターは、NNパラメーターに割り当てられた元の値 1 2 0 によって、又は元の値 1 2 0 に基づいて決定された量子化された値 1 3 0 によって表され得ることが明らかである。したがって、以下では、NNパラメーターを $1 2 0 / 1 3 0$ として示す。これは、例えば、NNパラメーターが元の値 1 2 0 と量子化された値 1 3 0 のどちらで表されても一般的に適用できる特徴を説明するものである。

30

【 0 0 6 2 】

デジタルデータ 2 0 0 は、NN表現 1 1 0 を定義し、NN表現 1 1 0 は、NNパラメーター 1 3 0 を表すために、量子化パラメーター 1 4 2 及び量子化値 1 5 2 を含んでおり、それにより、量子化パラメーター 1 4 2 から、量子化パラメーター 1 4 2 によって導出された被除数と精度パラメーター k 1 4 5 によって導出された除数との間の除算の剰余に基づいて乗数 1 4 4 が導出可能であり、かつ、量子化パラメーター 1 4 2 から、除算の商の丸めに基づいてビットシフト数 1 4 6 が導出可能である。NN表現 1 1 0 は、量子化パラメーター 1 4 2 及び量子化値 1 5 2 を含んでおり、それにより、NNパラメーター 1 3 0 は、量子化値 1 5 2 と乗数 1 4 4 に依存する因数 1 4 8 との間の積であって、ビットシフト数 1 4 6 に依存するビット数だけビットシフトされた積に対応する。

40

【 0 0 6 3 】

NN表現 1 1 0 から NNパラメーター 1 3 0 を導出する装置 3 0 0 は、例えば、量子化パラメーター導出ユニット 3 1 0 を使用して、NN表現 1 1 0 から量子化パラメーター 1 4 2 を導出し、例えば、量子化値導出ユニット 3 2 0 を使用して、NN表現 1 1 0 から量子化値 1 5 2 を導出するように構成される。加えて、装置 3 0 0 は、量子化パラメーター 1 4 2 から、乗数 1 4 4 及びビットシフト数 1 4 6 を導出するように構成される。装置 3

50

00は、量子化パラメーター142によって導出された被除数と精度パラメーター145によって導出された除数との間の除算の剰余に基づいて乗数144を導出し、除算の商の丸めに基づいてビットシフト数146を導出するように構成される。乗数144の導出は、乗数導出ユニット330を使用して実行されてもよく、ビットシフト数146の導出は、ビットシフト数導出ユニット340を使用して実行されてもよい。NNパラメーター130は、量子化値152と乗数144に依存する因数148との間の積であって、ビットシフト数146に依存するビット数だけビットシフトされた積に対応する(図2の装置100及びユニット150についての上記の対応する説明を参照)。NNパラメーター130は、例えば、NNパラメーター導出ユニット350を用いて導出されてもよい。NNパラメーター導出ユニット350は、装置100のオプションのユニット150と同じ特徴及び/又は機能を備えてもよい。

10

【0064】

以下では、装置100及び装置300の両方に適用可能な実施形態及び例が提示される。

【0065】

一実施形態によれば、NNパラメーター120/130は、重みパラメーター、バッチノルムパラメーター、及びバイアスのうちの1つである。重みパラメーター、例えば、 W の成分 w は、一对のニューロンの間のニューロン間活性化フィードフォワードを重み付けするために使用可能であり得るか、又は代替的に言えば、第1のニューロンと第2のニューロンとを接続するエッジに関係し、第2のニューロンのためのインバウンド活性化の和において第1のニューロンの活性化のフォーディングを重み付けする重みを表し得る。バッチノルムパラメーター、例えば、 μ 、 σ^2 、 γ は、ニューラルネットワーク層のアフィン変換をパラメーター化するために使用可能とすることができ、バイアス、例えば、 b_i の成分は、所定のニューラルネットワークニューロンのためのインバウンドニューロン間活性化フィードフォワードの和にバイアスをかけるために使用可能であり得る。

20

【0066】

一実施形態によれば、NNパラメーター120/130は、例えば図1に示すように、NNの複数のニューロン間活性化フィードフォワード122のうちの単一のニューロン間活性化フィードフォワード12_i、例えば W の成分 w に関して、NN20をパラメーター化する。装置100/装置300は、複数のニューロン間活性化フィードフォワード122の各々について、対応するNNパラメーター120/130をNN表現110に符号化/NN表現110から導出するように構成される。対応するNNパラメーター130は、NN表現110に含まれる。この場合、装置100は、複数のニューロン間活性化フィードフォワード122の各々について、それぞれのニューロン間活性化フィードフォワード12_iに関連付けられた関連量子化パラメーター142と、それぞれのニューロン間活性化フィードフォワード12_iに関連付けられた関連量子化値152とを決定(140)することによって、対応するNNパラメーター120を対応する量子化された値130に量子化するように構成され得る。関連量子化パラメーター142の決定(140)は、関連量子化パラメーター142から、関連量子化パラメーター142によって導出された被除数と、それぞれのニューロン間活性化フィードフォワード12_iに関連付けられた関連精度パラメーター145によって導出された除数との間の除算の剰余に基づいてそれぞれのニューロン間活性化フィードフォワード12_iに関連付けられた関連乗数144と、除算の商の丸めに基づいてそれぞれのニューロン間活性化フィードフォワード12_iに関連付けられた関連ビットシフト数146と、を導出することができるように実行される。この場合の対応する装置300は、複数のニューロン間活性化フィードフォワード122の各々について、NN表現110からそれぞれのニューロン間活性化フィードフォワード12_iに関連付けられた関連量子化パラメーター142を導出(310)し、NN表現110からそれぞれのニューロン間活性化フィードフォワード12_iに関連付けられた関連量子化値152を導出(320)するように構成される。導出(310及び320)は、例えばNN表現110から復号化することによって実行されてもよく、すなわちエッジごとに1つが復号化されてもよい。加えて、装置300は、複数のニューロン間活性化フィード

30

40

50

フォワード122の各々について、関連量子化パラメーター142から、関連量子化パラメーター142によって導出された被除数と、それぞれのニューロン間活性化フィードフォワード12_iに関連付けられた関連精度パラメーター145によって導出された除数との間の除算の剰余に基づいて、それぞれのニューロン間活性化フィードフォワード12_iに関連付けられた関連乗数144と、除算の商の丸めに基づいて、それぞれのニューロン間活性化フィードフォワード12_iに関連付けられた関連ビットシフト数146と、を導出するように構成される(330及び340参照)。導出(330及び340)は、例えばNN表現110から復号化することによって実行されてもよく、すなわちエッジごとに1つが復号化されてもよい。

【0067】

別の実施形態によれば、装置100/装置300は、NN20の複数のニューロン間活性化フィードフォワード122をニューロン間活性化フィードフォワードのサブグループ122a、122bに細分するように構成され、それにより、各サブグループは、NNの関連する一対のNN層に関連付けられ、関連する一対のNN層の間のニューロン間活性化フィードフォワードを含み、関連する一対の層以外の更なる一対のNN層の間のニューロン間活性化フィードフォワードを除外し、2つ以上のサブグループが所定のNN層に関連付けられる(例えば図3を参照)。サブグループ122aは、例えば、NN20の関連する一対のNN層114及び116₁に関連付けられ、関連する一対のNN層114及び116₁の間のニューロン間活性化フィードフォワードを含み、関連する一対の層114及び116₁以外の更なる一対のNN層の間、例えば更なる一対のNN層116₁及び116₂の間のニューロン間活性化フィードフォワードを除外する。サブグループ122a及び122bは、層116₁に関連付けられる。NN20の複数のニューロン間活性化フィードフォワード122の細分化は、例えば、NN20内の各エッジ/重み12のインデックスによって、又は他の形態で各層の対の間のエッジ12をセグメント化することによって実行されてもよい。NNパラメーター120/130は、NN2の複数のニューロン間活性化フィードフォワード122のうちの単一のニューロン間活性化フィードフォワード12_iに関してNN20をパラメーター化する。複数のニューロン間活性化フィードフォワード122の各々について、対応するNNパラメーター120/130がNN表現110に含まれる。装置300は、例えば、NN表現から復号化することによって、すなわち、エッジサブグループごとに1つのサブグループを復号化することによって、複数のニューロン間活性化フィードフォワード122の各々について、NN表現110から対応するNNパラメーター120/130を導出するように構成される。装置100/装置300は、ニューロン間活性化フィードフォワードのサブグループ122a、122bごとに、それぞれのサブグループ122a又は122bに関連付けられた関連量子化パラメーター142を決定(140)/導出(310)するように構成される。量子化パラメーター142は、それぞれのサブグループ122a又は122bに関連付けられた関連乗数144が、関連量子化パラメーター142によって導出された被除数とそれぞれのサブグループに関連付けられた関連精度パラメーター145によって導出された除数との間の除算の剰余に基づいて量子化パラメーター142から導出可能であるように、装置100によって決定(140)され、量子化パラメーター142は、それぞれのサブグループ122a又は122bに関連付けられた関連ビットシフト数146が、除算の商の丸めに基づいて量子化パラメーター142から導出可能であるように、装置100によって決定(140)される。装置300は、NN表現110から関連乗数144及び関連ビットシフト数146を導出するように構成される。装置100/装置300は、複数のニューロン間活性化フィードフォワード122の各々について、NN表現110からそれぞれのニューロン間活性化フィードフォワード12_iに関連付けられた関連量子化値152を決定(150)/導出(320)する(例えば、NN表現110から復号化することによって、すなわち、エッジごとに1つが復号化されることによって、導出(320)する)ように構成される。それぞれのニューロン間活性化フィードフォワード12_iの対応するNNパラメーター120/130は、関連量子化値142と、それぞれのニューロン間活性化フィードフォ

10

20

30

40

50

ワード 12_i が含まれるサブグループ、例えば $122a$ 又は $122b$ に関連付けられた関連乗数 144 に依存する因数 148 との間の積であって、それぞれのニューロン間活性化フィードフォワード 12_i が含まれるサブグループ、例えば $122a$ 又は $122b$ の関連ビットシフト数 146 に依存するビット数だけビットシフトされた積に対応する。

【0068】

関連精度パラメータ 145 は、例えば、 $NN20$ にわたって、又は各 NN 層 114 、 116_1 、及び 116_2 内でグローバルに等しく評価される。任意選択で、装置 100 / 装置 300 は、関連精度パラメータ 145 を NN 表現 110 に符号化 / NN 表現 110 から導出するように構成される。

【0069】

一実施形態によれば、装置 100 / 装置 300 は、コンテキスト適応型二値算術符号化 / 復号化を使用することによって、又は量子化パラメータ 142 を表すビットを NN 表現 110 に直接書き込む / NN 表現 110 から読み出すことによって、又は装置 100 / 装置 300 のコンテキスト適応型二値エンコーダ / デコーダの等確率バイパスモードを介して量子化パラメータ 142 を表すビットを NN 表現 110 から符号化 / 導出することによって、量子化パラメータ 142 を NN 表現 110 に符号化 / NN 表現 110 から導出するように構成される。装置 100 / 装置 300 は、二値化方式を使用してピンストリングを二値化 / 非二値化することによって NN 表現 110 から量子化パラメータ 142 を導出するように構成され得る。二値化方式は、例えば、指数ゴロム符号である。

【0070】

一実施形態によれば、装置 100 は、量子化パラメータ 142 を決定 (140) し、それを固定小数点表現、例えば 2 の補数表現の形態で NN 表現 110 に符号化するように構成される。装置 300 は、固定小数点表現、例えば 2 の補数表現の形態で NN 表現 110 から量子化パラメータ 142 を導出 (310) するように構成されてもよい。任意選択で、精度パラメータ 145 は 2^t であり、固定小数点表現、例えば 2 の補数表現のビット長は、 $NN20$ に対して一定になるように設定されるか、又は $NN20$ に対して一定である基底ビット長と t との和になるように設定される。

【0071】

一実施形態によれば、装置 100 / 装置 300 は、整数値シンタックス要素として量子化パラメータ 142 を NN 表現 110 に符号化 / NN 表現 110 から導出するように構成される。

【0072】

一実施形態によれば、装置 100 は、量子化値 152 を決定し、それを固定小数点表現、例えば 2 の補数表現の形態で NN 表現 110 に符号化するように構成される。装置 300 は、固定小数点表現、例えば 2 の補数表現の形態で NN 表現 110 から量子化値 152 を導出 (320) するように構成されてもよい。

【0073】

一実施形態によれば、装置 100 / 装置 300 は、二値化方式に従って量子化値 152 をピンストリングに二値化 / ピンストリングから非二値化し、コンテキスト適応算術符号化 / 復号化を使用してピンストリングのビットを符号化 / 復号化することによって、量子化値 152 を NN 表現 110 に符号化 / NN 表現 110 から導出するように構成される。

【0074】

一実施形態によれば、装置 100 / 装置 300 は、二値化方式に従って量子化値 152 をピンストリングに二値化 / ピンストリングから非二値化し、コンテキスト適応算術符号化 / 復号化を使用してピンストリングの第1のビットを符号化 / 復号化し、等確率バイパスモードを使用してピンストリングの第2のビットを符号化 / 復号化することによって、量子化値 152 を NN 表現 110 に符号化 / NN 表現 110 から復号化するように構成される。

【0075】

一実施形態によれば、量子化ステップサイズ 149 は、装置 100 及び / 又は装置 3

10

20

30

40

50

00によって、量子化パラメータ QP 142で示される符号付き整数及び正の整数パラメータ k 、すなわち精度パラメータ145から、以下の式に従って導出することができる。

$$mul = k + QP \% k$$

【数38】

$$shift = \lfloor QP/k \rfloor$$

$$= (mul / k) \cdot 2^{shift}$$

【0076】

乗数144は mul で表され、ビットシフト数146は $shift$ で表され、因数148は mul / k で表される。

【0077】

NNパラメータ130は、 $(mul / k) \cdot 2^{shift} \cdot P$ であり、ここで、 P は量子化値152である。

【0078】

フロア演算子

【数39】

$\lfloor \cdot \rfloor$

及びモジュロ演算子 $\%$ は、以下のように定義される。

【数40】

$\lfloor x \rfloor$

は、 x 以下の最大の整数である。

$x \% y$ は、

【数41】

$$x - y \cdot \lfloor x/y \rfloor$$

として定義されるモジュロ演算子である。

【0079】

任意選択で、装置100及び/又は装置300は、精度パラメータ k 145をデフォルト値に設定するように構成することができる。

【0080】

あるいは、装置100は、任意選択で、自然数又は2の累乗等の精度パラメータ k 145の複数の異なる整数値をテストすることができる。異なる整数値は、例えば、NN全体に対して、又は各層等のNNの各部分に対してテストされ、量子化誤差及びビットレートに関して、例えばそのラグランジュ和に関して最良の精度パラメータ k 145が選択される。装置100は、例えば、精度パラメータ k 145を決定して、例えば決定(140)において、乗数144及びビットシフト数146が量子化パラメータ142から導出可能であるかどうかをチェックするように構成されてもよい。任意選択で、装置100によって選択された精度パラメータ k 145は、NN表現110においてシグナリングされ、例えば、NN表現110に符号化される。装置300は、例えば、NN表現11

10

20

30

40

50

0 から精度パラメーター k_{145} を導出するように構成される。

【0081】

一実施形態によれば、精度パラメーター k_{145} は、2 の累乗である。

【0082】

一実施形態によれば、装置 100 / 装置 300 は、精度パラメーター k_{145} を表すビットを直接 NN 表現 110 に書き込む / NN 表現 110 から読み出すことによって、又は装置 100 / 装置 300 のコンテキスト適応型二値エンコーダー / デコーダーの等確率バイパスモードを介して精度パラメーター k_{145} を表すビットを NN 表現 110 に / NN 表現 110 から導出することによって、精度パラメーター k_{145} を NN 表現 110 に符号化 / NN 表現 110 から導出するように構成される。

10

【0083】

ビットストリーム、例えばデジタルデータ 200 において 32 ビット浮動小数点値をシグナリングする代わりに、パラメーター QP_{142} 及び k_{145} のみがシグナリングされる必要がある。一部の適用例では、ビットストリーム中で QP_{142} をシグナリングし、 k_{145} を何らかの固定値に設定することで十分な場合もある。

【0084】

好ましい実施形態において、パラメーター $QP' = QP - QP_0$ が、 QP_{142} の代わりにビットストリーム内でシグナリングされ、ここで、パラメーター QP_0 は、所定の定数値である。換言すれば、一実施形態によれば、装置 100 / 装置 300 は、関連量子化パラメーター QP_{142} を、参照量子化パラメーター QP_0 との差分の形態で、NN 表現 110 に符号化 / NN 表現 110 から導出するように構成される。

20

【0085】

別の好ましい実施形態において、 k_{145} は 2^t に設定される。このようにして、149 の算出は、除算を含まずに以下のように実行することができる。

$$= \text{mul} \cdot 2^{\text{shift} - t}$$

【0086】

これにより、一部の計算を、以下に例示されるように、浮動小数点領域の代わりに整数領域で実行することが可能となる。

【0087】

図 4 は、NN20 を使用して推論を実行するデバイス 400 を概略的に示している。デバイス 400 は、NN20 をパラメーター化するように構成された NN パラメトライザー 410 を備える。NN パラメトライザー 410 は、NN 表現 110 から NN パラメーター 130 を導出する装置 300 を備える。NN パラメーター 130 を導出する装置 300 は、図 2 の装置 300 に関して説明したものと同一又は同様の特徴を備えることができる。装置 300 は、NN パラメーター導出ユニットとして理解され得る。加えて、デバイス 400 は、NN20 を使用して、例えば、NN パラメトライザー 410 によって決定された NN20 のパラメーター化 450 を使用して、NN 入力 440 に基づいて推論出力 430 を計算するように構成された計算ユニット 420 を備える。

30

【0088】

例 1 :

一実施形態によれば、NN パラメトライザー 410 は、装置 300 を介して、第 1 の NN パラメーター及び第 2 の NN パラメーターのうちの少なくとも 1 つを導出するように構成され、それにより、第 1 の NN パラメーターは、第 1 の量子化値と第 1 の因数との間の積であって、第 1 のビット数だけビットシフトされた積に対応し、第 2 の NN パラメーターは、第 2 の量子化値と第 2 の因数との間の積であって、第 2 のビット数だけビットシフトされた積に対応する。

40

【0089】

第 1 の量子化値及び第 2 の量子化値は両方とも、図 2 において 152 で示される量子化値を表す。第 1 の因数及び第 2 の因数は両方とも、図 2 において 148 で示される因数を表す。

50

【 0 0 9 0 】

例えば、 $t = 2$ とし、 $k = 2^t$ とし、 $Q P_a$ で示される第1のQP、すなわち第1の量子化パラメータ142、関連する $shift_a$ 、すなわち第1のビットシフト数146、 mul_a 、すなわち第1の乗数144、及び Δ_a 、すなわち第1の量子化ステップサイズ149を定義する。

【 0 0 9 1 】

さらに、 $Q P_b$ で示される第2のQP、すなわち第2の量子化パラメータ142、関連する $shift_b$ 、すなわち第2のビットシフト数146、 mul_b 、すなわち第2の乗数144、及び Δ_b 、すなわち第2の量子化ステップサイズ149を定義する。

【 0 0 9 2 】

「第1の」パラメータ及び「第2の」パラメータは、この文脈では同じ参照番号で示されているが、それらが異なる値を有し得ることは明らかである。それらは、それらが図2に示されるどの特徴に属するかを明確にするために、同じ参照番号で示されているのみである。

【 0 0 9 3 】

$C = \Delta_a \cdot C_a$ が成り立つ第1の量子化行列 C_a を考える。

$D = \Delta_b \cdot D_b$ が成り立つ第2の量子化行列 D_b を考える。

すなわち、 C_a は、 $Q P_a$ を使用して量子化され、 D_b は、 $Q P_b$ を使用して量子化されている。

両方の行列は、同じ次元を有する。図2で説明した量子化値152は、 C_a の1つの成分又は D_b の1つの成分を表すことができる。例えば、 C_a は、複数の第1の量子化値152を含むことができ、 D_b は、複数の第2の量子化値152を含むことができる。

【 0 0 9 4 】

さらに、和 $C + D$ が以下のように算出されると仮定する。

【 数 4 2 】

$$\begin{aligned} C + D &= \Delta_a \cdot C_a + \Delta_b \cdot D_b = 2^{shift_a-2} \cdot mul_a \cdot C_a + 2^{shift_b-2} \cdot mul_b \cdot D_b \\ &= 2^{shift_a-2} \cdot (mul_a \cdot C_a + 2^{shift_b-shift_a} \cdot mul_b \cdot D_b) \end{aligned}$$

【 0 0 9 5 】

デバイス400は、第1のNNパラメータCの第1の量子化値 C_a によって形成され、第1の乗数 mul_a で重み付けされた第1の加数、例えば $mul_a \cdot C_a$ と、第2のNNパラメータDの第2の量子化値 D_b によって形成され、第2の乗数 mul_b で重み付けされ、第1のビット数及び第2のビット数の差、例えば、

【 数 4 3 】

$$2^{shift_b-shift_a}$$

だけビットシフトされた第2の加数、例えば、

【 数 4 4 】

$$2^{shift_b-shift_a} \cdot mul_b \cdot D_b$$

との間の和を形成することと、第1の加数及び第2の加数の和に、第1のビット数及び第2のビット数の一方に依存する、例えば第1のビットシフト数 $shift_a$ 又は第2のビットシフト数 $shift_b$ に依存するビット数だけビットシフト

10

20

30

40

50

【数 4 5】

$$2^{shift_a-2}$$

を施すことと、によって、第 1 の NN パラメータ C 及び第 2 の NN パラメータ D を加算して NN 2 0 の最終 NN パラメータを生成するように構成される。

【0096】

任意選択で、この算出 / 計算は、計算ユニット 4 2 0 によって実行することができる。この場合、計算ユニット 4 2 0 は、上述したように、第 1 の NN パラメータ C と第 2 の NN パラメータ D とを加算して NN 2 0 の最終的な NN パラメータを生成するように構成されている。

【0097】

式から分かるように、浮動小数点演算を必要とし得る C 及び D を導出する必要はない。その代わりに、C_a の要素、すなわち第 1 の量子化値 1 5 2 は、単に mul_a、すなわち第 1 の乗数 1 4 4 と乗算され、D_b の要素、すなわち第 2 の量子化値 1 5 2 は、mul_b、すなわち第 2 の乗数 1 4 4 と乗算され、因数

【数 4 6】

$$2^{shift_b-shift_a}$$

は、C_a の第 1 の量子化値 1 5 2、すなわち C_a の成分に関連する第 1 のビットシフト数 shift_a 1 4 6 と、D_b の第 2 の量子化値 1 5 2、すなわち D_b の成分に関連する第 2 のビットシフト数 shift_b 1 4 6 とに依存する単純なビットシフト演算として実施される。t = 2 であるので、整数変数 mul_a 及び mul_b は両方とも値 4、5、6、及び 7 のうちの 1 つであることに留意されたい。かかる小さい数を有する整数乗算は、ハードウェア又はソフトウェア実施態様において非常に効率的に実施され得る。

【0098】

一実施形態によれば、第 1 の NN パラメータは NN 2 0 のベース層表現を表し、第 2 の NN パラメータは NN 2 0 のエンハンスメント層表現を表す。代替的に、第 1 の NN パラメータは、例えば、NN 2 0 の現在の表現を表し、第 2 の NN パラメータは、現在の NN 表現の更新、すなわち、NN 2 0 の現在の表現の更新を表す。代替的に、例えば、第 1 の NN パラメータは、所定のニューラルネットワークニューロン 1 0 のインバウンドニューロン間活性化フィードフォワードの和にバイアスをかけるバイアス、すなわち b_i の成分を表し、第 2 の NN パラメータは、ニューラルネットワーク層 1 1 4、1 1 6₁ 又は 1 1 6₂ のアフィン変換をパラメータ化するバッチノルムパラメータ、すなわち μ、²、又は、例えば b + μ を表す。

【0099】

例 2 :

一実施形態によれば、NN パラメライザ 4 1 0 は、装置 3 0 0 を介して、第 3 の NN パラメータ及び第 4 の NN パラメータのうちの少なくとも 1 つを導出するように構成され、それにより、第 3 の NN パラメータは、第 3 の量子化値と第 3 の因数との間の積であって、第 3 のビット数だけビットシフトされた積に対応し、第 4 の NN パラメータは、第 4 の量子化値と第 4 の因数との間の積であって、第 4 のビット数だけビットシフトされた積に対応する。

【0100】

第 3 の量子化値及び第 4 の量子化値はともに、図 2 において 1 5 2 で示される量子化値を表す。第 3 の因数及び第 4 の因数は、両方とも、図 2 において 1 4 8 で示される因数を表す。

10

20

30

40

50

【 0 1 0 1 】

例えば、 $t = 2$ とし、 $k = 2^t$ とし、第1のQP、例えば、 QP_a で示される第3の量子化パラメータ142、関連する $shift_a$ 、すなわち第3のビットシフト数146、 mul_a 、すなわち第3の乗数144、及び γ_a 、すなわち第3の量子化ステップサイズ149を定義する。

【 0 1 0 2 】

さらに、第2のQP、例えば、 QP_b で示される第4の量子化パラメータ142、関連する $shift_b$ 、すなわち第4のビットシフト数146、 mul_b 、すなわち第4の乗数144、及び γ_b 、すなわち第4の量子化ステップサイズ149を定義する。

【 0 1 0 3 】

「第3の」パラメータ及び「第4の」パラメータは、この文脈では同じ参照番号で示されているが、それらが異なる値を有し得ることは明らかである。それらは、それらが図2に示されるどの特徴に属するかを明確にするために、同じ参照番号で示されているのみである。デバイス400は、第3のパラメータ及び/又は第4のパラメータのみ、又は、上記の例1で説明したように、更に第1のパラメータ及び/又は第2のパラメータを導出するように構成されてもよい。

【 0 1 0 4 】

$W = \gamma_a \cdot W_a$ が成り立つ量子化行列 W_a を考える。

$W = \gamma_b \cdot W_b$ が成り立つ量子化された転置ベクトル W_b を考える。

すなわち、 W_a は、 QP_a を使用して量子化され、 W_b は、 QP_b を使用して量子化されている。

図2で説明した量子化値152は、 W_a の1つの成分又は W_b の1つの成分を表すことができる。例えば、 W_a は複数の量子化値152を含むことができ、 W_b は複数の量子化値152を含むことができる。

【 0 1 0 5 】

さらに、要素ごとの積 $W \cdot \gamma$ が以下のように算出されると仮定する。

【 数 4 7 】

$$W \cdot \gamma = \Delta_a \cdot W_a \cdot \Delta_b \cdot \gamma_b = 2^{shift_a-2} \cdot mul_a \cdot W_a \cdot 2^{shift_b-2} \cdot mul_b \cdot \gamma_b$$

$$= 2^{shift_a+shift_b-4} \cdot mul_a \cdot mul_b \cdot W_a \cdot \gamma_b$$

【 0 1 0 6 】

この算出/計算は、例えば、第3のNNパラメータ W 及び第4のNNパラメータを乗算に施して、第3のNNパラメータ W の第3の量子化値 W_a によって形成される第1の因数と、第3の乗数 mul_a によって形成される第2の因数と、第4のNNパラメータの第4の量子化値 W_b によって形成される第3の因数と、第4の乗数 mul_b によって形成される第4の因数との積であって、第3のビット数 $shift_a$ によって形成される第1の加数と第4のビット数 $shift_b$ によって形成される第2の加数とを含む和に対応するビット数、例えば

【 数 4 8 】

$$2^{shift_a+shift_b-4}$$

だけビットシフトされた積を形成することによって積を得ることによって、演算ユニット420によって実行され得る。

10

20

30

40

50

【0107】

式から分かるように、浮動小数点演算を必要とし得る W 及び μ を導出する必要はない。代わりに、計算 $mul_a \cdot mul_b \cdot W_a \cdot \mu_b$ は、整数乗算のみを伴い、

【数49】

$$2^{shift_a + shift_b - 4}$$

との後続の乗算は、ビットシフトとして実施することができる。 $t = 2$ であるので、整数変数 mul_a 及び mul_b は両方とも値4、5、6、及び7のうちの1つであることに留意されたい。かかる小さい数を有する整数乗算は、ハードウェア又はソフトウェア実施態様において非常に効率的に実施され得る。

10

【0108】

一実施形態によれば、第3のNNパラメータは、第1のNN層114の第1のニューロン10₁から第2のNN層116₂の第2のニューロン10₂へのニューロン間活性化フィードフォワードを重み付けする重みパラメータ、例えば、 W の成分 w を表すか、又は代替的に、第3のNNパラメータは、第1のニューロン10₁と第2のニューロン10₂とを接続するエッジ12₁に関連し、第2のニューロン10₂のインバウンド活性化の和における第1のニューロン10₁の活性化の転送を重み付けする重みを表す。

【0109】

第4のNNパラメータは、例えば、バッチノルムパラメータ、例えば、 μ 、 μ^2 、又は μ を表す。バッチノルムパラメータは、例えば、第2のNN層116₁に対する第1のニューロン10₁の活性化フィードフォワード増幅を調整するものであり、例えばである。

20

【0110】

入力 X の量子化

一実施形態によれば、デバイス400は、例えば装置300を使用して、活性化を量子化された値、例えば X' に量子化することによって、活性化について第5の量子化パラメータ QP 、すなわち量子化パラメータ142、及び第5の量子化値、例えば X' 、すなわち量子化値152を決定することによって、NN入力 X 440を量子化するように構成され、それにより、第5の量子化パラメータ QP からの第5の乗数 mul 、すなわち乗数144の導出は、第5の量子化パラメータによって導出された被除数と、活性化に関連する精度パラメータ k 、すなわち精度パラメータ145によって導出された除数との間の除算の剰余と、除算の商の丸めに基づく第5のビットシフト数 $shift$ 、すなわちビットシフト数146とに基づいて、第5の量子化値と第5の乗数に依存する因数 mul/k 、すなわち因数148との間の積に対応する量子化された値が、第5のビットシフト数に依存する第5のビット数だけビットシフトされることをもたらす。

30

【0111】

好ましい実施形態において、バイアス層又はバッチ正規化層の入力 X 440も、本発明の量子化方法を使用して量子化される。図2の装置100の説明を参照。すなわち、 $X' = \mu \cdot X = mul \cdot 2^{shift-t} \cdot X$ が成り立つように、量子化パラメータ QP 及び関連する変数、関連する $shift$ 、 mul 、及び t ($t = 2$ 及び $k = 2^t$)が選択され、 X が X' に量子化される。次に、バイアス層又はバッチノルム層を実行するために X を使用する代わりに、 X' が入力として使用される。 X' は、通常、 X よりも極めて少ないビット/要素で表すことができ、これは、効率的なハードウェア又はソフトウェア実施態様のための別の利点であることに留意されたい。

40

【0112】

一実施形態によれば、NNパラメトラライザ410は、装置300を介して第6のNNパラメータを導出するように構成され、それにより、第6のNNパラメータは、第6の量子化値と第6の因数 mul/k との積であって、第6のビット数だけビットシフトさ

50

れた積に対する。デバイス400は、第6のNNパラメータ及び活性化に乗算を施して、第6のNNパラメータのための第6の量子化値によって形成される第1の因数と、第6の乗数によって形成される第2の因数と、第5の量子化値によって形成される第3の因数と、第5の乗数によって形成される第4の因数との積を形成することによって積であって、第6のビット数によって形成される第1の加数と第4のビット数によって形成される第2の加数とを含む和に対応するビット数だけビットシフトされた積を生成するように構成される。

【0113】

一実施形態によれば、第6のNNパラメータは、入力440を重み付けする重みパラメータWを表し、それにより、積 $W * X$ を算出/計算することができる。

10

【0114】

図2に戻って参照すると、以下では、装置100及び/又は装置300の更なる任意選択の特徴が説明される。

【0115】

パラメータQPの効率的な符号化及び復号化

好ましい実施形態において、パラメータQP、すなわち量子化パラメータ142は、以下の定義に従って、K次の符号付き指数ゴロム符号を使用して、装置100/装置300によってビットストリーム200内に符号化/ビットストリーム200から復号化される。

【0116】

別の好ましい実施形態は、次数Kが0に設定された先の好ましい実施形態と同じである。

20

【0117】

符号なし整数の指数ゴロム符号

符号なし整数の符号なし指数ゴロム符号は、高効率ビデオ符号化(HEVC: High Efficiency Video Coding)規格において定義されているシンタックス要素 $ue(v)$ の復号化仕様に従うものとする。

【0118】

この仕様を以下に簡単に説明する。

【0119】

次数Kの符号なし指数ゴロム符号で符号化された二値表現からの符号なし整数変数「decNum」の復号化は、以下の擬似符号に従って定義される。

30

$leadingZeroBits = -1$

$for (b = 0; !b; leadingZeroBits++)$

$b = read_bits(1)$

【0120】

次に、変数codeNumが以下のように割り当てられる。

$decNum = (2^{leadingZeroBits - 1}) * 2^K + read_bits(leadingZeroBits + K)$

【0121】

関数 $read_bits(x)$ は、ビットストリームからxビットを読み出し、それらを符号なし整数として返す。読み出されたビットは、最上位ビット(MSB)から最下位ビット(LSB)に順序付けられる。

40

【0122】

符号付き整数の指数ゴロム符号

符号付き整数の符号なし指数ゴロム符号は、高効率ビデオ符号化(HEVC)規格において定義されているシンタックス要素 $se(v)$ の復号化仕様に従うものとする。

【0123】

この仕様を以下に簡単に説明する。

【0124】

符号付き指数ゴロム符号で符号化された二値表現からの符号付き整数「signedD

50

「decNum」の復号化は、以下の通りである。最初に、符号なし整数が、上記で説明したようにHEVCのue(v)シンタックス要素復号化プロセスに従って復号化される。次に、符号なし整数は、以下の式に従って符号付き整数に変換される。

【数50】

$$\text{signedDecNum} = (-1)^{\text{decNum}+1} \cdot \lfloor \text{decNum}/2 \rfloor$$

シーリング演算子

10

【数51】

[x]

は、x以上の最小整数を返す。

【0125】

更に好ましい実施形態

好ましい実施形態において、パラメータk、すなわち精度パラメータ145は、 2^t に設定され、パラメータtは、bits_tビットを有する(例えば、bits_t = 3又はbits_t = 4を有する)符号なし整数表現を使用して符号化される。

20

【0126】

別の好ましい実施形態において、パラメータk、すなわち精度パラメータ145は、 2^t に設定され、パラメータtは、符号なし整数用の指数ゴロム符号を使用して符号化される。

【0127】

別の好ましい実施形態において、パラメータQP、すなわち量子化パラメータ142は、符号付き整数用の指数ゴロム符号を使用して符号化される。

【0128】

別の好ましい実施形態において、パラメータk、すなわち精度パラメータ145は、 2^t に設定され、パラメータQPは、bits_qpビットを使用して2の補数表現の符号付き整数を使用して符号化される。bits_qpは、例えば、12又は13のような一定値に設定されるか、又はbits_qpは、bits_qp0 + tに設定され、bits_qp0は、非ゼロの一定の整数値(例えば、bits_qp0 = 6)である。

30

【0129】

CABAC符号化ビットストリーム200の場合、パラメータt及び/又はQP142を表すビットは、(CABACのバイパスモードを使用して)バイパスピンとして符号化されるか、又はビットストリーム200に直接書き込まれるかのいずれかであり得る。

【0130】

別の好ましい実施形態において、パラメータW、b、 μ 、 2^2 、及びの各々は、パラメータの符号化の直前に符号化される個々のQP142の値で量子化される。

40

【0131】

別の好ましい実施形態において、第1のQP142がビットストリーム200に符号化され、モデルのパラメータのサブセットに関連付けられる。このサブセットの各パラメータxについて、1つのQPオフセットQP_xがパラメータごとに符号化され、パラメータを逆量子化するために使用される有効QP142、すなわちNNパラメータ120は、QP + QP_xとして与えられる。QP_xの二値表現は、好ましくは、QPの二値表現よりも少ないビットを使用する。例えば、QP_xは、符号付き整数又は(2の補数表現での)固定数のビットのための指数ゴロムコードを使用して符号化される。

【0132】

50

重みパラメータの符号化に関する更なる実施形態

図5に示される更なる好ましい実施形態は、重みパラメータ W_{545} の表現に関する。すなわち、それらをベクトル s_{546} と行列 W'_{544} との合成 $W = s \cdot W'$ として因数分解する。 W 及び W' 、すなわち重み行列 W_{544} は、次元 $n \times m$ の行列であり、 s は、長さ n の転置ベクトル s_{546} である。ベクトル s_{546} の各要素は、重み行列 W'_{544} の行方向のスケーリングファクターとして使用される。換言すれば、 s_{546} は、 W'_{544} の各列と要素ごとに乗算される。 s_{546} をローカルスケーリングファクター又はローカルスケール適応(LSA: local scale adaptation)と称する。

【0133】

図5は、NN20を使用して推論を実行するデバイス500を示している。デバイス500は、NN20を使用してNN入力440に基づいて推論出力430を計算するように構成される。NN20は、一对のNN層114及び116と、一对のNN層の第1の層114からNN層の第2の層116へのニューロン間活性化フィードフォワード122とを含む。デバイス500は、例えば、デバイス500の行列形成ユニット530を使用して、第1のNN層114のニューラルネットワークニューロン 10_1 の活性化520から行列 X_{532} を形成することによって、第1のNN層114のニューラルネットワークニューロン 10_1 の活性化520に基づいて第2のNN層116のニューラルネットワークニューロン 10_2 の活性化510を計算するように構成される。加えて、デバイス500は、 $s \cdot W' * X$ を計算(542)することによって、第1のNN層114のニューラルネットワークニューロン 10_1 の活性化520に基づいて、第2のNN層116のニューラルネットワークニューロン 10_2 の活性化510を計算するように構成され、ここで、 $*$ は、行列乗算を表し、 W' は、 n 及び m N である次元 $n \times m$ の重み行列 W_{544} であり、 s は、長さ n の転置ベクトル s_{546} であり、 \cdot は、 \cdot の一方の側の行列と \cdot の他方の側の転置ベクトルとの間の列に関するアダマール乗算を示す、デバイス500は、計算542を実行するように構成された計算ユニット540を備えることができる。

【0134】

一実施形態によれば、転置ベクトル s_{546} は、 W'_{544} を符号化するためのより高い圧縮及び/又はより高い推論忠実度に関して W'_{544} を最適化した結果である。

【0135】

その論理的根拠は、LSAが重み行列 W_{544} をスケーリングし、それにより、算術符号化方法がより高い符号化利得をもたらす、及び/又はニューラルネットワーク性能結果を増加させ、例えば、より高い精度を達成することである。例えば、 W の量子化後、 s_{546} は、入力データ440、例えば X_{532} を使用して又は使用せずに、量子化誤差を低減し、それにより量子化されたニューラルネットワークの予測性能を向上させるために適応され得る。

【0136】

したがって、 s_{546} 及び W'_{544} は、異なる量子化パラメータ、すなわち異なるQPを有することができる。これは、性能の観点からだけでなく、ハードウェア効率の観点からも有益であり得る。例えば、 W'_{544} は、入力 X_{532} との内積が8ビット表現で実行されてもよいが、スケーリングファクター s_{546} との後続の乗算が16ビットで実行され得るように量子化されてもよい。デバイス500は、例えば、内積を得るために n ビット固定小数点演算を使用して行列乗算 $W' * X$ を計算し、 $m > n$ である m ビット固定小数点演算を使用して内積を s_{546} と乗算するように構成される。

【0137】

しかしながら、 W'_{544} 及び s_{546} が両方とも n ビット表現に量子化される場合であっても、同じ推論精度を得るために W_{545} を量子化するのに必要な n よりも小さい n で十分な場合がある。同様に、 s_{546} が W'_{544} よりも少ないビットの表現に量子化された場合、表現の効率に関する利点を更に達成できる場合がある。

【0138】

一実施形態によれば、デバイス500は、NN表現110から W'_{544} を導出するよう

10

20

30

40

50

に構成されたNNパラメトライザー、例えば、図4に示すNNパラメトライザー410を備える。NNパラメトライザーは、NN表現110からNNパラメーターを導出する装置、例えば、図4又は図2に示される装置300を備える。重み行列W'544は、装置300によって導出されたNNパラメーターであり得る。任意選択で、NNパラメトライザー410は、W'544に関連するNNパラメーターと比較して異なる量子化パラメーター142を使用して、NN表現110からs546を導出するように更に構成される。

【0139】

好ましい実施形態において、重み行列W544の符号化は以下の通りである。第1に、LSAが使用されるかどうかを示すフラグが符号化される。フラグが1である場合、パラメーターs546及びW'544は、DeepCABACのような現行技術水準のパラメーター符号化方式を使用して符号化される。フラグが0である場合、W545が代わりに符号化される。

10

【0140】

別の好ましい実施形態において、前の好ましい実施形態による、異なるQP値がW'544及びs546に使用される。

【0141】

バッチノルム圧縮

図6に示す一実施形態は、バッチノルム圧縮を改善することに関する。図6は、NNのバッチノルム演算子710のNNパラメーター610、例えば、 μ 、 σ^2 、 γ 、 β 、及び任意選択でbをNN表現110に符号化する装置600と、NN表現110からNNのバッチノルム演算子710のNNパラメーター610、例えば、 μ 、 σ^2 、 γ 、 β 、及び任意選択でbを復号化する装置700とを示している。4つの実施形態が示されており、第1の実施形態は一般的な場合を説明し、他の実施形態は特別な場合を対象とする。

20

【0142】

概して、バッチノルム演算子710₁は、

【数52】

$$\frac{W * X + b - \mu}{\sqrt{\sigma^2 + \epsilon}} \cdot \gamma + \beta$$

30

として定義することができ、式中、

μ 、 σ^2 、 γ 、 β 、及びbは、バッチノルムパラメーター、例えば、各出力ノードについて1つの成分を含む転置ベクトルであり、

Wは、重み行列であり、例えば、その各行は1つの出力ノードに対するものであり、それぞれの行の各成分はXの1つの行に関連付けられており、

Xは、NN層の活性化から導出される入力行列であり、

bは、バイアスを形成する転置ベクトル、例えば、各出力ノードに対して1つの成分を含む転置ベクトルであり、

40

ϵ は、ゼロ除算回避のための定数であり、

\cdot は、 \cdot の一方の側の行列と他方の側の転置ベクトルとの間の列に関するアダマール乗算を示し、

*は、行列乗算を示す。

【0143】

第2の実施形態において、定数 ϵ は0であり、それにより、バッチノルム演算子710₂は、

【数53】

50

$$\frac{W \cdot X + b - \mu}{\sqrt{\sigma^2 + \epsilon}} \cdot \gamma + \beta$$

によって定義される。

【 0 1 4 4 】

第 3 の実施形態において、バイアス b は 0 であり、それにより、バッチノルム演算子 7 1 0 3 は、

【数 5 4】

10

$$\frac{W \cdot X - \mu}{\sqrt{\sigma^2 + \epsilon}} \cdot \gamma + \beta$$

によって定義される。

【 0 1 4 5 】

第 4 の実施形態において、バイアス b 及び定数 μ は 0 であり、それにより、バッチノルム演算子 7 1 0 4 は、

【数 5 5】

20

$$\frac{W \cdot X - \mu}{\sqrt{\sigma^2}} \cdot \gamma + \beta$$

によって定義される。

【 0 1 4 6 】

図 6 では、バッチノルム演算子 7 1 0 の一部のパラメータは、アポストロフィを有しており、アポストロフィなしのパラメータによって示される元のパラメータ 6 1 0 と、アポストロフィありのパラメータによって示される修正されたパラメータ 7 2 2、7 2 4、及び 7 3 2 との間の区別を可能にする。元のパラメータ 6 1 0 又は修正されたパラメータ 7 2 2、7 2 4 及び 7 3 2 のいずれかを、上記で定義されたバッチノルム演算子 7 1 0 のうちの 1 つのパラメータとして使用できることは明らかである。

30

【 0 1 4 7 】

装置 6 0 0 は、パラメータ μ 、 σ 、 ϵ 、及び γ を受信するように構成され（6 1 0 1 ~ 6 1 0 4 参照）、任意選択で b を受信するように構成される（6 1 0 1 及び 6 1 0 2 参照）。

【 0 1 4 8 】

第 1 の実施形態によれば、装置 6 0 0 は、

【数 5 6】

$$\beta' := \beta + \frac{(b - \mu) \cdot \gamma}{\sqrt{\sigma^2 + \epsilon}}$$

40

及び

【数 5 7】

$$\gamma' := \gamma \cdot \frac{\sqrt{\theta + \epsilon}}{\sqrt{\sigma^2 + \epsilon}}$$

を計算するように構成される。

50

【 0 1 4 9 】

代替の第 2 の実施形態によれば、装置 6 0 0 は、
【数 5 8】

$$\beta' := \beta + \frac{(b-\mu)*\gamma}{\sqrt{\sigma^2}} \text{ and}$$

及び

【数 5 9】

10

$$\gamma' := \gamma \cdot \frac{1}{\sqrt{\sigma^2}}.$$

を計算するように構成される。

【 0 1 5 0 】

代替の第 3 の実施形態によれば、装置 6 0 0 は、
【数 6 0】

20

$$\beta' := \beta - \frac{\mu*\gamma}{\sqrt{\sigma^2+\epsilon}}$$

及び

【数 6 1】

$$\gamma' := \gamma \cdot \frac{\sqrt{\theta+\epsilon}}{\sqrt{\sigma^2+\epsilon}}$$

30

を計算するように構成される。

【 0 1 5 1 】

代替の第 4 の実施形態によれば、装置 6 0 0 は、
【数 6 2】

$$\beta' := \beta - \frac{\mu*\gamma}{\sqrt{\sigma^2}}$$

40

及び

【数 6 3】

$$\gamma' := \gamma \cdot \frac{1}{\sqrt{\sigma^2}}.$$

を計算するように構成される。

【 0 1 5 2 】

計算されたパラメーター ' 及び ' は、バッチノルム演算子 7 1 0 の NN パラメーター

50

としてNN表現 1 1 0 に符号化され、例えば、同じ (' 及び ') が、各出力ノードに対して1つの成分を含む転置ベクトルでもあるようにする。

【 0 1 5 3 】

したがって、第 1 の実施形態のバッチノルム演算子 7 1 0 1 は、
【数 6 4 】

$$\frac{W \cdot X + b' - \mu'}{\sqrt{\sigma'^2 + \epsilon}} \cdot \gamma' + \beta'$$

10

として定義することができ、 $\sigma'^2 := \epsilon$ 、 $\mu' := 0$ 、及び $b' := 0$ であり、ここで、は所定のパラメータである。第 2 の実施形態のバッチノルム演算子 7 1 0 2 は、

【数 6 5 】

$$\frac{W \cdot X + b' - \mu'}{\sqrt{\sigma'^2}} \cdot \gamma' + \beta'$$

として定義することができ、 $\sigma'^2 := 1$ 、 $\mu' = 0$ 、及び $b' = 0$ である。第 3 の実施形態のバッチノルム演算子 7 1 0 3 は、

20

【数 6 6 】

$$\frac{W \cdot X - \mu'}{\sqrt{\sigma'^2 + \epsilon}} \cdot \gamma' + \beta'$$

として定義することができ、 $\sigma'^2 := \epsilon$ 及び $\mu' := 0$ であり、ここで、は所定のパラメータである。第 4 の実施形態のバッチノルム演算子 7 1 0 4 は、

【数 6 7 】

30

$$\frac{W \cdot X - \mu'}{\sqrt{\sigma'^2}} \cdot \gamma' + \beta'$$

として定義することができ、 $\sigma'^2 := 1$ 及び $\mu' := 0$ である。

【 0 1 5 4 】

所定のパラメータは 1 又は $1 - \epsilon$ であり、例えば、ここでも μ' 、 σ'^2 、 β' 、及び γ' は各出力ノードに対して1つの成分を含む転置ベクトルであり、Wは、重み行列であり、XはNN層の活性化から導出される入力行列であり、 b' はバイアスを形成する転置ベクトル、例えば、各出力ノードに対して1つの成分を含む転置ベクトルである。

40

【 0 1 5 5 】

装置 7 0 0 は、例えば、装置 7 0 0 に含まれ得る γ' 及び 導出ユニット 7 2 0 を使用することによって、NN表現から γ' 及び β' 、すなわち γ' 及び β' を導出するように構成される。

【 0 1 5 6 】

第 1 の実施形態によれば、装置 7 0 0 は、その全ての成分に適用される1つのシグナリング 7 3 4 を介して、 $\sigma'^2 := \epsilon$ 、 $\mu' := 0$ 、及び $b' := 0$ を推論又は導出するように構成され、ここで、は所定のパラメータである。

【 0 1 5 7 】

第 2 の実施形態によれば、装置 7 0 0 は、その全ての成分に適用される1つのシグナリ

50

ング734を介して、 $\sigma^2 := 1$ 、 $\mu' := 0$ 、及び $b' := 0$ であると推論又は導出するように構成される。

【0158】

第3の実施形態によれば、装置700は、その全ての成分に適用される1つのシグナリング734によって、 $\sigma^2 :=$ 及び $\mu' := 0$ を推論又は導出するように構成され、ここで、 σ^2 は所定のパラメータである。

【0159】

第4の実施形態によれば、装置700は、その全ての成分に適用される1つのシグナリング734によって、 $\sigma^2 := 1$ 及び $\mu' := 0$ であると推論又は導出するように構成される。

10

【0160】

パラメータ σ^2 、 μ' 、及び任意選択で b' のこの導出又は推論は、パラメータ推論/導出ユニット730を使用して実行され得る。

【0161】

所定のパラメータは1又は $1 -$ であり、例えば、ここでも μ' 、 σ^2 、 σ' 、及び σ'' は各出力ノードに対して1つの成分を含む転置ベクトルであり、 W は、重み行列であり、 X は NN 層の活性化から導出される入力行列であり、 b' はバイアスを形成する転置ベクトル、例えば、各出力ノードに対して1つの成分を含む転置ベクトルである。

【0162】

図6では、装置700によって導出又は推論されたパラメータはアポストロフィによって示されているが、装置700が元のパラメータ610を見ることがないという事実により、装置700によって導出又は推論されたパラメータはアポストロフィを使用せずに示されてもよい。装置700を考慮すると、導出又は推論されたパラメータは、唯一の既存のパラメータである。

20

【0163】

任意選択で、装置700は、例えば推論のために、導出又は推論されたパラメータ722、724、及び732とともにバッチノルム演算子を使用するように構成され得る。バッチノルム演算子計算ユニットは、バッチノルム演算子を使用するように構成され得る。代替的に、推論のデバイス、例えばデバイス400又はデバイス500は、バッチノルム演算子710のパラメータを取得する装置700を備えてもよい。

30

【0164】

定数スカラー値 ϵ 、すなわち、例えば、1又は $1 -$ に等しくなり得る所定のパラメータを導入すると、パラメータ b 、 μ 、 σ^2 、 σ' 、及び σ'' は、 $BN(X)$ 、すなわち、バッチノルム演算子710の結果を変更することなく、以下の順序付けられたステップによって修正することができる。

【数68】

1) $\beta := \beta + \frac{(b-\mu)*\gamma}{\sqrt{\sigma^2+\epsilon}}$

2) $\gamma := \gamma \cdot \frac{\sqrt{\theta+\epsilon}}{\sqrt{\sigma^2+\epsilon}}$

40

3) $\sigma^2 :=$

4) $\mu := 0$

5) $b := 0$

【0165】

各演算は、転置されたベクトルの要素に対する要素ごとの演算として解釈される。実施形態2~3に例示されるように、 $BN(X)$ を変更しない更なる修正も可能である。例えば、バイアス b 及び平均 μ は σ^2 に「積分」され、それにより、 b 及び μ はその後0に設定

50

される（第3の実施形態を参照）。又は、 μ^2 は、他のパラメーターがそれに応じて調整されるとき、BN(X)における分数の分母を1に等しく設定するために、 $1 - \mu$ （すなわち、 $\mu = 1 - \mu$ ）に設定され得る。

【0166】

それにより、全てのベクトル要素が同じ値を有するので、 b 、 μ^2 、 μ 、及び b を極めてより効率的に圧縮することができる。

【0167】

好ましい実施形態において、パラメーターの全ての要素が所定の定数値を有するかどうかを示すフラグ734が符号化される。パラメーターは、例えば、 b 、 μ 、 μ^2 、 μ 、又は μ^2 であってもよい。所定の値は、例えば、0、1、又は $1 - \mu$ であってもよい。フラグが1に等しい場合、パラメーターの全てのベクトル要素は、所定の値に設定される。そうでなければ、パラメーターは、例えばDeepCABACのような現行技術水準のパラメーター符号化方法の1つを使用して符号化される。

10

【0168】

別の好ましい実施形態において、全てのベクトル要素が同じ値を有するかどうかを示すフラグがパラメーターごとに符号化される。全てのベクトル要素が同じ値を有するとき、フラグは1に等しく、その値は、例えばDeepCABAC、又は及び指数ゴロム符号、又は固定長符号のような現行技術水準のパラメーター符号化方法を使用して符号化される。フラグが0である場合、パラメーターのベクトル要素は、例えばDeepCABACのような現行技術水準のパラメーター符号化方法の1つを使用して符号化される。

20

【0169】

一実施形態によれば、装置600/装置700は、表現110において/表現110から、 μ^2 の全ての成分（例えば、各成分は、対応する出力ノードを意味するWの対応する行に対するものである）が互いに等しいこと、及びその値を示す/導出するように構成される。加えて、又は代替として、装置600/装置700は、表現110において/表現110から、 μ' の全ての成分（例えば、各成分は、対応する出力ノードを意味するWの対応する行に対するものである）が互いに等しいこと、及びその値を示す/導出するように構成される。加えて、又は代替として、装置600/装置700は、存在する場合、例えば第1の実施形態及び第2の実施形態の場合であるが第3の実施形態及び第4の実施形態の場合ではない場合、表現119において/表現119から、 b' の全ての成分（例えば、各成分は、対応する出力ノードを意味するWの対応する行に対するものである）が互いに等しいこと、及びその値を示す/導出するように構成される。

30

【0170】

一実施形態によれば、装置600は、2つのバッチノルム符号化モードの間で切り替え可能であるように更に構成され、第1のバッチノルム符号化モードでは、装置600は、 μ' 及び μ^2 の計算及び符号化を実行するように構成され、第2のバッチノルム符号化モードでは、装置は、受信された μ 、 μ^2 又は μ 、 μ^2 、及び μ 、並びに存在する場合、 b を符号化するように構成される。換言すれば、受信されたパラメーター610は、第2のバッチノルムモードで表現110に直接符号化される。並行して、装置700も、2つのバッチノルム符号化モード間で切り替え可能であるように構成してもよく、第1のバッチノルム符号化モードでは、装置700は、導出することと、推論又は導出することとを実行するように構成され、第2の第1のバッチノルム符号化モードでは、装置700は、 μ 、 μ^2 又は μ 、 μ^2 、及び μ 、並びに存在する場合、 b を表現110から復号化するように構成される。換言すれば、パラメーター610は、第2のバッチノルムモードで表現110から直接復号化される。

40

【0171】

一実施形態によれば、装置600は、 μ' 及び μ^2 をNN表現110に量子化及び符号化する装置100（図2参照）を備える。例えば、装置600は、最初に計算620を実行し、得られたパラメーター μ' 及び μ^2 を、パラメーターの量子化のために装置100に渡す。一実施形態によれば、装置700は、NN表現110から μ' 及び μ^2 を導出する装置3

50

0 0 (図 2 参 照) を 備 え る。

【 0 1 7 2 】

理 解 を 容 易 に す る た め に、 X 及 び W と 層 の 対 と の 間 の 可 能 な 関 係 が 図 7 に 示 さ れ て お り、 左 は 全 結 合 層 $i + 1$ で あり、 右 は 畳 み 込 み 層 $i + 1$ で あり、 層 の ニ ュ ー ロ ン は 円 1 0 で 示 さ れ て い る。 各 層 の ニ ュ ー ロ ン は、 ア レ イ 位 置 (x, y) に 配 置 さ れ る。 各 層 i は、 q_i 列 の ニ ュ ー ロ ン 1 0 と p_i 行 の ニ ュ ー ロ ン 1 0 と を 有 す る。 全 結 合 の 場 合、 X_i は、 成 分

【 数 6 9 】

$$X_{1...p_i:q_i}$$

10

の ベ ク ト ル で あり、 各 X_g は、 位 置

【 数 7 0 】

$$\{[g/q_i];g\%q_i+1\}$$

に お け る ニ ュ ー ロ ン の 活 性 化 で 占 め ら れ、 W_i は、 成 分

【 数 7 1 】

$$W_{1...p_{i+1}:q_{i+1},1...p_i:q_i}$$

20

の 行 列 で あり、 各 $W_{g, h}$ は、 位 置

【 数 7 2 】

$$\{[g/q_{i+1}];g\%q_{i+1}+1\}$$

30

に お け る 層 $i + 1$ の ニ ュ ー ロ ン 1 0 と 位 置

【 数 7 3 】

$$\{[h/q_i];h\%q_i+1\}$$

40

に お け る 層 i の ニ ュ ー ロ ン 1 0 と の 間 の エ ッ ジ 1 2 の 重 み で 占 め ら れ る。 畳 み 込 み の 場 合、 X_i は、 成 分

【 数 7 4 】

$$X_{1...r:s,1...p_{i+1}:q_{i+1}}$$

の 行 列 で あり、 こ こ で、 各 $X_{g, h}$ は、 位 置

【 数 7 5 】

50

$$\{[(g + (h - 1) * q_i / (q_{i+1} + s - 1)) / s]; (g + (h - 1) * q_i / (q_{i+1} + s - 1)) \% s + 1\}$$

におけるニューロンの活性化で占められ、 W_i は、成分 $W_{1 \dots r \cdot s}$ のベクトルであり、ここで、各 $W_{g, h}$ は、層 i にわたって分散された $p_{i+1} \cdot q_{i+1}$ 位置のうちの1つに配置された層 i 中のサイズ $r \times s$ の矩形フィルターカーネル中のニューロンから、カーネル位置に対応する層 $i + 1$ 中のニューロン位置につながるエッジの重みで占められる。

【0173】

一部の態様を装置の文脈で説明してきたが、これらの態様に対応する方法の説明も表すことは明らかであり、ブロック又はデバイスは方法ステップ又は方法ステップの特徴に対応する。同様に、方法ステップの文脈で説明される態様はまた、対応する装置の対応するブロック又は項目又は特徴の説明を表す。方法ステップの一部又は全部は、例えばマイクロプロセッサ、プログラマブルコンピューター又は電子回路のようなハードウェア装置によって（又はそれを使用して）実行されてもよい。一部の実施形態において、最も重要な方法ステップのうちの1つ以上は、かかる装置によって実行されてもよい。

【0174】

本発明のNN表現を含む本発明のデジタルデータ、データストリーム、又はファイルは、デジタル記憶媒体上に記憶することができ、又は無線伝送媒体若しくはインターネット等の有線伝送媒体等の伝送媒体上で伝送することができる。

【0175】

特定の実装要件に応じて、本発明の実施形態は、ハードウェア又はソフトウェアで実行することができる。実施態様は、それぞれの方法が実行されるようにプログラム可能なコンピューターシステムと協働する（又は協働することができる）電子的可読制御信号が記憶されたデジタル記憶媒体、例えば、フロッピーディスク、DVD、Blu-ray（登録商標）、CD、ROM、PROM、EPROM、EEPROM又はフラッシュメモリを使用して実行することができる。したがって、デジタル記憶媒体はコンピューター可読であってもよい。

【0176】

本発明による一部の実施形態は、電子的可読制御信号を有するデータキャリアを備え、該電子的可読制御信号は、本明細書で説明される方法のうちの1つが実行されるように、プログラム可能なコンピューターシステムと協働することが可能である。

【0177】

概して、本発明の実施形態は、プログラムコードを有するコンピュータープログラム製品として実施することができ、プログラムコードは、コンピュータープログラム製品がコンピューター上で実行されるときに方法のうちの1つを実行するように動作可能である。プログラムコードは、例えば、機械可読キャリアに記憶されてもよい。

【0178】

他の実施形態は、機械可読キャリア上に記憶された、本明細書で説明される方法のうちの1つを実行するコンピュータープログラムを含む。

【0179】

したがって、換言すれば、本発明の方法の一実施形態は、コンピュータープログラムがコンピューター上で実行されるときに、本明細書で説明される方法のうちの1つを実行するプログラムコードを有するコンピュータープログラムである。

【0180】

したがって、本発明の方法の更なる実施形態は、本明細書に記載の方法のうちの1つを実行するコンピュータープログラムを記録したデータキャリア（又はデジタル記憶媒体、又はコンピューター可読媒体）である。データキャリア、デジタル記憶媒体又は記録された媒体は、典型的には有形及び/又は非一時的である。

10

20

30

40

50

【0181】

したがって、本発明の方法の更なる実施形態は、本明細書に記載された方法の1つを実行するコンピュータプログラムを表すデータストリーム又はシグナルのシーケンスである。データストリーム又はシグナルのシーケンスは、例えば、データ通信接続を介して、例えばインターネットを介して転送されるように構成されてもよい。

【0182】

更なる実施形態は、本明細書に記載の方法の1つを実行するように構成又は適合された処理手段、例えばコンピュータ又はプログラム可能な論理デバイスを含む。

【0183】

更なる実施形態は、本明細書に記載の方法のうちの1つを実行するコンピュータプログラムがインストールされたコンピュータを含む。

10

【0184】

本発明による更なる実施形態は、本明細書で説明される方法のうちの1つを実行するコンピュータプログラムを受信機に（例えば、電子的に又は光学的に）転送するように構成された装置又はシステムを含む。受信機は、例えば、コンピュータ、モバイルデバイス、メモリデバイス等であってもよい。装置又はシステムは、例えば、コンピュータプログラムを受信機に転送するためのファイルサーバを備えることができる。

【0185】

一部の実施形態において、プログラム可能論理デバイス（例えば、フィールドプログラマブルゲートアレイ）が、本明細書に説明される方法の機能性の一部又は全部を行うために使用されてもよい。一部の実施形態において、フィールドプログラマブルゲートアレイは、本明細書に説明される方法のうちの1つを行うために、マイクロプロセッサと協働してもよい。概して、方法は、任意のハードウェア装置によって実行されることが好ましい。

20

【0186】

本明細書で説明される装置は、ハードウェア装置を使用して、又はコンピュータを使用して、又はハードウェア装置とコンピュータとの組み合わせを使用して実施され得る。

【0187】

本明細書で説明される装置、又は本明細書で説明される装置の任意の構成要素は、少なくとも部分的にハードウェア及び/又はソフトウェアで実施され得る。

【0188】

本明細書で説明される方法は、ハードウェア装置を使用して、又はコンピュータを使用して、又はハードウェア装置とコンピュータとの組み合わせを使用して実行され得る。

30

【0189】

本明細書で説明される方法、又は本明細書で説明される装置の任意の構成要素は、少なくとも部分的にハードウェア及び/又はソフトウェアによって実行され得る。

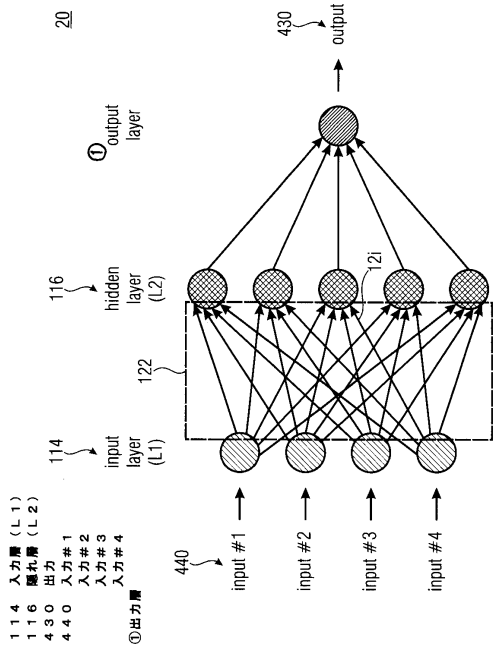
【0190】

上述の実施形態は、本発明の原理を単に例示するものである。本明細書に記載された構成及び詳細の変更及び変形が他の当業者に明らかであることが理解される。したがって、添付の特許請求の範囲によってのみ限定され、本明細書の実施形態の記述及び説明によって提示される特定の詳細によって限定されないことが意図される。

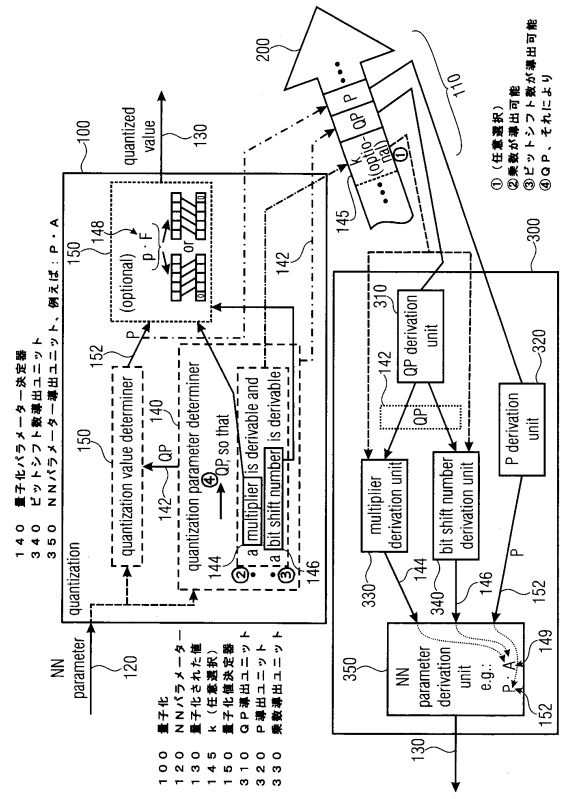
40

【図面】

【図 1】



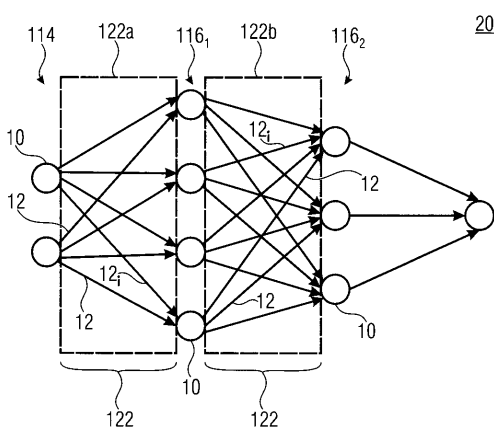
【図 2】



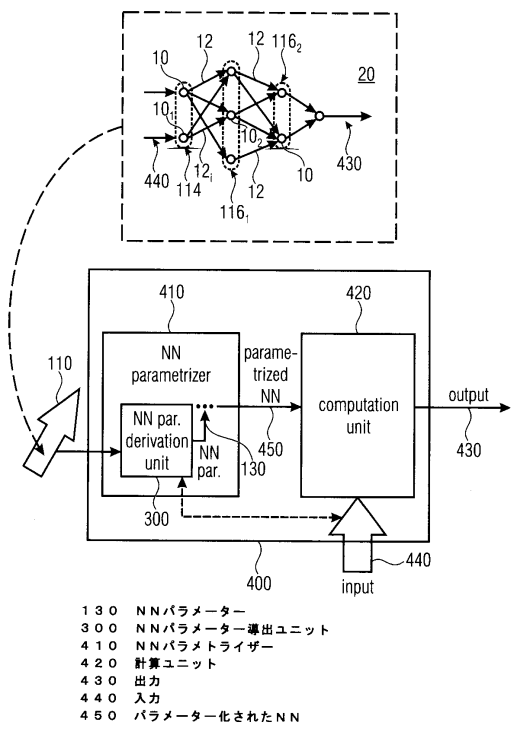
10

20

【図 3】



【図 4】

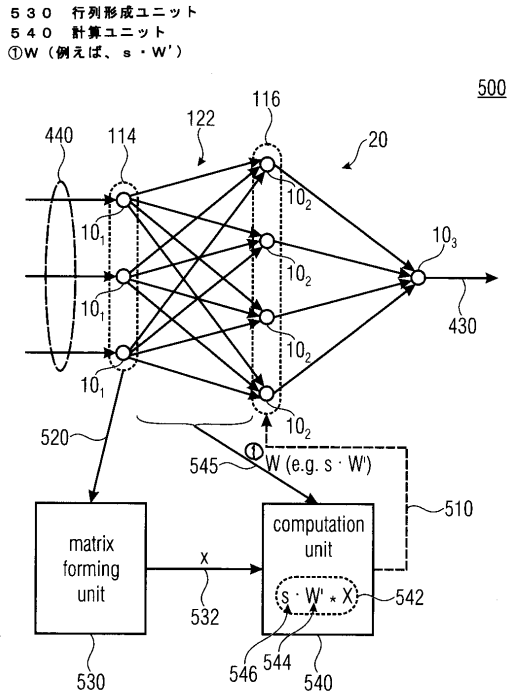


30

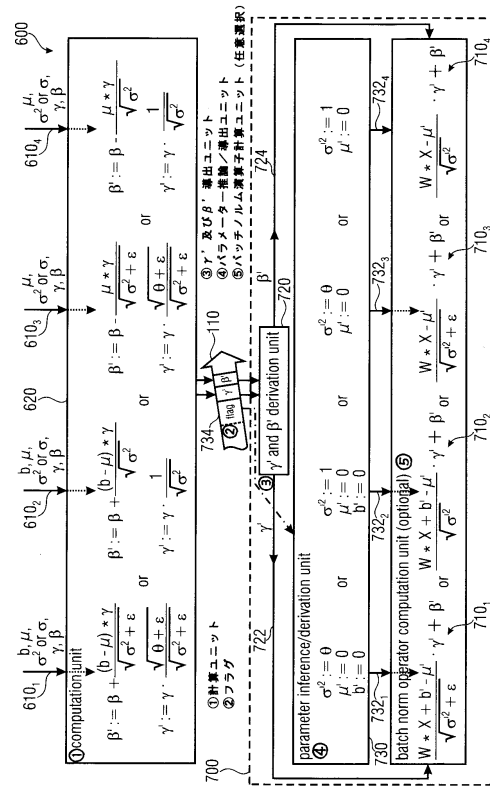
40

50

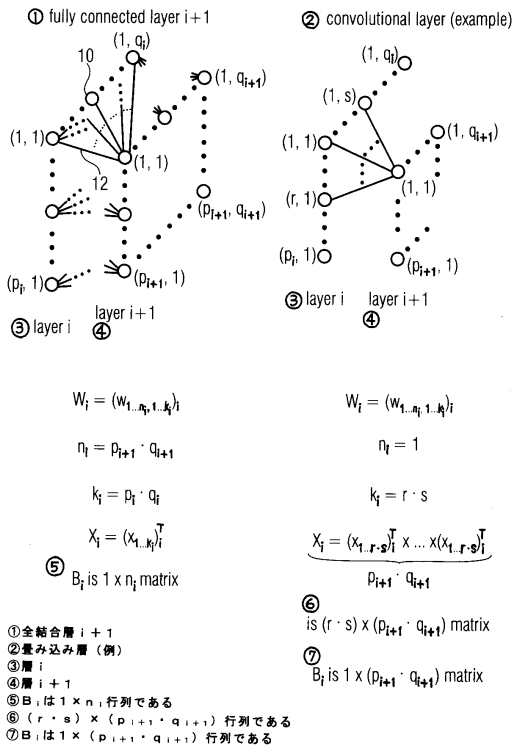
【 図 5 】



【 図 6 】



【 図 7 】



10

20

30

40

50

フロントページの続き

- スティテュート, ツェー/オー フラウンホーファー - インスティテュート フュア ナッハリヒ
テンテヒニーク
- (72)発明者 マーリンチ, タルマイ
ドイツ連邦共和国 ベルリン 10587 アインシュタインウーファー37 ハーハーイー, ハイ
ンリッヒ - ヘルツ - インスティテュート, ツェー/オー フラウンホーファー - インスティテュ
ート フュア ナッハリヒテンテヒニーク
- (72)発明者 サメック, ヴォイチェフ
ドイツ連邦共和国 ベルリン 10587 アインシュタインウーファー37 ハーハーイー, ハイ
ンリッヒ - ヘルツ - インスティテュート, ツェー/オー フラウンホーファー - インスティテュ
ート フュア ナッハリヒテンテヒニーク
- (72)発明者 ハーゼ, パウル
ドイツ連邦共和国 ベルリン 10587 アインシュタインウーファー37 ハーハーイー, ハイ
ンリッヒ - ヘルツ - インスティテュート, ツェー/オー フラウンホーファー - インスティテュ
ート フュア ナッハリヒテンテヒニーク
- (72)発明者 ミュラー, カーステン
ドイツ連邦共和国 ベルリン 10587 アインシュタインウーファー37 ハーハーイー, ハイ
ンリッヒ - ヘルツ - インスティテュート, ツェー/オー フラウンホーファー - インスティテュ
ート フュア ナッハリヒテンテヒニーク
- (72)発明者 キルヒホフファー, ハイナー
ドイツ連邦共和国 ベルリン 10587 アインシュタインウーファー37 ハーハーイー, ハイ
ンリッヒ - ヘルツ - インスティテュート, ツェー/オー フラウンホーファー - インスティテュ
ート フュア ナッハリヒテンテヒニーク
- (72)発明者 マーブ, デトレフ
ドイツ連邦共和国 ベルリン 10587 アインシュタインウーファー37 ハーハーイー, ハイ
ンリッヒ - ヘルツ - インスティテュート, ツェー/オー フラウンホーファー - インスティテュ
ート フュア ナッハリヒテンテヒニーク
- (72)発明者 シュヴァルツ, ハイコ
ドイツ連邦共和国 ベルリン 10587 アインシュタインウーファー37 ハーハーイー, ハイ
ンリッヒ - ヘルツ - インスティテュート, ツェー/オー フラウンホーファー - インスティテュ
ート フュア ナッハリヒテンテヒニーク
- (72)発明者 ヴィーガント, トーマス
ドイツ連邦共和国 ベルリン 10587 アインシュタインウーファー37 ハーハーイー, ハイ
ンリッヒ - ヘルツ - インスティテュート, ツェー/オー フラウンホーファー - インスティテュ
ート フュア ナッハリヒテンテヒニーク
- 審査官 大倉 峻吾
- (56)参考文献 国際公開第2019/008752 (WO, A1)
JACOB, Benoit et al., "Quantization and Training of Neural Networks for Efficient Integer-A
rithmetic-Only Inference", arXiv [online], 2017年12月, [2023年10月11日検索], イン
ターネット <URL: <https://arxiv.org/abs/1712.05877v1>>, 1712.05877v1
FLYNN, David et al., "G-PCC: Integer step sizes for in-tree geometry quantisation", MPEG
Point Cloud Compression [online], 2020年01月, [2023年10月11日検索], インターネ
ット <URL: [https://mpeg-pcc.org/index.php/public-contributions/m52522-g-pcc-intege
r-step-sizes-for-in-tree-geometry-quantisation/](https://mpeg-pcc.org/index.php/public-contributions/m52522-g-pcc-intege
r-step-sizes-for-in-tree-geometry-quantisation/)>, m52522
CAI, Yaohui et al., "ZeroQ: A Novel Zero Shot Quantization Framework", arXiv [online], 2
020年01月, [2023年10月11日検索], インターネット <URL: [https://arxiv.org/abs/
2001.00281v1](https://arxiv.org/abs/
2001.00281v1)>, 2001.00281v1
- (58)調査した分野 (Int.Cl., DB名)
G06N 3/02 - 3/10
G06N 20/00 - 99/00