

(19) 日本国特許庁(JP)

(12) 公開特許公報(A)

(11) 特許出願公開番号

特開2018-92377

(P2018-92377A)

(43) 公開日 平成30年6月14日(2018.6.14)

(51) Int.Cl.

G06N 3/063 (2006.01)

F I

G06N 3/063

テーマコード (参考)

審査請求 有 請求項の数 11 O L (全 23 頁)

(21) 出願番号 特願2016-235383 (P2016-235383)
 (22) 出願日 平成28年12月2日 (2016.12.2)
 (11) 特許番号 特許第6183980号 (P6183980)
 (45) 特許公報発行日 平成29年8月23日 (2017.8.23)

(71) 出願人 304021417
 国立大学法人東京工業大学
 東京都目黒区大岡山2丁目12番1号
 (74) 代理人 110001807
 特許業務法人磯野国際特許商標事務所
 (72) 発明者 中原 啓貴
 東京都目黒区大岡山2-12-1 国立大
 学法人東京工業大学内
 (72) 発明者 米川 晴義
 東京都目黒区大岡山2-12-1 国立大
 学法人東京工業大学内

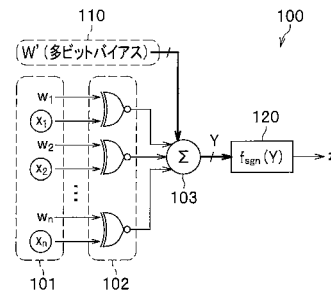
(54) 【発明の名称】 ニューラルネットワーク回路装置、ニューラルネットワーク、ニューラルネットワーク処理方法
 およびニューラルネットワークの実行プログラム

(57) 【要約】

【課題】 バッチ正規化回路が不要なニューラルネットワ
 ーク回路装置、ニューラルネットワーク、ニューラルネ
 ットワーク処理方法およびニューラルネットワークの実
 行プログラムを提供する。

【解決手段】 2値化ニューラルネットワーク回路100
 は、入力値 $x_1 \sim x_n$ (x_i) (2値) を入力する入力
 ノードおよび重み $w_1 \sim w_n$ (w_i) を入力する入力部
 101と、入力値 $x_1 \sim x_n$ および重み $w_1 \sim w_n$ を受
 け取り、XNOR論理を取るXNORゲート回路102
 と、多ビットバイアス W' を入力する多ビットバイアス
 W' 入力部110と、各XNOR論理値と多ビットバイ
 アス W' との総和を取る総和回路103と、総和を取っ
 た信号 Y に対して符号ビットのみを出力する活性化回路
 120と、を備える。

【選択図】 図9



【特許請求の範囲】

【請求項 1】

入力層、1 以上の中間層、および、出力層を少なくとも含むニューラルネットワークにおいて、前記中間層の中で、入力値に重みづけとバイアスを乗算するニューラルネットワーク回路装置であって、

入力値 x_i および重み w_i を受け取り、論理演算を行う論理回路部と、

多ビットバイアス W' を受け取り、前記論理回路部の出力と前記多ビットバイアス W' との総和を取る総和回路部と、

総和を取った多ビット信号 Y に対して符号ビットのみを出力する活性化回路部と、を備える

10

ことを特徴とする記載のニューラルネットワーク回路装置。

【請求項 2】

前記入力値 x_i および前記重み w_i を入力する入力部と、

前記多ビットバイアス W' を入力する多ビットバイアス入力部と、を備える

ことを特徴とする請求項 1 に記載のニューラルネットワーク回路装置。

【請求項 3】

前記入力値 x_i および前記重み w_i は、2 値信号である

ことを特徴とする請求項 1 または請求項 2 に記載のニューラルネットワーク回路装置。

【請求項 4】

前記多ビットバイアス W' は、学習後の多ビットバイアス値である

ことを特徴とする請求項 1 または請求項 2 に記載のニューラルネットワーク回路装置。

20

【請求項 5】

前記論理回路部は、否定排他的論理和または排他的論理和を含む

ことを特徴とする請求項 1 に記載のニューラルネットワーク回路装置。

【請求項 6】

前記論理回路部は、LUT (Look-Up Table) である

ことを特徴とする請求項 1 に記載のニューラルネットワーク回路装置。

【請求項 7】

前記符号ビットは、総和を取った前記多ビット信号 Y を活性化するかしないかで示す 2 値信号である

30

ことを特徴とする請求項 1 に記載のニューラルネットワーク回路装置。

【請求項 8】

前記多ビット信号 Y および多ビットバイアス W' は、下記式で示される

【数 4】

$$Y = \sum_{i=1}^n w_i x_i + W'$$

ことを特徴とする請求項 1 に記載のニューラルネットワーク回路装置。

【請求項 9】

請求項 1 乃至 8 のいずれか 1 項に記載のニューラルネットワーク回路装置を備えるニューラルネットワーク。

40

【請求項 10】

入力層、1 以上の中間層、および、出力層を少なくとも含むニューラルネットワークにおいて、前記中間層の中で、入力値に重みづけとバイアスを乗算するニューラルネットワーク処理方法であって、

入力値 x_i および重み w_i を受け取り、論理演算を行うステップと、

多ビットバイアス W' を受け取り、前記論理回路部の出力と前記多ビットバイアス W' との総和を取るステップと、

総和を取った多ビット信号 Y に対して符号ビットのみを出力するステップと、を有する

50

ことを特徴とするニューラルネットワーク処理方法。

【請求項 1 1】

入力層、1 以上の中間層、および、出力層を少なくとも含むニューラルネットワークにおいて、前記中間層の中で、入力値に重みづけとバイアスを乗算するニューラルネットワーク回路装置としてのコンピュータを、

入力値 x_i および重み w_i を受け取り、論理演算を行う論理回路手段、

多ビットバイアス W' を受け取り、前記論理回路部の出力と前記多ビットバイアス W' との総和を取る総和回路手段、

総和を取った多ビット信号 Y に対して符号ビットのみを出力する活性化回路手段、

として機能させるためのニューラルネットワークの実行プログラム。

10

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、ニューラルネットワーク回路装置、ニューラルネットワーク、ニューラルネットワーク処理方法およびニューラルネットワークの実行プログラムに関する。

【背景技術】

【0002】

古典的な順伝搬型ニューラルネットワーク (FFNN: Feedforward Neural Network)、RBF (Radial Basis Function) ネットワーク、正規化した RBF ネットワーク、自己組織化マップなどがある。RBFN は、誤差逆伝搬法に用いる活性化関数に放射基底関数を用いる。しかし、中間層が多く取れず高精度認識判定が難しかったり、HW 規模が大きく処理時間がかかる、などの問題があり手書き文字認識など応用分野が限定されていた。

20

近年、ADAS (advanced driver assistance system) 用の画像認識や自動翻訳などで注目を集める新方式として畳み込みニューラルネットワーク (CNN: Convolutional Neural Network) (層間が全結合でない NN) や再帰型ニューラルネットワーク (双方向伝搬) が登場している。CNN は、ディープニューラルネットワーク (DNN: Deep Neural Network) に畳込み演算を付加したものである。

【0003】

特許文献 1 には、誤り訂正符号の検査行列に基づいて、階層型ニューラルネットワークにおける疎結合のノード間で学習された重みの値と入力信号とを用いて、問題を解く処理部を備える処理装置が記載されている。

30

【0004】

既存の CNN は、短精度 (多ビット) による積和演算回路で構成されており、多数の乗算回路が必要である。このため、面積・消費電力が多くなる欠点があった。そこで、2 値化した精度、すなわち +1 と -1 のみ用いて CNN を構成する回路が提案されている (例えば、非特許文献 1 ~ 4 参照)。

【先行技術文献】

【特許文献】

【0005】

40

【特許文献 1】特開 2016 - 173843 号公報

【非特許文献】

【0006】

【非特許文献 1】M. Courbariaux, I. Hubara, D. Soudry, R.E. Yaniv, Y. Bengio, "Binarized neural networks: Training deep neural networks with weights and activations constrained to +1 or -1," Computer Research Repository (CoRR), 「2 値化 NN のアルゴリズム」、[online]、2016 年 3 月、[平成 28 年 10 月 5 日検索]、<URL: <http://arxiv.org/pdf/1602.02830v3.pdf> >

【非特許文献 2】Mohammad Rastegari, Vicente Ordonez, Joseph Redmon, Ali Farhadi, "XNOR-Net: ImageNet Classification Using Binary Convolutional Neural Networks,

50

" Computer Vision and Pattern recognition, 「2 値化 NN のアルゴリズム」、[online]、2016 年 3 月、[平成 28 年 10 月 5 日検索]、<URL: <https://arxiv.org/pdf/1603.05279v4> >

【非特許文献 3】Hiroki Nakahara, Haruyoshi Yonekawa, Tsutomu Sasao, Hisashi Iwamoto and Masato Motomura, " A Memory-Based Realization of a Binarized Deep Convolutional Neural Network, " Proc. of the 2016 International Conference on Field-Programmable Technology (FPT), Xi'an, China, Dec 2016 (To Appear).

【非特許文献 4】Eriko Nurvitadhi, David Sheffield, Jaewoong Sim, Asit Mishra, Ganes Venkatesh, Debbie Marr, " Accelerating Binarized Neural Networks: Comparison of FPGA, CPU, GPU, and ASIC, " Proc. of the 2016 International Conference on Field-Programmable Technology (FPT), Xi'an, China, Dec 2016 (To Appear).

10

【発明の概要】

【発明が解決しようとする課題】

【0007】

非特許文献 1 ~ 4 の技術では、精度を 2 値に落とすことで CNN の認識精度も落としてしまう。これを避けて 2 値化 CNN の精度を維持するためには、バッチ正規化回路が必要であるが、バッチ正規化回路は、複雑な回路であり、面積・消費電力が増大するという課題があった。

【0008】

本発明は、このような事情に鑑みてなされたものであり、バッチ正規化回路が不要なニューラルネットワーク回路装置、ニューラルネットワーク、ニューラルネットワーク処理方法およびニューラルネットワークの実行プログラムを提供することを課題とする。

20

【課題を解決するための手段】

【0009】

前記した課題を解決するため、本発明に係るニューラルネットワーク回路装置は、入力層、1 以上の中間層、および、出力層を少なくとも含むニューラルネットワークにおいて、前記中間層の中で、入力値に重みづけとバイアスを乗算するニューラルネットワーク回路装置であって、入力値 x_i および重み w_i を受け取り、論理演算を行う論理回路部と、多ビットバイアス W' を受け取り、前記論理回路部の出力と前記多ビットバイアス W' との総和を取る総和回路部と、総和を取った多ビット信号 Y に対して符号ビットのみを出力する活性化回路部と、を備えることを特徴とする。

30

【発明の効果】

【0010】

本発明によれば、バッチ正規化回路が不要なニューラルネットワーク回路装置、ニューラルネットワーク、ニューラルネットワーク処理方法およびニューラルネットワークの実行プログラムを提供することができる。

【図面の簡単な説明】

【0011】

【図 1】ディープニューラルネットワーク (DNN) の構造の一例を説明する図である。

【図 2】比較例のニューラルネットワークのニューラルネットワーク回路の構成の一例を示す図である。

40

【図 3】図 2 に示すニューラルネットワーク回路における活性化関数 $f_{act}(Y)$ を示す図である。

【図 4】図 2 に示すニューラルネットワーク回路の乗算回路を XNOR ゲート回路に置き換えた 2 値化ニューラルネットワーク回路の構成の一例を示す図である。

【図 5】図 4 に示す 2 値化ニューラルネットワーク回路における活性化関数 $f_{sgn}(B)$ を示す図である。

【図 6】比較例のバッチ正規化回路を備える 2 値化ニューラルネットワーク回路の構成の一例を示す図である。

【図 7】ニューラルネットワークの 2 値化ニューラルネットワーク回路のスケーリング (

50

)による正規化を示す図である。

【図8】ニューラルネットワークの2値化ニューラルネットワーク回路のシフト()による-1~+1の制限を示す図である。

【図9】本発明の実施形態に係るディープニューラルネットワークの2値化ニューラルネットワーク回路の構成を示す図である。

【図10】本発明の実施形態に係るディープニューラルネットワークの2値化ニューラルネットワーク回路の活性化回路を示す図である。

【図11】本発明の実施形態に係るディープニューラルネットワークの多ビット構成のニューラルネットワーク回路と2値化ニューラルネットワーク回路の認識精度を説明する図である。

10

【図12】本発明の実施形態に係るディープニューラルネットワークの2値化ニューラルネットワーク回路と既存の多ビット実装法との比較を行った結果を表にして示す図である。

【図13】本発明の実施形態に係るディープニューラルネットワークの2値化ニューラルネットワーク回路の実装例を説明する図である。

【図14】変形例のディープニューラルネットワークの2値化ニューラルネットワーク回路の構成を示す図である。

【図15】変形例の2値化ニューラルネットワーク回路のLUTの構成を示す図である。

【発明を実施するための形態】

【0012】

20

以下、図面を参照して本発明を実施するための形態(以下、「本実施形態」という)におけるディープニューラルネットワークについて説明する。

(背景説明)

図1は、ディープニューラルネットワーク(DNN)の構造の一例を説明する図である。

図1に示すように、ディープニューラルネットワーク(DNN)1は、入力層(input layer)11、任意の数の中間層である隠れ層(hidden layer)12、出力層(output layer)13から構成される。

入力層(input layer)11は、複数個(ここでは8)の入力ノード(ニューロン)を有する。隠れ層12は、複数(ここでは3層(hidden layer1, hidden layer2, hidden layer3))である。実際には、隠れ層12の層数nは、例えば20~100に達する。出力層13は、識別対象の数(ここでは4)の出力ノード(ニューロン)を有する。なお、層数およびノード数(ニューロン数)は、一例である。

30

ディープニューラルネットワーク1は、入力層11と隠れ層12のノード間が全て結合し、隠れ層12と出力層13のノード間が全て結合している。

【0013】

入力層11、隠れ層12および出力層13には、任意の数のノード(図1の印参照)が存在する。このノードは、入力を受け取り、値を出力する関数である。入力層11には、入力ノードとは別に独立した値を入れるバイアス(bias)ノードがある。構成は、複数のノードを持つ層を重ねることで構築される。伝播は、受け取った入力に対して重み(weight)をかけ、受け取った入力を次層に活性化関数(activation function)で変換して出力する。活性化関数は、sigmoid関数やtanh関数などの非線形関数、ReLU(Rectified Linear Unit function: 正規化線形関数)がある。ノード数を増やすことで、扱う変数を増やし、多数の要素を加味して値/境界を決定できる。層数を増やすことで、直線境界の組み合わせ、複雑な境界を表現できる。学習は、誤差を計算し、それを基に各層の重みを調整する。学習は、誤差を最小化する最適化問題を解くことであり、最適化問題の解法は誤差逆伝播法(Backpropagation)を使うのが一般的である。誤差は、二乗和誤差を使うのが一般的である。汎化能力を高めるために、誤差に正則化項を加算する。誤差逆伝播法は、誤差を出力層13から伝播させていき、各層の重みを調整する。

40

【0014】

50

図1のディープニューラルネットワーク1の構成を2次元に展開することで画像処理に適したCNNを構築できる。また、ディープニューラルネットワーク1にフィードバックを入れることで、双方向に信号が伝播するRNN(Recurrent Neural Network:再帰型ニューラルネットワーク)を構成することができる。

【0015】

図1の太破線三角部に示すように、ディープニューラルネットワーク1は、多層のニューラルネットワークを実現する回路(以下、ニューラルネットワーク回路という)2から構成されている。

本技術は、ニューラルネットワーク回路2を対象とする。ニューラルネットワーク回路2の適用箇所および適用数は限定されない。例えば、隠れ層12の層数 $n: 20 \sim 30$ の場合、これらの層のどの位置に適用してもよく、またどのノードを入出力ノードとするものでもよい。さらに、ディープニューラルネットワーク1に限らず、どのようなニューラルネットワークでもよい。ただし、入力層11または出力層13のノード出力には、2値化出力ではなく多ビット出力が求められるので、ニューラルネットワーク回路2は、対象外である。ただし、出力層13のノードを構成する回路に、乗算回路が残ったとしても面積的には問題にはならない。

なお、入力データに対し学習済のものを評価していくことを前提としている。したがって、学習結果として重み w_i は既に得られている。

【0016】

<ニューラルネットワーク回路>

図2は、比較例のニューラルネットワーク回路の構成の一例を示す図である。

比較例のニューラルネットワーク回路20は、図1のディープニューラルネットワーク1を構成するニューラルネットワーク回路2に適用できる。なお、以下の各図の表記において、値が多ビットである場合は太実線矢印、値が2値である場合は細太実線矢印で示す。

ニューラルネットワーク回路20は、入力値(判別データ) $X_1 \sim X_n$ (多ビット)を入力する入力ノード、重み $W_1 \sim W_n$ (多ビット)およびバイアス W_0 (多ビット)を入力する入力部21と、入力値 $X_1 \sim X_n$ および重み $W_1 \sim W_n$ を受け取り、入力値 $X_1 \sim X_n$ に重み $W_1 \sim W_n$ をそれぞれ乗算する複数の乗算回路22と、各乗算値とバイアス W_0 との総和を取る総和回路23と、総和を取った信号 Y を活性化関数 $f_{act}(Y)$ で変換する活性化関数回路24と、を備えて構成される。

以上の構成において、ニューラルネットワーク回路20は、入力値 $X_1 \sim X_n$ (多ビット)を受け取り、重み $W_1 \sim W_n$ を乗算した後に、バイアス W_0 を含めて総和を取った信号 Y を活性化関数回路24を通すことで人間のニューロンに模した処理を実現している。

【0017】

図3は、前記図2に示すニューラルネットワーク回路20における活性化関数 $f_{act}(Y)$ を示す図である。図3は、横軸に総和を取った信号 Y 、縦軸に活性化関数 $f_{act}(Y)$ の値をとる。図3の符号印は、 ± 1 の範囲の値をとる正側の活性化値(状態値)、図3の符号 \times 印は、 ± 1 の範囲の値をとる負側の活性化値である。

ニューラルネットワーク回路20(図2参照)は、多ビットで高い認識精度を実現している。このため、活性化関数回路24(図2参照)において、非線形な活性化関数 $f_{act}(Y)$ を用いることができる。すなわち、図4に示すように、非線形な活性化関数 $f_{act}(Y)$ は、傾きが非ゼロとなる部分(図4の破線囲み部分参照)に ± 1 の範囲の値をとる活性化値を設定できる。このため、ニューラルネットワーク回路20は、多様な活性を実現でき、認識精度は実用的な値になっていた。しかし、ニューラルネットワーク回路20は、大量の乗算回路22が必要になる。加えて、ニューラルネットワーク回路20は、入出力・重みが多ビットであることにより、大量のメモリが必要であり、読み書きの速度(メモリ容量・帯域)も問題である。

【0018】

<単に2値化した2値化ニューラルネットワーク回路>

図 2 に示す比較例のニューラルネットワーク回路 20 は、短精度（多ビット）による積和演算回路で構成されている。このため、多数の乗算回路 21 が必要であり、面積・消費電力が多くなる欠点があった。また、入出力・重みが多ビットであることで大量のメモリが必要であり、読み書きの速度（メモリ容量・帯域）が問題となっていた。

そこで、2 値化した精度、すなわち + 1 と - 1 のみ用いてニューラルネットワーク回路 2（図 1 参照）を構成する回路が提案された（非特許文献 1 ~ 4）。具体的には、図 2 に示すニューラルネットワーク回路 20 の乗算回路 21 を、論理ゲート（例えば X N O R ゲート回路）に置き換えることが考えられる。

【 0 0 1 9 】

図 4 は、図 2 に示すニューラルネットワーク回路 20 の乗算回路 21 を X N O R ゲート回路に置き換えた 2 値化ニューラルネットワーク回路の構成の一例を示す図である。

比較例の 2 値化ニューラルネットワーク回路 30 は、図 1 のニューラルネットワーク回路 2 に適用できる。

図 4 に示すように、比較例の 2 値化ニューラルネットワーク回路 30 は、入力値 $x_1 \sim x_n$ （2 値）を入力する入力ノード、重み $w_1 \sim w_n$ （2 値）およびバイアス w_0 （2 値）を入力する入力部 31 と、入力値 $x_1 \sim x_n$ および重み $w_1 \sim w_n$ を受け取り、X N O R（Exclusive NOR：否定排他的論理和）論理を取る複数の X N O R ゲート回路 32 と、X N O R ゲート回路 32 の各 X N O R 論理値とバイアス w_0 との総和を取る総和回路 33 と、総和を取った信号 Y のバッチ正規化した信号 B を活性化関数 $f \operatorname{sgn}(B)$ で変換する活性化関数回路 34 と、を備えて構成される。

2 値化ニューラルネットワーク回路 30 は、乗算回路 21（図 2 参照）が X N O R 論理を実現する X N O R ゲート回路 32 に置き換えられている。このため、乗算回路 21 を構成する際に必要であった面積を削減することができる。また、入力値 $x_1 \sim x_n$ 、出力値 z 、および重み $w_1 \sim w_n$ は、いずれも 2 値（- 1 と + 1）であるため、多値である場合と比較してメモリ量を大幅に削減でき、メモリ帯域を向上させることができる。

【 0 0 2 0 】

図 5 は、前記図 4 に示す 2 値化ニューラルネットワーク回路 30 における活性化関数 $f \operatorname{sgn}(B)$ を示す図である。図 5 は、横軸に総和を取った信号 Y、縦軸に活性化関数 $f \operatorname{sgn}(B)$ の値をとる。図 5 の符号 \oplus は、 ± 1 の範囲の値をとる正側の活性化値、図 5 の符号 \otimes は、 ± 1 の範囲の値をとる負側の活性化値である。

2 値化ニューラルネットワーク回路 30 は、入力値 $x_1 \sim x_n$ および重み $w_1 \sim w_n$ を単に 2 値化している。このため、図 5 の符号 a に示すように、 ± 1 のみ扱う活性化関数しか扱えないため、誤差が頻繁に生じてしまう。また、傾きが非ゼロとなる区間（図 5 の破線囲み部分参照）が不均等となり学習が上手く行われぬ。すなわち、図 6 の符号 b に示すように、不均等な幅により微分が定義できない。その結果、単に 2 値化した 2 値化ニューラルネットワーク回路 40 は、認識精度が大幅に落ち込んでしまう。

そこで、非特許文献 1 ~ 4 には、既存の 2 値化ニューラルネットワークの精度を維持するためにバッチ正規化を行う技術が記載されている。

【 0 0 2 1 】

<バッチ正規化回路を備える 2 値化ニューラルネットワーク回路>

図 6 は、2 値化した精度を是正して、C N N の認識精度を保つバッチ正規化回路を備える 2 値化ニューラルネットワーク回路 40 の構成の一例を示す図である。図 4 と同一構成部分には同一符号を付している。

図 6 に示すように、比較例の 2 値化ニューラルネットワーク回路 40 は、入力値 $x_1 \sim x_n$ （2 値）を入力する入力ノード $x_1 \sim x_n$ 、重み $w_1 \sim w_n$ （2 値）およびバイアス w_0 （2 値）を入力する入力部 31 と、入力値 $x_1 \sim x_n$ および重み $w_1 \sim w_n$ を受け取り、X N O R（Exclusive NOR：否定排他的論理和）論理を取る複数の X N O R ゲート回路 32 と、X N O R ゲート回路 32 の各 X N O R 論理値とバイアス w_0 （2 値）との総和を取る総和回路 33 と、2 値化によるバラツキの偏りを正規化範囲を広げ中心をシフトさせる処理で是正するバッチ正規化回路 41 と、総和を取った信号 Y のバッチ正規化した信

10

20

30

40

50

号 B を活性化関数 $f \operatorname{sgn}(B)$ で変換する活性化関数回路 3 4 と、を備えて構成される。

【 0 0 2 2 】

バッチ正規化回路 4 1 は、重み総和後、スケーリング () 値 (多ビット) による正規化を行う乗算回路 4 2 と、スケーリング () 値による正規化後、シフト () 値 (多ビット) によりシフトして 2 分類を行う加算器 4 3 と、からなる。スケーリング () 値およびシフト () 値の各パラメータは、事前に学習時に求めておく。

2 値化ニューラルネットワーク回路 4 0 は、バッチ正規化回路 4 1 を備えることで、2 値化した精度を是正して、CNN の認識精度を保つようにする。

なお、入力値 $x_1 \sim x_n$ と重み $w_1 \sim w_n$ との XNOR 論理を取る論理回路であれば、XNOR ゲートに限らずどのような論理ゲートでもよい。

10

【 0 0 2 3 】

しかしながら、図 6 に示すように、バッチ正規化回路 4 0 は、乗算回路 4 2 と加算器 4 3 が必要である。また、スケーリング () 値およびシフト () 値をメモリに格納しておく必要がある。メモリは、外付けであり大面積となり、読み出しの速度遅延となる。

2 値化ニューラルネットワーク回路 4 0 は、図 2 に示すニューラルネットワーク回路 2 0 のように、多数の乗算回路 2 1 を必要としないものの、バッチ正規化回路 4 1 において、面積が大きい乗算回路 4 2 と加算器 4 3 が必要となる。バッチ正規化回路 4 1 では、パラメータ格納用のメモリも必要となり、面積とメモリ帯域の削減が求められている。

【 0 0 2 4 】

<バッチ正規化回路が必要となる理由>

20

比較例の 2 値化ニューラルネットワーク回路 4 0 のバッチ正規化回路 4 1 が必要となる理由について説明する。

図 7 および図 8 は、比較例の 2 値化ニューラルネットワーク回路 4 0 のバッチ正規化による効果を説明する図である。図 7 は、スケーリング () による正規化を示す図、図 8 は、シフト () による $-1 \sim +1$ の制限を示す図である。

バッチ正規化とは、2 値化によるバラツキの偏りを是正する回路であり、重み総和後、スケーリング () 値による正規化を行った後、シフト () 値による適切な活性化による 2 分類を行う。これらのパラメータは事前に学習時に求めておく。具体的には、下記の通りである。

【 0 0 2 5 】

30

図 7 の白抜矢印および符号 c に示すように、バッチ正規化回路 4 1 の乗算回路 4 2 (図 6 参照) は、重み総和後の信号 (結果) Y を、スケーリング () 値により、幅「2」(図 7 の網掛け部参照) に正規化する。これにより、図 5 の幅 (図 5 の網掛け部参照) と比較して分かるように、単に 2 値化した 2 値化ニューラルネットワーク回路 3 0 では、不均等な幅により微分が定義できなかった不具合が、スケーリング () 値により幅「2」に正規化することで、不均等な幅が抑制される。

【 0 0 2 6 】

その上で、図 8 の白抜矢印および符号 d に示すように、バッチ正規化回路 4 1 の加算器 4 3 (図 6 参照) は、スケーリング () 値による正規化後の値を、シフト () 値により $-1 \sim +1$ の範囲になるよう制限する。すなわち、図 5 の幅 (図 5 の網掛け部参照) と比較して分かるように、図 5 の幅 (図 5 の網掛け部参照) が、 $+1$ 側により多くシフトしている場合には、シフト () 値により、スケーリング () 値による正規化後の値を $-1 \sim +1$ に制限することで、この幅の中心を 0 とする。図 5 の例では、負側の活性化値 (図 5 の破線囲み部の符号 x 印参照) が、本来あるべき負側に戻される。これにより、誤差の発生が減少し、認識精度を高めることができる。

40

このように、2 値化ニューラルネットワーク回路 4 0 には、バッチ正規化回路 4 1 が必要である。

【 0 0 2 7 】

<バッチ正規化回路を備える 2 値化ニューラルネットワーク回路の課題>

上記バッチ正規化回路 4 1 を導入することで、2 値化ニューラルネットワーク回路 4 0

50

の認識精度は、図2に示したニューラルネットワーク回路20とほぼ等しくなる。しかしながら、バッチ正規化回路41は、乗算回路42と加算器43が必要であり、多ビットのスケーリング()値およびシフト()値をメモリに格納しておく必要がある。このため、2値化ニューラルネットワーク回路40は、依然として複雑な回路であり、面積・消費電力を低減したいという切実な要請がある。

【0028】

2値化ニューラルネットワーク回路20では、例えば8,9ビットを1ビットに落とすので、計算精度が落ちる。NNに用いた場合、誤認識(認識失敗率)が80%に増加し、使用に耐えない。そこで、バッチ正規化で対処することになる。しかし、バッチ正規化回路41は、除算、若しくは浮動小数点の乗算と加算が必要であり、ハードウェア(HW)化して実装することが非常に困難である。また、外部メモリが必要となり、外部メモリとのアクセスのために遅延となる。

【0029】

(本発明の原理説明)

本発明者らは、バッチ正規化の操作を導入したNNに対して、これと等価なNNを解析的に求めると、バッチ正規化が不要なNNを得ることができることを発見したことが着眼点である。例えば、従来では、図3のような非線形な活性化関数 $f_{act}(Y)$ に対し、図3の符号印に示す状態値が得られた場合、不均等な幅を正規化するためスケーリングを行っていた。多ビットを2値して乗算する場合の演算精度を確保するためである。しかしながら、本発明者らは、ニューラルネットワーク回路における2値化の本質は、活性化するかしないかの(2値)だけに帰着することに着目した。スケーリングは、不要となりシフトのみで対応できることになる。

すなわち、重み積和後に2値化ニューラルネットワーク回路40のバッチ正規化回路41(図6参照)に入力される信号を Y とすると、バッチ正規化回路41から出力される信号(Y と等価となる信号) Y' は、次式(1)で示される。

【0030】

【数1】

$$Y' = \gamma \frac{Y - \mu_B}{\sqrt{\sigma_B^2 + \epsilon}} + \beta$$

$$= \frac{\gamma}{\sqrt{\sigma_B^2 + \epsilon}} \left(Y - \left(\mu_B - \frac{\sqrt{\sigma_B^2 + \epsilon}}{\gamma} \beta \right) \right) \quad \dots(1)$$

ただし、

: スケーリング値

: シフト値

μ_B : 平均値

σ_B^2 : 二乗和誤差

: パラメータ(調整用)

したがって、2値化活性化関数の値 $f'_{sgn}(Y)$ は、下記式(2)の条件で決まる。

【0031】

【数 2】

$$f'_{sgn}(Y) = \begin{cases} 1 & \left(\text{if } Y < -\mu_B + \frac{\sqrt{\sigma_B^2 + \epsilon}}{\gamma} \beta \right) \\ -1 & (\text{otherwise}) \end{cases} \quad \dots (2)$$

【0032】

よって、これらの解析的な操作から重み積和演算は、下記式(3)のように得られる。

10

【0033】

【数 3】

$$\begin{aligned} Y &= \sum_{i=0}^n w_i x_i - \mu_B + \frac{\sqrt{\sigma_B^2 + \epsilon}}{\gamma} \beta \\ &= \sum_{i=1}^n w_i x_i + \left(w_0 - \mu_B + \frac{\sqrt{\sigma_B^2 + \epsilon}}{\gamma} \beta \right) \\ &= \sum_{i=1}^n w_i x_i + W' \end{aligned} \quad \dots (3)$$

20

ただし、

W' : 多ビットバイアス

【0034】

バッチ正規化学習後に、これらの数学的な操作により、バッチ正規化と等価なNNの演算が得られる。

30

上記式(3)は、回路的にはバイアス値のみ多ビット構成にすればよいことを示している。回路は単純だが、単にバイアス値のみ多ビットにただけでは認識精度は向上せず、これらの解析的な洞察がなければ成り立たない。

【0035】

[実施形態の構成]

図9は、本発明の実施形態に係るニューラルネットワークの2値化ニューラルネットワーク回路の構成を示す図である。本実施形態の2値化ニューラルネットワーク回路は、ディープニューラルネットワークへの実装技術を提供する。

2値化ニューラルネットワーク回路100は、図1のニューラルネットワーク回路2に適用できる。

40

図9に示すように、2値化ニューラルネットワーク回路100(ニューラルネットワーク回路装置)は、入力値 $x_1 \sim x_n$ (x_i)(2値)を入力する入力ノードおよび重み $w_1 \sim w_n$ (w_i)(2値)を入力する入力部101と、入力値 $x_1 \sim x_n$ および重み $w_1 \sim w_n$ を受け取り、XNOR論理を取るXNORゲート回路102(論理回路部)と、多ビットバイアス W' (式(3)参照)を入力する多ビットバイアス W' 入力部110と、各XNOR論理値と多ビットバイアス W' との総和を取る総和回路103(総和回路部)と、総和を取った信号 Y に対して符号ビットのみを出力する活性化回路120(活性化回路部)と、を備えて構成される。

【0036】

上記入力値 x_i (2値)および重み w_i (2値)は、2値信号である。

50

上記多ビット信号 Y および多ビットバイアス W' は、前記式(3)で示される。

【0037】

2値化ニューラルネットワーク回路100は、ディープニューラルネットワーク1の隠れ層12(図1参照)に適用される。ここでは、ディープニューラルネットワーク1において、入力値に対し学習済のものを評価していくことを前提としている。したがって、学習結果として重みの多ビットバイアス W' は既に得られている。多ビットバイアス W' は、学習後の多ビットバイアス値である。ちなみに、図2のニューラルネットワーク回路20では、多ビットの重み $W_1 \sim W_n$ およびバイアス W_0 を用いるが、本実施形態の多ビットバイアス W' は、図2のニューラルネットワーク回路20における多ビットのバイアス W_0 とは異なるものである。

10

なお、NNでは、重みが、クライアントの認識物体毎に全て異なる。また学習により毎回異なることがある。画像処理では係数は、全て同じであり、この点でNNと画像処理では、HWが大きく異なる。

【0038】

XNORゲート回路102は、排他的論理和を含むどのような論理回路部でもよい。すなわち、入力値 $x_1 \sim x_n$ と重み $w_1 \sim w_n$ との論理を取る論理回路であれば、XNORゲートに限らずどのようなゲート回路でもよい。例えば、XORゲートにNOTゲートを組み合わせる、AND、ORゲートを組み合わせる、さらにはトランジスタスイッチを用いて作製するなど、論理的に等しいものであればどのようなものでもよい。

【0039】

活性化回路120は、総和を取った信号 Y に対して符号ビットのみを出力する活性化関数回路を模擬する回路となっている。符号ビットは、総和を取った多ビット信号 Y を活性化するかしないかで示す2値信号である。

20

【0040】

このように、2値化ニューラルネットワーク回路100は、バイアス値のみ多ビット構成にし、バイアス値を含めた総和から、符号ビットのみを出力する活性化回路120を備える。すなわち、2値化ニューラルネットワーク回路100は、図6に示す2値化ニューラルネットワーク回路40のバッチ正規化回路41および活性化関数回路34を、符号ビットのみを出力する活性化回路120で置き換えた構成となっている。このため、2値化ニューラルネットワーク回路100は、複雑なバッチ正規化回路41が不要なニューラルネットワーク回路となっている。

30

【0041】

図10は、2値化ニューラルネットワーク回路の活性化回路を示す図である。

図10に示すように、活性化回路120は、バイアス値を含めた総和の出力 Y から、符号ビットのみを出力する回路である。図10の回路では、符号ビットは、出力 $y[0]$ 、 $y[1]$ 、 \dots 、 $y[n-1]$ のうち、最上位ビット $y[n-1]$ であるとすると、符号ビットとして最上位ビット $y[n-1]$ のみを出力する。活性化回路120は、最上位ビット $y[n-1]$ のみが出力 z として出力される。図9では、活性化回路120は、活性化関数 $f \operatorname{sgn}(Y)$ と表記されているが、図6に示すスケールリング()による正規化およびシフト()による $-1 \sim +1$ の制限は行っておらず、最上位ビット $y[n-1]$ のみを出力する回路である。

40

【0042】

以下、上述のように構成された2値化ニューラルネットワーク回路100の動作について説明する。

2値化ニューラルネットワーク回路100は、図1に示すディープニューラルネットワーク1のニューラルネットワーク回路2に用いられる。この場合、2値化ニューラルネットワーク回路100の入力ノード $x_1 \sim x_n$ は、図1に示すディープニューラルネットワーク1のhidden layer1の入力ノードである。入力部101には、隠れ層12のhidden layer1の入力ノードの入力値 $x_1 \sim x_n$ (2値)および重み $w_1 \sim w_n$ (2値)が入力される。

XNORゲート回路102では、入力値 $x_1 \sim x_n$ および重み $w_1 \sim w_n$ を受け取り、

50

XNOR論理により2値(-1/+1)の乗算を行う。

2値化ニューラルネットワーク回路100は、多ビット構成の乗算回路21(図2参照)がXNOR論理を実現するXNORゲート回路102に置き換えられている。このため、乗算回路21を構成する際に必要であった面積を削減することができる。また、入力値 $x_1 \sim x_n$ および重み $w_1 \sim w_n$ は、いずれも2値(-1/+1)であるため、多ビット(多値)である場合と比較してメモリ容量を大幅に削減でき、メモリ帯域を向上させることができる。

【0043】

一方、前記式(3)に従った多ビットバイアス W' を入力する。多ビットバイアス W' は、2値化ニューラルネットワーク回路30, 40(図4および図6参照)のような2値のバイアス w_0 ではない。また、多ビットではあっても2値化ニューラルネットワーク回路20(図2参照)のようなバイアス W_0 とは異なる。多ビットバイアス W' は、前記式(3)に示すように、前記バイアス w_0 (2値)からバッチ正規化分を調整した、学習後のバイアス値である。

総和回路103には、バイアス値のみ多ビット構成にした多ビットバイアス W' が入力される。総和回路103は、XNORゲート回路102の各XNOR論理値と多ビットバイアス W' との総和を取り、総和の出力 Y (多ビット)を活性化回路120に出力する。

【0044】

図10に示すように、活性化回路120では、バイアス値を含めた総和の出力 Y (多ビット)から、符合ビットのみを出力する。図10の回路では、符号ビットは、出力 $y[0], y[1], \dots, y[n-1]$ のうち、最上位ビット $y[n-1]$ である。活性化回路120は、バイアス値を含めた総和の出力 Y から、最上位ビット $y[n-1]$ のみを出力 z として出力する。換言すれば、活性化回路120は、 $y[0], y[1], \dots, y[n-2]$ の数値を出力しない($y[0], y[1], \dots, y[n-2]$ の数値は使用しない)。

例えば、活性化回路120の入力 Y として4~5bitの信号が入力された場合、HWでは、通常最上位ビットを符号ビットとするので、最上位ビット(符合ビット)だけを出力する。すなわち、活性化回路120からは、活性化するかしないかの二通り(2値、すなわち+1か-1)が出力され、それが後段の中間層(隠れ層)のノードに伝達される。

【0045】

2値化ニューラルネットワーク回路100は、前記式(3)に示されるように、バッチ正規化の操作を導入したNNと等価なNNである。前記式(3)は、下記により実現される。すなわち、2値(1ビットのみ)にした入力値 x_i と重み w_i 、多ビットバイアス W' を入力とし、乗算の代わりとなるXNOR論理を取った後、バイアス値を含めたそれらの総和を取り(前記式(3)の第1項)、活性化回路120がバイアス値を含めた総和の出力 Y から、符合ビットのみを出力する(前記式(3)の第2項)。

このため、活性化回路120は、バイアス値を含めた総和の出力 Y から、符号ビットのみを出力する回路ではあるが、機能的には、活性化関数回路 $f \operatorname{sgn}(Y)$ と同様な機能、すなわち活性化関数回路 $f \operatorname{sgn}(Y)$ を模擬した回路となっている。

【0046】

本実施形態の効果を確認するため、VGG16(隠れ層が16層)ベンチマークNNを実装した。VGG16は、良く使われているベンチマークで再現性があるものである。

図11は、多ビット構成のニューラルネットワーク回路と2値化ニューラルネットワーク回路の認識精度を説明する図である。図11(a)は多ビット(32ビット浮動小数点)で構成したニューラルネットワーク回路20(図2参照)の認識精度、図11(b)は2値化ニューラルネットワーク回路100の認識精度を示す。図11の横軸は、利用した学習データに対して更新を終えたサイクルであるエポック(epoch)数、縦軸は誤認識(誤差)(Classification error)である。図11は、本実施形態をVGG16ベンチマークNNで実装し確認したものである。また、図11(a)は、ディープニューラルネットワーク用のフレームワークソフトウェアChainer(登録商標)のfloat32 CNNを用いている。また、図11(b)は、ディープニューラルネットワーク用のフレームワークソフトウェア

10

20

30

40

50

Chainer（登録商標）のfloat32 CNNを用いている。また、バッチ正規化なし、バッチ正規化ありを示している。

【0047】

図11(a)に示すように、多ビット構成のニューラルネットワーク回路20では、誤差(Classification error)が低く認識精度は高い。この多ビット構成のニューラルネットワーク回路20の認識精度を比較対象として、2値化ニューラルネットワーク回路の認識精度を検討する。

図11(b)の「バッチ正規化なし」に示すように、単に2値化した2値化ニューラルネットワーク回路30(図4参照)では、誤差率(Classification error)が大きく(約80%)認識精度は悪い。また、学習を続けても誤差率の改善は見られない(学習が収束しない)。

これに対して、図11(b)の「バッチ正規化あり」で示される本実施形態の2値化ニューラルネットワーク回路100は、多ビット構成のニューラルネットワーク回路20と比較して約6%の誤差(VGG-16を使用)に収まることが確認された。ただし、ニューロン数は同じ場合であるのでニューロン数を増やすとその差は縮まる。また、本実施形態の2値化ニューラルネットワーク回路100は、多ビット構成のニューラルネットワーク回路20と同様に、学習を続けるに従って収束していくことが確認された。

【0048】

本実施形態では、2値化ニューラルネットワーク回路40(図6参照)で必須であったバッチ正規化回路41(図6参照)自体を不要であり、それらのパラメータも不要であることから面積が削減でき、メモリ量も削減できる。また、図11(a)の「バッチ正規化あり」と図11(b)の「バッチ正規化あり」とを比較してわかるように、本実施形態の2値化ニューラルネットワーク回路100は、認識精度について、多ビット構成のニューラルネットワーク回路20(図2参照)と数%異なるだけであった。

【0049】

図12は、本実施形態の2値化ニューラルネットワーク回路100をFPGA(Digilent社NetFPGA-1G-CML)上に実装し、既存の多ビット実装法との比較を行った結果を表にして示す図である。

図12の表は、表下欄外に表記した[1]~[4]の学会発表者(論文発表年)のニューラルネットワークと本実施形態のニューラルネットワークをFPGA上に実現した場合に、各項目を対比して示したものである。「Platform」(プラットフォーム)、「Clock(MHz)」(同期化のための内部クロック)、「Bandwidth(GB/s)」(データ転送のバンド幅/外部にメモリを付けた場合の転送速度)、「Quantization Strategy」(量子化ビット数)、「Power(W)」(消費電力)、「Performance(GOP/s)」(チップ面積に対する性能)、「Resource Efficiency(GOP/s/Slices)」(リソース効率)、および「Power Efficiency(GOP/s/W)」(性能パワー効率)の各項目を対比して示した。この表において、特に注目すべき事項は下記の通りである。

【0050】

<消費電力>

本実施形態の2値化ニューラルネットワーク回路100は、表の従来例と比較して、電力のバランスが取れていることが挙げられる。従来例では、「Power(W)」に示すように、消費電力が大きい。消費電力が大きいので、これを回避する制御方法が複雑である。「Power(W)」に示すように、本実施形態では、従来例と比較して消費電力を1/2~1/3に低減することができた。

【0051】

<チップ面積>

本実施形態の2値化ニューラルネットワーク回路100は、バッチ正規化回路がなくメモリが不要であること、乗算回路が2値論理ゲートであること、活性化関数が単純であること(活性化関数回路ではなく活性化関数回路を模擬する活性化回路120であること)、から、表の「Performance(GOP/s)」に示すように、チップ面積に対する性能は、従来

10

20

30

40

50

例と比較して約30倍となる。すなわち、チップ面積が減る、外付けメモリが不要となる、メモリコントローラおよび活性化関数が単純になることなどの効果がある。チップ面積は価格に比例するので、価格も2桁程度安くなることが期待できる。

【0052】

<性能等価>

本実施形態の2値化ニューラルネットワーク回路100は、表の「Bandwidth (GB/s)」に示すように、従来例と比較してほぼ同等である。また、表の「Power (W)」に示すように、性能パワー効率は、面積を見ずにパワー効率だけを見たものでも約2倍となっている。さらに、表の「Power Efficiency (GOP/s/W)」に示すように、単位ワット数当たりの処理能力(基板全体のワット数)も約2倍となっている。

10

【0053】

[実装例]

図13は、本発明の実施形態に係る2値化ニューラルネットワーク回路の実装例を説明する図である。

<STEP1>

まず、与えられたデータセット(今回はImageNet、画像認識タスク用にデータ)を既存のディープニューラルネットワーク用のフレームワークソフトウェアであるChainer(登録商標)を用いてCPU(Central Processing Unit)101を有するコンピュータ上で学習を行った。このコンピュータは、ARMプロセッサなどのCPU101と、メモリと、ハードディスクなどの記憶手段(記憶部)と、ネットワークインタフェースを含むI/Oポートとを有する。このコンピュータは、CPU101が、メモリ上に読み込んだプログラム(2値化したニューラルネットワークの実行プログラム)を実行することにより、後記する各処理部により構成される制御部(制御手段)を動作させる。

20

【0054】

<STEP2>

次に、自動生成ツールを用いて、本実施形態の2値化ニューラルネットワーク回路100と等価なC++コードを自動生成し、C++コード102を得た。

<STEP3>

次に、FPGAベンダの高位合成ツール(Xilinx社SDSoC)(登録商標)を用いて、FPGA(field-programmable gate array)合成用にHDL(hardware description language)を生成した。

30

<STEP4>

次に、従来のFPGA合成ツールVivado(登録商標)を用いて、FPGA上に実現して画像認識タスクの検証を行った。

<STEP5>

検証後、基板103を完成させた。基板103には、2値化ニューラルネットワーク回路100がハードウェア化されて実装されている。

【0055】

以上説明したように、本実施形態に係る2値化ニューラルネットワーク回路100(図9参照)は、入力値 $x_1 \sim x_n$ (x_i)(2値)を入力する入力ノードおよび重み $w_1 \sim w_n$ (w_i)(2値)を入力する入力部101と、入力値 $x_1 \sim x_n$ および重み $w_1 \sim w_n$ を受け取り、XNOR論理を取るXNORゲート回路102と、多ビットバイアス W' (式(3)参照)を入力する多ビットバイアス W' 入力部110と、各XNOR論理値と多ビットバイアス W' との総和を取る総和回路103と、総和を取った信号Yに対して符号ビットのみを出力する活性化回路120と、を備える。

40

【0056】

この構成により、バッチ正規化回路自体が不要であり、それらのパラメータも不要であることから面積が削減でき、メモリ量も削減できる。また、本実施形態では、バッチ正規化回路がないにも拘わらず、バッチ正規化回路41を備える2値化ニューラルネットワーク回路40(図6参照)と性能的には等価な回路構成となっている。このように、バッチ

50

正規化回路の面積とパラメータを格納するメモリ面積・メモリ帯域を無くすことができ、かつ、性能的には等価な回路構成を実現することができる。例えば、図12の表に示すように、本実施形態に係る2値化ニューラルネットワーク回路100は、消費電力を半分に削減でき、面積を約30分の1に削減できた。

【0057】

本実施形態によれば、既存のバッチ正規化回路を備える2値化ニューラルネットワーク回路と比較して、面積を約30分の1に削減しつつ、認識精度はほぼ等価なCNNを構成できることが判明した。ディープラーニングを用いたADAS (Advanced Driver Assistance System: 先進運転支援システム) カメラ画像認識用のエッジ組み込み装置ハードウェア方式として実用化が期待される。特にADASでは、車載する上で高信頼性と低発熱が要求される。本実施形態に係る2値化ニューラルネットワーク回路100は、図12の表に示すように、消費電力が格段に低減していることに加え、外付けメモリが不要であるので、メモリを冷却する冷却ファンや冷却フィンも不要である。ADASカメラに搭載して好適である。

10

【0058】

[変形例]

図14は、変形例のディープニューラルネットワークの2値化ニューラルネットワーク回路の構成を示す図である。図9と同一構成部分には同一符号を付している。

本変形例は、乗算回路としての論理ゲートに代えて、LUT (Look-Up Table) を用いた例である。

20

2値化ニューラルネットワーク回路200は、図1のニューラルネットワーク回路2に適用できる。

図14に示すように、2値化ニューラルネットワーク回路200 (ニューラルネットワーク回路装置) は、入力値 $x_1 \sim x_n$ (x_i) (2値) を入力する入力ノード $x_1 \sim x_n$ および重み $w_1 \sim w_n$ (2値) を入力する入力部101と、入力値 $x_1 \sim x_n$ および重み $w_1 \sim w_n$ を受け取り、2値 (-1 / +1) の乗算を行うためのテーブル値を格納し演算時に参照されるLUT202 (論理回路部) と、多ビットバイアス W' (式(3)参照) を入力給する多ビットバイアス W' 入力手段110と、LUT202から参照した各テーブル値と多ビットバイアス W' との総和を取る総和回路103と、総和を取った信号 Y に対して符号ビットのみを出力する活性化関数回路を模擬する活性化回路120と、を備えて構成される。

30

【0059】

本変形例では、乗算回路としての論理ゲートに代えて、LUT (Look-Up Table) を用いた例である。

LUT202は、XNOR論理を行うXNORゲート回路102 (図9参照) に代えて、FPGAの基本構成要素であるルックアップテーブルを用いる。

【0060】

図15は、変形例の2値化ニューラルネットワーク回路200のLUT202の構成を示す図である。

図15に示すように、LUT202は、2入力 (x_1, w_1) に対する2値 (-1 / +1) のXNOR論理結果 Y を格納する。

40

【0061】

このように、変形例の2値化ニューラルネットワーク回路200は、図9のXNORゲート回路102を、LUT202に置き換えた構成となっている。変形例では、実施形態と同様に、バッチ正規化回路の面積とパラメータを格納するメモリ面積・メモリ帯域を無くすことができ、かつ、性能的には等価な回路構成を実現することができる。

また、本変形例では、XNOR演算を行う論理ゲートとして、LUT202を用いている。LUT202は、FPGAの基本構成要素であり、FPGA合成の際の親和性が高く、FPGAによる実装が容易である。

【0062】

50

本発明は上記の実施形態例に限定されるものではなく、特許請求の範囲に記載した本発明の要旨を逸脱しない限りにおいて、他の変形例、応用例を含む。

また、上記した実施形態例は本発明をわかりやすく説明するために詳細に説明したものであり、必ずしも説明した全ての構成を備えるものに限定されるものではない。また、ある実施形態例の構成の一部を他の実施形態例の構成に置き換えることが可能であり、また、ある実施形態例の構成に他の実施形態例の構成を加えることも可能である。また、実施形態例は、その他の様々な形態で実施されることが可能であり、発明の要旨を逸脱しない範囲で、種々の省略、置き換え、変更を行うことができる。これら実施形態やその変形例は、発明の範囲や要旨に含まれるとともに、特許請求の範囲に記載された発明とその均等の範囲に含まれる。

10

【0063】

また、上記実施形態において説明した各処理のうち、自動的に行われるものとして説明した処理の全部または一部を手動的に行うこともでき、あるいは、手動的に行われるものとして説明した処理の全部または一部を公知の方法で自動的に行うこともできる。この他、上述文書中や図面中に示した処理手順、制御手順、具体的名称、各種のデータやパラメータを含む情報については、特記する場合を除いて任意に変更することができる。

また、図示した各装置の各構成要素は機能概念的なものであり、必ずしも物理的に図示の如く構成されていることを要しない。すなわち、各装置の分散・統合の具体的形態は図示のものに限られず、その全部または一部を、各種の負荷や使用状況などに応じて、任意の単位で機能的または物理的に分散・統合して構成することができる。

20

【0064】

また、上記の各構成、機能、処理部、処理手段等は、それらの一部または全部を、例えば集積回路で設計する等によりハードウェアで実現してもよい。また、上記の各構成、機能等は、プロセッサがそれぞれの機能を実現するプログラムを解釈し、実行するためのソフトウェアで実現してもよい。各機能を実現するプログラム、テーブル、ファイル等の情報は、メモリや、ハードディスク、SSD (Solid State Drive) 等の記録装置、または、IC (Integrated Circuit) カード、SD (Secure Digital) カード、光ディスク等の記録媒体に保持することができる。

また、上記実施の形態では、装置は、ニューラルネットワーク回路装置という名称を用いたが、これは説明の便宜上であり、名称はディープニューラルネットワーク回路、ニューラルネットワーク装置、パーセプトロン等であってもよい。また、方法およびプログラムは、ニューラルネットワーク処理方法という名称を用いたが、ニューラルネットワーク演算方法、ニューラルネットワークプログラム等であってもよい。

30

【符号の説明】

【0065】

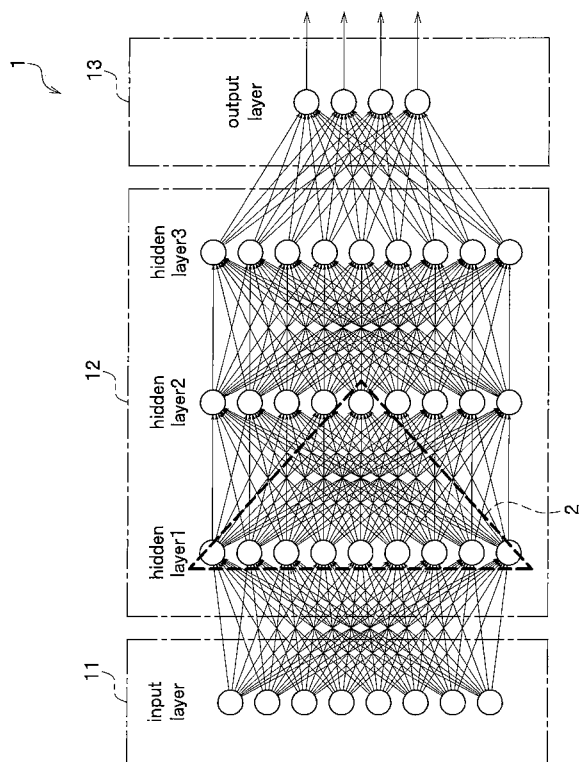
- 1 ディープニューラルネットワーク
- 2 ニューラルネットワーク回路
- 11 入力層
- 12 隠れ層 (中間層)
- 13 出力層
- 100, 200 2値化ニューラルネットワーク回路 (ニューラルネットワーク回路装置)
- 101 入力部
- 102 XNORゲート回路 (論理回路部, 論理回路手段)
- 103 総和回路 (総和回路部, 総和回路手段)
- 110 多ビットバイアス入力部
- 120 活性化回路 (活性化回路部, 活性化回路手段)
- 202 LUT (論理回路部)
- x1 ~ xn (xi) 入力値 (2値)
- w1 ~ wn (wi) 重み (2値)

40

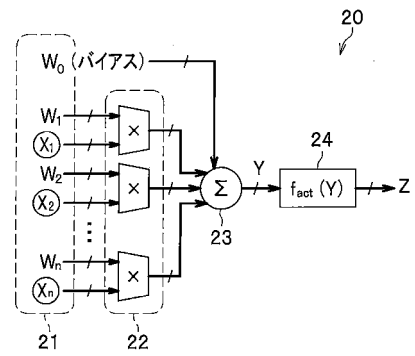
50

W' 多ビットバイアス

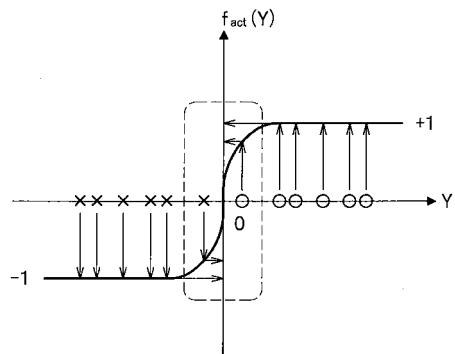
【図1】



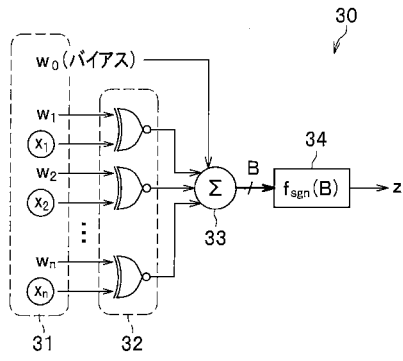
【図2】



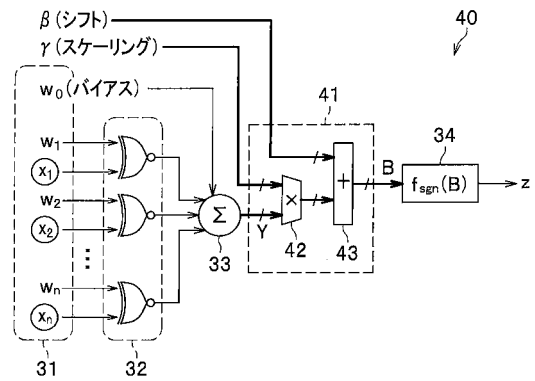
【図3】



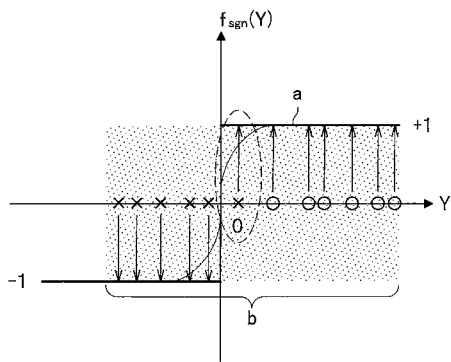
【 図 4 】



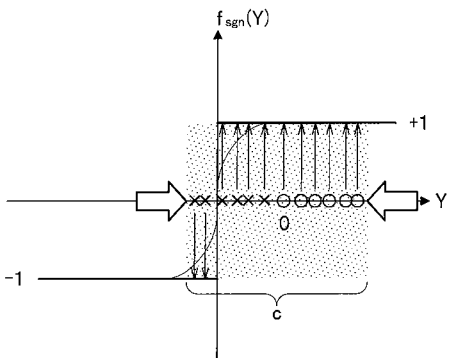
【 図 6 】



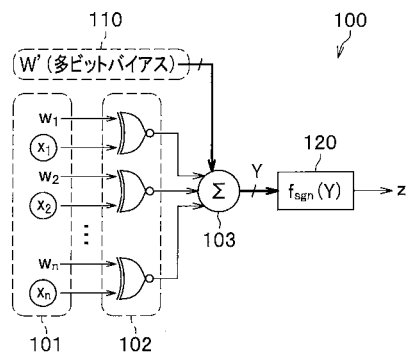
【 図 5 】



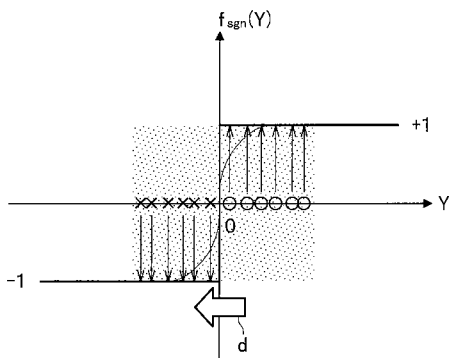
【 図 7 】



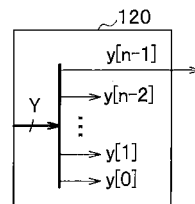
【 図 9 】



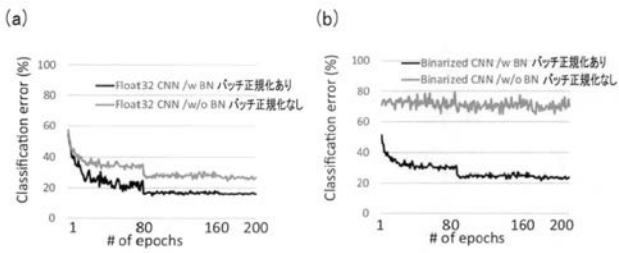
【 図 8 】



【 図 10 】



【 図 1 1 】

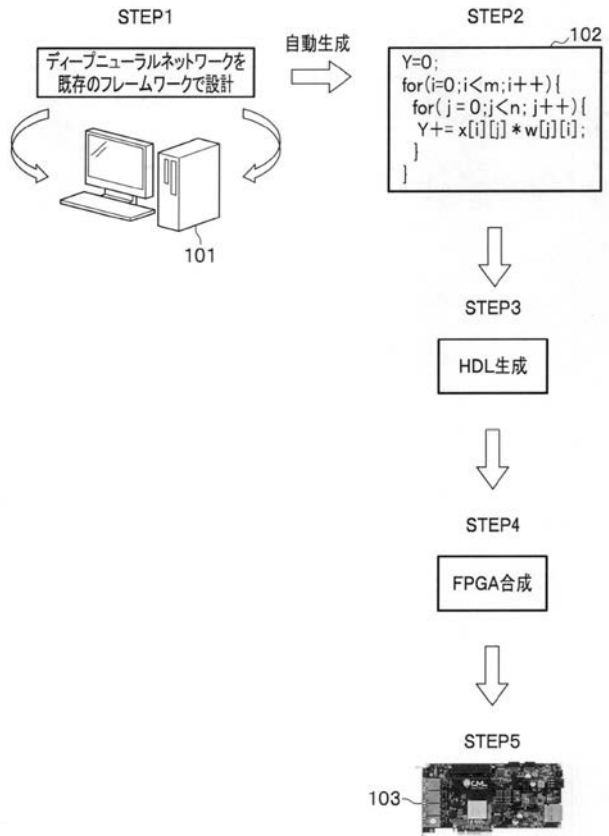


【 図 1 2 】

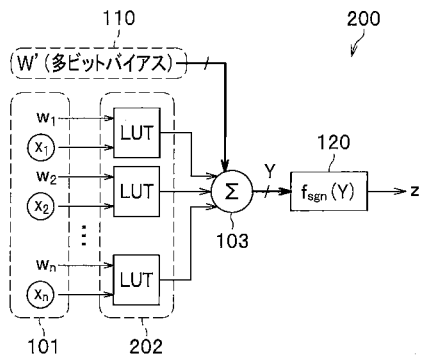
Year	2010 [1]	2014 [2]	2015 [3]	2016 [4]	提案手法
Platform	Virtex 5 SX240t	Zynq XC7Z045	Virtex7 VX485t	Zynq XC7Z045	Kintex7 325t
Clock(MHz)	120	150	100	150	450
Bandwidth(GB/s)	---	4.2	12.8	4.2	3.6
Quantization Strategy	48-bit fixed	16-bit fixed	32-bit float	16-bit fixed	binary (1bit)
Power(W)	14	8	18.61	9.63	4.48
Performance (GOP/s)	16	23.18	61.62	187.80	176.8
Resource Efficiency (GOP/s/Slices)	4.30x10 ⁻⁴	---	8.12x10 ⁻⁴	35.8x10 ⁻⁴	1389.9x10 ⁻⁴
Power Efficiency (GOP/s/W)	1.14	2.90	3.31	19.50	39.46

[1] S.Chakradhar et al., "A dynamically configurable coprocessor for convolutional neural networks," in ISCA.
 [2] V.Gokhale et al., "A 240 g-ops/s mobile coprocessor for deep neural networks," in CVPRW.
 [3] C. Zhang et al., "Optimizing fpga-based accelerator design for deep convolutional neural networks," in FPGA.
 [4] J. Qiu et al., "Going deeper with embedded FPGA platform for convolutional neural network," in FPGA.

【 図 1 3 】



【 図 1 4 】



【 図 1 5 】

LUT 202

x_1	w_1	Y
-1	-1	1
-1	+1	-1
+1	-1	-1
+1	+1	1

【手続補正書】

【提出日】平成29年6月23日(2017.6.23)

【手続補正1】

【補正対象書類名】特許請求の範囲

【補正対象項目名】全文

【補正方法】変更

【補正の内容】

【特許請求の範囲】

【請求項1】

入力層、1以上の中間層、および、出力層を少なくとも含むニューラルネットワーク回路装置であって、

前記中間層の中で、入力値 x_i および重み w_i を受け取り、論理演算を行う論理回路部と、

多ビットバイアス W' を受け取り、前記論理回路部の出力と前記多ビットバイアス W' との総和を取る総和回路部と、

総和を取った多ビット信号 Y に対して符号ビットのみを出力する活性化回路部と、を備え、

前記多ビット信号 Y および前記多ビットバイアス W' は、下記式で示される

【数3】

$$\begin{aligned}
 Y &= \sum_{i=0}^n w_i x_i - \mu_B + \frac{\sqrt{\sigma_B^2 + \epsilon}}{\gamma} \beta \\
 &= \sum_{i=1}^n w_i x_i + \left(w_0 - \mu_B + \frac{\sqrt{\sigma_B^2 + \epsilon}}{\gamma} \beta \right) \\
 &= \sum_{i=1}^n w_i x_i + W' \cdot
 \end{aligned}
 \tag{3}$$

ことを特徴とする記載のニューラルネットワーク回路装置。

【請求項2】

前記入力値 x_i および前記重み w_i を入力する入力部と、

前記多ビットバイアス W' を入力する多ビットバイアス入力部と、を備える

ことを特徴とする請求項1に記載のニューラルネットワーク回路装置。

【請求項3】

前記入力値 x_i および前記重み w_i は、2値信号である

ことを特徴とする請求項1または請求項2に記載のニューラルネットワーク回路装置。

【請求項4】

前記多ビットバイアス W' は、学習後の多ビットバイアス値である

ことを特徴とする請求項1または請求項2に記載のニューラルネットワーク回路装置。

【請求項5】

前記論理回路部は、否定排他的論理和または排他的論理和を含む

ことを特徴とする請求項1に記載のニューラルネットワーク回路装置。

【請求項6】

前記論理回路部は、LUT (Look-Up Table) である

ことを特徴とする請求項1に記載のニューラルネットワーク回路装置。

【請求項7】

前記符号ビットは、総和を取った前記多ビット信号 Y を活性化するかしないかで示す 2 値信号である

ことを特徴とする請求項 1 に記載のニューラルネットワーク回路装置。

【請求項 8】

前記多ビット信号 Y を、正規化範囲を広げ中心をシフトさせるバッチ正規化処理を行い出力される信号 Y' が、式 (1) で示される場合、

【数 1】

$$Y' = \gamma \frac{Y - \mu_B}{\sqrt{\sigma_B^2 + \epsilon}} + \beta$$

$$= \frac{\gamma}{\sqrt{\sigma_B^2 + \epsilon}} \left(Y - \left(\mu_B - \frac{\sqrt{\sigma_B^2 + \epsilon}}{\gamma} \beta \right) \right) \quad \dots(1)$$

前記多ビットバイアス W' は、

前記バッチ正規化処理による前記信号 Y' を含まない式 (3) で示される前記多ビット信号 Y で与えられる

【数 3】

$$Y = \sum_{i=0}^n w_i X_i - \mu_B + \frac{\sqrt{\sigma_B^2 + \epsilon}}{\gamma} \beta$$

$$= \sum_{i=1}^n w_i X_i + \left(w_0 - \mu_B + \frac{\sqrt{\sigma_B^2 + \epsilon}}{\gamma} \beta \right)$$

$$= \sum_{i=1}^n w_i X_i + W' \quad \dots(3)$$

ことを特徴とする請求項 1 に記載のニューラルネットワーク回路装置。

【請求項 9】

請求項 1 乃至 8 のいずれか 1 項に記載のニューラルネットワーク回路装置を備えるニューラルネットワーク。

【請求項 10】

入力層、1 以上の中間層、および、出力層を少なくとも含むニューラルネットワーク処理方法であって、

前記中間層の中で、入力値 x i および重み w i を受け取り、論理演算を行うステップと

、多ビットバイアス W' を受け取り、前記論理演算ステップの出力と前記多ビットバイアス W' との総和を取るステップと、

総和を取った多ビット信号 Y に対して符号ビットのみを出力するステップと、を有し、前記多ビット信号 Y および前記多ビットバイアス W' は、下記式で示される

【数 3】

$$\begin{aligned}
 Y &= \sum_{i=0}^n w_i X_i - \mu_B + \frac{\sqrt{\sigma_B^2 + \epsilon}}{\gamma} \beta \\
 &= \sum_{i=1}^n w_i X_i + \left(w_0 - \mu_B + \frac{\sqrt{\sigma_B^2 + \epsilon}}{\gamma} \beta \right) \\
 &= \sum_{i=1}^n w_i X_i + W' \cdot
 \end{aligned}
 \tag{3}$$

ことを特徴とするニューラルネットワーク処理方法。

【請求項 1 1】

入力層、1 以上の中間層、および、出力層を少なくとも含むニューラルネットワーク回路装置としてのコンピュータを、

前記中間層の中で、入力値 x_i および重み w_i を受け取り、論理演算を行う論理回路手段、

多ビットバイアス W' を受け取り、前記論理手段の出力と前記多ビットバイアス W' との総和を取る総和回路手段、

総和を取った多ビット信号 Y に対して符号ビットのみを出力する活性化回路手段、

ただし、前記多ビット信号 Y および前記多ビットバイアス W' は、下記式で示される、

【数 3】

$$\begin{aligned}
 Y &= \sum_{i=0}^n w_i X_i - \mu_B + \frac{\sqrt{\sigma_B^2 + \epsilon}}{\gamma} \beta \\
 &= \sum_{i=1}^n w_i X_i + \left(w_0 - \mu_B + \frac{\sqrt{\sigma_B^2 + \epsilon}}{\gamma} \beta \right) \\
 &= \sum_{i=1}^n w_i X_i + W' \cdot
 \end{aligned}
 \tag{3}$$

として機能させるためのニューラルネットワークの実行プログラム。

【手続補正 2】

【補正対象書類名】明細書

【補正対象項目名】0009

【補正方法】変更

【補正の内容】

【0009】

前記した課題を解決するため、本発明に係るニューラルネットワーク回路装置は、入力層、1 以上の中間層、および、出力層を少なくとも含むニューラルネットワーク回路装置であって、前記中間層の中で、入力値 x_i および重み w_i を受け取り、論理演算を行う論理回路部と、多ビットバイアス W' を受け取り、前記論理回路部の出力と前記多ビットバイアス W' との総和を取る総和回路部と、総和を取った多ビット信号 Y に対して符号ビッ

トのみを出力する活性化回路部と、を備え、前記多ビット信号 Y および前記多ビットバイアス W' は、下記式で示される

【数 3】

$$\begin{aligned}
 Y &= \sum_{i=0}^n w_i X_i - \mu_B + \frac{\sqrt{\sigma_B^2 + \epsilon}}{\gamma} \beta \\
 &= \sum_{i=1}^n w_i X_i + \left(w_0 - \mu_B + \frac{\sqrt{\sigma_B^2 + \epsilon}}{\gamma} \beta \right) \\
 &= \sum_{i=1}^n w_i X_i + W' .
 \end{aligned}
 \tag{3}$$

ことを特徴とする。