

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B1)

(11) 特許番号

特許第6054005号
(P6054005)

(45) 発行日 平成28年12月27日 (2016. 12. 27)

(24) 登録日 平成28年12月9日 (2016. 12. 9)

(51) Int. Cl. F 1
G 0 6 N 3 / 0 4 (2006. 01) G 0 6 N 3 / 0 4

請求項の数 18 (全 27 頁)

(21) 出願番号	特願2016-548332 (P2016-548332)	(73) 特許権者	000006013
(86) (22) 出願日	平成27年8月31日 (2015. 8. 31)		三菱電機株式会社
(86) 国際出願番号	PCT/JP2015/074720		東京都千代田区丸の内二丁目7番3号
審査請求日	平成28年7月25日 (2016. 7. 25)	(74) 代理人	100123434
(31) 優先権主張番号	特願2015-113440 (P2015-113440)		弁理士 田澤 英昭
(32) 優先日	平成27年6月3日 (2015. 6. 3)	(74) 代理人	100101133
(33) 優先権主張国	日本国 (JP)		弁理士 濱田 初音
早期審査対象出願		(74) 代理人	100199749
			弁理士 中島 成
		(74) 代理人	100188880
			弁理士 坂元 辰哉
		(74) 代理人	100197767
			弁理士 辻岡 将昭
		(74) 代理人	100201743
			弁理士 井上 和真

最終頁に続く

(54) 【発明の名称】 推論装置及び推論方法

(57) 【特許請求の範囲】

【請求項1】

ニューラルネットワークを構成している入力層の各ノードにデータが与えられると、前記データから前記入力層の各ノードでの活性度を算出する入力層活性度算出部と、

前記ニューラルネットワークを構成している中間層のノードと前記入力層のノードとを接続している各エッジの重みを記憶している中間層記憶部と、

前記入力層活性度算出部により算出された入力層の各ノードでの活性度及び前記中間層記憶部に記憶されている各エッジの重みの中から、前記中間層の各ノードと接続関係がある前記入力層の各ノードでの活性度と各エッジの重みとを取得し、前記取得した入力層の各ノードでの活性度と各エッジの重みとを用いて、前記中間層の各ノードでの活性度を算出する中間層活性度算出部と、

前記中間層活性度算出部により算出された中間層の各ノードでの活性度を用いて、前記ニューラルネットワークを構成している出力層の各ノードでの活性度を算出する出力層活性度算出部と

を備えており、かつ、中間層の各ノード当りの前記入力層との平均接続本数が50本以下であることを特徴とする推論装置。

【請求項2】

前記ニューラルネットワークを構成している中間層が複数存在しており、

前記中間層記憶部は、前記ニューラルネットワークを構成している中間層毎に、当該中間層のノードが前記入力層のノードと接続されていれば、当該中間層のノードと前記入力

層のノードとを接続している各エッジの重みを記憶し、当該中間層のノードが他の中間層のノードと接続されていれば、当該中間層のノードと他の中間層のノードとを接続している各エッジの重みを記憶しており、

前記中間層活性度算出部は、前記ニューラルネットワークを構成している中間層のノードと接続されているノードが前記入力層のノードであれば、前記入力層活性度算出部により算出された入力層の各ノードでの活性度及び前記中間層記憶部に記憶されている各エッジの重みの中から、前記ニューラルネットワークを構成している中間層の各ノードと接続関係がある前記入力層の各ノードでの活性度と各エッジの重みとを取得し、前記ニューラルネットワークを構成している中間層のノードと接続されているノードが他の中間層のノードであれば、他の中間層のノードでの活性度及び前記中間層記憶部に記憶されている各エッジの重みの中から、前記ニューラルネットワークを構成している中間層の各ノードと接続関係がある他の中間層のノードでの活性度と各エッジの重みとを取得し、前記取得した入力層又は他の中間層のノードでの活性度と各エッジの重みとを用いて、前記ニューラルネットワークを構成している中間層の各ノードでの活性度を算出し、かつ、全部または一部の中間層の各ノード当りの前記中間層との平均接続本数が50本以下であることを特徴とする請求項1記載の推論装置。

10

【請求項3】

前記出力層のノードと前記中間層のノードとを接続している各エッジの重みを記憶している出力層記憶部を備え、

前記出力層活性度算出部は、前記中間層活性度算出部により算出された中間層の各ノードでの活性度及び前記出力層記憶部に記憶されている各エッジの重みの中から、前記出力層の各ノードと接続関係がある前記中間層の各ノードでの活性度と各エッジの重みとを取得し、前記取得した中間層の各ノードでの活性度と各エッジの重みとを用いて、前記出力層の各ノードでの活性度を算出することを特徴とする請求項1あるいは請求項2記載の推論装置。

20

【請求項4】

前記出力層のノードが前記入力層のノードと接続されていれば、前記出力層のノードと前記入力層のノードとを接続している各エッジの重みを記憶し、前記出力層のノードが前記中間層のノードと接続されていれば、前記出力層のノードと前記中間層のノードとを接続している各エッジの重みを記憶している出力層記憶部を備え、

30

前記出力層活性度算出部は、前記出力層のノードと接続されているノードが前記入力層のノードであれば、前記入力層活性度算出部により算出された入力層の各ノードでの活性度及び前記出力層記憶部に記憶されている各エッジの重みの中から、前記出力層の各ノードと接続関係がある前記入力層の各ノードでの活性度と各エッジの重みとを取得し、前記出力層のノードと接続されているノードが前記中間層のノードであれば、前記中間層活性度算出部により算出された中間層のノードでの活性度及び前記出力層記憶部に記憶されている各エッジの重みの中から、前記出力層の各ノードと接続関係がある中間層のノードでの活性度と各エッジの重みとを取得し、前記取得した入力層又は中間層のノードでの活性度と各エッジの重みとを用いて、前記出力層の各ノードでの活性度を算出することを特徴とする請求項1あるいは請求項2記載の推論装置。

40

【請求項5】

前記中間層記憶部は、前記各エッジの重みのほかに、前記中間層の各ノードに与えられているバイアス値を記憶しており、

前記中間層活性度算出部は、前記入力層の各ノードでの活性度と各エッジの重みと前記バイアス値とを用いて、前記中間層の各ノードでの活性度を算出することを特徴とする請求項1あるいは請求項2記載の推論装置。

【請求項6】

前記中間層活性度算出部は、前記中間層のノード毎に、当該ノードと接続関係がある前記入力層の各ノードでの活性度と、前記中間層の当該ノードと前記入力層の各ノードを接続している各エッジの重みとの積和演算を実施して、前記積和演算の演算結果と前記中間

50

層の当該ノードのバイアス値を加算し、その加算した結果を前記ニューラルネットワークの活性化関数の引数として用いることで、前記中間層の当該ノードの活性度として、前記活性化関数の関数値を算出することを特徴とする請求項 5 記載の推論装置。

【請求項 7】

前記出力層記憶部は、前記各エッジの重みのほかに、前記出力層の各ノードに与えられているバイアス値を記憶しており、

前記出力層活性度算出部は、前記中間層の各ノードでの活性度と各エッジの重みと前記バイアス値とを用いて、前記出力層の各ノードでの活性度を算出することを特徴とする請求項 3 記載の推論装置。

【請求項 8】

前記出力層活性度算出部は、前記出力層のノード毎に、当該ノードと接続関係がある前記中間層の各ノードでの活性度と、前記出力層の当該ノードと前記中間層の各ノードを接続している各エッジの重みとの積和演算を実施して、前記積和演算の演算結果と前記出力層の当該ノードのバイアス値を加算し、その加算した結果を前記ニューラルネットワークの活性化関数の引数として用いることで、前記出力層の当該ノードの活性度として、前記活性化関数の関数値を算出することを特徴とする請求項 7 記載の推論装置。

【請求項 9】

前記ニューラルネットワークを構成している入力層及び中間層のノードを接続している複数のエッジ、あるいは、前記ニューラルネットワークを構成している中間層及び出力層のノードを接続している複数のエッジがループをなしており、前記ループが 6 本以上のエッジで形成されていることを特徴とする請求項 1 あるいは請求項 2 記載の推論装置。

【請求項 10】

前記ニューラルネットワークを構成している複数の中間層のノードを接続している複数のエッジがループをなしており、前記ループが 6 本以上のエッジで形成されていることを特徴とする請求項 2 記載の推論装置。

【請求項 11】

前記ニューラルネットワークを構成している中間層の各ノードは、前記入力層における全ノードのうち、ランダムに選択された一部のノードと接続されており、

前記ニューラルネットワークを構成している出力層の各ノードは、前記中間層における全ノードのうち、ランダムに選択された一部のノードと接続されていることを特徴とする請求項 1 あるいは請求項 2 記載の推論装置。

【請求項 12】

前記ニューラルネットワークを構成している中間層の各ノードは、前記入力層又は他の中間層における全ノードのうち、ランダムに選択された一部のノードと接続されていることを特徴とする請求項 2 記載の推論装置。

【請求項 13】

前記ニューラルネットワークを構成している中間層の各ノードは、前記入力層における全ノードのうち、隣接していない一部のノードと接続されており、

前記ニューラルネットワークを構成している出力層の各ノードは、前記中間層における全ノードのうち、隣接していない一部のノードと接続されていることを特徴とする請求項 1 あるいは請求項 2 記載の推論装置。

【請求項 14】

前記ニューラルネットワークを構成している複数の中間層の各ノードは、前記入力層又は他の中間層における全ノードのうち、隣接していない一部のノードと接続されていることを特徴とする請求項 2 記載の推論装置。

【請求項 15】

前記ニューラルネットワークを構成している中間層の各ノード当りの前記入力層のノードとの平均接続本数が前記入力層のノードの個数の 10 分の 1 以下であることを特徴とする請求項 1 あるいは請求項 2 記載の推論装置。

【請求項 16】

10

20

30

40

50

前記ニューラルネットワークを構成している複数の中間層の各ノード当りの前記入力層又は他の中間層のノードとの平均接続本数が前記入力層又は他の中間層のノードの個数の10分の1以下であることを特徴とする請求項2記載の推論装置。

【請求項17】

中間層記憶部が、ニューラルネットワークを構成している中間層のノードと入力層のノードとを接続している各エッジの重みを記憶しており、

入力層活性度算出部が、前記ニューラルネットワークを構成している入力層の各ノードにデータが与えられると、前記データから前記入力層の各ノードでの活性度を算出し、

中間層活性度算出部が、前記入力層活性度算出部により算出された入力層の各ノードでの活性度及び前記中間層記憶部に記憶されている各エッジの重みの中から、前記中間層の各ノードと接続関係がある前記入力層の各ノードでの活性度と各エッジの重みとを取得し、前記取得した入力層の各ノードでの活性度と各エッジの重みとを用いて、前記中間層の各ノードでの活性度を算出し、

出力層活性度算出部が、前記中間層活性度算出部により算出された中間層の各ノードでの活性度を用いて、前記ニューラルネットワークを構成している出力層の各ノードでの活性度を算出し、かつ、中間層の各ノード当りの前記入力層との平均接続本数が50本以下であることを特徴とする推論方法。

【請求項18】

前記ニューラルネットワークを構成している中間層が複数存在しており、

前記中間層記憶部は、前記ニューラルネットワークを構成している中間層毎に、当該中間層のノードが前記入力層のノードと接続されていれば、当該中間層のノードと前記入力層のノードとを接続している各エッジの重みを記憶し、当該中間層のノードが他の中間層のノードと接続されていれば、当該中間層のノードと他の中間層のノードとを接続している各エッジの重みを記憶しており、

前記中間層活性度算出部は、前記ニューラルネットワークを構成している中間層のノードと接続されているノードが前記入力層のノードであれば、前記入力層活性度算出部により算出された入力層の各ノードでの活性度及び前記中間層記憶部に記憶されている各エッジの重みの中から、前記ニューラルネットワークを構成している中間層の各ノードと接続関係がある前記入力層の各ノードでの活性度と各エッジの重みとを取得し、前記ニューラルネットワークを構成している中間層のノードと接続されているノードが他の中間層のノードであれば、他の中間層のノードでの活性度及び前記中間層記憶部に記憶されている各エッジの重みの中から、前記ニューラルネットワークを構成している中間層の各ノードと接続関係がある他の中間層のノードでの活性度と各エッジの重みとを取得し、前記取得した入力層又は他の中間層のノードでの活性度と各エッジの重みとを用いて、前記ニューラルネットワークを構成している中間層の各ノードでの活性度を算出することを特徴とし、かつ、全部または一部の中間層の各ノード当りの前記中間層との平均接続本数が50本以下であることを特徴とする請求項17記載の推論方法。

【発明の詳細な説明】

【技術分野】

【0001】

この発明は、ニューラルネットワークを用いた推論装置及び推論方法に関するものである。

【背景技術】

【0002】

機械学習手法の1つとして、ニューラルネットワークは高い問題解決能力を有し、画像認識、音声認識、異常検知、将来予測などの多くの処理に用いられている。

ニューラルネットワークの構造の1つとして、階層型ニューラルネットワークがあり、学習手法として、主に教師あり学習と、教師なし学習の二種類がある。

教師あり学習は、複数の学習例の入力データと目標出力を与え、実際の出力と目標出力が一致するように、ニューラルネットワークの結合状態を調整する手法である。また、教

10

20

30

40

50

教師なし学習は、目標出力を与えずに、学習例が有する本質的な特徴を抽出できるように、ニューラルネットワークの結合状態を調整する手法である。

【0003】

例えば、教師あり学習法に属する誤差逆伝播法（バックプロパゲーションアルゴリズム）には、ニューラルネットワークの階層の数が多くなると、学習結果が収束しなくなる問題を発生することがある。

上記の問題を解決するために、例えば、自己符号化器（Autoencoder）や制約ボルツマンマシン（Restricted Boltzmann Machine）などの教師なし学習を用いて層毎の事前学習（Pre-training）を実施することで、ニューラルネットワークの結合状態の初期値を決定し、その後、誤差逆伝播法を用いて、
10 ニューラルネットワークの結合状態を調整（Fine-tuning）するようにしているものがある。

これにより、学習結果が収束しなくなる問題の発生を招くことなく、実際の出力と目標出力が一致するように、ニューラルネットワークの結合状態を調整することができる。

【0004】

階層型ニューラルネットワークは、複数のノード（節点）及びノード間の結合するエッジ（枝）で構成されるグラフ構造で表すことができるが、例えば、4層のニューラルネットワークでは、複数のノードが入力層、第1中間層、第2中間層、出力層で階層化され、同一階層に属するノード間のエッジは存在せずに、隣接している層の間だけにエッジが存在する。中間層は隠れ層と呼ばれることがある。
20

各エッジには、繋いだ2つのノード間の結合度合を示すパラメータが存在し、そのパラメータは、エッジ重みと呼ばれている。

【0005】

階層型ニューラルネットワークを用いて、学習または推論を行う際、その計算量及びメモリ量はエッジ数に比例する。一般的に、各層に属するノードは、隣接している層に属する全てのノードとエッジで接続されているため、計算量及びメモリ量がノード数と直接に関係する。

例えば、入力層のノード数が N 、第1中間層のノード数が M_1 、第2中間層のノード数が M_2 、出力層のノード数が1である場合には、入力層と第1中間層の間エッジ数が $N \times M_1$ 、第1中間層と第2中間層の間エッジ数が $M_1 \times M_2$ 、第2中間層と出力層の間エッジ数が M_2 となるため、学習または推論を行う際の計算量及びメモリ量が、 $(N \times M_1 + M_1 \times M_2 + M_2)$ に比例する。
30

【0006】

特に、中間層のノード数が入力層のノード数に比例する場合、ノード数が N 個の入力層に対して、第1中間層のノード数が $M_1 = a \times N$ 個、第2中間層のノード数が $M_2 = b \times N$ 個となる。この場合、ニューラルネットワークにおけるエッジの総数が $N \times a \times N + a \times N \times b \times N + b \times N = (a + a \times b) \times N^2 + b \times N$ になり、学習または推論を行う際の計算量及びメモリ量が、 $(a + a \times b) \times N^2 + b \times N$ に比例する。

【0007】

階層型ニューラルネットワークは、以上のような構造を持つことが多く、計算量及びメモリ量が、入力データ数である N の2乗、即ち、入力層のノード数 N の2乗に比例して増加するため、入力データ数の増大と共に計算量及びメモリ量が飛躍的に増加し、計算機リソース不足、処理遅延、装置コスト増大などの問題が発生する。
40

以下の特許文献1には、複数の入力データの相関関係に基づいて、複数の入力データをグループ化することで、入力層と中間層の間エッジ数や、中間層と出力層の間エッジ数を削減している。

【先行技術文献】

【特許文献】

【0008】

【特許文献1】特開2011-54200号公報（図1）

10

20

30

40

50

【発明の概要】

【発明が解決しようとする課題】

【0009】

従来の推論装置は以上のように構成されているので、入力層と中間層の間のエッジ数や、中間層と出力層の間のエッジ数を削減することができる。しかし、同一グループに属する入力層と中間層の間では、入力層の各ノードが、中間層の全てのノードと接続されるため、エッジの削減数が限定的であり、依然として、推論を行う際の計算量及びメモリ量が大きくなってしまいう課題があった。

【0010】

この発明は上記のような課題を解決するためになされたもので、推論を行う際の計算量及びメモリ量を削減することができる推論装置及び推論方法を得ることを目的とする。また、推論精度が高い推論装置及び推論方法を得ることを目的とする。

【課題を解決するための手段】

【0011】

この発明に係る推論装置は、ニューラルネットワークを構成している入力層の各ノードにデータが与えられると、そのデータから入力層の各ノードでの活性度を算出する入力層活性度算出部と、ニューラルネットワークを構成している中間層のノードと入力層のノードとを接続している各エッジの重みを記憶している中間層記憶部と、入力層活性度算出部により算出された入力層の各ノードでの活性度及び中間層記憶部に記憶されている各エッジの重みの中から、中間層の各ノードと接続関係がある入力層の各ノードでの活性度と各エッジの重みとを取得し、その取得した入力層の各ノードでの活性度と各エッジの重みとを用いて、中間層の各ノードでの活性度を算出する中間層活性度算出部と、中間層活性度算出部により算出された中間層の各ノードでの活性度を用いて、ニューラルネットワークを構成している出力層の各ノードでの活性度を算出する出力層活性度算出部とを備えており、かつ、中間層の各ノード当りの前記入力層との平均接続本数が50本以下であることを特徴とするようにしたものである。

【発明の効果】

【0012】

この発明によれば、推論を行う際の計算量及びメモリ量を削減することができる効果がある。また、他の効果として、より高い推論精度を得ることができる。

【図面の簡単な説明】

【0013】

【図1】この発明の実施の形態1による推論装置を示す構成図である。

【図2】この発明の実施の形態1による推論装置を示すハードウェア構成図である。

【図3】推論装置がコンピュータで構成される場合のハードウェア構成図である。

【図4】この発明の実施の形態1による推論装置の処理内容である推論方法を示すフローチャートである。

【図5】入力層活性度算出部1、第1中間層活性度算出部5、第2中間層活性度算出部6、第3中間層活性度算出部7及び出力層活性度算出部9の処理内容を示すフローチャートである。

【図6】この発明の実施の形態1による推論装置が適用する階層型ニューラルネットワークを示す説明図である。

【図7】インデックス、エッジ重み及びバイアス値の一例を示す説明図である。

【図8】複数のエッジによって形成されるループの例を示す説明図である。

【図9】入力層に与えられる画像データを10種のクラスに識別するニューラルネットワークの一例を示す説明図である。

【図10】入力層から第1中間層、第1中間層から第2中間層（文脈層）、第2中間層（文脈層）から第1中間層、第1中間層から出力層に枝接続を持っているリカレントニューラルネットワーク（RNN）であるエルマンネット（Elman Network）を示す説明図である。

10

20

30

40

50

【図11】この発明の実施の形態7による推論装置を示す構成図である。

【図12】中間層内でノード同士の枝接続や自己接続が存在するとともに、入力層から中間層を飛ばして出力層へ接続する枝が存在するニューラルネットワークであるEcho State Networkの例を示す説明図である。

【図13】図12のEcho State Networkを層単位で表している説明図である。

【図14】この発明の実施の形態8による推論装置を示す構成図である。

【発明を実施するための形態】

【0014】

以下、この発明をより詳細に説明するために、この発明を実施するための形態について、添付の図面にしたがって説明する。 10

【0015】

実施の形態1.

図1はこの発明の実施の形態1による推論装置を示す構成図であり、図2はこの発明の実施の形態1による推論装置を示すハードウェア構成図である。

図1では、複数のノードが、入力層、第1中間層、第2中間層、第3中間層、出力層で階層化されている5層の階層型ニューラルネットワークを用いる推論装置の例を示している。また、図1の例では、入力層に与えられるデータが画像データである例を示している。

ここでは、5層の階層型ニューラルネットワークを用いる例を示しているが、5層の階層型ニューラルネットワークに限るものではなく、3層や4層、あるいは、6層以上の階層型ニューラルネットワークを用いるものであってもよい。 20

因みに、3層の階層型ニューラルネットワークを用いる場合、中間層は第1中間層だけになり、後述する第2中間層記憶部3、第3中間層記憶部4、第2中間層活性度算出部6及び第3中間層活性度算出部7が不要になる。

【0016】

図1及び図2において、入力層活性度算出部1は例えばCPU(Central Processing Unit)を実装している半導体集積回路、あるいは、ワンチップマイコンなどで構成されている入力層活性度算出回路11で実現されるものであり、階層型ニューラルネットワークを構成している入力層の各ノードに画像データが与えられると、その画像データから入力層の各ノードでの活性度を算出する処理を実施する。 30

【0017】

中間層記憶部である第1中間層記憶部2は例えばRAMやハードディスクなどの記憶媒体からなる中間層記憶装置12で実現されるものであり、第1中間層の各ノードと入力層の各ノードとの接続関係を示すインデックス(接続情報)と、第1中間層のノードと入力層のノードとを接続している各エッジの重みと、第1中間層の各ノードに与えられているバイアス値とを記憶している。

中間層記憶部である第2中間層記憶部3は例えばRAMやハードディスクなどの記憶媒体からなる中間層記憶装置12で実現されるものであり、第2中間層の各ノードと第1中間層の各ノードとの接続関係を示すインデックスと、第2中間層のノードと第1中間層のノードとを接続している各エッジの重みと、第2中間層の各ノードに与えられているバイアス値とを記憶している。 40

中間層記憶部である第3中間層記憶部4は例えばRAMやハードディスクなどの記憶媒体からなる中間層記憶装置12で実現されるものであり、第3中間層の各ノードと第2中間層の各ノードとの接続関係を示すインデックスと、第3中間層のノードと第2中間層のノードとを接続している各エッジの重みと、第3中間層の各ノードに与えられているバイアス値とを記憶している。

【0018】

中間層活性度算出部である第1中間層活性度算出部5は例えばCPUを実装している半導体集積回路、あるいは、ワンチップマイコンなどで構成されている中間層活性度算出回 50

路 1 3 で実現されるものであり、第 1 中間層記憶部 2 に記憶されているインデックスを参照して、入力層活性度算出部 1 により算出された入力層の各ノードでの活性度と第 1 中間層記憶部 2 に記憶されている各エッジの重み及びバイアス値の中から、第 1 中間層の各ノードと接続関係がある入力層の各ノードでの活性度と各エッジの重みとバイアス値とを取得し、その取得した入力層の各ノードでの活性度と各エッジの重みとバイアス値とを用いて、第 1 中間層の各ノードでの活性度を算出する処理を実施する。

【 0 0 1 9 】

中間層活性度算出部である第 2 中間層活性度算出部 6 は例えば CPU を実装している半導体集積回路、あるいは、ワンチップマイコンなどで構成されている中間層活性度算出回路 1 3 で実現されるものであり、第 2 中間層記憶部 3 に記憶されているインデックスを参照して、第 1 中間層活性度算出部 5 により算出された第 1 中間層の各ノードでの活性度と第 2 中間層記憶部 3 に記憶されている各エッジの重み及びバイアス値の中から、第 2 中間層の各ノードと接続関係がある第 1 中間層の各ノードでの活性度と各エッジの重みとバイアス値とを取得し、その取得した第 1 中間層の各ノードでの活性度と各エッジの重みとバイアス値とを用いて、第 2 中間層の各ノードでの活性度を算出する処理を実施する。

【 0 0 2 0 】

中間層活性度算出部である第 3 中間層活性度算出部 7 は例えば CPU を実装している半導体集積回路、あるいは、ワンチップマイコンなどで構成されている中間層活性度算出回路 1 3 で実現されるものであり、第 3 中間層記憶部 4 に記憶されているインデックスを参照して、第 2 中間層活性度算出部 6 により算出された第 2 中間層の各ノードでの活性度と第 3 中間層記憶部 4 に記憶されている各エッジの重み及びバイアス値の中から、第 3 中間層の各ノードと接続関係がある第 2 中間層の各ノードでの活性度と各エッジの重みとバイアス値とを取得し、その取得した第 2 中間層の各ノードでの活性度と各エッジの重みとバイアス値とを用いて、第 3 中間層の各ノードでの活性度を算出する処理を実施する。

【 0 0 2 1 】

出力層記憶部 8 は例えば RAM やハードディスクなどの記憶媒体からなる出力層記憶装置 1 4 で実現されるものであり、出力層の各ノードと第 3 中間層の各ノードとの接続関係を示すインデックス（接続情報）と、出力層のノードと第 3 中間層のノードとを接続している各エッジの重みと、出力層の各ノードに与えられているバイアス値とを記憶している。

出力層活性度算出部 9 は例えば CPU を実装している半導体集積回路、あるいは、ワンチップマイコンなどで構成されている出力層活性度算出回路 1 5 で実現されるものであり、出力層記憶部 8 に記憶されているインデックスを参照して、第 3 中間層活性度算出部 7 により算出された第 3 中間層の各ノードでの活性度と出力層記憶部 8 に記憶されている各エッジの重み及びバイアス値の中から、出力層の各ノードと接続関係がある第 3 中間層の各ノードでの活性度と各エッジの重みとバイアス値とを取得し、その取得した第 3 中間層の各ノードでの活性度と各エッジの重みとバイアス値とを用いて、出力層の各ノードでの活性度を算出する処理を実施する。

【 0 0 2 2 】

図 1 では、推論装置の構成要素である入力層活性度算出部 1、第 1 中間層記憶部 2、第 2 中間層記憶部 3、第 3 中間層記憶部 4、第 1 中間層活性度算出部 5、第 2 中間層活性度算出部 6、第 3 中間層活性度算出部 7、出力層記憶部 8 及び出力層活性度算出部 9 のそれぞれが専用のハードウェアで構成されているものを想定しているが、推論装置がコンピュータで構成されていてもよい。

図 3 は推論装置がコンピュータで構成される場合のハードウェア構成図である。

推論装置がコンピュータで構成される場合、第 1 中間層記憶部 2、第 2 中間層記憶部 3、第 3 中間層記憶部 4 及び出力層記憶部 8 をコンピュータのメモリ 2 1 上に構成するとともに、入力層活性度算出部 1、第 1 中間層活性度算出部 5、第 2 中間層活性度算出部 6、第 3 中間層活性度算出部 7 及び出力層活性度算出部 9 の処理内容を記述しているプログラムをコンピュータのメモリ 2 1 に格納し、当該コンピュータのプロセッサ 2 2 がメモリ 2

10

20

30

40

50

1 に格納されているプログラムを実行するようにすればよい。

図 4 はこの発明の実施の形態 1 による推論装置の処理内容である推論方法を示すフローチャートであり、図 5 は入力層活性度算出部 1、第 1 中間層活性度算出部 5、第 2 中間層活性度算出部 6、第 3 中間層活性度算出部 7 及び出力層活性度算出部 9 の処理内容を示すフローチャートである。

【 0 0 2 3 】

図 6 はこの発明の実施の形態 1 による推論装置が適用する階層型ニューラルネットワークを示す説明図である。

図 6 の階層型ニューラルネットワークは、複数のノードが、入力層、第 1 中間層、第 2 中間層、第 3 中間層、出力層で階層化されている 5 層の階層型ニューラルネットワークである。

10

図 7 はインデックス、エッジ重み及びバイアス値の一例を示す説明図である。

図 7 では、ノードの接続情報であるインデックスが、例えば、第 1 中間層における“ N ”のノードは、入力層における“ 0 ”，“ 3 ”，“ 5 ”のノードと接続されている旨を示している。

また、図 7 では、例えば、第 1 中間層における“ N ”のノードと、入力層における“ 0 ”のノードとを接続しているエッジの重みが“ 0 . 2 ”、第 1 中間層における“ N ”のノードと、入力層における“ 3 ”のノードとを接続しているエッジの重みが“ - 0 . 5 ”、第 1 中間層における“ N ”のノードと、入力層における“ 5 ”のノードとを接続しているエッジの重みが“ 0 . 1 ”である旨を示している。

20

さらに、図 7 では、例えば、第 1 中間層における“ N ”のノードのバイアス値が“ 1 . 8 ”である旨を示している。

【 0 0 2 4 】

次に動作について説明する。

入力層活性度算出部 1 は、階層型ニューラルネットワークを構成している入力層の各ノードに画像データが与えられると、その画像データから入力層の各ノードでの活性度 A_{IN} を算出する（図 4 のステップ S T 1）。

入力層活性度算出部 1 に与えられる画像データが、例えば、0 ~ 255 の画素値 P を有するピクセルからなる画像を示すデータであり、各ピクセルの画素値 P が入力層の各ノードに与えられる場合、入力層の各ノードでの活性度 A_{IN} は、下記の式（1）のように算出することができる。

30

$$A_{IN} = \frac{P}{255} \quad (1)$$

ここでは、画像データが入力される場合を想定し、各ピクセルの画素値 P を 255 で除算することで正規化して、浮動小数点値（0 . 0 ~ 1 . 0）を入力層の各ノードでの活性度 A_{IN} とする例を示しているが、単なる正規化のみでなく、入力されるデータの種類に応じて、データ間引き、量子化、変換等の処理を実施するようにしてもよい。

【 0 0 2 5 】

第 1 中間層活性度算出部 5 は、入力層活性度算出部 1 が入力層の各ノードでの活性度 A_{IN} を算出すると、第 1 中間層記憶部 2 に記憶されているインデックスを参照して、第 1 中間層のノード毎に、当該ノードに接続されている入力層の各ノードを確認して、その入力層の各ノードでの活性度 A_{IN} を取得する。

40

例えば、第 1 中間層における“ N ”のノードの場合、第 1 中間層記憶部 2 に記憶されているインデックスが、入力層における“ 0 ”，“ 3 ”，“ 5 ”のノードと接続されている旨を示しているため、入力層活性度算出部 1 により算出された入力層の各ノードでの活性度 A_{IN} のうち、入力層における“ 0 ”，“ 3 ”，“ 5 ”のノードでの活性度 A_{IN-0} 、 A_{IN-3} 、 A_{IN-5} を取得する。

【 0 0 2 6 】

また、第 1 中間層活性度算出部 5 は、第 1 中間層記憶部 2 に記憶されているインデックスを参照して、第 1 中間層のノード毎に、当該ノードに接続されているエッジを確認して

50

、第1中間層記憶部2からそのエッジの重み w を取得する。

例えば、第1中間層における“N”のノードの場合、第1中間層記憶部2に記憶されているインデックスが、入力層における“0”、“3”、“5”のノードと接続されている旨を示しているため、第1中間層における“N”のノードと、入力層における“0”のノードとを接続しているエッジの重み w_{N-0} として“0.2”を取得し、第1中間層における“N”のノードと、入力層における“3”のノードとを接続しているエッジの重み w_{N-3} として“-0.5”を取得する。また、第1中間層における“N”のノードと、入力層における“5”のノードとを接続しているエッジの重み w_{N-5} として“0.1”を取得する。

【0027】

また、第1中間層活性度算出部5は、第1中間層のノード毎に、第1中間層記憶部2から当該ノードのバイアス値 B_{1M} を取得する。

例えば、第1中間層における“N”のノードの場合、バイアス値 B_{1M-N} として“1.8”を取得する。

【0028】

第1中間層活性度算出部5は、第1中間層のノード毎に、入力層の各ノードでの活性度 A_{IN} 、エッジの重み w 、バイアス値 B_{1M} を取得すると、活性度 A_{IN} 、エッジの重み w 、バイアス値 B_{1M} を用いて、第1中間層のノード毎の活性度 A_{1M} を算出する(ステップST2)。

以下、第1中間層における“N”のノードの活性度 A_{1M-N} の算出例を具体的に説明する。

【0029】

まず、第1中間層活性度算出部5は、第1中間層記憶部2に記憶されているインデックスを読み込んで(図5のステップST11)、そのインデックスを参照することで、活性度の算出に用いるパラメータとして、入力層における“0”、“3”、“5”のノードでの活性度 A_{IN-0} 、 A_{IN-3} 、 A_{IN-5} と、エッジの重み w_{N-0} 、 w_{N-3} 、 w_{N-5} と、第1中間層における“N”のノードのバイアス値 B_{1M-N} とを取得する(ステップST12)。

次に、第1中間層活性度算出部5は、下記の式(2)に示すように、入力層における“0”、“3”、“5”のノードでの活性度 A_{IN-0} 、 A_{IN-3} 、 A_{IN-5} と、エッジの重み w_{N-0} 、 w_{N-3} 、 w_{N-5} との積和演算を実施する(ステップST13)。

$$MADD = A_{IN-0} \times w_{N-0} + A_{IN-3} \times w_{N-3} + A_{IN-5} \times w_{N-5} \quad (2)$$

次に、第1中間層活性度算出部5は、下記の式(3)に示すように、積和演算の演算結果MADDと第1中間層における“N”のノードのバイアス値 B_{1M-N} とを加算する(ステップST14)。

$$ADD = MADD + B_{1M-N} \quad (3)$$

【0030】

第1中間層活性度算出部5は、活性度の計算に用いるニューラルネットワークの活性化関数 $F(\text{activation function})$ として、線形関数、シグモイド関数、ソフトマックス関数、整流化線形関数(ReLU)などを事前に用意しており、式(3)の加算結果ADDを活性化関数 F の引数とし、下記の式(4)に示すように、第1中間層における“N”のノードの活性度 A_{1M-N} として、その活性化関数の関数値を算出する(ステップST15)。

$$A_{1M-N} = F(ADD) \quad (4)$$

ここでは、第1中間層における“N”のノードの活性度 A_{1M-N} の算出例を示したが、第1中間層における他のノードの活性度 A_{1M} についても同様に算出する。

【0031】

第2中間層活性度算出部6は、第1中間層活性度算出部5が第1中間層の各ノードでの活性度 A_{1M} を算出すると、第2中間層の各ノードでの活性度 A_{2M} を算出する(図4の

10

20

30

40

50

ステップST3)。

第2中間層活性度算出部6による第2中間層の各ノードでの活性度 A_{2M} の算出方法は、第1中間層活性度算出部5による第1中間層の各ノードでの活性度 A_{1M} の算出方法と同様である。

即ち、第2中間層活性度算出部6は、第2中間層記憶部3に記憶されているインデックスを参照して、第2中間層のノード毎に、当該ノードに接続されている第1中間層の各ノードを確認して、その第1中間層の各ノードでの活性度 A_{1M} を取得する。

また、第2中間層活性度算出部6は、第2中間層記憶部3に記憶されているインデックスを参照して、第2中間層のノード毎に、当該ノードに接続されているエッジを確認して、第2中間層記憶部3からそのエッジの重み w を取得する。

10

【0032】

また、第2中間層活性度算出部6は、第2中間層のノード毎に、第2中間層記憶部3から当該ノードのバイアス値 B_{2M} を取得する。

第2中間層活性度算出部6は、第2中間層のノード毎に、第1中間層の各ノードでの活性度 A_{1M} 、エッジの重み w 、バイアス値 B_{2M} を取得すると、第1中間層活性度算出部5と同様の計算方法で、活性度 A_{1M} 、エッジの重み w 、バイアス値 B_{2M} を用いて、第2中間層のノード毎の活性度 A_{2M} を算出する。

【0033】

第3中間層活性度算出部7は、第2中間層活性度算出部6が第2中間層の各ノードでの活性度 A_{2M} を算出すると、第3中間層の各ノードでの活性度 A_{3M} を算出する(ステップST4)。

20

第3中間層活性度算出部7による第3中間層の各ノードでの活性度 A_{3M} の算出方法は、第1中間層活性度算出部5による第1中間層の各ノードでの活性度 A_{1M} の算出方法と同様である。

即ち、第3中間層活性度算出部7は、第3中間層記憶部4に記憶されているインデックスを参照して、第3中間層のノード毎に、当該ノードに接続されている第2中間層の各ノードを確認して、その第2中間層の各ノードでの活性度 A_{2M} を取得する。

また、第3中間層活性度算出部7は、第3中間層記憶部4に記憶されているインデックスを参照して、第3中間層のノード毎に、当該ノードに接続されているエッジを確認して、第3中間層記憶部4からそのエッジの重み w を取得する。

30

【0034】

また、第3中間層活性度算出部7は、第3中間層のノード毎に、第3中間層記憶部4から当該ノードのバイアス値 B_{3M} を取得する。

第3中間層活性度算出部7は、第3中間層のノード毎に、第2中間層の各ノードでの活性度 A_{2M} 、エッジの重み w 、バイアス値 B_{3M} を取得すると、第1中間層活性度算出部5と同様の計算方法で、活性度 A_{2M} 、エッジの重み w 、バイアス値 B_{3M} を用いて、第3中間層のノード毎の活性度 A_{3M} を算出する。

【0035】

出力層活性度算出部9は、第3中間層活性度算出部7が第3中間層の各ノードでの活性度 A_{3M} を算出すると、出力層の各ノードでの活性度 A_{OUT} を算出する(ステップST5)。

40

出力層活性度算出部9による出力層の各ノードでの活性度 A_{OUT} の算出方法は、第1中間層活性度算出部5による第1中間層の各ノードでの活性度 A_{1M} の算出方法と同様である。

即ち、出力層活性度算出部9は、出力層記憶部8に記憶されているインデックスを参照して、出力層のノード毎に、当該ノードに接続されている第3中間層の各ノードを確認して、その第3中間層の各ノードでの活性度 A_{3M} を取得する。

また、出力層活性度算出部9は、出力層記憶部8に記憶されているインデックスを参照して、出力層のノード毎に、当該ノードに接続されているエッジを確認して、出力層記憶部8からそのエッジの重み w を取得する。

50

【 0 0 3 6 】

また、出力層活性度算出部 9 は、出力層のノード毎に、出力層記憶部 8 から当該ノードのバイアス値 B_{OUT} を取得する。

出力層活性度算出部 9 は、出力層のノード毎に、第 3 中間層の各ノードでの活性度 A_{3M} 、エッジの重み w 、バイアス値 B_{OUT} を取得すると、第 1 中間層活性度算出部 5 と同様の計算方法で、活性度 A_{3M} 、エッジの重み w 、バイアス値 B_{OUT} を用いて、出力層のノード毎の活性度 A_{OUT} を算出する。

【 0 0 3 7 】

出力層活性度算出部 9 により算出された出力層のノード毎の活性度 A_{OUT} は、推論装置の推論結果として出力される。

例えば、画像に映っているものが人、犬、猫、自動車のいずれであるかを識別する場合、出力層は、4 つのノードから構成され、各ノードの活性度が、それぞれ人、犬、猫、自動車である可能性を示す値になるように学習される。

推論時は、出力層の中で一番活性度が大きいノードを選び、例えば、それが猫である可能性を出力するノードであれば、猫という推論結果を出力する。単なる識別結果のみでなく、活性度を用いた信頼度の算出や回帰予測値出力等の処理を実施してもよい。

【 0 0 3 8 】

以上で明らかのように、この実施の形態 1 によれば、第 1 中間層活性度算出部 5 が、第 1 中間層記憶部 2 に記憶されているインデックスを参照して、入力層活性度算出部 1 により算出された入力層の各ノードでの活性度と第 1 中間層記憶部 2 に記憶されている各エッジの重み及びバイアス値の中から、第 1 中間層の各ノードと接続関係がある入力層の各ノードでの活性度と各エッジの重みとバイアス値とを取得し、その取得した入力層の各ノードでの活性度と各エッジの重みとバイアス値とを用いて、第 1 中間層の各ノードでの活性度を算出するように構成したので、推論を行う際の計算量及びメモリ量を削減することができる効果を奏する。

即ち、第 1 中間層活性度算出部 5 では、第 1 中間層の各ノードと接続関係がある入力層の各ノードについてだけ計算を行えばよいため、推論を行う際の計算量及びメモリ量を大幅に削減することができる。

【 0 0 3 9 】

また、第 2 中間層活性度算出部 6 についても、第 2 中間層の各ノードと接続関係がある第 1 中間層の各ノードについてだけ計算を行えばよいため、第 1 中間層活性度算出部 5 と同様に、推論を行う際の計算量及びメモリ量を大幅に削減することができる。

また、第 3 中間層活性度算出部 7 についても、第 3 中間層の各ノードと接続関係がある第 2 中間層の各ノードについてだけ計算を行えばよいため、第 1 中間層活性度算出部 5 と同様に、推論を行う際の計算量及びメモリ量を大幅に削減することができる。

さらに、出力層活性度算出部 9 についても、出力層の各ノードと接続関係がある第 3 中間層の各ノードについてだけ計算を行えばよいため、第 1 中間層活性度算出部 5 と同様に、推論を行う際の計算量及びメモリ量を大幅に削減することができる。

【 0 0 4 0 】

この実施の形態 1 では、第 1 中間層活性度算出部 5、第 2 中間層活性度算出部 6、第 3 中間層活性度算出部 7 及び出力層活性度算出部 9 が、各ノードでの活性度を算出する際、前段の層の各ノードでの活性度とエッジの重みとの積和演算を実施するものを示したが、前段の層の各ノードでの活性度の最大値又は平均値を求め、その最大値又は平均値を式 (3) の加算結果 ADD の代わりに用いてもよい。

【 0 0 4 1 】

実施の形態 2 .

上記実施の形態 1 では、推論装置が適用するニューラルネットワークを構成している各層の各ノードが、前段又は後段の層の全てのノードとは接続されておらず、一部のノードと接続されているものを示している。

各層の各ノードが、前段又は後段の層の一部のノードと接続される場合でも、ノード間

10

20

30

40

50

の接続形態によっては、ノード間を接続する複数のエッジによってループが形成されることがある。

ここで、ニューラルネットワークにおいて、或るノードから同じエッジを一度しか通らないように辿って元のノードに戻る経路をループと称し、ループを構成するエッジの本数をループの長さと呼ぶ。

【0042】

図8は複数のエッジによって形成されるループの例を示す説明図である。

図8(a)(b)は4本のエッジによって形成されるループの例を示し、図8(c)は6本のエッジによって形成されるループの例を示し、図8(d)は8本のエッジによって形成されるループの例を示している。

例えば、階層型ニューラルネットワークでは、最短で長さ4のループが生じ得るが、特に、長さ4のループは学習時の逆誤差伝播法で伝播する勾配計算情報が容易に循環するため、推論精度低下の要因になる。また、ベイジアンネットの確率伝播法のように双方向に情報伝播することで推論するモデルでも、短いループの存在によって伝播情報が循環してしまっ

【0043】

そこで、この実施の形態2では、ニューラルネットワークを構成している各層のノード間を接続しているエッジ、即ち、第1中間層のノードと入力層のノードとを接続しているエッジ、第2中間層のノードと第1中間層のノードとを接続しているエッジ、第3中間層のノードと第2中間層のノードとを接続しているエッジ、出力層のノードと第3中間層のノードとを接続しているエッジがループを形成する場合、推論装置が適用するニューラルネットワークは、ループが6本以上のエッジで形成されているニューラルネットワークに限るものとする。

したがって、この実施の形態2では、図8(a)(b)に示すように4本のエッジによってループが形成されているニューラルネットワークは、推論装置の適用対象とならないが、図8(c)(d)に示すように6本又は8本のエッジによってループが形成されているニューラルネットワークは、推論装置の適用対象となる。

これにより、ニューラルネットワーク内に存在する長さの短いループに伴う推論精度の低下を抑制することができる効果を奏する。即ち、推論精度を維持したまま計算量及びメモリ量を削減することができる。

【0044】

実施の形態3

上記実施の形態1では、推論装置が適用するニューラルネットワークを構成している各層の各ノードが、前段又は後段の層の全てのノードとは接続されておらず、一部のノードと接続されているものを示している。

このとき、第1中間層の各ノードは、入力層における全ノードのうち、ランダムに選択された一部のノードと接続されているものであってもよい。

同様に、第2中間層の各ノードは、第1中間層における全ノードのうち、ランダムに選択された一部のノードと接続されているものであってもよく、また、第3中間層の各ノードは、第2中間層における全ノードのうち、ランダムに選択された一部のノードと接続されているものであってもよい。

また、同様に、出力層の各ノードは、第3中間層における全ノードのうち、ランダムに選択された一部のノードと接続されているものであってもよい。

【0045】

ランダムに選択される条件として、各層(出力層、第3中間層、第2中間層、第1中間層)の各ノード当りの前段の層(第3中間層、第2中間層、第1中間層、入力層)のノードとの平均接続本数が50本以下であるという条件を課してもよい。

あるいは、各層(出力層、第3中間層、第2中間層、第1中間層)の各ノード当りの前段の層(第3中間層、第2中間層、第1中間層、入力層)のノードとの平均接続本数が前段の層のノードの個数の10分の1以下であるという条件を課してもよい。

10

20

30

40

50

【 0 0 4 6 】

各層の各ノードが前段の層における全ノードと接続される形態では、各層のノード数が M で、前段の層のノード数が N である場合、各層での活性度の計算量とメモリ量が $N \times M$ のオーダーになるが、各層の各ノード当りの前段の層のノードの平均接続本数 n ($n < N$) が、50本以下であるという条件、あるいは、前段の層のノードの個数 N の10分の1以下であるという条件を課することで、長さの短いループの発生確率を低下させて、推論精度の低下を抑制することができるとともに、計算量及びメモリ量を削減することができる。

【 0 0 4 7 】

実施の形態 4 .

上記実施の形態 1 では、推論装置が適用するニューラルネットワークを構成している各層の各ノードが、前段又は後段の層の全てのノードとは接続されておらず、一部のノードと接続されているものを示している。

このとき、第 1 中間層の各ノードは、入力層における全ノードのうち、隣接していない一部のノードと接続されているようにしてもよい。

同様に、第 2 中間層の各ノードは、第 1 中間層における全ノードのうち、隣接していない一部のノードと接続されているようにしてもよく、また、第 3 中間層の各ノードは、第 2 中間層における全ノードのうち、隣接していない一部のノードと接続されているようにしてもよい。

また、同様に、出力層の各ノードは、第 3 中間層における全ノードのうち、隣接していない一部のノードと接続されているようにしてもよい。

例えば、第 1 中間層における“ N ”のノードの場合、入力層における“ 0 ”のノードと、“ 3 ”のノードとは隣接していないため、第 1 中間層における“ N ”のノードは、入力層における“ 0 ”、“ 3 ”のノードと接続される形態が許容されるが、入力層における“ 0 ”のノードと、“ 1 ”のノードとは隣接しているため、第 1 中間層における“ N ”のノードは、入力層における“ 0 ”、“ 1 ”のノードと接続される形態が許容されない。

【 0 0 4 8 】

この実施の形態 4 において、接続される形態が許容される条件として、各層（出力層、第 3 中間層、第 2 中間層、第 1 中間層）の各ノード当りの前段の層（第 3 中間層、第 2 中間層、第 1 中間層、入力層）のノードとの平均接続本数が 50 本以下であるという条件を課してもよい。

あるいは、各層（出力層、第 3 中間層、第 2 中間層、第 1 中間層）の各ノード当りの前段の層（第 3 中間層、第 2 中間層、第 1 中間層、入力層）のノードとの平均接続本数が前段の層のノードの個数の10分の1以下であるという条件を課してもよい。

上記の条件を課することで、長さの短いループの発生確率を低下させて、推論精度の低下を抑制することができるとともに、計算量及びメモリ量を削減することができる。

【 0 0 4 9 】

実施の形態 5 .

上記実施の形態 1 ~ 4 では、推論装置が適用するニューラルネットワークとして、3つの中間層をもつ階層型フィードフォワードのニューラルネットワークを例に挙げているが、中間層は3つより多くても少なくともよい。また、中間層がないロジスティック回帰モデルのような構造であってもよい。

また、層間のノードを全結合する層や、畳み込みニューラルネットワークのような畳み込み層およびプーリング層や、リカレントニューラルネットワークにおける LSTM（長期短期記憶）ブロックなど、ニューラルネットワークの従来手法と組み合わせてもよい。

ここで、畳み込みニューラルネットワークは、畳み込み層とプーリング層が繰り返された構造になっている。例えば、畳み込み層は、画像の局所的な特徴抽出を担う層であり、プーリング層は、局所毎に特徴をまとめあげる層である。

【 0 0 5 0 】

上記実施の形態 1 ~ 4 では、推論装置が適用するニューラルネットワークとして、階層

10

20

30

40

50

型フィードフォワードのニューラルネットワークを例に挙げているが、層を飛ばした接続があってもよいし、同じ層に属するノード同士で接続があってもよいし、接続先と接続元が同一の自己接続があってもよいし、エッジがループを形成するような循環接続があってもよい(リカレントニューラルネットワーク)。

また、自己組織化マップ(SOM)、連想記憶モデル、ホップフィールドネットワーク、ボルツマンマシンなど、他のグラフを用いて推論するニューラルネットワークであってもよい。さらに、ニューラルネットワークに限らず、ベイジアンネットワークなど、他のグラフを用いて推論するモデルでもよい。

【0051】

上記実施の形態1~4では、入力層のノードが $0, 1, \dots, N-1$ 、第1中間層のノードが $N, N+1, \dots, N+M-1$ といったように一次元のインデックスを付けているが、入力層のノードが $(0, 0), (0, 1), \dots, (0, N-1)$ 、第1中間層のノードが $(1, 0), (1, 1), \dots, (1, M-1)$ といったように二次元のインデックスを付けてもよいし、メモリのアドレスをインデックスとして用いてもよいし、他のインデックスを付けてもよい。

10

【0052】

上記実施の形態1~4では、推論装置が適用するニューラルネットワークとして、エッジの数とエッジ重みの数が一致する例を挙げているが、畳み込みネットワークにおける畳み込みフィルタ係数のように、複数のエッジ重みを共有化するようにしてもよい。

上記実施の形態1~4では、各ノードにおける活性度の計算過程を順に記載しているが、互いに依存しない計算を複数のCPUやGPUを用いて並列化し、さらに高速化することも可能である。

20

【0053】

上記実施の形態1~4では、画像データを入力して、画像を分類する画像分類システムを例に挙げているが、データと対応する教師信号が準備できており、教師あり学習を行うことができるならば、データの入力に対して何らかの推論結果を出力する推論システム全般に適用可能である。

例えば、画像を入力して検知したい物体領域の位置や大きさを出力してもよいし、画像を入力して、その画像を説明するテキストを出力してもよいし、ノイズが入った画像を入力して、ノイズを除去した画像を出力してもよいし、画像とテキストを入力して、画像をテキストに従って変換してもよい。

30

また、音声を入力して音素や単語を出力してもよいし、音声を入力して次に発話される単語を予測してもよいし、音声を入力して、それに対する適切な応答音声も出力してもよいし、テキストを入力して別の言語のテキストを出力してもよいし、時系列を入力して将来の時系列を予測してもよいし、時系列を入力して時系列の状態を推定してもよい。

【0054】

上記実施の形態1~4では、データと対応する教師信号を用いる教師あり学習によって学習したモデルで推論するシステム例を挙げているが、教師信号のないデータを用いる教師なし学習や半教師あり学習によって学習したモデルで推論するシステムでもよい。

上記実施の形態1~4では、推論装置が、図示せぬデータ入力装置から画像データを与えられて、第1中間層の各ノードでの活性度を算出する例を示したが、図示せぬデータ入力装置が、第1中間層の各ノードでの活性度を算出し、推論装置が、第2及び第3中間層及び出力層の各ノードでの活性度を算出するようにしてもよい。データ入力装置の出力の次元数が入力の次元数よりも少ない場合、データ入力装置がデータ圧縮の機能も併せ持つことになる。

40

上記実施の形態1~4では、各ノードに対して一度だけ活性度を算出する例を挙げているが、ベイジアンネットワークの確率伝播法のようにノード間で繰り返し何度も情報交換して推論精度を向上させてもよい。

【0055】

実施の形態6 .

50

上記実施の形態 1 ~ 4 では、推論装置が適用するニューラルネットワークとして、入力層を除く全ての層で枝接続のインデックスを保持している例を挙げているが、一部の層だけが枝接続のインデックスを保持し、他の層では通常のニューラルネットワークと同様の枝接続であるものでもあってもよい。

ここで、枝接続のインデックスとは、図 7 に示すようなインデックスであり、エッジ重みやバイアス値を含む概念である。

また、通常のニューラルネットワークと同様の枝接続とは、接続先の層における全てのノードと接続されてる枝接続（全接続層の枝接続）を意味するほか、接続先の層における或るノード及び当該ノードの周辺ノードと接続されている畳み込み層やプーリング層などの公知のニューラルネットワークの枝接続を意味する。

【 0 0 5 6 】

図 9 は入力層に与えられる画像データを 10 種のクラスに識別するニューラルネットワークの一例を示す説明図である。

図 9 は例では、入力層と出力層の間に 5 つの中間層、即ち、第 1 中間層、第 2 中間層、第 3 中間層、第 4 中間層及び第 5 中間層が接続されている。

また、図 9 は例では、入力層から第 1 中間層が畳み込み層 3 1、第 1 中間層から第 2 中間層がプーリング層 3 2、第 2 中間層から第 3 中間層が畳み込み層 3 3、第 3 中間層から第 4 中間層がプーリング層 3 4、第 4 中間層から第 5 中間層が上記実施の形態 1 ~ 4 で示しているインデックスを保持する層 3 5、第 5 中間層から出力層が全接続層 3 6 である。

このため、第 5 中間層における各ノードは、図 7 に示している第 1 中間層と同様に、第 4 中間層における接続元のノードを示すインデックスと、その接続に対応するエッジ重み及びバイアス値を保持している。

【 0 0 5 7 】

例えば、入力層に与えられる画像データが、縦 60 × 横 60 画素の画像データであれば、図 9 のニューラルネットワークでは、3600 個（= 60 × 60 × 1 個）のノードを有する入力層が必要となる。

このとき、例えば、入力層から第 1 中間層への畳み込み層 3 1 のフィルタサイズが 5 × 5 × 1、この畳み込み層 3 1 でのマップ数が 100、第 1 中間層から第 2 中間層へのプーリング層 3 2 及び第 3 中間層から第 4 中間層へのプーリング層 3 4 がフィルタサイズ 2 × 2 × 1 の最大値プーリングである場合、第 1 中間層のサイズが 56 × 56 × 100（= (60 - 5 + 1) × (60 - 5 + 1) × 100）、第 2 中間層のサイズが 28 × 28 × 100（= (56 / 2) × (56 / 2) × 100）となる。

また、第 3 中間層のサイズが 24 × 24 × 200（= (28 - 5 + 1) × (28 - 5 + 1) × 200）、第 4 中間層のサイズが 12 × 12 × 200（= (24 / 2) × (24 / 2) × 200）、第 5 中間層のサイズが 1 × 1 × 1000、出力層のノード数が 1 × 1 × 10 となる。

【 0 0 5 8 】

なお、入力層から第 1 中間層に情報を伝播する際に伝播値を計算する活性化関数、第 2 中間層から第 3 中間層に情報を伝播する際に伝播値を計算する活性化関数や、第 4 中間層から第 5 中間層に情報を伝播する際に伝播値を計算する活性化関数として、例えば、ReLU (Rectified Linear Unit) が用いられ、第 5 中間層から出力層に情報を伝播する際に伝播値を計算する活性化関数として、例えば、正規化指数関数であるソフトマックス関数 (Softmax 関数) が用いられる。

【 0 0 5 9 】

図 9 のニューラルネットワークでは、入力層から第 4 中間層までの畳み込み層 3 1, 3 3 とプーリング層 3 2, 3 4 によって、入力された画像の位置変化に対してロバストに画像データの特徴量を抽出することができる。

また、第 4 中間層から第 5 中間層へのインデックスを保持する層によって、上記実施の形態 1 ~ 4 と同様に、推論を行う際の計算量及びメモリ量を大幅に削減することができる。

10

20

30

40

50

【 0 0 6 0 】

この実施の形態 6 では、画像データが入力層に与えられる例を示しているが、入力層に与えられるデータは画像データに限るものではなく、例えば、センサにより観測されたデータであるセンサ信号、音声やテキストなどのデータなどであってもよい。

また、この実施の形態 6 では、入力層に与えられる画像データを 10 種のクラスに識別する例を示しているが、ニューラルネットワークを構成している出力層を変更することで、画像データのクラスを識別する推論以外の推論を行うようにしてもよい。

例えば、画像データのノイズを取り除くデノイジング、回帰予測や尤度算出などの推論を行うようにしてもよい。

また、推論の目的に合わせて、各層のノード数やフィルタサイズを変えてもよい。

10

【 0 0 6 1 】

図 9 のニューラルネットワークでは、畳み込み層 3 1、プーリング層 3 2、畳み込み層 3 3、プーリング層 3 4、インデックスを保持する層 3 5、全接続層 3 6 の順序で枝接続している例を示しているが、上記実施の形態 1 ~ 4 に示すインデックスを保持する層が 1 つでも接続されていればよく、上記の順序以外の順序で枝接続しているものであってもよい。また、プーリング層 3 2、3 4 が接続されていないものであってもよい。

【 0 0 6 2 】

実施の形態 7 .

実施の形態 1 ~ 6 では、推論装置が適用しているニューラルネットワークが有向ループをなしていないフィードフォワードニューラルネットワーク (F F N N) である例を挙げているが、推論装置が、ネットワークの一部が有向ループをなしているリカレントニューラルネットワーク (R N N) を適用し、そのリカレントニューラルネットワークの一部の層が、上記実施の形態 1 ~ 4 に示すインデックスを保持する層であるものであってもよい。

20

【 0 0 6 3 】

図 1 0 は入力層から第 1 中間層、第 1 中間層から第 2 中間層 (文脈層)、第 2 中間層 (文脈層) から第 1 中間層、第 1 中間層から出力層に枝接続を持っているリカレントニューラルネットワーク (R N N) であるエルマンネットワーク (E l m a n N e t w o r k) を示す説明図である。

図 1 0 のエルマンネットワークでは、第 2 中間層 (文脈層) のノード数が、第 1 中間層のノード数と等しいものとする。

30

【 0 0 6 4 】

図 1 1 はこの発明の実施の形態 7 による推論装置を示す構成図であり、図 1 1 において、図 1 と同一符号は同一または相当部分を示すので説明を省略する。

中間層活性度算出部である第 2 中間層活性度算出部 4 1 は例えば CPU を実装している半導体集積回路、あるいは、ワンチップマイコンなどで構成されている中間層活性度算出回路 1 3 で実現されるものであり、第 1 中間層の各ノードの活性度を第 2 中間層 (文脈層) の各ノードの活性度としてコピーする処理を実施する。

この実施の形態 7 では、第 2 中間層活性度算出部 4 1 が、第 1 中間層の各ノードの活性度を第 2 中間層 (文脈層) の各ノードの活性度としてコピーする例を想定しているが、これは一例に過ぎず、例えば、入力層活性度算出部 1 と同様に、式 (1) によって第 2 中間層 (文脈層) の各ノードの活性度を算出するようにしてもよい。

40

【 0 0 6 5 】

中間層記憶部である第 1 中間層記憶部 4 2 は例えば RAM やハードディスクなどの記憶媒体からなる中間層記憶装置 1 2 で実現されるものであり、第 1 中間層の各ノードと第 2 中間層 (文脈層) の各ノードとの接続関係を示すインデックスと、第 1 中間層のノードと第 2 中間層のノードとを接続している各エッジの重みと、第 1 中間層の各ノードに与えられているバイアス値とを記憶している。

中間層活性度算出部である第 1 中間層活性度算出部 4 3 は例えば CPU を実装している半導体集積回路、あるいは、ワンチップマイコンなどで構成されている中間層活性度算出

50

回路 13 で実現されるものであり、第 1 中間層記憶部 42 に記憶されているインデックスを参照して、第 2 中間層活性度算出部 41 により求められた第 2 中間層（文脈層）の各ノードでの活性度と第 1 中間層記憶部 42 に記憶されている各エッジの重み及びバイアス値の中から、第 1 中間層の各ノードと接続関係がある第 2 中間層（文脈層）の各ノードでの活性度と各エッジの重みとバイアス値とを取得し、その取得した第 2 中間層（文脈層）の各ノードでの活性度と各エッジの重みとバイアス値とを用いて、第 1 中間層の各ノードでの活性度を算出する処理を実施する。

【 0066 】

図 11 では、推論装置の構成要素である入力層活性度算出部 1、第 1 中間層記憶部 2、第 1 中間層活性度算出部 5、第 2 中間層活性度算出部 41、第 1 中間層記憶部 42、第 1 中間層活性度算出部 43、出力層記憶部 8 及び出力層活性度算出部 9 のそれぞれが専用のハードウェアで構成されているものを想定しているが、推論装置がコンピュータで構成されていてもよい。

10

推論装置がコンピュータで構成される場合、第 1 中間層記憶部 2、第 1 中間層記憶部 42 及び出力層記憶部 8 を図 3 に示すコンピュータのメモリ 21 上に構成するとともに、入力層活性度算出部 1、第 1 中間層活性度算出部 5、第 2 中間層活性度算出部 41、第 1 中間層活性度算出部 43 及び出力層活性度算出部 9 の処理内容を記述しているプログラムを図 3 に示すコンピュータのメモリ 21 に格納し、当該コンピュータのプロセッサ 22 がメモリ 21 に格納されているプログラムを実行するようにすればよい。

【 0067 】

20

次に動作について説明する。

ただし、第 2 中間層活性度算出部 41、第 1 中間層記憶部 42 及び第 1 中間層活性度算出部 43 以外は、上記実施の形態 1 と同様であるため、ここでは、第 2 中間層活性度算出部 41、第 1 中間層記憶部 42 及び第 1 中間層活性度算出部 43 についてのみ説明する。

第 2 中間層活性度算出部 41 は、第 1 中間層活性度算出部 5 が上記実施の形態 1 と同様に、入力層活性度算出部 1 により算出された入力層の各ノードの活性度を用いて、第 1 中間層の各ノードの活性度 A_{1M} を算出すると、第 1 中間層の各ノードの活性度 A_{1M} を第 2 中間層（文脈層）の各ノードの活性度 A_{2M} としてコピーする。

これにより、ある時刻 t における第 2 中間層の各ノードの活性度 A_{2M} は、時刻 t における第 1 中間層の各ノードの活性度 A_{1M} と同じになる。

30

【 0068 】

第 1 中間層活性度算出部 43 は、第 2 中間層活性度算出部 41 が第 2 中間層（文脈層）の各ノードの活性度 A_{2M} を求めると、第 1 中間層の各ノードの活性度 A'_{1M} を算出する。

第 1 中間層活性度算出部 43 による第 1 中間層の各ノードでの活性度 A'_{1M} の算出方法は、第 1 中間層活性度算出部 5 による第 1 中間層の各ノードでの活性度 A_{1M} の算出方法と同様である。

即ち、第 1 中間層活性度算出部 43 は、第 1 中間層記憶部 42 に記憶されているインデックスを参照して、第 1 中間層のノード毎に、当該ノードに接続されている第 2 中間層（文脈層）の各ノードを確認して、第 2 中間層（文脈層）の各ノードでの活性度 A_{2M} を取得する。

40

また、第 1 中間層活性度算出部 43 は、第 1 中間層記憶部 42 に記憶されているインデックスを参照して、第 1 中間層のノード毎に、当該ノードに接続されているエッジ（第 2 中間層のノードと接続されているエッジ）を確認して、第 1 中間層記憶部 42 からそのエッジの重み w を取得する。

【 0069 】

また、第 1 中間層活性度算出部 43 は、第 1 中間層のノード毎に、第 1 中間層記憶部 42 から当該ノードのバイアス値 B_{1M} を取得する。

第 1 中間層活性度算出部 43 は、第 1 中間層のノード毎に、第 2 中間層（文脈層）の各ノードでの活性度 A_{2M} 、エッジの重み w 、バイアス値 B_{1M} を取得すると、第 1 中間層

50

活性度算出部 5 と同様の計算方法で、活性度 A_{2M} 、エッジの重み w 、バイアス値 B_{1M} を用いて、第 1 中間層のノード毎の活性度 A'_{1M} を算出する。

出力層活性度算出部 9 は、第 1 中間層活性度算出部 4 3 が第 1 中間層の各ノードでの活性度 A'_{1M} を算出すると、第 1 中間層の各ノードでの活性度 A'_{1M} を用いて、出力層の各ノードでの活性度 A_{OUT} を算出する。

出力層活性度算出部 9 による出力層の各ノードでの活性度 A_{OUT} の算出方法は上記実施の形態 1 と同様である。

【0070】

以上で明らかのように、この実施の形態 7 によれば、第 1 中間層活性度算出部 4 3 が、第 1 中間層記憶部 4 2 に記憶されているインデックスを参照して、第 2 中間層活性度算出部 4 1 により求められた第 2 中間層（文脈層）の各ノードでの活性度と第 1 中間層記憶部 4 2 に記憶されている各エッジの重み及びバイアス値の中から、第 1 中間層の各ノードと接続関係がある第 2 中間層（文脈層）の各ノードでの活性度と各エッジの重みとバイアス値とを取得し、その取得した第 2 中間層（文脈層）の各ノードでの活性度と各エッジの重みとバイアス値とを用いて、第 1 中間層の各ノードでの活性度を算出するように構成したので、ネットワークの一部が有向ループをなしているリカレントニューラルネットワーク（RNN）を適用する場合であっても、推論を行う際の計算量及びメモリ量を削減することができる効果を奏する。

即ち、第 1 中間層活性度算出部 4 3 では、第 1 中間層の各ノードと接続関係がある第 2 中間層（文脈層）の各ノードについてだけ計算を行えばよいため、推論を行う際の計算量及びメモリ量を大幅に削減することができる。

【0071】

この実施の形態 7 では、推論装置が適用するリカレントニューラルネットワーク（RNN）がエルマンネットワークである例に挙げているが、これは一例に過ぎず、例えば、Jordan Network や、LSTM（Long Short Term Memory）ブロックを持つリカレントニューラルネットワーク、階層型のリカレントニューラルネットワーク、双方向のリカレントニューラルネットワーク、連続時間のリカレントニューラルネットワークなどを適用するものであってもよい。

ただし、この実施の形態 7 では、いずれのリカレントニューラルネットワークを適用する場合でも、一部の層で枝接続のインデックスを保持するものとする。

【0072】

実施の形態 8 .

上記実施の形態 1 ~ 7 では、推論装置が適用するニューラルネットワークとして、同じ層のノード同士の枝接続や自己接続が無いフィードフォワードニューラルネットワーク（FFNN）やリカレントニューラルネットワーク（RNN）である例に挙げているが、同じ層のノード同士の枝接続や、接続元ノードと接続先ノードが同一のノードである自己接続が有るフィードフォワードニューラルネットワーク（FFNN）やリカレントニューラルネットワーク（RNN）を適用するようにしてもよい。また、層を飛ばした接続があるフィードフォワードニューラルネットワーク（FFNN）やリカレントニューラルネットワーク（RNN）を適用するようにしてもよい。

【0073】

図 1 2 は中間層内でノード同士の枝接続や自己接続が存在するとともに、入力層から中間層を飛ばして出力層へ接続する枝が存在するニューラルネットワークである Echo State Network の例を示す説明図である。

図 1 2 において、中間層のノード同士の枝接続や自己接続は、中間層から中間層への枝接続とみなすことができるため、Echo State Network は、層単位では図 1 3 のように表すことができる。

【0074】

図 1 4 はこの発明の実施の形態 8 による推論装置を示す構成図であり、図 1 4 において、図 1 と同一符号は同一または相当部分を示すので説明を省略する。

中間層記憶部 5 1 は例えば R A M やハードディスクなどの記憶媒体からなる中間層記憶装置 1 2 で実現されるものであり、中間層の各ノードと入力層又は出力層の各ノードとの接続関係を示すインデックスと、中間層のノードと入力層又は出力層のノードとを接続している各エッジの重みと、中間層の各ノードに与えられているバイアス値とを記憶している。

また、中間層記憶部 5 1 は中間層のノード同士の枝接続や自己接続の関係を示すインデックスと、中間層のノード同士の枝接続や自己接続している各エッジの重みと、中間層の各ノードに与えられているバイアス値とを記憶している。

【 0 0 7 5 】

中間層活性度算出部 5 2 は例えば C P U を実装している半導体集積回路、あるいは、ワンチップマイコンなどで構成されている中間層活性度算出回路 1 3 で実現されるものであり、図 1 の第 1 中間層活性度算出部 5 と同様に、中間層記憶部 5 1 に記憶されているインデックスを参照して、入力層活性度算出部 1 により算出された入力層の各ノードでの活性度又は出力層活性度算出部 5 4 により算出された出力層の各ノードでの活性度と中間層記憶部 5 1 に記憶されている各エッジの重み及びバイアス値の中から、中間層の各ノードと接続関係がある入力層又は出力層の各ノードでの活性度と各エッジの重みとバイアス値とを取得し、その取得した入力層又は出力層の各ノードでの活性度と各エッジの重みとバイアス値とを用いて、中間層の各ノードでの活性度を算出する処理を実施する。

また、中間層活性度算出部 5 2 は中間層記憶部 5 1 に記憶されているインデックスを参照して、中間層における算出済みの各ノードでの活性度と中間層記憶部 5 1 に記憶されている各エッジの重み及びバイアス値の中から、中間層における接続先の各ノードと接続関係がある中間層における接続元の各ノードでの活性度と各エッジの重みとバイアス値とを取得し、その取得した中間層における接続元の各ノードでの活性度と各エッジの重みとバイアス値とを用いて、中間層における接続先の各ノードでの活性度を算出する処理を実施する。

ここで、中間層における接続先のノードとは、中間層における他のノードと接続されているノード、あるいは、中間層において、自己のノードと接続されているノードを意味する。

また、中間層における接続元のノードとは、接続先のノードと接続されている中間層における他のノード、あるいは、中間層における自己接続のノードを意味する。

【 0 0 7 6 】

出力層記憶部 5 3 は例えば R A M やハードディスクなどの記憶媒体からなる出力層記憶装置 1 4 で実現されるものであり、出力層の各ノードと入力層又は中間層の各ノードとの接続関係を示すインデックス（接続情報）を記憶している。

また、出力層記憶部 5 3 は出力層のノードが入力層のノードと接続されていれば、出力層のノードと入力層のノードとを接続している各エッジの重みと、入力層のノードと接続されている出力層のノードに与えられているバイアス値とを記憶している。

また、出力層記憶部 5 3 は出力層のノードが中間層のノードと接続されていれば、出力層のノードと中間層のノードとを接続している各エッジの重みと、中間層のノードと接続されている出力層のノードに与えられているバイアス値とを記憶している。

【 0 0 7 7 】

出力層活性度算出部 5 4 は例えば C P U を実装している半導体集積回路、あるいは、ワンチップマイコンなどで構成されている出力層活性度算出回路 1 5 で実現されるものであり、出力層のノードと接続されているノードが入力層のノードであれば、入力層活性度算出部 1 により算出された入力層の各ノードでの活性度と出力層記憶部 5 3 に記憶されている各エッジの重み及びバイアス値の中から、出力層の各ノードと接続関係がある入力層の各ノードでの活性度と各エッジの重みとバイアス値とを取得し、出力層のノードと接続されているノードが中間層のノードであれば、中間層活性度算出部 5 2 により算出された中間層のノードでの活性度と出力層記憶部 5 3 に記憶されている各エッジの重み及びバイアス値の中から、出力層の各ノードと接続関係がある中間層のノードでの活性度と各エッジ

10

20

30

40

50

の重みとバイアス値とを取得する処理を実施する。

また、出力層活性度算出部 5 4 は取得した入力層又は中間層のノードでの活性度と各エッジの重みとバイアス値とを用いて、出力層の各ノードでの活性度を算出する処理を実施する。

【 0 0 7 8 】

図 1 4 では、推論装置の構成要素である入力層活性度算出部 1、中間層記憶部 5 1、中間層活性度算出部 5 2、出力層記憶部 5 3 及び出力層活性度算出部 5 4 のそれぞれが専用のハードウェアで構成されているものを想定しているが、推論装置がコンピュータで構成されていてもよい。

推論装置がコンピュータで構成される場合、中間層記憶部 5 1 及び出力層記憶部 5 3 を図 3 に示すコンピュータのメモリ 2 1 上に構成するとともに、入力層活性度算出部 1、中間層活性度算出部 5 2 及び出力層活性度算出部 5 4 の処理内容を記述しているプログラムを図 3 に示すコンピュータのメモリ 2 1 に格納し、当該コンピュータのプロセッサ 2 2 がメモリ 2 1 に格納されているプログラムを実行するようにすればよい。

【 0 0 7 9 】

次に動作について説明する。

中間層活性度算出部 5 2 は、入力層活性度算出部 1 が上記実施の形態 1 と同様に入力層の各ノードの活性度を算出すると、中間層記憶部 5 1 に記憶されているインデックスを参照して、中間層の各ノードのうち、入力層のノードと接続関係があるノードを確認するとともに、出力層のノードと接続関係があるノードを確認する。

中間層活性度算出部 5 2 は、中間層の各ノードのうち、入力層のノードと接続関係があるノードの活性度については、図 1 の第 1 中間層活性度算出部 5 と同様に、入力層の各ノードの活性度を用いて算出する。

また、中間層の各ノードのうち、出力層のノードと接続関係があるノードの活性度については、出力層活性度算出部 5 4 により算出された出力層の活性度を用いて算出する。

中間層活性度算出部 5 2 による出力層のノードと接続関係があるノードの活性度の算出方法は、入力層のノードと接続関係があるノードの活性度の算出方法と同様であり、活性度算出対象のノードが接続されているノードが、入力層のノードではなく、出力層のノードである点だけが相違している。

【 0 0 8 0 】

中間層活性度算出部 5 2 は、入力層又は出力層のノードと接続関係がある中間層のノードの活性度を算出すると、中間層記憶部 5 1 に記憶されているインデックスを参照して、中間層の各ノードのうち、接続先のノード（中間層における他のノードと接続されているノード、あるいは、中間層において、自己のノードと接続されているノード）と接続関係がある接続元のノードを確認する。

中間層活性度算出部 5 2 は、接続先のノードと接続関係がある接続元のノードを確認すると、中間層における算出済みの各ノードでの活性度と中間層記憶部 5 1 に記憶されている各エッジの重み及びバイアス値の中から、中間層における接続先の各ノードと接続関係がある中間層における接続元の各ノードでの活性度と各エッジの重みとバイアス値とを取得する。

中間層における接続元のノードが、入力層又は出力層のノードと接続関係がある中間層のノードであれば、先に説明したように既に算出済みである。このため、入力層又は出力層のノードと接続関係がある中間層のノードに近いノードから順番に、活性度算出対象のノード（接続先のノード）とすればよい。

【 0 0 8 1 】

中間層活性度算出部 5 2 は、その取得した中間層における接続元の各ノードでの活性度と各エッジの重みとバイアス値とを用いて、中間層における接続先の各ノードでの活性度を算出する。

中間層活性度算出部 5 2 による中間層における接続先のノードでの活性度の算出方法は、入力層のノードと接続関係があるノードの活性度の算出方法と同様であり、接続元のノ

10

20

30

40

50

ードが、入力層のノードではなく、中間層のノードである点だけが相違している。

【 0 0 8 2 】

出力層活性度算出部 5 4 は、出力層記憶部 5 3 に記憶されているインデックスを参照して、出力層の各ノードが接続されている入力層又は中間層のノードを確認する。

出力層活性度算出部 5 4 は、出力層のノードと接続されているノードが入力層のノードであれば、入力層活性度算出部 1 により算出された入力層の各ノードでの活性度と出力層記憶部 5 3 に記憶されている各エッジの重み及びバイアス値の中から、出力層の各ノードと接続関係がある入力層の各ノードでの活性度と各エッジの重みとバイアス値とを取得する。一方、出力層のノードと接続されているノードが中間層のノードであれば、中間層活性度算出部 5 2 により算出された中間層のノードでの活性度と出力層記憶部 5 3 に記憶されている各エッジの重み及びバイアス値の中から、出力層の各ノードと接続関係がある中間層のノードでの活性度と各エッジの重みとバイアス値とを取得する。

出力層活性度算出部 5 4 は、入力層又は中間層のノードでの活性度と各エッジの重みとバイアス値とを取得すると、その取得した入力層又は中間層のノードでの活性度と各エッジの重みとバイアス値とを用いて、出力層の各ノードでの活性度を算出する。

【 0 0 8 3 】

以上で明らかなように、この実施の形態 8 によれば、中間層活性度算出部 5 2 が、中間層記憶部 5 1 に記憶されているインデックスを参照して、中間層における算出済みの各ノードでの活性度と中間層記憶部 5 1 に記憶されている各エッジの重み及びバイアス値の中から、中間層における接続先の各ノードと接続関係がある中間層における接続元の各ノードでの活性度と各エッジの重みとバイアス値とを取得し、その取得した中間層における接続元の各ノードでの活性度と各エッジの重みとバイアス値とを用いて、中間層における接続先の各ノードでの活性度を算出するように構成したので、中間層内でノード同士の枝接続や自己接続が存在するニューラルネットワークを適用する場合であっても、推論を行う際の計算量及びメモリ量を削減することができる効果を奏する。

【 0 0 8 4 】

また、この実施の形態 8 によれば、出力層活性度算出部 5 4 が、出力層のノードと接続されているノードが入力層のノードであれば、入力層活性度算出部 1 により算出された入力層の各ノードでの活性度と出力層記憶部 5 3 に記憶されている各エッジの重み及びバイアス値の中から、出力層の各ノードと接続関係がある入力層の各ノードでの活性度と各エッジの重みとバイアス値とを取得し、その取得した入力層のノードでの活性度と各エッジの重みとバイアス値とを用いて、出力層の各ノードでの活性度を算出するように構成したので、入力層から中間層を飛ばして出力層へ接続する枝が存在するニューラルネットワークを適用する場合であっても、推論を行う際の計算量及びメモリ量を削減することができる効果を奏する。

【 0 0 8 5 】

この実施の形態 8 では、推論装置が適用するニューラルネットワークが Echo State Network である例に挙げているが、これは一例に過ぎず、例えば、全結合のリカレントニューラルネットワーク、ホップフィールドネットワーク、ボルツマンマシンなどを適用するものであってもよい。

ただし、この実施の形態 8 では、いずれのニューラルネットワークを適用する場合でも、一部の層で枝接続のインデックスを保持するものとする。

【 0 0 8 6 】

実施の形態 9 .

上記実施の形態 1 ~ 6 では、推論装置が適用しているニューラルネットワークがフィードフォワードニューラルネットワーク (F F N N) である例を挙げ、上記実施の形態 7 , 8 では、推論装置が適用しているニューラルネットワークがリカレントニューラルネットワーク (R N N) である例を挙げているが、これは一例に過ぎず、推論装置が下記に示すようなニューラルネットワークを適用するものであってもよい。ただし、いずれのニューラルネットワークを適用する場合でも、一部の層で枝接続のインデックスを保持するもの

10

20

30

40

50

とする。

【0087】

例えば、推論装置が適用するニューラルネットワークとして、放射基底関数(RBF)ネットワーク、自己組織化マップ(SOM)、学習ベクトル量子化法(LVQ)、モジュールニューラルネットワーク、スパイクニューラルネットワーク、動的ニューラルネットワーク、カスケードニューラルネットワーク、階層型時間記憶(HTM)等のFFNNやRNN以外のニューラルネットワークが考えられる。

【0088】

実施の形態10.

上記実施の形態1~9では、推論装置の学習手法が、データと対応する教師信号を用いる教師あり学習、教師信号のないデータを用いる教師なし学習、あるいは、半教師あり学習である例を挙げているが、推論装置の学習手法が、強化学習であってもよい。

ここで、強化学習とは、ある環境下におけるエージェントが、現在の状態を観測し、取るべき行動を決定するためのモデルを学習する手法である。エージェントは、コンピュータのユーザが連続した操作をしなくても、自律的に情報収集や状況判断を行って適切な処理動作を実行する機能を意味する。

エージェントが行動を選択すると、環境から報酬を得るが、強化学習では、一連の行動で報酬を最大化できるようなポリシーを学習する。

【0089】

強化学習では、現在の状態、あるいは、行動がどのくらい良いかを計る指標として、状態 s の価値を表す状態価値関数 $V(s)$ や、状態 s のときに行動 a を選択することで環境から得られる報酬を表す行動価値関数 $Q(s, a)$ を用いる。強化学習のアルゴリズムとしては、 $Sarsa$ や Q 学習($Q-learning$)等のTD(Temporal Difference: 時間差分)学習が用いられる。

推論装置の学習手法が強化学習である場合、状態 s を入力として、状態価値関数 $V(s)$ や行動価値関数 $Q(s, a)$ を出力するニューラルネットワークを学習し、これらを用いてTD学習することになる。即ち、一部の層で枝接続のインデックスを保持するようなニューラルネットワークを用いて、状態価値関数 $V(s)$ や行動価値関数 $Q(s, a)$ を計算して、強化学習を行う。

【0090】

なお、本願発明はその発明の範囲内において、各実施の形態の自由な組み合わせ、あるいは各実施の形態の任意の構成要素の変形、もしくは各実施の形態において任意の構成要素の省略が可能である。

【産業上の利用可能性】

【0091】

この発明に係る推論装置は、推論を行う際の計算量やメモリ量を削減する必要が高いものに適している。

【符号の説明】

【0092】

1 入力層活性度算出部、2 第1中間層記憶部(中間層記憶部)、3 第2中間層記憶部(中間層記憶部)、4 第3中間層記憶部(中間層記憶部)、5 第1中間層活性度算出部(中間層活性度算出部)、6 第2中間層活性度算出部(中間層活性度算出部)、7 第3中間層活性度算出部(中間層活性度算出部)、8 出力層記憶部、9 出力層活性度算出部、11 入力層活性度算出回路、12 中間層記憶装置、13 中間層活性度算出回路、14 出力層記憶装置、15 出力層活性度算出回路、21 メモリ、22 プロセッサ、31 畳み込み層、32 プーリング層、33 畳み込み層、34 プーリング層、35 インデックスを保持する層、36 全接続層、41 第2中間層活性度算出部(中間層活性度算出部)、42 第1中間層記憶部(中間層記憶部)、43 第1中間層活性度算出部(中間層活性度算出部)、51 中間層記憶部、52 中間層活性度算出部、53 出力層記憶部、54 出力層活性度算出部。

10

20

30

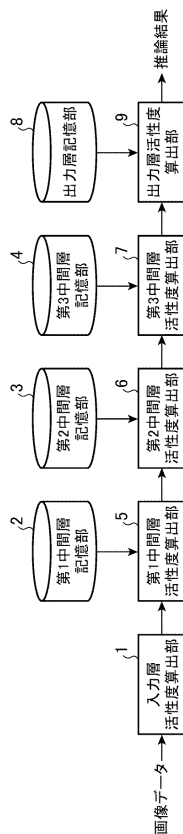
40

50

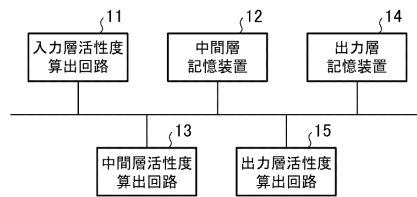
【要約】

第1中間層活性度算出部(5)が、第1中間層記憶部(2)に記憶されているインデックスを参照して、入力層活性度算出部(1)により算出された入力層の各ノードでの活性度と第1中間層記憶部(2)に記憶されている各エッジの重み及びバイアス値の中から、第1中間層の各ノードと接続関係がある入力層の各ノードでの活性度と各エッジの重みとバイアス値とを取得し、その取得した入力層の各ノードでの活性度と各エッジの重みとバイアス値とを用いて、第1中間層の各ノードでの活性度を算出する。これにより、推論を行う際の計算量及びメモリ量を削減することができる。また、より高い推論精度を得ることができる。

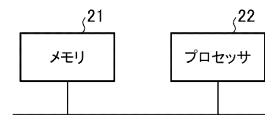
【図1】



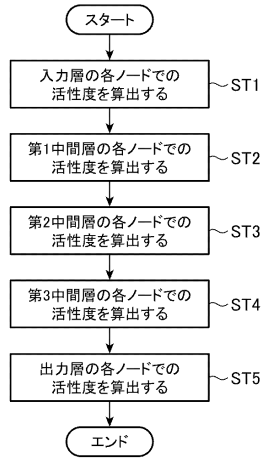
【図2】



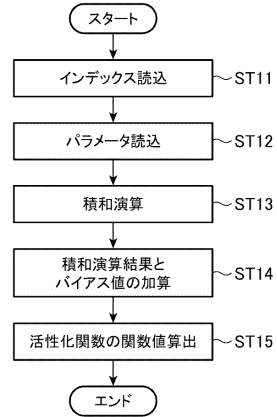
【図3】



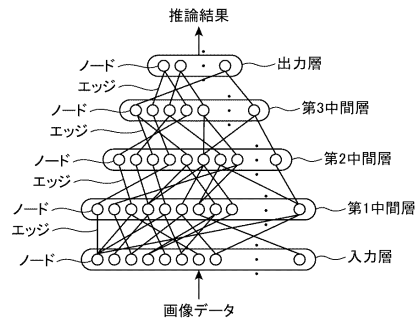
【 図 4 】



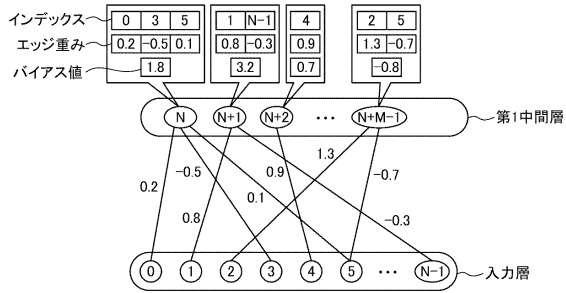
【 図 5 】



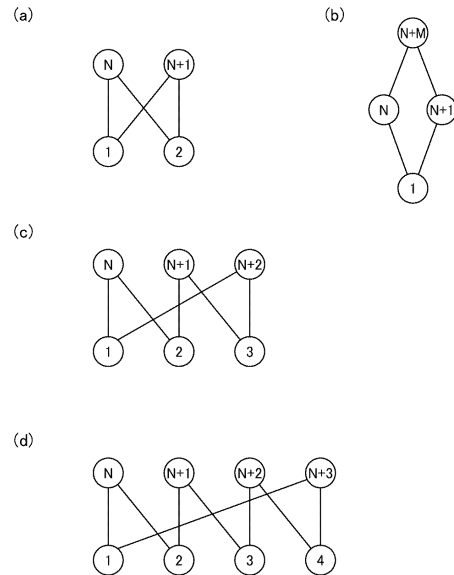
【 図 6 】



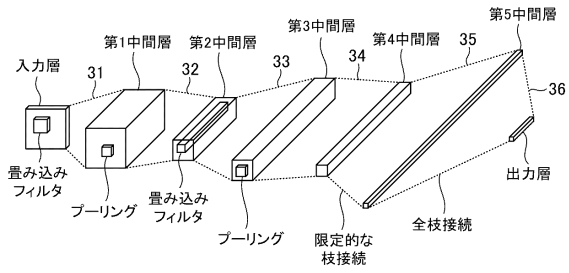
【 図 7 】



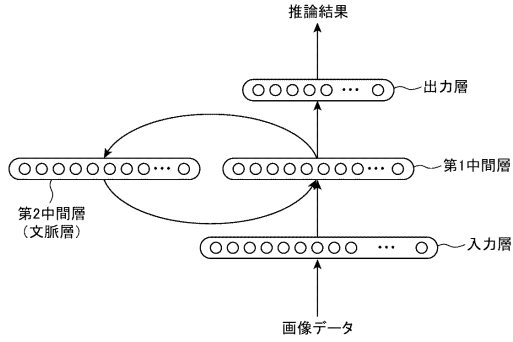
【 図 8 】



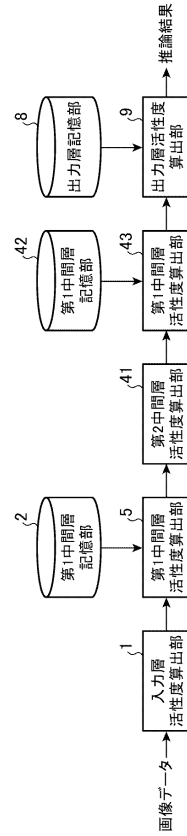
【図9】



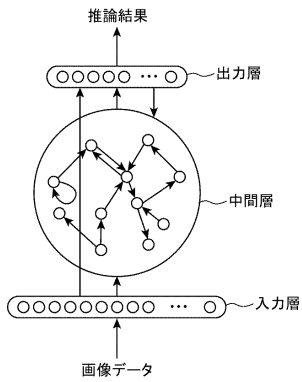
【図10】



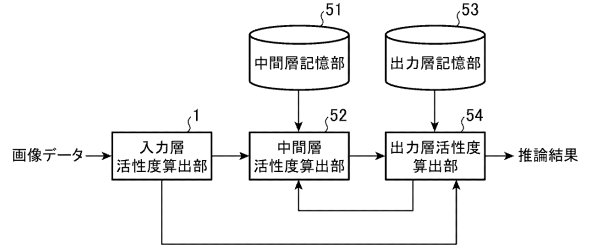
【図11】



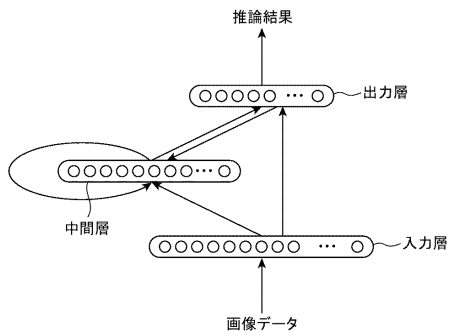
【図12】



【図14】



【図13】



フロントページの続き

- (72)発明者 松本 渉
東京都千代田区丸の内二丁目7番3号 三菱電機株式会社内
- (72)発明者 吉村 玄太
東京都千代田区丸の内二丁目7番3号 三菱電機株式会社内
- (72)発明者 趙 雄心
東京都千代田区丸の内二丁目7番3号 三菱電機株式会社内

審査官 多賀 実

- (56)参考文献 特開平3 - 55658 (JP, A)
特開平2 - 236659 (JP, A)
特開2002 - 251601 (JP, A)
特開平4 - 355889 (JP, A)

- (58)調査した分野(Int.Cl., DB名)
G06N3/02 - 3/10
G06T7/00