

(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2021-22050
(P2021-22050A)

(43) 公開日 令和3年2月18日(2021.2.18)

(51) Int.Cl.
G06N 3/08 (2006.01)

F I
G06N 3/08 120

テーマコード (参考)

審査請求 未請求 請求項の数 10 O L (全 14 頁)

(21) 出願番号 特願2019-137019 (P2019-137019)
(22) 出願日 令和1年7月25日 (2019.7.25)

(71) 出願人 504145283
国立大学法人 和歌山大学
和歌山県和歌山市栄谷930番地
(74) 代理人 100111567
弁理士 坂本 寛
(72) 発明者 和田 俊和
和歌山県和歌山市栄谷930番地 国立大
学法人和歌山大学内
(72) 発明者 菅間 幸司
和歌山県和歌山市栄谷930番地 国立大
学法人和歌山大学内

(54) 【発明の名称】 ニューラルネットワークの圧縮方法、ニューラルネットワーク圧縮装置、コンピュータプログラム、及び圧縮されたニューラルネットワークデータの製造方法

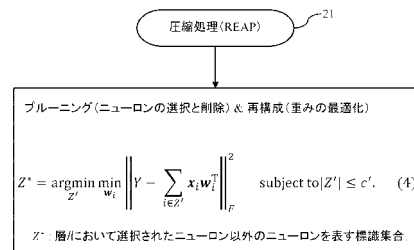
(57) 【要約】

【課題】再構成誤差が最小になるようにプルーニングが行われるようにする。

【解決手段】開示のニューラルネットワークの圧縮方法は、選択されたプルーニング対象に対するプルーニングと、プルーニングされたニューラルネットワークの再構成と、が行われる工程を備える。前記プルーニング対象は、プルーニング及び再構成によって生じる再構成誤差が最小になるように選択される。

【選択図】 図4

図4



【特許請求の範囲】**【請求項 1】**

選択されたブルーニング対象に対するブルーニングと、ブルーニングされたニューラルネットワークの再構成と、が行われる工程を備え、

前記ブルーニング対象は、ブルーニング及び再構成によって生じる再構成誤差が最小になるように選択される

ニューラルネットワークの圧縮方法。

【請求項 2】

前記ブルーニングは、全結合層におけるニューロンのブルーニングである

請求項 1 に記載のニューラルネットワークの圧縮方法。

10

【請求項 3】

前記ブルーニングは、畳み込み層におけるチャンネルのブルーニングである

請求項 1 に記載のニューラルネットワークの圧縮方法。

【請求項 4】

前記再構成誤差は、ブルーニング対象の挙動ベクトルを、前記ブルーニング対象以外の他のブルーニング単位の挙動ベクトルが張る部分空間に射影したときの射影残差に基づいて計算される

請求項 1 から 3 のいずれか 1 項に記載のニューラルネットワークの圧縮方法。

【請求項 5】

前記射影残差は、ブルーニング対象の挙動ベクトルの双直交基底に基づいて計算される

請求項 4 に記載のニューラルネットワークの圧縮方法。

20

【請求項 6】

前記射影残差は、グラン・シュミットの直交化計算の反復適用により計算される

請求項 4 に記載のニューラルネットワークの圧縮方法。

【請求項 7】

前記再構成誤差は、並列計算される

請求項 1 から 6 のいずれか 1 項に記載のニューラルネットワークの圧縮方法。

【請求項 8】

選択されたブルーニング対象に対するブルーニングと、ブルーニングされたニューラルネットワークの再構成と、を行うニューラルネットワーク圧縮装置であって、

前記ブルーニング対象は、ブルーニング及び再構成によって生じる再構成誤差が最小になるように選択されるよう構成されている

ニューラルネットワーク圧縮装置。

30

【請求項 9】

コンピュータを、ニューラルネットワーク圧縮装置として機能させるためのコンピュータプログラムであって、

前記ニューラルネットワーク圧縮装置は、選択されたブルーニング対象に対するブルーニングと、ブルーニングされたニューラルネットワークの再構成と、を行うよう構成され、

前記ニューラルネットワーク圧縮装置において、前記ブルーニング対象は、ブルーニング及び再構成によって生じる再構成誤差が最小になるように選択される

コンピュータプログラム。

40

【請求項 10】

選択されたブルーニング対象に対するブルーニングと、ブルーニングされたニューラルネットワークの再構成と、が行われる圧縮工程と、

前記圧縮工程により圧縮されたニューラルネットワークのデータを出力する工程と、を有し、

前記圧縮工程において、前記ブルーニング対象は、ブルーニング及び再構成によって生じる再構成誤差が最小になるように選択される

圧縮されたニューラルネットワークデータの製造方法。

50

【発明の詳細な説明】

【技術分野】

【0001】

本開示は、ニューラルネットワークの圧縮に関する。

【背景技術】

【0002】

ディープニューラルネットワーク(DNN)のようなニューラルネットワークの圧縮手法として、プルーニング(Pruning; 枝刈り)と、再構成(Reconstruction)と、を行う手法がある。

【0003】

プルーニングは、全結合型ニューラルネットワーク(FCN)においては、ニューロン(とそのニューロンに接続された重み)の削除として行われ、畳み込み型ニューラルネットワーク(CNN)においては、チャンネルの削除として行われる(非特許文献1参照)。チャンネルの削除は、削除されるチャンネルに属する重み全体の削除として行われる。

【0004】

再構成はプルーニング後に行われる。再構成では、プルーニング前の出力に近づくように、ニューラルネットワークの重みパラメータの調整が行われる。例えば、FCNにおいては、再構成として、ニューロン間の結合の重みパラメータの調整が行われ、CNNにおいては、再構成として、フィルタ(カーネル)における重みパラメータの調整が行われる。

【先行技術文献】

【非特許文献】

【0005】

【非特許文献1】Yihui He, Xiangyu Zhang, Jian Sun, "Channel Pruning for Accelerating Very Deep Neural Networks," Proc. of ICCV2017, 2017

【発明の概要】

【0006】

プルーニングと再構成とを行う従来の圧縮手法においては、再構成後における誤差(再構成誤差)が最小になるようにプルーニングが行われるわけではない、という課題が存在することを、本発明者らは見出した。以下では、この課題について説明する。なお、以下の説明では、簡単化のため、FCNを前提として説明する。

【0007】

プルーニングをする際には、削除されるニューロン(プルーニング対象)を選択する必要がある。削除されるニューロンは、削除されるニューロンが存在する層の次層に与える誤差に着目して、選択される。具体的には、ニューロンの削除によって次層に与える誤差が最小となるニューロンが、削除されるニューロンとして選択される。例えば、ニューロン A_1 を削除した場合に、次層に与える誤差が E_1 であり、ニューロン A_2 を削除した場合に、次層に与える誤差が E_2 である場合、誤差 E_1 が誤差 E_2 よりも小さければ、ニューロン A_1 が、ニューロン A_2 よりも優先して、削除されるニューロンとして選択されることになる。

【0008】

プルーニング後の再構成では、削除されずに残ったニューロンから次層のニューロンへ向かう結合における重みパラメータが、調整される。重みの調整は、ニューロンの削除により次層に与える誤差が最小になるように実行される。例えば、プルーニングにおいてニューロン A_1 を削除することで次層に与える誤差が E_1 である場合、再構成では、誤差 E_1 ができるだけ小さくなるように、重みの調整が行われる。重みの調整により最小化された誤差 E_1 は、再構成誤差 E_{1r} と呼ばれる。

【0009】

以上のような従来の圧縮手法では、再構成誤差 E_{1r} が最小になるようにプルーニングが行われるわけではない。例えば、前述のように、ニューロン A_1 を削除した場合に次層

10

20

30

40

50

に与える誤差が E_1 であり、再構成誤差が E_{1r} であるとする。また、ニューロン A_2 を削除した場合に次層に与えられる誤差が E_2 であり、再構成誤差が E_{2r} であるとする。この場合において、誤差 E_1 が誤差 E_2 よりも小さいとしても、再構成誤差 E_{1r} が再構成誤差 E_{2r} よりも大きいことがある。すなわち、削除により生じる誤差が最小であっても、再構成誤差が最小になるとの保証はない。

【0010】

したがって、上記の課題の解決が望まれる。本開示において、上記の課題は、再構成誤差が最小になるようにブルーニングすることにより解決される。更なる詳細は、後述の実施形態として説明される。

【図面の簡単な説明】

10

【0011】

【図1】図1は、ニューラルネットワーク圧縮装置及びニューラルネットワーク利用装置の構成図である。

【図2】図2は、ニューラルネットワークの構成及び伝播量 Y の定式化の説明図である。

【図3】図3は、比較例に係る圧縮処理のフローチャートである。

【図4】図4は、実施形態に係る圧縮処理のフローチャートである。

【図5】図5は、実施形態に係る圧縮処理のための Greedy Algorithm である。

【図6】図6は、部分空間 U への x_j の射影を示す図である。

【図7】図7は、残差 r_j の計算方法を示す。

20

【図8】図8は、残差 r_j の計算方法を示す。

【図9】図9は、別のニューロンを削除する際の再構成誤差の計算方法を示す。

【図10】図10は、別のニューロンを削除する際の再構成誤差の計算方法を示す。

【図11】図11は、REAPの畳み込み層への適用を示す。

【図12】図12は、グラン・シュミットの直交化計算の適用による射影残差 r_j の計算方法を示す。

【図13】図13は、実験結果を示す図である。

【発明を実施するための形態】

【0012】

< 1. ニューラルネットワークの圧縮方法、ニューラルネットワーク圧縮装置、コンピュータプログラム、及び圧縮されたニューラルネットワークデータの製造方法の概要 >

30

【0013】

(1) 実施形態に係るニューラルネットワークの圧縮方法は、選択されたブルーニング対象に対するブルーニングと、ブルーニングされたニューラルネットワークの再構成と、が行われる工程を備える。前記ブルーニング対象は、ブルーニング及び再構成によって生じる再構成誤差が最小になるように選択される。これにより、再構成誤差が最小になるようにブルーニングされる。

【0014】

(2) 前記ブルーニングは全結合層におけるニューロンのブルーニングであってもよい。

【0015】

(3) 前記ブルーニングは畳み込み層におけるチャンネルのブルーニングであってもよい。

40

【0016】

(4) 前記再構成誤差は、ブルーニング対象の挙動ベクトルを、前記ブルーニング対象以外の他のブルーニング単位の挙動ベクトルが張る部分空間に射影したときの射影残差に基づいて計算されるのが好ましい。

【0017】

(5) 前記射影残差は、ブルーニング対象の挙動ベクトルの双直交基底に基づいて計算されるのが好ましい。

【0018】

(6) 前記射影残差は、グラン・シュミットの直交化計算の反復適用により計算されてもよい。

50

【 0 0 1 9 】

(7) 前記再構成誤差は、並列計算等で高速化されるのが好ましい。

【 0 0 2 0 】

(8) 実施形態に係るニューラルネットワーク圧縮装置は、選択されたブルーニング対象に対するブルーニングと、ブルーニングされたニューラルネットワークの再構成と、を行うよう構成されている。ニューラルネットワーク圧縮装置は、前記ブルーニング対象は、ブルーニング及び再構成によって生じる再構成誤差が最小になるように選択されるよう構成されている。

【 0 0 2 1 】

(9) 実施形態に係るコンピュータプログラムは、コンピュータを、ニューラルネットワーク圧縮装置として機能させるためのコンピュータプログラムである。前記ニューラルネットワーク圧縮装置は、選択されたブルーニング対象に対するブルーニングと、ブルーニングされたニューラルネットワークの再構成と、を行うよう構成され、前記ニューラルネットワーク圧縮装置において、前記ブルーニング対象は、ブルーニング及び再構成によって生じる再構成誤差が最小になるように選択される。

10

【 0 0 2 2 】

(1 0) 実施形態に係る圧縮されたニューラルネットワークデータの製造方法は、選択されたブルーニング対象に対するブルーニングと、ブルーニングされたニューラルネットワークの再構成と、が行われる圧縮工程と、前記圧縮工程により圧縮されたニューラルネットワークのデータを出力する工程と、を有し、前記圧縮工程において、前記ブルーニング対象は、ブルーニング及び再構成によって生じる再構成誤差が最小になるように選択される。

20

【 0 0 2 3 】

< 2 . ニューラルネットワークの圧縮方法、ニューラルネットワーク圧縮装置、コンピュータプログラム、及び圧縮されたニューラルネットワークデータの製造方法の例 >

【 0 0 2 4 】

図 1 は、実施形態に係るニューラルネットワーク圧縮装置（以下、「圧縮装置」という）1 0 とニューラルネットワーク利用装置（以下、「利用装置」という）1 0 0 とを示している。実施形態に係る圧縮装置 1 0 は、ニューラルネットワーク N 1 を圧縮して小規模化するための圧縮処理 2 1 を実行する。圧縮処理 2 1 を実行することにより実施される方法は、圧縮されたニューラルネットワークの製造方法又は圧縮されたニューラルネットワークデータの製造方法でもある。

30

【 0 0 2 5 】

ニューラルネットワークは、複数の人工ニューロン（「ノード」ともいう）が結合した人工的な計算機構である。ニューラルネットワークは、例えば、ディープニューラルネットワーク（DNN）である。DNNは、例えば、全結合型ニューラルネットワーク（FCN）であってもよいし、畳み込み型ニューラルネットワーク（CNN）であってもよい。以下では、圧縮処理の対象となるニューラルネットワーク N 1 を、「原ニューラルネットワーク」といい、圧縮されたニューラルネットワーク N 2 を「圧縮ニューラルネットワーク」という。なお、実施形態に係る圧縮装置 1 0 は、原ニューラルネットワーク N 1 の機械学習（深層学習）のための処理も実行可能である。圧縮装置 1 0 は、学習済の原ニューラルネットワーク N 1 を圧縮する。

40

【 0 0 2 6 】

圧縮装置 1 0 は、1 又は複数のプロセッサ 2 0 及び記憶装置 3 0 を有するコンピュータによって構成されている。1 又は複数のプロセッサ 2 0 は、例えば、グラフィックプロセッシングユニット（GPU）を含む。1 又は複数のプロセッサ 2 0 は、さらに CPU を含んでもよい。GPU のような大規模並列計算機構は、大規模なニューラルネットワークに関する処理を実行するための大量の計算に適している。

【 0 0 2 7 】

記憶装置 3 0 は、プロセッサ 2 0 によって実行されるコンピュータプログラム 3 1 を記

50

憶している。プロセッサ 20 は、コンピュータプログラム 31 を実行することで、圧縮処理 21 を行う。圧縮処理 21 は、プルーニング (Pruning ; 枝刈り) と再構成 (Reconstruction) とを含む。

【0028】

記憶装置 30 は、圧縮処理 21 によって製造された圧縮ニューラルネットワーク N2 を表すデータ (圧縮ニューラルネットワークデータ) N20 を記憶することができる。圧縮ニューラルネットワークデータ N20 は、圧縮ニューラルネットワーク N2 を表現する各種のパラメータ (重み、結合関係など) からなるデータである。圧縮装置 10 は、圧縮ニューラルネットワークデータ N20 を、ニューラルネットワークエンジン等へ、出力することができる。圧縮ニューラルネットワークデータ N20 は、ニューラルネットワークエンジンに読み込まれることで、そのニューラルネットワークエンジンを圧縮ニューラルネットワーク N2 として機能させる。

10

【0029】

利用装置 100 は、圧縮ニューラルネットワークデータ N20 を読み込んで、圧縮ニューラルネットワーク N2 として機能するニューラルネットワークエンジンを有する。ニューラルネットワークエンジンは、例えば、プロセッサ 200 と記憶装置 300 とを備える。プロセッサ 200 は、例えば、組み込み系システムにおける低消費電力の CPU でよい。圧縮ニューラルネットワークデータ N20 は、原ニューラルネットワーク N1 のデータに比べて、サイズが小さいため、低消費電力の CPU による処理が可能である。

【0030】

組み込み系システムは、汎用的なコンピュータシステムではなく、特定の用途に向けられたコンピュータシステムであり、例えば、スマートフォン・家電などの家庭用機器、産業用ロボットなどの産業用機器、各種の医療用機器、自動車・ドローンなどのビークル、及びその他の機器におけるコンピュータシステムである。組み込み系システムでは、プロセッサとして、低消費電力の CPU が使われることが多いが、圧縮ニューラルネットワークデータ N20 は、データサイズが小さいため、実行が容易である。

20

【0031】

圧縮ニューラルネットワーク N2 は、例えば、画像・音声の変換、セグメンテーション、識別などの用途に用いられる。より具体的には、例えば、店舗等の客数計測、男女・年齢層分析、車両計数、車種分析など、対象物の画像から必要な情報を抽出するために用いることができる。原ニューラルネットワーク N1 は大規模であり、計算コストが大きいため、組み込み系システムでの実行が困難であるが、圧縮ニューラルネットワーク N2 は、小規模化されているため、組み込み系システムでの実行が容易である。

30

【0032】

以下、圧縮処理 21 について説明する。以下では、理解の容易のため、まず、全結合型ニューラルネットワーク (FCN) を前提に、圧縮処理 21 を説明し、その後、同様の圧縮処理 21 を、畳み込み型ニューラルネットワーク (CNN) に適用できることを説明する。

【0033】

図 2 は、原ニューラルネットワーク N1 である全結合型ニューラルネットワーク (FCN) における層 1 と、層 1 の次の層である層 1 + 1 と、を示している。図 2 では、2 つの層 (1 層, 1 + 1 層) が代表的に示されている。FCN における各層は、層状に並べられた人工ニューロン (以下、単に「ニューロン」という) が層間で結合されている全結合層である。各層中における丸印がニューロンである。層 1 に含まれるニューロン数は c であり、層 1 + 1 に含まれるニューロン数は C である。

40

【0034】

図 2 中の式 (1) は、ニューラルネットワークにデータが与えられた時における、層 1 から次の層 1 + 1 への伝播量 Y を定式化している。ここでは、 Y は、層 1 + 1 の C 個のニューロンの内部活性度を表す $N \times C$ 行列とする。換言すると、 Y は、層 1 から与えられる、層 1 + 1 への入力でもある。

50

【 0 0 3 5 】

N個のデータ（例えば、N個の画像データ）を、層1のc個のニューロンに与えた場合、層1の各ニューロンからはN個の出力が生じる。層1におけるi番目のニューロンの出力が x_i で表される。 x_i は、N次元のベクトルである。 x_i は、i番目のニューロンにN個のデータが与えられた時のi番目のニューロンの出力（挙動）を示す。すなわち、 x_i はi番目のニューロンの挙動ベクトル（ニューロン挙動ベクトル）でもある。なお、ニューラルネットワークに与えられるデータは、画像以外の他のデータ、例えば、音声データ等であってもよい。画像データ等のデータは、各ニューロンの挙動を把握するために、ニューラルネットワークに与えられる。

【 0 0 3 6 】

図2中の w_i は、1層のi番目のニューロンから、1+1層のC個のニューロンへ向かう結合の重み（重み係数）からなるC次元の重みベクトルである。

【 0 0 3 7 】

この場合、次層1+1への伝播量Yは、層1におけるニューロンの出力 x_i と、層1から次層1+1への重みベクトル w_i と、によって、図2中の式(1)に示すように定式化される。

【 0 0 3 8 】

実施形態に係る圧縮処理21の目的は、上記のYをできるだけ変化させることなく、ニューロンの数を、所望の数ほど減少させることである。ニューロンを減少させても、Yの変化が少なければ、原ニューラルネットワークN1の性能を維持することができる。つまり、ニューラルネットワークを圧縮しても、精度低下を防止できる。なお、上記のように、FCNでは、ニューロンがプルーニング単位であるが、CNNでは、チャンネルがプルーニング単位である。なお、プルーニング単位は、削除の単位である。

【 0 0 3 9 】

実施形態に係る圧縮処理21の説明に先立ち、比較例に係る圧縮処理121を説明する。図3は、比較例に係る圧縮処理121を示している。図3に示す圧縮処理121は、プルーニング工程122と、再構成工程123と、を有している。比較例においては、プルーニング工程122と再構成工程123とは完全に分離している。

【 0 0 4 0 】

プルーニング工程122では、ある層1に含まれる複数のニューロン（複数のプルーニング単位）から削除されるニューロン（プルーニング対象）が選択され、選択されたニューロンの削除が行われる。削除されるニューロンが層1の中から選択される場合、次層1+1への伝播量Yに与える影響が最も小さくなるニューロンが、削除されるニューロン（プルーニング対象）として選択される。この選択の際には、Lasso回帰を用いて、図3の式(2)に従ってニューロンが選択される（非特許文献1参照。非特許文献1ではチャンネルが選択される）。比較例においては、ニューロンの選択の際には、重みが調整されることはない。

【 0 0 4 1 】

比較例では、次層1+1での内部活性度の誤差に関するペナルティ項に重要度ベクトルの L_1 ノルムを加えている。重要度ベクトルの L_1 ノルムを最小化する重要度ベクトルを求めると、次層1+1の活性度の誤差を低く抑えつつゼロ要素の多いが得られ、削除すべきニューロンを決定できる。すなわち、式(2)の最適化の結果、最適化されたニューロンの重要度ベクトル w^* が得られるが、そのベクトルのi番目の要素 w_i^* が0ならば、i番目のニューロンは不要であり、削除されるニューロンとして選択される。

【 0 0 4 2 】

比較例において、削除されるニューロンの数は、ハイパーパラメータの微調整によってコントロールされる。 λ を増加させれば、削除されるニューロンの数が増え、 λ を減少させれば、削除されるニューロンの数が減る。比較例においては、削除されるニューロンの数は、 λ によってコントロールされるため、削除されるニューロンの数の制御は難しい。

10

20

30

40

50

【0043】

再構成工程123では、層1において、プルーニング後に残ったニューロンが、次層 $l+1$ に与える Y が、プルーニング前における Y （本来の Y ）に近づくように、重みが調整（最適化）される。重みの調整は、図3中の式(3)に従って行われる。式(3)は、再構成誤差を最小化する重みベクトルを求める。ここでの再構成誤差は、プルーニング前の Y と、プルーニング後に重みを調整したときの Y と、の差に基づく。

【0044】

比較例においては、プルーニング対象は、再構成を行う前の誤差を最小化するように選択されており、再構成後に最も誤差が小さくなるように選択されているわけではない。つまり、比較例では、プルーニング対象の選択は、再構成前の誤差に基づいて行われており、再構成は、再構成後の誤差に基づいて行われており、プルーニングと再構成とが、異なる基準で行われている。また、比較例においては、Lassoを用いており、削除されるニューロンの数をコントロールするには、の人手による微調整が必要となる。つまり、比較例では、削除されるニューロンの数のコントロールは困難である。

10

【0045】

図4は、実施形態に係る圧縮処理21を示している。以下では、実施形態に係る圧縮処理21を、REAP (Reconstruction Error Aware Pruning) という。

【0046】

比較例では、プルーニングと再構成とが異なる基準で行われていたのに対して、REAPでは、プルーニングと再構成とを同じ基準で行う。すなわち、REAPでは、再構成誤差が最小になるようにプルーニングされるとともに再構成される。REAPでは、図4中の式(4)に従って、削除されるニューロンが決定される。なお、式(4)の Z^* は、層1において、削除して残ったニューロンを示す。式(4)においては、重みベクトル w_j は、 Z' を Z^* に固定する前に最適化される。したがって、式(4)によれば、次層 $l+1$ への伝播量 Y の再構成誤差を最小化するニューロンの集合が求まる。

20

【0047】

REAPでは、再構成誤差が最小になるようにプルーニング対象であるニューロンが選択されるため、REAPは、再構成誤差が最小になるとは限らない比較例に比べて、有利である。

【0048】

式(4)は、組み合わせ最適化問題であり、グリーディ法 (Greedy Algorithm) によって解かれる。図5は、式(4)を解くためのアルゴリズムを示している。まず、ステップS1において、層1における j 番目のニューロンを消してみる。ステップS2において、層1において残ったニューロンのみで Y を再構成して誤差（再構成誤差）を計算する。再構成誤差は、図5中において式(5)として示されるコスト関数 $P(Z')$ を計算することで求まる。再構成誤差が求まると、一旦、削除した j 番目のニューロンを元に戻す。ステップS3で示される繰り返しループにおいては、ステップS1及びステップS2がすべての j ($j \in Z$) について繰り返され、最も $P(Z')$ の値が小さくなるニューロンが、プルーニング対象として選択され、最終的に削除される。

30

【0049】

ステップS3の繰り返しループによって、一つのニューロンが削除される。ステップS4で示される繰り返しループにおいては、残ったニューロンのみで、再度、ステップS3の繰り返しループが実行される。再度、ステップ3の繰り返しループが実行されると、別のニューロンが削除される。

40

【0050】

層1からいくつのニューロンを削除するかは、ステップS4の繰り返しループを何回実行するかによって決まる。したがって、所望数 D のニューロンを削除したい場合、ステップS4の繰り返しループを D 回実行すればよい。したがって、REAPでは、削除されるニューロンの数のコントロールは容易である。

【0051】

50

REAPでは、比較例に比べて、計算量が増加する。すなわち、REAPでは、再構成誤差を求める際に最小二乗法を適用するため、連立方程式を解く必要がある。そして、解く必要のある連立方程式は層内のニューロン数分存在する。例えば、層 l のニューロン数が c であり、次層 $l+1$ のニューロン数が C である場合、重みパラメータ数は $c \times C$ になる。1個のニューロンを削除する場合、係数行列のサイズが $(c-1)C \times (c-1)C$ となる。したがって、一つの連立方程式を解くための時間計算量は、 $O(c^3 C^3)$ である。この連立方程式を c 回解く必要があることから、 $O(c^4 C^3)$ の時間計算量となる。

【0052】

ここで、本来解こうとしている最小二乗問題は、あるニューロンの挙動を表す挙動ベクトル x_j がなくなったときに、残りのニューロン集合 Z' の挙動ベクトル $x_i (i \in Z')$ の線形和で、次層 $l+1$ への本来の伝播量 Y を近似する、という問題である。この近似による誤差は、図6に示すように、 x_j を $x_i (i \in Z')$ の線形和で表現した際の誤差 r_j に起因している。したがって、この誤差 r_j を最小化すれば、誤差 r_j を示すベクトルに、次層 $l+1$ への重みベクトル w_i^T を掛けるだけで、次層 $l+1$ での活性度(伝播量 Y)の誤差が計算できる。

10

【0053】

すなわち、 j 番目のニューロンの削除により生じる再構成誤差を計算するためには、図6中の式(6-1)に示すコスト関数を計算する必要がある。後述するように、式(6-1)のコスト関数は、式(6-2)のように表される。したがって、 x_j を $x_i (i \in Z')$ が張る部分空間に射影した際の射影残差 r_j を、巨大な係数行列を用いることなく、計算することができれば、再構成誤差を効率よく計算することができる。一例として、残差 r_j は、式(6-3)のように計算される。

20

【0054】

残差 r_j の求め方の一例は、図7及び図8に詳しく説明されている。図7及び図8に示す計算方法は、残差 r_j が、 x_j の双直交基底と線形従属であることを利用したものである。すなわち、残差 r_j は、 x_j の双直交基底に対する、 x_j の射影である。残差 r_j は、 x_j の双直交基底に基づいて計算される。図7及び図8に示す計算方法によれば、連立方程式の係数行列を使用せずに、残差 r_j を効率的に計算できる。また、図9及び図10は、あるニューロンを削除した後に、別のニューロンを削除するための再構成誤差を効率的に計算する方法を説明している。

30

【0055】

高速化のため、再構成誤差の計算は、並列処理で行うのが好ましい。再構成誤差を並列計算することで、高速に再構成誤差を求めることができる。ただし、最小二乗法を解く際に必要となる係数行列を格納するメモリ量は、並列化によって増大する。

【0056】

したがって、再構成誤差の計算を効率的に行うには、消費メモリ量に相当する空間計算量を削減することが好ましい。ここで、空間計算量(消費メモリ)は一つの連立方程式当たり、 $O(c^2 C^2)$ であり、 c 個並列計算で同時に計算すると $O(c^3 C^2)$ になる。空間計算量を削減するには、連立方程式の係数行列を使用せずに最小二乗法の計算を行うのが好ましい。

40

【0057】

図12は、並列計算で残差 r_j を求める例を示している。図12では、グラム・シュミット(Gram-Schmit)の直交化計算を反復適用することで、残差 r_j を効率的に計算できる。図11に示す計算方法では、連立方程式の係数行列を使用せずに、残差 r_j を効率的に並列計算できる。ニューロン数が非常に多い(実行環境によるが、目安として、4096個以上)の場合は、グラム・シュミットを用いる解法の方が高速である。

【0058】

図11は、REAPが畳み込み層におけるプルーニング及び再構成に適用できることを説明している。図11中の式(19)に示すように、畳み込み層におけるスライディング

50

ウィンドウ操作は、行列乗算の和によって表される。式(19)は、全結合層のための式(1)と同様の形式であることから、畳み込み層においても、全結合層と同様に、REAPを適用できることがわかる。

【0059】

図13は、ニューラルネットワークの圧縮をREAP及び比較例によって行った実験結果を示している。実験では、ImageNetデータセットによってトレーニングしたVGG16を、原ニューラルネットワークN1として用いた。原ニューラルネットワークN1に対する圧縮処理21としてREAPを適用した場合及び比較例を適用した場合それぞれについて、画像の認識精度(正解率)を求めた。

【0060】

図13の横軸は、FLOPs(浮動小数点演算数)を示し、縦軸は、正解率を示す。FLOPsが小さいほど、圧縮ニューラルネットワークN2の演算数が少なく、圧縮率が大きいことを示す。図13に示すように、比較例では、圧縮率を増加(削除されるニューロン数を増加)させると、正解率が0.7(70%)程度まで下がるのに対して、REAPでは、圧縮率を増加させても、正解率は0.8(80%)程度までしか下がらなかった。したがって、REAPの方が、圧縮による精度低下を抑制できていることがわかる。

【0061】

<3.付記>

本発明は、上記実施形態に限定されるものではなく、様々な変形が可能である。

【符号の説明】

【0062】

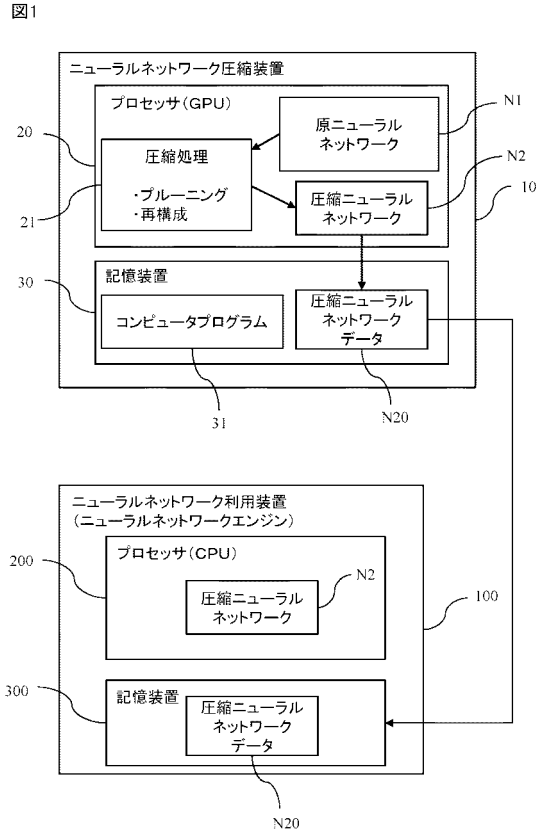
10 : ニューラルネットワーク圧縮装置
 20 : プロセッサ
 21 : 圧縮処理
 30 : 記憶装置
 31 : コンピュータプログラム
 100 : ニューラルネットワーク利用装置
 121 : 圧縮処理
 122 : プルーニング工程
 123 : 再構成工程
 200 : プロセッサ
 300 : 記憶装置
 N1 : 原ニューラルネットワーク
 N2 : 圧縮ニューラルネットワーク
 N20 : 圧縮ニューラルネットワークデータ

10

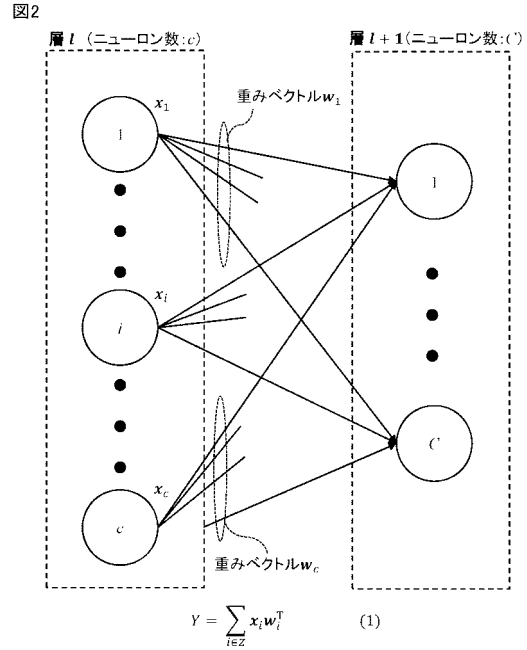
20

30

【 図 1 】

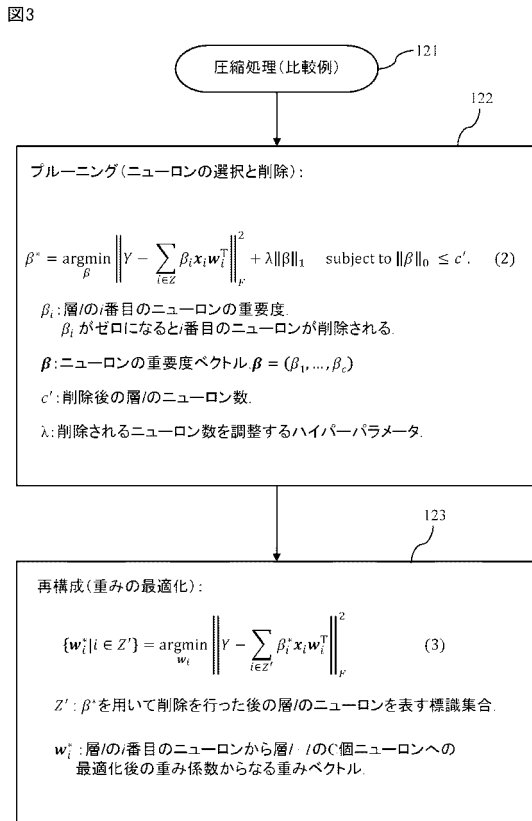


【 図 2 】

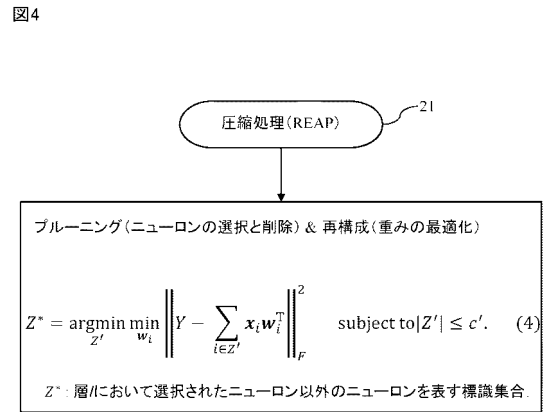


Z : 層 l のニューロンを表す標識集合, $Z = \{1, \dots, c\}$
 Y : N 個のデータを層 l の c 個のニューロンに与えたときに、層 l の c' 個のニューロンの内部活性化度を表す $N \times c'$ の行列 (層 $l+1$ への伝播量).
 x_i : N 個のデータを層 l の i 番目のニューロンに与えたときに、そのニューロンの出力からなる N 次元ベクトル (挙動ベクトル).
 w_i : 層 l の i 番目のニューロンから層 l の c' 個ニューロンへの重み係数からなる重みベクトル
 $x_i w_i^T$: 層 l の i 番目のニューロンによって発生させられる層 l の c' 個のニューロンの内部活性化度.

【 図 3 】

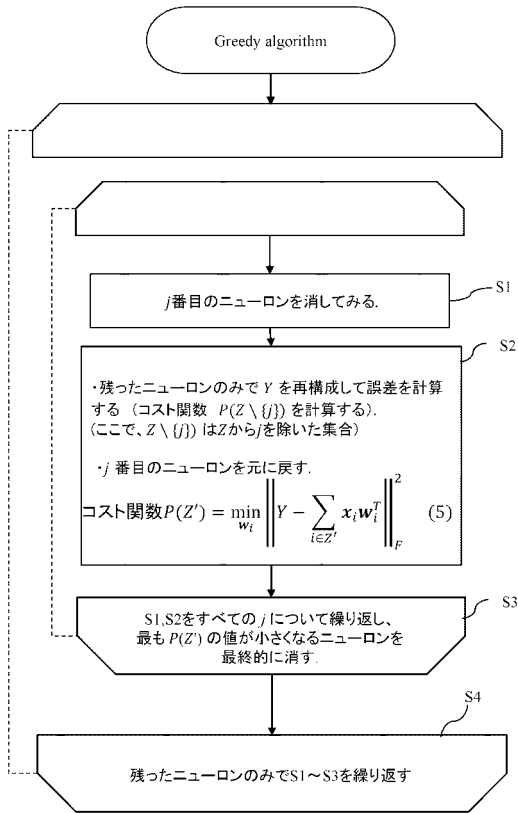


【 図 4 】



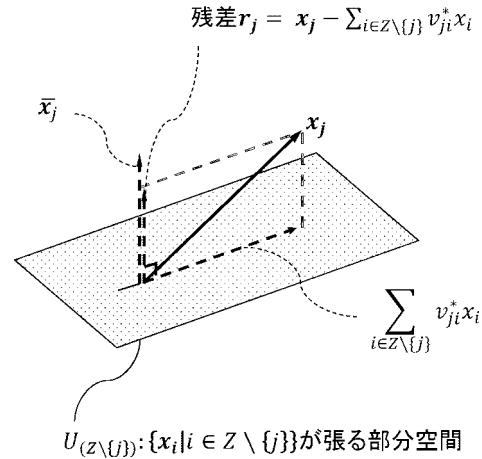
【 図 5 】

図5



【 図 6 】

図6



$U_{(Z \setminus \{j\})}$ への x_j の直交射影の計算は、 $\{x_i | i \in Z \setminus \{j\}\}$ から x_j を再構成する問題を最小二乗法により解くことと等価である。残差 r_j は、 x_j の双対基底である \bar{x}_j に線形従属である。

$$P(Z \setminus \{j\}) = \min_{w_i} \left\| Y - \sum_{i \in Z \setminus \{j\}} x_i w_i^T \right\|_F^2 \quad (6-1)$$

$$= \|r_j w_j^T\|_F^2 \quad (6-2)$$

r_j は、 x_j の \bar{x}_j に対する射影であり、以下のように計算される。

$$r_j = \frac{x_j^T \bar{x}_j}{\|\bar{x}_j\|^2} \bar{x}_j = \frac{\bar{x}_j}{\|\bar{x}_j\|^2} \quad (6-3)$$

【 図 7 】

図7

残差 r_j の計算 (1/2)

j 番目のニューロンの削除により生じる再構成誤差を計算するため、以下のコスト関数 $P(Z \setminus \{j\})$ を計算する。

$$P(Z \setminus \{j\}) = \min_{w_i} \left\| Y - \sum_{i \in Z \setminus \{j\}} x_i w_i^T \right\|_F^2 \quad (6-1)$$

部分問題の定義:

$$\{v_{ji}^* | i \in Z \setminus \{j\}\} = \underset{v_{ji}}{\operatorname{argmin}} \left\| x_j - \sum_{i \in Z \setminus \{j\}} v_{ji} x_i \right\|^2 \quad (7)$$

v_{ji} : x_i から x_j を再構成するための係数

式(7)を解くと、式(6-1)の解は、式(1)より、以下のように得られる。

$$P(Z \setminus \{j\}) = \left\| \left(x_j - \sum_{i \in Z \setminus \{j\}} v_{ji}^* x_i \right) w_j^T \right\|_F^2 \quad (8)$$

部分問題の効率的な解法:

式(7)は、 $\{x_i | i \in Z \setminus \{j\}\}$ から再構成された x_j の誤差を最小化するための問題である。 $\{x_i | i \in Z \setminus \{j\}\}$ から再構成された x_j の計算は、図6に示す部分空間 $U_{(Z \setminus \{j\})}$ への x_j の直交射影の計算と等価である。

x_j の残差を r_j で示すと、以下の式が得られる。

$$r_j = x_j - \sum_{i \in Z \setminus \{j\}} v_{ji}^* x_i \quad (9)$$

明らかに、 r_j は $x_i (i \in Z \setminus \{j\})$ に直交し、 x_j には直交しない。すなわち、 r_j は x_j の双対基底に線形従属である。この r_j の性質を利用し、計算を効率化する。

【 図 8 】

図8

残差 r_j の計算 (2/2)

ここで、 $\{\bar{x}_i | i \in Z\}$ を、 $\{x_i | i \in Z\}$ の双対基底とする。 $X = [x_1 \dots x_n], \bar{X} = [\bar{x}_1 \dots \bar{x}_n]$ とすると、双対基底の定義より、 \bar{X} は次のように与えられる。

$$\bar{X} = (X^{\theta})^T \quad (10)$$

X^{θ} : X の一般化逆行列

r_j は x_j の双対基底に線形従属であるため、 r_j を得るには、図6に示すように、 x_j の \bar{x}_j への直交射影を計算すればよい。したがって、 r_j は、以下のように計算される。

$$r_j = \frac{x_j^T \bar{x}_j}{\|\bar{x}_j\|^2} \bar{x}_j = \frac{\bar{x}_j}{\|\bar{x}_j\|^2} \quad (6-3)$$

ここで、 D を対角行列とする。この対角行列の (i,i) 成分は $1/\|\bar{x}_i\|^2$ で与えられる。すると、 $R = [r_1 \dots r_n]$ は、以下のシンプルな行列乗算によって得られる。

$$R = \bar{X} D \quad (11)$$

ここで、 V^* を行列とする。この行列の (j,i) 成分は $v_{ji}^* (i \neq j)$ である。 $v_{ji}^* (i \neq j)$ は式(7)を解くことで得られる。 $v_{ji}^* = 0$ である。

式(9)より以下の方程式が成り立つ。

$$X = R + X V^{*T} \quad (12)$$

したがって、 $V^{*T} = I - X^{\theta} R$ が得られる。ここで、 I は単位行列である。

上記のように、とり得る j それぞれについて、式(7)を解く代わりに、いくつかの簡単な行列演算を行うことにより、等価な解が得られる。

【 図 9 】

図9

別のニューロンを削除する際の再構成誤差の計算 (1/2)

ここで、 j 番目のニューロンが削除されたものとする。
式(4)をgreedy流に解くには、次のステップとして、次に削除されるニューロンを選択するために、 $k(k \in Z \setminus \{j\})$ について $P(Z \setminus \{j, k\})$ を計算する。

単純なアプローチは、とり得る k それぞれについて $Z' = Z \setminus \{j, k\}$ とし、式(5)で定義される最小二乗問題を解くことである。しかし、このアプローチは計算コストが大きい。

以下に、 $P(Z \setminus \{j, k\})$ をより効率的に計算する方法を示す。

ここで、 r'_j を x_j の残差とし、 r'_k を x_k の残差とする。残差は、 $\{x_i | i \in Z \setminus \{j, k\}\}$ から計算される。

$Z' = Z \setminus \{j, k\}$ の場合、式(5)の解は、以下のように与えられる。

$$P(Z \setminus \{j, k\}) = \|r'_j w_j^T + r'_k w_k^T\|_F^2 \quad (13)$$

したがって、 r'_j 及び r'_k が得られれば $P(Z \setminus \{j, k\})$ を計算できる。

r'_j 及び r'_k は、いくつかのシンプルな線形代数のトリックによって得られる。以下では、 r'_k の計算方法だけを示す。 r'_j は r'_k と同様の方法で計算される。

以下の式(14)及び(15)が得られる。

$$x_j = r_j + v_{jk}^* x_k - \sum_{i \in Z \setminus \{j, k\}} v_{ji}^* x_i \quad (14)$$

$$x_k = r_k + v_{kj}^* x_j - \sum_{i \in Z \setminus \{j, k\}} v_{ki}^* x_i \quad (15)$$

【 図 1 0 】

図10

別のニューロンを削除する際の再構成誤差の計算 (2/2)

j 番目及び k 番目のニューロンを削除した後では、 x_k の再構成のために x_j を利用することができない。そこで、式(14)を式(15)に代入し、以下の式を得る。

$$x_k = \frac{r_k + v_{kj}^* r_j}{1 - v_{jk}^* v_{kj}^*} + \sum_{i \in Z \setminus \{j, k\}} \frac{v_{ki}^* - v_{kj}^* v_{ji}^*}{1 - v_{jk}^* v_{kj}^*} x_i \quad (16)$$

明らかに、 $x_q^T r_j = 0$ ($q \in Z \setminus \{j\}$)
 $x_p^T r_k = 0$ ($p \in Z \setminus \{k\}$)
であるから
 $x_i^T \frac{r_k + v_{kj}^* r_j}{1 - v_{jk}^* v_{kj}^*} = 0$ が、 $i \in Z \setminus \{j, k\}$ について成り立つ。

このことは、式(16)の右辺の第1項が、 $\{x_i | i \in Z \setminus \{j, k\}\}$ から再構成された x_k の残差を示すことを意味する。したがって、残差 r'_k は以下のように表される。

$$r'_k = \frac{r_k + v_{kj}^* r_j}{1 - v_{jk}^* v_{kj}^*} \quad (17)$$

残差 r'_j は残差 r'_k と同様に計算される。
 r'_j 及び r'_k を式(13)に代入すると、 $P(Z \setminus \{j, k\})$ が計算される。
 $P(Z \setminus \{j, k\})$ は、 $\{x_i | i \in Z \setminus \{j, k\}\}$ から再構成された Y の誤差を示す。

同時に、再構成のための係数を更新する。
 $\{v_{ki}^* | i \in Z \setminus \{j, k\}\}$ は、 $\{x_i | i \in Z \setminus \{j, k\}\}$ から x_j を再構成するための係数を示すと、以下の式が得られる。

$$v'_{ki} = \frac{v_{ki}^* - v_{kj}^* v_{ji}^*}{1 - v_{jk}^* v_{kj}^*} \quad (18)$$

上記のように、式(5)で定義される最小二乗問題を直接解くことなく、 $P(Z \setminus \{j, k\})$ を計算することができる。したがって、 $P(Z \setminus \{j, k\})$ を最小化する k を決定し、 k 番目のニューロンを削除することができる。上記と同じ手順を繰り返せば、さらに別のニューロンを削除できる。

【 図 1 1 】

図11

REAPの畳み込み層への適用

a_l は、入力チャネル数を示し、
 A_l は、出力チャネル数を示し、
 t_w は、特徴マップの幅を示し、
 t_h は、特徴マップの高さを示し、
 s_w は、重みテンソルの幅を示し、
 s_h は、重みテンソルの高さを示すものとする。

畳み込み層における演算に関して、
 N 入力画像に対応する特徴マップを示す $N \times a \times t_w \times t_h$ テンソルでのスライディングウィンドウ操作、及び
重みを示す $A \times a \times s_w \times s_h$ テンソルは、次のように表現できる。

$$Y = \sum_{i \in B} \phi_i \psi_i^T = \sum_{i \in B} \sum_{m \in C} \phi_{i(m)} \psi_{i(m)}^T \quad (19)$$

ここで、
 $B = \{1, \dots, a\}$ は、入力チャネルの標識集合を示す。
行列 $\phi_i \in \mathbb{R}^{N \times t_w \times s_w \times s_h}$ は、 i 番目の reshape された入力特徴マップを示す。
 ϕ_i の各行は元の特徴マップの部分テンソルを示す。
行列 $\psi_i \in \mathbb{R}^{A \times s_w \times s_h}$ は、 i 番目の reshape された重みテンソルを示す。
 $T = \{1, \dots, s_w, s_h\}$ は、 ϕ_i 及び ψ_i の列標識集合を示す。
 $\phi_{i(m)}$ 及び $\psi_{i(m)}$ は、 ϕ_i 及び ψ_i の m 番目の列を示す。

畳み込み層に関しては、チャネルレベルでのブルーニング(削除)が行われる。換言すると、 i 番目のチャネルの削除のために、 $\phi_{i(m)}$ 及び $\psi_{i(m)}$ ($m \in T$) が同時に削除される。他の手順については、全結合層における手順と同様である。

【 図 1 2 】

図12

グラム・シュミット(Gram-Schmit)の直交化計算の適用による射影残差 r_j の計算

グラム・シュミットの直交化計算で用いられる以下の射影計算を用いる。

$$F(r, x) = r - \frac{(r, x)}{\|x\|^2} x \quad (20)$$

ここでは、 x_1 を x_j 以外のベクトルに順次射影する $F(F(x_i, x_1), x_2), \dots, x_n)$ を計算すれば、これが r_j に漸近するという性質を用いる。この性質を用いると、並列計算のための共有メモリ領域に x_1, \dots, x_n を格納しておけば、以下のように、各射影残差 r_j を並列計算できる。

$$r_1 = F(F(x_1, x_2), x_3), \dots, x_n)$$

$$r_2 = F(F(x_2, x_1), x_3), \dots, x_n)$$

$$\vdots$$

$$r_n = F(F(x_n, x_1), x_n), \dots, x_{n-1})$$

【 図 1 3 】

図13

