2C1-OS-06a-5

級数展開に基づく表層非線形ネットワーク

Shallow Nonlinear Network Based on Fourier Series

| 窪澤駿平 *1*2 | 渡辺太郎*1 | 隅田英一郎 *1 | 岡田将吾* ² | 新田克己*2 |
|------------------|---------------|-----------------|--------------------|---------------|
| Shumpei KUBOSAWA | Taro WATANABE | Eiichiro SUMITA | Shogo OKADA | Katsumi NITTA |

*1情報通信研究機構 先進的音声翻訳研究開発推進センター

ASTREC, National Institute of Information and Communications Technology

*2東京工業大学大学院総合理工学研究科

Interdisciplinary Graduate School of Science and Engineering, Tokyo Institute of Technology

One of the most important characteristics of an artificial neural network is the ability of capturing nonlinearity inherent in the dataset. Existing deep architectures take the advantage of capturing global nonlinearity by stacking many hidden layers in order to obtain better representations. We propose an alternative method to handle nonlinearity under a shallow network, i.e., without a deep architecture, by fitting coefficients in a decomposed formula of a multivariate function motivated by Fourier series expansion. Analysis of this network on MNIST datasets reveals that this network enables efficient representation with respect to amount of total parameters. On this network, ℓ_2 regularization greatly helps to avoid overfitting, and L-BFGS method is effective for updating the parameters.

1. はじめに

言語モデルや音声認識,画像認識など様々な識別タスクにおいて、多層パーセプトロン(multi-layered perceptron; MLP) すなわちフィードフォワード・ニューラルネットワーク(feedforward neural network; FFNN)が広く用いられている. 識 別タスクにおいて MLP が利用される理由は、並列計算が容易 に実現できる単純な計算構造(計算グラフ)によって、データ セットに関する線形分離可能性などの性質や、これに基づく カーネルの選び方などの設定によらず、任意の問題を表現可能 だからである.任意の問題、すなわちデータセットの入出力関 係を表す連続関数は、少なくとも 2 層の MLP によって近似可 能であることが証明されている(万能性定理 *1)[1].

また,近年の深層学習(deep learning)では,多様な識別タ スクにおいて,2層を上回る多層化による汎化性能の向上が示 されている.なお,2層で十分とする万能性定理は,あくまで も「所与である」任意の入出力対に関する表現の万能性を示し たものであり,汎化性能その他については触れていない.一方 で,深層学習により汎化性能が向上したものの,計算グラフの 大規模化により全体のパラメータ数は増大した.また,各層が 多数直列に結合されることにより,訓練時の前向き計算の待ち 時間が増大し,データ量の増加とは別に計算時間も増加した.

そこで本稿では、少ない層数のニューラルネットワークにお ける、汎化性能の向上を目的とした新たなアプローチを提案す る.連続関数の大域的な傾向は、関数をフーリエ変換した周波 数領域における低周波成分に存在する.そこで本アプローチで は、フーリエ級数展開に基づく計算グラフによってデータセッ トを表現し、さらに低周波成分を優先的に捉えることで大域的 な特徴を獲得する.本アプローチの性質については、MNIST 手書き数字識別タスクを用いた実験結果を基に考察する.

連絡先: NICT ASTREC 先進的翻訳技術研究室,
〒 619-0289 京都府相楽郡精華町光台3-5, e-mail: shumpei.kubosawa@nict.go.jp

*1 一般には入力層も1層とみなして「3層パーセプトロンの表現万 能性定理」等呼ばれるが、本稿では活性化関数を含む層のみを数え ることとする。



図 1:2 次元空間における1次元へのアフィン変換の効果.格 子上の点を変換した際の様子を、各点の色と大きさで表した. 色は変換後の符号を、大きさは変換後の値の大小を表す.黒色 の直線が識別超平面(変換後の値が0)であり、これと直交す る方向(紫色の直線方向)にのみ識別が有効であることを表す.

2. 多層化の理由

多くの識別タスクにおける入出力は実数ベクトルである. こ のため、入出力の次元数をそれぞれ $d_i, d_o \in \mathbb{N}$ とすると、識 別タスクは、データセットの入力データ空間からラベル空間へ の写像 $f: \mathbb{R}^{d_i} \to \mathbb{R}^{d_o}$ を得る学習タスクとして定式化される. MLP を含むニューラルネットワークは一般に、予めパラメー タと計算順序および入出力により定義される計算グラフ (関数 定義)を用意し、そのパラメータを訓練データセットにフィッ ティングすることにより識別器を得る枠組みである. MLP は、 複数のパーセプトロンを入力側から出力側へ直列に接続した 識別モデルである. 一方で、1 層のパーセプトロンであって も、識別可能な問題は多く存在する. 例えば、シグモイド関 数 $\sigma(x) = \frac{1}{1+e^{-x}}$ を活性化関数とするパーセプトロンはロジ スティック回帰と等価である. では、なぜ多層化が必要なのだ ろうか?

多層化の理由の第一は、線形分離不可能な問題を含む任意の問題を表現するためである。Minsky[5] は、パーセプトロンでは線形分離不可能な問題の表現と訓練の両方が不可能であることを示した。そもそも、MLPの各層を構成するパーセプトロンは、入出力の次元数を $l_i, l_o \in \mathbb{N}$ とすると、入力ベク



図 2: flatabs 関数を活性化関数とした場合の 1 層パーセプト ロンが XOR 問題を識別する様子. 左:学習結果 (メッシュ) と教師信号点 (小球). 右: flatabs 関数の形状.

トル $\mathbf{x} \in \mathbb{R}^{l_i}$ について $\mathbf{y} = \sigma(W\mathbf{x} + \mathbf{b})$ を出力する. ここで, $W \in \mathbb{R}^{l_o} \times \mathbb{R}^{l_i}$ は重み行列, $\mathbf{b} \in \mathbb{R}^{l_o}$ はバイアスベクトル, $\sigma: \mathbb{R} \to \mathbb{R}$ はベクトルの要素毎に適用される活性化関数であ り,フィッティング対象のパラメータは W と b である.線形 分離可能性という制約は、このパーセプトロンの計算過程にお ける2要素に起因している.1個目はアフィン変換Wx+b であり、2個目は(広義)単調な活性化関数 σ である. アフィ ン変換は、入力空間上の各データ点を、ある1直線(あるいは 超平面)からの符号付きの距離に変換する。アフィン変換の効 果を図1に示す。活性化関数σは、入力ベクトルにアフィン 変換を適用したベクトル Wx+bの各要素に適用される.つ まり、アフィン変換後に単調関数を適用する構造のため、いか なる超平面でも二分できない場合は識別不可能である。一方 で、活性化関数を非単調な関数、例えば $\sigma(x) = 1 - e^{-x^2}$ (以 後, flatabs 関数と呼ぶ) に置き換えると,線形分離不可能な XOR 問題であっても、1 層パーセプトロンによって識別可能 である.この様子を図2に示す.ただし、識別性能及び訓練効 率について最適な活性化関数を、事前に設計することは困難で ある. そこで、少なくとも1層パーセプトロンの各素子には入 力空間を超平面で二分する識別能力があることを利用し、まず 1層目(隠れ層)で入力空間を様々な超平面で二分し、次の層 (出力層) ではさらにその層の入力空間を様々な超平面で二分 する, すなわち1層目で様々に二分された部分入力空間の組 み合わせによって,任意の関数を表現可能にしたものが2層 MLP である。2層 MLP が任意の問題を表現可能であること は万能性定理によって裏付けられている。一方で、本質的には 2層のパーセプトロンで十分のはずが、深層学習ではそれを上 回る層数が用いられている。なぜなのだろうか?

多層化の理由の第二は,過適合の抑制である.2層パーセプ トロンの場合,訓練データセットの識別誤差を減らすために1 層目の素子数を増やすと、過適合する傾向がある。これは1層 目で入力空間を様々に二分する際に、各訓練データ点を個別に 表現する様に超平面が構成され、訓練データが存在しない領域 における補間作用が働かなくなるためである。一方,深層学習 では、入力側の層から出力側の層に向かって、入力空間を様々 に二分(1層目)して出来る部分入力空間の組み合わせ(2層 目)により出来る部分入力空間の組み合わせ(3層目)…とい う様に、組み合わせによる表現が強制される。このため、個々 のデータ点という局所への最適化が起こりづらくなる。つま り、接続された各2層は、十分な素子数があれば万能な表現 が可能だが、「各層における部分入力空間の組み合わせで訓練 データ全体を表現せよ」という制約がネットワーク全体に与え られることで、訓練データセットの大域的な傾向を捉えやすく なり,汎化性能が向上するものと考えられる.



図 3: 様々な活性化関数の形状. 左: 図中左上の凡例に示す各 関数の形状. 中央・右:引数がスカラーおよび2次元ベクトル の場合の maxout 関数形状 (紫色)の例.

3. 関連研究

誤差関数の勾配と誤差逆伝播に基づいて学習されるニュー ラルネットワークにおいて,活性化関数に求められる条件は (劣) 微分可能性である.そこで,この条件を満たす様々な関 数が活性化関数として使用されてきた.これまで取り上げて きたシグモイド関数に加え,ソフトマックス関数やハイパーボ リックタンジェントなどは,勾配計算が容易であることから広 く用いられている.他にも,(広義)単調関数としては softplus $\sigma(x) = \log(1 + e^x)$ や ReLU $\sigma(x) = \max(0, x)$ (rectified linear unit) [2] などが提案されており,非単調関数としては放 射基底関数 (radial basis function; RBF)や maxout[3],さ らに sin, cos (以後,これらをまとめてシヌソイドと呼ぶ)な どが使用されてきた.図3左に,スカラーを入力としたとき のこれらの関数の形状を示す.ソフトマックス関数は,基本的 にシグモイド関数と同様の形状である.

ReLUは、softplus を区分線形化したものであり、 ℓ_1 正則化 と併用することでスパースな内部表現を得やすい特徴がある. ReLU と ℓ_1 正則化を組み合わせると、値が 0 となるパラメー タが増えることが期待される.これは必要最低限度の素子数で 識別を行う制約として働くため、汎化性能が向上したものと考 えられる.一方で、ReLU は多層化を前提としており、また機 能しなくなる素子の発生を見込んでネットワークを設計する必 要があり、パラメータ数は増加する.

RBF ネットワークは、2 層 MLP の隠れ層における活性化 関数が放射基底関数であるものを指し、これも任意の連続関 数を近似可能であることが示されている[7].一方で、個々の RBF 素子は入力空間の局所的な特徴を捉えるため、隠れ層の 素子数を増やすことにより訓練データでの識別性能は向上する が、汎化性能の向上については必ずしも期待できない.これは 2 層パーセプトロンと同様である.加えて、訓練データが複雑 な分布である場合、素子数の増加が避けられない.

maxout は, $\sigma(\mathbf{x}) = \max(\mathbf{w}_1^\mathsf{T}\mathbf{x} + b_1, \mathbf{w}_2^\mathsf{T}\mathbf{x} + b_2, \cdots, \mathbf{w}_m^\mathsf{T}\mathbf{x} + b_m)$ ただし $\mathbf{w}_* \in \mathbb{R}^{l_i}$, $b_* \in \mathbb{R}$, $m \in \mathbb{N}$ という関数であり, \mathcal{P} ラメータは $\mathbf{w}_* \ge b_*$ である. つまり, 各 maxout 素子の内部 には m 個のサブモデル (アフィン変換)があり, それらの最 大値を関数出力として採用するものである. maxout の入力が 1 次元の場合と 2 次元の場合の関数形状を図 3 中央および右 に示す. maxout は, 個々のサブモデルの最大値を用いること で,素子の入力空間における区分線形化された凸関数を構成 する. これも,入力空間の局所的な特徴を捉える関数である. め, 本質的に多層化を前提とした活性化関数である.

シヌソイドを活性化関数として用いる 2 層 MLP は,フー リエ級数展開による関数近似と等価であると捉えられる。例え ば,隠れ層の活性化関数が sin である場合,このネットワーク の n 番目の出力は次式により表される:

$$y_n = \mathbf{a}_n^{\dagger} \sin \left(W \mathbf{x} + \mathbf{b} \right)$$

ただし $\mathbf{a}_{\star} \in \mathbb{R}^{h}$, $h \in \mathbb{N}$ (隠れ層の sin 素子数) であり,出力 層の活性化関数は恒等写像としている(出力層のバイアスも省 略). これは、フーリエ正弦級数展開と等価であるため,任意 の奇関数を表現可能である.ある訓練データセットについてパ ラメータ \mathbf{a}_{\star} , W, \mathbf{b} を最適化した場合, \mathbf{a}_{\star} , W, \mathbf{b} は,訓練 データセット全体の周波数領域における振幅,周波数,位相に 対応する.なお, sin と cos の違いは位相のみだが,位相もパ ラメータであるため,仮に sin または cos のどちらか一方のみ を活性化関数として採用したとしても,任意の連続関数を表現 可能である.

4. 表層非線形ネットワーク

少ない層数のニューラルネットワークによって汎化性能の向 上を図るためには、データセットの大域的な特徴を効率的に表 現する必要がある。そこで本稿では、フーリエ級数展開に基づ いているシヌソイドを活性化関数として用いる2層 MLP に 関連して、次に示す3個の手法を組み合わせたアプローチを 提案する:

- 1. 隠れ層の各出力は、各出力に個別のフーリエ級数展開表現を用いる.
- 2. パラメータ最適化に L-BFGS 法 [6] を用いる.
- 3. 周波数パラメータ最適化の際に ℓ2 正則化を行う.

提案手法で用いるネットワークの計算グラフを図4に示す 深 層学習の場合は、入力に近い層から順に、画像で言うエッジな ど学習対象の概念を構成する基本部品(局所的な分布)が学習 され、層が出力に近づくにつれて学習対象の概念(大域的な分 布)が構成的に学習される、と説明される、一方、本アプロー チのネットワークはフーリエ展開に基づき、入力について複数 の重み付けをされた cos 素子が並列していることにより、訓練 データセットにおける局所的な分布は cos の入力の時点で大 きく重みづけられ(高周波成分),大域的な分布は同様に小さ く重みづけられること(低周波成分)で表現される.また、2 層目が全結合ではないことにより、2層目の各素子出力(図中 hidden unit)は個別の非線形関数として機能する。以後、こ の非線形関数の単位を非線形ユニットと呼ぶ。非線形ユニット は、識別器という観点からは特徴量化を担うカーネルの様な役 割を持つ. 従来のシヌソイドを用いた 2 層 MLP では出力層も 全結合であったが、これを排することによって、非線形ユニッ トがそれぞれ個別に最適化される設計とした.



図 4: 本アプローチで用いるネットワーク構造(計算グラフ)



図 5: 活性化関数の周期性と ℓ_2 正則化が, 誤差関数の形状に 与える影響の例. 左: flatabs 関数による活性化. 中央: cos 活 性化・ ℓ_2 正則化なし. 右: cos 活性化・ ℓ_2 正則化あり.

パラメータ最適化については、本アプローチも誤差逆伝播に 基づくが、ニューラルネットワークで一般に用いられる最急降 下法は用いず,L-BFGS 法を用いる.誤差関数は従来のネット ワークと同じ2乗誤差またはクロスエントロピー誤差を用い るが、シヌソイドを活性化関数として用いるネットワークでは 誤差関数の形状が問題となる。1層パーセプトロンににおける XOR 問題を例に、活性化関数として flatabs 関数を用いた場 合と cos を用いた場合の重み行列に関する2 乗誤差関数の形状 を図5(左,中央)に示す(ただし, cosは値域を[0,1]とする ため活性化関数を $\sigma(x) = \frac{\cos(x)+1}{2}$ とし、バイアスはそれぞれ 適当な値に固定した). cosの周期性が誤差関数にも現れるこ とが図に示されている。この例では、どの局所最小値であって も同じ誤差値が得られるが、より複雑な問題では至る所に極小 値が存在する.このため,誤差関数の現在のパラメータ点にお ける最急勾配方向ヘパラメータを移動させる最急降下法では, 妥当な解に収束することが期待できない。一方で、この凹凸が 激しくなるのは cos の引数側のパラメータに関してのみであ る. そこで,前向き計算における cos より前のパラメータ,す なわち周波数パラメータにℓ2正則化項を加えることで、ある 程度誤差関数を滑らかにする。 ℓ2 正則化項を追加した場合の 誤差関数の形状を図5(右)に示す。ℓ2正則化の影響により、 誤差関数は大域的には凸となる. ここで準ニュートン法である L-BFGS 法を用いることで、妥当な解への収束が期待される.

周波数パラメータの ℓ2 正則化には、収束させること以外に も重要な目的がある.画像を例にすると、自然画像をフーリエ 変換した場合、大域的な情報は低周波領域に存在する、高周波 成分は、主にエッジや孤立した点などの、極端に値が変化する 箇所の存在を表している。一般のデータセットにおけるこれら 高周波成分は、外れ値の存在を表すものと考えられる。一方 で、本アプローチでフーリエ級数を利用するのは、大域的な特 徴を効率よく捉えることが目的である。そこで、周波数パラ メータについて lo 正則化することにより、低周波領域でデー タセット全体を近似する. 画像処理における標準画像を2値 化したものを例に、周波数パラメータの l₂ 正則化の効果を図 6に示す.ここでは画像を、各ピクセルの位置座標が入力デー タであり、各座標の輝度が教師信号であるデータセットとみな している. ℓ2 正則化を加えない場合は局所的な特徴を捉えた ことによるノイズが見られるが、 62 正則化を行った場合は大 域的な特徴のみを獲得していることが、図6に示されている。

5. 実験結果

提案手法による識別タスクでの汎化性能とパラメータ数に ついて, MNIST 手書き数字データを用いて, maxout および ReLU (rectifier) と比較した. MNIST は,入力が28×28 画 素のグレースケール画像の各行を並べた784 次元ベクトルで あり,入力ベクトルを数字10 クラスに分類するタスクである.



図 6: 画像をデータセットとして用いて可視化した周波数パラ メータへの ℓ₂ 正則化の効果. 左:元画像(教師信号). 中央: ℓ₂ 正則化なしの学習結果. 右:ℓ₂ 正則化ありの学習結果.

全 70.000 データのうち、訓練データセットは 60.000 データで あり、テストデータセットは残りの10,000 データである。本 アプローチでは事前学習 (pre-training) を行わないため、事 前学習を行わないモデルのみを比較対象として引用する。比較 対象のモデルにおけるパラメータ数は、各論文および公開され ているソースコードより算出した.全ての手法において,出力 層の活性化関数はソフトマックス関数である。この設定におけ る maxout は全3層のモデルであり, rectifier は全4層であ る。また、maxout では正則化に dropout が用いられている。 テストデータセットにおける識別誤り率の比較を表1に示す. 表中、提案手法の名称の後の括弧内の数値は、順に非線形ユ ニット数と各ユニット内の cos 素子数を表す。本手法は、従来 提案されてきた手法と比べて、汎化性能は及ばなかった。一方 で,識別性能に対するパラメータ数の観点からは,各パラメー タが他のモデルよりも識別に大きく寄与していることが裏付け られた.

表 1: MNIST データセットを用いた識別誤り率の比較

| 識別モデル | 識別誤り率 | パラメータ数 |
|--------------|-------|--------|
| maxout[3] | 0.94% | 1,233K |
| rectifier[2] | 1.43% | 3,798K |
| 提案手法(30/10) | 1.79% | 472K |
| 提案手法(20/10) | 1.81% | 157K |
| 提案手法(10/15) | 1.84% | 118K |

6. 考察

ここで、提案手法が汎化性能の点で従来手法に及ばなかった ことについて考察する.提案手法では,周波数パラメータの ℓ2 正則化により汎化性能を得ようとしている. そこで、 ℓ2 正則 化を行う場合と行わない場合とを比較する. これらの2条件に おける識別誤り率とイテレーション回数の関係を図7に示す. この実験では、MNIST の 50,000 件を訓練データセットとし、 残りの 10,000 件ずつをバリデーションデータセットとテスト データとして評価に用いた.図より、 ℓ2 正則化は確かに効果的 だが、学習の中盤(400 イテレーション)以降で l2 正則化だ けでは抑制しきれない過学習が発生していることが判明した. また、学習の進行に伴い、出力層の重み行列と振幅パラメータ のノルムが、訓練データの識別誤り率にのみ関係して増加する 現象を確認している。このため、これらのパラメータについて ℓ2 正則化やノルム正規化 [4] を行ったが,汎化性能は向上しな かった.従って、振幅パラメータあるいは周波数パラメータの 分布に関する他の制約や、周波数領域におけるサンプリング等 の方法を検討する余地があると考察される.



図 7: 訓練時のイテレーション回数(横軸)と識別誤り率(縦 軸・対数目盛)の関係に対する,周波数パラメータへの ℓ_2 正 則化が与える影響. 左: ℓ_2 正則化あり. 右: ℓ_2 正則化なし. ℓ_2 正則化を行わないと汎化作用がほぼ無いことが示されている.

7. おわりに

本稿では、少ない層数のネットワーク構造により、パラメー タ数を減らして前向き計算の効率を上げることと、汎化性能の 向上を目的とした、ニューラルネットワークの新たなアプロー チを提案した.提案手法のネットワークはフーリエ級数展開に 基づく構造であり、データセットをフーリエ変換した際の低周 波成分で近似するという制約を与えることで、データセット の大域的な特徴を得て汎化性能を上げるアプローチをとった. MNIST データセットによる評価実験では、汎化性能は従来手 法を上回らなかったが、パラメータ数という観点からは効率的 に識別されていることを確認した.一方で、低周波成分による 近似だけでは、従来手法による汎化性能を上回ることが難しい ことが判明した.このため、汎化性能向上のために更なる制約 を追加する等の方法について、検討を進める予定である.

参考文献

- Cybenko, G. (1989). Approximation by superpositions of a sigmoidal function. Mathematics of control, signals and systems, 2(4), 303-314.
- [2] Glorot, X., Bordes, A., and Bengio, Y. (2011). Deep sparse rectifier networks. Proceedings of the 14th International Conference on Artificial Intelligence and Statistics. JMLR W&CP Volume. Vol. 15.
- [3] Goodfellow, I., Warde-farley, D., Mirza, M., Courville, A., and Bengio, Y. (2013). Maxout Networks. In Proceedings of the 30th International Conference on Machine Learning (ICML-13), pages 1319-1327.
- [4] Hinton, G. E., Srivastava, N., Krizhevsky, A., Sutskever, I., & Salakhutdinov, R. R. (2012). Improving neural networks by preventing co-adaptation of feature detectors. arXiv preprint arXiv:1207.0580.
- [5] Minsky, M., Papert, S. (1972). Perceptrons: An Introduction to Computational Geometry, The MIT Press, Cambridge MA.
- [6] Liu, D. C., and Nocedal, J. (1989). "On the Limited Memory Method for Large Scale Optimization". Mathematical Programming B 45 (3): 503–528. doi:10.1007/BF01589116.
- [7] Park, J., and Sandberg, I. W. (1991). Universal approximation using radial-basis-function networks. Neural computation, 3(2), 246-257.