

判別分析とパターン認識(2) – ニューラルネットワーク

ニューラルネットワークは、人間の脳細胞と神経回路を模した演算によって、人間の行なう複雑なパターン認識などを実現しようというものです。ニューラルネットワークの特徴のひとつが、その学習能力です。学習とは、ネットワークの入出力例（パターン認識でいえば、文字の画像と読み取り結果など）をネットワークに呈示し、その入出力関係をネットワークが行なうようにネットワークを調節することです。これは、統計学でいえば、標本から入出力関係を推定することと同じです。そこで、今日の講義では、ニューラルネットワークの学習について、その基本的部分を説明します。

ニューロンとニューラルネットワーク

ニューラルネットの研究は、まず神経素子すなわち**ニューロン (neuron)** をモデル化することから始まりました。生理学的研究から、

1. ニューロンは他のニューロンから電氣的・化学的の刺激を受けとる。
2. 受け取る刺激の総和が一定量を超えると「発火」して、他のニューロンに刺激を与える。

というモデルが提示されました。これを図1のような記号で表します。単純化するため、 i 番のニューロンの x_i の状態 (status) を、 $x_i = 1 \rightarrow$ 発火している・ $x_i = 0 \rightarrow$ 発火していないとし、発火しているとき1の量の刺激が他のニューロンへ送られます。また、刺激を受け取る際、神経素子間の結合には重みが設定されており、ひとつのニューロンには、各結合から（重み $\times 1$ または 0 ）の刺激が届きます。すなわち、 j 番のニューロンの状態を x_j とし、 j 番と i 番のニューロン間の結合重みを w_{ji} とすると、 j 番のニューロンから i 番のニューロンに届く刺激の量は $w_{ji}x_j$ となります。そして、その総和 $\sum_j w_{ji}x_j$ が一定量（しきい値 (threshold)）を超えると、この i 番のニューロンが発火します。すなわち、しきい値を T とし、しきい値関数 f を

$$f(x) = \begin{cases} 1 & x \geq T, \\ 0 & x < T. \end{cases} \quad (1)$$

とするとき、 $f(\sum_j w_{ji}x_j) = 1$ ならば i 番のニューロンが発火、すなわち $x_i = 1$ とし、それ以外は $x_i = 0$ とします。

ニューラルネットによる情報処理では、このようなニューロンを多数結合し、ニューロンの初期状態を与えることで情報を入力し、上のような刺激の伝達によってニューロンの状態を変化させて、変化後のニューロンの状態を出力とします。ニューロンの結合の形（トポロジー）には、おおまかにいって**階層型**と**相互結合型**の2種類があります。

図2の(a)が階層型ネットワークで、視神経における情報処理をモデルとして生まれたものです。1つの層にあるニューロン全体で1つの情報を保持し（例えば、ニューロンを画素と考えて視覚的パターンとする）、これがある層から次の層へ一定方向に刺激となって流れてゆきます。そして、最終層に刺激が届いたときの最終層のニューロンの状態が、このニューラルネットワークの出力結果となります。階層型ネットワークは、画像処理やパターン認識に多く用いられます。

一方、図2の(b)が相互結合型ネットワークで、こちらは層はなく、各ニューロンが互いに各ニューロンと結合しています。こちらの場合は、結合重みを一種の「プログラム」として問題に合わせて適切

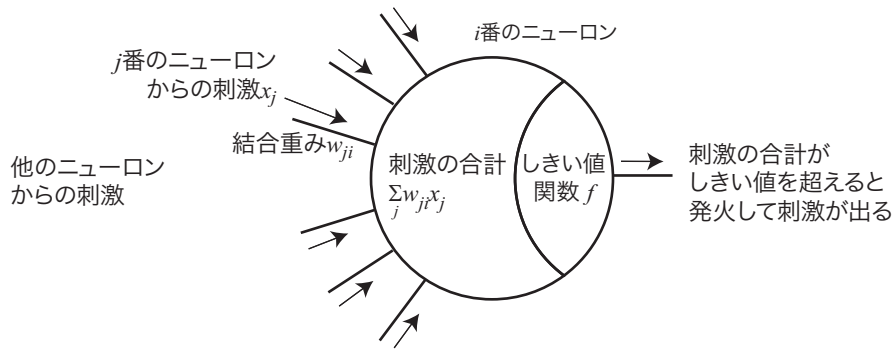


図 1: ニューロンのモデル

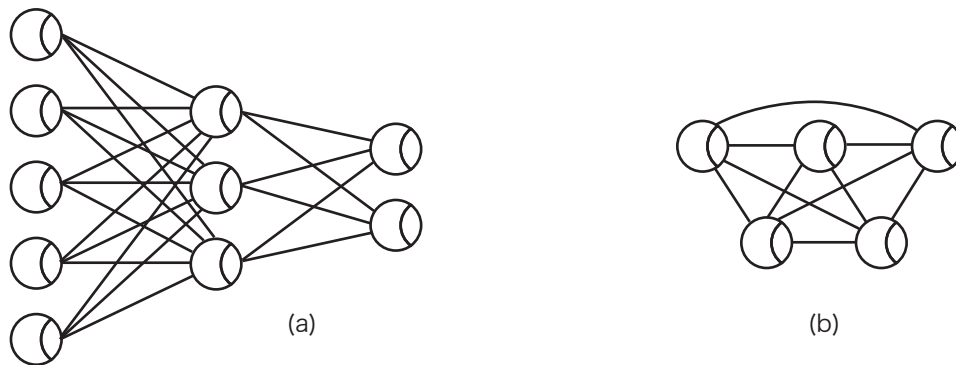


図 2: ニューラルネットワークのトポロジー. (a) 階層型 (b) 相互結合型

に設定し、ニューロンの初期状態を設定した後互いに刺激を交換しあいます。その動作が収束してそれ以上刺激のやりとりが起こらなくなったとき、その時のニューロンの状態を出力結果（あるいは問題の解）とします。相互結合型ネットワークは、最適配置問題を解いたり、連想記憶（記憶したいニューロンの状態に応じて結合重みを設定しておき、ニューロンの初期状態を与えて刺激を交換しあうと、記憶した状態の中で初期状態にもっとも近い状態に収束する。すなわち、初期状態からそれにもっとも近い記憶が「想起」される）として用いられます。

以下、この講義では、パターン認識に関連の深い、階層型ニューラルネットワークを取り扱います。

パーセプトロン学習とその限界

パーセプトロン (perceptron) とは、ニューラルネットワーク研究の初期に提案された、もっとも簡単な階層型ニューラルネットワークで、図3のように入力層、中間層、出力層の3層からなっています。その利用法は、例えば入力層に「A」という文字パターンが入力されると出力層で「A」に対応するニューロンが発火し、以下、「B」のパターンにはBのニューロンが、「C」のパターンにはCのニューロンが、... というようにしてパターン認識を行うというものです。

パーセプトロンをはじめとする階層型ニューラルネットワークの特徴は、結合重みを学習によって定

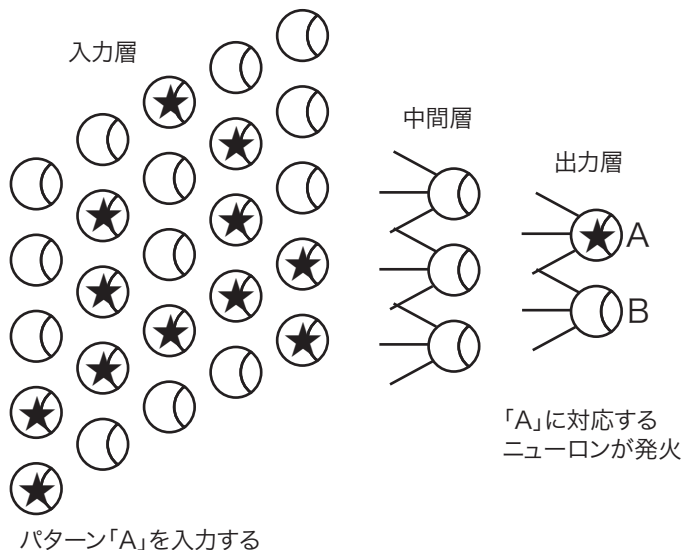


図 3: 階層型ニューラルネットワークによるパターン認識

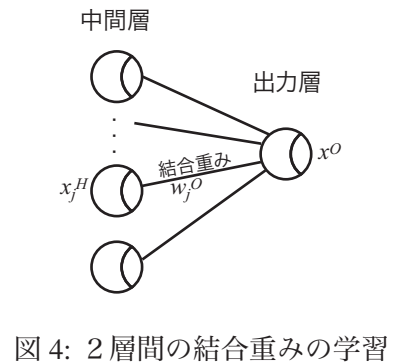


図 4: 2 層間の結合重みの学習

めることです。学習とは、ニューラルネットワークが望みの出力をするようにするための結合重みの調整を、入出力の例のみを用いて行うことです。パーセプトロンにおいては、学習は中間層と出力層の間だけで、以下のような手順で実現されます。以下、簡単のため出力層のニューロンは1つだけでその状態を x^O とします (図 4)。

1. 所望の入出力の例をいくつか用意します。これは例えば、入力層に「A」というパターンを入力したとき出力層のニューロンの状態が「1」になってほしい、といったものです。各入出力例において、入力層に入力例を与えたときの、中間層の j 番目のニューロンの状態を x_j^H とします。
2. 中間層と出力層間の初期結合重み w_j^O を適当に (例えばランダムに) 決めておきます。
3. 入力例を入力層に与え、その時の出力 x^O を調べます。
4. 出力層での所望の出力を y^O とします。もし、所望の出力が $y^O = 1$ なのに現実の出力が $x^O = 0$ である ($x^O - y^O = -1$) ならば、それは出力層のニューロンにやって来る刺激の総和 $\sum_j w_j^O x_j^H$ が足りないからです。そこで、中間層のうち状態 x_j^H が 1 であるようなニューロンとの結合重み w_j^O を増やします。また、もし、所望の出力が $y^O = 0$ なのに現実の出力が $x^O = 1$ である ($x^O - y^O = +1$) ならば、それは出力層のニューロンにやって来る刺激の総和 $\sum_j w_j^O x_j^H$ が多すぎるからなので、中間層のうち状態が 1 であるようなニューロンとの結合重みを減らします。以上のことは、中間層の j 番のニューロンと出力層のニューロンとの結合の新しい重み $w_j'^O$ を

$$w_j'^O = w_j^O - (x^O - y^O)x_j^H. \quad (2)$$

にしたがって更新することで実現できます。

5. 所望の出力が得られるようになるまで、3、4を繰り返します。

この方法をパーセプトロン学習といいます。この方法は簡単ですが、学習アルゴリズムが収束せずに結合重み w_j^O がいくつかの組合せの間を循環してしまう場合があることが知られています。これは、1

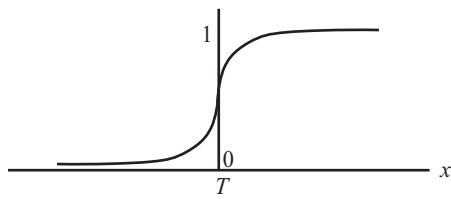


図 5: シグモイド関数

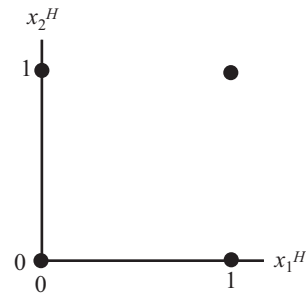


図 6: ニューロンの状態. (x_1^H, x_2^H) の組は図の 4 つの ● のいずれかにある.

回の学習で w_j^O から $w_j'^O$ へ修正する修正量が大きすぎ、急激に結合重みを修正しすぎるからです。そこで、しきい関数 f を、例えば次のシグモイド関数

$$f(x) = \frac{1}{1 + \exp(-\lambda(x - T))}, \quad (3)$$

のような連続関数に置き換えます (図 5, λ はパラメータで, $\lambda \rightarrow \infty$ のとき (3) 式の f はしきい関数になります)。そして、ニューロンの状態 x_j^H や x^O を連続値とし、(2) 式に小さな正の実数 ε (学習係数といいます) を導入して

$$w_j'^O = w_j^O + \varepsilon(y^O - x^O)x_j^H. \quad (4)$$

のように「徐々に」学習する方法があります。これは δ -ルールとよばれています。

さて、パーセプトロンでは、出力層のニューロンの状態が (1) 式で述べたとおり

$$x^O = f\left(\sum_j w_j^O x_j^H\right). \quad (5)$$

となります。この操作は、中間層のニューロンの状態の重み付き線形和を求め、それがしきい値 T より大きいかわりに小さいかで出力層のニューロンの状態を決めることを意味しています。

そこで、中間層のニューロンが x_1^H, x_2^H の 2 つしかない場合を考え、両ニューロンの状態を図 6 の座標平面で表します。このとき、(5) 式の計算は図 7(a) の座標平面上で

$$w_1^O x_1^H + w_2^O x_2^H = T, \quad (6)$$

という直線をひき、現在の中間層のニューロンの状態 x_1^H, x_2^H がこの直線のどちらにあるかによって x^O が 1 か 0 かを決定していることになります。

ところが、例えば「 $x^O = [x_1^H$ と x_2^H の排他的論理和 (XOR)]」という演算では、図 7(b) のように、 x^O が 1 になるような x_1^H と x_2^H の値の領域と、 x^O が 0 になるような x_1^H と x_2^H の値の領域とを、直線で区切ることができません (このような演算を線形非分離といいます)。したがって、どのように学習しても中間層と出力層の 2 層では「 x_1^H と x_2^H の排他的論理和」という演算を得ることはできません¹。パーセプトロンでは入力層と中間層との間の結合重みは固定されていますから、入力層と中間層との間で行う演算の設定によっては、どう学習しても求めるパターン認識ができない、ということが起こります。

¹ 「パーセプトロン学習が収束しない場合」とは、所望の入出力例が線形非分離であるときです。δ-ルールでは、この場合、入出力例との誤差の 2 乗が最小である結合重みに収束します。

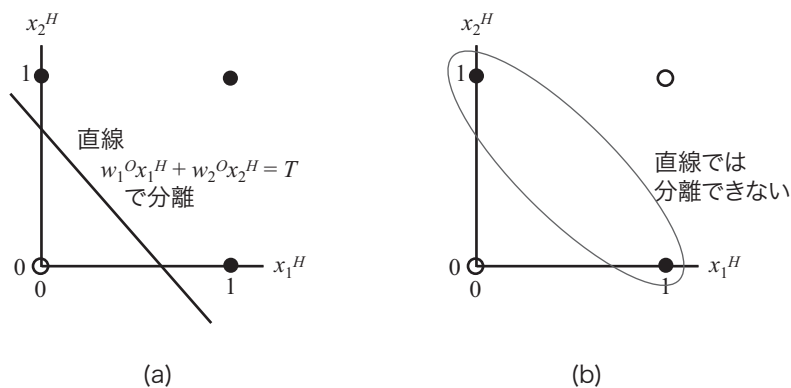


図 7: 線形分離性. (a) ●と○は線形分離. (b) 排他的論理和. ●と○は線形非分離.

これは、パーセプトロンの限界の1つです。これを指摘した Minsky と Papert は、さらに入力層と中間層の間の構成に現実的な制限を課すと表現できない問題がたくさんあることを述べ、中間層と出力層の間の学習だけでは、大したことは学習できないことを示しました。その結果、ニューラルネットワークの研究は一気に下火になってしまったのです。この限界を超えるには、中間層と出力層間だけでなく、ネットワークの層間の結合全体を学習する方法を考える必要があります。それを実現したのが、次章で述べる誤差逆伝播法です。

誤差逆伝播法

前節のパーセプトロン学習や δ -ルールにおいては、「所望の入出力例」は入力層と出力層の状態としてしか与えられていません。ですから、学習によって最小化する「入力例に対する、所望の出力と現実の出力との誤差」は、出力層についてしか得られません。パーセプトロン学習や δ -ルールにおいては、この誤差を使って出力層とその1つ手前の層（全体で n 層あるとすれば、第 $(n-1)$ 層）との結合重みを調整しました。

この方法を使ってさらにもう1つ手前の層（第 $(n-2)$ 層）との結合重みを調整するには、第 $(n-1)$ 層の各ニューロンでの誤差を知る必要があります。そのためには、出力層の1つのニューロンでの誤差を第 $(n-1)$ 層の各ニューロンに割り振らなければなりません。この割り振りを、次の考え方で行います。

1. 第 $(n-1)$ 層の各ニューロンのうち、出力層のこのニューロンとの結合重みが大きいニューロンにより多く誤差を割り振るほうが、割り振られた誤差が修正されたとき、その修正の出力層での効果が大きくなって好都合です。
2. 第 $(n-1)$ 層の各ニューロンの状態は、第 $(n-2)$ 層から送られてくる刺激の総和にシグモイド関数のような非線形関数を適用して決まります。このとき、刺激の総和の変化に対するニューロンの状態の変化率が大きい、つまり「感度のよい」ニューロンにより多く誤差を割り振るほうが、やはり割り振られた誤差の修正の効果が大きくなって好都合です。

これを、一般の第 k 層と第 $(k+1)$ 層について、記号と式を使って書くと、つぎのようになります（図 8 を参照してください）。

いま、第 k 層の j 番目のニューロンが受け持つ誤差の量を δ_j^k とし、第 $(k+1)$ 層の l 番目のニューロンが受け持っている誤差の量を δ_l^{k+1} とします。また、第 k 層の j 番目のニューロンから $(k+1)$ 層の l 番目

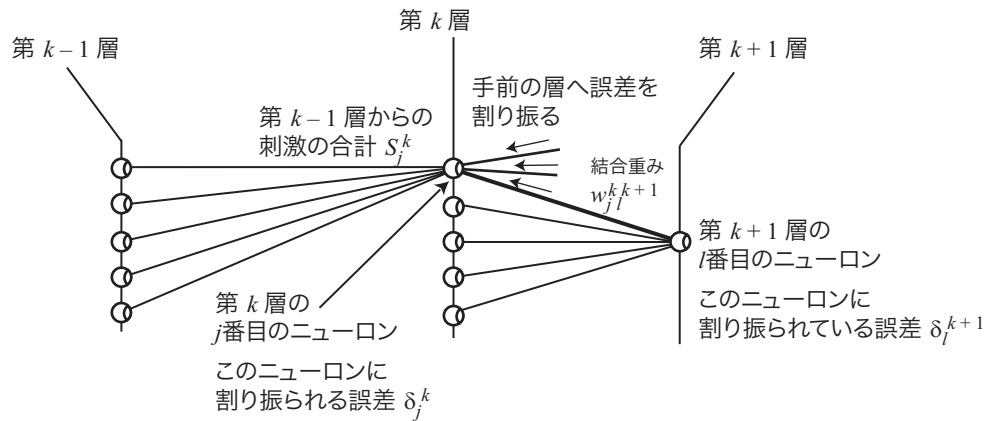


図 8: 誤差逆伝播法

のニューロンへの結合の重みを $w_{jl}^{(k+1)}$ で表します。さらに、第 k 層の j 番目のニューロンの状態を x_j^k で表します。

さて、上の「考え方」1. の条件は、「 δ_j^k のうち第 $(k+1)$ 層の l 番目のニューロンから割り振られる量」が $w_{jl}^{(k+1)}$ に比例することを意味します。また、「第 k 層の j 番目のニューロンに第 $(k-1)$ 層から送られてくる刺激の総和」 S_j^k は

$$S_j^k = \sum_i w_j^{(k-1)k} x_i^{k-1}, \quad (7)$$

と表せますから、「考え方」2. の条件は、「 δ_j^k のうち第 $(k+1)$ 層の l 番目のニューロンから割り振られる量」が、 dx_j^k/dS_j^k にも比例することを意味します。これらの誤差が第 $(k+1)$ 層の各ニューロンから割り振られて、その合計が δ_j^k となるわけですから

$$\delta_j^k = \sum_l \left(w_{jl}^{k+1} \frac{dx_j^k}{dS_j^k} \right) \delta_l^{k+1}. \quad (8)$$

となります。ここで、(しきい関数の代用のシグモイド関数などの) 非線形関数を $f()$ とすると

$$x_j^k = f(S_j^k), \quad (9)$$

ですから、

$$\frac{dx_j^k}{dS_j^k} = f'(S_j^k), \quad (10)$$

となり、(8) 式は

$$\begin{aligned} \delta_j^k &= \sum_l (w_{jl}^{k+1} f'(S_j^k)) \delta_l^{k+1} \\ &= \left[\sum_l (w_{jl}^{k+1} \delta_l^{k+1}) \right] f'(S_j^k). \end{aligned} \quad (11)$$

と表されます。この式によって、第 k 層の j 番目のニューロンが受け持つ誤差の量が、第 k 層の各ニューロンが持つ誤差から決まります。 $k=n$ のとき、すなわち出力層については第 $(k+1)$ 層がありませんが、

出力層での誤差とは所望の出力と現実の出力との差ですから、出力層の l 番目のニューロンにおける所望の出力を y_l とすると

$$\delta_l^n = (x_l^n - y_l) f'(S_l^n), \quad (12)$$

となります。

(11)(12) 式で、各層の各ニューロンが受け持つ誤差が決まりました。そこで、これらの誤差に (4) 式の δ -ルールを適用して、新しい結合重み $w'_{j_l}{}^{k,k+1}$ を

$$w'_{j_l}{}^{k-1,k} = w_{j_l}{}^{k-1,k} - \varepsilon \delta_l^k x_j^{k-1} \quad (13)$$

のように結合の重みを修正します。 ε は (4) 式と同じ学習係数です。

以上の方法で 3 層以上の階層型ニューラルネットワークの学習を行う方法を、**誤差逆伝播法 (error back propagation method, EBP)** といいます。この名前は、出力層での所望の出力と現実の出力との差をもとに、手前の層に誤差を順に伝えていくことから来ています。

さて、ここまでの説明では誤差逆伝播法を「考え方」1., 2. から導きましたが、実はこのようにして導かれた誤差に基づいた (13) 式の結合重みの修正法は、最終層での誤差の 2 乗和がその時点でもっとも多く減る方向に結合重みを修正することが証明できます。このことを、この修正法が**最急降下 (steepest descent)** になっているといえます。最急降下であることの証明は付録に回すとして、ここでは最急降下とはどういう意味かを説明します。

結合重みのある組合せに対して、誤差 2 乗和は 1 通りに定まります。したがって、各結合重みを座標軸とする多次元空間の各点に、誤差 2 乗和の値が決まっているような「場」を考えることができます。結合重みの修正は、誤差 2 乗和がもっとも小さい点をめざして、多次元空間の中のある経路をたどってゆくことに相当します。最急降下とは、現在いる点から誤差 2 乗和がもっとも大きく減少する方向へ進むことを意味しています。

これがどのような方向かを求めてみましょう。 \mathbf{x} を多次元空間の一点とし、 $f(\mathbf{x})$ を各点の関数（上の例なら誤差 2 乗和）とします。ここで、 s をパラメータとして、 $\mathbf{x}(s)$ で、多次元空間中のある経路を表すことにします。

このとき、経路上の $\mathbf{x}(s)$ における勾配を表すベクトルは

$$\mathbf{b} = \frac{d\mathbf{x}}{ds}. \quad (14)$$

で表されます。そこで、この経路の方向に f を微分すると、

$$\frac{\partial f}{\partial s} = \frac{\partial f}{\partial \mathbf{x}} \cdot \frac{d\mathbf{x}}{ds} = \mathbf{b} \cdot \text{grad}f, \quad (15)$$

が得られます。 $\text{grad}f$ は $f(\mathbf{x})$ を各座標軸の方向に微分したベクトルで、たとえば $\mathbf{x} = (x, y, z)$ という 3 次元空間なら、 $\text{grad}f$ は

$$\left(\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z} \right). \quad (16)$$

というベクトルになります。(15) 式によると、経路上の勾配 \mathbf{b} が $\text{grad}f$ に平行で逆方向であるとき、 f の変化率が負で絶対値が最大になります。つまり、もっとも f を減少させる方向、すなわち最急降下の方向は $-\text{grad}f$ であることがわかります。

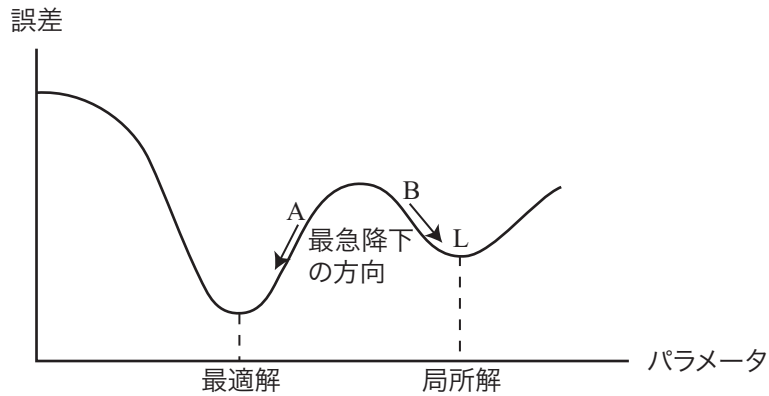


図 9: 最急降下と局所解

最急降下の方向は、 f が現時点でもっとも減少する方向を示しているだけで、**最小値に向かう方向を示しているのではない**ことに注意してください。このため、最急降下法では、「 f の値が、周囲よりは小さいが、遠く離れた場所の値よりは大きい」という**局所解**にしか到達できないことがあります。これを解決するために、最急降下するだけでなく「時々降下しない方向へも進む」などの方法が研究されています。

付録：誤差逆伝播法が最急降下になっていることの証明

ある入出力例における、所望の出力例と現実の出力（第 n 層でのニューロンの状態）の誤差の 2 乗和を

$$E = \frac{1}{2} \sum_l (x_l^n - y_l)^2. \quad (\text{A1})$$

と定義します。このとき、 E は各結合重みの関数です。いま、結合重み $w_j^{k-1 k}$ 方向の E の微分を求めると

$$\frac{\partial E}{\partial w_j^{k-1 k}} = \frac{\partial E}{\partial S_l^k} \cdot \frac{\partial S_l^k}{\partial w_j^{k-1 k}}, \quad (\text{A2})$$

となり、また本文 (7) 式から

$$\frac{\partial S_l^k}{\partial w_j^{k-1 k}} = \frac{\partial}{\partial w_j^{k-1 k}} \sum_i w_i^{k-1 k} x_i^{k-1} = x_j^{k-1}. \quad (\text{A3})$$

となります。そこで

$$\delta_l^k = \frac{\partial E}{\partial S_l^k}, \quad (\text{A4})$$

とにおいて、(A3)(A4) 式を (A2) 式に代入すると

$$\frac{\partial E}{\partial w_j^{k-1 k}} = \delta_l^k x_j^{k-1}. \quad (\text{A5})$$

となります。この式を、結合重みを修正する方法を表す (13) 式、すなわち

$$w_j^{k-1 k} = w_j^{k-1 k} - \epsilon \delta_l^k x_j^{k-1},$$

と比較すると、(13)式は、結合重みを $\partial E / \partial w_{jl}^{k-1}$ 方向、すなわち $-\text{grad}E$ の方向に変化させています。すなわち、(13)式による結合重みの修正は、 E の最急降下になっていることがわかります。

さて、(A4)式の δ は、本文の「考え方」1., 2. から導いた δ と同じものであることを示しましょう。出力層以外、すなわち $k \neq n$ のとき、式 (A4) を

$$\frac{\partial E}{\partial S_j^k} = \sum_l \frac{\partial E}{\partial S_l^{k+1}} \cdot \frac{\partial S_l^{k+1}}{\partial x_j^k} \cdot \frac{\partial x_j^k}{\partial S_j^k} \quad (\text{A6})$$

と変形します。一方、本文の (10) 式から

$$\frac{dx_j^k}{dS_j^k} = f'(S_j^k), \quad (\text{A7})$$

です。また、(7)式の添え字を付け替えると

$$S_l^{k+1} = \sum_i w_{il}^{k+1} x_i^k, \quad (\text{A8})$$

ですから、

$$\frac{\partial S_l^{k+1}}{\partial x_j^k} = w_{jl}^{k+1}. \quad (\text{A9})$$

となります。(A7)式、(A9)式を(A6)式に代入すると

$$\frac{\partial E}{\partial S_j^k} = \sum_l \frac{\partial E}{\partial S_l^{k+1}} \cdot w_{jl}^{k+1} \cdot f'(S_j^k). \quad (\text{A10})$$

となり、(A4)式の δ の定義から

$$\delta_j^k = \frac{\partial E}{\partial S_j^k}, \quad \delta_j^{k+1} = \frac{\partial E}{\partial S_l^{k+1}}. \quad (\text{A11})$$

です。これを(A10)式に代入すると

$$\begin{aligned} \delta_j^k &= \sum_l \delta_j^{k+1} \cdot w_{jl}^{k+1} \cdot f'(S_j^k) \\ &= \left[\sum_l (w_{jl}^{k+1} \delta_j^{k+1}) \right] f'(S_j^k), \end{aligned} \quad (\text{A12})$$

となり、本文の δ の定義と同じになります。

出力層では、すなわち $k = n$ のときは、(A4)式を

$$\frac{\partial E}{\partial S_l^n} = \frac{\partial E}{\partial x_l^n} \cdot \frac{dx_l^n}{dS_l^n}. \quad (\text{A13})$$

と変形します。誤差 2 乗和の定義である (A1) 式から

$$\begin{aligned} \frac{\partial E}{\partial x_l^n} &= \frac{\partial}{\partial x_l^n} \left\{ \frac{1}{2} \sum_l (x_l^n - y_l)^2 \right\} \\ &= x_l^n - y_l \end{aligned} \quad (\text{A14})$$

となります。これを (A13) 式に代入すると

$$\frac{\partial E}{\partial S_l^n} = (x_l^n - y_l) \cdot \frac{dx_l^n}{dS_l^n}. \quad (\text{A15})$$

となり、(A4) 式の δ の定義から

$$\delta_l^n = \frac{\partial E}{\partial S_l^n}, \quad (\text{A16})$$

で、また (A7) 式から

$$\frac{dx_l^n}{dS_l^n} = f'(S_l^n), \quad (\text{A17})$$

ですから、(A16) 式、(A17) 式を (A15) 式に代入すると

$$\delta_l^n = (x_l^n - y_l) f'(S_l^n), \quad (\text{A18})$$

となり、本文の (12) 式の δ の定義と同じになります。