

目次

- ★ 12.3 ロジスティック回帰
- ★ 13.1 ニューラルネットワークって？
- ★ 13.2 多層パーセプトロン

★ 12 パターン認識と機械学習 (3) — 識別のための教師あり学習 (承前)

★ 12.3 ロジスティック回帰

識別のための手法の別の例として、ロジスティック回帰 (☆1) を紹介する。名前に「回帰」とあって紛らわしいが、「識別」のための手法である (☆2)。

★ 12.3.1 2 クラス識別問題の場合の定式化

簡単のため、まずは識別すべきクラスの数に 2 に限定された場合を考える。2 つのクラスのうち一方を ‘positive’ (正) クラス、他方を ‘negative’ (負) クラスと呼ぶことにする。学習データは $\{(\mathbf{x}_n, y_n) | n = 1, 2, \dots, N\}$ という形とする。 \mathbf{x}_n は D 次元のデータ (特徴ベクトル) である。 y_n は、 \mathbf{x}_n の所属クラスの正解を表し、 \mathbf{x}_n が positive クラスに属すべきものなら $y_n = 1$ 、さもなければ $y_n = 0$ とする。

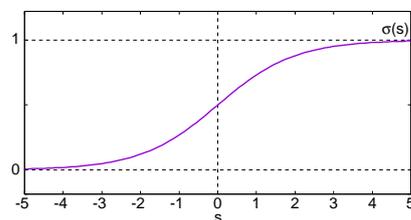
ここで、 $(D + 1)$ 個のパラメータ w_0, w_1, \dots, w_D を持つ次のような関数を考える。

$$f(\mathbf{x}; w_0, w_1, \dots, w_D) = \frac{1}{1 + \exp\left(-\left(w_0 + \sum_{d=1}^D w_d x_d\right)\right)} \quad (1)$$

$s = w_0 + \sum_{d=1}^D w_d x_d = w_0 + \mathbf{w}^\top \mathbf{x}$ とおくと、式 (1) は

$$\sigma(s) = \frac{1}{1 + \exp(-s)} \quad (2)$$

という形をしている。この関数 $\sigma(s)$ のことを **シグモイド関数** (☆3) という。式の形と下のグラフからわかるように、この式は任意の s に対して 0 より大きく 1 より小さい値をとる。そこで、式 (1) の $f(\mathbf{x})$ の値が、「 \mathbf{x} は positive クラスに所属するものである」ということの意味を確信度を表すと考えることにする。



このように考えると、学習データのうち $y_n = 1$ であるものについては positive クラス所属なのだから $f(\mathbf{x}_n)$ が 1 に近くなるように、 $y_n = 0$ であるものについては逆に $f(\mathbf{x}_n)$ が 0 に近くなるようにうまくパラメータを調節すれば、式 (1) の値で 2 クラスの識別ができそうである。未知のデータ \mathbf{x} に対しても、例えば $f(\mathbf{x}) > \frac{1}{2}$ なら positive クラス、さもなければ negative クラスと判断すればよい。

☆1) ロジスティック回帰: logistic regression.

☆2) この手法の背景には、その他の多くの機械学習手法と同様に、データ等の確率的・統計的性質を考慮した問題設定や議論があるのだが、この授業では説明を省略する。

☆3) シグモイド関数: sigmoid function. 次回話題であるニューラルネットワークでも登場します。

そこで、学習データからパラメータ w_0, w_1, \dots, w_D を定めるために、パラメータの「悪さ」を表す関数を定義する。ロジスティック回帰では、次式で表される**交差エントロピー** (☆4) を用いる。

$$H(w_0, w_1, \dots, w_D) = \sum_{n=1}^N h_n \quad (3)$$

$$h_n = -y_n \log z_n - (1 - y_n) \log(1 - z_n) \quad (4)$$

ただし、 $z_n = f(\mathbf{x}_n)$ である。この交差エントロピー H がなるべく小さくなるようなパラメータが「良い」と考えて、 H を最小化するパラメータを探す。これが、2クラス問題の場合のロジスティック回帰の考え方である。

★ 12.3.2 勾配法によるパラメータの逐次修正

上記のような問題設定は、最小二乗法による回帰の場合にも考えた。最小二乗法の場合、学習データに対する二乗誤差の和で誤差関数を定義し、それを最小にするパラメータを求めたのだった。このように、最小二乗法もロジスティック回帰も、目的とする関数（二乗誤差や交差エントロピー）を最小化する解（特定のパラメータの値）を探す問題である、というところは共通である。このような問題は、**最適化問題** (☆5) と呼ばれる。

最小二乗法の場合、最適化問題を解くことは比較的容易である。誤差関数をパラメータについて微分して $\mathbf{0}$ とおくと線形の連立方程式が得られるので、それを解けばよいのだった。しかし、ロジスティック回帰の場合、 H の微分を $\mathbf{0}$ とおいて得られる連立方程式は非線形であり、簡単に解くことはできない。それでも、ロジスティック回帰の場合、 H のパラメータに関する微分が求まるので、それを利用した最適化手法である**勾配法** (☆6) が使える。ここでは、その最も単純な方法の一つである**最急降下法** (☆7) を説明し、これによってパラメータを逐次修正していく学習アルゴリズムを構成できることを示す。

2クラス問題のロジスティック回帰における $(D+1)$ 個のパラメータのうちの一つを w と表して考える。 H を最小にするパラメータ w の値を求めるための最急降下法の手順は次のようになる。

1. $w(t)$ の初期値 ($t=0$ における値) を適当に定める。
2. ステップ t における値 $w(t)$ を用いて微係数 $\left. \frac{\partial H}{\partial w} \right|_{w=w(t)}$ を計算し、次のステップ $t+1$ における値 $w(t+1)$ を次式により求める。

$$\begin{cases} w(t+1) = w(t) + \Delta w \\ \Delta w = -\eta \left. \frac{\partial H}{\partial w} \right|_{w=w(t)} \end{cases} \quad (5)$$

Δw は w の修正量を表す。 η は正の小さな定数である。

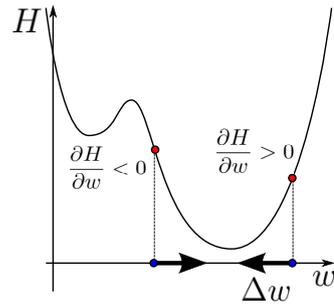
3. H の値が十分小さくなっていれば終了、さもなければ 2. を繰り返す。

☆4) 交差エントロピー: cross entropy. これが何者かは、この授業では説明しません。パターン認識や情報理論を勉強するとわかるかも。

☆5) 最適化: optimization. 負号を付けるだけなので、最大化でも同じこと。機械学習、パターン認識、コンピュータビジョンではよく出てくるが、それ以外の幅広い分野で頻出の問題である。

☆6) 勾配法: gradient method.

☆7) 最急降下法: steepest descent method.



上図からわかるように、最急降下法は「現在地での傾きを調べ、下り方向にちょっとだけ進む」ことを繰り返す手法である。初期値によっては極小解に陥って最小解にたどり着けないこともある。

上記のロジスティック回帰の場合、微係数を具体的に計算してみると、

$$\frac{\partial H}{\partial w_d} = \sum_{n=1}^N (z_n - y_n) x_{n,d} \quad (d = 0, 1, \dots, D) \quad (6)$$

となる。ただし、 $x_{n,d}$ はベクトル x_n の d 番目の要素 ($d = 0$ のときは 1 とみなす) である。これを式 (5) に当てはめれば、学習アルゴリズムが得られる。

以下の Q は、式 (6) を導出するためのヒントである。

Q1. シグモイド関数の微分は、シグモイド関数自身を用いて

$$\frac{d\sigma(s)}{ds} = \sigma(s) \times \text{hoge} \quad (7)$$

と表せる。ただし、hoge は $\sigma(s)$ を用いた式である。hoge を求めなさい。

Q2. 上記の結果を利用して、 $d = 1, 2, \dots, D$ の場合の

$$\frac{\partial z_n}{\partial w_d} = \frac{\partial}{\partial w_d} \sigma \left(- \left(w_0 + \sum_{d=1}^D w_d x_{n,d} \right) \right) \quad (8)$$

を $x_{n,d}$ と z_n のみを用いた式で表しなさい ($z_n = \sigma(\dots)$, 合成関数の微分…).

Q3. 上記の結果を利用して、 $d = 1, 2, \dots, D$ の場合の

$$\frac{\partial}{\partial w_d} (y_n \log z_n) = y_n \frac{\partial}{\partial w_d} \log z_n \quad (9)$$

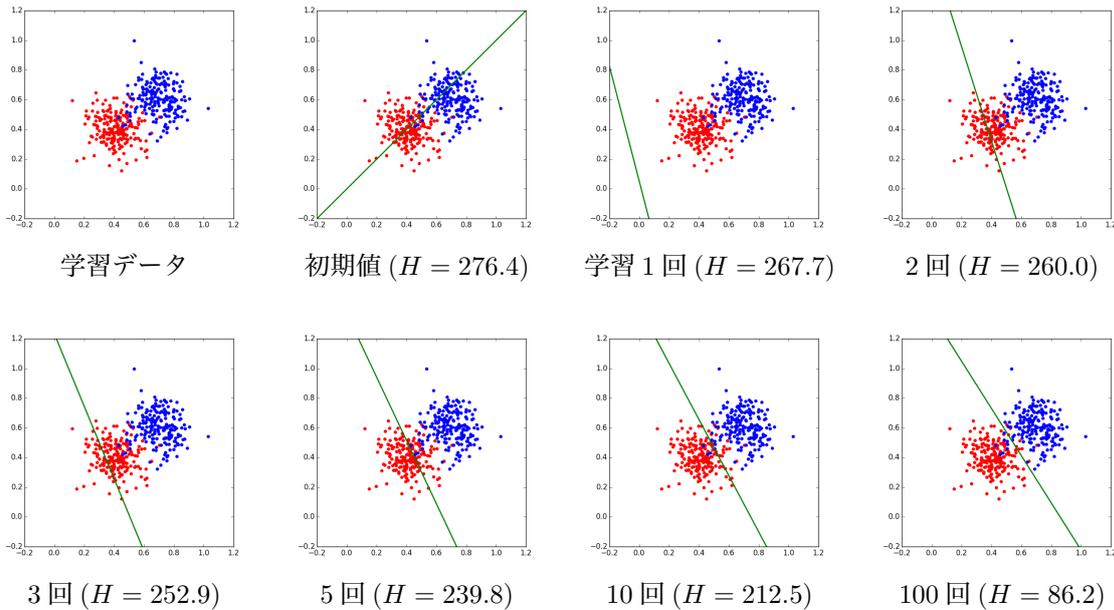
を $x_{n,d}, y_n$ と z_n のみを用いた式で表しなさい (やばし合成関数の微分…).

Q4. 上記の結果を利用して、 $d = 1, 2, \dots, D$ の場合の $\frac{\partial h_n}{\partial w_d}$ を $x_{n,d}, y_n$ と z_n のみを用いた式で表しなさい。また、 $x_{n,0} = 1$ と考えると、その式が $d = 0$ の場合にも当てはまることを確かめなさい。

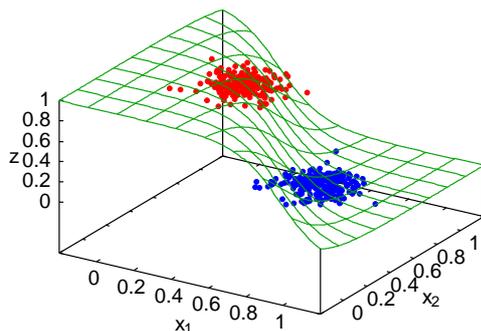
Q5. 上記の結果を利用して、式 (6) が成り立つことを示しなさい。

★ 12.3.3 2 クラスデータのロジスティック回帰の例

乱数を使って 2 次元の人工データを生成し、ロジスティック回帰によってそれらを 2 クラスに識別する実験を行った結果を示す。以下の図の赤と青の点は 2 つのクラスの学習データを表している。緑色は、 $f(x_1, x_2; w_0, w_1, w_2) = \frac{1}{2}$ を満たす点 (x_1, x_2) の集合、すなわち赤クラスと青クラスの識別境界を表している。式 (1) からわかるようにこれらの点は $w_0 + w_1x_1 + w_2x_2 = 0$ を満たすので、境界は直線である。



以下の図は、100 回学習後の $z = f(x_1, x_2)$ の表す曲面を可視化したものである。赤クラス青クラスのデータを、それぞれ $z = 1$ および $z = 0$ の平面上に重ねて表示してある。



★ 12.3.4 クラス数が 3 以上の場合のロジスティック回帰

上述の定式化では 2 クラス問題しか扱うことができないが、3 クラス以上の識別ができるように拡張することは容易である。ただし、紙面と時間の都合で、具体的な定式化や実験結果については省略する。

★ 13 パターン認識と機械学習 (4) — ニューラルネットワーク

ニューラルネットワークは、元々は脳における情報処理のモデルとして提案されたものであるが、現在では機械学習の手法として幅広く応用されている。

★ 13.1 ニューラルネットワークって？

★ 13.1.1 神経細胞とそのモデル

ヒトの脳には数百億のニューロン（神経細胞, neuron）が存在している。ニューロンは信号を伝達する機能をもっており、多数のニューロンが相互につながりあって多様な情報処理を行なっている。生命維持から感情、思考にいたるまでの脳機能は、主にこれらニューロンの集団が担っているものと考えられている。

ニューロンのふるまいを思い切って単純化すると、次のようにまとめることができる。

1. 他のニューロンの出力を受け取る。ただし、そのまま受け取るのではなく、ニューロン間のつながりの強さに応じて重みづけされた値を受け取る。
2. それらの和を求める
3. その値に応じて自身の出力を決める

このようなニューロンのふるまいは、次式のようにモデル化できる (☆ 8)。

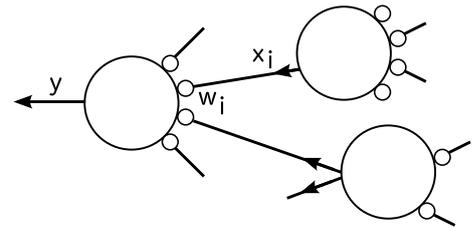
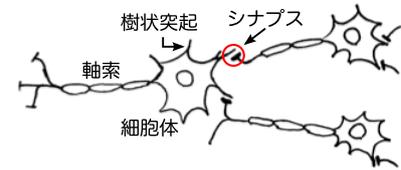
$$y = \sigma \left(\sum_{i=1}^n w_i x_i + \theta \right) \quad (10)$$

ただし、 x_i はこのニューロンに信号を伝達する n 個のニューロンのうち i 番目のものの出力であり、 w_i は x_i の結合重み (☆ 9) を表すパラメータである。また、 θ はしきい値 (☆ 10) と呼ばれるパラメータである。関数 $\sigma(s)$ は活性化関数 (☆ 11) と呼ばれるものであり、以下に示すステップ関数 (☆ 12) やシグモイド関数 (☆ 13)、Rectified Linear 関数 (☆ 14) などがよく用いられる。ニューロンの入力と出力の関係は一般に非線形であり、これらはその性質をモデル化したものとなっている。

$$\text{ステップ関数} \quad \sigma(s) = \begin{cases} 1 & s \geq 0 \\ 0 & \text{otherwise} \end{cases} \quad (11)$$

$$\text{シグモイド関数} \quad \sigma(s) = \frac{1}{1 + e^{-s}} \quad (12)$$

$$\text{Rectified Linear} \quad \sigma(s) = \begin{cases} s & s \geq 0 \\ 0 & \text{otherwise} \end{cases} \quad (13)$$



☆ 8) 1940 年代に McCulloch と Pitts が提案した。ただし、彼らのモデルではニューロンの出力は 0, 1 の二値である。

☆ 9) 結合重み: connection weight.

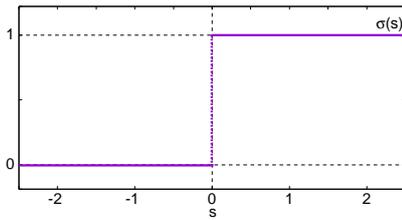
☆ 10) しきい値: threshold value.

☆ 11) 活性化関数: activation function.

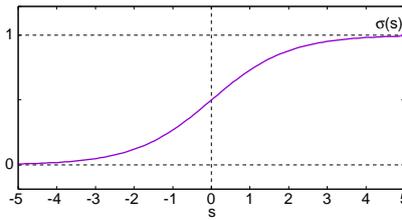
☆ 12) ステップ関数: step function.

☆ 13) シグモイド関数: sigmoid function.

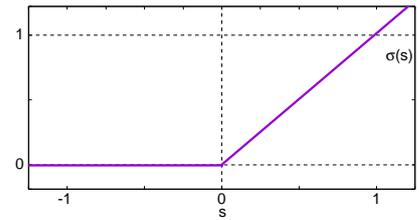
☆ 14) ランプ関数ということもある。



ステップ



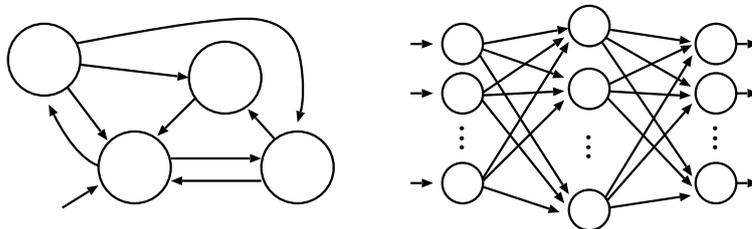
シグモイド



Rectified Linear

★ 13.1.2 ニューラルネットワーク, 多層パーセプトロン

上記のようなニューロンモデルを複数つなぎあわせたものを**ニューラルネットワーク** (☆ 15) という。様々なつなぎ方ものが考えられるが、図右のようにニューロンが層を成し、層間のニューロンのみに単方向のつながりがあるタイプものを**多層パーセプトロン** (☆ 16) と呼ぶ。



☆ 15) ニューラルネットワーク: neural network, 神経回路網とも。

☆ 16) 多層パーセプトロン: Multi-Layer Perceptron, MLP.

★ 13.1.3 深層学習

ニューラルネットワークに関しては数十年前から地道な研究が続けられてきたが、近年、その進歩とコンピュータの高性能化が相まって、いわゆる**人工知能** (☆ 17) の一種として急速に発展し、様々な分野で応用されるようになってきている (☆ 18)。多くの場合、上述の多層パーセプトロンのように階層的な構造をしたニューラルネットワークで、たくさんの層を積み重ねた (したがって学習するパラメータの数も非常に多い) ものが用いられる。このような多層のニューラルネットワークの学習は、**深層学習** (☆ 19) と呼ばれている。次頁以降でその概略を説明する。

☆ 17) 人工知能: Artificial Intelligence.

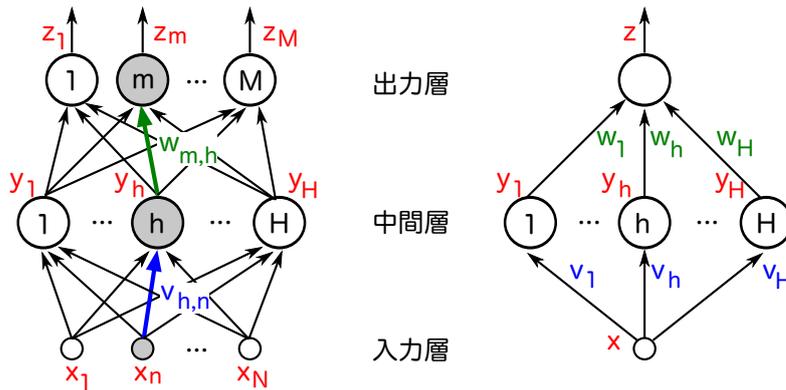
☆ 18) 画像認識, 自動運転, 音声認識, 自動翻訳, 将棋や囲碁でプロに勝つ AI, etc.

☆ 19) 深層学習: deep learning.

★ 13.2 多層パーセプトロン

★ 13.2.1 多層パーセプトロンとその学習

下図に、2種類のニューラルネットを示す。左右の図のいずれも、大きめの丸一つが一つのニューロンを表している。左のものは入力と出力とつながった構造をしている。一方、右の方は、入力と出力の間に「隠れた」ニューロンの層（これを「隠れ層」または「中間層」という）を有している。一般に、「多層」パーセプトロンという時は、このような隠れ層を1層以上有するものを指す。



図左の多層パーセプトロン（以下 MLP と略記する）の入出力は、次のような式で表される。

$$y_h = \sigma \left(v_{h,0} + \sum_{d=1}^D v_{h,d} x_d \right) \quad (h = 1, 2, \dots, H) \quad (14)$$

$$z_m = \sigma \left(w_{m,0} + \sum_{h=1}^H w_{m,h} y_h \right) \quad (m = 1, 2, \dots, M) \quad (15)$$

ここで、 y_h は隠れ層の h 番目のニューロンの値であり、 $v_{h,d}$ はこのニューロンと入力の d 番目の要素との間の結合重みである。同様に、 z_m は出力層の m 番目のニューロンの値であり、 $w_{m,h}$ はこのニューロンと隠れ層の h 番目のニューロンとの間の結合重みである。関数 $\sigma(s)$ は前述の活性化関数である。

MLP の特徴は、上述のように隠れ層を有することである。隠れ層ニューロンの活性化関数にシグモイドのように非線形なものを用いると、MLP 全体の入出力も非線形関数となる。そのため、入力と出力の間に複雑な関係があるようなデータの場合でも、うまく学習できると期待される。

MLP においては、ニューロン間の結合の強さを表す値が、学習すべきパラメータとなる。式 (14) と (15) で表される MLP の場合、入力-隠れ層のニューロン間の結合を表す値 $v_{h,d}$ と、隠れ層-出力層のニューロン間の結合を表す値 $w_{m,h}$ がパラメータである。それらの学習には、ロジスティック回帰の場合と同様に、勾配法/最急降下法を用いることが多い。目的関数としては、出力の正解と実際の出力との間の二乗誤差（最小二乗法で説明した誤差関数と同じ形のもの、次節も参照）や交差エントロピー（前回資料参照）などを用いることができる。いずれの場合も、上記パラメータに関する微分が計算できるので、パラメータを適当な初期値から逐次修正していく形の学習アルゴリズムを構成することができる。

★ 13.2.2 多層パーセプトロンの応用例

多層パーセプトロン (MLP) は、教師データと目的関数の与え方次第で、回帰／識別どちらの問題にも適用することができる。回帰の場合、通常は二乗誤差を目的関数とする。一方、識別の場合、交差エントロピーを用いることが多い。以下、回帰と識別それぞれの応用例を示す。

非線形回帰 最も単純な形の MLP として、上の図および式 (14) と (15) で $D = M = 1$ とした場合、つまり入力も出力も 1 つの場合を考える。出力層ニューロンの活性化関数は恒等関数 ($\sigma(s) = s$) とする (☆ 20) このとき、入力 x に対するこの MLP の出力 z は、次式のようになる。

$$y_h = s(v_{h,0} + v_{h,1}x) \quad (h = 1, 2, \dots, H) \quad (16)$$

$$z = \sum_{h=1}^H w_h y_h \quad (17)$$

☆ 20) シグモイドと違って出力が (0, 1) に限定されないの
で、この問題のような回帰／
関数近似では出力の活性化関
数としてよく用いられる。

この MLP に対して、入力と出力の正解のペア N 個から成る学習データ $\{(x_n, \tilde{z}_n) | n = 1, 2, \dots, N\}$ を与え、 x_n に対する MLP 出力 z_n と正解 \tilde{z}_n との間の二乗誤差の和

$$E = \frac{1}{2} \sum_{n=1}^N (\tilde{z}_n - z_n)^2 \quad (18)$$

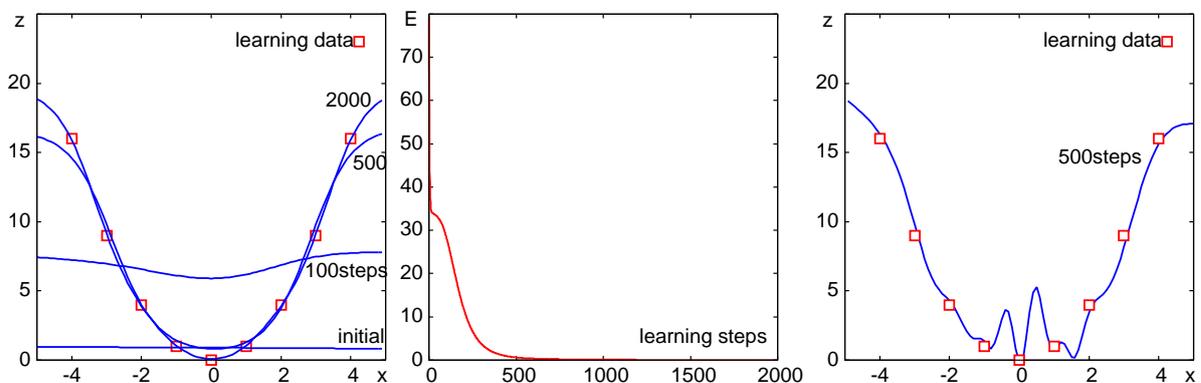
を目的関数として学習させる。これは、入出力がともに 1 変数の場合の回帰問題への MLP の適用例となっている。MLP を用いると、非線形な活性化関数をもった隠れ層のおかげで複雑な曲線を当てはめることができる。

以下に、実際に学習を行なった実験の結果を示す。二次関数 $z = x^2$ を近似することを目標に、学習データを $(-4, 16), (-3, 9), \dots, (4, 16)$ の $N = 9$ 個とした。また、中間層のニューロン数は $H = 10$ と $H = 100$ の二通りとした。

下図左は、 $H = 10$ の場合の学習過程における MLP 出力の変化を示している。学習が進むにつれてうまく学習データを近似できるようになっていることがわかる。このことは、下図中に示した二乗誤差 E の変化の様子からもわかる。

一方、下図右は、 $H = 100$ とした場合の学習結果の一例を示している。この場合、学習データに対する二乗誤差 E はほぼ 0 になっており、MLP 出力の曲線は全ての学習データ点を通っているが、その形は近似対象である放物線とはほど遠いものとなっている。このような現象は**過適合・過学習** (☆ 21) と呼ばれている。

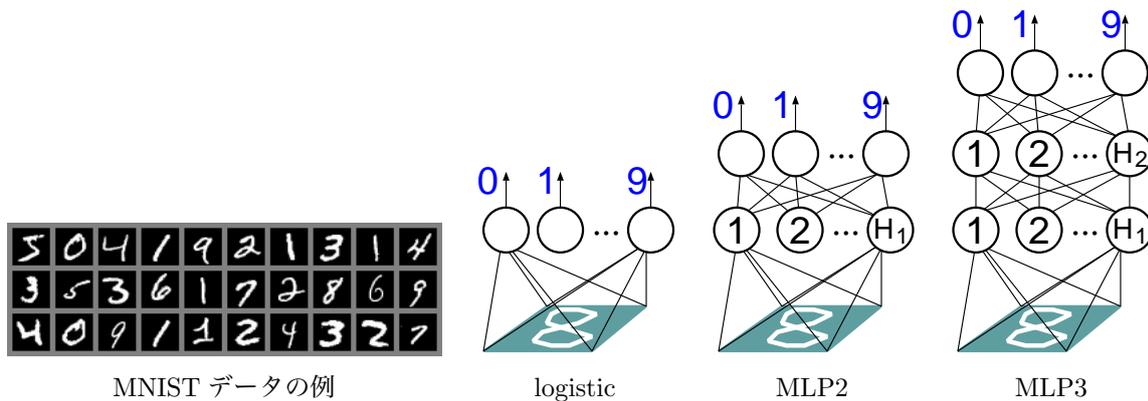
☆ 21) 過適合: overfitting.



手書き数字の識別 多層パーセプトロン (MLP) を識別問題に用いる例として, MNIST 手書き数字データセット で実験を行ってみる. このデータセットは, 28×28 画素の '0' から '9' までの手書き数字のグレイスケール画像 7 万枚 (学習用 6 万枚+テスト用 1 万枚) から成るものであり, 機械学習の練習問題としてよく用いられている. 入力の次元数 $D = 28 \times 28 = 784$ であり, クラス数 $K = 10$ である.

☆ 22) <http://yann.lecun.com/exdb/mnist/>

ここでは, ロジスティック回帰 (logistic と表記, 前回登場のものをクラス数が 3 以上に一般化したもの), それに隠れ層を 1 つ追加した MLP (MLP2 と表記), さらにもう 1 つ追加した MLP (MLP3) の 3 通り (下図参照) で実験を行った. 隠れ層のニューロン数は $H_1 = 500$ (MLP2) および $H_1 = H_2 = 500$ (MLP3) とし, 隠れ層ニューロンの活性化関数は全て Rectified Linear とした.



左下の表に, 3 万回の学習 (詳しい条件は説明を省略する) を行った後の MLP を用いて測った誤識別率 (MLP が出力したクラスが誤りだった割合) を示す. ロジスティック回帰も MLP も, パラメータの初期値が異なれば異なる結果が得られるので, ここでは初期値を 5 通りずつ変えて得られた誤識別率の平均を示している. また, 右下のグラフに学習の進行の様子を示す. 横軸は学習回数, 縦軸は学習データに対する交差エントロピーの平均である.

	logistic	MLP2	MLP3
学習データの誤識別率 [%]	6.9	0.036	0.0
テストデータの誤識別率 [%]	7.5	1.9	1.9

