

# A I 基礎セミナー

## 第8回

### 機械学習（2回目：機械学習の概要と教師あり学習（分類））

#### 改訂履歴

日付	担当者	内容
2021/05/08	M. Takeda	Git 公開

#### 目次

(1) はじめに

(2) 機械学習の手法

(2.1) 教師あり学習

(2.2) 教師なし学習

(2.3) 強化学習

(2.4) 転移学習

### (3) 脳神経系のモデル化

- (3.1) 神経細胞間の信号伝達の様子
- (3.2) 神経細胞間の信号伝達のモデル化
- (3.3) ニューラルネットワークと適用モデルの選択

### (4) 前処理

- (4.1) データを機械学習モデルに入力するまでの流れ
- (4.2) データの種類について
- (4.3) 何を入力データとするか
- (4.4) 入力データのクリーニング
- (4.5) 入力データの欠損値の扱い
- (4.6) 入力データの標準化
- (4.7) 入力データの加工による作成
- (4.8) 入力データの符号化

### (5) 活性化関数

### (6) 損失関数

- (6.1) 損失関数・誤差関数・コスト関数・目的関数
- (6.2) 尤度と損失関数
- (6.3) 損失関数一覧

### (7) モデルの構造

- (7.1) 順伝播型ニューラルネットワーク
- (7.2) 畳み込みニューラルネットワーク
  - (7.2.1) 畳み込み層
  - (7.2.2) プールング層

## (8) モデルの学習

### (8.1) 勾配降下法と学習率

### (8.2) 勾配降下法の改良版

#### (8.2.1) 確率的勾配降下法

#### (8.2.2) AdaGrad

#### (8.2.3) Adam

### (8.3) 誤差逆伝搬法

### (8.4) 「順伝播型ニューラルネットワーク」の誤差逆伝搬法

#### (8.4.1) 各層のユニットの誤差

#### (8.4.2) ユニットの誤差とコスト関数の偏微分との関係式

#### (8.4.3) 出力層のユニットの誤差

#### (8.4.4) 隣接する層間のユニットの誤差についての逆漸化式

#### (8.4.5) 全ての中間層のユニットの誤差 $\delta$ を計算する

#### (8.4.6) ユニットの誤差 $\delta$ を用いて勾配降下法で学習する流れ（例）

### (8.5) モデルの学習の検証

#### (8.5.1) モデルの学習の検証方法

#### (8.5.2) 過学習の対策とドロップアウト

## (9) まとめ

## (10) 確認問題

## (11) 確認問題回答用紙

## (1) はじめに

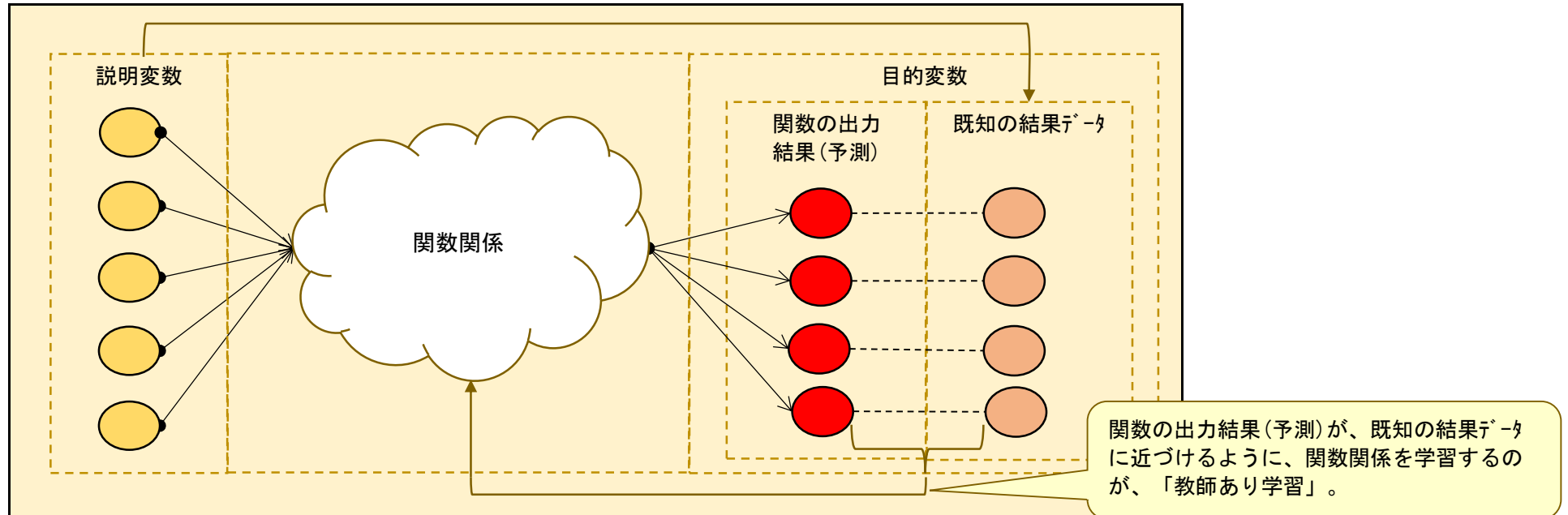
- ・「セミナー7回目（機械学習 1回目）」では、機械学習による画像認識を実装例を示して紹介しましたが、「セミナー8回目（機械学習 2回目）」では、機械学習の基本概念と「教師あり学習（分類）」についての解説を行います。
- ・今回のセミナーは、「セミナー7回目」で扱ったモデルとその学習方法の実装や他の「教師あり学習（分類）」の文献を見て、内容が手に取って読めるようになる、ということを目指します。
- ・確認問題は本文を要約したような問題となっております。  
解いてみて理解度を確認してみてください。

## (2) 機械学習の手法

- ・「機械学習（マシンラーニング, machine learning, ML）」とは、人間の学習の仕組みを機械（特にコンピュータ）で実現する技術・手法の総称です。機械学習では、大量のデータ間の統計的関係性をモデルとして定式化して、そのパラメータを訓練用データで学習することにより推定します。そしてそのモデルを予測や分類といった問題に適用します。機械学習の手法は、「教師あり学習」「教師なし学習」「強化学習」の3種類に分類できます。

### (2.1) 教師あり学習

- ・ 何かの因果関係を探る場合、統計モデルでは、原因となる変数群を元に結果となる変数群を説明しようと試みます。原因となる変数のことを「説明変数（セツメイヘンズウ、explanatory variable）」と言い、結果となる変数のことを「目的変数（モクテキヘンズウ、response variable）」と言います。
- ・ 「教師あり学習（キョウシアリガクシュウ, Supervised learning）」は 異なるデータ間の関係性を、一方を入力（説明変数）、他方を出力（目的変数）とする関数関係として記述し、関数の出力を既知の結果データ（目的変数の既知データ）に近づけるように、関数関係を学習する方法です。この場合、関数の出力結果（予測）に対して、既知の結果データが教師データの役割をします（下図参照）。



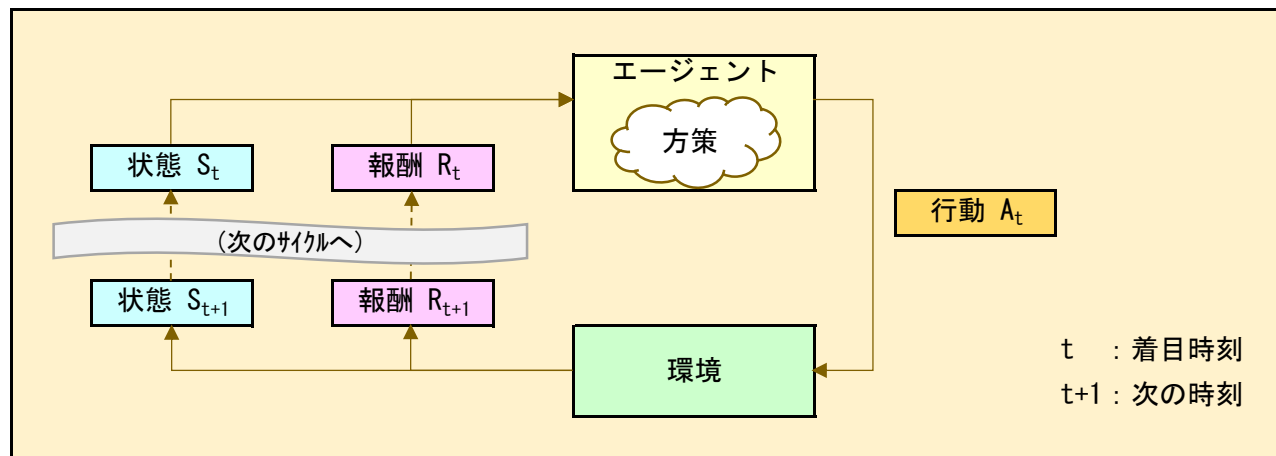
- ・教師あり学習において、目的変数に対応する確率変数が連続確率変数の場合、関数関係を推定する学習を「回帰 (カク, Regression)」と呼びます。  
「回帰」により分析することを「回帰分析 (カクブンセキ, Regression analysis)」と呼びます。  
「回帰」では、観測データをもとに入力データ (説明変数) と出力データ (目的変数) の数値間の関連性を導き出し、入力データに対する出力データを数値で予測します。
- ・教師あり学習において、目的変数に対応する確率変数が離散確率変数の場合、関数関係を推定する学習を「分類 (ブンレイ, Classification)」と呼びます。  
「分類」での関数の出力結果は、各ラベルが付与される確率を与えるものであり、確率が最大となるラベルが正解ラベル (教師データ) と一致するように学習します。

## (2.2) 教師なし学習

- ・ 入力変数（説明変数）に対する、出力変数（目的変数）が与えられていない状態で、入力データ自体の背後にある性質や特徴や構造などを分析することを、「教師なし学習（キョウシナシガクシュ、Unsupervised Learning）」と言います。
- ・ 「教師なし学習」の例として、「クラスタリング（Clustering）」によるデータ分類や、「次元削減（ジゲンサクゲン、Dimensionality Reduction）」による特徴量抽出などが挙げられます。
- ・ 「クラスタリング」は入力データを分類するものですが、「教師あり学習」の「分類」では、分類済みの入力データで分類方法を学習するのに対し、「教師なし学習」の「クラスタリング」では、未分類の似たデータを集めてデータ構造を発見します。「クラスタリング」のアルゴリズムとしては、「k平均法（ケイヘキンホウ、k-means clustering）」や「混合ガウス分布（コンゴウガウスブンブ、Gaussian mixture models、GMM）」などがあります。
- ・ 「次元削減」は入力データの次元を落とし、できるだけ「情報」を保ちながらより小さな次元の情報へ対応させることです。「次元削減」の例として、多変量データ解析の分野で提案された「主成分分析（シュセイブンブンセキ、principal component analysis、PCA）」や「因子分析（インブンセキ、factor analysis、FA）」がよく知られています。

## (2.3) 強化学習

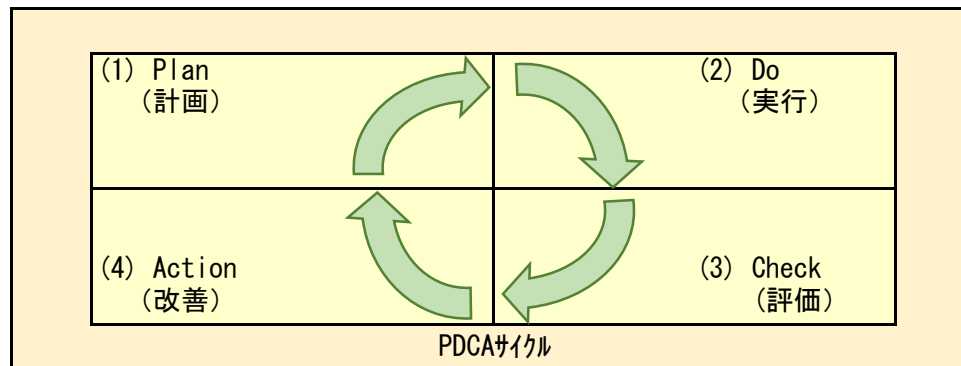
- ・「強化学習 (きょうかがくしゅう、Reinforcement learning)」は、  
「認知」した状況の下で最適な「行動」を「判断」すること、すなわち「制御」を学習する方法です。
- ・ゲームなどでは、『その時の「状態」を認知し、最終的な「報酬」を最大化するように、次にどんな「方策」を講じるか判断して「行動」に移す』ということを繰り返します。「強化学習」では、この流れを学習して、制御を行います。
- ・強化学習では、「エージェント (agent、制御器)」が「環境 (environment)」との相互作用を通じて情報収集しながら、「環境」を制御する方法を学習します。  
「環境」は、状態とその変化を規定する制御対象の系(システム)のことです。  
「エージェント」は、環境の置かれた状態にもとづいて行動を選択しながら環境を制御する主体で、制御プログラムにあたります。
- ・「エージェント」がある時刻  $t$  において「行動(action)  $A_t$ 」を通して「環境」に影響を及ぼす際に、  
「エージェント」は、時刻  $t$  における「状態 (state)  $S_t$ 」をみて、  
何らかの「方策 (policy)」に基づいて「行動  $A_t$ 」を決定します。  
一方、「環境」は、「エージェント」の「行動  $A_t$ 」がもたらす影響に応じて「報酬 (reward)  $R_{t+1}$ 」を「エージェント」に返すと同時に、  
「状態 (state)」を更新( $S_t \Rightarrow S_{t+1}$ )します。  
こうした「エージェント」と「環境」のやりとりを繰り返しつつ、  
最終的な「報酬」を最大化するように、「エージェント」の「方策」を最適化することが、強化学習の目的です。





- ・「強化学習」は、「PDCAサイクル（ピーデーシーエーサイクル、PDCA cycle、plan-do-check-act cycle）」と似ています。これは、「Plan（計画）」「Do（実行）」「Check（評価）」「Action（改善）」の四段階を繰り返すことにより、業務を継続的に改善する手法です。

- (1) Plan（計画）：従来の実績や将来の予測などをもとにして業務計画を作成する。
- (2) Do（実行）：計画に沿って業務を行う。
- (3) Check（評価）：業務の実施が計画に沿っているかどうかを評価する。
- (4) Action（改善）：実施が計画に沿っていない部分を調べて改善をする。



## (2.4) 転移学習

- ・上記の分類とは異なりますが、学習済みモデルの一部を学習しなおして使う「転移学習（テンイガクシュウ、Transfer Learning）」という学習方法があります。

転移学習は、ある領域で学習させたモデルを、別の領域に適応させる技術です。

具体的には、広くデータが手に入る領域で学習させたモデルを少ないデータしかない領域に適応させたり、シミュレート環境で学習させたモデルを現実に適応させたりする技術です。

### 【出典・参考】

機械学習⇒ <https://ja.wikipedia.org/wiki/機械学習>

教師あり学習⇒ <https://ja.wikipedia.org/wiki/教師あり学習>

教師あり学習⇒「現場で使える！強化学習・深層学習による探索と制御入門」（2019年07月 翔泳社）

教師なし学習⇒「現場で使える！Python 機械学習入門」（2019年05月 翔泳社 大曾根圭輔 他著）

k平均法⇒ <https://ja.wikipedia.org/wiki/K平均法>

次元削減⇒ [https://qiita.com/aya\\_taka/items/4d3996b3f15aa712a54f](https://qiita.com/aya_taka/items/4d3996b3f15aa712a54f)

強化学習⇒「現場で使える！強化学習・深層学習による探索と制御入門」（2019年07月 翔泳社）

PDCAサイクル⇒ <https://flhouse.co.jp/article/178/pdca>

PDCAサイクル⇒ <https://ja.wikipedia.org/wiki/PDCAサイクル>

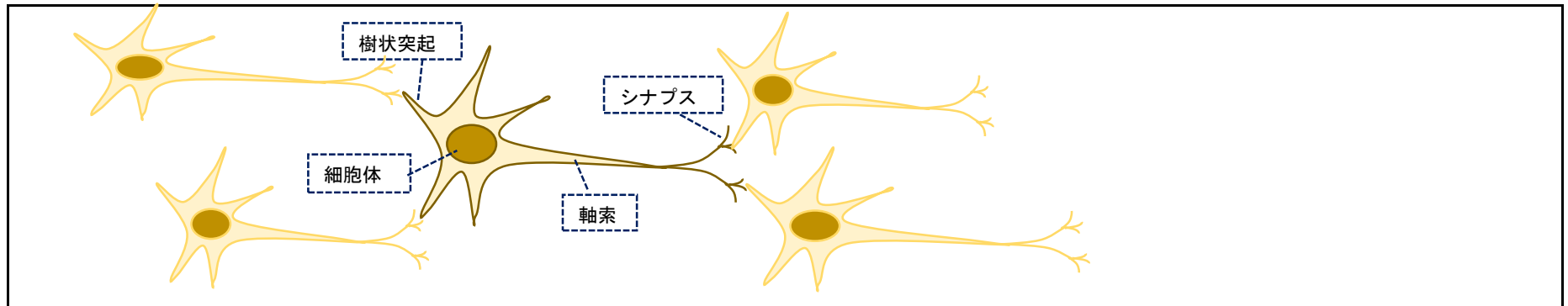
転移学習⇒ <https://qiita.com/icoxfog417/items/48cbf087dd22f1f8c6f4>

### (3) 脳神経系のモデル化

(※ この章は、「セナ-2回目」の内容を含みます。)

#### (3.1) 神経細胞間の信号伝達の様子

- ・ 20世紀の初頭に、人間の脳は「ニューロン (neuron)」と呼ばれる神経細胞が、「シナプス (synapse)」と呼ばれる結合部位を介して、多数結合してできているネットワークであることが示されました (下図)。  
神経細胞は主に、細胞核のある細胞体、他の細胞からの入力を受ける樹状突起、他の細胞に出力する軸索から構成されています。  
樹状突起が受け取った信号は細胞体で処理され、軸索を通して、次の神経細胞に伝達されます。



- 細胞体 : 細胞核とそれを取り巻く細胞質からなり、細胞としての機能はほとんどここで行われます。  
樹状突起と軸索が会合する部位でもあります。
- 樹状突起 : 細胞体から木の枝のように分岐しながら広がる構造であり、他の神経細胞などから信号を受け取る働きをします。  
一つの神経細胞に、樹状突起は複数あります。
- 軸索 : 細胞体から延びている突起状の構造で、神経細胞において信号の出力を担います。  
一つの神経細胞に、軸索は基本的には一本だけあります。
- シナプス : 前の細胞の軸索終末と後ろの細胞の樹状突起の間の情報を伝達する部分にある、伝達構造です。

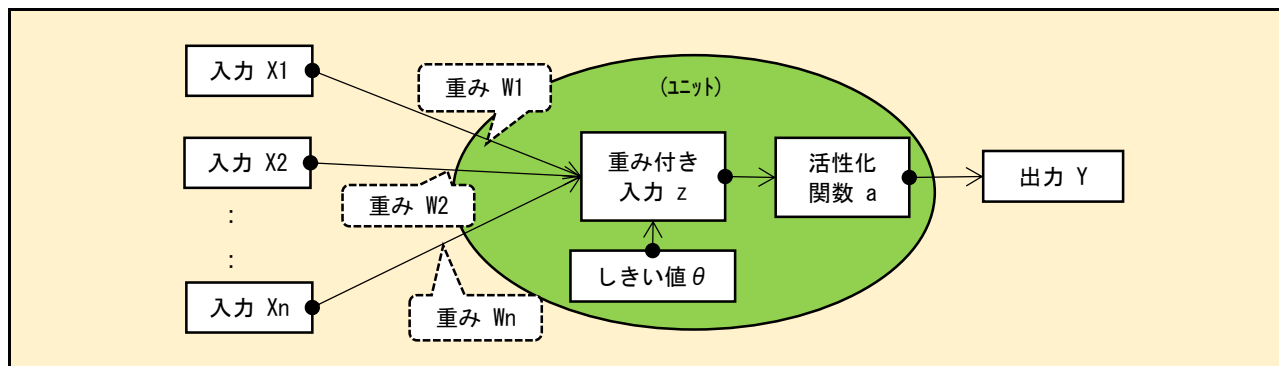
#### 【出典・参考】

神経細胞⇒ <https://ja.wikipedia.org/wiki/神経細胞>

「深層学習 Deep Learning」人工知能学会監修 (近代科学社 2015年11月)

### (3.2) 神経細胞間の信号伝達のモデル化

- ・ 個々の神経細胞間の信号伝達の様子を以下のようにモデル化します。これを「形式ニューロン」「ユニット」「ノード」または「セル」と呼びます。



- ・ 関連する変数名・関数名は、以下のとおり：

$X_i$  : ニューロン  $i$  ( $1 \sim n$ ) からの入力信号

$W_i$  : ニューロン  $i$  ( $1 \sim n$ ) からの入力信号  $X_i$  の重み

$\theta$  : 入力信号の重み付き線形和 ( $\sum_{i=1}^n (X_i * W_i)$ ) に対し、出力信号の有無を判別する為のしきい値

$b$  : バイアス (しきい値  $\theta$  の符号を変えたもの  $b = -\theta$ )

$z$  : 重み付き入力

$a$  : 活性化関数 (カセバカンス、activation function、伝達関数(デントツカンス、transfer function)ともいいます)

$Y$  : 出力信号

と表現した時、重み付き入力  $z$  は以下の式で与えられます：

$$\begin{aligned} z &= \{ X_1 * W_1 + X_2 * W_2 + \dots + X_n * W_n \} - \theta \\ &= \{ \sum_{i=1}^n (X_i * W_i) \} - \theta \\ &= \{ \sum_{i=1}^n (X_i * W_i) \} + b \end{aligned} \quad \left\{ \begin{array}{l} \text{入力信号の重み付き線形和(括弧 \{ \} 内の値)} < \text{しきい値 } \theta \text{ ならば、} z < 0 \\ \text{入力信号の重み付き線形和(括弧 \{ \} 内の値)} \geq \text{しきい値 } \theta \text{ ならば、} z \geq 0 \end{array} \right.$$

- ・ 出力信号  $Y$  は、重み付き入力  $z$  を入力変数として、活性化関数  $a$  で、以下の式で与えられます：

活性化関数  $a$  は、重み付き入力  $z$  に応じて、出力信号  $Y$  の大きさを規定する関数です：

$$\begin{aligned} Y &= a(z) \\ &= a\left(\left\{\sum_{i=1}^n (X_i * W_i)\right\} + b\right) \end{aligned}$$

### (3.3) ニューラルネットワークと適用モデルの選択

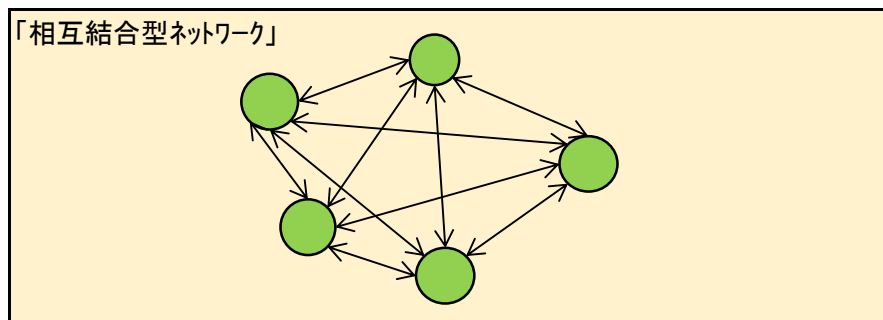
- ・「ユニット（形式ニューロン）」を複数つなげて「数理的な神経回路網」を構成したものが「ニューラルネットワーク（Neural Network, NN）」で、「階層型ニューラルネットワーク」と「相互結合型ネットワーク」があります。

- ・「階層型ニューラルネットワーク」は、右図のように数層から構成されます。

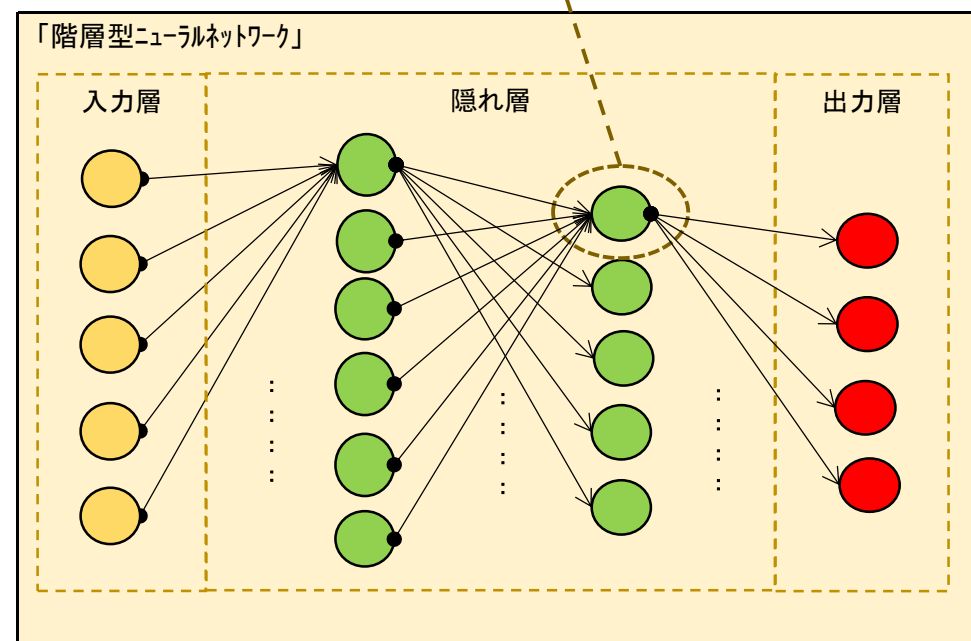
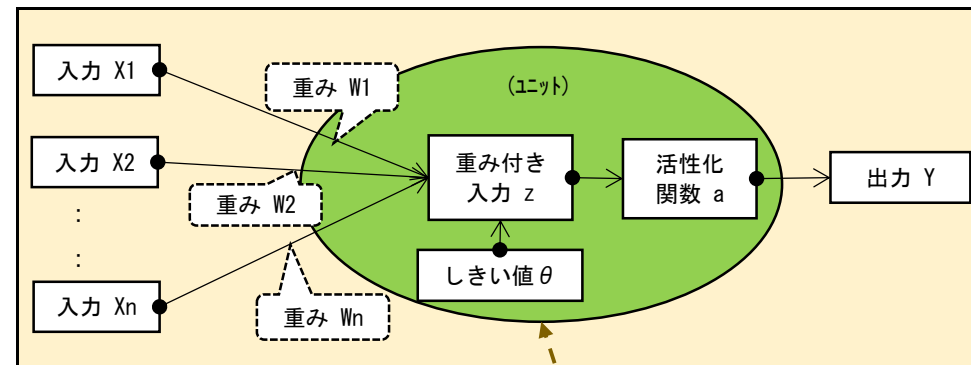
入力層	ニューラルネットワークへの情報を取り込む層
隠れ層	入力側からの信号に対して、重みづけ評価を行いながら、出力側へ信号伝達する層（中間層ともいいます）
出力層	ニューラルネットワークが算出した結果を出力する層

- ・「相互結合型ネットワーク」は、下図のように、層がなくユニット同士が互いに結合しています。このユニットの持つ値により、ある状態を記憶することができ、これを連想記憶と言います。

相互結合型ネットワークは、制約ボルツマンマシンを起源としており、それを深層化した深層ボルツマンマシン、深層信念ネットワークなどがあります。



- ・ 目的の為にどのような「ニューラルネットワーク」のモデルを適用するかについては、「セミナ2回目 機械学習のモデルとアルゴリズム」を参考にしてください。

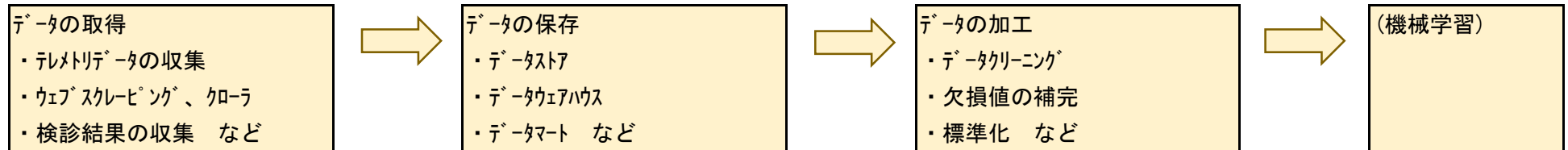


- : 入力層のユニット
- : 中間層のユニット
- : 出力層のユニット

## (4) 前処理

### (4.1) データを機械学習モデルに入力するまでの流れ

- ・データを機械学習モデルに入力するまでの流れは以下のようになります：



- ・「テレメトリデータ (telemetry data)」とは、「テレメタ(遠隔計測装置)」を使って、遠隔地の測定結果をセンタに送信して収集したデータのことです。  
「アメダス (AMeDAS: Automated Meteorological Data Acquisition System: 自動気象データ収集システム)」で収集されるデータや、  
ガスや電気のメーター、自動販売機の売り上げ管理システムなどで収集されるデータがそれにあたります。
- ・「ウェブスクレーピング (Web scraping)」は、ウェブサイトから情報を抽出するコンピュータソフトウェア技術のことです。  
「クローラ (Crawler)」とは、ウェブ上の文書や画像などを周期的に取得し、自動的にデータベース化するプログラムのことです。
- ・「データストア (data store)」とは、データを決められた形式で記憶装置などに永続的に保存・蓄積するソフトウェア、および、そのような機能・処理のことです。データを保存する機能一般を抽象的に表す用語で、ファイルシステムやデータベースもデータストアの具体的な実装の一つということができます。
- ・「データウェアハウス (Data Ware House, DWH)」とは、意思決定のために、基幹系などの複数システムから、必要なデータを収集し、  
目的別に再構成して時系列に蓄積した統合データベースで、データ分析や意思決定に役立てます。
- ・「データマート (data mart)」とは、企業などで情報システムに記録・蓄積されたデータから、利用部門や用途、目的などに応じて  
必要なものだけを抽出、集計し、利用しやすい形に格納したデータベースのことです。「マート (mart)」は「小売店」の意味です。

## (4.2) データの種類について

- ・データの種類について、表現する情報の性質に基づき、数学・統計学的に分類するのが「尺度水準（スケール・レベル、Level of measurement）」です。  
以下に、一覧表にしておきます：

尺度水準	説明
名義尺度 (メイギンシャクド、 nominal scale)	「名義尺度」とは、分類して名称を付与することに意味があり、順序には意味がない尺度のことを言います。 (例) 性別、血液型、電話番号、カードのID番号
順序尺度 (ジュンシヨウシャクド、 ordinal scale)	「順序尺度」とは、間隔には意味がないが、大小や順序に意味がある尺度のことを言います。 「順位尺度（ジュンイシャクド）」とも言います。 (例) 売上げランキングの順位、成績の5段階評価
間隔尺度 (カンカクシャクド、 interval scale)	「間隔尺度」とは、大小関係に加えて数値の差が意味をもつ尺度で、間隔に意味があります。 原点は相対的な意味しか持ちません。「距離尺度」とも言います。 「順序尺度」の性質も備えています。 (例) 温度(°C)
比率尺度 (ヒリツシャクド、 ratio scale)	「比率尺度」とは、数値の大小関係・差・比に意味がある尺度で、絶対原点をもつ間隔尺度をいいます。 「比例尺度（ヒレイシャクド）」「比尺度」とも言います。 「順序尺度」、「間隔尺度」の性質も備えています。 (例) 絶対温度(K)、体重、身長、年齢

### (4.3) 何を入力データとするか

(※ この節は、「セナ-6回目」からの再掲を含みます。)

- ・ 何を入力データとするかについては、その適用業務の専門知識や経験などに基づいて決めることになります。

- ・ 入力データとなる説明変数を選ぶ際に、変数間の多重共線性に気を付ける必要があります。

「多重共線性 (タジユキョウセンセイ, multicollinearity)」(略称「マルチコ」)とは、統計モデルの説明変数間の相関関係が高いときに解析計算が不安定になる、といった性質です。これにより、正しく推計できなくなるといった現象が起こります。

(行列の逆行列計算では、行(または列)の間に線形の関係(一方の行が他方の行の定数倍であるなど)がある場合、行列式が0になり計算できませんが、それに近い現象が発生しているとみて良いでしょう。)

- ・ 「多重共線性」の問題の最も一般的な解消法は、「相関関係が高いと考えられる説明変数を外すこと」です。

(例) コンビニの月間のアイスクリームの売り上げ

- ・ 目的変数: (1) コンビニの月間のアイスクリームの売り上げ
- ・ 説明変数: (1) 来客数  
(2) 最高気温が30℃以上の日数  
(3) 降雨日数  
(4) 月間降水量  
(5) 平均気温

という分析モデルを考えたとき、説明変数で「(3) 降雨日数」と「(4) 月間降水量」は相関が高く、多重共線性が発生する可能性が高いとみられます。

この場合、(3)か(4)の何れか一方を説明変数として不採用とします。



#### (4.4) 入力データのクリーニング

- ・入力データ（説明変数）を統計処理するに先立って、  
入力データの中から重複や誤記、表記のぶれなどを探し出し、削除や修正などを行い、入力データの品質を高める必要があります。  
この作業を「データクリーニング（data cleaning）」または「データクレンジング（data cleansing）」と言います。

- ・「データクリーニング」は、使用する単位や書式や定義域（入力値の範囲）等を基準にして行います。

以下に様々な例を示します：

- （例1） 電話番号の文字列を、半角のみにし、ハイフンを取り除く。（“0 1 － 2 3 4 5 － 6 7 8 9” ⇒ “0123456789”）
- （例2） 住所の誤字を直す。（“東京都太田区” ⇒ “東京都大田区”）
- （例3） 数字文字列の3桁区切り表記を、数字列に直す。（“123,456,789” ⇒ “123456789”）
- （例4） 欠けているデータを補完する。（（“123”, “789”） ⇒ （“123”, “456”, “789”））
- （例5） 名簿の重複レコードを削除する。（（{“1”, “山田”}, {“1”, “山田”}, {“2”, “佐藤”}） ⇒ （{“1”, “山田”}, {“2”, “佐藤”}））
- （例6） 長さの単位を[cm]から[mm]へ統一する。（“175” ⇒ “1750”）
- （例7） 異常な地上気圧を（誤記と見做して）補正する。（“10240” ⇒ “1024”）
- （例8） 異常な地上気圧を（計測機器故障と見做して）削除する。（“10240, 998, 30, 1013” ⇒ “998, 1013”）

#### (4.5) 入力データの欠損値の扱い

- ・ 入力データ（説明変数）が数種類からなる場合、各入力レコードに必ずしも全ての種類のデータが揃っている訳ではありません。欠けているデータのことを「欠損値（ケツシツ、missing value）」と言います。
- ・ 「欠損値」の例を以下に示します：
  - （例1） アンケートデータの未回答項目
  - （例2） 健康診断時の未検査項目
- ・ 欠損値がある場合、統計処理に支障があるため、以下の様な対応をします：
  - （1） 欠損値を含む行（レコード）を削除する。
  - （2） 欠損値を含む列（説明変数、目的変数）を削除する。
  - （3） 欠損値を何らかの値（指定値、平均値、内挿値、外挿値など）で補完する。
- ・ Python のライブラリ「Pandas」の「DataFrame」クラスで欠損値対応のメソッドが提供されているので、実装時に参照すると良いでしょう。（「Missing data handling」の機能）

## (4.6) 入力データの標準化

- ・平均と分散についての公式を用いて、平均  $\mu$ 、標準偏差  $\sigma$  の確率変数  $X$  を、平均 0、標準偏差 1 の確率変数  $Z$  へ変換することができます。

この操作を「標準化 (ヒョウジュンカ, standardization)」と言い、

$Z$  を「標準化確率変数 (ヒョウジュンカカリツハンスウ, standardized random variable)」

と言います。確率変数  $X$  から次式により標準化確率変数  $Z$  を求めます：

標準化確率変数

$$Z = (X - \mu) / \sigma$$

$\mu$  : 平均値

$\sigma$  : 標準偏差

- ・入力データ (説明変数) が数種類からなる場合で、入力データ毎の平均や分布が大きく異なる場合などに標準化を行います。

これには、以下の様な効果があります。

- (1) 特定の特徴に偏った結果が出ることを防ぐことができます。
- (2) 異なる項目のデータであってもその大小を比較できるようになります。

## (4.7) 入力データの加工による作成

・モデルの学習時用の入力データが少ない場合、データに以下の様な加工を加えることにより、入力データを増やすことができます。

- (1) 画像の場合、回転移動、拡大・縮小、平行移動、輝度変換およびそれらの組合せで別の画像を作成する。  
(これは、Keras の ImageDataGenerator クラスを用いて実装可能です)
- (2) 画像の場合、元の画像にノイズを加えることにより別の画像を作成する。  
(これは、numpy を用いて、マスクノイズデータやガウシアンノイズデータを加える、といった実装で可能です)

## (4.8) 入力データの符号化

・入力データを符号化してモデルに取り込んだ方が良い場合があります。

・出力が K 種類への分類モデルの場合、目的変数を K 次元のベクトルとします。

ある入力データが第 k 番目 ( $k=1\sim K$ ) のクラスへ分類されることが分かっている場合、

目的変数ベクトルの k 番目の要素だけ 1 でそれ以外を 0 とする符号化を行うと便利です。

これを「1-of-K 符号化法 (ワンオブケフコウカク、One-of-K encoding)」と言います (「one-hot encoding」とも言います)。

(例)

5つのクラスに分類するモデルの場合、n 番目の入力データ  $X_n$  が3番目のクラスに分類された場合、

目的変数ベクトル  $T_n$  を「3」と表現するのではなく、次のように表現します：

$$T_n = (0, 0, 1, 0, 0)$$

## 【出典・参考】

全般⇒「現場で使える！Python 機械学習入門」（2019年05月 翔泳社 大曾根圭輔 他著）

テレメトリデータ⇒ <https://kotobank.jp/word/テレメトリー-577329>

ウェブスクレイピング⇒ <https://ja.wikipedia.org/wiki/ウェブスクレイピング>

データウェアハウス⇒ [http://www.advanlink.co.jp/crmconsultation/crm\\_dwh.html](http://www.advanlink.co.jp/crmconsultation/crm_dwh.html)

データマート⇒ <http://e-words.jp/w/データマート.html>

データストア⇒ <http://e-words.jp/w/データストア.html>

全般⇒「現場で使える！Python 機械学習入門」（2019年05月 翔泳社 大曾根圭輔 他著）

尺度水準⇒ <https://ja.wikipedia.org/wiki/尺度水準>

尺度水準⇒ <https://to-kei.net/basic/glossary/scale/>

多重共線性⇒ <https://xica.net/vno4ul5p/>

多重共線性⇒ <https://support.minitab.com/ja-jp/minitab/18/help-and-how-to/modeling-statistics/regression/supporting-topics/model-assumptions/multicollinearity-in-regression/>

行列式の計算⇒「ベクトルと行列」（1976年05月 現代数学社 横田一郎、他著）

データクリーニング⇒ <https://fukaichi.jp/web-marketing-glossary-list/data-cleaning>

データクリーニング⇒ <http://e-words.jp/w/データクレンジング.html>

欠損値⇒ <https://www.weblio.jp/content/欠損値>

欠損値, Pandas⇒ <https://techacademy.jp/magazine/17697>

Pandas⇒ <https://pandas.pydata.org/pandas-docs/stable/reference/frame.html>

標準化確率変数⇒ <https://k-san.link/standardized/>

標準化⇒ <https://bellcurve.jp/statistics/course/19647.html>

入力データの加工⇒

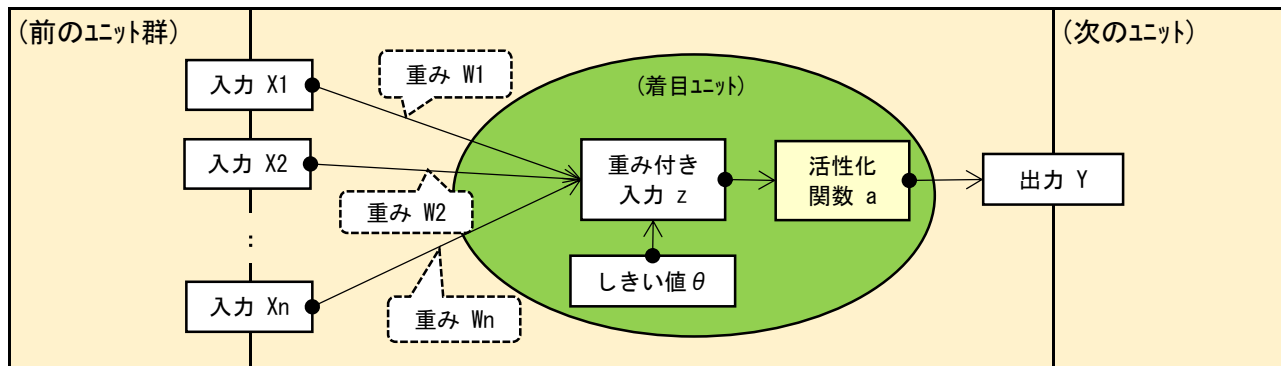
「現場で使える！TensorFlow開発入門 Kerasによる深層学習モデル構築手法」（2018年04月 翔泳社 太田満久、須藤広大、黒澤匠雅、小田大輔 共著）

## (5) 活性化関数

(※ この章は、「セミナー2回目」の内容を含みます。)

- あるユニットに着目した時、そのユニットへの入力信号の重み付き線形和にバイアスを加えたもの(下図で「重み付き入力  $z$ 」)を、どのように活性化して次のユニットに渡すかを、「活性化関数 (カセイクナシ、Activation function)」で指定します。

「活性化関数」は「出力関数 (シュツヨウカナシ、Output function)」とも言います。



- 「活性化関数」には様々なものが考案されており、主なものを一覧表に示します。  
このうち「セミナー7回目」で使用しているものは、「ReLU」「sigmoid」「ソフトマックス関数」の3つです。

- 「ソフトマックス関数 (softmax function)」は、 $k$ 個の入力値  $X_i$  ( $i=1\sim k$ )の順序関係を保ちながら、 $X_i$  を確率としての値  $Y_i$  ( $i=1\sim k$ ,  $0\leq Y_i\leq 1$ ,  $\sum_{i=1}^k Y_i=1$ )に変換する関数で、 $k$ 個のクラスへ分類する際の所属の確率を示すのによく用いられます。

「シグモイド関数 (sigmoid function)」は「ソフトマックス関数」で分類するクラス数を2個に限定したものに相当します。

- (1)  $k$  個のクラスへ分類する時のソフトマックス関数

$$Y_i = \exp(X_i) / \sum_{j=1}^k \exp(X_j) \quad (i=1\sim k)$$

- (2) 2個のクラスへ分類する時のソフトマックス関数

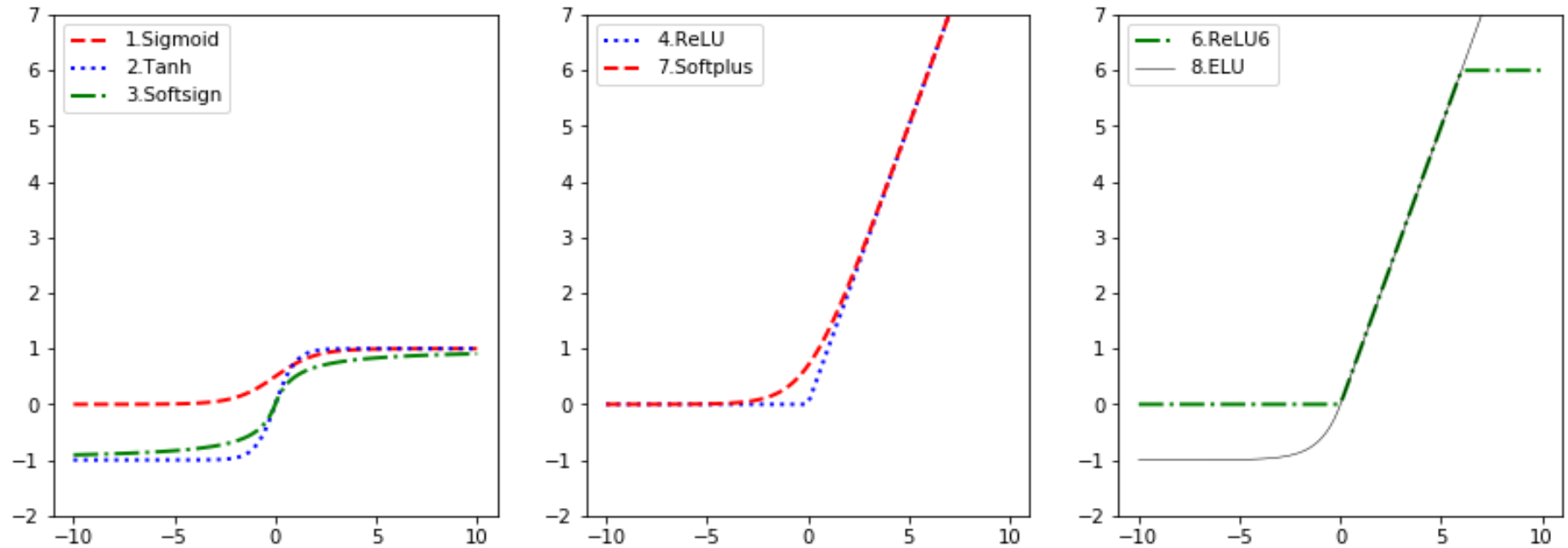
$$Y_1 = \exp(X_1) / \{\exp(X_1) + \exp(X_2)\} = 1 / \{1 + \exp(X_2 - X_1)\} = 1 / \{1 + \exp(-Z)\} \quad (Z = X_1 - X_2) \quad \dots \text{シグモイド関数}$$

$$Y_2 = \exp(X_2) / \{\exp(X_1) + \exp(X_2)\} = 1 - Y_1$$

- 「ReLU (正規化線形関数)」は計算が容易で(正(負)の領域では微分が 1(0))、「sigmoid関数」で問題となる学習時の「誤差逆伝播法」での「勾配消失問題 (コハバ イショウツモンダ イ、vanishing gradient problem)」もなく、一般的に使用されています。

No.	活性化関数名	和名	英語	式「 $Y = a(z)$ 」	値域	tensorflow の関数名	備考
1	sigmoid	シグモイド関数 (ロジスティック関数)	sigmoid function (logistic function)	$1 / \{ 1 + \exp(-z) \}$	0 ~ 1	sigmoid	微分可能で解析的に扱いが容易だが、誤差逆伝播法の誤差項が、学習中に0に近づいていく傾向にある為(勾配消失問題 Vanishing Gradient problem)、あまり使用されていません。 「ソフトマックス関数」で分類するクラス数を2個に限定したものに相当します(1986年)。
2	tanh	ハイパボリックタンジェント (双曲線正接関数)	hyperbolic tangent	$\{ \exp(z) - \exp(-z) \} / \{ \exp(z) + \exp(-z) \}$	-1 ~ +1	tanh	活性化関数は原点を通すべきと言う考えから、標準シグモイド関数よりもそれを線形変換した tanh の方が良いと提案されました。(1998年)
3	softsign	ソフトサイン関数	softsign function	$z / \{ \text{abs}(z) + 1 \}$	-1 ~ +1	softsign	正弦関数に対する連続近似として位置づけられます。 (2010年 Xavier Glorot)
4	ReLU	正規化線形関数 (ランプ関数)	Rectified Linear Unit	$\max(0, z)$	0 ~ $\infty$	relu	ReLU (Rectified Linear Units) は、入力値が負の値の場合は0、正の値の場合は入力値をそのまま出力します。 最も一般的で基本的な方法です。(2011年 Xavier Glorot)
5	Clipped ReLU	Clipped ReLU	Clipped Rectified Linear Unit	$\min(\max(0, z), \text{clipVal})$ clipVal : 0以上の浮動小数点数。正の部分の最大値。	0 ~ clipVal	---	Clipped ReLU は、ReLUの特別なバージョンで、出力値が一定の大きさ以上にはならないように変更されたものです。ReLUでは、入力が正の値のときに線形な変換が適用されますが、Clipped ReLU では入力に対して、出力は頭打ち状態 (Clipping) になります。
6	ReLU6	正規化線形関数 (ランプ関数)	Rectified Linear Unit 6	$\min(\max(0, z), 6)$	0 ~ +6	relu6	ReLU関数で上限を6としたもの。 Clipped RELUの一つで、計算速度が速く、値の消失や発散の問題にも悩まされません。
7	softplus	ソフトプラス関数	softplus function	$\log(\exp(z) + 1)$	0 ~ $\infty$	softplus	ReLU関数を滑らかにしたものです。
8	ELU		Exponential Linear Unit	$\exp(z) - 1$ for $z < 0$ $z$ for $z \geq 0$	-1 ~ $\infty$	elu	ELUは、ReLUの特別なバージョンで、負の入力に対して指数関数 (Exp) から1を引いた値を適用しています。
9	LeakyReLU	LeakyReLU	Leaky Rectified Linear Unit	$\alpha * z$ for $z < 0$ $z$ for $z \geq 0$  $\alpha$ : 0以上の浮動小数点数。負の部分の傾き。	$-\infty \sim \infty$	leaky_relu	LeakyReLUは、ReLUの特別なバージョンで、ユニットがアクティブでないとき ( $z < 0$ ) でも微少な勾配を可能とします。これにより勾配消失を防ぎ、学習速度を向上させる効果が期待できます。DCGAN (Deep Convolutional Generative Adversarial Network) と呼ばれるモデルなどで利用されます。
10	softmax	ソフトマックス関数	softmax function	$Y_i = \exp(X_i) / \sum_{j=1}^k \exp(X_j)$  ( $i=1 \sim k$ )	0 ~ 1	softmax	出力が k 個ある場合に、k個の入力値 $X_i$ ( $i=1 \sim k$ ) の順序関係を保ちながら、確率としての値 $Y_i$ ( $i=1 \sim k$ , $0 \leq Y_i \leq 1$ , $\sum_{i=1}^k Y_i = 1$ ) に変換する関数。k個のクラスへ分類する際の所属の確率を示すのによく用いられます。

- ・上記で紹介した「活性化関数」のうち、幾つかを図示します。



#### (活性化関数とその描画の実装例)

- ・活性化関数はライブラリで提供されているか、提供されているライブラリのモジュールの引数として活性化関数の名前を指定するだけでも実装できますが、ここでは、敢えて関数の定義のまま実装します。

#### (リスト08-(05)-1\_様々な活性化関数)

```
#####/
# リスト08-(05)-1_様々な活性化関数
#####/
import numpy as np
import matplotlib.pyplot as plt
%matplotlib inline
```



```

#(1) シグモイド関数
def f_sigmoid(x):
    return 1 / ( 1 + np.exp(-x) )

#(2) ハイパーボリックタンジェント (双曲線正接関数)
def f_tanh(x):
    return (np.exp(x) - np.exp(-x)) / (np.exp(x) + np.exp(-x))

#(3) ソフトサイン関数
def f_softsign(x):
    return x / (np.abs(x) + 1)

#(4) 正規化線形関数 (ランプ関数)
def f_ReLU(x):
    xcnt = len(x)
    ylist = np.zeros(xcnt)
    for ii in range(xcnt):
        if x[ii] < 0:
            ylist[ii] = 0
        else:
            ylist[ii] = x[ii]
    return ylist

#(6) 正規化線形関数 (ランプ関数、上限=6)
def f_ReLU6(x):
    xcnt = len(x)
    ylist = np.zeros(xcnt)
    for ii in range(xcnt):
        if x[ii] < 0:
            ylist[ii] = 0
        elif x[ii] > 6:
            ylist[ii] = 6
        else:
            ylist[ii] = xlist[ii]
    return ylist

```

```

#(7) ソフトプラス関数
def f_softplus(x):
    return np.log(np.exp(x) + 1)

#(8) Exponential Linear Unit
def f_ELU(x):
    xcnt = len(x)
    ylist = np.zeros(xcnt)
    for ii in range(xcnt):
        if x[ii] < 0:
            ylist[ii] = np.exp(x[ii]) -1
        else:
            ylist[ii] = x[ii]
    return ylist

# x データ列
xlist = np.linspace(-10, 10, 100)

# y データ列
y_sigmoid = f_sigmoid(xlist)
y_tanh = f_tanh(xlist)
y_softsign = f_softsign(xlist)
y_ReLU = f_ReLU(xlist)
y_ReLU6 = f_ReLU6(xlist)
y_softplus = f_softplus(xlist)
y_ELU = f_ELU(xlist)

# グラフ表示
plt.figure(figsize=(14, 5))
plt.subplots_adjust(wspace=0.2, hspace=0.5)

plt.subplot(1, 3, 1)
plt.plot(xlist, y_sigmoid, 'r--', label='1. Sigmoid', linewidth=2)
plt.plot(xlist, y_tanh, 'b:', label='2. Tanh', linewidth=2)
plt.plot(xlist, y_softsign, 'g-.', label='3. Softsign', linewidth=2)
plt.ylim([-2, 7])
plt.legend(loc='upper left')

```

```
plt.subplot(1, 3, 2)
plt.plot(xlist, y_ReLU, 'b:', label='4. ReLU', linewidth=2)
plt.plot(xlist, y_softplus, 'r--', label='7. Softplus', linewidth=2)
plt.ylim([-2, 7])
plt.legend(loc='upper left')

plt.subplot(1, 3, 3)
plt.plot(xlist, y_ReLU6, 'g-.', label='6. ReLU6', linewidth=2)
plt.plot(xlist, y_ELU, 'k-', label='8. ELU', linewidth=0.5)
plt.ylim([-2, 7])
plt.legend(loc='upper left')

plt.show()
```

#### 【出典・参考】

活性化関数全般 ⇒ <http://yagami12.hatenablog.com/entry/2017/09/17/111935>

活性化関数全般 ⇒ [https://www.renom.jp/ja/notebooks/tutorial/basic\\_algorithm/activation\\_types/notebook.html](https://www.renom.jp/ja/notebooks/tutorial/basic_algorithm/activation_types/notebook.html)

keras ライブラリ ⇒ <https://keras.io/ja/activations/>

keras ライブラリ ⇒ <https://keras.io/ja/layers/advanced-activations/>

ReLU6 ⇒ [https://www.tensorflow.org/api\\_docs/python/tf/nn/relu6](https://www.tensorflow.org/api_docs/python/tf/nn/relu6)

ReLU6 ⇒ <http://www.cs.utoronto.ca/~kriz/conv-cifar10-aug2010.pdf>

LeakyReLU ⇒ [https://web.stanford.edu/~awni/papers/relu\\_hybrid\\_icml2013\\_final.pdf](https://web.stanford.edu/~awni/papers/relu_hybrid_icml2013_final.pdf)

LeakyReLU ⇒ [https://www.tensorflow.org/api\\_docs/python/tf/nn/leaky\\_relu](https://www.tensorflow.org/api_docs/python/tf/nn/leaky_relu)

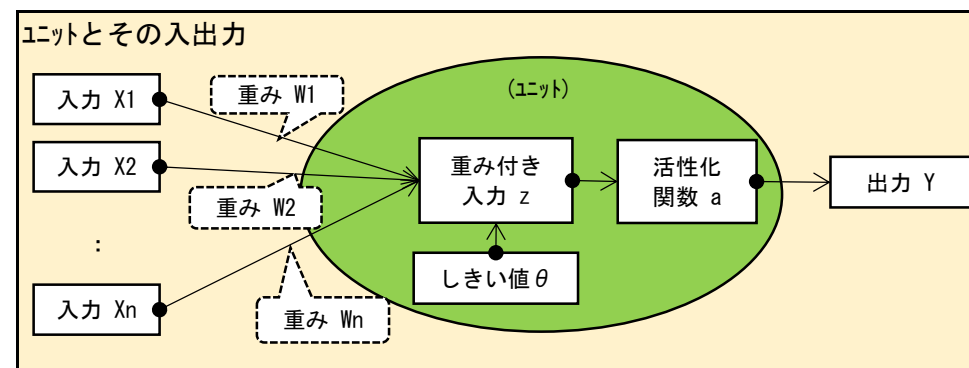
ELU ⇒ <https://arxiv.org/pdf/1511.07289v1.pdf>

sigmoid関数とソフトマックス関数の関係 ⇒ <https://mathtrain.jp/softmax>

## (6) 損失関数

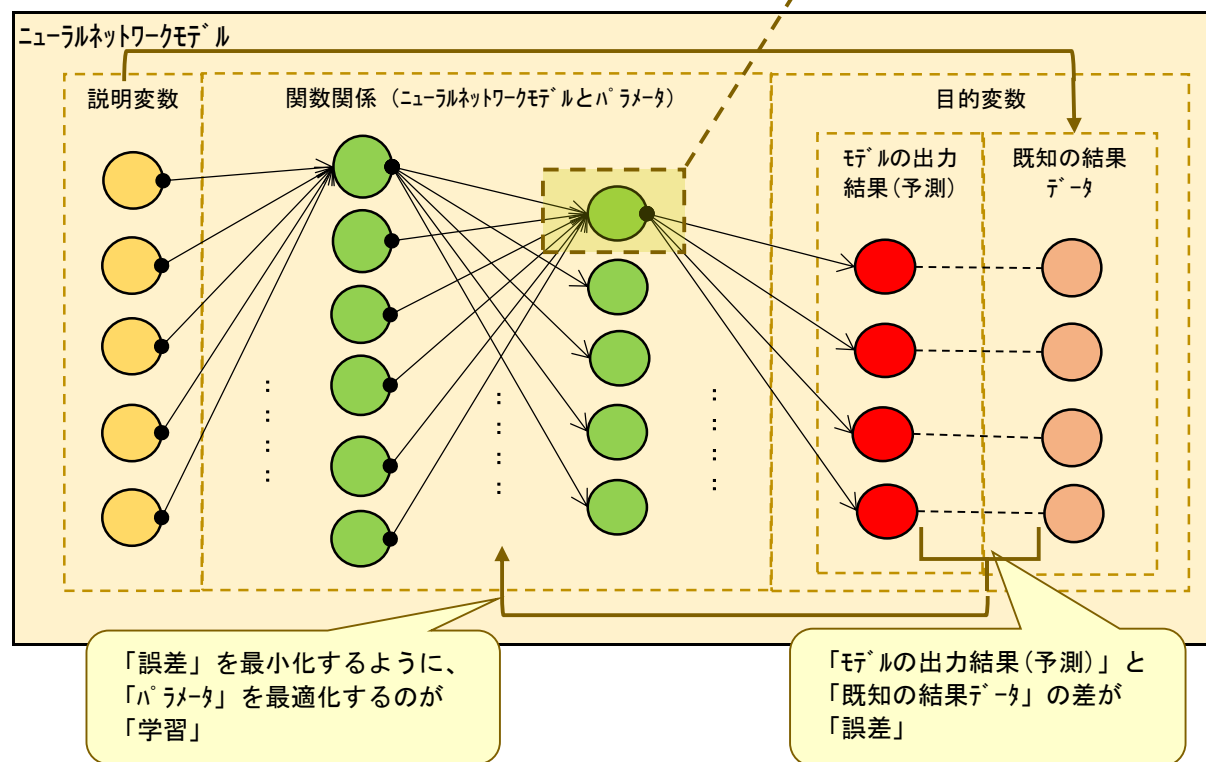
### (6.1) 損失関数・誤差関数・コスト関数・目的関数

- ・「教師あり学習」では、一方を入力（説明変数）、他方を出力（目的変数）とする関数関係として記述し、関数の出力を観測データ（目的変数）に近づけるように、関数関係を学習します。
- ・この関数関係は「ニューラルネットワークモデル」とその「パラメータ」に相当します。「ニューラルネットワークモデル」は使用する層やユニットの数、ユニット間の結合方法、使用する活性化関数などのネットワーク構成で、「パラメータ」は、各ユニットの重み $W_i$ としきい値 $\theta$ などの調整可能な値のことです。



- ・説明変数と目的変数の関係を良く説明できるモデルが「良いモデル」であり、「モデルの出力結果(予測)」と「既知の結果データ」との差（「誤差 (ゴサ, error)」）が最小になるように、「パラメータ」を最適化するのが「学習」です。
- ・「学習」を通して最小化したいのはこの誤差であり、誤差を表す関数のことを、「誤差関数 (ゴサ関数, error function)」、又は「損失関数 (ロジツクス関数, loss function)」、又は「コスト関数 (コスト関数, cost function)」と言います。

これらは「最小化する」という目的があるので、「目的関数 (オブジェクティブ関数, objective function)」とも言います。



## (6.2) 尤度と損失関数

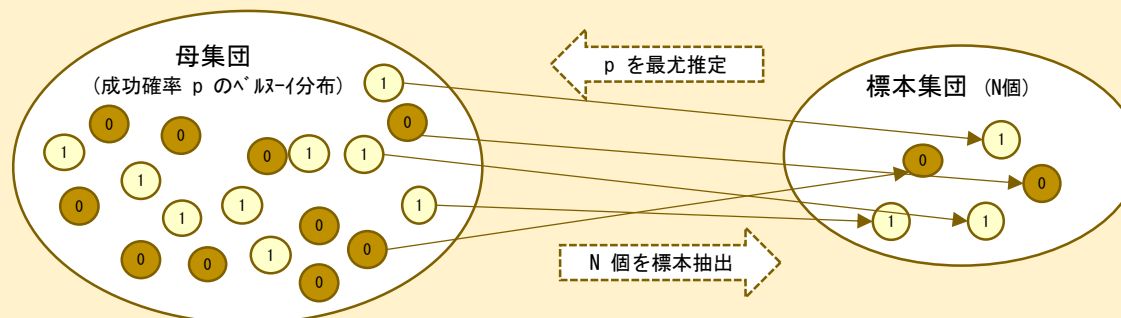
- ・ 統計学において、未知または既知のある確率分布の下で、確率変数の結果が出現する場合に、確率変数の出現結果からみて、元々の確率分布を推測する時の「尤もらしさ (モットモリサ)」を表す数値を、「尤度 (ユド、likelihood)」といい、「尤度」を計算する関数を「尤度関数 (ユドカンス、likelihood function)」と言います。
- ・ この「尤度」を最大化するように元の確率分布を推測することを「最尤推定 (サイクスタイ、maximum likelihood estimation、MLE)」と言います。「離散型確率分布」の場合は「確率質量関数」を、「連続型確率分布」の場合は「確率密度関数」を推定することになります。
- ・ 例として、ベルヌーイ分布の最尤推定方法を示します。 (ベルヌーイ分布については、「セミナー6回目」を参照のこと)

### ベルヌーイ分布の最尤推定方法

標本  $\{X_1, X_2, \dots, X_N\}$  はある母集団からの  $N$  個の無作為標本であるとしてます。

その母集団が従う確率分布がベルヌーイ分布で、1回のベルヌーイ試行で成功する(1となる)確率を  $p$  (未知数)とします。

この時、失敗する(0となる)確率は  $1-p$  となります。



この試行で標本  $\{X_1, X_2, \dots, X_N\}$  の観測値  $\{t_1, t_2, \dots, t_N\}$  の各  $t_n$  は  $\{1, 0\}$  の何れかを各々  $\{p, 1-p\}$  の確率で取ります。各  $X_n$  が値  $t_n$  を取るというベルヌーイ分布の確率質量関数  $P(X_n=t_n)$  は次式となります：

$$P(X_n=t_n) = p^{t_n} (1-p)^{1-t_n} \quad \text{for } t_n \in \{0, 1\}, n=1 \sim N \quad \dots (\text{式6.2-1})$$

この式を元に、標本  $\{X_1, X_2, \dots, X_N\}$  の観測値が  $\{t_1, t_2, \dots, t_N\}$  となる確率(同時確率)は、各観測結果  $X_n$  についての確率(式6.2-1)の直積で以下の式になります：

$$\begin{aligned} P(X_1=t_1, X_2=t_2, \dots, X_N=t_N) &= \prod_{n=1}^N P(X_n=t_n) = \prod_{n=1}^N p^{t_n} (1-p)^{1-t_n} \\ &= p^z (1-p)^{N-z} \quad (z \text{ は } t_n=1 \text{ となる } t_n \text{ の総和 : } z = \sum_{n=1}^N t_n) \end{aligned}$$

これが観測値が  $\{ t_1, t_2, \dots, t_N \}$  の下での成功確率  $p$  を推定する為の尤度関数  $L$  となります。

$$L(p|t_1, t_2, \dots, t_N) = \prod_{n=1}^N p^{t_n} (1-p)^{1-t_n} \quad \dots (式6.2-2)$$

$$= p^z (1-p)^{N-z} \quad (zはt_n=1となるt_nの総和: z = \sum_{n=1}^N t_n) \quad \dots (式6.2-3)$$

$p$  の最尤推定値は、尤度関数  $L$  を最大にする  $p$  の値になります。

ここで自然対数関数  $\log(x)$  が  $x$  について単調増加であり、関数値の大小関係が  $x$  の大小関係を保つという性質から、尤度関数  $L$  を最大にする  $p$  の値は、この対数をとった「対数尤度関数」の値も最大にすることが分かります。

対数尤度関数は以下の様になります（ $\log$ 関数の性質を用いて式を変形します）：

$$\log\{L(p|t_1, t_2, \dots, t_N)\} = \log\left\{\prod_{n=1}^N p^{t_n} (1-p)^{1-t_n}\right\} \quad \dots (式6.2-2) \text{ より}$$

$$= \sum_{n=1}^N \{t_n * \log p + (1-t_n) * \log(1-p)\} \quad \dots (式6.2-4)$$

$$= \log\{p^z (1-p)^{N-z}\} \quad \dots (式6.2-3) \text{ より}$$

$$= z * \log p + (N-z) * \log(1-p) \quad \dots (式6.2-5)$$

対数尤度関数  $\log(L)$  の  $p$  についての偏微分が 0 を与える  $p$  が最尤推定値であることの条件となります：

$$\partial \log(L) / \partial p = z/p - (N-z)/(1-p) = 0$$

$$\therefore p = z/N$$

$$\therefore p = 1/N \sum_{n=1}^N t_n \quad \dots (式6.2-7)$$

これは、標本平均に一致しています。

例えば、標本  $\{ t_1, t_2, \dots, t_N \}$  のうち、50% が 1（成功）ならば  $p=0.5$  が最尤推定値となります。

- ・「尤度」を最大化することは、モデルの学習時の評価指数である「損失関数」の値を最小化することになります。

この関係から、上記の例のようなクラス分類では「対数尤度関数（ $\log$  likelihood function）」にマイナス(-1)をかけて更に標本サイズ  $N$  で平均を取った関数を「損失関数」として扱います。

「対数尤度関数」(式6.2-4)に  $-1/N$  をかけて得られた関数は、

次に述べる損失関数一覧の「2クラス分類の場合の交差エントピー-誤差」CE(式6.3-4)になります。

$$CE = -1/N \sum_{n=1}^N \{t_n \log(y_n) + (1-t_n) \log(1-y_n)\} \quad \dots (式6.3-4)$$

## (6.3) 損失関数一覧

・「損失関数」には様々なものが提案されており、次表に代表的なものを幾つか掲載します。

No.	損失関数名	英語	略称	式	用途・特徴など
1	平均二乗誤差 (ヘイキンニジョウゴサ)	Mean Squared Error	MSE	$MSE = 1/N \sum_{n=1}^N (t_n - y_n)^2 \quad \dots (式6.3-1)$ $t_n$ : 実際の値 $y_n$ : モデルの予測値 $N$ : データ個数 ( $n=1 \sim N$ )	回帰問題において平均二乗誤差は最も一般的で、線形回帰モデルやニューラルネットワーク、決定木といった様々なモデルにおいて用いられます。平均二乗誤差の性質として外れ値に対して敏感であることが挙げられます。
2	平均絶対誤差 (ヘイキンゼツタイゴサ)	Mean Absolute Error	MAE	$MSE = 1/N \sum_{n=1}^N  t_n - y_n  \quad \dots (式6.3-2)$ $t_n$ : 実際の値 $y_n$ : モデルの予測値 $N$ : データ個数 ( $n=1 \sim N$ )	平均絶対誤差の性質として外れ値に強いことが挙げられます。
3	平均二乗対数誤差 (ヘイキンニジョウイスクゴサ)	Mean Squared Logarithmic Error	MSLE	$MSLE = 1/N \sum_{n=1}^N \{\log(1+t_n) - \log(1+y_n)\}^2 \quad \dots (式6.3-3)$ $t_n$ : 実際の値 $y_n$ : モデルの予測値 $N$ : データ個数 ( $n=1 \sim N$ )	平均二乗対数誤差を用いたモデルは予測が実値を上回りやすくなるという傾向があります。
4	交差エントロピー誤差 (コウサイエントロピーゴサ)	Cross Entropy Error	CE	<p>2クラス分類の場合</p> $CE = -1/N \sum_{n=1}^N \{t_n \log(y_n) + (1-t_n) \log(1-y_n)\} \quad \dots (式6.3-4)$ $t_n$ : 入力 $x_n$ のクラス1への所属確率の実際の値 (0, 1) $y_n$ : 入力 $x_n$ のクラス1への所属確率の予測値 (0~1) $N$ : データ個数 <p>Kクラス分類の場合</p> $CE = -1/N \sum_{n=1}^N \sum_{k=1}^K t_{nk} \log \{q(Y=k X=x_n)\} \quad \dots (式6.3-5)$ $t_{nk}$ : 入力 $x_n$ のクラスkへの所属確率の実際の値 (0, 1) $q(Y=k X=x_n)$ : 入力 $x_n$ のクラスkへの所属確率の予測値 (0~1) $N$ : データ個数 ( $n=1 \sim N$ ) $K$ : 分類のクラス数 ( $k=1 \sim K$ )	クラス分類に一般的に用いられます。

【出典・参考】

損失関数⇒ [https://to-kei.net/neural-network/nn\\_loss\\_function/](https://to-kei.net/neural-network/nn_loss_function/)

損失関数⇒ [https://www.renom.jp/ja/notebooks/tutorial/basic\\_algorithm/lossfunction/notebook.html](https://www.renom.jp/ja/notebooks/tutorial/basic_algorithm/lossfunction/notebook.html)

尤度⇒「統計学」森棟, 照井, 中川, 西埜, 黒住 共著 有斐閣 2017年12月 改訂版第3刷

目的関数⇒[http://research.nii.ac.jp/~uno/opt\\_1.htm](http://research.nii.ac.jp/~uno/opt_1.htm)



## (7) モデルの構造

- ここでは「教師あり学習（分類）」のモデルについて、「階層型ニューラルネットワーク」のうち以下の2つを用いた分類モデルを紹介します。

- (1) 「順伝播型ニューラルネットワーク（ジョンペンガタニューラルネットワーク、Feed forward neural network、FFNN）」
- (2) 「畳み込みニューラルネットワーク（タミコニューラルネットワーク、Convolutional neural network、CNN）」

- 「階層型ニューラルネットワーク」は、基本的に以下の層から構成されます。

層名	説明
入力層	ニューラルネットワークへの情報を取り込む層
隠れ層	入力側からの信号に対して、重みづけ評価を行いながら、出力側へ信号伝達する層（中間層とも言います）
出力層	ニューラルネットワークが算出した結果を出力する層

## (7.1) 順伝播型ニューラルネットワーク

(※ この節の一部は、「セミナー2回目、7回目」からの再掲を含みます。)

- ・「順伝播型ニューラルネットワーク (フィードフォワードニューラルネットワーク、Feed forward neural network、FFNN)」は、以下の特徴を持つニューラルネットワークです。

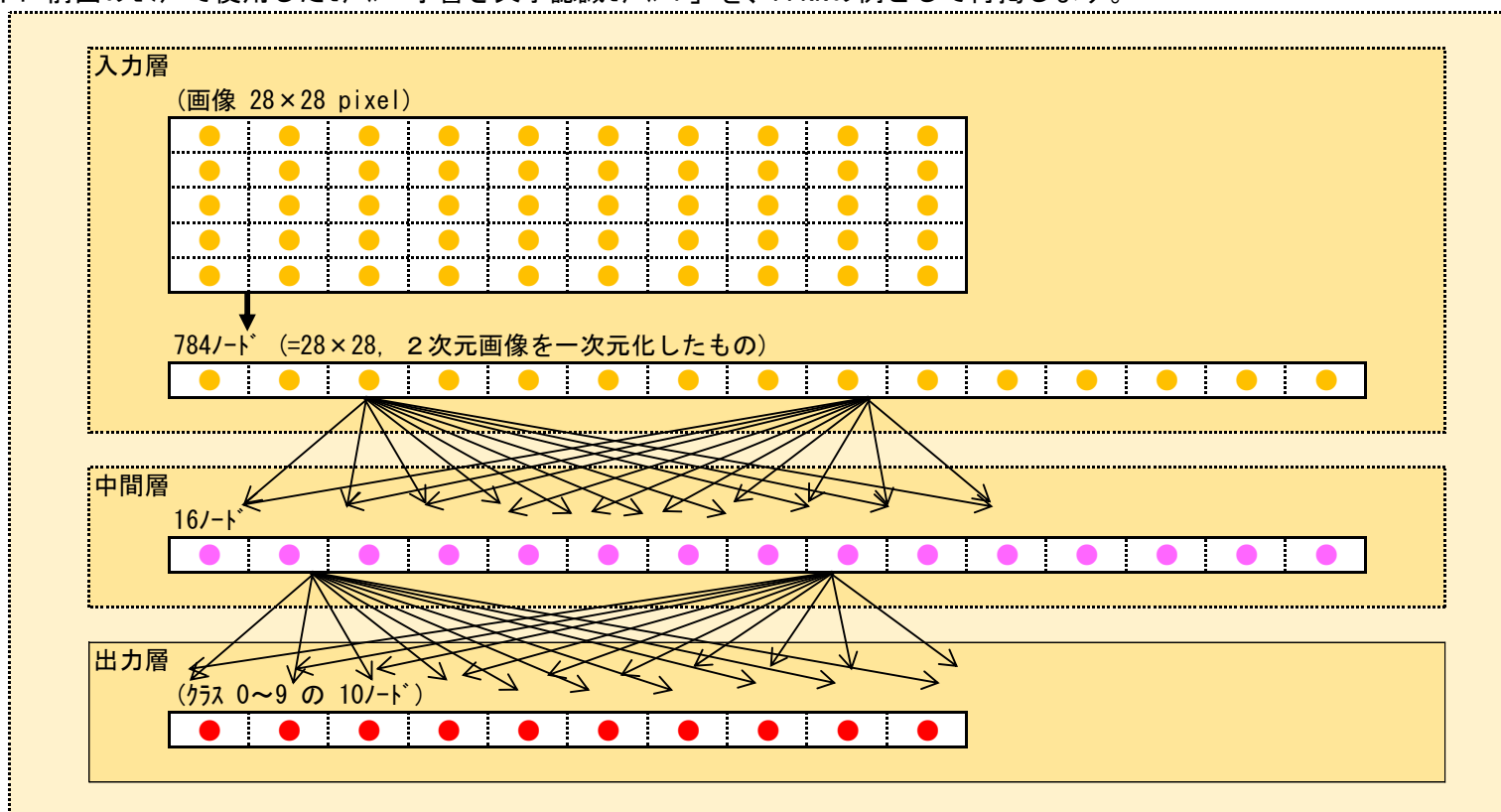
- (1) 隣接する層間のすべてのユニットが完全に接続されている (完全結合)
- (2) 入力層から出力層に向けて一方通行で流れる
- (3) 入力層と出力層の間に、1つ以上の隠れ層がある

- ・「順伝播型ニューラルネットワーク」で、隠れ層の層数を多くしたものが

「深層順伝播型ニューラルネットワーク (ディープフィードフォワードニューラルネットワーク、Deep Feed forward neural network、DFF)」です。

- ・ このタイプのネットワークは「誤差逆伝播法 (バックプロパゲーション、Backpropagation)」で学習します。

- ・ 以下に前回のセミナーで使用したモデル「手書き文字認識モデル1」を、FFNNの例として再掲します。



(入力画像)



(28×28 モノクロ)

(出力: クラスへの所属確率)

0: 0.45, 1: 0.01, 2: 0.06  
3: 0.05, 4: 0.02, 5: 0.18  
6: 0.01, 7: 0.01, 8: 0.21  
9: 0.01

・モデル「手書き文字認識モデル1」とその学習方法は以下のようなものです：

- (1) モデルは「順伝播型ニューラルネットワーク (FFNN)」であり、入力層から出力層に向けて一方通行で流れる (Sequential)。
- (2) 入力データは、(Y方向サイズ, X方向サイズ, 入力チャネル数) = (28, 28, 1) のモノクロ画像 (0~9の手書き数字1文字) であるが、二次元画像を一次元化した784 (=28×28)個の各データを、入力層の各ノードとする。  
各入力値は、0~255 の 256 階調のところを、255 で割って 0~1 に正規化したものを用いる。  
入力層の各ノードは中間層の各ノードとの間で、完全に接続させる (全結合)。
- (3) 入力と出力の間に、ノード数が 16 から構成される1つの中間層 (隠れ層) を置く。  
中間層の活性化関数として「ReLU (正規化線形関数, RectifiedLinear Unit)」を用いる。  
中間層の各ノードは出力層の各ノードとの間で、完全に接続させる (全結合)。
- (4) 出力層のノード数は、10 (入力画像の文字の種類 (0~9) の数) である。  
データ型は「float32 (入力画像枚数, 10)」で、入力画像の各クラスへの所属確率 [0~1] を表す。  
出力層の活性化関数として「ソフトマックス関数 (softmax)」を用いる。
- (5) 損失関数 (誤差を計算する関数) として、「交差エントロピー誤差 (cross entropy)」を用いる。
- (6) 学習方法 (モデルの最適化方法) として、勾配降下法的一种である「Adam」を用いる。  
学習時には「エポック数」として10、「バッチサイズ」として「1000」で学習を行う。

- ・モデル「手書き文字認識モデル1」の定義と学習の実装部分は以下のようなものです：

(リスト07-(03)-1\_手書き文字認識モデル1)からの抜粋と要約

```
IMG_SIZE_Y = 28          # 画像サイズ(Y方向)
IMG_SIZE_X = 28          # 画像サイズ(X方向)
IMG_CLASS_NO = 10        # 画像の分類クラス数(0~9 の10クラス)

#=====
# 「手書き文字認識モデル1」を定義する
#=====
np.random.seed(1)
predModel = Sequential()
predModel.add(Dense(16, input_dim=inTrain1.shape[1], activation=actvFunc))
predModel.add(Dense(IMG_CLASS_NO, activation='softmax'))
predModel.compile(loss='categorical_crossentropy',
                  optimizer=Adam(), metrics=['accuracy'])
predModel.summary()

#=====
# 「手書き文字認識モデル1」で学習する
#=====
epochNo = 10
batchSize = 1000

history = predModel.fit(inTrain1, outTrain1,
                        epochs=epochNo, batch_size=batchSize,
                        verbose=0, validation_data=(inTest1, outTest1))
score = predModel.evaluate(inTest1, outTest1, verbose=0)
```

- ・モデル「手書き文字認識モデル1」はモデルクラスの「summary()」メソッドによりそのサマリを出力でき、以下のようになります：

(リスト07-(03)-1\_手書き文字認識モデル1)の実行結果の抜粋と要約

Model: "sequential\_2"

Layer (type)	Output Shape	Param #
=====		
dense_3 (Dense)	(None, 16)	12560
=====		
dense_4 (Dense)	(None, 10)	170
=====		

Total params: 12,730

Trainable params: 12,730

Non-trainable params: 0

#### 【出典・参考】

⇒ <https://qiita.com/ma-oshita/items/99b2cf313494adbb964d>

⇒ <https://ja.wikipedia.org/wiki/ニューラルネットワーク#順伝播型ニューラルネットワーク>

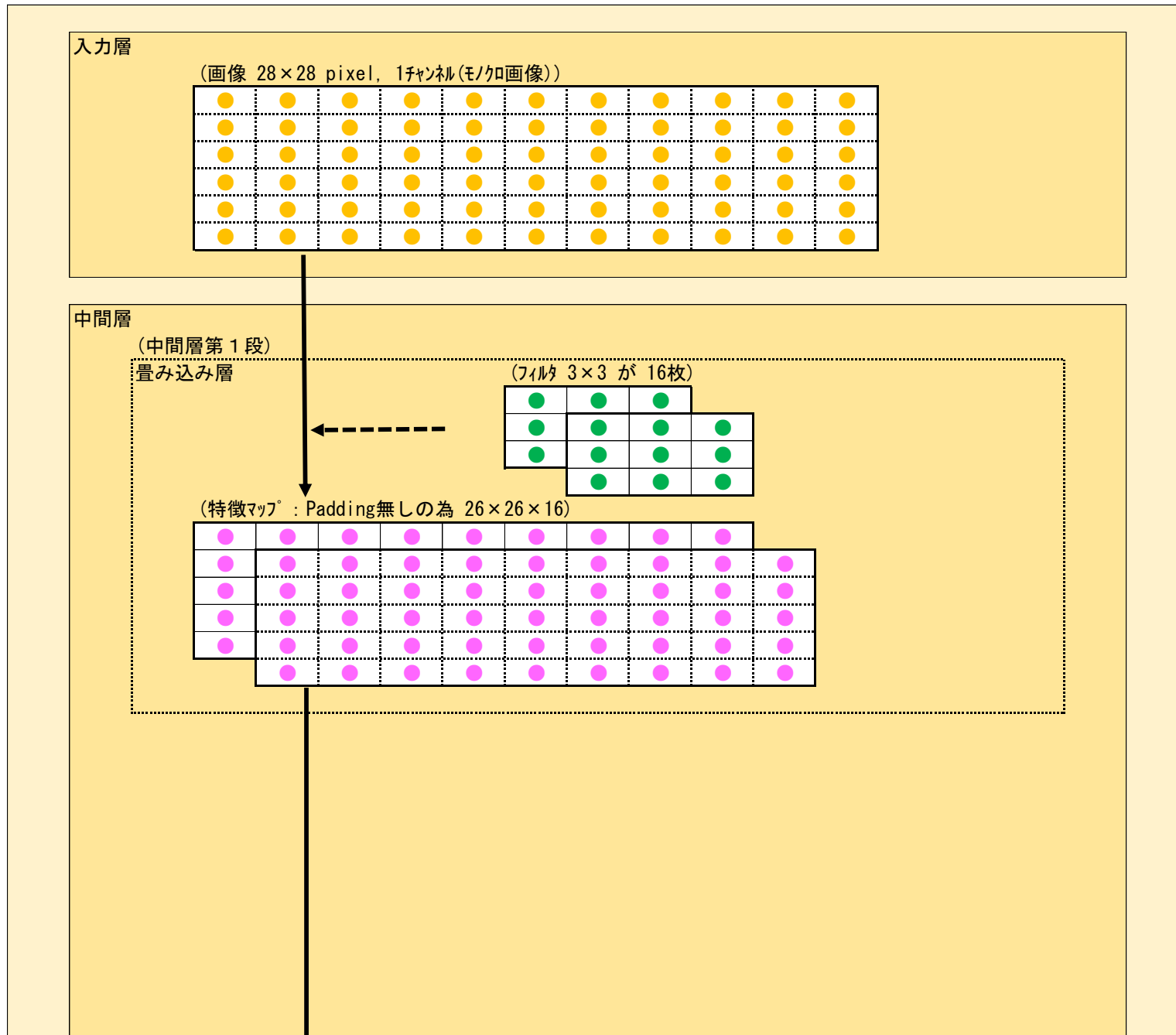
## (7.2) 畳み込みニューラルネットワーク

(※ この節の一部は、「セミナー2回目、7回目」からの再掲を含みます。)

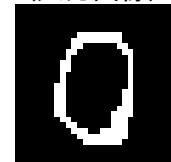
- ・「畳み込みニューラルネットワーク（タタミニューラルネットワーク、Convolutional neural network、CNN）」は、階層型ニューラルネットワークに、畳み込み構造を入れたモデルです。「深層畳み込みニューラルネットワーク（シンソウタタミニューラルネットワーク、Deep Convolutional neural network、DCN）」は、CNN の中間層（隠れ層）を深くしたものです。  
これは、人の視覚をモデルにした構成になっており、画像認識分野を中心に幅広く利用されているディープラーニング（深層学習）手法の一つです。
- ・「畳み込みニューラルネットワーク」には、以下のような特徴があります：
  - (1) 複雑なパターン認識問題にも、簡潔なネットワークで対応できます。
  - (2) 全体としてユニット数が少なくて済むので、計算が容易です。
- ・このモデルは、中間層（隠れ層）が、畳み込み層（コンボリューション層）、正規化層、プーリング層、全結合層で構成された複数の層からなります。  
正規化層とプーリング層は必須ではありません。また、畳み込み層でフィルタの「移動量（ドワリヨウ、ストライド、stride）」を調整することにより情報縮約を行って、プーリング層を代行することもあります。
- ・各層は以下のようなものです。

層名	説明
畳み込み層（タタミソウ、Convolution Layer）	入力側からの信号に対してフィルタ処理を行い「特徴マップ（トクショマップ、feature map）」を出力する層です。 コンボリューション層ともいいます。使用するフィルタは「カーネル (kernel)」「オペレータ」「マスク」とも言います。
プーリング層（プーリングソウ、Pooling Layer）	畳み込み層で作成された「特徴マップ」に対し情報縮約操作を行う層です。 プーリングには、以下のようなものがあります： <ul style="list-style-type: none"><li>・最大値プーリング（Max Pooling） 着目領域での最大値で縮約する</li><li>・平均値プーリング（Average Pooling） 着目領域での平均値で縮約する</li><li>・Lpプーリング（Lp Pooling） 着目領域での中央値を強調した値で縮約する</li></ul>
全結合層（ゼンカツゴウソウ、Fully Connected Layer）	隣接する層の各ユニット同士が、全ての組合せで結合することを全結合と言います。 通常、出力層と前段の数層を全結合とし、これを全結合層としてモデル化します。
正規化層（セイヤカソウ、Normalization Layer）	各サンプルの平均と分散をそれぞれ0と1になるように、入力側からの信号を変換する層です。 これで、入力データ間のばらつきを補正します。
Flatten層（フラットソウ、Flatten Layer）	前段の層から出てきた多次元のデータを一次元に変換する層です。 次段の層の各ユニットと全結合する準備などの為に入れます。

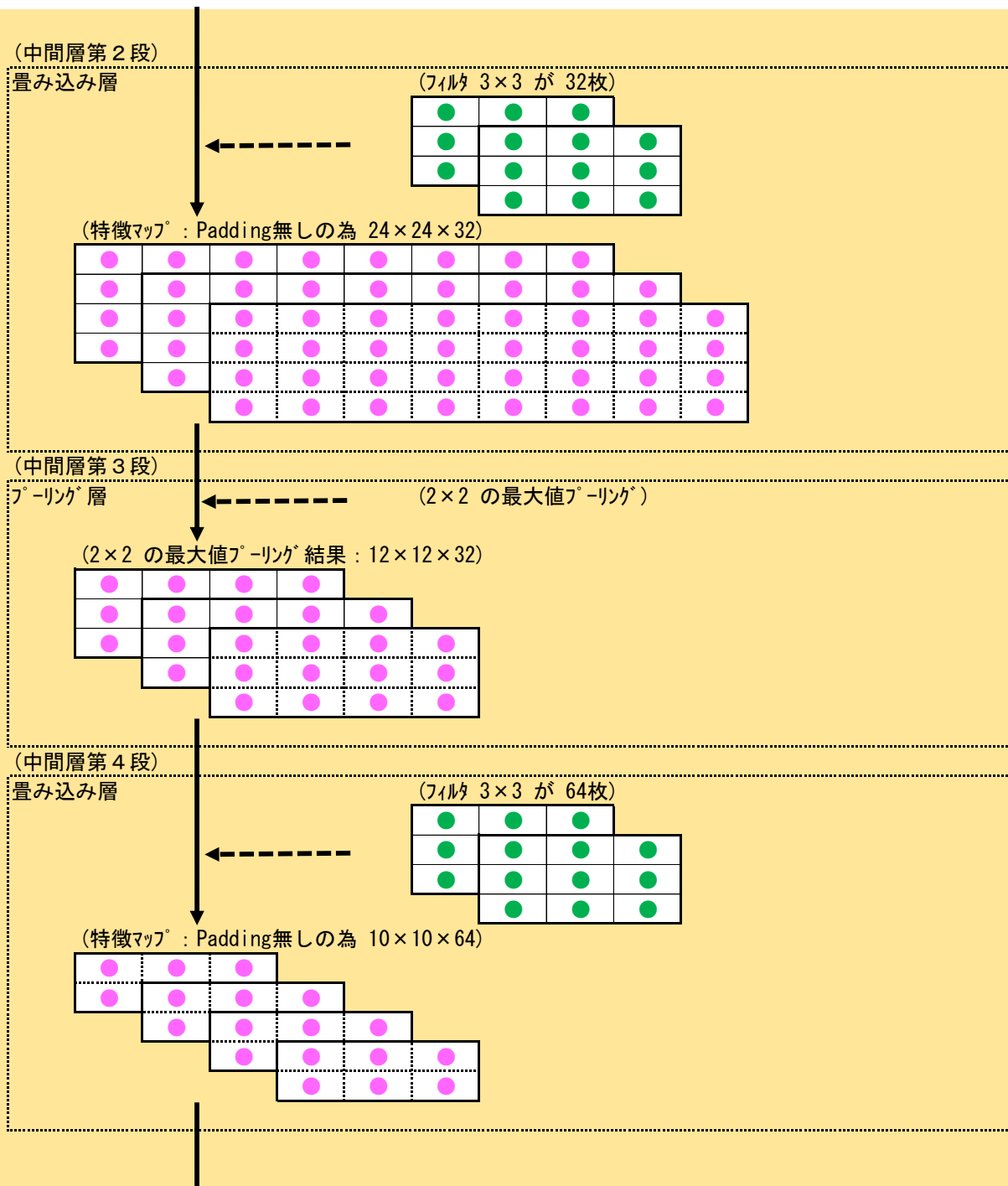
- ・ 以下に前回のセミナーで使用したモデル「手書き文字認識モデル3」を、CNNの例として再掲します。



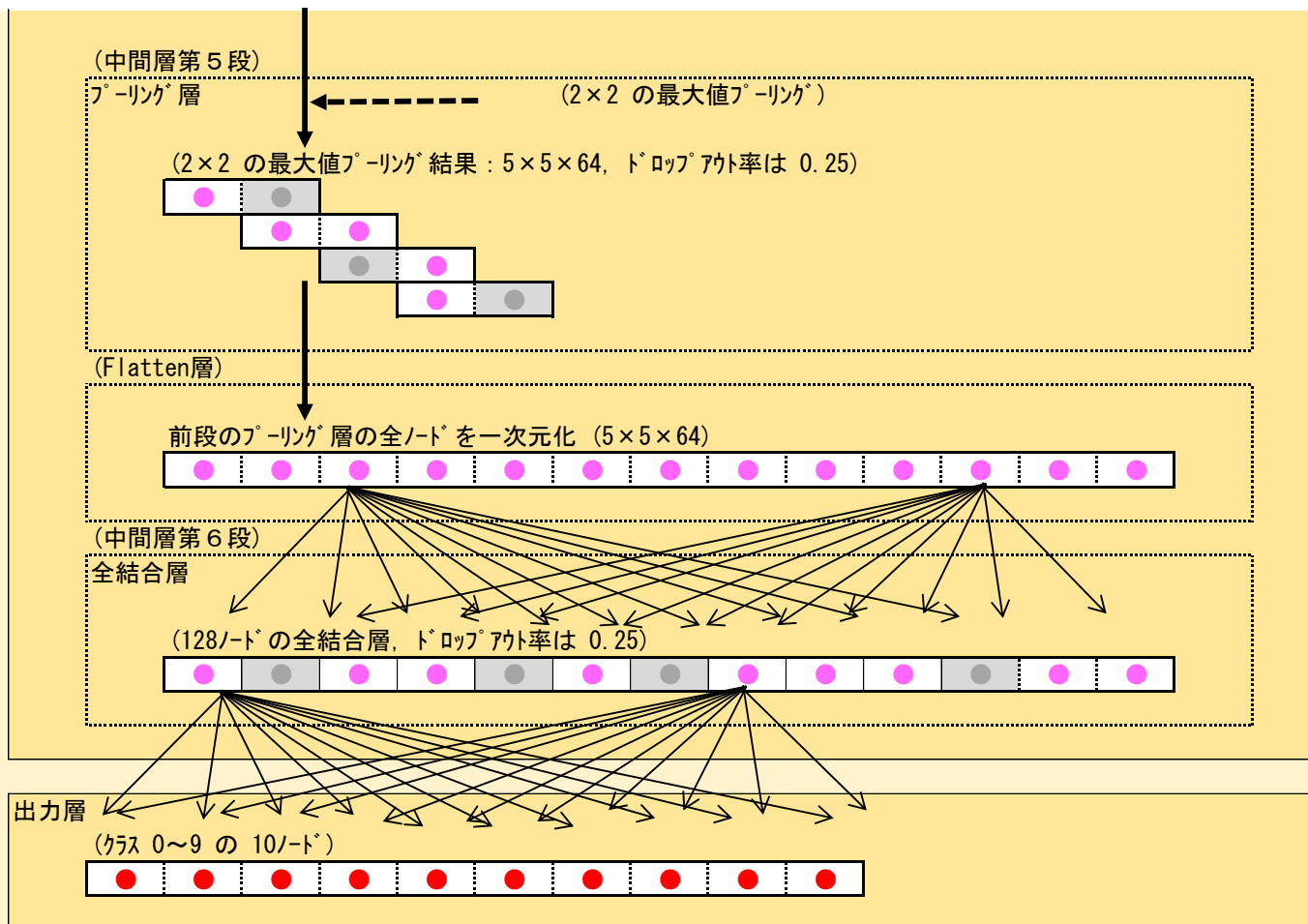
(入力画像)



( $28 \times 28$  モノクロ)







(出力 : クラスへの所属確率)

0: 0.99, 1: 0.00, 2: 0.00  
 3: 0.00, 4: 0.00, 5: 0.00  
 6: 0.00, 7: 0.00, 8: 0.00  
 9: 0.01

・モデル「手書き文字認識モデル3」とその学習方法は以下のようなものです：

- (1) モデルは「畳み込みニューラルネットワーク」であり、入力層から出力層に向けて一方通行で流れる(Sequential)。
- (2) 入力データは、(Y方向サイズ, X方向サイズ, 入力チャネル数) = (28, 28, 1) のモノクロ画像 (0~9の手書き数字1文字) であり、これをそのまま入力層の入力データとする。  
各入力値は、0~255 の 256 階調のところを、255 で割って 0~1 に正規化したものを用いる。
- (3) 入力と出力の間に、以下の6段で構成する中間層(隠れ層)を置く。
  - (3.1) 中間層第1段は、畳み込み層で、フィルタ  $3 \times 3$  を 16枚で、Padding無しで、畳み込む。
  - (3.2) 中間層第2段は、畳み込み層で、フィルタ  $3 \times 3$  を 32枚で、Padding無しで、畳み込む。
  - (3.3) 中間層第3段は、 $2 \times 2$ の最大値プーリング層とする。
  - (3.4) 中間層第4段は、畳み込み層で、フィルタ  $3 \times 3$  を 64枚で、Padding無しで、畳み込む。
  - (3.5) 中間層第5段は、 $2 \times 2$ の最大値プーリング層で、ドロップアウト率を 0.25 とする。  
この出力を、Flatten層を通して一次元化し、その各ノードを、中間層第5段の各ノードとの間で全結合させる。
  - (3.6) 中間層第6段は、全結合層で、ドロップアウト率を 0.25 とする。  
各ノードを、出力層の各ノードとの間で全結合させる。
  - (3.7) 中間層の活性化関数として「ReLU (正規化線形関数, Rectified Linear Unit)」を用いる。
- (4) 出力層のノード数は、10 (入力画像の文字の種類 (0~9) の数) である。  
データ型は「float32 (入力画像枚数, 10)」で、入力画像の各クラスへの所属確率 [0~1] を表す。  
中間層の出力を「Flatten」層を通すことで、「 $28 \times 28 \times 8 \text{枚} \times \text{入力チャネル数}$ 」という4次元のデータを一次元化して、出力層との間で全結合させる。  
出力層の活性化関数として「ソフトマックス関数 (softmax)」を用いる。
- (5) 損失関数(誤差を計算する関数)として、「交差エントロピー誤差 (cross entropy)」を用いる。
- (6) 学習方法(モデルの最適化方法)として、勾配降下法的一种である「Adam」を用いる。  
学習時には「エポック数」として10、「バッチサイズ」として「1000」で学習を行う。

- ・モデル「手書き文字認識モデル3」の定義と学習の実装部分は以下のようなものです：

(リスト07-(03)-3\_手書き文字認識モデル3)からの抜粋と要約

```
IMG_SIZE_Y = 28          # 画像サイズ(Y方向)
IMG_SIZE_X = 28          # 画像サイズ(X方向)
IMG_CLASS_NO = 10        # 画像の分類クラス数(0~9 の10クラス)

#=====
# 「手書き文字認識モデル3」を定義する
#=====
np.random.seed(1)
predModel = Sequential()
predModel.add(Conv2D(16, (3, 3), input_shape=(IMG_SIZE_Y, IMG_SIZE_X, 1), activation='relu'))
predModel.add(Conv2D(32, (3, 3), activation='relu'))
predModel.add(MaxPooling2D(pool_size=(2, 2)))
predModel.add(Conv2D(64, (3, 3), activation='relu'))
predModel.add(MaxPooling2D(pool_size=(2, 2)))
predModel.add(Dropout(0.25))
predModel.add(Flatten())
predModel.add(Dense(128, activation='relu'))
predModel.add(Dropout(0.25))
predModel.add(Dense(IMG_CLASS_NO, activation='softmax'))
predModel.compile(loss='categorical_crossentropy', optimizer=Adam(), metrics=['accuracy'])
predModel.summary()

#=====
# 「手書き文字認識モデル3」で学習する
#=====
epochNo = 10
batchSize = 1000
history = predModel.fit(inTrain3, outTrain3,
                        epochs=epochNo, batch_size=batchSize,
                        verbose=0, validation_data=(inTest3, outTest3))
score = predModel.evaluate(inTest3, outTest3, verbose=0)
```

- ・モデル「手書き文字認識モデル3」はモデルクラスの「summary()」メソッドによりそのサリを出力でき、以下のようになります：

( リスト07-(03)-3\_手書き文字認識モデル3 ) の実行結果の抜粋と要約

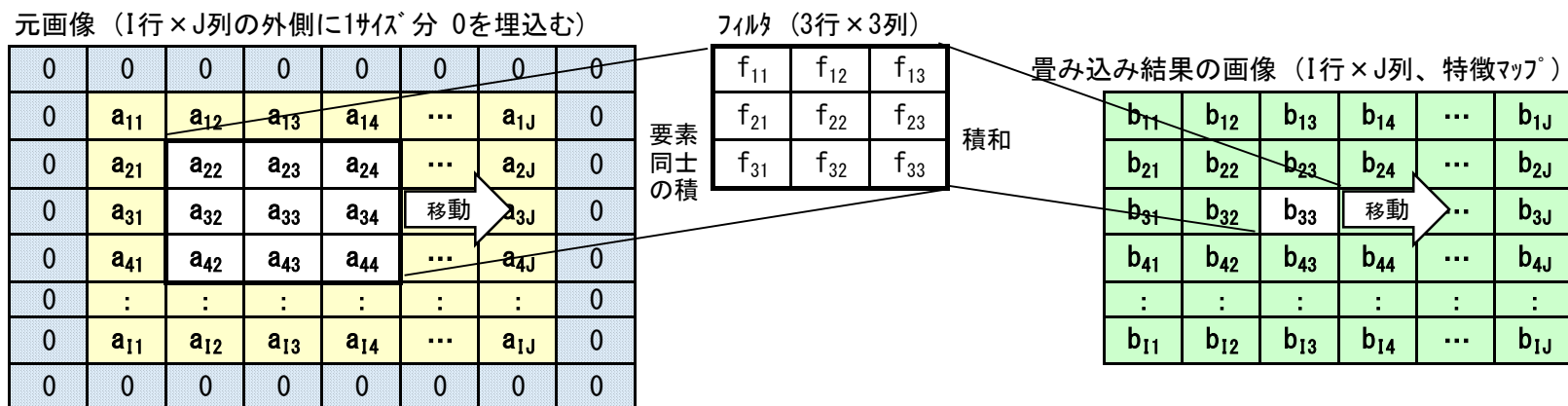
Model: "sequential_4"		
Layer (type)	Output Shape	Param #
=====		
conv2d_2 (Conv2D)	(None, 26, 26, 16)	160
conv2d_3 (Conv2D)	(None, 24, 24, 32)	4640
max_pooling2d_1 (MaxPooling2D)	(None, 12, 12, 32)	0
conv2d_4 (Conv2D)	(None, 10, 10, 64)	18496
max_pooling2d_2 (MaxPooling2D)	(None, 5, 5, 64)	0
dropout_1 (Dropout)	(None, 5, 5, 64)	0
flatten_2 (Flatten)	(None, 1600)	0
dense_6 (Dense)	(None, 128)	204928
dropout_2 (Dropout)	(None, 128)	0
dense_7 (Dense)	(None, 10)	1290
=====		
Total params: 229,514		
Trainable params: 229,514		
Non-trainable params: 0		
=====		

## (7.2.1) 畳み込み層

- ・「畳み込み層 (タピコミツ、Convolution Layer)」は、入力側からの信号に対して「畳み込み」によるフィルタ処理を行い、「特徴マップ (トクショウマップ, feature map)」を作成する層です。
- ・「畳み込み」によるフィルタ処理は、具体的には以下のように、元画像の領域にフィルタを重ねて、要素同士の積和を畳み込み結果の領域の中心位置の値とする演算処理で、これを全領域に渡って移動しながら、畳み込み結果画像を得るものです。

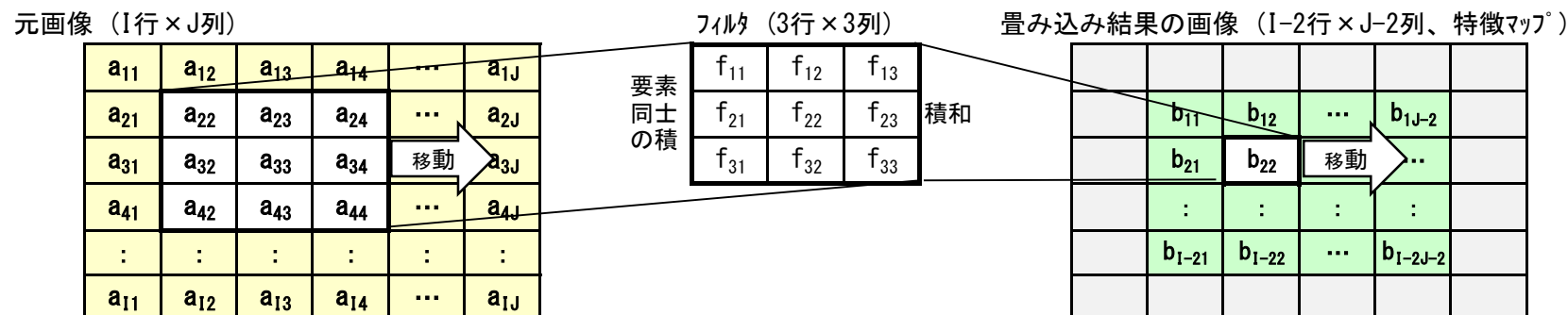
(1) padding 有りの場合 (元画像がI行×J列、フィルタが3行×3列、ストライドのサイズは1)

$$b_{ij} = \sum_{m=1}^3 \sum_{n=1}^3 a_{i+m-2, j+n-2} \times f_{mn} \quad i=1 \sim I, j=1 \sim J, \text{ 拡張部分 } (a_{0,*} / a_{*,0} / a_{I+1,*} / a_{*,J+1}) \text{ は } 0$$



(2) padding 無しの場合 (元画像がI行×J列、フィルタが3行×3列、ストライドのサイズは1)

$$b_{ij} = \sum_{m=1}^3 \sum_{n=1}^3 a_{i+m-1, j+n-1} \times f_{mn} \quad i=1 \sim I-2, j=1 \sim J-2$$



- ・「畳み込み」演算では、元画像の縁部分にフィルタをかけると、フィルタの一部が元画像の領域からはみ出してしまいます。  
この為、フィルタ演算結果として採用できるのは、元画像の領域からはみ出していない領域に限られ、  
フィルタ演算結果のサイズが元画像よりも小さくなってしまいます。

- ・これを避けるために、元画像の縁から外側に、フィルタのはみ出し部分がないサイズ分だけ0などを埋め込みます。  
この操作を「パディング (padding)」と言い、サイズを「パディングサイズ (padding size)」と言います。  
「パディング」によりフィルタ適用後画像（特徴マップ）のサイズが小さくなるのを防ぐことができます。

元画像 (I行×J列)

0	0	0				
0	$a_{11}$	$a_{12}$	$a_{13}$	$a_{14}$	...	$a_{1J}$
0	$a_{21}$	$a_{22}$	$a_{23}$	$a_{24}$	...	$a_{2J}$
	$a_{31}$	$a_{32}$	$a_{33}$	$a_{34}$	...	$a_{3J}$
	$a_{41}$	$a_{42}$	$a_{43}$	$a_{44}$	...	$a_{4J}$
	:	:	:	:	:	:
	$a_{I1}$	$a_{I2}$	$a_{I3}$	$a_{I4}$	...	$a_{IJ}$

フィルタ (3行×3列)

$f_{11}$	$f_{12}$	$f_{13}$
$f_{21}$	$f_{22}$	$f_{23}$
$f_{31}$	$f_{32}$	$f_{33}$

要素  
同士の  
積

- ...  $a_{11}$  を中心とした領域に  
フィルタ (3行×3列) をかける際に、  
フィルタは元画像の領域から一部がはみだしてしまいます。  
1サイズ分 (フィルタの半分のサイズ) 外側に0を埋めて対応します。

元画像 (I行×J列)

$a_{11}$	$a_{12}$	$a_{13}$	$a_{14}$	...	$a_{1J}$
$a_{21}$	$a_{22}$	$a_{23}$	$a_{24}$	...	$a_{2J}$
$a_{31}$	$a_{32}$	$a_{33}$	$a_{34}$	...	$a_{3J}$
$a_{41}$	$a_{42}$	$a_{43}$	$a_{44}$	...	$a_{4J}$
:	:	:	:	:	:
$a_{I1}$	$a_{I2}$	$a_{I3}$	$a_{I4}$	...	$a_{IJ}$

フィルタ (3行×3列)

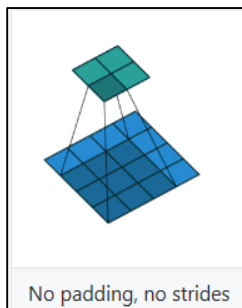
$f_{11}$	$f_{12}$	$f_{13}$
$f_{21}$	$f_{22}$	$f_{23}$
$f_{31}$	$f_{32}$	$f_{33}$

要素  
同士の  
積

- ...  $a_{22}$  を中心とした領域に  
フィルタ (3行×3列) をかける際に、  
フィルタは元画像の領域内にあるので、  
外側に0を埋める必要はありません。

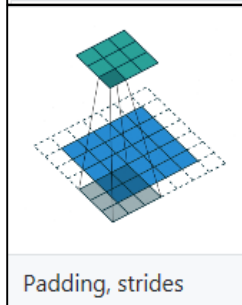
- ・「畳み込み」演算時に、全領域に渡って移動しながら元画像の領域にフィルタを重ねて演算を行います、この時のフィルタの「移動量（ドキュメント、ストライド、stride）」を指定できます。上記の説明では「ストライド」のサイズを「1」にしています。
- ・「移動量」を調整することにより結果的に情報縮約を行うことになり、これでプーリング層を代行することもあります。例えば、「ストライド」として「2」を指定すると、フィルタ適用後画像（特徴マップ）が、元画像の2分の1に縮んだものとなります。
- ・以下に畳み込み時のフィルタ・ストライド・パディングのパターンを幾つか、引用で示します。  
（「[https://github.com/vdumoulin/conv\\_arithmetic](https://github.com/vdumoulin/conv_arithmetic)」からの引用。参照するとアニメーションになっていて解り易いです）

例1



- ・（例1）は、 $4 \times 4$ の画像（下図）に対して、フィルタのサイズを $3 \times 3$ 、ストライドのサイズを $1 \times 1$ 、パディング無しで畳み込みを行った結果、 $2 \times 2$ の特徴マップが作成されたものです（上図）。

例2

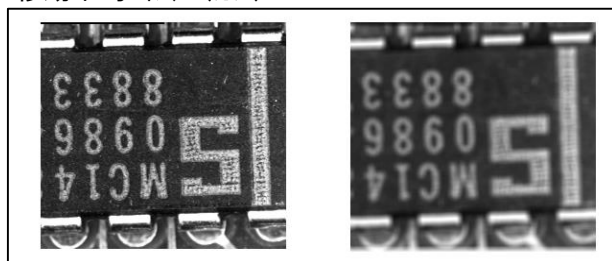


- ・（例2）は、 $5 \times 5$ の画像（下図）に対して、フィルタのサイズを $3 \times 3$ 、ストライドのサイズを $2 \times 2$ 、パディングあり（片側サイズ1）で畳み込みを行った結果、 $3 \times 3$ の特徴マップが作成されたものです（上図）。

- ・画像処理では様々なフィルタがあり、これによって画像解析の為の前処理を行っています。  
CNN では、学習を通してこのフィルタを自動作成します。

- ・「畳み込み」は、画像編集ソフトでは「効果」で「ぼかし」「エッジ検出」などの機能として組み込まれているもの多く見受けられます。  
これも画像処理フィルタであり、以下に幾つか紹介します。「<https://imaging-solution.net/imaging/filter-algorithm/>」からの引用です。

移動平均フィルタ (※1)



(3 × 3)

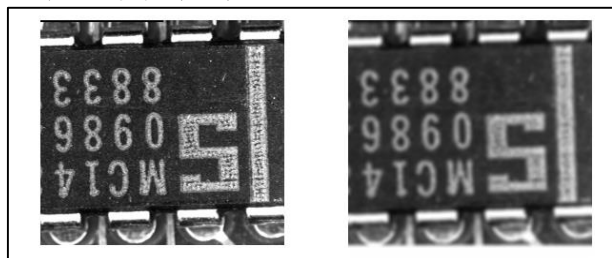
$\frac{1}{9}$	$\frac{1}{9}$	$\frac{1}{9}$
$\frac{1}{9}$	$\frac{1}{9}$	$\frac{1}{9}$
$\frac{1}{9}$	$\frac{1}{9}$	$\frac{1}{9}$

(5 × 5)

$\frac{1}{25}$	$\frac{1}{25}$	$\frac{1}{25}$	$\frac{1}{25}$	$\frac{1}{25}$
$\frac{1}{25}$	$\frac{1}{25}$	$\frac{1}{25}$	$\frac{1}{25}$	$\frac{1}{25}$
$\frac{1}{25}$	$\frac{1}{25}$	$\frac{1}{25}$	$\frac{1}{25}$	$\frac{1}{25}$
$\frac{1}{25}$	$\frac{1}{25}$	$\frac{1}{25}$	$\frac{1}{25}$	$\frac{1}{25}$
$\frac{1}{25}$	$\frac{1}{25}$	$\frac{1}{25}$	$\frac{1}{25}$	$\frac{1}{25}$

- ・移動平均フィルタ  
(平均化フィルタ、単に平滑化フィルタともいう)では、  
注目画素のその周辺の輝度値を用いて、  
輝度値の平均を求め、  
処理後画像の輝度値とするフィルタです。

ガウシアンフィルタ (※1)



(3 × 3)

$\frac{1}{16}$	$\frac{2}{16}$	$\frac{1}{16}$
$\frac{2}{16}$	$\frac{4}{16}$	$\frac{2}{16}$
$\frac{1}{16}$	$\frac{2}{16}$	$\frac{1}{16}$

(5 × 5)

$\frac{1}{256}$	$\frac{4}{256}$	$\frac{6}{256}$	$\frac{4}{256}$	$\frac{1}{256}$
$\frac{4}{256}$	$\frac{16}{256}$	$\frac{24}{256}$	$\frac{16}{256}$	$\frac{4}{256}$
$\frac{6}{256}$	$\frac{24}{256}$	$\frac{36}{256}$	$\frac{24}{256}$	$\frac{6}{256}$
$\frac{4}{256}$	$\frac{16}{256}$	$\frac{24}{256}$	$\frac{16}{256}$	$\frac{4}{256}$
$\frac{1}{256}$	$\frac{4}{256}$	$\frac{6}{256}$	$\frac{4}{256}$	$\frac{1}{256}$

- ・ガウシアンフィルタ  
は、注目画素に近いほど、  
平均値を計算するときの重みを大きくし、  
遠くなるほど重みを小さくなるようにガウス分布の  
関数を用いて割合を計算しているフィルタです。  
人間の目の焦点の当て方に近いと言われています。

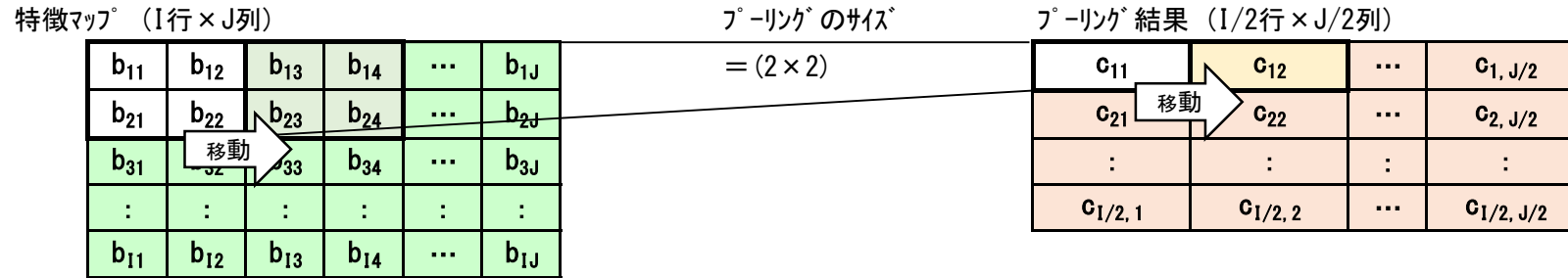
(※1) 左が原画で右がフィルタ処理後の画像



## (7.2.2) プーリング層

- 「プーリング層（プーリングゆ、Pooling Layer）」は、畳み込み層で作成された「特徴マップ」に対し情報縮約操作を行う層です。  
「プーリング」には、「最大値プーリング」「平均値プーリング」「Lpプーリング」といったものがあります。  
何れも元画像の全領域を小領域（ $2 \times 2$ 、 $3 \times 3$ などサイズを任意に指定）に分けて、各小領域での特徴量（最大値、平均値など）で縮約します。  
縮約した結果のサイズは、元画像のサイズを小領域のサイズで割ったものになります。  
これにより、特徴として重要な情報を残しながら元の画像を縮小し、画像の位置ずれに対する頑強さも持つことができます。

- 以下にプーリングのサイズが $2 \times 2$ の場合について示します。



- 以下に主なプーリングを一覧にして示します。

No.	プーリング名	英語	式（プーリングのサイズ $=2 \times 2$ の場合）	説明
1	最大値プーリング (サイズイチプーリング)	Max Pooling	$c_{ij} = \max \{ b_{m,n} \}_{m=2(i-1)+1 \sim 2(i-1)+2, n=2(j-1)+1 \sim 2(j-1)+2}$ $i=1 \sim I/2, j=1 \sim J/2$	元画像の全領域を、各小領域での最大値で縮約するものです。
2	平均値プーリング (エイベンチプーリング)	Average Pooling	$c_{ij} = 1/4 \times \sum_{m=2(i-1)+1}^{2(i-1)+2} \sum_{n=2(j-1)+1}^{2(j-1)+2} b_{m,n}$ $i=1 \sim I/2, j=1 \sim J/2$	元画像の全領域を、各小領域での平均値で縮約するものです。
3	Lpプーリング (エルピープーリング)	Lp Pooling	$c_{ij} = \{ 1/4 \times \sum_{m=2(i-1)+1}^{2(i-1)+2} \sum_{n=2(j-1)+1}^{2(j-1)+2} (b_{m,n})^p \}^{1/p}$ $i=1 \sim I/2, j=1 \sim J/2$ <p><math>P</math> はべき指数</p>	元画像の全領域を、各小領域での中央値を強調した値で縮約するものです。

## 【出典・参考】

畳み込み⇒ [https://github.com/vdumoulin/conv\\_arithmetic](https://github.com/vdumoulin/conv_arithmetic)

CNN⇒ [https://deeppage.net/deep\\_learning/2016/11/07/convolutional\\_neural\\_network.html](https://deeppage.net/deep_learning/2016/11/07/convolutional_neural_network.html)

DCNN⇒ [http://www.nlab.ci.i.u-tokyo.ac.jp/pdf/CNN\\_survey.pdf](http://www.nlab.ci.i.u-tokyo.ac.jp/pdf/CNN_survey.pdf)

DCNN⇒ <https://papers.nips.cc/paper/4824-imagenet-classification-with-deep-convolutional-neural-networks.pdf>

空間フィルタ⇒ [http://www.clg.niigata-u.ac.jp/~lee/jyugyou/img\\_processing/medical\\_image\\_processing\\_03\\_press.pdf](http://www.clg.niigata-u.ac.jp/~lee/jyugyou/img_processing/medical_image_processing_03_press.pdf)

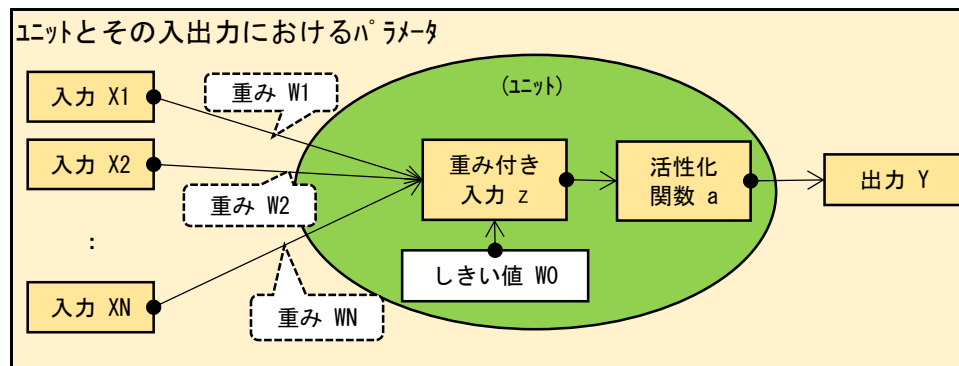
フィルタ⇒ <https://codezine.jp/article/detail/129>

プーリング⇒ <http://pynote.hatenablog.com/entry/dl-pooling>

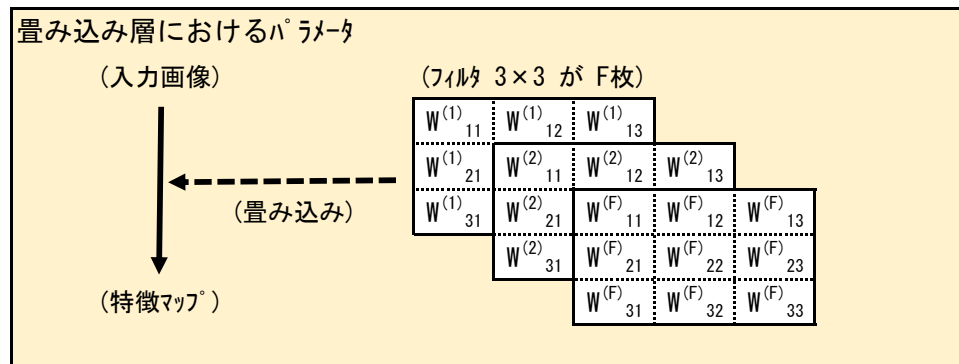
プーリング⇒ 「深層学習 Deep Learning」 人工知能学会監修（近代科学社 2015年11月）

## (8) モデルの学習

- 教師あり学習では、説明変数と目的変数の関係を説明するモデルを作成し、「モデルの出力結果(予測)」と「既知の結果データ」との差が最小化する(コスト関数(損失関数)の値が最小になる)ように、モデルの「パラメータ」を最適化します。これが「学習」です。
- モデルの「パラメータ」は、各ユニットの重み $W_i$ ／しきい値 $W_0$ ／フィルタの各値 $W_{ij}^{(f)}$ などの調整可能な値です。



- 左図で、ユニットとその入出力におけるパラメータは、ユニットの重み  $W_i$ 、しきい値  $W_0$  です。  
( $i=1\sim N$ : 入力の数分)



- 左図で、畳み込み層におけるパラメータは、フィルタの各値  $W_{ij}^{(f)}$  です。  
( $f=1\sim F$ : フィルタ枚数分、 $i=1\sim 3$ ,  $j=1\sim 3$ : フィルタのサイズ分)

- では、どのようにモデルの「パラメータ」を最適化するのでしょうか？  
以下に、この最適化の方法を紹介します。

## (8.1) 勾配降下法と学習率

(※ この節の一部は、「セミナー5回目」からの再掲を含みます。)

- 最適化対象のパラメータを  $\theta$  で表します。パラメータ  $\theta$  の関数であるコスト関数  $C(\theta)$  が与えられた時、これを最小にする  $\theta^*$  を求める最適化を、「勾配降下法 (コバニウカキ, gradient descent method)」で行います。

- 「勾配降下法」は、  
現ステップ (t) のパラメータ  $\theta^{(t)}$  における勾配と学習率  $\eta_t$  から、  
次ステップ (t+1) のパラメータ  $\theta^{(t+1)}$  を次式により更新して、  
コスト関数  $C(\theta)$  が最小になる  $\theta^*$  を逐次近似により求める方法です。

勾配降下法によるパラメータの最適化

$$\theta^{(t+1)} = \theta^{(t)} - \eta_t \times \partial C(\theta) / \partial \theta |_{\theta = \theta^{(t)}}$$

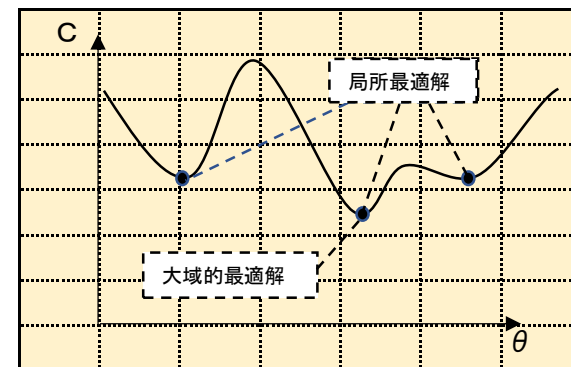
…(式8.1-1)

- $C(\theta)$  : コスト関数 (パラメータ  $\theta$  の関数)
- $\theta^{(t)}$  : ステップ t におけるパラメータ  $\theta$  の値
- $\eta_t$  : ステップ t における学習率  $\eta$  の値 ( $>0$ )
- $\partial C(\theta) / \partial \theta |_{\theta = \theta^{(t)}}$  :  $\theta^{(t)}$  におけるコスト関数の勾配

(「勾配降下法」の詳細な説明は「セミナー5回目」を参照してください。)

- 勾配降下法によるパラメータの最適化の(式8.1-1)では、「学習率 (ガクシュリツ, learning rate)」と呼ばれるパラメータ  $\eta$  (イタ) があります。  
これは、勾配降下法で、現ステップ (t) のパラメータ  $\theta^{(t)}$  におけるコスト関数  $C(\theta)$  の勾配から、次ステップ (t+1) のパラメータ  $\theta^{(t+1)}$  を計算する時に、勾配をどの程度の割合で変化分として反映させるか? という割合を示しています。  
「学習率」は「学習係数」とも言います。

- 「学習率」を大きくすると、最適解を通り過ぎて見落としてしまう可能性が高くなり、  
「学習率」を小さくすると、最適解に辿り着くのに時間がかかったり、  
最小値でない極小値にはまってしまうと抜け出られない場合があるという欠点があります。  
(局所的に最小値を与える解を「局所最適解 (キョクショサイケイ, local optimal solution)」と言いますが、極小値が最小値とは限りません。)  
これを「局所解問題 (キョクショカイレンダイ, local mininum problem)」と言います。  
尚、全領域にわたり最小値となる解を「大域的最適解 (タイキョクサイケイ, global optimal soution)」と言います。



- ・以降に述べる様々な改良版の「勾配降下法」と区別する場合、上記の「勾配降下法」を「バッチ勾配降下法 (バッチ勾配降下法, Batch gradient descent method)」と呼ぶこともあります。
- ・「バッチ勾配降下法」は、学習データが多いと計算コストがとて大きくなってしまいます。また、学習データが増えるたびに全ての学習データで再学習が必要になってしまうという欠点もあります。

( 勾配降下法による最小値探索の例 )

- ・以下の図は多変数 ( $x_1, x_2$ ) の関数  $y$  と、その勾配降下法による最小値探索の例です。

$$y(x_1, x_2) = x_1^2 + x_2^2 + 3x_2 + 4$$

この関数が構成する面の勾配ベクトル  $\nabla y$  は、次式で与られます。

$$\nabla y = (\partial y / \partial x_1, \partial y / \partial x_2) = (2x_1, 2x_2 + 3)$$

下図では、丸点とそれをつなぐ矢印は、初期位置  $r_0 = (x_{10}, x_{20})$  から始まって、

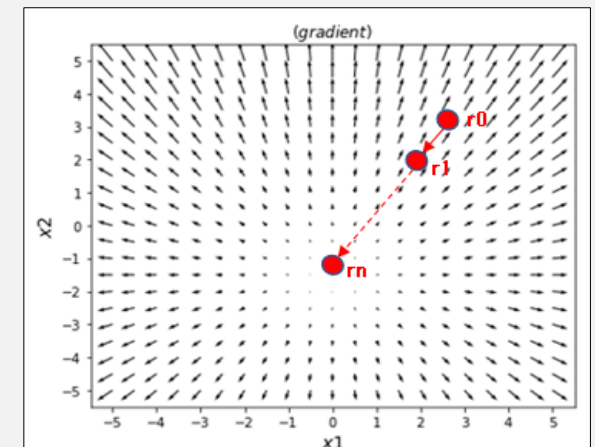
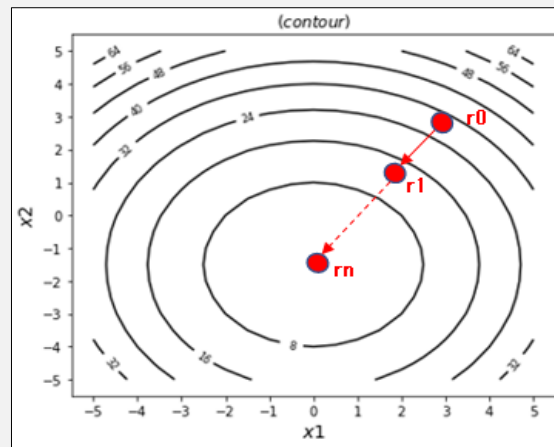
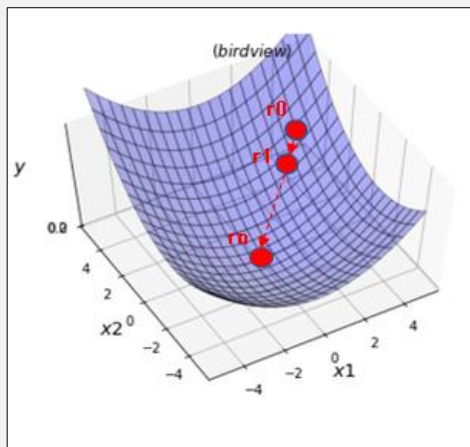
$$r_{i+1} = r_i - \eta (\partial y / \partial x_1, \partial y / \partial x_2) \quad \eta (>0) : \text{学習率}$$

という移動を繰り返して、関数値が最小の点  $r_n = (x_{1n}, x_{2n})$  に落ち着く様子を示しています。

・・・ (birdview)、(contour) の図

・・・ (gradient) の図

・・・ 各図

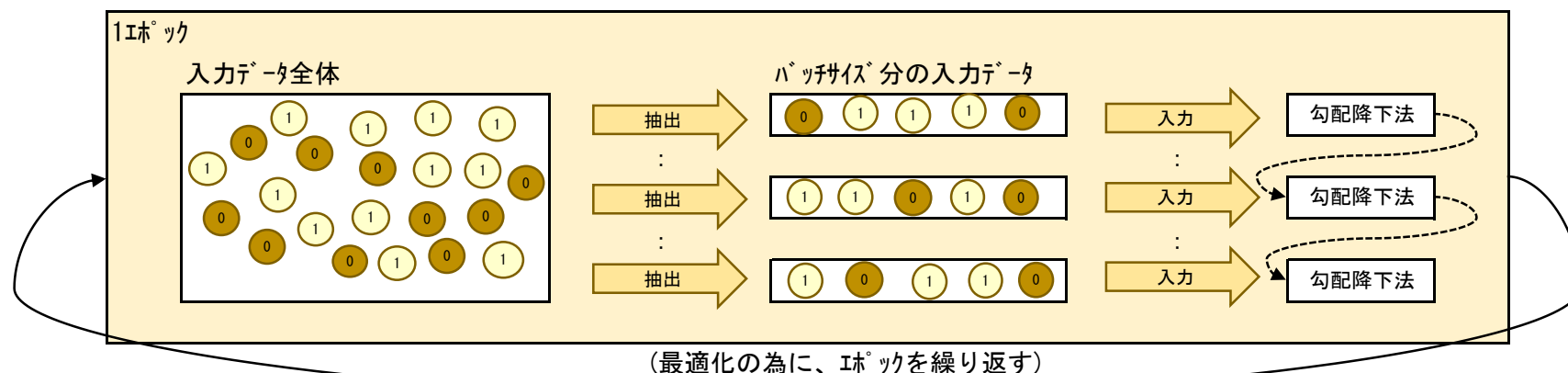


## (8.2) 勾配降下法の改良版

- ・「バッチ勾配降下法」の欠点を改善するために、様々な改良版が考案されていますので幾つか紹介します。

### (8.2.1) 確率的勾配降下法

- ・入力データ全てに対する誤差関数の勾配を、更新の1ステップごとに計算する場合、データが大きいとその計算に時間がかかってしまいます。そのような場合には、データの一部で誤差関数の勾配を計算する、「確率的勾配降下法 (カリツテキコバ イコウホウ, Stochastic gradient descent method, SGD)」という方法が使われます。「確率的勾配降下法」は「確率的勾配法」とも言います。
- ・「確率的勾配法」では入力データ全体の一部で計算します。計算1回分に使用するデータの個数を「バッチサイズ (batch size)」と言います。「バッチサイズ」分の入力データを使用した計算で、パラメータが1回更新されることになります。入力データ全体を用いて全体の誤差が最小になる方向へ進む「勾配降下法」とは異なり、「バッチサイズ」分の入力データを使用した「確率的勾配法」は、訓練データを少なくする為、勾配の推定が不安定となり、ふらつきながら徐々に誤差の少ない方向へ進むことになります。逆にこの不安定さのお陰で、局所解を抜け出せることもありうる (局所解問題の解消)、という長所があります。
- ・入力データ全てを使いパラメータ更新を行うとき、この更新の単位を「エポック (epochs)」で表します。「確率的勾配法」では、1「エポック」のパラメータ更新を行うのに「バッチサイズ」分の入力データを使用した更新を何度か繰り返します。そしてパラメータが最適化されるまで「エポック」を何度か繰り返します (モデルによっては収束しない場合もあります)。



(例：MNIST データセット・・・セミナ7回目を参照)

- ・ 訓練データの件数は 60000 であり、バッチサイズを 1000 とすると、  
訓練データを全て使うのに、 $60000/1000=60$  回のパラメータ更新が行われます。
- ・ これが 1エポックでのパラメータ更新です。  
このエポックを何度か繰り返して、パラメータの最適値を計算します。

### (8.2.2) AdaGrad

- ・ 「AdaGrad (エイダグランド)」は、2011年にJohn Duchiらが提唱した手法であり、学習率を自動で調整するアルゴリズムです。  
ステップ (t)が進むにつれて、学習率  $\eta_t$  が徐々に小さくなるように調整します。

AdaGradで、ステップ (t)における学習率  $\eta_t$

$$\eta_t = \eta_0 / \left\{ \sqrt{\sum_{j=1}^t \left( \partial C(\theta) / \partial \theta \mid_{\theta = \theta(j)} \right)^2} \right\}$$

…(式8.2-2)

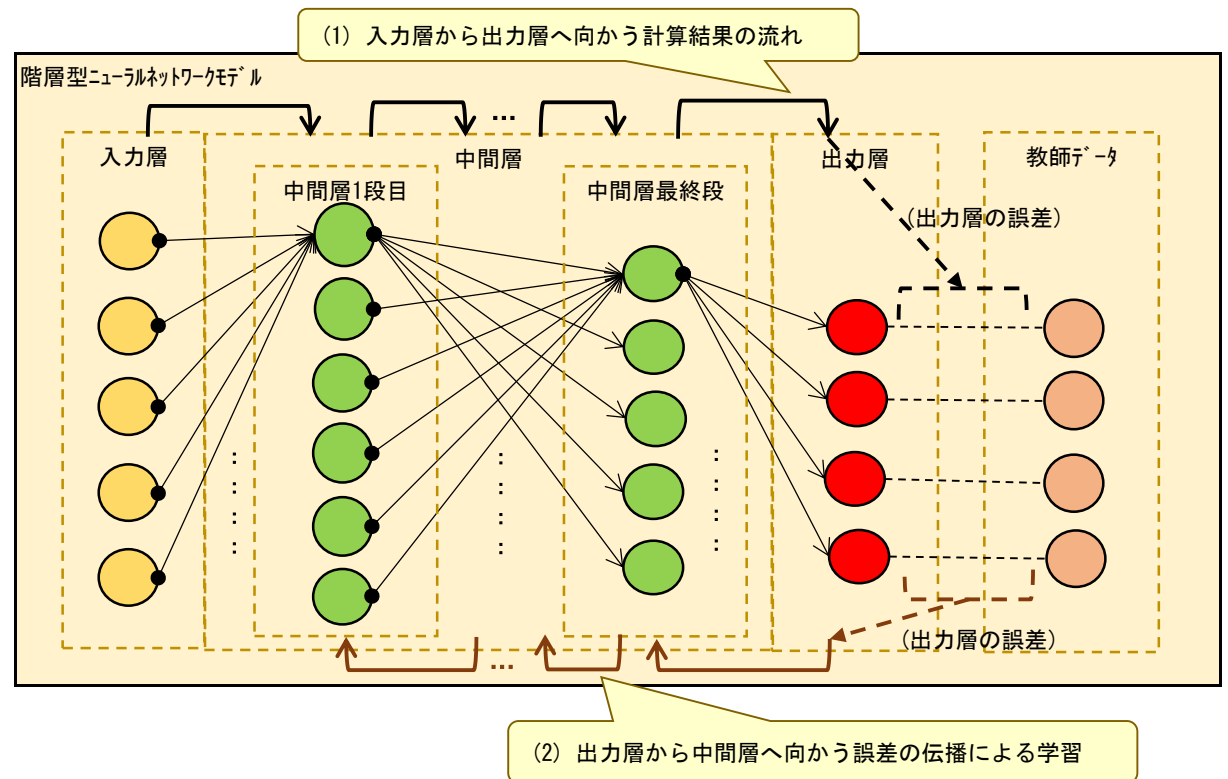
- $C(\theta)$  : コスト関数 (パラメータ  $\theta$  の関数)
- $\theta$  : 最適化対象のパラメータ
- $\eta_0$  : 初期学習率
- $\partial C(\theta) / \partial \theta \mid_{\theta = \theta(j)}$  :  $\theta(j)$ におけるコスト関数の勾配
- $j$  : ステップ 参照用添え字

### (8.2.3) Adam

- ・ 「Adam (アダム, Adaptive moment estimation)」は、2015年にDiederik P. Kingma 氏が提唱した手法で、  
AdaGrad や RMSProp, AdaDelta を改良したものです。  
アルゴリズムは『「深層学習 Deep Learning」近代科学社 人工知能学会監修 2015年11月』などに解説がありますので、参照して下さい。
- ・ 「Adam」は、セミナ7回目の実装例でも使用しています。

### (8.3) 誤差逆伝搬法

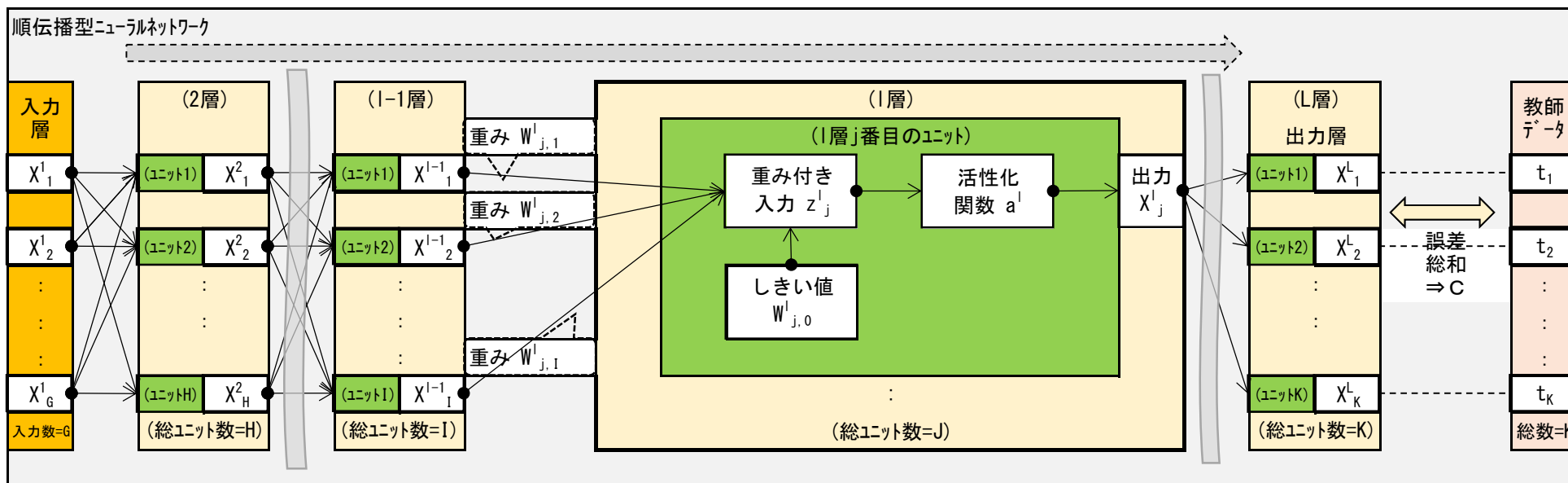
- ・「階層型ニューラルネットワーク」の学習方法として「誤差逆伝播法 (コシヤゲンパキ、Backpropagation method)」が有名です。  
手法そのものは以前からありましたが、1986年に「backwards propagation of errors (後方への誤差伝播)」の略からデビッド・ラumlハートらによって命名されたものです。
- ・「誤差逆伝播法」では、出力層での誤差 (階層型ネットワークの出力で生じる教師データとの差) を使って、  
モデルの出力の計算方向が中間層⇒出力層の方向なのとは逆向きに、  
出力層⇒中間層の方向へと、勾配降下法を用いてパラメータを更新してゆきます (下図)。
- ・勾配降下法では、誤差関数の偏微分を計算しますが、これを偏微分の「連鎖律」を用いて解いてゆきます。  
「連鎖律」の説明は「セナ-5回目」を御参照下さい。





## (8.4) 「順伝播型ニューラルネットワーク」の誤差逆伝搬法

ここでは「順伝播型ニューラルネットワーク (Feed forward neural networks、FF, FFNN)」(下図)について、誤差逆伝搬法をより詳しく見てゆきます。



・これより以降は、以下のような記法を用いることにします。添え字が多いですがご了承下さい。

- L : 層の総数 (入力層を第1層、出力層を第L層とします)
- l : 層番号 (l=2~L、入力層は層の厚さとしては数えませんが、説明上第1層とします)
- I : l-1層の総ユニット数
- J : l層の総ユニット数
- i : l-1層でのユニット番号 (i=1~I)
- j : l層でのユニット番号 (j=1~J)
- $X^{l-1}_i$  : l-1層i番目のユニットの出力値 (=l層への入力信号)
- $W^l_{j,i}$  : l層j番目のユニットについて、l-1層i番目のユニットからの入力信号の重み
- $W^l_{j,0}$  : l層j番目のユニットについて、l-1層からの入力信号に付加するバイアス
- $z^l_j$  : l層j番目のユニットの重み付き入力 (  $z^l_j = \sum_{i=0}^I W^l_{j,i} * X^{l-1}_i$  (但し  $X^{l-1}_0=1$ : バイアス項) )
- $a^l$  : l層の活性化関数 (activation function)
- $X^l_j$  : l層j番目のユニットの出力値 (  $X^l_j = a^l(z^l_j)$  )
- C : コスト関数 (平均二乗誤差(MSE)の場合  $C = (1/K) \sum_{k=1}^K (t_k^L - X_k^L)^2$  (K: は出力層のユニット数) )

## (8.4.1) 各層のユニットの誤差

- ・はじめに「 $l$ 層 $j$ 番目のユニットの誤差」 $\delta^l_j$ なる概念を導入します。

これは、次の定義で与えられるもので、 $l$ 層 $j$ 番目のユニットの重み付き入力 $z^l_j$ が、コスト関数 $C$ に与える変化率を意味しています。

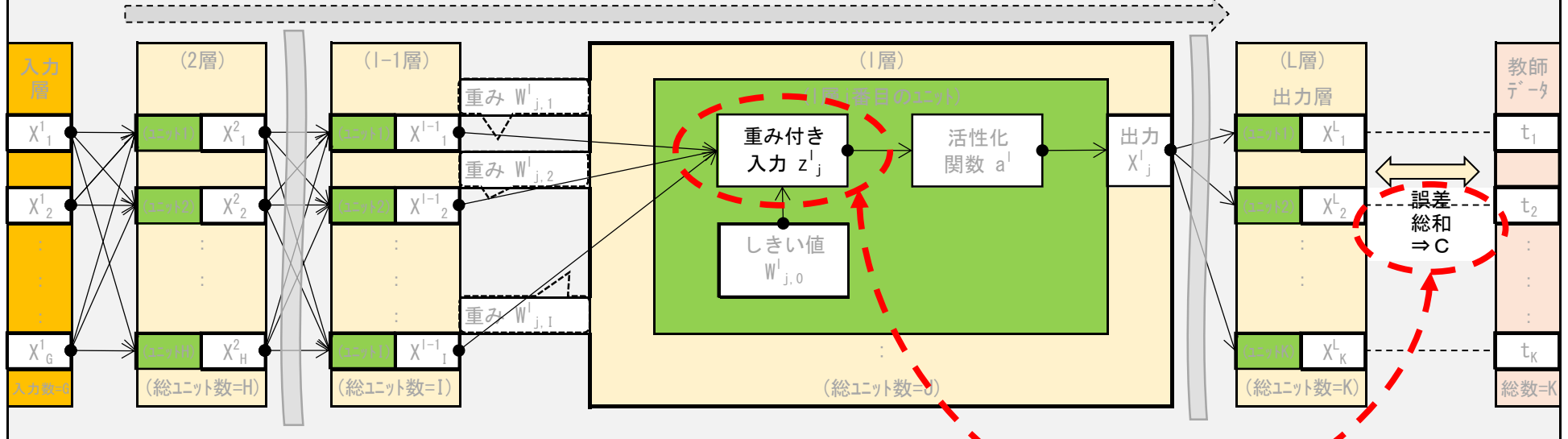
「誤差逆伝播法」では、このユニットの誤差「 $\delta^l_j$ 」についての逆漸化式を解くことにより、パラメータの最適化を行います。

「 $l$ 層 $j$ 番目のユニットの誤差」の定義

$$\delta^l_j = \partial C / \partial z^l_j \quad \dots(\text{式8.4-1})$$

- $l$  : 層番号 ( $l=2 \sim L$ 。説明上入力層を第1層とします)
- $j$  :  $l$ 層でのユニット番号 ( $j=1 \sim J$ )
- $\delta^l_j$  :  $l$ 層 $j$ 番目のユニットの誤差
- $C$  : コスト関数
- $z^l_j$  :  $l$ 層 $j$ 番目のユニットの重み付き入力

順伝播型ニューラルネットワーク



$$\delta^l_j = \partial C / \partial z^l_j$$

## (8.4.2) ユニットの誤差とコスト関数の偏微分との関係式

- 「 $l$ 層 $j$ 番目のユニットの誤差 $\delta_j^l$ 」と、コスト関数 $C$ の重み $W_{j,i}^l$ に関する偏微分について、以下の関係があります：

「 $l$ 層 $j$ 番目のユニットの誤差 $\delta_j^l$ 」と、コスト関数 $C$ の重み $W_{j,i}^l$ に関する偏微分との関係式

$$\partial C / \partial W_{j,i}^l = \delta_j^l X^{l-1}_i \quad \dots(\text{式8.4-2})$$

特に  $i=0$  (バイアス用のユニット番号) の場合、

$$\partial C / \partial W_{j,0}^l = \delta_j^l \quad \dots(\text{式8.4-3})$$

$C$  : コスト関数  
 $l$  : 層番号 ( $l=2 \sim L$ 、説明上、入力層を第1層とします)  
 $j$  :  $l$ 層でのユニット番号 ( $j=1 \sim J$ )  
 $i$  :  $l-1$ 層でのユニット番号 ( $i=1 \sim I$ ,  $i=0$  はバイアス用のユニット番号)  
 $W_{j,i}^l$  :  $l$ 層 $j$ 番目のユニットについて、 $l-1$ 層 $i$ 番目のユニットからの入力信号の重み  
 $W_{j,0}^l$  :  $l$ 層 $j$ 番目のユニットについて、 $l-1$ 層からの入力信号に付加するバイアス  
 $X^{l-1}_i$  :  $l-1$ 層 $i$ 番目のユニットの出力値 ( $l$ 層への入力信号)  
 $X^{l-1}_0$  :  $l-1$ 層からのユニットの出力のバイアス用の出力値 (1固定)  
 $\delta_j^l$  :  $l$ 層 $j$ 番目のユニットの誤差

(説明)

$l$ 層 $j$ 番目のユニットの重み付き入力  $z_j^l$  は、次式で与えられます

$$z_j^l = \sum_{i=0}^I W_{j,i}^l * X^{l-1}_i$$

これより

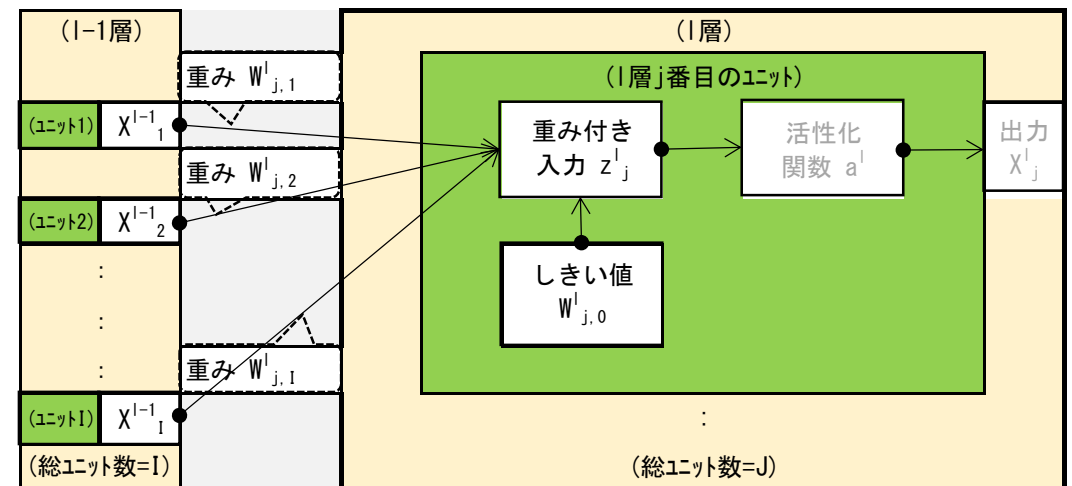
$$\partial z_j^l / \partial W_{j,i}^l = X^{l-1}_i$$

従って

$$\begin{aligned} \partial C / \partial W_{j,i}^l &= (\partial C / \partial z_j^l) (\partial z_j^l / \partial W_{j,i}^l) \\ &= \delta_j^l X^{l-1}_i \end{aligned}$$

特に  $i=0$  (バイアス用のユニット番号) の場合、 $X^{l-1}_0=1$  より、

$$\partial C / \partial W_{j,0}^l = \delta_j^l$$



### (8.4.3) 出力層のユニットの誤差

- ・ 出力層(L層) j番目のユニットの誤差  $\delta_j^L$  は、以下のように具体的に計算できます：

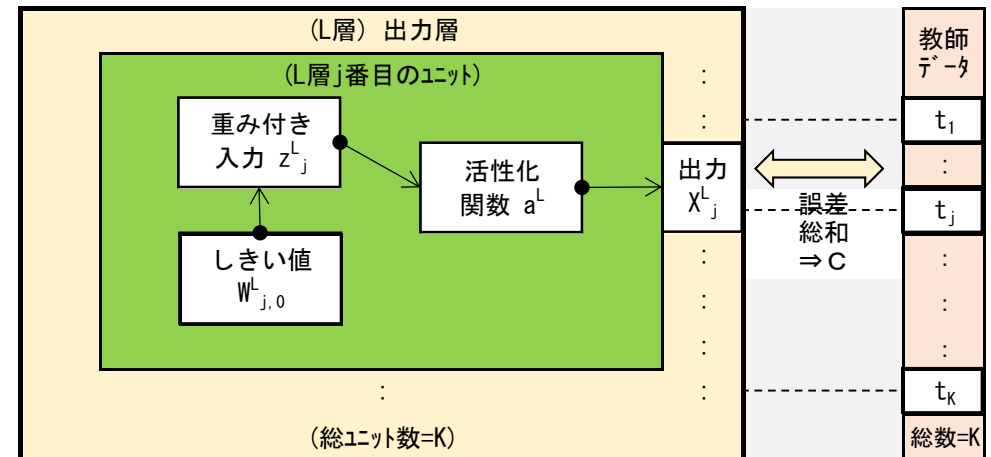
出力層(L層) j番目のユニットの誤差  $\delta_j^L$  の計算

$$\begin{aligned}\delta_j^L &=_{(*)0} \partial C / \partial z_j^L \\ &=_{(*)1} (\partial C / \partial X_j^L) (\partial X_j^L / \partial z_j^L) \quad \dots (\text{式8.4-4})\end{aligned}$$

(\*0) : ユニットの誤差の定義

(\*1) :  $X_j^L = a^L(z_j^L)$  の関係式から偏微分の連鎖律を適用

$\delta_j^L$  : 出力層j番目のユニットの誤差  
L : 出力層の層番号  
j : 出力層でのユニット番号 (j=1~K)  
C : コスト関数  
 $z_j^L$  : 出力層j番目のユニットの重み付き入力値  
 $X_j^L$  : 出力層j番目のユニットの出力値  $X_j^L = a^L(z_j^L)$   
 $a^L$  : 出力層の活性化関数



( 出力層のユニットの誤差の計算例 )

・ 出力層(L層) に関して、コスト関数Cとして平均二乗誤差(MSE)を適用する場合、

コスト関数Cは、「学習データ」の入力値から算出した予測値  $X_j^L$  と正解値  $t_j$  の平均二乗誤差の総和に相当します。

$$C = 1/(N*K) \sum_{n=1}^N \sum_{j=1}^K (X_j^L[n] - t_j[n])^2$$

$$\therefore \partial C / \partial X_j^L = 2/(N*K) \times (X_j^L - t_j)$$

活性化関数  $a$  として、

シグモイド関数  $\sigma(z) = 1/(1+\exp(-z))$  を持ちいる場合、

$$X_j^L = \sigma(z_j^L)$$

$$\therefore \partial X_j^L / \partial z_j^L = \partial \sigma(z_j^L) / \partial z_j^L$$

$$= \sigma(z_j^L) \{1 - \sigma(z_j^L)\}$$

従って、この時の「出力層( $l=L$ )  $j$  番目のユニットの誤差」  $\delta_j^L$  は、

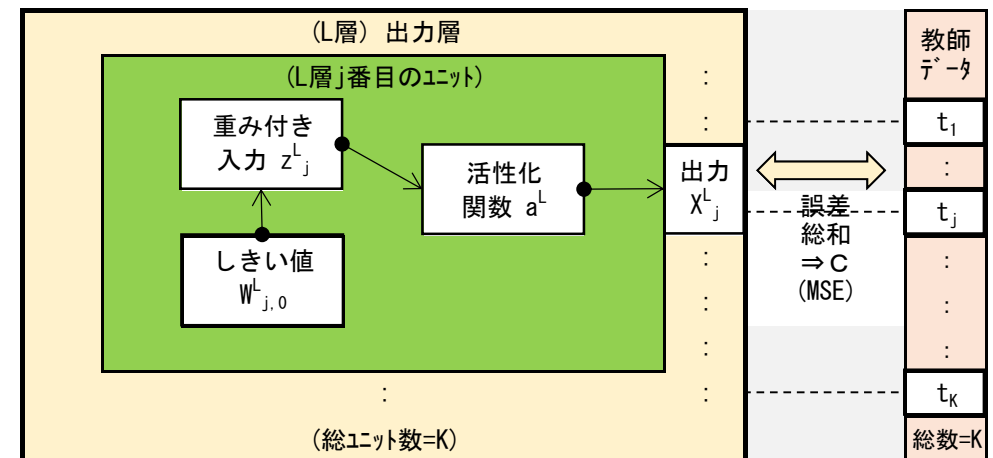
$$\delta_j^L =_{(*1)} (\partial C / \partial X_j^L) (\partial X_j^L / \partial z_j^L)$$

$$=_{(*2)} 2/(N*K) \times (X_j^L - t_j) \times \sigma(z_j^L) \{1 - \sigma(z_j^L)\}$$

(\*1) : (式8.4-4)

(\*2) : 上の関係式

$N, n$  : 学習データ数とその参照添え字  
 $K, j$  : 出力層のユニット数とその参照添え字  
 $X_j^L[n]$  :  $n$  番目のデータについての出力層  $j$  番目ユニットの予測値  
 $t_j[n]$  :  $n$  番目のデータについての出力層  $j$  番目ユニットの教師データ  
 $X_j^L$  :  $\sum_{n=1}^N X_j^L[n]$   
 $t_j$  :  $\sum_{n=1}^N t_j[n]$   
 $z_j^L$  : 出力層  $j$  番目のユニットの重み付き入力  
 $= \sum_{i=0}^I W_{j,i}^L * X_i^{L-1}$  (但し  $X_0^{L-1}=1$ )



#### (8.4.4) 隣接する層間のエットの誤差についての逆漸化式

- ・ 中間層である l 層 j 番目のエットに関して、l-1 層 i 番目のエットとの間に、次の逆漸化式が成り立ちます：

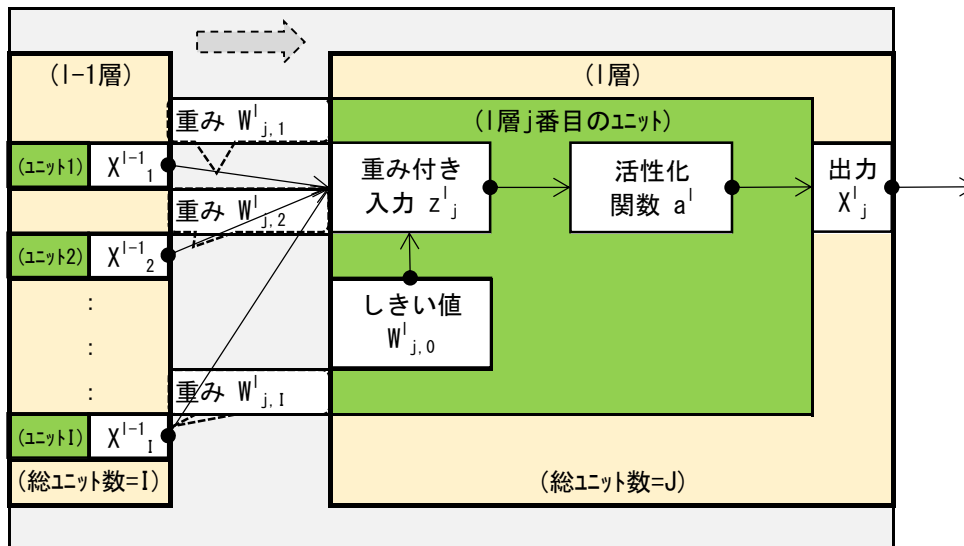
l 層 j 番目のエットと、l-1 層 i 番目のエットとの間に成立する逆漸化式

$$\delta^{l-1}_i = \left\{ \sum_{j=1}^J \delta^l_j * W^l_{j,i} \right\} * \partial a^{l-1} / \partial z^{l-1}_i$$

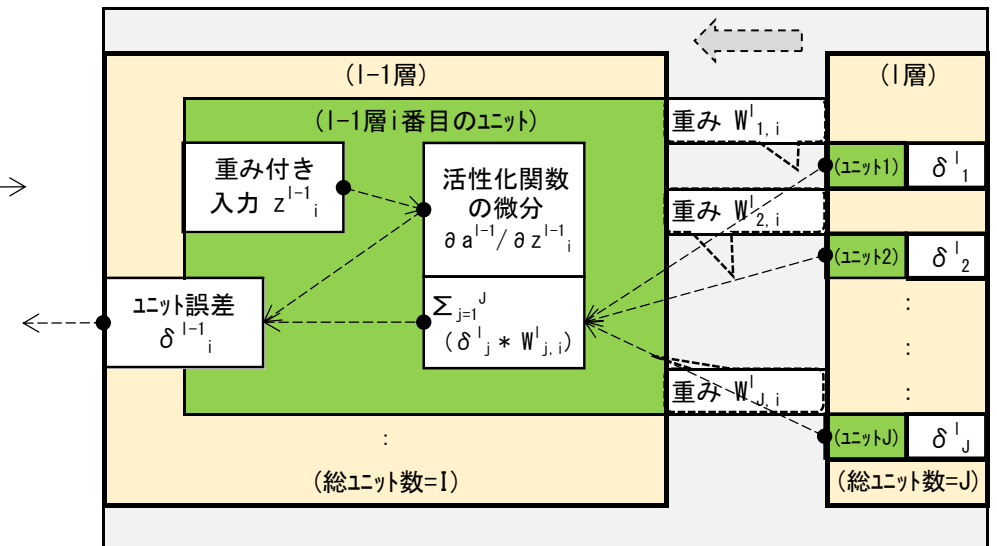
…(式8.4-5)

- l : 層番号 (l=2~L、説明上入力層を第1層とします)
- $\delta^{l-1}_i$  : l-1 層 i 番目のエットの誤差
- $\delta^l_j$  : l 層 j 番目のエットの誤差
- I, i : l-1 層でのエット数, エット番号 (i=1~I)
- J, j : l 層でのエット数, エット番号 (j=1~J)
- $W^l_{j,0}$  : l 層 j 番目のエットについてのバイアス
- $W^l_{j,i}$  : l 層 j 番目のエットへの l-1 層 i 番目のエットからの入力信号の重み
- $a^{l-1}$  : l-1 層の活性化関数
- $z^{l-1}_i$  : l-1 層 i 番目のエットの重み付き入力値

(出力計算時の流れ)



(誤差逆伝搬時の流れ)



- ・ (式8. 4-5) の導出について、補足説明を加えておきます。

$$\begin{aligned}
 \delta^{l-1}_i &=_{(\ast 0)} \partial C / \partial z^{l-1}_i \\
 &=_{(\ast 1)} \sum_{j=1}^J (\partial C / \partial z^l_j) (\partial z^l_j / \partial X^{l-1}_i) (\partial X^{l-1}_i / \partial z^{l-1}_i) \\
 &=_{(\ast 0, 2, 3)} \{ \sum_{j=1}^J \delta^l_j \ast W^l_{j,i} \} \ast \partial a^{l-1} / \partial z^{l-1}_i
 \end{aligned}$$

(※0) : ニットの誤差の定義

$$\begin{aligned}
 \delta^{l-1}_i &= \partial C / \partial z^{l-1}_i \\
 \delta^l_j &= \partial C / \partial z^l_j
 \end{aligned}$$

(※1) : 以下の関数関係に基づいて偏微分の連鎖律を適用

$$\begin{aligned}
 X^{l-1}_i &= a^{l-1}(z^{l-1}_i) \\
 z^l_j &= \sum_{i=0}^I W^l_{j,i} \ast X^{l-1}_i \\
 &= \sum_{i=0}^I W^l_{j,i} \ast a^{l-1}(z^{l-1}_i)
 \end{aligned}$$

(※2) : l層の第jニットの重み付き入力

$$\begin{aligned}
 z^l_j &= \sum_{i=0}^I W^l_{j,i} \ast X^{l-1}_i \\
 \text{より、} X^{l-1}_i \text{ についての偏微分をとると、} \\
 \partial z^l_j / \partial X^{l-1}_i &= W^l_{j,i}
 \end{aligned}$$

(※3) : l-1層の第iニットの出力は、活性化関数  $a^{l-1}$  を用いて

$$\begin{aligned}
 X^{l-1}_i &= a^{l-1}(z^{l-1}_i) \\
 \text{より、その偏微分をとると、} \\
 \partial X^{l-1}_i / \partial z^{l-1}_i &= \partial a^{l-1} / \partial z^{l-1}_i
 \end{aligned}$$

C : コスト関数

l : 層番号 (l=2~L、説明上入力層を第1層とします)

I, i : l-1層でのユニット数, ユニット番号 (i=1~I)

$\delta^{l-1}_i$  : l-1層i番目のユニットの誤差

$z^{l-1}_i$  : l-1層i番目のユニットの重み付き入力値

$X^{l-1}_i$  : l-1層i番目のユニットの出力値

$a^{l-1}$  : l-1層の活性化関数

J, j : l層でのユニット数, ユニット番号 (j=1~J)

$\delta^l_j$  : l層j番目のユニットの誤差

$z^l_j$  : l層j番目のユニットの重み付き入力値

$X^l_j$  : l層j番目のユニットの出力値

$a^l$  : l層の活性化関数

$W^l_{j,0}$  : l層j番目のユニットについてのバイアス

$W^l_{j,i}$  : l層j番目のユニットへのl-1層i番目のユニットからの入力信号の重み

#### (8.4.5) 全ての中間層のユニットの誤差 $\delta$ を計算する

- ・「漸化式 (ゼンガシ、recurrence relation)」の場合は、初項と、第n項で第n+1項を表現する関係式を与えると、全ての項が計算できます。  
これに対し「逆漸化式」の場合は、末項と、第n項で第n-1項を表現する関係式を与えると、全ての項が計算できます。
- ・「(8.4.3) 出力層のユニットの誤差」で末項「 $\delta^L_j$ 」を与え、  
「(8.4.4) 隣接する層間のユニットの誤差についての逆漸化式」で、  
第l-1層i番目のユニットの誤差「 $\delta^{l-1}_i$ 」を、第l層j番目のユニットの誤差「 $\delta^l_j$ 」で表現する関係式を与えることにより、  
ユニットの誤差  $\delta$  は、全ての中間層のユニットについて計算することができます。

出力層j番目のユニットの誤差  $\delta^L_j$

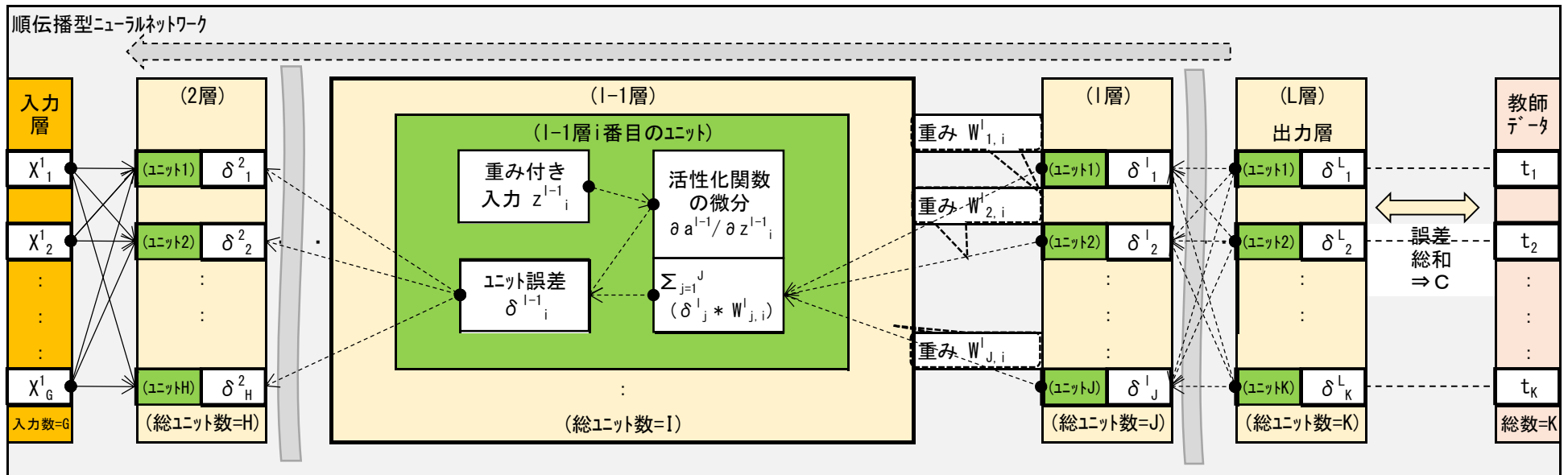
$$\delta^L_j = (\partial C / \partial X^L_j) (\partial X^L_j / \partial z^L_j)$$

…(式8.4-4) 再掲

l層j番目のユニットと、l-1層i番目のユニットとの間に成立する逆漸化式

$$\delta^{l-1}_i = \left\{ \sum_{j=1}^J \delta^l_j * W^l_{j,i} \right\} * \partial a^{l-1} / \partial z^{l-1}_i$$

…(式8.4-5) 再掲



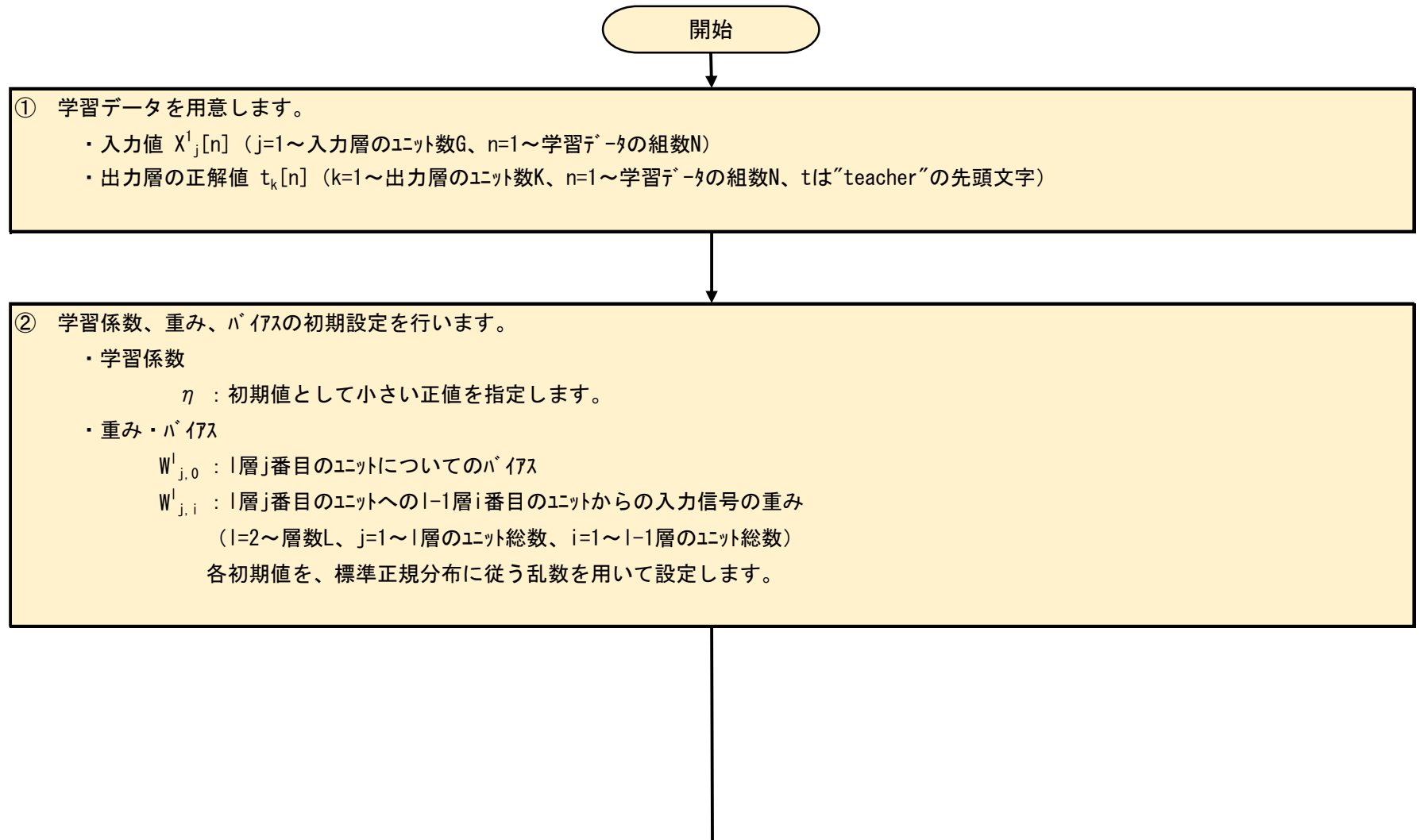


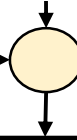
#### (8.4.6) ニットの誤差 $\delta$ を用いて勾配降下法で学習する流れ (例)

・ ニットの誤差  $\delta$  を導入した「誤差逆伝播法」による学習(モデルのパラメータ最適化)のアルゴリズム例は、以下のようになります。

以下の例では、コスト関数  $C$ 、第  $l$  層の活性化関数  $a^l$  として、次のものを使用します。

- (1) コスト関数  $C$  として二乗誤差(SE)
- (2) 第  $l$  層の活性化関数  $a^l$  としてシグモイド関数





③ 現在の重みとバイアスを用いて、各層のユニットの出力値を第2層（中間層の第1層）から出力層に向かって順番に計算し（フォワードプロパゲーション）、最後に出力層の出力値と学習データとの二乗誤差(SE)を算出し、これをコスト関数Cとします。  
 (l=2~層数L分、l=1:入力層、l=L:出力層、j=1~各層l毎のユニット数、a<sup>l</sup>は第l層の活性化関数)

- ・ 入力層の出力値  
 $X_j^1$  : 入力層j番目ユニットの出力値
- ・ 中間層・出力層の出力値  
 $X_j^l$  : 第l層j番目ユニットの出力値  

$$= a^l \left( \sum_{i=0}^{\text{l-1層のユニット数}} (W_{j,i}^l * X_i^{l-1}[n]) \right)$$
- ・ コスト関数(二乗誤差(SE))  

$$C = \sum_{j=1}^{\text{出力層のユニット数}} 1/2 (X_j^L[n] - t_j[n])^2$$

n : 学習データの組番号。n=1~学習データの組数N

$W_{j,0}^l$  : l層j番目のユニットについてのバイアス

$W_{j,i}^l$  : l層j番目のユニットへのl-1層i番目のユニットからの入力信号の重み

$X_0^{l-1}$  : l-1層からのユニットの出力のバイアス用の出力値 (1固定)

$X_i^{l-1}$  : l-1層i番目のユニットの出力値

④ 誤差逆伝播法から、出力層から中間層に向かって逆漸化式を用いて、各層のユニットの誤差δを計算します(バックプロパゲーション)。  
 (例では、第l層の活性化関数 a<sup>l</sup> としてシグモイド関数「 $\sigma(z) = 1/(1+\exp(-z))$ 」を使用)

- ・ 出力層のユニットの誤差  
 $\delta_j^L$  : 出力層j番目のユニットの誤差 (j=1~出力層のユニット総数)  

$$= (\partial C / \partial X_j^L) (\partial X_j^L / \partial z_j^L)$$

$$= (X_j^L - t_j) \exp(-z_j^L) / (1 + \exp(-z_j^L))$$
- ・ 中間層のユニットの誤差  
 $\delta_i^{l-1}$  : 第l-1層i番目のユニットの誤差 (l=3~L-1、i=1~l層のユニット総数)  

$$= \left\{ \sum_{j=1}^{\text{l層のユニット数}} (\delta_j^l \times W_{j,i}^l) \right\} * \partial a^{l-1} / \partial z_i^{l-1}$$

$$= \left\{ \sum_{j=1}^{\text{l層のユニット数}} (\delta_j^l \times W_{j,i}^l) \right\} * \exp(-z_i^{l-1}) / (1 + \exp(-z_i^{l-1}))$$

$z_j^l$  : l層j番目のユニットの重み付き入力

$$= \sum_{i=0}^{\text{l-1層のユニット総数}} W_{j,i}^l * X_i^{l-1} \quad (\text{但し } X_0^{l-1}=1)$$

⑤ 各層のユニットの誤差  $\delta$  からコスト関数  $C$  の、重み・バイアスに関する偏微分を計算します (バックプロパゲーションの続き)。  
( $l=2 \sim$  層数  $L$  分、 $j=1 \sim l$  層の各層毎のユニット総数、 $i=1 \sim l-1$  層の各層毎のユニット総数)

・ 重みに関するコスト関数の偏微分

$l$  層  $j$  番目のユニットについて、 $l-1$  層  $i$  番目のユニットからの入力信号の重み  $W_{j,i}^l$  に関するコスト関数  $C$  の偏微分

$$\partial C / \partial W_{j,i}^l = \delta_j^l X_{i,i}^{l-1}$$

・ バイアスに関するコスト関数の偏微分

$l$  層  $j$  番目のユニットについてのバイアスに関するコスト関数の偏微分

$$\partial C / \partial W_{j,0}^l = \delta_j^l$$

$W_{j,i}^l$  :  $l$  層  $j$  番目のユニットの、 $l-1$  層  $i$  番目のユニットからの入力信号の重み  
 $W_{j,0}^l$  :  $l$  層  $j$  番目のユニットの、 $l-1$  層からの入力信号に付加するバイアス  
 $X_{i,i}^{l-1}$  :  $l-1$  層  $i$  番目のユニットの出力値 ( $l$  層への入力信号)

⑥ 上記③～⑤の結果を全入力データについて加え合わせ、コスト関数  $C_T$  とその勾配  $\nabla C_T$  を算出します。

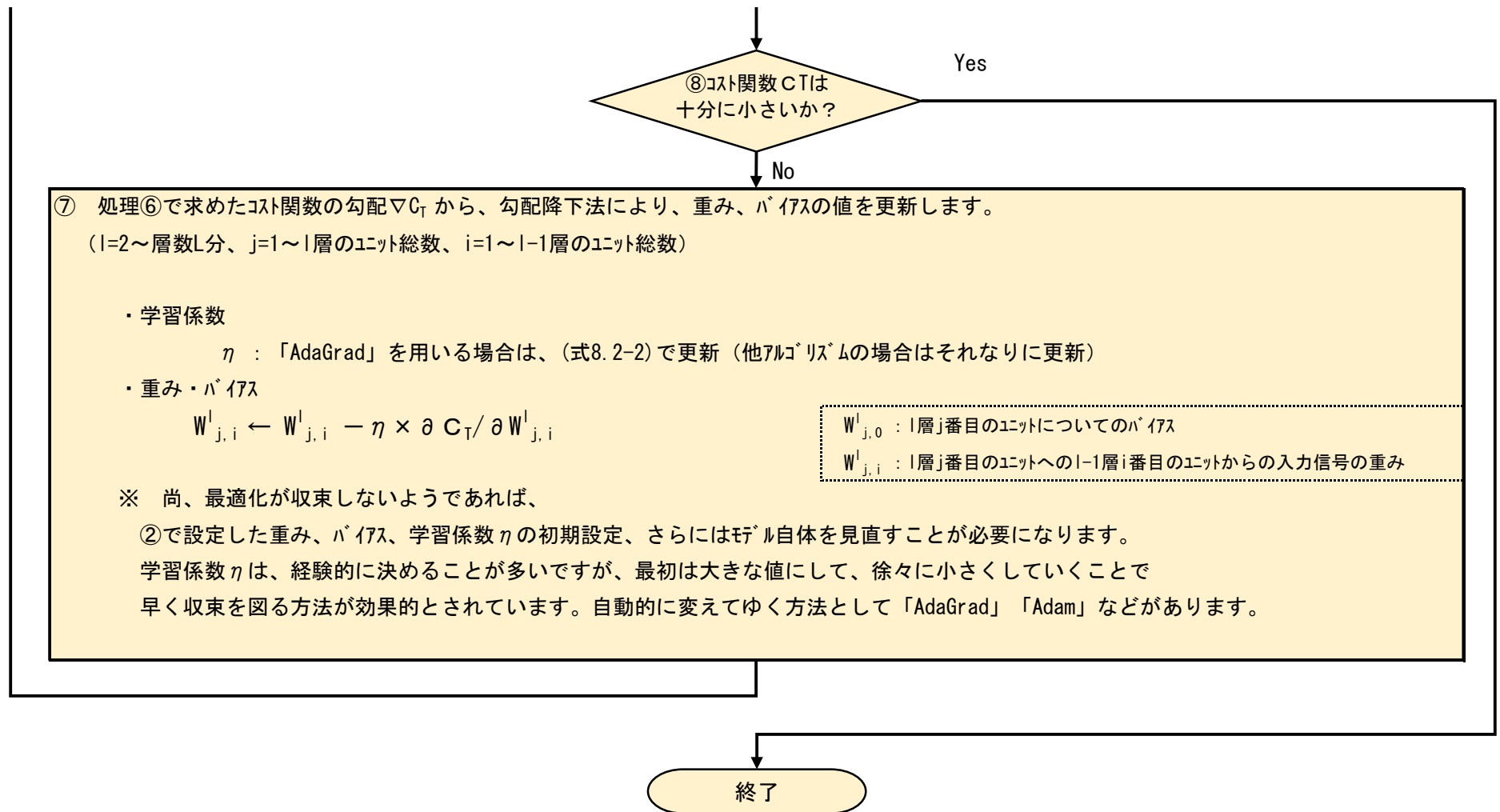
・ 全入力データについて加え合わせたコスト関数

$$C_T = \sum_{\text{学習データ数}} \text{処理③で算出された } C$$

・ 全入力データについて加え合わせたコスト関数の勾配

$$\begin{aligned} \nabla C_T &= \sum_{\text{学習データ数}} \text{処理⑤で算出された } C \text{ の偏微分を成分とするベクトル} \\ &= (\partial C_T / \partial W_{j,0}^l, \partial C_T / \partial W_{j,1}^l, \dots, \partial C_T / \partial W_{j,l-1}^l \text{ 層のユニット総数}) \end{aligned}$$

$l$  : 層番号 ( $l=2 \sim$  層数  $L$  分)  
 $j$  :  $l$  層でのユニット番号 ( $j=1 \sim l$  層のユニット総数)  
 $i$  :  $l-1$  層でのユニット番号 ( $i=1 \sim l-1$  層のユニット総数)  
 $n$  : 学習データの組番号。  $n=1 \sim$  学習データの組数  $N$



## (8.5) モデルの学習の検証

- ・モデルを作成し、学習を通してモデルのパラメータの最適化を行います。

学習時のデータではコスト関数  $C$  が最小になるものの、学習時以外のデータでは誤差が大きくなってしまいます。

これは、学習時のデータに特化した学習をしてしまって、汎用になっていないということであり、

これを「過学習（かがくじゅ、over fitting）」と言います。

### (8.5.1) モデルの学習の検証方法

- ・使用しているモデルと、そのパラメータの学習状況を把握する為の検証方法を幾つか紹介します。

以下は、データを、モデルの学習用と、評価用に分けて評価する方法です。

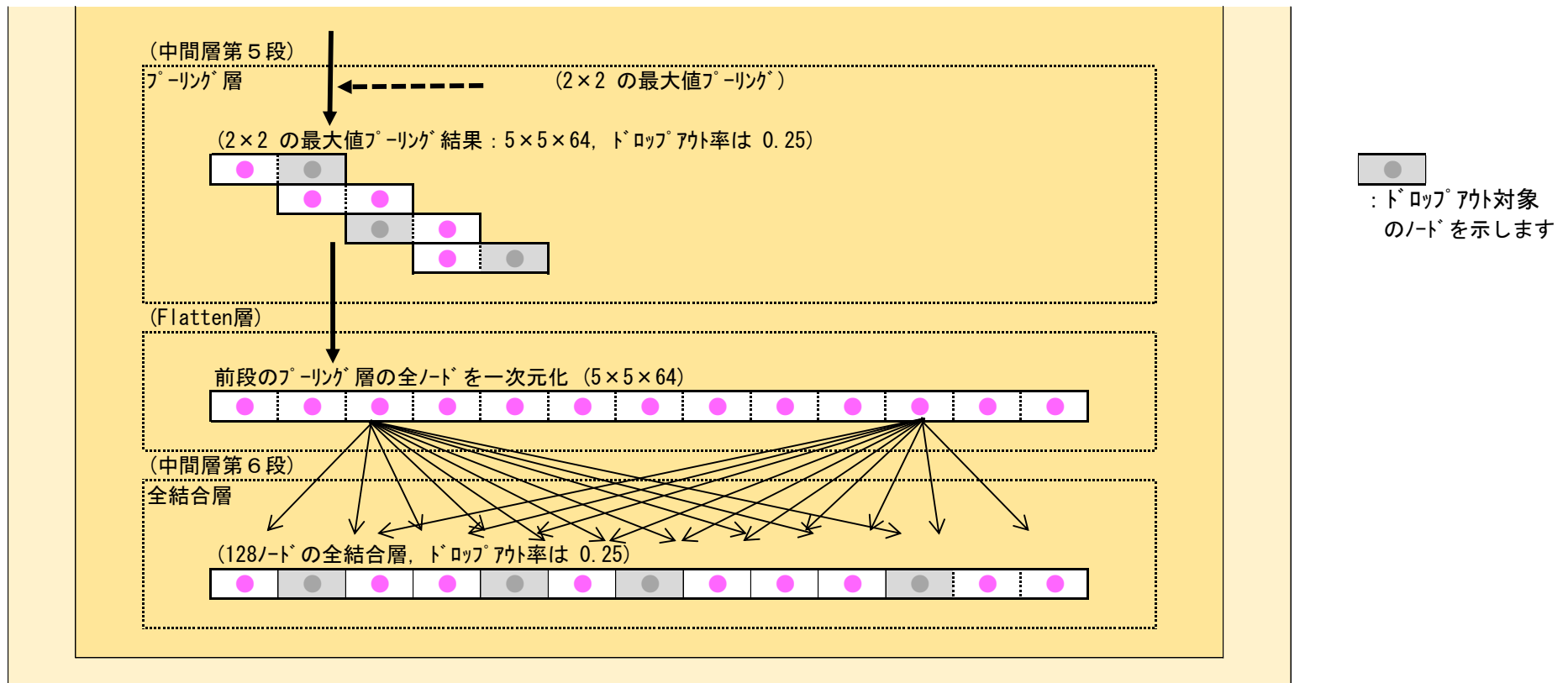
学習用データは、「訓練データ（クレンデータ、training data）」「訓練事例集合（クレンジレイシュゴウ、training set）」とも言います。

評価用データは、「テストデータ（test data）」「テスト事例集合（テストレイシュゴウ、testing set）」とも言います。

検証方法	説明	長所・短所
ホールドアウト検証 (hold out validation)	<ul style="list-style-type: none"><li>・集められたデータのうち、学習専用と評価専用の2つのグループに分けて、学習データでモデルを作成し、評価データで作成されたモデルの評価を行う。</li><li>これにより作成したモデルがどの位未知データを判別できるかが評価でき、過学習を防ぐことにつながる。</li><li>・尚、評価データは学習データに含めてはいけない。</li></ul>	<ul style="list-style-type: none"><li>・集めたデータを学習専用と評価専用に分けるため、学習に使えるデータの量が減るという欠点がある。</li><li>・データを学習用と評価用に分ける際、分け方に偶然偏りができた場合にはうまく評価ができないという欠点がある。</li><li>・データ数が少ない時は、上記の欠点が顕著に表れる。</li><li>・データ量が多くなると「K-分割交差検証」の結果と殆ど変わらなくなる。</li></ul>
K-分割交差検証 (K-fold cross-validation)	<ul style="list-style-type: none"><li>・集められたデータ（説明変数と目的変数）をK個に分割する。そのうちの1個を評価データとし、残りのK-1個を学習データとして、モデルのパラメータを算出して、評価データでの平均二乗誤差を計算する。</li><li>これをK個分繰り返して、平均二乗誤差の平均をこのモデルの評価値とする。</li></ul>	<ul style="list-style-type: none"><li>・データ量が多くなると計算に時間がかかるという欠点がある。</li></ul>
leave-one-out 交差検証 (一個抜き交差検証) (leave-one-out cross-validation (LOOCV))	<ul style="list-style-type: none"><li>・K-分割交差検証で、分割数Kをデータ数Nとした場合の検証である。この時、評価データ数は1、学習データ数はN-1になる。データ数が特に少ない時にこの手法が使われる。</li></ul>	<ul style="list-style-type: none"><li>・データ量が多くなると計算に時間がかかるという欠点がある。</li></ul>

## (8.5.2) 過学習の対策とドロップアウト

- ・「過学習」の対策として、入力にノイズを与えたり、学習時に正則化を適用したり、様々な方法が提案されてきましたが、「ドロップアウト (Dropout)」という方法が有効であることが知られています。
- ・「ドロップアウト」は、学習時に入力層や中間層のユニットを、各層毎に一定の割合  $\alpha$  ( $0 \leq \alpha < 1$ ) でユニットを選出し、それ以外は無効な状態で学習を行います（選出されなかったユニットは存在しないものとして学習を行います）。これをミニバッチ毎にユニットを選出しながら、学習を繰り返します。
- ・学習後に予測をする場合は、全てのユニットを使用しますが、学習時は  $\alpha$  の割合のユニットしか存在しなかったため、予測時にはドロップアウトを施した層の出力の重みを  $\alpha$  倍にして出力が大きくなることを防ぎます。
- ・既述のモデル「手書き文字認識モデル3」では、プーリング層と全結合層について、ドロップアウト率  $\alpha = 0.25$  で作成しています。



## 【出典・参考】

勾配降下法⇒ <https://qiita.com/YudaiSadakuni/items/ece07b04c685a64eaa01>

勾配降下法⇒ <https://postd.cc/optimizing-gradient-descent/#batchgradientdescent>

AdaGrad、Adam⇒ <https://qiita.com/tokkuman/items/1944c00415d129ca0ee9>

AdaGrad、Adam⇒ 「深層学習 Deep Learning」人工知能学会監修（近代科学社 2015年11月）

モデルの学習⇒[https://ml4a.github.io/ml4a/jp/how\\_neural\\_networks\\_are\\_trained/](https://ml4a.github.io/ml4a/jp/how_neural_networks_are_trained/)

本章全般⇒「Pythonで動かして学ぶ！ あたらしい機械学習の教科書」伊藤真著（翔泳社 2018年1月）

確率的勾配降下法⇒ <http://yaju3d.hatenablog.jp/entry/2017/08/27/233459>

誤差逆伝搬法⇒ 「ディープラーニングがわかる数学入門」涌井良幸・貞美著（技術評論社 2017.04）

トロッポアウト⇒ 「深層学習 Deep Learning」人工知能学会監修（近代科学社 2015年11月）

## (9) まとめ

- ・ 以上、「教師あり学習（分類）」を中心に機械学習の基本的な技術的要素を解説してきましたが、如何でしたでしょうか？ 今回の解説で、「セミナー7回目」の「モデルとその学習方法」が手に取るように読めるようになったのでは？と期待しています。
- ・ 誤差逆伝搬法については、「順伝播型ニューラルネットワーク」だけ取り上げましたが、「畳み込みニューラルネットワーク」について詳しく知りたい方は、以下の書籍などに詳しく書かれているので、参照してみてください。

「Pythonで動かして学ぶ！ あたらしい機械学習の教科書」伊藤真著（翔泳社 2018年1月）

「ディープラーニングがわかる数学入門」涌井良幸・貞美著（技術評論社 2017.04）

- ・ 実装レベルではすでに提供されているAPIを使いこなす・・・といったことが中心となります。  
APIを使いこなすことができれば、仕事をこなすことは十分できるような時代になってはいますが、今回のセミナーで取り上げたような基本的な背景を押さえておくことは、無駄にはならないと思います。
- ・ 今回取り上げた内容は機械学習の全体のほんの一部にすぎませんので、これを踏み台にして、自分も含めて、どんどん取り組んで行っていただきたいと思います。

以上です。



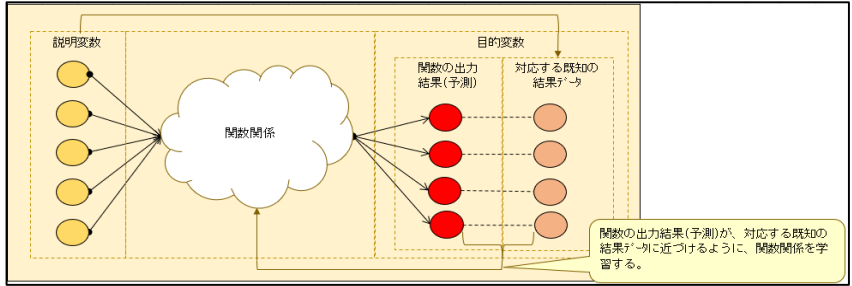
(10) 確認問題

(1) 学習の種類

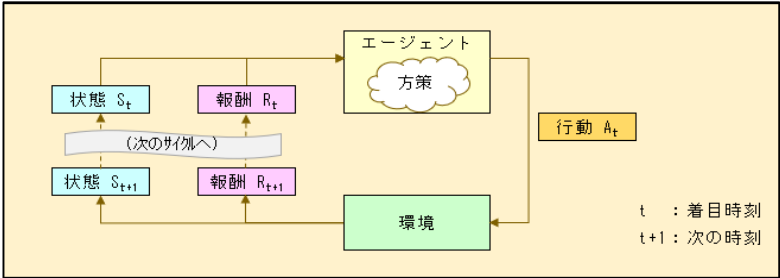
以下の空欄に最もあてはまる語句を選択肢から選び、その記号を回答欄に記入してください。

- ・ (1.1) \_\_\_\_\_ は 異なるデータ間の関係性を、一方を入力（説明変数）、他方を出力（目的変数）とする関数関係として記述し、関数の出力を既知の結果データ（目的変数の既知データ）に近づけるように、関数関係を学習する方法です。この場合、関数の出力結果（予測）に対して、既知の結果データが教師データの役割をします（図10.1.1 参照）。
- ・ (1.2) \_\_\_\_\_ は、入力変数（説明変数）に対する、出力変数（目的変数）が与えられておらず、入力データ自体の背後にある性質や特徴や構造などを分析・発見することと言えます。例として、「クラスタリング（Clustering）」によるデータ分類や、「次元削減（ジゲンサクゲン、Dimensionality Reduction）」による特徴量抽出などが挙げられます。
- ・ (1.3) \_\_\_\_\_ は、「認知」した状況の下で最適な「行動」を「判断」すること、すなわち「制御」を学習する方法です（図10.1.2 参照）。強化学習では、「エージェント（agent、制御器）」が「環境（environment）」との相互作用を通じて情報収集しながら、環境を制御する方法を学習します。「環境」は、状態とその変化を規定する制御対象の系（システム）のことです。「エージェント」は、環境の置かれた状態にもとづいて行動を選択しながら環境を制御する主体で、制御プログラムにあたります。

(図10.1.1)



(図10.1.2)



(選択肢)

- (a) 「強化学習（キョウカガクシユ、Reinforcement learning）」
- (b) 「教師あり学習（キョウシアリガクシユ、Supervised learning）」
- (c) 「教師なし学習（キョウナシガクシユ、Unsupervised Learning）」
- (d) 「転移学習（テンイガクシユ、Transfer Learning）」

(回答)

(1.1)	
(1.2)	
(1.3)	

(2) 入力準備

以下の空欄に最もあてはまる語句を選択肢から選び、その記号を回答欄に記入してください。

- ・入力データ（説明変数）を統計処理するに先立って、  
入力データの中から重複や誤記、表記の揺れなどを探し出し、削除や修正などを行い、入力データの品質を高める必要があります。  
この作業を(2.1)\_\_\_\_\_または「データクレンジング（data cleansing）」と言います。
- ・入力データ（説明変数）が数種類からなる場合、各入力モードに必ずしも全ての種類のデータが揃っている訳ではありません。  
欠けているデータのことを(2.2)\_\_\_\_\_と言います。  
欠損値がある場合、統計処理に支障があるため、当該項目を削除したり、何らかの値で補完したり、といった対応をします。
- ・平均と分散についての公式を用いて、平均 $\mu$ 、標準偏差 $\sigma$ の確率変数  $X$ を、平均0、標準偏差1の確率変数  $Z$ へ変換することが出来ます。  
この操作を(2.3)\_\_\_\_\_と言います。これには、以下の様な効果があります。
  - (1) 特定の特徴に偏った結果が出ることを防ぐことが出来ます。
  - (2) 異なる項目のデータであってもその大小を比較できるようになります。
- ・入力データを符号化してベクトルに取り込んだ方が良い場合があります。  
出力が  $K$ 種類への分類ベクトルの場合、目的変数を $K$ 次元のベクトルとします。  
ある入力データが第 $k$ 番目のクラスへ分類されることが分かっている場合、目的変数ベクトルの  $k$ 番目の要素だけ 1でそれ以外を 0とする符号化を  
(2.4)\_\_\_\_\_（または「one-hot encoding」）と言います。

(選択肢)

- (a) 「1-of-K 符号化法（ワンオブケーフオフカーク、One-of-K encoding）」
- (b) 「データクリーニング（data cleaning）」
- (c) 「ハフマンブロック符号化法（Huffman block coding method）」
- (d) 「欠損値（ケツソンチ、missing value）」
- (e) 「標準化（ヒョウジュンカ、standardized）」

(回答)

(2.1)	
(2.2)	
(2.3)	
(2.4)	

### (3) 活性化関数

以下の空欄に最もあてはまる語句を選択肢から選び、その記号を回答欄に記入してください。

- あるユニットに着目した時、そのユニットへの入力信号の重み付き線形和にバイアスを加えたもの(下図で重み付き入力  $z$ ) を、どのように活性化して次のユニットに渡すかを、(3.1) \_\_\_\_\_ で指定します。

「活性化関数」は「出力関数 (シュツヨウカンズ、Output function)」とも言います。

- 活性化関数の一つ(3.2) \_\_\_\_\_ は、 $k$ 個の入力値  $X_i$  ( $i=1\sim k$ ) の順序関係を保ちながら、 $X_i$  を確率としての値  $Y_i$  ( $i=1\sim k$ 、各値は0~1の範囲で総和が1になる) に変換する関数で、 $k$ 個のクラスへ分類する際の所属の確率を示すのによく用いられます。

$k$  個のクラスへ分類する時のソフトマックス関数

$$Y_i = \exp(X_i) / \sum_{j=1}^k \exp(X_j) \quad (i=1\sim k)$$

- 活性化関数の一つ(3.3) \_\_\_\_\_ は「ソフトマックス関数」で分類するクラス数を2個に限定したものに相当します。

2個のクラスへ分類する時のソフトマックス関数

$$\begin{aligned} Y_1 &= \exp(X_1) / \{\exp(X_1) + \exp(X_2)\} = 1 / \{1 + \exp(X_2 - X_1)\} = 1 / \{1 + \exp(-Z)\} & (Z = X_1 - X_2) & \text{= シグモイド関数} \\ Y_2 &= \exp(X_2) / \{\exp(X_1) + \exp(X_2)\} = 1 - Y_1 \end{aligned}$$

- 活性化関数の一つ(3.4) \_\_\_\_\_ は計算が容易で(正の領域では微分が1)、  
「sigmoid関数」で問題となる学習時の「誤差逆伝播法」での(3.5) \_\_\_\_\_ もなく、  
一般的に使用されています。

(選択肢)

- (a) 「ReLU (正規化線形関数)」
- (b) 「シグモイド関数 (sigmoid function)」
- (c) 「ソフトマックス関数 (softmax function)」
- (d) 「過学習 (ガクシュウ、over fitting)」
- (e) 「活性化関数 (カッセイカンズ、Activation function)」
- (f) 「勾配消失問題 (コウバイショウツツモンダイ、vanishing gradient problem)」

(回答)

(3.1)	
(3.2)	
(3.3)	
(3.4)	
(3.5)	

#### (4) 損失関数・誤差関数・コスト関数・目的関数

以下の空欄に最もあてはまる語句を選択肢から選び、その記号を回答欄に記入してください。

- ・説明変数と目的変数の関係を良く説明できるモデルが「良いモデル」であり、「モデルの出力結果(予測)」と「既知の結果データ」との差(「誤差 (ゴサ、error)」) が最小になるように、モデルの「パラメータ」を最適化するのが「学習」です。
- ・「学習」を通して最小化したいのはこの誤差であり、誤差を表す関数のことを、「誤差関数 (ゴサかんすう、error function)」、又は(4.1) \_\_\_\_\_、又は「コスト関数 (コストかんすう、cost function)」と言います。これらは「最小化する」という目的があるので、(4.2) \_\_\_\_\_とも言います。
- ・統計学において、未知または既知のある確率分布の下で、確率変数の結果が出現する場合に、確率変数の出現結果からみて、元々の確率分布を推測する時の「尤もらしさ (モットモラシサ)」を表す数値を、(4.3) \_\_\_\_\_または「尤度関数 (ユウドかんすう、likelihood function)」と言います。
- ・この「尤度」を最大化するように元々の確率分布を推測することを(4.4) \_\_\_\_\_と言います。「離散型確率分布」の場合は「確率質量関数」を、「連続型確率分布」の場合は「確率密度関数」を推定することになります。「尤度」を最大化することは、モデルの学習時の評価指数である「損失関数」の値を最小化することになります。
- ・自然対数関数  $\log(x)$  が  $x$  について単調増加であり、関数値の大小関係が  $x$  の大小関係を保つという性質から、「尤度」を最大にする  $x$  の値は、この対数をとった(4.5) \_\_\_\_\_の値も最大にします。

(選択肢)

- (a) 「損失関数 (ソンシツかんすう、loss function)」
- (b) 「目的関数 (モクテキかんすう、objective function)」
- (c) 「尤度 (ユウド、likelihood)」
- (d) 「最尤推定 (サイユウスタイイ、maximum likelihood estimation、MLE)」
- (f) 「対数尤度 (タイスウユウド、Log likelihood)」

(回答)

(4.1)	
(4.2)	
(4.3)	
(4.4)	
(4.5)	

(5) 順伝播型ニューラルネットワーク (FFNN)、畳み込みニューラルネットワーク (CNN)

以下の空欄に最もあてはまる語句を選択肢から選び、その記号を回答欄に記入してください。

・ (5. 1) \_\_\_\_\_ は、以下の特徴を持つニューラルネットワークです。

(5. 2) \_\_\_\_\_ は、

FFNN の中間層（隠れ層）を深くしたものです。

(1) 隣接する層間のすべてのユニットが完全に接続されている（完全結合）

(2) 入力層から出力層に向けて一方通行で流れる

(3) 入力と出力の間に、1 つ以上の隠れ層がある

・ (5. 3) \_\_\_\_\_ は、階層型ニューラルネットワークに、畳み込み構造を

入れ込んだモデルで、人の視覚をモデルにした構成になっており、画像認識分野を中心に幅広く利用されているディープラーニング（深層学習）手法の一つです。以下のような特徴があります。

(5. 4) \_\_\_\_\_ は、

CNN の中間層（隠れ層）を深くしたものです。

(1) 複雑なパターン認識問題にも、簡潔なネットワークで対応できます。

(2) 全体としてユニット数が少なくて済むので、計算が容易です。

・ 上記何れのタイプのネットワークも (5. 5) \_\_\_\_\_ で学習します。

（選択肢）

(a) 「誤差逆伝播法（Backpropagation）」

(b) 「順伝播型ニューラルネットワーク（Feed forward neural network、FFNN）」

(c) 「畳み込みニューラルネットワーク（Convolutional neural network、CNN）」

(d) 「深層順伝播型ニューラルネットワーク（Deep Feed forward neural network、DFF）」

(f) 「深層畳み込みニューラルネットワーク（Deep Convolutional neural network、DCN）」

（回答）

(5. 1)	
(5. 2)	
(5. 3)	
(5. 4)	
(5. 5)	

(6) 畳み込みニューラルネットワーク (CNN) を構成する層

以下の空欄に最もあてはまる語句を選択肢から選び、その記号を回答欄に記入してください。

・「畳み込みニューラルネットワーク (タミコミニューラルネットワーク、Convolutional neural network、CNN)」は、中間層 (隠れ層) が、畳み込み層 (コンボリューション層)、正規化層、プーリング層、全結合層で構成された複数の層からなります。正規化層とプーリング層は必須ではありません。また、畳み込み層でフィルタの「移動量 (ドウリョウ、ストライド、stride)」を調整することにより情報縮約を行って、プーリング層を代行することもあります。

・各層は以下のようなものです。

層名	説明
(6. 1)	入力側からの信号に対してフィルタ処理を行い、「特徴マップ (トクショウマップ、feature map)」を作成する層です。コンボリューション層ともいいます。フィルタは「カーネル (kernel)」「オペレータ」「マスク」とも言います。
(6. 2)	畳み込み層で作成された「特徴マップ」に対し情報縮約操作を行う層。 プーリングには、以下のようなものがあります： ・最大値プーリング (Max Pooling) 着目領域での最大値で縮約する ・平均値プーリング (Average Pooling) 着目領域での平均値で縮約する ・Lpプーリング (Lp Pooling) 着目領域での中央値を強調した値で縮約する
(6. 3)	隣接する層の各ユニット同士が、全ての組合せで結合することを全結合と言います。 通常、出力層と前段の数層を全結合とし、これを全結合層としてモデル化します。
(6. 4)	各サンプルの平均と分散をそれぞれ0と1になるように、入力側からの信号を変換する層。 これで、入力データ間のばらつきを補正します。
(6. 5)	前段の層から出てきた多次元のデータを一次元に変換する層。 次段の層の各ユニットと全結合する準備などの為に入れます。

(選択肢)

- (a) 「Flatten層 (フラットソウ、Flatten Layer)」
- (b) 「プーリング層 (プーリングソウ、Pooling Layer)」
- (c) 「畳み込み層 (タミコミソウ、Convolution Layer)」
- (d) 「正規化層 (セイギカソウ、Normalization Layer)」
- (f) 「全結合層 (ゼンケツゴソウ、Fully Connected Layer)」

(回答)

(6. 1)	
(6. 2)	
(6. 3)	
(6. 4)	
(6. 5)	

## (7) モデルの学習

以下の空欄に最もあてはまる語句を選択肢から選び、その記号を回答欄に記入してください。

- ・教師あり学習では、説明変数と目的変数の関係を説明するモデルを作成し、  
「モデルの出力結果(予測)」と「既知の結果データ」との差が最小化する(コスト関数(損失関数)の値が最小になる)ように、  
モデルの「パラメータ」(各ユニットの重み $W_i$ /しきい値 $W_0$ /フィルタの各値 $W_{ij}^{(f)}$ などの調整可能な値)を最適化します。これが「学習」です。
- ・最適化対象のパラメータを一般化して $\theta$ で表します。パラメータ $\theta$ の関数であるコスト関数 $C(\theta)$ が与えられた時、  
これを最小にする $\theta^*$ を求める最適化を、(7.1) \_\_\_\_\_ で行います。
- ・「勾配降下法」は、現ステップ $t$ のパラメータ $\theta^{(t)}$ における勾配と学習率 $\eta_t$ から、  
次ステップ $t+1$ のパラメータ $\theta^{(t+1)}$ を次式により更新して、コスト関数 $C(\theta)$ が最小になる $\theta^*$ を逐次近似により求める方法です。

勾配降下法によるパラメータの最適化

$$\theta^{(t+1)} = \theta^{(t)} - \eta_t \times \partial C(\theta) / \partial \theta \big|_{\theta = \theta^{(t)}}$$

$C(\theta)$ : コスト関数 (パラメータ $\theta$ の関数)

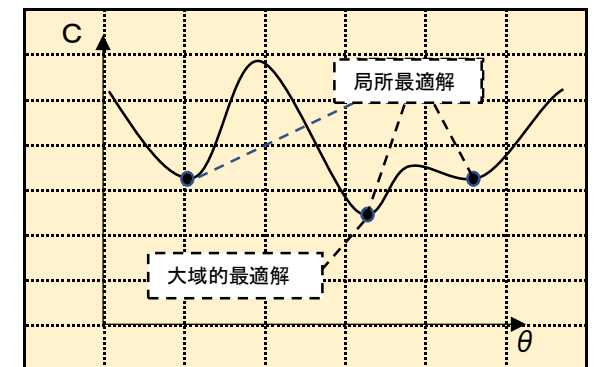
$\theta^{(t)}$ : ステップ $t$ におけるパラメータ $\theta$ の値

$\eta_t$ : ステップ $t$ における学習率 $\eta$ の値 ( $>0$ )

$\partial C(\theta) / \partial \theta \big|_{\theta = \theta^{(t)}}$ :  $\theta^{(t)}$ におけるコスト関数の勾配

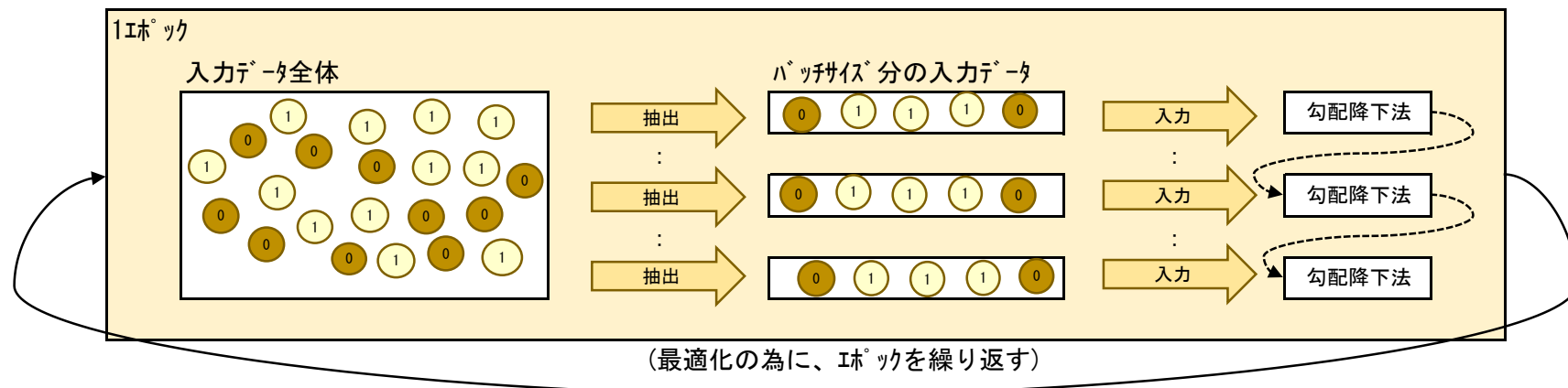
- ・勾配降下法によるパラメータの最適化の式では、(7.2) \_\_\_\_\_ と呼ばれるパラメータ $\eta$ (イータ)があります。  
これは、勾配降下法で、現ステップ $t$ のパラメータ $\theta^{(t)}$ におけるコスト関数 $C(\theta)$ の勾配から、次ステップ $t+1$ のパラメータ $\theta^{(t+1)}$ を計算する時に、  
勾配をどの程度の割合で変化分として反映させるか? という割合を示しています。  
「学習率」は「学習係数」とも言います。

- ・「学習率」を大きくすると、最適解を通り過ぎて見落としてしまう可能性が高くなり、  
「学習率」を小さくすると、最適解に辿り着くのに時間がかかったり、  
極小値が最小値とは限らないのに、最小値でない極小値  
(これを「局所最適解 (キョクショサイケイ、local optimal solution)」と言います)  
にはまってしまうと抜け出られない場合があるという欠点があります。  
これを(7.3) \_\_\_\_\_ と言います。  
尚、全領域にわたり最小値となる解を「大域的最適解 (タイケイサイケイ、global optimal solution)」  
と言います。





- ・入力データ全てに対する誤差関数の勾配を、更新の1ステップごとに計算する場合、データが大きいとその計算に時間がかかってしまいます。そのような場合には、データの一部で誤差関数の勾配を計算する、  
(7.4) \_\_\_\_\_ という方法が使われます。  
「確率的勾配降下法」では入力データ全体の一部で計算しますが、計算1回分に使用するデータの個数を(7.5) \_\_\_\_\_ と言います。  
「バッチサイズ」分の入力データを使用した計算で、パラメータが1回更新されることになります。
- ・入力データ全てを使いパラメータ更新を行うとき、この更新の単位を(7.6) \_\_\_\_\_ で表します。  
「確率的勾配法」では、1「エポック」のパラメータ更新を行うのに「バッチサイズ」分の入力データを使用した更新を何度も繰り返します。  
そしてパラメータが最適化されるまで「エポック」を何度も繰り返します(モデルによっては収束しない場合もあります)。



(選択肢)

- 「勾配降下法 (コハバ イウカホ, gradient descent method)」
- 「学習率 (ガクシュリツ, learning rate)」
- 「局所解問題 (キョクショカイモンダイ, local mininum problem)」
- 「確率的勾配降下法 (カクリツキコハ イウカホ, Stochastic gradient descent method, SGD)」
- 「バッチサイズ (batch size)」
- 「エポック (epochs)」

(回答)

(7.1)	
(7.2)	
(7.3)	
(7.4)	
(7.5)	
(7.6)	



## (8) 「順伝播型ニューラルネットワーク」の誤差逆伝搬法

以下の空欄に最もあてはまる語句を選択肢から選び、その記号を回答欄に記入してください。

- ・「階層型ニューラルネットワーク」の学習方法として(8.1) \_\_\_\_\_ が有名です。  
手法そのものは以前からありましたが、1986年に「backwards propagation of errors（後方への誤差伝播）」の略からデビッド・ラムハートらによって命名されたものです。
- ・「誤差逆伝播法」では、出力層での誤差(階層型ネットワークの出力で生じる教師データとの差)を使って、モデルの出力の計算方向が中間層⇒出力層の方向なのとは逆向きに、出力層⇒中間層の方向へと、勾配降下法を用いてパラメータを更新してゆきます(下図)。
- ・下図では、以下のような記法を用います。

$L$  : 層の総数（入力層を第1層、出力層を第 $L$ 層とします）

$J$  :  $l$ 層のユニット数

$l$  : 層番号（ $l=2\sim L$ 、入力層は層の厚さとしては数えませんが、説明上第1層とします）

$j$  :  $l$ 層でのユニット番号（ $j=1\sim J$ ）

$i$  :  $l-1$ 層でのユニット番号（ $i=1\sim I$ ）

$I$  :  $l-1$ 層のユニット数

$x^{l-1}_i$  :  $l-1$ 層 $i$ 番目のユニットの出力値（ $l$ 層への入力信号）

$w_{j,i}^l$  :  $l$ 層 $j$ 番目のユニットについて、 $l-1$ 層 $i$ 番目のユニットからの入力信号の重み

$w_{j,0}^l$  :  $l$ 層 $j$ 番目のユニットについて、 $l-1$ 層からの入力信号に付加するバイアス

$z_j^l$  :  $l$ 層 $j$ 番目のユニットの重み付き入力（ $z_j^l = \sum_{i=0}^I w_{j,i}^l * x^{l-1}_i$ （但し  $x^{l-1}_0=1$  でバイアス用））

$a^l$  :  $l$ 層の活性化関数（activation function）

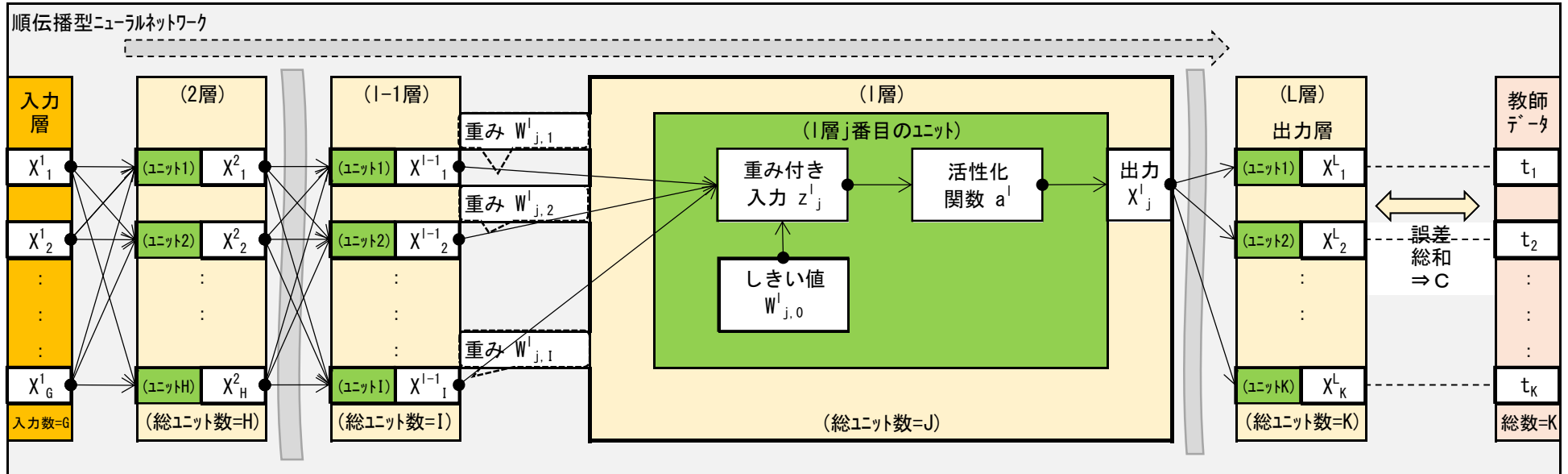
$x_j^l$  :  $l$ 層 $j$ 番目のユニットの出力値（ $x_j^l = a^l(z_j^l)$ ）

$C$  : コスト関数（平均二乗誤差(MSE)の場合  $C = (1/K) \sum_{k=1}^K (t_k^L - x_k^L)^2$ （ $K$  : は出力層のユニット数））

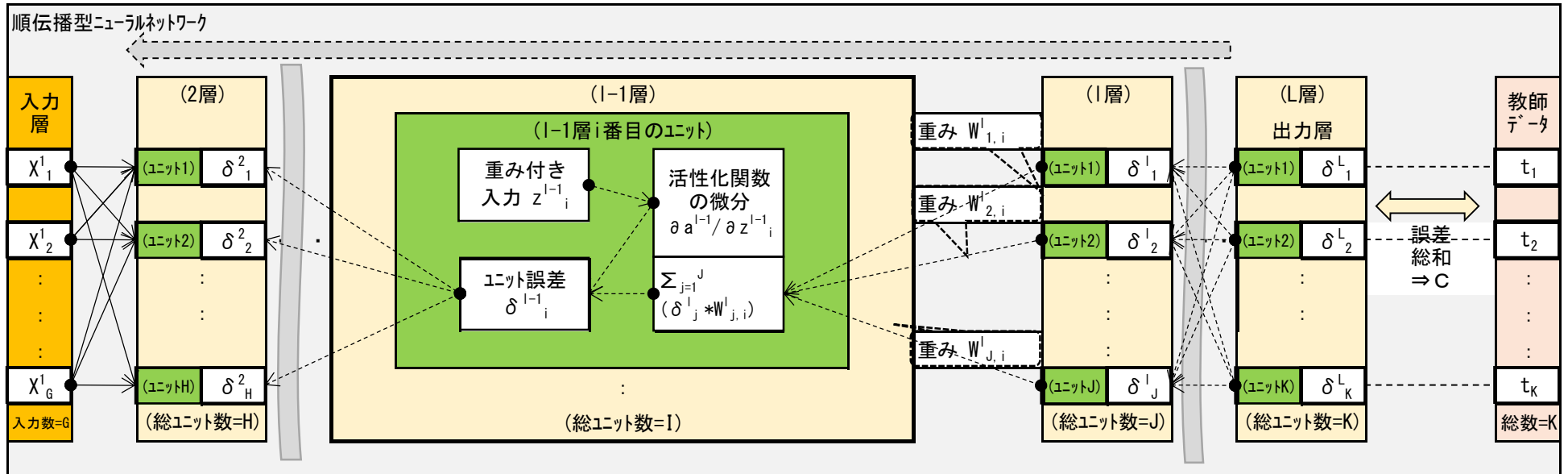
$\delta^{l-1}_i$  :  $l-1$ 層 $i$ 番目のユニットの誤差

$\delta_j^l$  :  $l$ 層 $j$ 番目のユニットの誤差

(出力計算時の流れ)



(誤差逆伝搬時の流れ)



- ・ 漸化式の場合は、初項と、第n+1項を第n項で表現する関係式を与えると、全ての項が計算できます。  
これに対し、逆漸化式の場合は、末項と第n-1項を第n項で表現する関係式を与えると、全ての項が計算できます。

- ・ 出力層のユニットの誤差で末項「 $\delta^L_*$ 」を与え、  
隣接する層間のユニットの誤差についての逆漸化式で、  
第l-1層のユニットの誤差「 $\delta^{l-1}_*$ 」を、第l層のユニットの誤差「 $\delta^l_*$ 」で表現する関係式を与えることにより、  
ユニットの誤差  $\delta$  は、全ての中間層のユニットについて計算することができます。

出力層j番目のユニットの誤差  $\delta^L_j$

$$\delta^L_j = (\partial C / \partial X^L_j) (\partial X^L_j / \partial z^L_j)$$

l層j番目のユニットと、l-1層の各ユニットとの間に成立する逆漸化式

$$\delta^{l-1}_i = \left\{ \sum_{j=1}^J \delta^l_j * W^l_{j,i} \right\} * \partial a^{l-1} / \partial z^{l-1}_i$$

- ・ これを元に、重みに関するコスト関数の偏微分とバイアスに関するコスト関数の偏微分を計算して、  
更に(8.2)を用いて、重み、バイアスの値を更新します。

$$\partial C / \partial W^l_{j,i} = \delta^l_j X^{l-1}_i \quad : \text{l層j番目のユニットについて、l-1層i番目のユニットからの入力信号の重みに関するコスト関数の偏微分}$$

$$\partial C / \partial W^l_{j,0} = \delta^l_j \quad : \text{l層j番目のユニットについてのバイアスに関するコスト関数の偏微分}$$

(選択肢)

- 「誤差逆伝播法 (コサギヤクデンパホウ, Backpropagation method)」
- 「勾配降下法 (コウバイカクホウ, gradient descent method)」

(回答)

(8. 1)	
(8. 2)	

(9) モデルの学習の検証

以下の空欄に最もあてはまる語句を選択肢から選び、その記号を回答欄に記入してください。

- ・モデルを作成し、学習を通してモデルのパラメータの最適化を行います。  
学習時のデータでは、コスト関数Cが最小になるが、学習時以外のデータでは、誤差が大きく出てしまうことがあります。  
これは、学習時のデータに特化した学習をしてしまって、汎用になっていないということであり、これを(9.1)\_\_\_\_\_と言います。
- ・「過学習」の対策として、入力にノイズを与えたり、学習時に正則化を適用したり、様々な方法が提案されてきましたが、(9.2)\_\_\_\_\_という方法が有効であることが知られています。
- ・「ドロップアウト」は、学習時に入力層や中間層のユニットを、各層毎に一定の割合  $\alpha$  ( $0 \leq \alpha < 1$ ) でユニットを選出し、それ以外は無効な状態で学習を行います（選出されなかったユニットはないものとして学習を行います）。これをミニバッチ毎にユニットを選出しながら、学習を繰り返します。
- ・学習後に予測をする場合は、全てのユニットを使用しますが、学習時は  $\alpha$  の割合のユニットしか存在しなかったため、予測時にはドロップアウトを施した層の出力の重みを  $\alpha$  倍にして出力が大きくなることを防ぎます。

(選択肢)

- (a) 「過学習 (かぶつき、over fitting)」
- (b) 「ドロップアウト (Dropout)」
- (c) 「ドロップイン (Dropin)」

(回答)

(9.1)	
(9.2)	

以上。

(11) 確認問題回答用紙

提出者

:

提出日

:

年

月

日

回答

(全 37 問)

No.	回答
(1. 1)	
(1. 2)	
(1. 3)	
(2. 1)	
(2. 2)	
(2. 3)	
(2. 4)	
(3. 1)	
(3. 2)	
(3. 3)	
(3. 4)	
(3. 5)	
(4. 1)	
(4. 2)	
(4. 3)	
(4. 4)	
(4. 5)	

No.	回答
(5. 1)	
(5. 2)	
(5. 3)	
(5. 4)	
(5. 5)	
(6. 1)	
(6. 2)	
(6. 3)	
(6. 4)	
(6. 5)	
(7. 1)	
(7. 2)	
(7. 3)	
(7. 4)	
(7. 5)	
(7. 6)	

No.	回答
(8. 1)	
(8. 2)	
(9. 1)	
(9. 2)	

(※黄色の枠のみ記入をお願いします)

※ ご意見・ご要望などありましたら、下欄に記してください。