

## 機械学習・深層学習の概観

人間の知能と機械が同レベルに達するか超えるものを「人工知能」と呼ぶが、ビッグデータが取り扱える時代になり機械学習の手法が身近になって、その実現に近づいている。

### 1 機械学習の要素

大量のデータから必要な情報を抽出する「機械学習」は、分類・回帰・クラスタリングなどに応用される。類似度を明示的に扱うか（いわゆる機械学習）、暗黙的に評価するか（ニューラルネットワーク neural network; NN, 多層にした NN を深層学習 deep learning と呼ぶ）によって大きく 2 種類に分かれる。

#### 1.1 3つの応用

- **分類 (Classification)**

データの類似性に基づいて、データをカテゴリー（あるいはクラス）に分ける機能である。郵便番号に記載された文字の判読、画像中の猫の判定などが代表例である。

データの中でどこに境界線（境界面）を引くか、という数学になる。

⇒ 数学的手法の例：教科書 §4.6 判別分析

⇒ 明示的方法の例： $k$  近傍法, サポートベクターマシン

- **クラスタリング (Clustering)**

類似しているデータどうしを集めて集合（クラスター）をつくり、データを仕分けする機能である。類似性分類の前処理でもある。正解か不正解かが明確ならば、「教師あり学習」と言える。

データ間の距離を定義する、という数学になる。

⇒ 数学的手法の例：教科書 §4.7 クラスタ分析

⇒ 明示的方法の例：階層クラスタリング法, 多次元尺度構成法, カーネル  $k$  平均法

- **回帰 (Regression)**

データの類似性に基づいて、データの数値を予測する機能である。株価予測や製造業における工程管理などに応用される。

データ  $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$  から、全体の傾向を知るための関数  $y = f(x)$  を導く数学になる。

⇒ 数学的手法の例：教科書 §4.3 回帰分析

⇒ 明示的方法の例：カーネル線形回帰

#### 1.2 類似度関数

データが類似していることを示す方法の 1 つは、データ間の距離関数を定義することである。距離が小さいほど類似している、という考えである。データ  $(x_1, \dots, x_n)$  と  $(y_1, \dots, y_n)$  のユークリッド距離（あるいはノルム）は、

$$\|\mathbf{x} - \mathbf{y}\| = \sqrt{(x_1 - y_1)^2 + \dots + (x_n - y_n)^2} \quad (1)$$

である。一般に、 $p$  ノルムを

$$\|\mathbf{x} - \mathbf{y}\|_p = \sqrt[p]{(x_1 - y_1)^p + \dots + (x_n - y_n)^p} \quad (2)$$

とする。また、重みをつけるさまざまな距離が提案されている。

別の方法として、データの内積をとる（カーネル関数  $k(x, y)$  を定義する）考えがある。値が大きいほど類似していることになる。もっとも簡単なものは

$$k(x, y) = \sum_i x_i y_i \quad (3)$$

とするものだ。より、一般には、カーネル関数は、グラム行列

$$G(x_1, \dots, x_n; k) = \begin{bmatrix} k(x_1, x_1) & \dots & k(x_1, x_n) \\ \vdots & \ddots & \vdots \\ k(x_n, x_1) & \dots & k(x_n, x_n) \end{bmatrix} \quad (4)$$

が正定値となればよい。（ $G$  が正定値とは、 $G = V^T V$  ( $V$  は行列) と分解できること、あるいは  $G$  が対称行列かつ固有値が非負であること）。

#### 1.3 目的関数 objective func.

データ  $(x_1, \dots, x_m)$  と  $(y_1, \dots, y_n)$  があり、両者の関係が線形であるとして、その関係式（目的関数）

$$y = ax + b \quad (5)$$

を導きたい。データからは

$$\begin{aligned} y_1 &= a_{11}x_1 + a_{12}x_2 + \dots + a_{1m}x_m + b_1 \\ y_2 &= a_{21}x_1 + a_{22}x_2 + \dots + a_{2m}x_m + b_2 \\ &\vdots \\ y_n &= a_{n1}x_1 + a_{n2}x_2 + \dots + a_{nm}x_m + b_n \end{aligned} \quad (6)$$

と書ける。これをベクトルと行列を用いて、

$$\mathbf{y} = \mathbf{A}\mathbf{x} + \mathbf{b} \quad (7)$$

と表す。以下のように定義した。

$$\mathbf{y} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix}, \quad \mathbf{x} = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_m \end{pmatrix}, \quad \mathbf{b} = \begin{pmatrix} b_1 \\ b_2 \\ \vdots \\ b_n \end{pmatrix},$$

$$\mathbf{A} = \begin{pmatrix} a_{11} & a_{12} & \dots & a_{1m} \\ a_{21} & a_{22} & \dots & a_{2m} \\ \vdots & \vdots & \ddots & \vdots \\ a_{n1} & a_{n2} & \dots & a_{nm} \end{pmatrix}.$$

## 2 目的関数最適化の例：最小二乗法

### 2.1 目的関数

例えば、最小二乗法で、データ  $\mathbf{x}$  と  $\mathbf{y}$  の関係式  $\mathbf{y} = \mathbf{A}\mathbf{x} + \mathbf{b}$  を求めたいとき、誤差

$$e = \sum_i \{y_i - (ax_i + b)\}^2$$

を最小にするような  $(a, b)$  を求める問題になる。ここでは誤差（距離）最小を原理としている。

これを一般化して、ある関数  $S$  を最小化（あるいは最大化でもよい）するような問題を考えよう。例えばノルム

$$S = \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2 \quad (8)$$

を最小化する問題を考える。 $\mathbf{A}$  を既知と考えると  $S$  を最小とする解  $\mathbf{x}_0$  を求めたい。式 (8) は

$$\begin{aligned} S &= (\mathbf{y} - \mathbf{A}\mathbf{x})^T (\mathbf{y} - \mathbf{A}\mathbf{x}) \\ &= \mathbf{y}^T (\mathbf{y} - \mathbf{A}\mathbf{x}) - \mathbf{x}^T \mathbf{A}^T (\mathbf{y} - \mathbf{A}\mathbf{x}) \\ &= \mathbf{y}^T \mathbf{y} - \mathbf{y}^T \mathbf{A}\mathbf{x} - \mathbf{x}^T \mathbf{A}^T \mathbf{y} + \mathbf{x}^T \mathbf{A}^T \mathbf{A}\mathbf{x} \\ &= \mathbf{y}^T \mathbf{y} - 2\mathbf{y}^T \mathbf{A}\mathbf{x} + \mathbf{x}^T \mathbf{A}^T \mathbf{A}\mathbf{x} \end{aligned}$$

$S$  を  $\mathbf{x}$  で微分して、(ベクトルの微分は成分ごとの微分をまとめたもの)

$$\frac{\partial S}{\partial \mathbf{x}} = -2\mathbf{A}^T \mathbf{y} + 2\mathbf{A}^T \mathbf{A}\mathbf{x}$$

となるので、 $S$  を最小にする  $\mathbf{x}$  は、この微分がゼロとなることから

$$\mathbf{x} = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{y} \quad (9)$$

などと求められる。

### 2.2 最小二乗解が不偏推定量であること

式 (7) の  $\mathbf{b}$  が誤差  $\mathbf{e}$  とする。 $\mathbf{e}$  は、平均ゼロ、分散  $\sigma_0$  の正規分布にしたがうとする。最小二乗解を導いた式 (9) で、真の値  $\mathbf{x}$  を  $\mathbf{x}_0$  とすると、

$$\begin{aligned} \hat{\mathbf{x}} &= (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T (\mathbf{A}\mathbf{x}_0 + \mathbf{e}) \\ &= \mathbf{x}_0 + (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{e}. \end{aligned} \quad (10)$$

誤差  $\mathbf{e}$  は正規分布にしたがう確率変数とすれば、 $\hat{\mathbf{x}}$  の各成分も正規分布にしたがう。期待値をとると、

$$E[\hat{\mathbf{x}}] = \mathbf{x}_0 + (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T E[\mathbf{e}] = \mathbf{x}_0 \quad (11)$$

最後の部分は、 $E[\mathbf{e}] = 0$  を用いた。この式は、最小二乗解  $\hat{\mathbf{x}}$  の期待値が真の値に一致するという保証（不偏推定量であること）を与えてくれる。

### 2.3 最小二乗解の分散

最適化された解を得たとしても、その不定性を添えないと不毛な議論になる。最小二乗解  $\hat{\mathbf{x}}$  の分散を計算しよう。

$$\begin{aligned} &E[(\hat{\mathbf{x}} - \mathbf{x}_0)(\hat{\mathbf{x}} - \mathbf{x}_0)^T] \\ &= E\{[(\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T E[\mathbf{e}]][(\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T E[\mathbf{e}]]^T\} \\ &= (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T E[\mathbf{e}\mathbf{e}^T] \mathbf{A} (\mathbf{A}^T \mathbf{A})^{-1} \\ &\stackrel{*}{=} (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T (\sigma_0^2 \mathbf{I}) \mathbf{A} (\mathbf{A}^T \mathbf{A})^{-1} = \sigma_0^2 (\mathbf{A}^T \mathbf{A})^{-1} \end{aligned}$$

ここで、\*部分では、 $V[\mathbf{e}] = \sigma_0^2 \mathbf{I}$  を用いた。この式で、 $(\mathbf{A}^T \mathbf{A})^{-1}$  の  $(i, i)$  成分を  $a_{ii}$  とする。

$\sigma_0^2$  の不偏推定量は、データ数を  $N$ 、モデルのパラメータ数を  $K$  とすれば、

$$\hat{\sigma}^2 = \frac{\|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2}{N - K} \quad (12)$$

となるので、パラメータ  $\hat{x}_i$  の分散は、 $\hat{\sigma}^2 a_{ii}$  となる。この平方根

$$\hat{s}(\hat{x}_i) = \hat{\sigma} \sqrt{a_{ii}} \quad (13)$$

は  $\hat{x}_i$  の標準誤差と呼ばれる。

## 3 ニューラルネット (NN) の考え

NN は暗黙的に目的関数を探る方法である。正解・不正解の判定を促すために、神経の閾値に見立てた活性化関数 (activation func.)

$$z = h(u), \quad \text{関数形は任意,} \quad \begin{cases} h(u) \rightarrow 1 & u \rightarrow +\infty \\ h(u) \rightarrow 0 & u \rightarrow -\infty \end{cases} \quad (14)$$

を定義する。そして式 (5) を得るために、式 (6) のように、 $n$  本の連立 (ユニット) を考え、

$$\mathbf{z}^{(1)} = h(\mathbf{A}^{(1)} \mathbf{x} + \mathbf{b}^{(1)}) \quad (15)$$

とする。これをさらに多段階の中間層を考え (ユニット数は各層で異なってよい)、

$$\mathbf{z}^{(1)} = h(\mathbf{A}^{(1)} \mathbf{x} + \mathbf{b}^{(1)}) \quad (16)$$

$$\mathbf{z}^{(2)} = h(\mathbf{A}^{(2)} \mathbf{z}^{(1)} + \mathbf{b}^{(2)}) \quad (17)$$

$$\vdots$$

$$\mathbf{z}^{(N)} = h(\mathbf{A}^{(N)} \mathbf{z}^{(N-1)} + \mathbf{b}^{(N)}) \quad (18)$$

とする。これらをまとめて、入力  $\mathbf{x}$  から出力  $\mathbf{y} = \mathbf{z}^{(N)}$  への対応式とみなすのである。決めるべきパラメータが多いので、学習させる量も多く必要になるが、多層にすれば予測精度がよくなると期待される。

ユニット間のつながり方にもさまざまな提案がなされている。すべてのユニットを次のユニットにつなぐ **全結合** (fully-connected) か、入力の一部だけを次の層のユニットにつなぐ **畳み込み NN** (convolutional NN; CNN) か。(CNN は画像判別問題でよく使われる)。多層を順にこなす **順伝播型** (feedforward) か、一度出力された層の結果をもう一度使う **再帰型** (recurrent NN; RNN) か。(RNN は時系列データでよく使われる)。

## 参考文献

- [1] 植村誠 『物理のためのデータサイエンス入門』 (講談社, 2023)
- [2] 申吉浩 監修 『機械学習アルゴリズム入門 類似性の科学』 (工学社, 2022)
- [3] 真貝寿明 『徹底攻略 確率統計』 (共立出版, 2012)