

## 最小二乗法によるビッグデータ解析

東京工業大学 杉山将

本節では、最小二乗法を用いた教師付きのビッグデータ解析手法を概説する<sup>1)</sup>。

### 1. 教師付き学習

教師付き学習とは、入力 $\mathbf{x}$ と出力 $y$ とをつなぐ関数 $y = f(\mathbf{x})$ を、入出力が対になった $n$ 個の訓練標本 $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$ から学習する問題である。入力 $\mathbf{x}$ は一般に $d$ 次元の実数値ベクトルであり、出力 $y$ は実数値、あるいは、 $\pm 1$ のようなカテゴリ値を取る。出力 $y$ が実数値を取る場合は回帰問題とよび (図1)、カテゴリ値を取る場合は分類問題とよぶ (図2)。

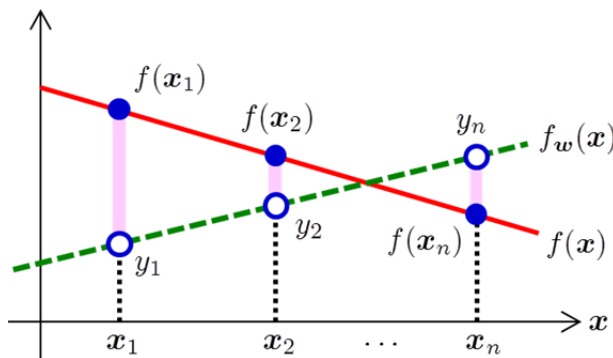


図1：教師付き回帰.

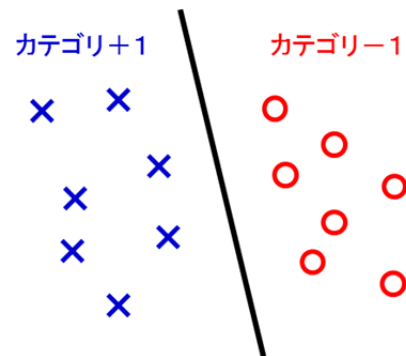


図2：教師付き分類

ビッグデータとは、入力 $\mathbf{x}$ の次元数 $d$ と訓練標本数 $n$ が非常に大きい状況を指す。ビッグデータ解析においては、以下の単純な線形モデル $f_w$ を用いて関数 $f$ をモデル化する事が多い。

$$f_w(\mathbf{x}) = \mathbf{x}^T \mathbf{w}$$

ここで $\mathbf{w}$ はモデルのパラメータであり $d$ 次元の実数値ベクトルである。モデル $f_w$ が真の関数 $f$ にできるだけ近くなるようにパラメータ $\mathbf{w}$ を学習することが、教師付き学習の目的である。

以下では、まずは出力 $y$ が実数値を取る回帰問題に対する学習法を紹介し、続いて出力 $y$ がカテゴリ値を取る分類問題に対する学習法を紹介する。最後に、最小二乗法の更なる応用について述べる。

### 2. 最小二乗回帰とその確率的勾配法による実装

回帰問題では、パラメータ $\mathbf{w}$ を次の二乗損失を最小化するように学習する。

$$\min_{\mathbf{w}} \sum_i^n (y_i - \mathbf{x}_i^T \mathbf{w})^2$$

この最小解 $\hat{\mathbf{w}}$ は以下のように解析的に求められる。

$$\hat{\mathbf{w}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

ただし、 $\mathbf{X}$ は訓練入力標本を並べた $n \times d$ 行列であり、 $\mathbf{y}$ は訓練出力標本を並べた $n$ 次元の縦ベクトルである。

$$\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_n)^T, \quad \mathbf{y} = (y_1, \dots, y_n)^T$$

しかし、入力 $\mathbf{x}$ の次元数 $d$ と訓練標本数 $n$ が非常に大きいとき、訓練入力標本行列 $\mathbf{X}$ をメモリに保持することは困難である。また、 $d \times d$ 行列である $\mathbf{X}^T \mathbf{X}$ の逆を計算することも困難である。

そこで、次の確率的勾配法を用いて逐次的に解を求めることにする (図3)。

- (1) パラメータ $\mathbf{w}$ を適当に初期化する。
- (2)  $n$ 個の訓練標本 $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$ からランダムに一つ選ぶ。
- (3) 選んだ訓練標本 $(\mathbf{x}, y)$ に対する二乗誤差 $(y - \mathbf{x}^T \mathbf{w})^2$ を小さくするようにパラメータ $\mathbf{w}$ を更新する。

$$\mathbf{w} \leftarrow \mathbf{w} + \epsilon \cdot (y - \mathbf{x}^T \mathbf{w}) \mathbf{x}$$

ただし、 $\epsilon$ は勾配降下の歩幅を表す小さい正のスカラールである。

- (4) 収束するまで (2) ~ (3) を繰り返す。

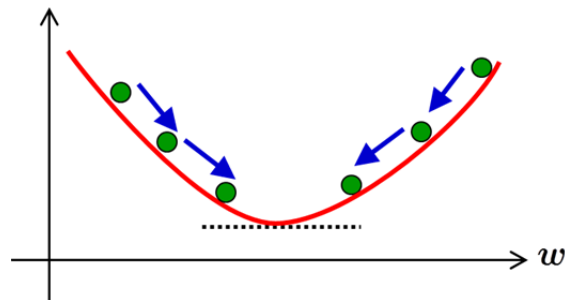


図3：確率的勾配法。誤差規準を小さくするように勾配を降下していく。

### 3. 受動攻撃回帰

確率的勾配法はビッグデータ解析の有効な学習アルゴリズムであるが、勾配降下の歩幅 $\epsilon$ の決め方が難しいという問題がある。すなわち、 $\epsilon$ が大きすぎると収束せず (図4)、 $\epsilon$ が小さすぎると収束に時間がかかる (図5)。最初は $\epsilon$ を大きめに設定し、学習が進むにつれて徐々に $\epsilon$ を小さくしていくのが妥当と考えられるが、どのくらいの速さで $\epsilon$ を減らすかを定めるのもまた困難である。

一方、最小二乗回帰では一気に谷底までジャンプするようにパラメータを更新することが可能である。そうすると勾配降下の歩幅 $\epsilon$ を決める必要がなく、実用上非常に便利である。

しかし、このような急激な学習を行うとパラメータの値が大きく変化するため、これまでの学習で得られていた知識がたった一つのデータの追加によって失われてしまう恐れがある。

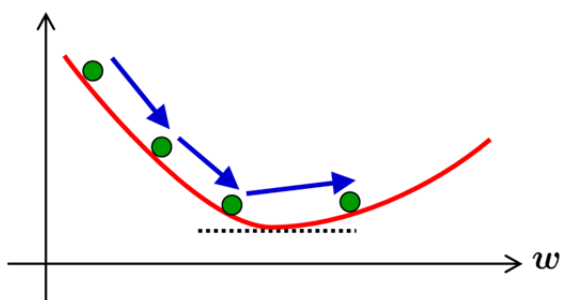


図4：確率的勾配法の問題点1.  
歩幅が大きすぎると収束しない。

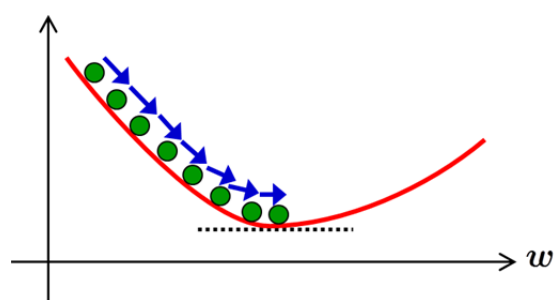


図5：確率的勾配法の問題点2.  
歩幅が小さすぎると収束が遅い。

この問題を避けるためには、現在の解からの移動量を考慮しながら勾配降下を行う受動攻撃回帰が有効である。受動攻撃回帰では、次式を最小にするようにパラメータ $\mathbf{w}$ を更新する。

$$\min_{\mathbf{w}} [(y - \mathbf{x}^T \mathbf{w})^2 + \gamma \|\mathbf{w} - \hat{\mathbf{w}}\|^2]$$

ただし、 $\hat{\mathbf{w}}$ は現在のパラメータ値、 $\|\cdot\|$ はユークリッドノルム、 $\gamma$ は現在の解からの移動量を調整する正のスカラをそれぞれ表す。この最小化問題は解析的に解くことができ、次のパラメータ更新式が得られる。

$$\mathbf{w} \leftarrow \mathbf{w} + \frac{y - \mathbf{x}^T \mathbf{w}}{\|\mathbf{x}\|^2 + \gamma} \mathbf{x}$$

これは、確率的勾配法において歩幅を $\epsilon = 1/(\|\mathbf{x}\|^2 + \gamma)$ とおくことに対応しており、受動攻撃回帰では $\mathbf{x}$ のノルムに合わせて適応的に歩幅を調整していることがわかる。

#### 4. 適応正規化回帰

統計的な枠組みのもとでパラメータ $\mathbf{w}$ をデータから学習する限りは、推定に対する不確定性が必ず伴う。すなわち、パラメータ $\mathbf{w}$ も確率分布を持つ。受動攻撃回帰では、現在のパラメータ値 $\hat{\mathbf{w}}$ からの移動量をユークリッドノルムで測ったが、パラメータが確率分布を持つ場合は、例えばカルバック・ライブラー距離など、確率分布間の距離尺度を用いる方が自然である。確率密度関数 $p$ から $q$ へのカルバック・ライブラー距離は、次式で定義される。

$$\text{KL}(p||q) = \int p(\mathbf{x}) \log \frac{p(\mathbf{x})}{q(\mathbf{x})} d\mathbf{x}$$

パラメータの分布として、ガウス分布を用いることにする。ガウス分布の確率密度関数は次式で定義される (図6)。

$$N_{\mathbf{w}}(\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{d/2} \det(\boldsymbol{\Sigma})^{1/2}} \exp\left(-\frac{1}{2}(\mathbf{w} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{w} - \boldsymbol{\mu})\right)$$

ただし、 $\boldsymbol{\mu}, \boldsymbol{\Sigma}$ はガウス分布の期待値ベクトルと分散共分散行列を表し、 $\det(\cdot)$ は行列式を表す。

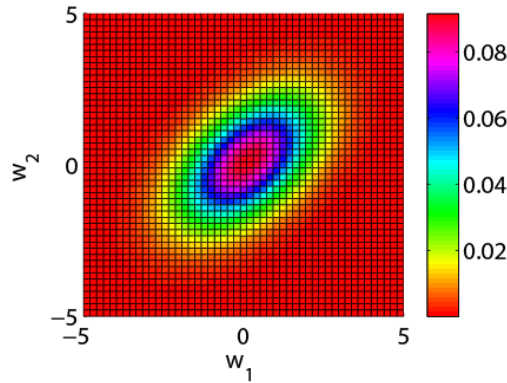


図6：ガウス分布。楕円形の等高線を持つ。

適応正則化回帰では、パラメータ $\mathbf{w}$ そのものでなくその期待値と分散共分散行列 $\boldsymbol{\mu}, \boldsymbol{\Sigma}$ を考え、次式を最小にするように $\boldsymbol{\mu}, \boldsymbol{\Sigma}$ を更新する。

$$\min_{\boldsymbol{\mu}, \boldsymbol{\Sigma}} \left[ (y - \mathbf{x}^T \boldsymbol{\mu})^2 + 2\gamma \cdot \text{KL} \left( N(\boldsymbol{\mu}, \boldsymbol{\Sigma}) \parallel N(\tilde{\boldsymbol{\mu}}, \tilde{\boldsymbol{\Sigma}}) \right) + \mathbf{x}^T \boldsymbol{\Sigma} \mathbf{x} \right]$$

ただし、 $\tilde{\boldsymbol{\mu}}, \tilde{\boldsymbol{\Sigma}}$ は現在のパラメータ $\tilde{\mathbf{w}}$ の期待値と分散共分散行列を表す。第3項目の $\mathbf{x}^T \boldsymbol{\Sigma} \mathbf{x}$ は、分散共分散行列 $\boldsymbol{\Sigma}$ を標本 $\mathbf{x}$ に合わせて適応的に正則化する役割を果たす。この最小化問題は解析的に解くことができ、次の更新式が得られる。

$$\boldsymbol{\mu} \leftarrow \boldsymbol{\mu} + \frac{y - \mathbf{x}^T \boldsymbol{\mu}}{\mathbf{x}^T \boldsymbol{\Sigma} \mathbf{x} + \gamma} \boldsymbol{\Sigma} \mathbf{x}, \quad \boldsymbol{\Sigma} \leftarrow \boldsymbol{\Sigma} - \frac{1}{\mathbf{x}^T \boldsymbol{\Sigma} \mathbf{x} + \gamma} \boldsymbol{\Sigma} \mathbf{x} \mathbf{x}^T \boldsymbol{\Sigma}$$

$\boldsymbol{\mu}$ の更新式を受動攻撃回帰の $\mathbf{w}$ の更新式と比べると、分散共分散行列 $\boldsymbol{\Sigma}$ が表す楕円形状（図6）を考慮した形になっていることがわかる。 $\boldsymbol{\Sigma}$ が単位行列のとき（すなわち楕円が単位円になるとき）、 $\boldsymbol{\mu}$ の更新式は受動攻撃回帰の $\mathbf{w}$ の更新式と一致する。

分散共分散行列 $\boldsymbol{\Sigma}$ の計算が煩雑な場合は、対角成分のみを残して非対角成分をゼロにするという近似を行うことにより、計算時間とメモリ消費量を大幅に削減できる。

## 5. 二値分類問題への拡張

ここまで、出力 $y$ が実数値を取る回帰問題に対する学習アルゴリズムを紹介してきた。二値分類問題では出力 $y$ は $\pm 1$ のカテゴリ値を取るが、これを実数とみなせば上記の回帰アルゴリズムをそのまま分類学習にも適用できる。そして、テスト入力 $\mathbf{x}$ に対するモデル $f_{\mathbf{w}}(\mathbf{x})$ の出力をその正負に従って $\pm 1$ のカテゴリ値に変換すれば、テスト入力の分類が行なえる。

しかし、分類問題では学習した関数の符号のみが重要であることから、分類学習に二乗

誤差を用いるのはやや不自然である。分類学習では、誤分類数に対応する0/1損失の最小化が理想的である。

$$\min_{\mathbf{w}} \sum_i^n I(y_i \mathbf{x}_i^T \mathbf{w} > 0)$$

ここで、 $I(\cdot)$ は条件が真なら0ば、偽ならば1を出力する指示関数である。しかし、0/1損失の最小化は $n$ 個の標本のカテゴリ割り当てを求める離散最適化問題に帰着されるため、 $2^n$ の組み合わせから最適なカテゴリの割り当てを求めることは一般に困難である。

そこで、次の二乗ヒンジ損失を代理損失として用いることにする。

$$\min_{\mathbf{w}} \sum_i^n (\max(0, 1 - y_i \mathbf{x}_i^T \mathbf{w}))^2$$

図7に示すように、二乗ヒンジ損失は0/1損失の滑らかな上界となっている。一方、出力 $y$ が $\pm 1$ のカテゴリ値を取るとき、二乗損失の最小化は

$$\min_{\mathbf{w}} \sum_i^n (1 - y_i \mathbf{x}_i^T \mathbf{w})^2$$

と表現できる。図7に示すように、二乗損失も0/1損失の上界となっていることがわかるが、右側で損失が増加するため0/1損失と形状が大きく異なる。

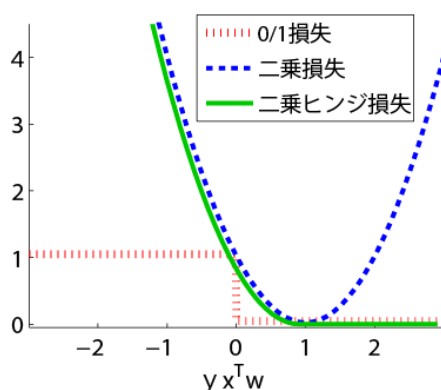


図7：分類問題の損失関数。

この二乗ヒンジ損失を用いた分類学習規準に対しても、回帰の場合と同様にして学習アルゴリズムを導出できる。

- 確率的勾配分類のパラメータ更新式：

$$\mathbf{w} \leftarrow \mathbf{w} + \epsilon \cdot \max(0, 1 - y \mathbf{x}^T \mathbf{w}) y \mathbf{x}$$

- 受動攻撃分類のパラメータ更新式：

$$\mathbf{w} \leftarrow \mathbf{w} + \frac{\max(0, 1 - y \mathbf{x}^T \mathbf{w}) y}{\|\mathbf{x}\|^2 + \gamma} \mathbf{x}$$

- 適応正則化分類のパラメータ更新式：

$$\boldsymbol{\mu} \leftarrow \boldsymbol{\mu} + \frac{\max(0, 1 - \mathbf{y}\mathbf{x}^T\boldsymbol{\mu})\mathbf{y}}{\mathbf{x}^T\boldsymbol{\Sigma}\mathbf{x} + \gamma}\boldsymbol{\Sigma}\mathbf{x}, \quad \boldsymbol{\Sigma} \leftarrow \boldsymbol{\Sigma} - \frac{1}{\mathbf{x}^T\boldsymbol{\Sigma}\mathbf{x} + \gamma}\boldsymbol{\Sigma}\mathbf{x}\mathbf{x}^T\boldsymbol{\Sigma}$$

これらは、回帰学習のパラメータ更新式の $\mathbf{y} - \mathbf{x}^T\mathbf{w}$ を $\max(0, 1 - \mathbf{y}\mathbf{x}^T\mathbf{w})$ に置き換えた形になっている。

## 6. まとめ

本節では、最小二乗法に基づく教師付きのデータ解析手法を紹介した。これらの学習法はオンライン学習とよばれる手法であり、わずか数行のプログラムによって簡単に実装できるという優れた特徴を持っている。ビッグデータ解析にはこのような単純な学習法が特に有効であり、今後益々重要性が高まってくるものと思われる。

このように回帰と分類に対してはビッグデータ解析に適した優れた学習アルゴリズムが開発されているが、回帰と分類以外にも、クラスタリング、異常検知、特徴選択など様々なデータタスクが存在する。機械学習の最も汎用的なアプローチは、データを生成する確率分布を推定することである<sup>2)</sup>。なぜならば、データの生成分布を知ることは、そのデータに関する全ての知識を得ることと本質的に等価だからである。従って、データの生成分布がうまく推定できれば、あらゆるデータ解析タスクを精度良く解決できる。しかし、高次元の確率分布をうまく推定することは一般に困難であり、データの生成分布に基づくアプローチでは必ずしも精度良くデータ解析を行えない。

一方、クラスタリング、異常検知、特徴選択などの様々なデータ解析タスクそれぞれに対して、ビッグデータ解析に適した優れたアルゴリズムを個別に開発することも現実的には困難である。実際、回帰や分類の基礎理論は既に1960年ころには研究されており、半世紀にも及ぶ更なる研究開発を経て、様々な実データ解析に応用されるようになった。ビッグデータ時代には大量のデータを前に迅速な意思決定を行うことが望まれており、様々なデータ解析タスクに対して長期間に及ぶ研究開発を行うことは現実的でない。

この状況を打破するために、生成分布推定アプローチとタスク特化アプローチの中間的なアプローチが近年提案された<sup>3,4)</sup>。それは、ある条件を満たすデータ解析タスクのクラスに対して学習アルゴリズムを開発するというアプローチである(図8)。具体的には、複数の確率分布が含まれるデータ解析タスクのうち、それらの確率分布そのものは必要なく、確率密度関数の比さえわかればデータ解析を行えるというクラスを考える。このクラスには、非定常環境適応学習<sup>5,6)</sup>、特徴選択、クラスタリング<sup>7)</sup>、条件付き確率推定、独立成分分析、異常値検出、変化検知<sup>8)</sup>など、多くの重要なデータ解析タスクが含まれる。そして、この確率密度関数の比を、それぞれの確率密度関数を推定することなく直接推定する。密度比の推定は確率密度の推定よりも一般に容易であるため(図9)、密度比推定によって上記の全てのデータ解析タスクを統一的かつ優れた精度で解決できる。そして、密度比は最小二乗法によって推定できるため、本節で紹介した学習アルゴリズムと同等の考え方を密度比推定にも適用できる。

密度比推定によって様々な機械学習タスクを統一的に解決できるため、密度比推定の精度や計算効率を更に向上させることにより、様々な機械学習アルゴリズムの性能を一挙に

改善できる。今後は、密度比推定の基礎技術を更に発展させていくとともに、密度比推定によって解決できる新たなデータ解析タスクを開拓し、それらの機械学習技術を様々な実世界問題の解決に活用していくことが期待される。



図8：機械学習のアプローチ。各●印は回帰、分類などのデータ解析タスクを表す。

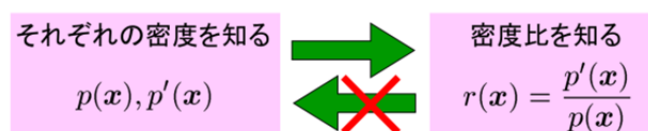


図9：密度比推定。それぞれの密度を推定するよりも密度比を推定するほうが容易である。

## 参考文献

- 1) 杉山 将. イラストで学ぶ機械学習：最小二乗法による識別モデル学習を中心に，講談社，2013.
- 2) 杉山 将. 統計的機械学習：生成モデルに基づくパターン認識，オーム社，2009.
- 3) M. Sugiyama et al. Density Ratio Estimation in Machine Learning, Cambridge University Press, 2012.
- 4) 杉山 将. 機械学習入門. オペレーションズ・リサーチ, vol.57, no.7, pp.353-359, 2012.
- 5) M. Sugiyama et al. Machine Learning in Non-Stationary Environments: Introduction to Covariate Shift Adaptation, MIT Press, 2012.
- 6) 杉山 将. 非定常環境下での教師付き学習：データの入力分布が変化する場合. 画像ラボ, vol.18, no.10, pp.1-6, 2007.
- 7) 杉山 将. 機械学習によるデータの自動クラスタリング. シミュレーション, vol.31, no.2, pp.36-40, 2012.
- 8) 杉山 将. 確率分布間の距離推定：機械学習分野における最新動向. 日本応用数理学会論文誌, vol.23, no.3, pp.439-452, 2013.