

線形回帰とロジスティック回帰

中田和秀

東京工業大学 工学院 経営工学系

機械学習入門

<http://www.iee.e.titech.ac.jp/~nakatalab/text/lecture.html>

概要

ここでは線形モデルによって回帰や判別を行う方法について説明する。これらは昔から使われてきた単純な方法であるが、計算が容易である、解析しやすい、人間が理解しやすいなどの優れた特徴を持っている。

目次：

1. 線形回帰

1.1 教師あり学習からの導出

1.2 統計的な解析

1.3 正則化項

1.4 学習アルゴリズム

2. ロジスティック回帰（線形モデルによる判別）

2.1 教師あり学習からの導出

2.2 統計的な解析

2.3 学習アルゴリズム

記号の使い方：

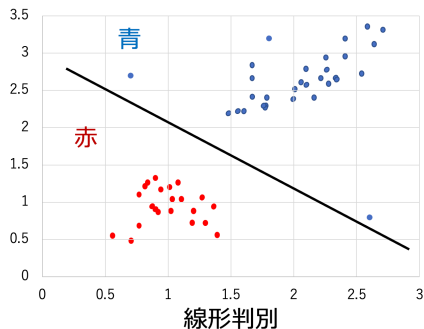
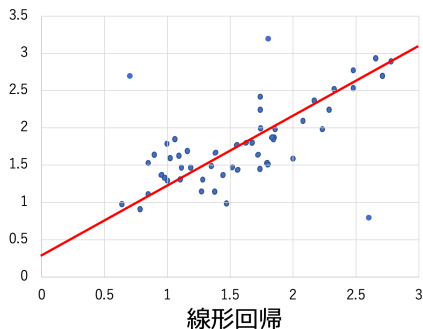
- $A := B$ は、 B で A を定義する、 B を A に代入することを意味する
- $[n]$ は n までのインデックスの集合を表し $[n] := \{1, 2, \dots, n\}$

線形モデル

線形モデルによる教師あり学習

- 回帰：線形回帰
- 判別・分類：ロジスティック回帰

線形性： $f(\alpha_1 \mathbf{x}_1 + \alpha_2 \mathbf{x}_2) = \alpha_1 f(\mathbf{x}_1) + \alpha_2 f(\mathbf{x}_2)$



線形モデルの特徴

特徴

- 計算（学習）や理論的解析が容易
- 人間が理解しやすい
- 複雑な現象（非線形性を持つもの）を説明できない

データ： $\{(\mathbf{x}_d, y_d)\}_{d \in [D]}$, $\mathbf{x}_d \in \mathbb{R}^n$, $y_d \in \mathbb{R}$

線形関数（アフィン関数）： $y = \mathbf{w}^T \mathbf{x} + b$

$$\tilde{\mathbf{x}} := \begin{pmatrix} 1 \\ \mathbf{x} \end{pmatrix} \in \mathbb{R}^{1+n}, \quad \tilde{\mathbf{w}} := \begin{pmatrix} b \\ \mathbf{w} \end{pmatrix} \in \mathbb{R}^{1+n} \quad \text{とすると、} \quad \mathbf{w}^T \mathbf{x} + b = \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}$$

以下では、1次元を増やした $\tilde{\mathbf{x}}_d, \tilde{\mathbf{w}}$ を $\mathbf{x}_d, \mathbf{w} \in \mathbb{R}^n$ として考える。

線形回帰

予測器

$$y \simeq \hat{y} = \mathbf{w}^T \mathbf{x} \quad (\text{パラメタは } \mathbf{w})$$

誤差関数

予測誤差の2乗

$$L(y, \hat{y}) := \frac{1}{2}(y - \hat{y})^2 \quad y: \text{目標} \quad \hat{y}: \text{予測値}$$

学習

$$\min_{\mathbf{w}} \sum_{d \in [D]} (y_d - \mathbf{w}^T \mathbf{x}_d)^2$$

目的関数の $1/2D$ は省略

線形回帰は最小2乗法や重回帰とも呼ばれる

学習アルゴリズム

$$\mathbf{X} := \begin{pmatrix} \mathbf{x}_1^T \\ \mathbf{x}_2^T \\ \vdots \\ \mathbf{x}_D^T \end{pmatrix} \in \mathbb{R}^{D \times n}, \quad \mathbf{y} := \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_D \end{pmatrix} \in \mathbb{R}^D \quad \text{とする。}$$

学習の行列表現

$$\min_{\mathbf{w}} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|^2$$

\mathbf{X} が列フルランク $\text{rank}(\mathbf{X}) = n$ と仮定

学習アルゴリズム

最適解は解析的に与えられる

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

$$\begin{aligned} F(\mathbf{w}) &:= \|\mathbf{y} - \mathbf{X}\mathbf{w}\|^2 = (\mathbf{y} - \mathbf{X}\mathbf{w})^T (\mathbf{y} - \mathbf{X}\mathbf{w}) \\ &= \mathbf{y}^T \mathbf{y} - 2\mathbf{y}^T \mathbf{X}\mathbf{w} + \mathbf{w}^T \mathbf{X}^T \mathbf{X}\mathbf{w} \end{aligned}$$

最適解の必要条件 $\nabla F(\mathbf{w}^*) = \mathbf{0}$ を考える。

$$\nabla F(\mathbf{w}^*) = -2\mathbf{X}^T \mathbf{y} + 2\mathbf{X}^T \mathbf{X}\mathbf{w}^* = \mathbf{0}$$

$$\iff \mathbf{X}^T \mathbf{X}\mathbf{w}^* = \mathbf{X}^T \mathbf{y}$$

$$\iff \mathbf{w}^* = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \quad (\because \text{rank}(\mathbf{X}) = n \text{ より逆行列は存在})$$

また、 $\nabla^2 F(\mathbf{w})$ は半正定値である。なぜなら、任意の $\mathbf{v} \in \mathbb{R}^n$ に対して

$$\mathbf{v}^T \nabla^2 F(\mathbf{w}) \mathbf{v} = 2\mathbf{v}^T \mathbf{X}^T \mathbf{X}\mathbf{v} = 2(\mathbf{X}\mathbf{v})^T (\mathbf{X}\mathbf{v}) \geq 0$$

よって、 $F(\mathbf{w})$ は凸関数となり、 $\nabla F(\mathbf{w}^*) = \mathbf{0}$ は最適解の必要十分条件

以上より、 $\mathbf{w}^* = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$ が最適解。

確率モデルに基づいた解析

データ $\{(\mathbf{x}_d, y_d)\}_{d \in [D]}$ が次の確率モデルから生成されたサンプル（標本）とする。

$$y_d = \mathbf{w}^T \mathbf{x}_d + \epsilon_d \quad (d \in [D])$$

$$\iff \mathbf{y} = \mathbf{X}\mathbf{w} + \boldsymbol{\epsilon}$$

$\boldsymbol{\epsilon}$ はノイズ（測定誤差など）を表す確率変数

ここで確率変数 $\boldsymbol{\epsilon} \in \mathbb{R}^D$ が従う仮定として、次のものを考える。

誤差項に対する仮定

- 仮定 1: $\mathbb{E}[\boldsymbol{\epsilon}] = \mathbf{0}$
- 仮定 2: $\mathbb{V}[\boldsymbol{\epsilon}] = \sigma^2 \mathbf{I}$
- 仮定 3: $\boldsymbol{\epsilon} \sim N(\mathbf{0}, \sigma^2 \mathbf{I})$

\mathbf{I} は単位行列

※ 仮定 3 は仮定 1, 2 を含むより強い仮定

不偏推定量

不偏推定量

仮定 1 のもとで、推定量 $\mathbf{w}^* := (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$ は不偏推定量である。

$$\mathbb{E}[\mathbf{w}^*] = \mathbf{w}$$

証明：

$$\begin{aligned}\mathbf{w}^* &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \\ &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T (\mathbf{X} \mathbf{w} + \boldsymbol{\epsilon}) \\ &= \mathbf{w} + (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \boldsymbol{\epsilon}\end{aligned}$$

$$\begin{aligned}\mathbb{E}[\mathbf{w}^*] &= \mathbb{E}[\mathbf{w} + (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \boldsymbol{\epsilon}] \\ &= \mathbf{w} + (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbb{E}[\boldsymbol{\epsilon}] \\ &= \mathbf{w} \quad (\because \text{仮定 1})\end{aligned}$$

この意味でこの計算に正当性があるといえる

最良線形不偏推定量

仮定 2 のもとで

$$\begin{aligned}\mathbb{V}[\mathbf{w}^*] &= \mathbb{V}[\mathbf{w} + (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \boldsymbol{\epsilon}] \\ &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbb{V}[\boldsymbol{\epsilon}] \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \\ &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \sigma^2 \mathbf{I} \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \quad (\because \text{仮定 2}) \\ &= \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1}\end{aligned}$$

ガウス・マルコフの定理

仮定 1,2 のもとで、 $\mathbf{w}^* := (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$ は、 \mathbf{y} の線形変換で表される不偏推定量 $\tilde{\mathbf{w}}$ の中で最も分散が小さい。

$$\mathbb{V}[\tilde{\mathbf{w}}] \succeq \mathbb{V}[\mathbf{w}^*]$$

※ $\mathbf{A} \succeq \mathbf{B}$ は、 $\mathbf{A} - \mathbf{B}$ が半正定値行列であることを意味する

この性質を最良線形不偏推定量や BLUE(best linear unbiased estimator) と呼ぶ

証明

推定量 \tilde{w} として、 y の線形変換で表せるものを考える。

$$\tilde{w} = Cy$$

任意の w に対して不偏推定量であるためには、

$$w = \mathbb{E}[\tilde{w}] = \mathbb{E}[C(Xw + \epsilon)] = CXw$$

である必要があるため、 $CX = I$.

$$\mathbb{V}[\hat{w}] = \mathbb{V}[C(Xw + \epsilon)] = C\mathbb{V}[\epsilon]C^T = C\sigma^2 IC^T = \sigma^2 CC^T$$

$$\mathbb{V}[w^*] = \sigma^2 (X^T X)^{-1}$$

$P := C - (X^T X)^{-1} X^T$ とすると、

$$\begin{aligned} PP^T &= (C - (X^T X)^{-1} X^T)(C - (X^T X)^{-1} X^T)^T \\ &= CC^T - (X^T X)^{-1} X^T C^T - CX(X^T X)^{-1} + (X^T X)^{-1} X^T X (X^T X)^{-1} \\ &= CC^T - (X^T X)^{-1} - (X^T X)^{-1} + (X^T X)^{-1} \\ &= CC^T - (X^T X)^{-1} \end{aligned}$$

証明 (続き)

つまり、

$$\mathbb{V}[\tilde{\mathbf{w}}] - \mathbb{V}[\mathbf{w}^*] = \sigma^2 \mathbf{C}\mathbf{C}^T - \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1} = \sigma^2 \mathbf{P}\mathbf{P}^T$$

となるため、

$$\mathbb{V}[\tilde{\mathbf{w}}] \succeq \mathbb{V}[\mathbf{w}^*]$$

この意味で \mathbf{w}^* は線形不偏推定量 $\hat{\mathbf{w}}$ の中で最も分散が小さいといえる。

また、 $\text{Tr}(\mathbb{V}[\tilde{\mathbf{w}}]) \geq \text{Tr}(\mathbb{V}[\mathbf{w}^*])$ が成り立つ。そして、

$$\begin{aligned} \text{Tr}(\mathbb{V}[\tilde{\mathbf{w}}]) &= \text{Tr}(\mathbb{E}[(\tilde{\mathbf{w}} - \mathbf{w})(\tilde{\mathbf{w}} - \mathbf{w})^T]) = \mathbb{E}[\text{Tr}((\tilde{\mathbf{w}} - \mathbf{w})(\tilde{\mathbf{w}} - \mathbf{w})^T)] \\ &= \mathbb{E}[\text{Tr}((\tilde{\mathbf{w}} - \mathbf{w})^T(\tilde{\mathbf{w}} - \mathbf{w}))] = \mathbb{E}[\|\tilde{\mathbf{w}} - \mathbf{w}\|^2]. \end{aligned}$$

同様に、 $\text{Tr}(\mathbb{V}[\mathbf{w}^*]) = \mathbb{E}[\|\mathbf{w}^* - \mathbf{w}\|^2]$ である。よって、

$$\mathbb{E}[\|\tilde{\mathbf{w}} - \mathbf{w}\|^2] \geq \mathbb{E}[\|\mathbf{w}^* - \mathbf{w}\|^2].$$

これは \mathbf{w}^* は真の値 \mathbf{w} からの距離が小さいことを意味する。

最尤推定法

仮定 3 では次の式に従う。

$$\mathbf{y} = \mathbf{X}\mathbf{w} + \boldsymbol{\epsilon}, \quad \boldsymbol{\epsilon} \sim N(\mathbf{0}, \sigma^2 \mathbf{I})$$

このとき、

$$\mathbf{y} \sim N(\mathbf{X}\mathbf{w}, \sigma^2 \mathbf{I})$$

尤度関数

$$l(\mathbf{y}; \mathbf{w}) = \frac{1}{(2\pi\sigma^2)^{D/2}} \exp \left\{ -\frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{X}\mathbf{w})^T (\mathbf{y} - \mathbf{X}\mathbf{w}) \right\}$$

対数尤度関数

$$\log l(\mathbf{y}; \mathbf{w}) = -\frac{D}{2} \log 2\pi\sigma^2 - \frac{1}{2\sigma^2} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|^2$$

対数尤度関数の最大化は線形回帰と等価

$$\max_{\mathbf{w}} \log l(\mathbf{y}; \mathbf{w})$$

つまり、線形回帰は最尤推定法とみなせる

決定係数

決定係数 R^2

- 線形回帰のデータへの当てはまりの良さの尺度
- $0 \leq R^2 \leq 1$ で、大きいほど当てはまりがよい。
- R^2 は寄与率とも言う。また、その平方根を重相関係数と言う。

定義：

- $\mathbb{V}[y]$ ： データ $\{y_d\}_{d \in [D]}$ の標本分散
- $\mathbb{V}[\hat{y}]$ ： 予測値 $\{\hat{y}_d\}_{d \in [D]}$ の標本分散
- E^2 ： 平均 2 乗予測誤差 $\frac{1}{D} \sum_{d \in [D]} (y_d - \hat{y}_d)^2$
- $\rho(y, \hat{y})$ ： $\{(y, \hat{y}_d)\}_{d \in [D]}$ の標本相関係数

決定係数

$$R^2 := \frac{\mathbb{V}[\hat{y}]}{\mathbb{V}[y]} = 1 - \frac{E^2}{\mathbb{V}[y]} = \rho(y, \hat{y})^2$$

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

$$\hat{\mathbf{y}} := \mathbf{X} \mathbf{w}^*$$

$$\boldsymbol{\epsilon} := \mathbf{y} - \hat{\mathbf{y}}$$

まず、

$$\mathbf{X}^T \boldsymbol{\epsilon} = \mathbf{X}^T (\mathbf{y} - \hat{\mathbf{y}}) = \mathbf{X}^T \mathbf{y} - \mathbf{X}^T \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} = \mathbf{0}$$

が成り立つ。このとき、

$$\hat{\mathbf{y}}^T \boldsymbol{\epsilon} = (\mathbf{X} \mathbf{w}^*)^T \boldsymbol{\epsilon} = \mathbf{w}^{*T} \mathbf{X}^T \boldsymbol{\epsilon} = 0 \quad (1)$$

である。また、 \mathbf{X}^T の1行目は $\mathbf{e}^T := (1, 1, \dots, 1)$ なので、 $\mathbf{X}^T \boldsymbol{\epsilon} = \mathbf{0}$ の一番目は

$$\mathbf{e}^T \boldsymbol{\epsilon} = 0 \quad (2)$$

となる。このとき、次の計算より y と \hat{y} の標本平均は等しくなる。

$$\mathbb{E}[y] = \frac{1}{D} \mathbf{e}^T \mathbf{y} = \frac{1}{D} \mathbf{e}^T (\hat{\mathbf{y}} + \boldsymbol{\epsilon}) = \frac{1}{D} \mathbf{e}^T \hat{\mathbf{y}} = \mathbb{E}[\hat{y}]$$

証明 (続き)

$$\begin{aligned}\mathbb{V}[y] &= \frac{1}{D}(\mathbf{y} - \mathbb{E}[y]\mathbf{e})^T(\mathbf{y} - \mathbb{E}[y]\mathbf{e}) \\ &= \frac{1}{D}(\hat{\mathbf{y}} + \boldsymbol{\epsilon} - \mathbb{E}[\hat{\mathbf{y}}]\mathbf{e})^T(\hat{\mathbf{y}} + \boldsymbol{\epsilon} - \mathbb{E}[\hat{\mathbf{y}}]\mathbf{e}) \\ &= \frac{1}{D}(\hat{\mathbf{y}} - \mathbb{E}[\hat{\mathbf{y}}]\mathbf{e})^T(\hat{\mathbf{y}} - \mathbb{E}[\hat{\mathbf{y}}]\mathbf{e}) + \frac{2}{D}(\hat{\mathbf{y}} - \mathbb{E}[\hat{\mathbf{y}}]\mathbf{e})^T\boldsymbol{\epsilon} + \frac{1}{D}\boldsymbol{\epsilon}^T\boldsymbol{\epsilon} \\ &= \mathbb{V}[\hat{\mathbf{y}}] + E^2 \quad (\because (1)(2))\end{aligned}$$

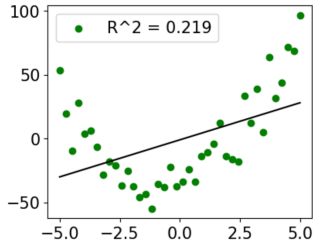
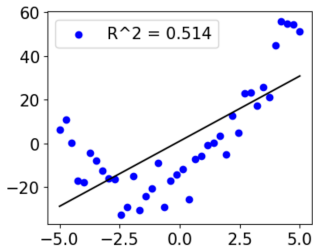
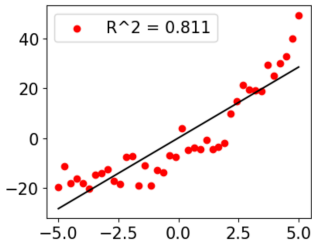
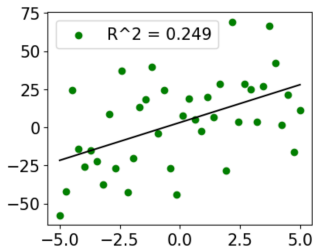
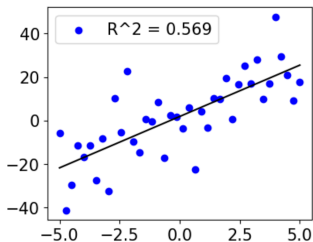
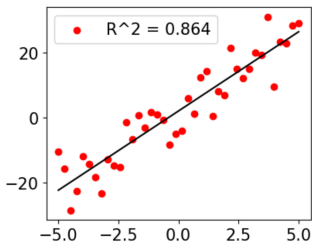
よって、
$$\frac{\mathbb{V}[\hat{\mathbf{y}}]}{\mathbb{V}[y]} = \frac{\mathbb{V}[y] - E^2}{\mathbb{V}[y]} = 1 - \frac{E^2}{\mathbb{V}[y]}$$

$$\begin{aligned}\text{Cov}[y, \hat{\mathbf{y}}] &:= \frac{1}{D}(\mathbf{y} - \mathbb{E}[y]\mathbf{e})^T(\hat{\mathbf{y}} - \mathbb{E}[\hat{\mathbf{y}}]\mathbf{e}) \\ &= \frac{1}{D}(\hat{\mathbf{y}} + \boldsymbol{\epsilon} - \mathbb{E}[\hat{\mathbf{y}}]\mathbf{e})^T(\hat{\mathbf{y}} - \mathbb{E}[\hat{\mathbf{y}}]\mathbf{e}) \\ &= \frac{1}{D}(\hat{\mathbf{y}} - \mathbb{E}[\hat{\mathbf{y}}]\mathbf{e})^T(\hat{\mathbf{y}} - \mathbb{E}[\hat{\mathbf{y}}]\mathbf{e}) + \frac{1}{D}\boldsymbol{\epsilon}^T(\hat{\mathbf{y}} - \mathbb{E}[\hat{\mathbf{y}}]\mathbf{e}) \\ &= \mathbb{V}[\hat{\mathbf{y}}] \quad (\because (1)(2))\end{aligned}$$

よって、
$$\rho(y, \hat{\mathbf{y}})^2 := \frac{\text{Cov}[y, \hat{\mathbf{y}}]^2}{\mathbb{V}[y]\mathbb{V}[\hat{\mathbf{y}}]} = \frac{\mathbb{V}[\hat{\mathbf{y}}]^2}{\mathbb{V}[y]\mathbb{V}[\hat{\mathbf{y}}]} = \frac{\mathbb{V}[\hat{\mathbf{y}}]}{\mathbb{V}[y]}$$

決定係数

上段は一次関数、下段は二次関数の予測。 R^2 だけでは全てを判断できない。



正則化項の導入

線形回帰を行う上での問題

- 多重共線性：

$\text{rank}(\mathbf{X}) < n$ のとき、 \mathbf{w}^* が計算できない。

- 数値的不安定：

$\mathbf{X}^T \mathbf{X}$ に 0 に近い固有値があると、 \mathbf{w}^* の分散が大きくなる。

$$\mathbb{V}[\mathbf{w}^*] = \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1}$$

- 過学習：

テストデータで予測精度が悪化することもある

これらの問題を解決するために、正則化項を加えることが有効

正則化付き学習

正則化付き学習

- Ridge 回帰

$$\min_{\mathbf{w}} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|^2 + \lambda \|\mathbf{w}\|^2$$

- LASSO 回帰

$$\min_{\mathbf{w}} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|^2 + \lambda \|\mathbf{w}\|_1$$

- Elastic Net 回帰

$$\min_{\mathbf{w}} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|^2 + \lambda_1 \|\mathbf{w}\|_1 + \lambda_2 \|\mathbf{w}\|^2$$

$$\|\mathbf{w}\|_1 := \sum_{i=1}^n |w_i|$$

- ハイパーパラメータ $\lambda \geq 0$ は正則化の強さを表す定数
- LASSO 回帰は最適解 \mathbf{w}^* の各要素が 0 になりやすい。このため、予測器の解釈がし易くなる。
- Elastic net 回帰は、Ridge 回帰と LASSO 回帰の中間の性質がある。

学習アルゴリズム

Ridge 回帰： $\min_{\mathbf{w}} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|^2 + \lambda\|\mathbf{w}\|^2$

解析解が得られ、 $\mathbf{w}^* = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y}$.

LASSO 回帰： $\min_{\mathbf{w}} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|^2 + \lambda\|\mathbf{w}\|_1$

この問題は $\mathbf{w}, \mathbf{s} \in \mathbb{R}^n$ を変数とする次のような 2 次計画問題となる。

$$\begin{aligned} \min_{\mathbf{w}, \mathbf{s}} \quad & \mathbf{w}^T \mathbf{X}^T \mathbf{X} \mathbf{w} - 2\mathbf{y}^T \mathbf{X} \mathbf{w} + \lambda \mathbf{e}^T \mathbf{s} \\ \text{s.t.} \quad & -\mathbf{s} \leq \mathbf{w} \leq \mathbf{s}. \end{aligned}$$

Elastic Net 回帰： $\min_{\mathbf{w}} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|^2 + \lambda_1\|\mathbf{w}\|_1 + \lambda_2\|\mathbf{w}\|^2$

この問題は $\mathbf{w}, \mathbf{s} \in \mathbb{R}^n$ を変数とする次のような 2 次計画問題となる。

$$\begin{aligned} \min_{\mathbf{w}, \mathbf{s}} \quad & \mathbf{w}^T (\mathbf{X}^T \mathbf{X} + \lambda_2 \mathbf{I}) \mathbf{w} - 2\mathbf{y}^T \mathbf{X} \mathbf{w} + \lambda_1 \mathbf{e}^T \mathbf{s} \\ \text{s.t.} \quad & -\mathbf{s} \leq \mathbf{w} \leq \mathbf{s}. \end{aligned}$$

2 次計画問題の汎用解法で解くことも可能だが、専用解法の方が高速

LASSO の学習アルゴリズム

反復解法で解く。 k 回目の反復点を \mathbf{w}_k とする。

目的関数 $f(\mathbf{w}) := \|\mathbf{y} - \mathbf{X}\mathbf{w}\|^2 + \lambda\|\mathbf{w}\|_1$ に対し、次を満たす関数 $g_k(x)$ を考える。

- $f(\mathbf{w}_k) = g_k(\mathbf{w}_k)$
- $\nabla f(\mathbf{w}_k) = \nabla g_k(\mathbf{w}_k)$
- $f(\mathbf{w}) \leq g_k(\mathbf{w}) \quad (\forall \mathbf{w} \in \mathbb{R}^n)$

そして、 $f(\mathbf{w})$ を最小化する代わりに、次のように反復点を更新する。

$$\mathbf{w}_{k+1} := \underset{\mathbf{w}}{\operatorname{argmin}} g_k(\mathbf{w})$$

- $f(\mathbf{w}_k)$ ($k = 1, 2, \dots$) は非増加 (減少 or 変わらない)

$$\begin{aligned} f(\mathbf{w}_k) &= g_k(\mathbf{w}_k) \\ &\geq g(\mathbf{w}_{k+1}) \\ &\geq f(\mathbf{w}_{k+1}) \end{aligned}$$

LASSO の学習アルゴリズム

- $f(\mathbf{w}_k)$ が変わらない時、最適解が得られる

$$f(\mathbf{w}_{k+1}) = f(\mathbf{w}_k)$$

$\implies \mathbf{w}_k$ が $g_k(\mathbf{w})$ の最適解

$$\implies \nabla g_k(\mathbf{w}_k) = \mathbf{0}$$

$$\implies \nabla f(\mathbf{w}_k) = \mathbf{0}$$

$\implies \mathbf{w}_k$ は $f(\mathbf{w})$ の最適解

$g_k(\mathbf{x})$ には次の性質も必要

- $\min_{\mathbf{w}} g_k(\mathbf{w})$ が効率よく計算できる。
1 反復辺りの計算時間を減らすため
- $g_k(\mathbf{w})$ は $f(\mathbf{w})$ のタイトなバウンドである。
最適解を得るまでの反復回数を減らすため

$g_k(\mathbf{w})$ の構築

$$\begin{aligned} f(\mathbf{w}) &:= \|\mathbf{y} - \mathbf{X}\mathbf{w}\|^2 + \lambda\|\mathbf{w}\|_1 \\ &= \|(\mathbf{y} - \mathbf{X}\mathbf{w}_k) - \mathbf{X}(\mathbf{w} - \mathbf{w}_k)\| + \lambda\|\mathbf{w}\|_1 \\ &= (\mathbf{w} - \mathbf{w}_k)^T \mathbf{X}^T \mathbf{X} (\mathbf{w} - \mathbf{w}_k) - 2(\mathbf{y} - \mathbf{X}\mathbf{w}_k)^T \mathbf{X} (\mathbf{w} - \mathbf{w}_k) \\ &\quad + (\mathbf{y} - \mathbf{X}\mathbf{w}_k)^T (\mathbf{y} - \mathbf{X}\mathbf{w}_k) + \lambda\|\mathbf{w}\|_1 \end{aligned}$$

と変形できる。ここで、次のように $g_k(\mathbf{w})$ を定義する (ρ は $\mathbf{X}^T \mathbf{X}$ の最大固有値)。

$$\begin{aligned} g_k(\mathbf{w}) &:= \rho(\mathbf{w} - \mathbf{w}_k)^T (\mathbf{w} - \mathbf{w}_k) - 2(\mathbf{y} - \mathbf{X}\mathbf{w}_k)^T \mathbf{X} (\mathbf{w} - \mathbf{w}_k) \\ &\quad + (\mathbf{y} - \mathbf{X}\mathbf{w}_k)^T (\mathbf{y} - \mathbf{X}\mathbf{w}_k) + \lambda\|\mathbf{w}\|_1 \end{aligned}$$

このとき、次の関係が成り立つ。

- $f(\mathbf{w}_k) = g_k(\mathbf{w}_k)$
- $\nabla f(\mathbf{w}_k) = \nabla g_k(\mathbf{w}_k)$
- $f(\mathbf{w}) \leq g_k(\mathbf{w}) \quad (\forall \mathbf{w} \in \mathbb{R}^n)$

証明

$$\begin{aligned}f(\mathbf{w}_k) &= g_k(\mathbf{w}_k) \\ &= -2(\mathbf{y} - \mathbf{X}\mathbf{w}_k)^T \mathbf{X}(\mathbf{w} - \mathbf{w}_k) + (\mathbf{y} - \mathbf{X}\mathbf{w}_k)^T (\mathbf{y} - \mathbf{X}\mathbf{w}_k) + \lambda \|\mathbf{w}\|_1\end{aligned}$$

$$\begin{aligned}\nabla f(\mathbf{w}_k) &= -2\mathbf{X}^T (\mathbf{y} - \mathbf{X}\mathbf{w}_k) + \lambda \nabla \|\mathbf{w}\|_1 \\ &= \nabla g_k(\mathbf{w}_k)\end{aligned}$$

$\rho \mathbf{I} \succeq \mathbf{X}^T \mathbf{X}$ より、

$$\begin{aligned}\rho(\mathbf{w} - \mathbf{w}_k)^T (\mathbf{w} - \mathbf{w}_k) - (\mathbf{w} - \mathbf{w}_k)^T \mathbf{X}^T \mathbf{X} (\mathbf{w} - \mathbf{w}_k) \\ &= (\mathbf{w} - \mathbf{w}_k)^T (\rho \mathbf{I} - \mathbf{X}^T \mathbf{X}) (\mathbf{w} - \mathbf{w}_k) \\ &\geq 0\end{aligned}$$

$g_k(\mathbf{w})$ の変形

$\mathbf{c} := -\mathbf{X}^T(\mathbf{y} - \mathbf{X}\mathbf{w}_k)$ とすると、次のように変形できる。

$$\begin{aligned}g_k(\mathbf{w}) &= \rho(\mathbf{w} - \mathbf{w}_k)^T(\mathbf{w} - \mathbf{w}_k) + \mathbf{c}^T(\mathbf{w} - \mathbf{w}_k) + \lambda\|\mathbf{w}\|_1 + \text{定数} \\&= \sum_{i \in [n]} \rho(w_i - [\mathbf{w}_k]_i)^2 + c_i(w_i - [\mathbf{w}_k]_i) + \lambda|w_i| + \text{定数} \\&= \sum_{i \in [n]} \rho\left(w_i - [\mathbf{w}_k]_i + \frac{1}{\rho}c_i\right)^2 + \lambda|w_i| + \text{定数} \\&= \sum_{i \in [n]} \rho\left(w_i - \left([\mathbf{w}_k]_i - \frac{1}{\rho}c_i\right)\right)^2 + \lambda|w_i| + \text{定数}\end{aligned}$$

要素ごとに分割できるため、簡単に最適解が得られる。

反復点の更新

$$\min_{w_i} \rho \left(w_i - \left([\mathbf{w}_k]_i - \frac{1}{\rho} c_i \right) \right)^2 + \lambda |w_i|$$

の最適解は $w_i^* = S_{\lambda/2\rho} \left([\mathbf{w}_k]_i - \frac{1}{\rho} c_i \right)$ となる (次スライド)。

ただし、 $S_\alpha(x)$ はソフト閾値関数で

$$S_\alpha(x) := \begin{cases} x - \alpha & (\alpha \leq x) \\ 0 & (-\alpha < x < \alpha) \\ x + \alpha & (x \leq -\alpha) \end{cases}$$

よって、次の式によって反復点を更新する。

$$\forall i \in [n], \quad [\mathbf{w}_{k+1}]_i := S_{\lambda/2\rho} \left([\mathbf{w}_k]_i - \frac{1}{\rho} c_i \right)$$

この解法を ISTA (iterative shrinkage threshold algorithm) と呼ぶ。

これに改良を加えた FISTA という解法もある。

証明

$y \in \mathbb{R}$, $\alpha \geq 0$ を定数としたとき、関数 $f(x) := \frac{1}{2}(x - y)^2 + \alpha|x|$ の最小解は、

$$x^* = S_\alpha(y) := \begin{cases} y - \alpha & (\alpha \leq y) \\ 0 & (-\alpha < y < \alpha) \\ y + \alpha & (y \leq -\alpha) \end{cases} \quad \text{となる。}$$

証明：

- $-\alpha < y < \alpha$ の場合、

$$\begin{aligned} f(x) &= \frac{1}{2}(x - y)^2 + \alpha|x| \geq \frac{1}{2}(x^2 - 2|x||y| + y^2) + \alpha|x| \\ &\geq \frac{1}{2}(x^2 - 2\alpha|x| + y^2) + \alpha|x| = \frac{1}{2}(x^2 + y^2) \geq \frac{1}{2}y^2 = f(0) \end{aligned}$$

最後の不等号は $x = 0$ のときのみ等号が成立する。

- $\alpha \leq y$ の場合、 $\tilde{f}(x) := \frac{1}{2}(x - y)^2 + \alpha x$ とすると、 $\forall x, f(x) \geq \tilde{f}(x)$ より、

$$f(y - \alpha) \geq \min_x f(x) \geq \min_x \tilde{f}(x) = \tilde{f}(y - \alpha) = f(y - \alpha)$$

$\min_x \tilde{f}(x)$ は $x^* = y - \alpha$ が唯一の最小点。

- $y \leq -\alpha$ の場合も、 $\tilde{f}(x) := \frac{1}{2}(x - y)^2 - \alpha x$ とすれば、上の証明と同様

Elastic Net 回帰の学習アルゴリズム

$$f(\mathbf{w}) := \|\mathbf{y} - \mathbf{X}\mathbf{w}\|^2 + \lambda_1 \|\mathbf{w}\|_1 + \lambda_2 \|\mathbf{w}\|^2$$

次の関数を考える。

$$g_k(\mathbf{w}) := (\rho + \lambda_2)(\mathbf{w} - \mathbf{w}_k)^T(\mathbf{w} - \mathbf{w}_k) - 2(\mathbf{y} - \mathbf{X}\mathbf{w}_k)^T \mathbf{X}(\mathbf{w} - \mathbf{w}_k) \\ + (\mathbf{y} - \mathbf{X}\mathbf{w}_k)^T(\mathbf{y} - \mathbf{X}\mathbf{w}_k) + \lambda_1 \|\mathbf{w}\|_1$$

後は LASSO と同じように考えることができる。

最適解は次の計算による反復で求めることができる。

$$\forall i, \quad [\mathbf{w}_{k+1}]_i := S_{\lambda_1/2(\rho+\lambda_2)} \left([\mathbf{w}_k]_i - \frac{1}{\rho + \lambda_2} c_i \right)$$

確率モデルに基づいた解析

ベイズの定理を使うと、データ $\mathcal{D} = \{(\mathbf{x}_d, y_d)\}_{d \in [D]}$ のもとでのパラメタ \mathbf{w} の確率は、次のようになる。

$$P(\mathbf{w}|\mathcal{D}) = \frac{P(\mathcal{D}|\mathbf{w})P(\mathbf{w})}{P(\mathcal{D})}$$

- $P(\mathbf{w})$: パラメタ \mathbf{w} の事前確率
- $P(\mathbf{w}|\mathcal{D})$: データ \mathcal{D} が与えられた状況でのパラメタ \mathbf{w} の事後確率

MAP(maximum a posteriori) 推定

事後確率 $P(\mathbf{w}|\mathcal{D})$ が最大となるパラメタ \mathbf{w} を採用する

$$\max_{\mathbf{w}} \log P(\mathbf{w}|\mathcal{D}) \iff \max_{\mathbf{w}} \log P(\mathcal{D}|\mathbf{w}) + \log P(\mathbf{w})$$

Ridge 回帰と MAP 推定

仮定

- パラメタの事前分布： $\mathbf{w} \sim N(\mathbf{0}, \tau^2 \mathbf{I})$
- 仮定 3： $\mathbf{y} = \mathbf{X}\mathbf{w} + \epsilon$, $\epsilon \sim N(\mathbf{0}, \sigma^2 \mathbf{I})$

$$\begin{aligned} \log P(\mathcal{D}|\mathbf{w}) + \log P(\mathbf{w}) &= -\frac{D}{2} \log 2\pi\sigma^2 - \frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{X}\mathbf{w})^T (\mathbf{y} - \mathbf{X}\mathbf{w}) \\ &\quad - \frac{n}{2} \log 2\pi\tau^2 - \frac{1}{2\tau^2} \mathbf{w}^T \mathbf{w} \\ &= -\frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{X}\mathbf{w})^T (\mathbf{y} - \mathbf{X}\mathbf{w}) - \frac{1}{2\tau^2} \mathbf{w}^T \mathbf{w} + \text{定数} \end{aligned}$$

この関数の最大化は、 λ を適当に定めた $\|\mathbf{y} - \mathbf{X}\mathbf{w}\|^2 + \lambda\|\mathbf{w}\|^2$ の最小化と等価
つまり、Ridge 回帰は上記の仮定での MAP 推定とみなせる

LASSO 回帰と MAP 推定

別の仮定

- パラメタの事前分布： w の各要素が平均 0 のラプラス分布に従う
※ 平均 μ , 分散 $2b^2$ であるラプラス分布の密度関数

$$f(x; \mu, b) = \frac{1}{2b} \exp \left\{ -\frac{|x - \mu|}{b} \right\}$$

- 仮定 3： $y = Xw + \epsilon$, $\epsilon \sim N(\mathbf{0}, \sigma^2 I)$

この仮定のもとでの、事後確率の最大化は LASSO 回帰における学習と等価
つまり、LASSO 回帰は MAP 推定とみなせる。

Elastic Net 回帰と等価な事前分布を作ることも可能だが、名前が付いているようなありふれた分布ではない

判別・分類

ここでは2クラスの判別・分類を考える。

$$\text{データ: } \{(\mathbf{x}_d, y_d)\}_{d \in [D]}, \quad \mathbf{x}_d \in \mathbb{R}^n, y_d \in \{0, 1\}$$

ロジスティック回帰

- ロジスティック回帰では確率的判別モデルを考える

$y = 1$ の確率を計算する。 $f: \mathbb{R}^n \rightarrow [0, 1]$

- 確率（実数）を予測するので回帰
- 新しいデータ \mathbf{x} に対し、予測した確率が閾値（例えば 0.5）以上のとき $y = 1$ 、そうでないとき $y = 0$ と予測
- 多クラス判別に拡張した多項ロジスティック回帰というものもある

ロジスティック回帰

予測器

線形関数 $w^T x$ とシグモイド関数を使った回帰。パラメタは w

$$\hat{y} = \sigma(w^T x)$$

$\hat{y} \in [0, 1]$: 予測値 (確率)

シグモイド関数:

$$\sigma(a) := \frac{1}{1 + \exp(-a)} = \frac{\exp a}{1 + \exp a}$$

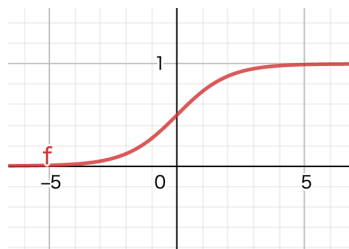
$$0 < \sigma(a) < 1$$

$$\lim_{a \rightarrow -\infty} \sigma(a) = 0, \quad \lim_{a \rightarrow \infty} \sigma(a) = 1$$

$$a_1 < a_2 \implies \sigma(a_1) < \sigma(a_2)$$

$$\sigma(-a) = 1 - \sigma(a)$$

$$\sigma'(a) = \sigma(a)(1 - \sigma(a))$$



判別面

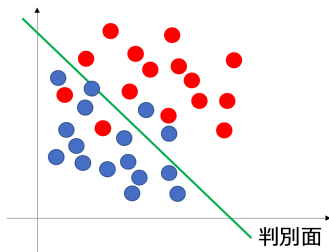
\mathbf{x} に対する判別 (α は閾値)

$$\text{判別結果} = \begin{cases} 1 & (\sigma(\mathbf{w}^T \mathbf{x}) \geq \alpha) \\ 0 & (\sigma(\mathbf{w}^T \mathbf{x}) \leq \alpha) \end{cases}$$

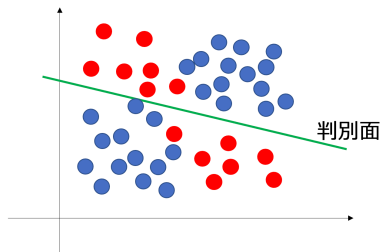
$\sigma(a)$ は単調増加関数なので、判別面は超平面となる。

$$C_1 := \{\mathbf{x} \mid \mathbf{w}^T \mathbf{x} \geq \sigma^{-1}(\alpha)\}, \quad C_0 := \{\mathbf{x} \mid \mathbf{w}^T \mathbf{x} < \sigma^{-1}(\alpha)\}$$

線形分離可能



線形分離不可能



学習

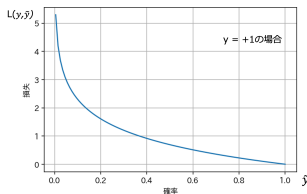
誤差関数

(2クラス判別における) クロスエントロピー誤差

$$L(y, \hat{y}) := -y \log \hat{y} - (1 - y) \log(1 - \hat{y})$$

$y \in \{0, 1\}$: 目標, $\hat{y} \in [0, 1]$: 予測値

$y = 1$ のとき



ロジスティック回帰における学習

$$\min_{\mathbf{w}} \sum_{d \in [D]} L(y_d, \sigma(\mathbf{w}^T \mathbf{x}_d))$$

確率モデルに基づいた解析 1

確率 \hat{y} で $y = 1$ 、それ以外（確率 $1 - \hat{y}$ ）で $y = 0$ となるベルヌーイ試行で、データが生成されたとする。

$$\text{ただし } \hat{y} = \sigma(\mathbf{w}^T \mathbf{x})$$

尤度関数

$$l(\mathbf{y}; \mathbf{w}) := \prod_{d \in [D]} (\hat{y}_d)^{y_d} \times (1 - \hat{y}_d)^{1 - y_d}$$

対数尤度関数

$$\begin{aligned} \log l(\mathbf{y}; \mathbf{w}) &:= \sum_{d \in [D]} \{y_d \log \hat{y}_d + (1 - y_d) \log(1 - \hat{y}_d)\} \\ &= - \sum_{d \in [D]} L(y_d, \hat{y}_d) \end{aligned}$$

クロスエントロピー誤差を用いた学習は最尤推定に相当する

$$\max_{\mathbf{w}} \log l(\mathbf{y}; \mathbf{w}) \iff \min_{\mathbf{w}} \sum_{d \in [D]} L(y_d, \sigma(\mathbf{w}^T \mathbf{x}_d))$$

確率モデルに基づいた解析 2

データ y_d , 予測値 \hat{y}_d とする

経験分布 P

$$p(y) := \begin{cases} 1 & (y = y_d) \\ 0 & (y \neq y_d) \end{cases}$$

予測した分布 Q

$$q(y) := \begin{cases} \hat{y}_d & (y = 1) \\ 1 - \hat{y}_d & (y = 0) \end{cases}$$

とすると、

$$\begin{aligned} KL(P||Q) &:= \sum_{x_i \in X} p(x_i) \log \frac{p(x_i)}{q(x_i)} \\ &= -y_d \log \hat{y}_d - (1 - y_d) \log(1 - \hat{y}_d) \\ &= L(y_d, \hat{y}_d) \end{aligned}$$

クロスエントロピー誤差は、経験分布と予測分布の距離 (KL ダイバージェンス)

別の見方

予測器

線形関数 $w^T x$ を使った回帰。パラメタは w

$$\hat{y} := w^T x$$

誤差関数

ロジスティック損失： $y \in \{1, -1\}$ とする

$$L(y, \hat{y}) := \log(1 + \exp(-y\hat{y}))$$

この予測器と誤差関数を用いた学習は次のようになる。

学習

$$\min_w \sum_{d \in [D]} L(y_d, w^T x_d)$$

これは、ロジスティック回帰における学習と一致する。

等価性の証明

$$\hat{y} := \mathbf{w}^T \mathbf{x}, \quad L(y, \hat{y}) := \log(1 + \exp(-y\hat{y}))$$

より、次の関係が成り立つ。

$$L(y_d, \mathbf{w}^T \mathbf{x}_d) = \begin{cases} \log \{1 + \exp(+\mathbf{w}^T \mathbf{x}_d)\} & (y_d = +1) \\ \log \{1 + \exp(-\mathbf{w}^T \mathbf{x}_d)\} & (y_d = -1) \end{cases}$$

一方、

$$\hat{y} := \sigma(\mathbf{w}^T \mathbf{x}_d), \quad \sigma(a) := \frac{1}{1 + \exp(-a)},$$
$$L(y, \hat{y}) := -y \log \hat{y} - (1 - y) \log(1 - \hat{y})$$

より、次の関係が成り立つ。

$$L(y_d, \sigma(\mathbf{w}^T \mathbf{x}_d)) = -y_d \log \frac{1}{1 + \exp(-\mathbf{w}^T \mathbf{x}_d)} - (1 - y_d) \log \frac{1}{1 + \exp(\mathbf{w}^T \mathbf{x}_d)}$$
$$= \begin{cases} \log \{1 + \exp(+\mathbf{w}^T \mathbf{x}_d)\} & (y_d = 1) \\ \log \{1 + \exp(-\mathbf{w}^T \mathbf{x}_d)\} & (y_d = 0) \end{cases}$$

よって、両者は等しい。

学習アルゴリズム

目的関数を $F(\mathbf{w})$ とする。このとき、

$$\begin{aligned}\nabla F(\mathbf{w}) &= \mathbf{X}^T (\hat{\mathbf{y}} - \mathbf{y}) \\ \nabla^2 F(\mathbf{w}) &= \mathbf{X}^T \mathbf{D} \mathbf{X}\end{aligned}$$

ここで、 $\hat{y}_d = \sigma(\mathbf{w}^T \mathbf{x}_d)$, $D_{dd} = \hat{y}_d(1 - \hat{y}_d) > 0$ ($d \in [D]$)

式変形は次のスライド

性質など

- $\nabla^2 F(\mathbf{w})$ は半正定値行列より、 $F(\mathbf{w})$ は凸関数
- $\nabla F(\mathbf{w}) = \mathbf{0}$ が最適解の必要十分条件だが、解析的に解くことはできない
- $F(\mathbf{w})$ の最小点を Newton 法で求める

$$\hat{y}_d = \sigma(\mathbf{w}^T \mathbf{x}_d)$$

$$\frac{\partial \hat{y}_d}{\partial \mathbf{w}} = \hat{y}_d(1 - \hat{y}_d)\mathbf{x}_d^T$$

$$F(\mathbf{w}) = \sum_{d \in [D]} \{-y_d \log \hat{y}_d - (1 - y_d) \log(1 - \hat{y}_d)\}$$

$$\nabla F(\mathbf{w}) = \left(\frac{\partial F(\mathbf{x})}{\partial \mathbf{w}} \right)^T = \sum_{d \in [D]} \left\{ -\frac{y_d}{\hat{y}_d} \hat{y}_d(1 - \hat{y}_d)\mathbf{x}_d + \frac{1 - y_d}{(1 - \hat{y}_d)} \hat{y}_d(1 - \hat{y}_d)\mathbf{x}_d \right\}$$

$$= \sum_{d \in [D]} (-y_d + \hat{y}_d)\mathbf{x}_d$$

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

$$\nabla^2 F(\mathbf{w}) = \sum_{d \in [D]} \hat{y}_d(1 - \hat{y}_d)\mathbf{x}_d\mathbf{x}_d^T$$

$$= \mathbf{X}^T \mathbf{D} \mathbf{X}$$

Newton 法

Newton 法の適用

$$\begin{aligned}\mathbf{w}_{k+1} &:= \mathbf{w}_k - (\nabla^2 F(\mathbf{w}_k))^{-1} \nabla F(\mathbf{w}_k) \\ &= \mathbf{w}_k - (\mathbf{X}^T \mathbf{D} \mathbf{X})^{-1} \mathbf{X}^T (\hat{\mathbf{y}} - \mathbf{y}) \\ &= (\mathbf{X}^T \mathbf{D} \mathbf{X})^{-1} \left(\mathbf{X}^T \mathbf{D} \mathbf{X} \mathbf{w}_k - \mathbf{X}^T (\hat{\mathbf{y}} - \mathbf{y}) \right) \\ &= (\mathbf{X}^T \mathbf{D} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{D} (\mathbf{X} \mathbf{w}_k - \mathbf{D}^{-1} (\hat{\mathbf{y}} - \mathbf{y}))\end{aligned}$$

$\mathbf{z} := \mathbf{X} \mathbf{w}_k - \mathbf{D}^{-1} (\hat{\mathbf{y}} - \mathbf{y})$ とおくと、

$$\mathbf{w}_{k+1} := (\mathbf{X}^T \mathbf{D} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{D} \mathbf{z}.$$

これは、内積が $(\mathbf{x} \cdot \mathbf{y}) := \mathbf{x}^T \mathbf{D} \mathbf{y}$ で定義される内積空間上での直交射影

学習アルゴリズム

学習アルゴリズム

反復重み付き最小2乗法 (iterative reweighted least squares method; IRLS)

- 1 $k := 1$ として、 \mathbf{w}_1 の初期値を定める。
- 2 $\hat{y}_d := \sigma(\mathbf{w}_k^T \mathbf{x}_d) \quad d \in [D]$
 $D_{dd} := \hat{y}_d(1 - \hat{y}_d) \quad d \in [D]$
- 3 $\mathbf{z} := \mathbf{X}\mathbf{w}_k - \mathbf{D}^{-1}(\hat{\mathbf{y}} - \mathbf{y})$
- 4 $\mathbf{w}_{k+1} := (\mathbf{X}^T \mathbf{D} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{D} \mathbf{z}$
- 5 終了条件を満たしていれば終了。
そうでなければ、 $k := k + 1$ として (2) に戻る