

Table of Contents

- 1 Breast Cancer Wisconsin (Diagnostic) Data Set
 - 1.1 Attribute Information:
 - 1.2 分類器
 - 1.3 仮説クラス
- 2 最急降下法
 - 2.1 print_w
 - 2.2 データの読み込みと初期化
 - 2.3 最急降下法によるw探索(steepest descent)
- 3 結果
- 4 QR decomposition

Breast Cancer Wisconsin (Diagnostic) Data Set

[https://archive.ics.uci.edu/ml/datasets/breast+cancer+wisconsin+\(diagnostic\)](https://archive.ics.uci.edu/ml/datasets/breast+cancer+wisconsin+(diagnostic))

Attribute Information:

1. ID number
2. Diagnosis (M = malignant:-1, B = benign:1) M:悪性, B:良性
3. 3-32

Ten real-valued features are computed for each cell nucleus:

- 半径radius (mean of distances from center to points on the perimeter)
- テクスチャtexture (standard deviation of gray-scale values)
- 境界の長さperimeter
- 面積area
- なめらかさsmoothness (local variation in radius lengths)
- コンパクトさcompactness (perimeter^2 / area - 1.0)
- くぼみ度合いconcavity (severity of concave portions of the contour)
- くぼみの数concave points (number of concave portions of the contour)
- 対称性symmetry
- フラクタル次元fractal dimension ("coastline approximation" - 1)

のそれぞれのmean, stderr, worst数値を保持している.

<http://people.idsia.ch/~juergen/deeplearningwinsMICCAIgrandchallenge.html>

分類器

用意された訓練(training)データには, \mathbf{A} に上に記した特徴量が, \mathbf{b} に悪性(-1)か良性(1)かを示す数値が入っている.

与えられた特徴ベクトル \mathbf{y} に対し, 細胞組織が悪性か良性かを分類する関数 $C(\mathbf{y})$ を選び出すプログラムを作成しよう.

仮説クラス

分類器は可能な分類器の集合(仮説クラス)から選ばれる. この場合, 仮説クラスとは特徴ベクトルの空間 \mathbb{R}^D から \mathbb{R} への線形関数 $h(\cdot)$ である. すると分類器は次のような関数として定義される.

$$C(\mathbf{y}) = \begin{cases} +1 & \text{when } h(\mathbf{y}) \geq 0 \\ -1 & \text{when } h(\mathbf{y}) < 0 \end{cases}$$

$h(\mathbf{y}) = \mathbf{w} \cdot \mathbf{y}$

各線形関数 $h: \mathbb{R}^D \rightarrow \mathbb{R}$ に対して, 次のような D ベクトル \mathbf{w} が存在する.

したがって, そのような線形関数を選ぶことは, 結局 D ベクトル \mathbf{w} を選ぶことに等しい. 特に, \mathbf{w} を選ぶことは, 仮説クラス h を選ぶことと等価なので, \mathbf{w} を仮説ベクトルと呼ぶ.

問題を単純化すると, 分類器を単なるベクトルとみなして, データとの掛け算で予測がつきます. 本来は予測を-1,1とかに投影しないといけないんですが, 単純化のためにそのままの値を使います. 問題は どうやってこの仮説ベクトル \mathbf{w} の各要素の値を決定するか? ですよ.

最急降下法

損失関数に

$$L(\mathbf{w}) = \sum_{i=1}^n (\mathbf{A}_i \cdot \mathbf{w} - b_i)^2$$

最小二乗法

を選ぶと, ベクトル \mathbf{w} の j 偏微分は,

$$\begin{aligned} \frac{\partial L}{\partial w_j} &= \sum_{i=1}^n \frac{\partial}{\partial w_j} (\mathbf{A}_i \cdot \mathbf{w} - b_i)^2 \\ &= \sum_{i=1}^n 2(\mathbf{A}_i \cdot \mathbf{w} - b_i) \mathbf{A}_{ij} \end{aligned}$$

連立方程式

となる. ここで, \mathbf{A}_{ij} は \mathbf{A}_i の j 番目の要素を意味する. この偏微分 $\frac{\partial L}{\partial w_j}$ を \mathbf{w}_j の勾配(slope)として, $L(\mathbf{w})$ の極小値(local minimum)を求める.

このような探索方法を最急降下法(steepest descent method)と呼ぶ.

print_w

出てきた \mathbf{w} の j 要素をきれいに表示する関数を用意しておきます.

```
In [ ]: def print_w(w):
```