

▼ 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](#)
- 3 [code\(python\)](#)
 - 3.1 [print w](#)
 - 3.2 [データの読み込みと初期化](#)
 - 3.3 [最急降下法によるw探索\(steepest descent\)](#)
- 4 [結果](#)
- 5 [QR decomposition](#)

▼ 1 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))
([https://archive.ics.uci.edu/ml/datasets/breast+cancer+wisconsin+\(diagnostic\)](https://archive.ics.uci.edu/ml/datasets/breast+cancer+wisconsin+(diagnostic)))

1.1 Attribute Information:

1. ID number
2. Diagnosis (M = malignant, B = benign) 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 ($\text{perimeter}^2 / \text{area} - 1.0$)
- 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)

<http://people.idsia.ch/~juergen/deeplearningwinsMICCAIgrandchallenge.htm>
([http://people.idsia.ch/~juergen/deeplearningwinsMICCAIgrandchallenge.h](http://people.idsia.ch/~juergen/deeplearningwinsMICCAIgrandchallenge.htm))

1.2 分類器

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

▼ 1.3 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。

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

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 w が存在する。

$$h(y) = w \cdot y$$

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 w を選ぶことに等しい。特に、 w を選ぶことは、仮説クラス h を選ぶことと等価なので、 w を**仮説ベクトル**と呼ぶ。

単に、ベクトルの掛け算で分類器はできそう。問題は どうやってこの仮説ベクトル w の各要素の値を決定するか？ですね。

▼ 2 最急降下法

損失関数に

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

を選ぶと、ベクトル w の j 偏微分は、

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

となる。ここで、 A_{ij} は A_i の j 番目の要素です。この偏微分 $\frac{\partial L}{\partial w_j}$ を w_j の勾配(slope)として、 $L(w)$ 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。

▼ 3 code(python)

- /Users/bob/python/doing_math_with_python/numerical_calc/breast_c

▼ 3.1 print_w

出てきた 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 [1]: 1 def print_w(w):
2         params = ["radius",
3                 "texture", "perimeter", "area",
4                 "smoothness", "compactness", "concavity", "concave
5                 points",
6                 "symmetry", "fractal dimension"]
7         print("      (params)      :      ", end="")
8         print("      (mean)      (stderr)      (worst)")
9         for i, param in enumerate(params):
10            print("%18s:" % param, end="")
11            for j in range(3):
12                print("%13.9f" % w[i*3+j], end="")
13            print()
```

▼ 3.2 データの読み込みと初期化

```
In [2]: 1 import numpy as np
2         tmp = np.fromfile('./codes/train_A.data',
3                 np.float64, -1, " ")
4         A = tmp.reshape(300,30)
5         tmp = np.fromfile('./codes/train_b.data',
6                 np.float64, -1, " ")
7         b = tmp.reshape(300,1)
8         w = np.zeros(30).reshape(30,1)
9         for i in range(30):
10            w[i] = 0
```

▼ 3.3 最急降下法によるw探索(steepest descent)

In [3]:

```

1 loop, sigma = 300, 3.0*10**(-9)
2 for i in range(loop):
3     dLw = A.dot(w)-b
4     w = w -
      (dLw.transpose().dot(A)).transpose()*sigma
5
6 print_w(w)

      (params)      :      (mean)      (stderr)      (w
orst)
      radius:  0.000426997  0.000741817  0.0025
48876
      texture:  0.001687946  0.000004707  0.0000
00127
      perimeter: -0.000003968 -0.000002078  0.0000
08954
      area:  0.000003595  0.000002569  0.0000
70324
      smoothness:  0.000001139 -0.000881778  0.0000
00430
      compactness:  0.000000441  0.000000723  0.0000
00267
      concavity:  0.000001200  0.000000191  0.0004
11499
      concave points:  0.000921972  0.002395138 -0.0019
32789
      symmetry:  0.000005930 -0.000003750 -0.0000
08147
      fractal dimension: -0.000002341  0.000011565  0.0000
03523

```

▼ 4 結果

In [4]:

```

1 def show_accuracy(mA, vb, vw):
2     # M:悪性(-1), B:良性(1)
3
4     correct,safe_error,critical_error=0,0,0
5     predict = mA.dot(vw)
6     n = vb.size
7     for i in range(n):
8         if predict[i]*vb[i]>0:
9             correct += 1
10        elif (predict[i]<0 and vb[i]>0):
11            safe_error += 1
12        elif (predict[i]>0 and vb[i]<0):
13            critical_error += 1
14    print("      correct: %4d/%4d" % (correct,n))
15    print("      safe error: %4d" % safe_error)
16    print("critical error: %4d" % critical_error)
17

```

In [5]:

```

1 show_accuracy(A, b, w)

      correct:  274/ 300
      safe error:    5
      critical error:  21

```

In [6]:

```

1 tmp = np.fromfile('./codes/validate_A.data',
2 np.float64, -1, " ")
3 A = tmp.reshape(260,30)
4 tmp = np.fromfile('./codes/validate_b.data',
5 np.float64, -1, " ")
6 b = tmp.reshape(260,1)
7
8 show_accuracy(A, b, w)

      correct:  240/ 260
      safe error:    10
      critical error:   10

```

▼ 5 QR decomposition

QR分解を使うとより簡単に最小値を求めることができる。行列 A は正方形でないので、逆行列をもとめることができない。しかし、その場合でも $\|A \cdot w - b\|^2$ を最小にする w を求めることができる。

QR分解によって、 $n \times m$ 行列は

$$A = QR$$

と分解される。ここで、 Q は $n \times m$ 行列、 R は $m \times m$ の正方形行列。逆行列を求めることができる。

$\|A w - b\|$ がzeroとなるのはQRを使って、

$$Q \cdot R \cdot w = b$$

$$R \cdot w = Q^t \cdot b$$

$$R^{-1} \cdot R \cdot w = R^{-1} \cdot Q^t \cdot b$$

となりそう。

In [7]:

```

1 import numpy as np
2
3 tmp = np.fromfile('./codes/train_A.data',
4 np.float64, -1, " ")
5 A = tmp.reshape(300,30)
6 tmp = np.fromfile('./codes/train_b.data',
7 np.float64, -1, " ")
8 b = tmp.reshape(300,1)
9
10 q, r = np.linalg.qr(A)

```

```
In [8]: 1 ww = np.linalg.inv(r).dot(np.transpose(q).dot(b))
```

```
In [9]: 1 q.shape
```

```
Out[9]: (300, 30)
```

```
In [10]: 1 print(r[0,0:5])
```

```
[ -2.57579883e+02 -3.32324268e+02 -1.68607899e+03
 -1.29450676e+04
 -1.65446346e+00]
```

```
In [11]: 1 show_accuracy(A, b, ww)
```

```
correct: 286/ 300
safe error: 1
critical error: 13
```

```
In [12]: 1 print_w(ww)
```

```
(params)      :      (mean)      (stderr)      (w
orst)
      radius:  0.869921844 -0.024313948 -0.0626
79561
      texture: -0.003274619 -8.790300861  1.7471
47500
      perimeter: -0.202849407 -6.506451098  5.0617
60446
      area: 49.167541566 -0.956591421 -0.0820
52658
      smoothness: -0.007943157  0.004976908-27.8419
44367
      compactness:  3.301527110  4.985959134-16.3188
86295
      concavity: 10.316289081-21.332232171 -0.4086
05816
      concave points: -0.003345722 -0.000677873  0.0025
10735
      symmetry:  4.531369718  0.590110016 -0.7193
68704
      fractal dimension: -2.158965299 -3.803467225-12.2984
17038
```

```
In [13]: 1 tmp = np.fromfile('./codes/validate_A.data',
      np.float64, -1, " ")
2 A = tmp.reshape(260,30)
3 tmp = np.fromfile('./codes/validate_b.data',
      np.float64, -1, " ")
4 b = tmp.reshape(260,1)
5
6 show_accuracy(A, b, ww)
```

```
correct: 252/ 260
safe error: 6
critical error: 2
```

```
In [ ]: 1
```