

先端データ解析論
第十一回 レポート

37-196360 森田涼介

2019 年 7 月 8 日

宿題 1

フィッシャー判別分析を実装する。

クラス内散布行列は,

$$\mathbf{S}^{(w)} = \sum_{y=1}^c \sum_{i: y_i=y} (\mathbf{x}_i - \boldsymbol{\mu}_y)(\mathbf{x}_i - \boldsymbol{\mu}_y)^T \quad (1)$$

クラス間散布行列は,

$$\mathbf{S}^{(b)} = \sum_{y=1}^c n_y \boldsymbol{\mu}_y \boldsymbol{\mu}_y^T \quad (2)$$

である。クラス内散布を小さく, クラス間散布を大きくするような写像を求める。つまり, 次を求めればよい。

$$\mathbf{T}_{\text{FDA}} = \arg \max_{\mathbf{T}} \left((\mathbf{T} \mathbf{S}^{(w)} \mathbf{T}^T)^{-1} \mathbf{T} \mathbf{S}^{(b)} \mathbf{T}^T \right) \quad (3)$$

これは, 次の一般化固有値問題を解くことにより, 以下のように求まる。

$$\mathbf{S}^{(b)} \boldsymbol{\xi} = \lambda \mathbf{S}^{(w)} \boldsymbol{\xi} \quad (4)$$

$$\boldsymbol{\xi}_j^T \mathbf{S}^{(w)} \boldsymbol{\xi}_j = 1 \quad (5)$$

$$\mathbf{T}_{\text{FDA}} = (\boldsymbol{\xi}_1, \dots, \boldsymbol{\xi}_m)^T \quad (6)$$

7 ページの Listing 1 にプログラムを示した。結果は図 1, 2 に示した通りである。

これらより, フィッシャー判別分析によって異なるクラスのデータをうまく分離できるような写像が求まるが, クラス内にクラスタ構造があるとうまくいかないことがわかる。

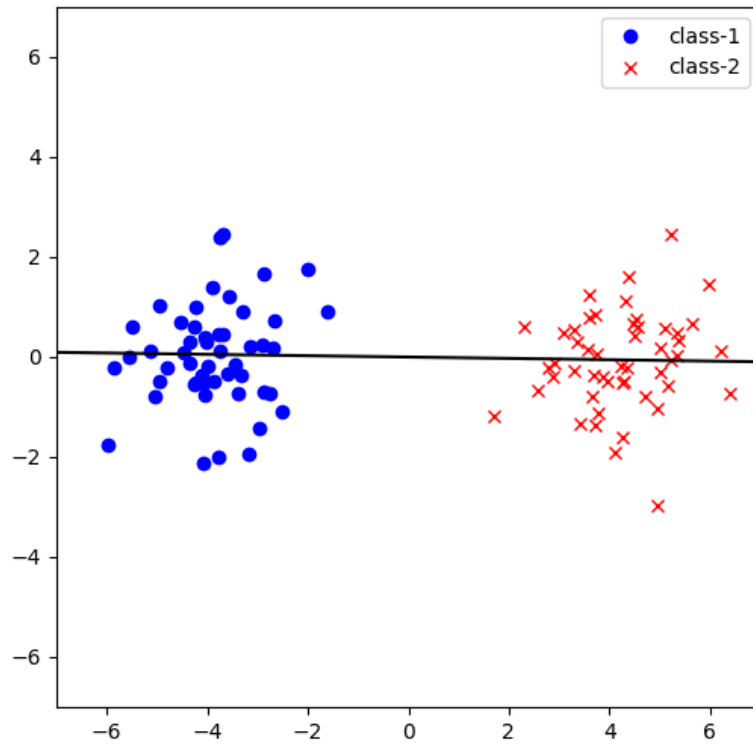


図 1: 2 クラスタのデータに対する結果

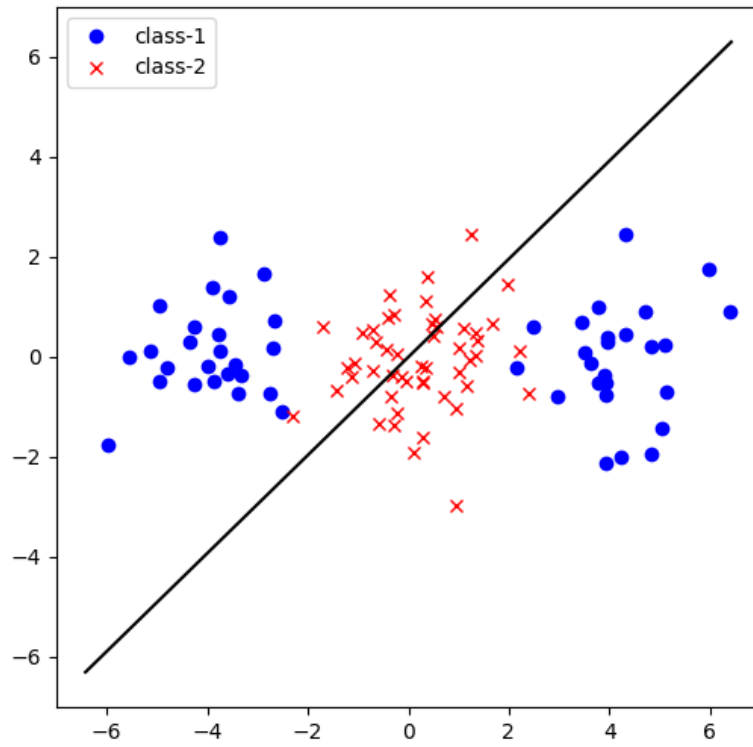


図 2: 3 クラスタのデータに対する結果

宿題 2

標本 $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$ を考える。 \mathbf{x} は中心化されているとする。つまり,

$$\boldsymbol{\mu} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i = \mathbf{0} \quad (7)$$

とする。このとき、散布行列は,

$$\mathbf{C} = \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^T \quad (8)$$

クラス内散布行列は,

$$\mathbf{S}^{(w)} = \sum_{y=1}^c \sum_{i: y_i=y} (\mathbf{x}_i - \boldsymbol{\mu}_y) (\mathbf{x}_i - \boldsymbol{\mu}_y)^T \quad (9)$$

クラス間散布行列は,

$$\mathbf{S}^{(b)} = \sum_{y=1}^c n_y \boldsymbol{\mu}_y \boldsymbol{\mu}_y^T \quad (10)$$

である。ここで、 n_y はクラス y の標本数であり,

$$\boldsymbol{\mu}_y = \frac{1}{n_y} \sum_{i: y_i=y} \mathbf{x}_i \quad (11)$$

である。

いま、散布行列について,

$$\mathbf{C} = \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^T = \frac{1}{2} \left(\sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^T + \sum_{i'=1}^n \mathbf{x}_{i'} \mathbf{x}_{i'}^T \right)$$

と変形できる。

$$\sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^T = \sum_{i'=1}^n \frac{1}{n} \left(\sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^T \right) = \frac{1}{n} \sum_{i, i'=1}^n \mathbf{x}_i \mathbf{x}_i^T$$

であることから,

$$\mathbf{C} = \frac{1}{2} \cdot \frac{1}{n} \sum_{i, i'}^n (\mathbf{x}_i \mathbf{x}_i^T + \mathbf{x}_{i'} \mathbf{x}_{i'}^T)$$

となる。ここで、式 (7) から,

$$\begin{aligned} \sum_{i, i'=1}^n \mathbf{x}_i \mathbf{x}_{i'}^T &= \sum_{i=1}^n \mathbf{x}_i \left(\sum_{i'=1}^n \mathbf{x}_{i'} \right)^T + \sum_{i'=1}^n \left(\sum_{i=1}^n \mathbf{x}_i \right) \mathbf{x}_{i'}^T \\ &= 2 \sum_{i=1}^n \mathbf{x}_i \sum_{i'=1}^n \mathbf{x}_{i'}^T \\ &= 2n \boldsymbol{\mu} \cdot n \boldsymbol{\mu}^T \\ &= \mathbf{0} \end{aligned}$$

が成立することを用いると、結局、

$$\begin{aligned}\mathbf{C} &= \frac{1}{2} \cdot \frac{1}{n} \sum_{i,i'=1}^n (\mathbf{x}_i \mathbf{x}_i^T - 2\mathbf{x}_i \mathbf{x}_{i'}^T + \mathbf{x}_{i'} \mathbf{x}_{i'}^T) \\ &= \frac{1}{2} \sum_{i,i'=1}^n \frac{1}{n} (\mathbf{x}_i - \mathbf{x}_{i'}) (\mathbf{x}_i - \mathbf{x}_{i'})^T\end{aligned}\tag{12}$$

得る。

同様に、クラス内散布行列は、

$$\begin{aligned}\mathbf{S}^{(w)} &= \sum_{y=1}^c \sum_{i:y_i=y} (\mathbf{x}_i - \boldsymbol{\mu}_y) (\mathbf{x}_i - \boldsymbol{\mu}_y)^T \\ &= \frac{1}{2} \sum_{y=1}^c \frac{1}{n_y} \sum_{i,i':y_{i,i'}=y} \left\{ (\mathbf{x}_i - \boldsymbol{\mu}_y) (\mathbf{x}_i - \boldsymbol{\mu}_y)^T + (\mathbf{x}_{i'} - \boldsymbol{\mu}_y) (\mathbf{x}_{i'} - \boldsymbol{\mu}_y)^T \right\}\end{aligned}$$

と変形でき、ここで、式 (11) から、

$$\begin{aligned}\sum_{i,i':y_{i,i'}=y} (\mathbf{x}_i - \boldsymbol{\mu}_y) (\mathbf{x}_{i'} - \boldsymbol{\mu}_y)^T &= 2 \sum_{i:y_i=y} (\mathbf{x}_i - \boldsymbol{\mu}_y) \cdot \sum_{i':y_{i'}=y} (\mathbf{x}_{i'} - \boldsymbol{\mu}_y)^T \\ &= 2(n_y \boldsymbol{\mu}_y - n_y \boldsymbol{\mu}_y) (n_y \boldsymbol{\mu}_y - n_y \boldsymbol{\mu}_y)^T \\ &= \mathbf{0}\end{aligned}$$

となることから、結局、

$$\begin{aligned}\mathbf{S}^{(w)} &= \frac{1}{2} \sum_{y=1}^c \frac{1}{n_y} \sum_{i,i':y_{i,i'}=y} \left\{ (\mathbf{x}_i - \boldsymbol{\mu}_y) (\mathbf{x}_i - \boldsymbol{\mu}_y)^T - 2(\mathbf{x}_i - \boldsymbol{\mu}_y) (\mathbf{x}_{i'} - \boldsymbol{\mu}_y)^T + (\mathbf{x}_{i'} - \boldsymbol{\mu}_y) (\mathbf{x}_{i'} - \boldsymbol{\mu}_y)^T \right\} \\ &= \frac{1}{2} \sum_{y=1}^c \frac{1}{n_y} \sum_{i,i':y_{i,i'}=y} \left\{ (\mathbf{x}_i - \boldsymbol{\mu}_y) - (\mathbf{x}_{i'} - \boldsymbol{\mu}_y) \right\} \left\{ (\mathbf{x}_i - \boldsymbol{\mu}_y) - (\mathbf{x}_{i'} - \boldsymbol{\mu}_y) \right\}^T \\ &= \frac{1}{2} \sum_{y=1}^c \frac{1}{n_y} \sum_{i,i':y_{i,i'}=y} (\mathbf{x}_i - \mathbf{x}_{i'}) (\mathbf{x}_i - \mathbf{x}_{i'})^T\end{aligned}$$

となる。これは、

$$Q_{i,i'}^{(w)} = \begin{cases} 1/n_y & (y_i = y_{i'} = y) \\ 0 & (y_i \neq y_{i'}) \end{cases}\tag{13}$$

なる $Q_{i,i'}^{(w)}$ を用いて、

$$\mathbf{S}^{(w)} = \frac{1}{2} \sum_{i,i'=1}^n Q_{i,i'}^{(w)} (\mathbf{x}_i - \mathbf{x}_{i'}) (\mathbf{x}_i - \mathbf{x}_{i'})^T\tag{14}$$

と表せる。

いま, $\mathbf{S}^{(w)}$ について変形すると,

$$\begin{aligned}
\mathbf{S}^{(w)} &= \sum_{y=1}^c \sum_{i:y_i=y} (\mathbf{x}_i - \boldsymbol{\mu}_y)(\mathbf{x}_i - \boldsymbol{\mu}_y)^T \\
&= \sum_{y=1}^c \sum_{i:y_i=y} (\mathbf{x}_i \mathbf{x}_i^T - 2\boldsymbol{\mu}_y \mathbf{x}_i^T + \boldsymbol{\mu}_y \boldsymbol{\mu}_y^T) \\
&= \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^T - \sum_{y=1}^c 2\boldsymbol{\mu}_y \sum_{i:y_i=y} \mathbf{x}_i^T + \sum_{y=1}^c n_y \boldsymbol{\mu}_y \boldsymbol{\mu}_y^T \\
&= \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^T - \sum_{y=1}^c 2\boldsymbol{\mu}_y n_y \boldsymbol{\mu}_y^T + \sum_{y=1}^c n_y \boldsymbol{\mu}_y \boldsymbol{\mu}_y^T \\
&= \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^T - \sum_{y=1}^c n_y \boldsymbol{\mu}_y \boldsymbol{\mu}_y^T
\end{aligned}$$

となり, 式 (8), (10) とから,

$$\mathbf{C} = \mathbf{S}^{(w)} + \mathbf{S}^{(b)} \quad (15)$$

が成立する。

式 (12), (14), (15) から,

$$\begin{aligned}
\mathbf{S}^{(b)} &= \mathbf{C} - \mathbf{S}^{(w)} \\
&= \frac{1}{2} \sum_{i,i'=1}^n \left(\frac{1}{n} - \mathcal{Q}_{i,i'}^{(w)} \right) (\mathbf{x}_i - \mathbf{x}_{i'}) (\mathbf{x}_i - \mathbf{x}_{i'})^T \\
&= \frac{1}{2} \sum_{i,i'=1}^n \mathcal{Q}_{i,i'}^{(b)} (\mathbf{x}_i - \mathbf{x}_{i'}) (\mathbf{x}_i - \mathbf{x}_{i'})^T \quad (16)
\end{aligned}$$

$$\begin{aligned}
\mathcal{Q}_{i,i'}^{(b)} &= \frac{1}{n} - \mathcal{Q}_{i,i'}^{(w)} \\
&= \begin{cases} 1/n - 1/n_y & (y_i = y_{i'} = y) \\ 1/n & (y_i \neq y_{i'}) \end{cases} \quad (17)
\end{aligned}$$

となる。

プログラム

実行環境と用いた言語・ライブラリを以下の表 1 に示す。

表 1: プログラムの実行環境

OS	: Microsoft Windows 10 Pro (64bit)
CPU	: Intel(R) Core(TM) i5-4300U
RAM	: 4.00 GB
使用言語	: Python3.6
可視化	: matplotlib ライブラリ

Listings 1: assignment1.py

```
1  # -*- coding: utf-8 -*-
2
3
4  import numpy as np
5  import matplotlib.pyplot as plt
6  import scipy.linalg
7
8
9  def generate_data(sample_size=100, pattern='two_cluster'):
10     if pattern not in ['two_cluster', 'three_cluster']:
11         raise ValueError('Dataset pattern must be one of '
12                             '[two_cluster, three_cluster].')
13     x = np.random.normal(size=(sample_size, 2))
14     if pattern == 'two_cluster':
15         x[:sample_size // 2, 0] -= 4
16         x[sample_size // 2:, 0] += 4
17     else:
18         x[:sample_size // 4, 0] -= 4
19         x[sample_size // 4:sample_size // 2, 0] += 4
20     y = np.ones(sample_size, dtype=np.int64)
21     y[sample_size // 2:] = 2
22     return x, y
23
24
25 def scatter_matrices(x, y):
26     n = x.shape[0]
27     d = x.shape[1]
28
29     labels = np.unique(y)
30     C, Sw, Sb = np.zeros((d, d)), np.zeros((d, d)), np.zeros((d, d))
31     for label in labels:
32         x_y = x[(y == label), :]
```



```

33     n_y = x_y.shape[0]
34     mu_y = x_y.mean(axis=0)[: , np.newaxis]
35     Sb += n_y * mu_y.dot(mu_y.T)
36
37     for i in range(n_y):
38         x_i = x_y[i][: , np.newaxis]
39         diff = (x_i - mu_y)
40         Sw += diff.dot(diff.T)
41
42         C += x_i.dot(x_i.T)
43
44     return C, Sw, Sb
45
46
47 def train(x, y, n_components):
48     """Fisher Discriminant Analysis.
49     Implement this function
50
51     Returns
52     -----
53     T : (1, 2) ndarray
54         The embedding matrix.
55     """
56
57     C, Sw, Sb = scatter_matrices(x, y)
58     eigen_values, eigen_vectors = scipy.linalg.eig(Sb, Sw)
59
60     # normalize
61     for i in range(len(eigen_vectors)):
62         eigen_vectors[i] = eigen_vectors[i]/np.linalg.norm(eigen_vectors[i])
63
64     T = eigen_vectors[:n_components]
65     return T
66
67
68 def visualize(x, y, T, path=None):
69     plt.figure(1, (6, 6))
70     plt.xlim(-7., 7.)
71     plt.ylim(-7., 7.)
72     plt.plot(x[y == 1, 0], x[y == 1, 1], 'bo', label='class-1')
73     plt.plot(x[y == 2, 0], x[y == 2, 1], 'rx', label='class-2')
74     plt.plot(
75         np.array([-T[:, 0], T[:, 0]]) * 9,
76         np.array([-T[:, 1], T[:, 1]]) * 9,
77         'k-',
78     )
79     plt.legend()

```

```

80     if path:
81         plt.savefig(str(path))
82     plt.show()
83
84
85 def main():
86     # settings
87     n = 100
88     n_components = 1
89     #mode = 'two_cluster'
90     mode = 'three_cluster'
91     fig_path = f'../figures/assignment1_result_{mode}.png'
92     np.random.seed(10)
93
94
95     # generate data
96     x, y = generate_data(sample_size=n, pattern=mode)
97     #print(x.shape, y.shape)
98
99
100    # train
101    T = train(x, y, n_components)
102
103
104    # result
105    print(f'data: {mode}  (#sample = {n})')
106    print(f'T = {T}')
107
108    visualize(x, y, T, path=fig_path)
109
110
111
112 if __name__ == '__main__':
113     main()

```