

# 資料科學

## Data Science

### 作業五 HW5

電機所 R11921038 江讀晉

2022/12/9

## Problem 1. Decision Tree split by Gini Index

### 1. Every Gini index of candidates

Depth	Feature	Feature value		Gini index
0	<b>Car Type</b>	<b>Family</b>	<b>Sports, Luxury</b>	<b>0.2769</b>
	Shirt Size	Small, Medium	Extra Large, Large	0.3542
	Car Type	Sports	Family, Luxury	0.3690
	Shirt Size	Extra Large	Small, Medium, Large	0.3938
	Shirt Size	Medium	Small, Large, Extra Large	0.4198
	Gender	M	F	0.45
	Shirt Size	Small, Extra Large	Medium, Large	0.4727
	Shirt Size	Large	Small, Medium, Extra Large	0.4750
	Car Type	Luxury	Family, Sports	0.4791
	Shirt Size	Small	Medium, Large, Extra Large	0.48
	Shirt Size	Small, Large	Medium, Extra Large	0.4949
1	<b>Shirt Size</b>	<b>Large, Extra Large</b>	<b>Small, Medium</b>	<b>0.1231</b>
	Shirt Size	Extra Large	Small, Medium, Large	0.2517
	Shirt Size	Medium	Small, Large, Extra Large	0.3077
	Shirt Size	Large	Small, Medium, Extra Large	0.3487
	Shirt size	Small	Medium, Large, Extra Large	0.3692
	Car Type	Sports	Luxury	0.3919
	Gender	M	F	0.4154
	Shirt Size	Small, Extra Large	Medium, Large	0.4154
	Shirt Size	Small, Large	Medium, Extra Large	0.4249
2	<b>Car Type</b>	<b>Sports</b>	<b>Luxury</b>	<b>0.2</b>
	Gender	M	F	0.2667
	Shirt Size	Large	Extra Large	0.2667
3	<b>Shirt Size</b>	<b>Large</b>	<b>Extra Large</b>	<b>0</b>

Table 1. Gini index of each candidate. The bold texts are the selected rules at each step.

## 2. Decision tree

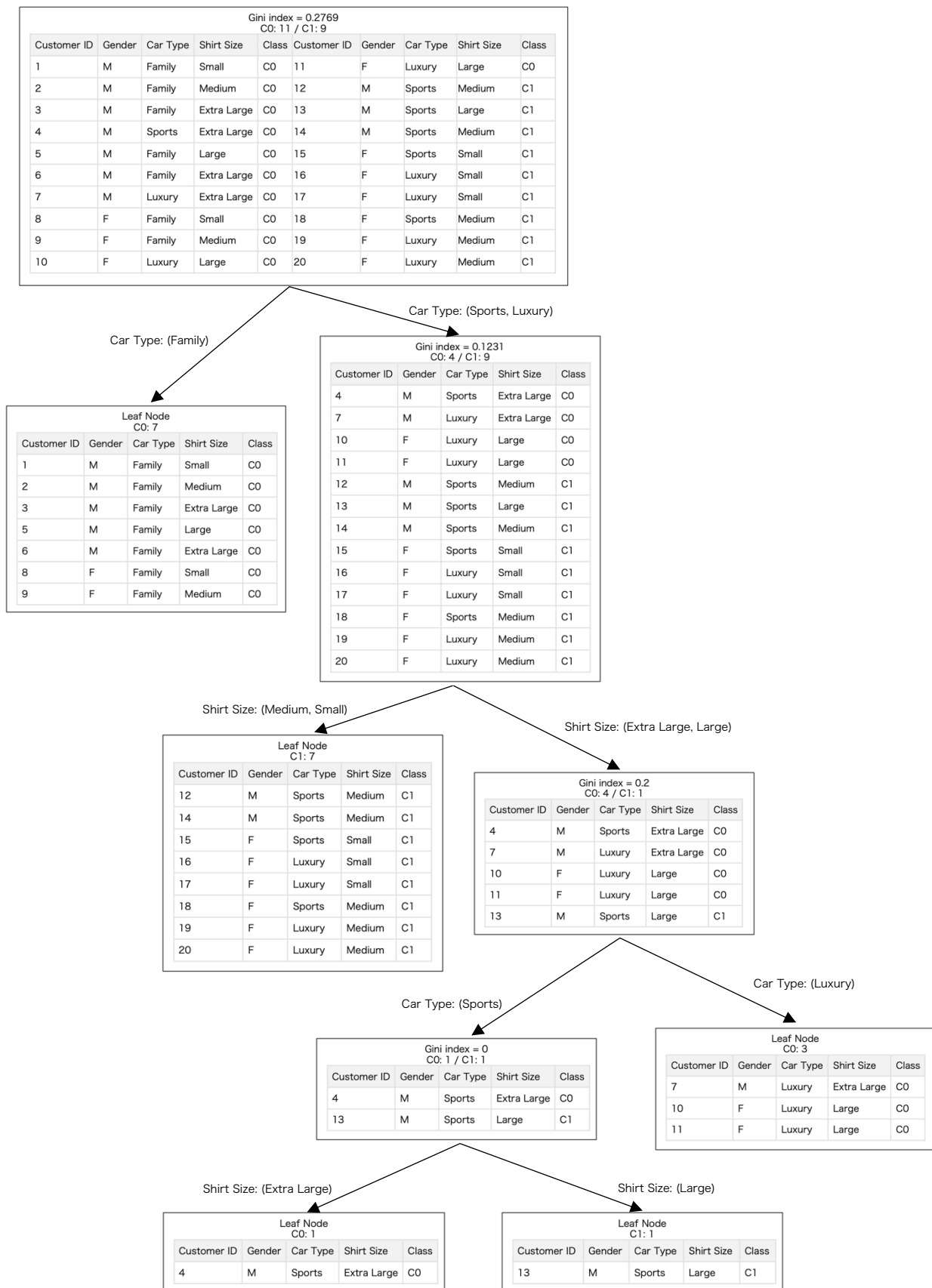


Fig. 1 Decision tree with node information

### 3. The programming implementation

```
1  class DecisionTree:
2      def __init__(self, data, columns = None, depth = 0, max_depth = 10):
3          self.data = data
4          self.columns = columns
5          self.depth = depth
6          self.max_depth = max_depth
7          self.left = None
8          self.right = None
9          self.gini_candidates = {'Feature': [], 'Value': [], 'Gini Index': []}
10         self.best_gini = None
11         self.split_feature = None
12         self.split_value = None
13         self.target = None
14         self.build_tree()
15
16     def build_tree(self):
17         # If the data is empty, return
18         if self.data.empty:
19             return
20
21         # If the data is pure, return
22         if len(self.data[self.columns[-1]].unique()) == 1:
23             self.target = self.data[self.columns[-1]].unique()[0]
24             return
25
26         # If the depth is greater than the max depth, return
27         if self.depth >= self.max_depth:
28             self.target = self.data[self.columns[-1]].value_counts().idxmax()
29             return
30
31         # Get the best split feature and value
32         self.best_gini, self.split_feature, self.split_value = self.get_best_split()
33
34         # Split the data
35         left_data = self.data[self.data[self.split_feature].isin(self.split_value)]
36         right_data = self.data[~self.data[self.split_feature].isin(self.split_value)]
37
38         # Display the decision tree
39         self.display_tree(left_data, right_data)
40
41         # Save the children nodes to csv file
42         left_data.to_csv(f'left_node_{self.depth}.csv', index = False)
43         right_data.to_csv(f'right_node_{self.depth}.csv', index = False)
44
45         # Build the left and right subtrees
46         self.left = DecisionTree(left_data, self.columns, self.depth + 1, self.max_depth)
47         self.right = DecisionTree(right_data, self.columns, self.depth + 1, self.max_depth)
48
49
50     def display_tree(self, left_data, right_data):
51         print('Depth:', self.depth)
52
53         candidates = pd.DataFrame(self.gini_candidates)
54         candidates = candidates.sort_values(by = 'Gini Index', ascending = True)
55         print(f'Gini candidates:\n{candidates}')
56
57         split = pd.DataFrame()
58         split['Split Feature'] = [self.split_feature]
59         split['Split Value'] = [self.split_value]
60         split['Gini Index'] = [self.best_gini]
61
62         print(f'Split parameter:\n{split}')
63         print(f'Left node:\n{left_data}')
64         print(f'Right node:\n{right_data}\n')
65
66     def get_best_split(self):
67         # Get the best split feature and value
```

```

68     best_split_feature = None
69     best_split_value = None
70     best_gini = 1.0
71
72     # Loop through all the features
73     for feature in self.columns[1:-1]:
74         # Get the unique values of the feature
75         values = self.data[feature].unique()
76         # print('ori', values)
77
78         # Build up the combination of the values
79         new_values = []
80         for i in range(len(values) // 2):
81             new_values.extend(list(combinations(values, i + 1)))
82
83         # Loop through all the values
84         for value in new_values:
85             # Split the data
86             left_data = self.data[self.data[feature].isin(value)]
87             right_data = self.data[~(self.data[feature].isin(value))]
88
89             # Calculate the gini index
90             gini = self.get_gini(left_data, right_data)
91             self.gini_candidates['Feature'].append(feature)
92             self.gini_candidates['Value'].append(value)
93             self.gini_candidates['Gini Index'].append(round(gini, 4))
94
95             # If the gini index is less than the best gini index, update the best gini
96             index if gini < best_gini:
97                 best_gini = round(gini, 4)
98                 best_split_feature = feature
99                 best_split_value = value
100
101         # Return the best split feature and value
102         return best_gini, best_split_feature, best_split_value
103
104     def get_gini(self, left_data, right_data):
105         # Calculate the gini index
106         gini_left = 0.0
107         gini_right = 0.0
108
109         # Get the target categories
110         target = self.data[self.columns[-1]].unique()
111
112         # Calculate the gini index
113         if len(left_data) > 0:
114             for t in target:
115                 gini_left += (len(left_data[left_data[self.columns[-1]] == t]) /
116 len(left_data)) ** 2
117             gini_left = 1 - gini_left
118
119         if len(right_data) > 0:
120             for t in target:
121                 gini_right += (len(right_data[right_data[self.columns[-1]] == t]) /
122 len(right_data)) ** 2
123             gini_right = 1 - gini_right
124
125         gini = (len(left_data) / len(self.data)) * gini_left + (len(right_data) /
126 len(self.data)) * gini_right
127
128         return gini
129
130 if __name__ == '__main__':
131     # Read in the data
132     df = pd.read_csv('data.csv')
133     print(df)

```

```

132 # Get the column names
133 columns = df.columns
134 print(columns)
135
136 # Get the unique feature categories and the target categories
137 gender = df[columns[1]].unique()
138 car = df[columns[2]].unique()
139 shirt = df[columns[3]].unique()
140 target = df[columns[4]].unique()
141
142 # Build the decision tree
143 tree = DecisionTree(df, columns)

```

## Problem 2. Naïve Bayes Classifier

Given tuple: (Gender=M, Car Type=Sports, Shirt Size=Medium)

*A: features of Gender, Car Type and Shirt Size*

$$P(A|C0) = P(\text{Gender}|C0) \cdot P(\text{Car Type}|C0) \cdot P(\text{Shirt Size}|C0) = \frac{7}{11} \cdot \frac{1}{11} \cdot \frac{2}{11} = \frac{14}{1331} = 0.0105$$

$$P(A|C1) = P(\text{Gender}|C1) \cdot P(\text{Car Type}|C1) \cdot P(\text{Shirt Size}|C1) = \frac{3}{9} \cdot \frac{5}{9} \cdot \frac{5}{9} = \frac{25}{243} = 0.1029$$

$$P(A|C0)P(C0) = \frac{14}{1331} \cdot \frac{11}{20} = 0.0058$$

$$P(A|C1)P(C1) = \frac{25}{243} \cdot \frac{9}{20} = 0.0463$$

As  $P(A|C0)P(C0) < P(A|C1)P(C1)$ , the tuple (Gender=M, Car Type=Sports, Shirt Size=Medium) should be classified into class **C1**.

## Problem 3. SVM (Support Vector Machine)

Positive samples,  $y = 1$ : (4, 3), (4, 8), (7, 2)

Negative samples,  $y = -1$ : (-1, -2), (-1, 3), (2, -1), (2, 1)

### 1. Objectives and constraints

The objective with the constraint is to

$$\begin{aligned} & \text{maximize } \frac{2}{\|w\|^2} \\ & \text{subject to } y_i(w^T x_i + b) - 1 \geq 0, \forall x_i \end{aligned}$$

Hence, the hinge loss with the regularization of the weights is used.

$$J = \lambda \|w\|^2 + \frac{1}{n} \sum_{i=1}^n \max(0, 1 - y_i(w \cdot x_i + b))$$

## 2. Support vectors

The two support vectors are **(4, 3)** and **(2, 1)**. These two support vectors are highlighted in the Fig. 3.

## 3. Computing progress

Based on the loss function, the gradient of weight and bias can be computed.

$$\begin{aligned} \text{if } y_i(w \cdot x_i + b) \geq 1, \text{ then } J_i = \lambda \|w\|^2 & \quad \therefore \frac{\partial J_i}{\partial w_k} = 2\lambda w_k \\ \text{else } J_i = \lambda \|w\|^2 + 1 - y_i(w \cdot x_i + b) & \quad \therefore \begin{cases} \frac{\partial J_i}{\partial w_k} = 2\lambda w_k - y_i \cdot x_i \\ \frac{\partial J_i}{\partial b} = -y_i \end{cases} \end{aligned}$$

The following is my programming implementation.

```
1 class LinearSVM:
2     def __init__(self, learning_rate=0.001, lambda_param=0.1, n_iters=1000):
3         self.lr = learning_rate
4         self.lambda_param = lambda_param
5         self.n_iters = n_iters
6         self.w = None
7         self.b = None
8         self.w_list = []
9         self.b_list = []
10
11     def fit(self, x, y):
12         self.w = np.zeros(x.shape[1])
13         self.b = 0
14
15         self.w_list = [list(self.w)]
16         self.b_list = [self.b]
17         for _ in range(self.n_iters):
18             for i in range(x.shape[0]):
19                 if y[i] * (np.dot(x[i], self.w) + self.b) >= 1:
20                     self.w -= self.lr * (2 * self.lambda_param * self.w)
21                 else:
22                     self.w -= self.lr * (2 * self.lambda_param * self.w - np.dot(x[i], y[i]))
23                     self.b -= self.lr * (-y[i])
24
25             if np.linalg.norm(self.w - self.w_list[-1]) < 1e-4:
26                 break
27             else:
28                 self.w_list.append(list(self.w))
29                 self.b_list.append(self.b)
30
31 if __name__ == '__main__':
32     # Construct the data
33     data = np.array([[4, 3], [4, 8], [7, 2], [-1, -2], [-1, 3], [2, -1], [2, 1]])
34     target = np.array([1, 1, 1, -1, -1, -1, -1])
35
36     # Train the model
37     svm = LinearSVM(n_iters=10000)
38     svm.fit(data, target)
```

## 4. Weight, bias and hyperplane

Based the codes above, the weight and the bias are converged to

$$w = \begin{bmatrix} 0.4999023 \\ 0.5003976 \end{bmatrix}, b = -2.5039999$$

If the weight and the bias are round to 1 decimal place, then

$$\mathbf{w} = \begin{bmatrix} 0.5 \\ 0.5 \end{bmatrix}, b = -2.5$$

The curves of the weight and the bias in every iteration are shown in the Fig. 2.

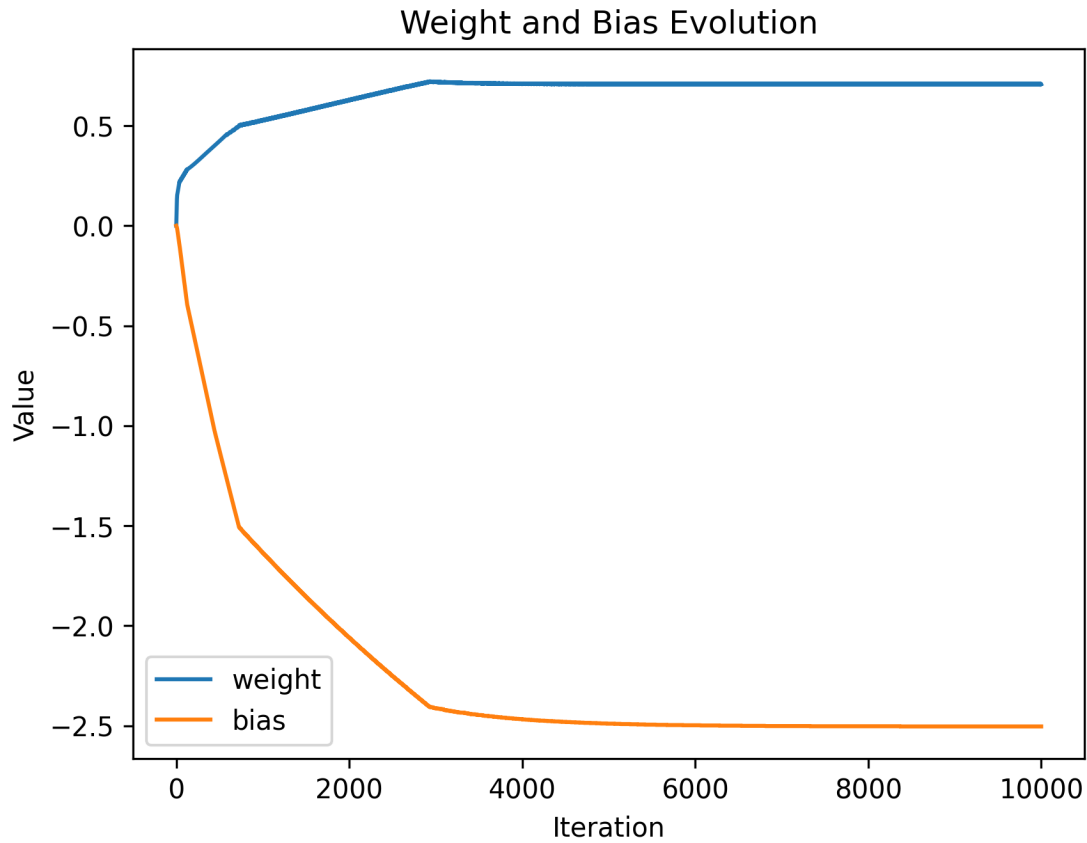


Fig. 2 The evolution of the weight and the bias

The decision boundary or hyperplane is

$$\begin{aligned} \mathbf{y} &= \mathbf{w}^T \mathbf{x} + b \\ &= w_1 x_1 + w_2 x_2 + b \\ &= 0.5x_1 + 0.5x_2 - 2.5 \end{aligned}$$

To illustrate the results, the data points with labels, hyperplane (decision boundary), two boundaries of margin and support vectors are manifested in the Fig. 3.

Note that HW5-P3 is a hard SVM and there are exactly two support vectors on the boundaries of margin.

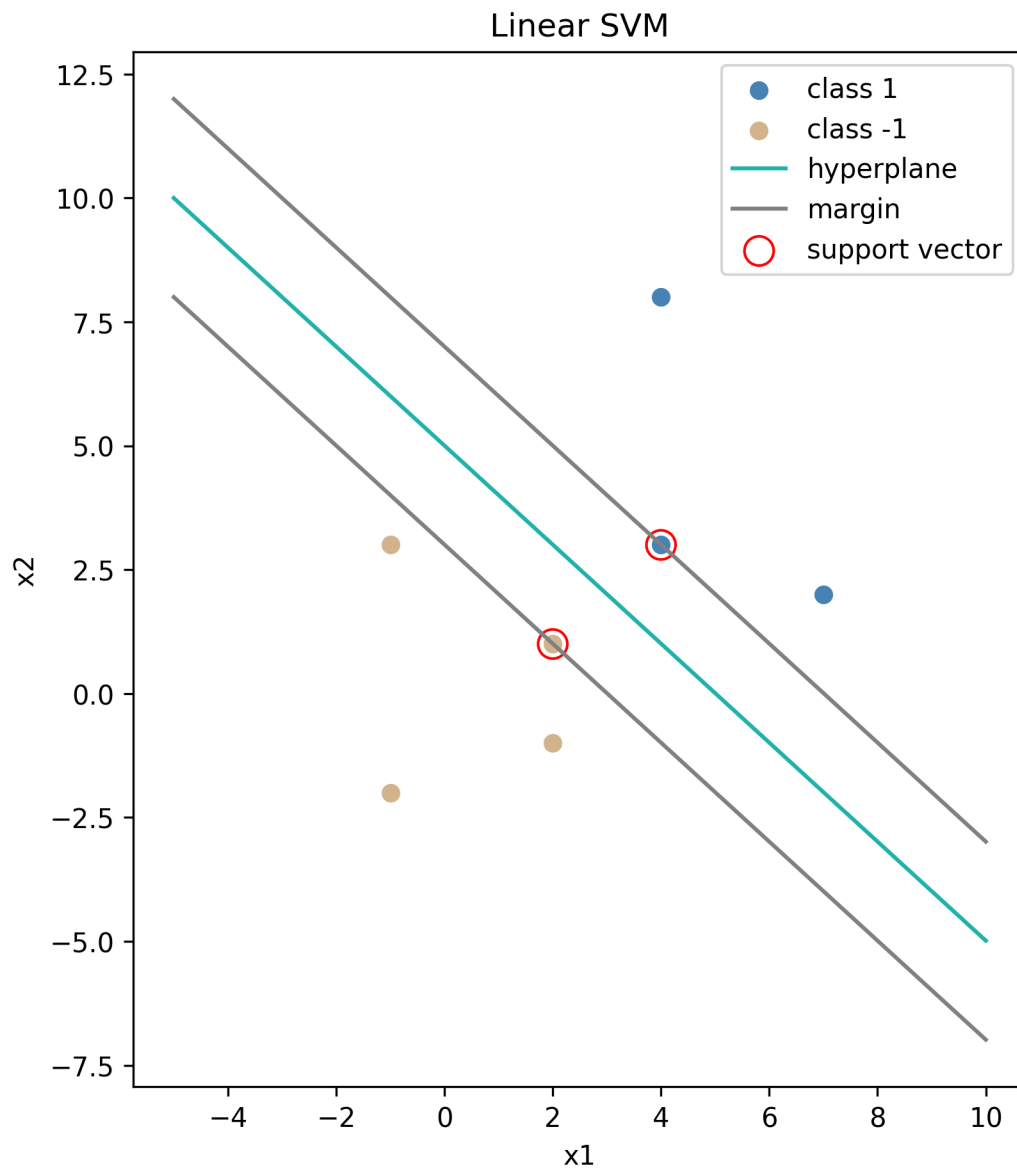


Fig. 3 The result of the SVM classifier