

Лабораторна робота №2. Методи класифікації

Виконав студент групи КМ-91мп

Галета М.С.

Завдання на лабораторну роботу

1. Обрати відповідний файл з даними.
2. Збудувати дерево прийняття рішень для визначення значення цільової характеристики (останній стовпчик) на основі вхідних аргументів (інші стовпчики). Записи у файлі з даними потрібно розділити на навчальну та відкладену вибірки.
3. Параметри дерева та алгоритм його побудови обрати самостійно. Дерево не повинно бути занадто гіллястим, за потреби використовуйте процедуру обрізання. Використовувати бібліотечні функції забороняється.
4. Оцінити якість збудованого дерева за допомогою методу відкладеної вибірки. Надати графічне порівняння результатів, отриманих за допомогою застосування збудованого дерева прийняття рішень до відкладеної вибірки, із фактичними значеннями цільової характеристики.
5. Застосувати для визначення цільової характеристики на даних відкладеної вибірки метод k найближчих сусідів. Визначити кількість сусідів, при якій метод дає найкращу точність.
6. Порівняти результати з п.п. 4 та 5.

In [1]:

```
1 import numpy as np
2 import math
3 import pandas as pd
4 import matplotlib.pyplot as plt
5 from sklearn.model_selection import train_test_split
6
7 %matplotlib inline
```

Зчитування датасету та його розбиття на тренувальну та відкладену вибірки

In [2]:

```
1 df = pd.read_csv('MP-04-Galeta.csv', sep=';', names=['x1', 'x2', 'x3', 'x4', 'x5', 'x6']
2
3 X_train, X_test = train_test_split(df, test_size=0.2, random_state=42)
```

Побудова дерева рішень

In [3]:

```
1 class Tree:
2     def __init__(self, parent=None):
3         self.parent = parent
4         self.children = []
5         self.splitFeature = None
6         self.splitFeatureValue = None
7         self.label = None
```

In [4]:

```
1 def dataToDistribution(data):
2     ''' Функція перетворює набір даних, який має n можливих
3         класифікаційних міток, в розподіл ймовірностей з n записами. '''
4
5     allLabels = [label for (point, label) in data]
6     numEntries = len(allLabels)
7     possibleLabels = set(allLabels)
8     dist = []
9     for aLabel in possibleLabels:
10         dist.append(float(allLabels.count(aLabel)) / numEntries)
11     return dist
```

In [5]:

```
1 def entropy(dist):
2     ''' Функція обчислює ентропію Шеннона для заданого розподілу ймовірностей. '''
3     return -sum([p * math.log(p, 2) for p in dist])
```

In [6]:

```
1 def splitData(data, featureIndex):
2     ''' Функція виконує ітерацію над підмножинами даних,
3         що відповідають кожному значенню функції в індексі featureIndex. '''
4
5     # отримання можливих значень заданої ознаки
6     attrValues= [point[featureIndex] for (point, label) in data]
7     for aValue in set(attrValues):
8         # обчислення частини розбиття, що відповідає обраному значенню
9         dataSubset= [(point, label) for (point, label) in data if point[featureIndex] == aValue]
10        yield dataSubset
```

In [7]:

```
1 def gain(data, featureIndex):
2     ''' Функція обчислює очікуваний приріст інформації від
3         розбиття даних на всі можливі значення функції '''
4
5     entropyGain = entropy(dataToDistribution(data))
6     for dataSubset in splitData(data, featureIndex):
7         entropyGain -= entropy(dataToDistribution(dataSubset))
8     return entropyGain
```

In [8]:

```
1 def homogeneous(data):
2     ''' Функція повертає True, якщо дані мають однакову мітку,
3         і False в іншому випадку '''
4     return len(set([label for (point, label) in data])) <= 1
```

In [9]:

```
1 def majorityVote(data, node):
2     ''' Функція повертає вузол мітки з більшістю міток класу в даному наборі даних '''
3
4     labels = [label for (point, label) in data]
5     choice = max(set(labels), key=labels.count)
6     node.label= choice
7     return node
```

In [10]:

```

1 def buildDecisionTree(data, root, remainingFeatures, max_depth):
2     ''' Функція будує дерево рішень з даних даних,
3         додаючи дітей до кореневого вузла (який може бути піддеревом) '''
4
5     # Глобальна змінна, яка позначає поточну глибину дерева
6     global current_depth
7
8     # Зупиняємось, якщо поточна глибина більша за максимально допустиму глибину дерева
9     if current_depth >= max_depth:
10         remainingFeatures = []
11
12     # Якщо всі елементи мають однаковий клас, то класифікуємо листок дерева міткою цього класу
13     if homogeneous(data):
14         root.label = data[0][1]
15         root.classCounts = {root.label: len(data)}
16         return root
17
18     # Якщо немає більше ознак для подальшого розбиття,
19     # то класифікуємо листок дерева міткою класу якого більшість
20     if len(remainingFeatures) == 0:
21         return majorityVote(data, root)
22
23     # Знаходження індексу найкращої функції для поділу
24     bestFeature = max(remainingFeatures, key=lambda index: gain(data, index))
25     # Якщо приріст інформації дорівнює нулю,
26     # то класифікуємо листок дерева міткою класу якого більшість
27     if gain(data, bestFeature) == 0:
28         return majorityVote(data, root)
29
30     root.splitFeature = bestFeature
31
32     #Додавання дочірніх вузлів для вже існуючих
33     for dataSubset in splitData(data, bestFeature):
34         aChild = Tree(parent=root)
35         aChild.splitFeatureValue = dataSubset[0][0][bestFeature]
36         root.children.append(aChild)
37
38     #Запуск рекурсивного процесу, де дочірні вузли виступають в якості батьківських
39     buildDecisionTree(dataSubset, aChild, remainingFeatures - set([bestFeature]),
40
41     current_depth += 1
42
43     return root

```

In [11]:

```

1 def decisionTree(data, max_depth):
2     ''' Функція повертає збудоване дерево '''
3     return buildDecisionTree(data, Tree(), set(range(len(data[0][0]))), max_depth)

```

In [12]:

```
1 def classify(tree, point):
2     ''' Функція класифікації даних шляхом проходження даного дерева рішень. '''
3     if tree.children== []:
4         return tree.label
5     else:
6         matchingChildren = []
7         for child in tree.children:
8             if child.splitFeatureValue == point[tree.splitFeature]:
9                 matchingChildren.append(child)
10
11         try:
12             return classify(matchingChildren[0], point)
13         except Exception:
14             return child.label
```

Оцінювання якості збудованого дерева

In [13]:

```

1 train_data = [(x[:-1], x[-1]) for x in X_train.values.tolist()]
2 test_data = [(x[:-1], x[-1]) for x in X_test.values.tolist()]
3
4 y_train = np.array([label for point, label in train_data])
5 y_test = np.array([label for point, label in test_data])
6
7 for max_depth in range(2, 11):
8     current_depth = 0
9     tree = decisionTree(train_data, max_depth=max_depth)
10
11     y_predicted_train = [classify(tree, point) for point, label in train_data]
12     y_predicted_test = [classify(tree, point) for point, label in test_data]
13
14     accuracy_test = (y_predicted_test == y_test).mean()
15     accuracy_train = (y_predicted_train == y_train).mean()
16
17     print('==== max_depth - ', max_depth, '====')
18     print('Train accuracy:', np.round(accuracy_train*100,2), '%')
19     print('Test accuracy:', accuracy_test*100, '%')
20     print('')

```

```

==== max_depth - 2 ====
Train accuracy: 44.23 %
Test accuracy: 23.076923076923077 %

```

```

==== max_depth - 3 ====
Train accuracy: 44.23 %
Test accuracy: 23.076923076923077 %

```

```

==== max_depth - 4 ====
Train accuracy: 44.23 %
Test accuracy: 23.076923076923077 %

```

```

==== max_depth - 5 ====
Train accuracy: 44.23 %
Test accuracy: 23.076923076923077 %

```

```

==== max_depth - 6 ====
Train accuracy: 46.15 %
Test accuracy: 15.384615384615385 %

```

```

==== max_depth - 7 ====
Train accuracy: 46.15 %
Test accuracy: 15.384615384615385 %

```

```

==== max_depth - 8 ====
Train accuracy: 57.69 %
Test accuracy: 0.0 %

```

```

==== max_depth - 9 ====
Train accuracy: 63.46 %
Test accuracy: 0.0 %

```

```

==== max_depth - 10 ====
Train accuracy: 63.46 %
Test accuracy: 0.0 %

```

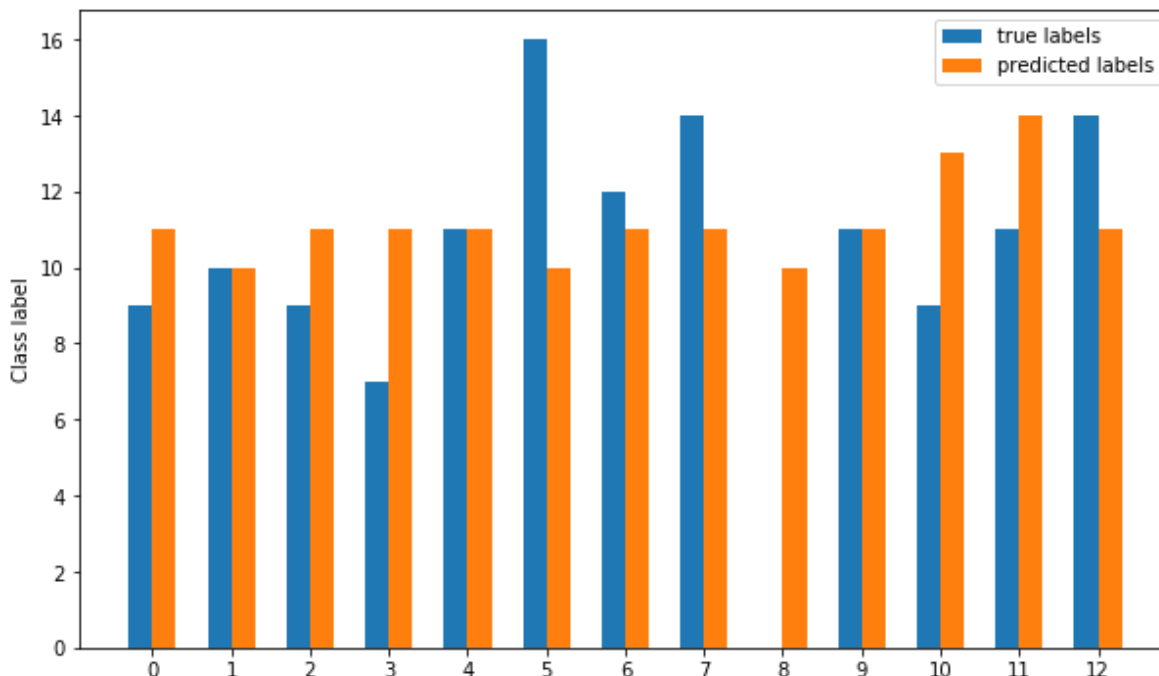
Отже, при глибині дерева 2, 3, 4 або 5 воно найкраще класифікує відкладену вибірку

In [14]:

```
1 train_data = [(x[:-1], x[-1]) for x in X_train.values.tolist()]
2 test_data = [(x[:-1], x[-1]) for x in X_test.values.tolist()]
3
4 y_train = np.array([label for point, label in train_data])
5 y_test = np.array([label for point, label in test_data])
6
7 current_depth = 0
8 tree = decisionTree(train_data, max_depth=4)
9
10 y_predicted_train = [classify(tree, point) for point, label in train_data]
11 y_predicted_test = [classify(tree, point) for point, label in test_data]
```

In [15]:

```
1 x = np.arange(len(y_test))
2 width = 0.3
3
4 fig, ax = plt.subplots(figsize=(10,6))
5 ax.bar(x - width/2, y_test, width, label='true labels')
6 ax.bar(x + width/2, y_predicted_test, width, label='predicted labels')
7
8 ax.set_ylabel('Class label')
9 ax.set_xticks(x)
10 ax.legend()
11 plt.show()
```



Метод KNN

In [16]:

```
1 x_train = X_train[X_train.columns[0:-1]].values
2 y_train = X_train[X_train.columns[-1]].values
3 x_test = X_test[X_test.columns[0:-1]].values
4 y_test = X_test[X_test.columns[-1]].values
```

In [17]:

```
1 def compute_distances(x_known, x_unknown):
2
3     num_pred = x_unknown.shape[0]
4     num_data = x_known.shape[0]
5
6     dists = np.zeros((num_pred, num_data))
7
8     for i in range(num_pred):
9         for j in range(num_data):
10             dists[i,j] = np.sum(x_known[j] == x_unknown[i])
11
12     return dists
```

In [18]:

```
1 def k_nearest_labels(dists, y_known, k):
2
3     num_pred = dists.shape[0]
4     n_nearest = []
5
6     for j in range(num_pred):
7         dst = dists[j]
8         closest_y = y_known[np.argpartition(dst, k-1)[:k]]
9
10         n_nearest.append(closest_y)
11
12     return np.asarray(n_nearest)
```


In [19]:

```

1  class KNearest_Neighbours:
2
3      def __init__(self, k):
4          self.k = k
5          self.test_set_x = None
6          self.train_set_x = None
7          self.train_set_y = None
8
9
10     def fit(self, train_set_x, train_set_y):
11         self.train_set_x = train_set_x
12         self.train_set_y = train_set_y
13
14     def predict(self, test_set_x):
15
16         y_labels = k_nearest_labels(compute_distances(self.train_set_x, test_set_x), self.k)
17         y_predictions = []
18         for i in range(y_labels.shape[0]):
19             bc = np.bincount(y_labels[i])
20             y_predictions.append(np.arange(len(bc))[bc == bc.max()].min())
21
22     return y_predictions

```

In [20]:

```

1  for k in range(1, 10, 2):
2      knn = KNearest_Neighbours(k)
3      knn.fit(x_train, y_train)
4      y_pred_knn = knn.predict(x_test)
5      accuracy = (y_pred_knn == y_test).mean() * 100
6      print("k = {}; accuracy = {}".format(k, accuracy))

```

```

k = 1; accuracy = 7.6923076923076925%
k = 3; accuracy = 15.384615384615385%
k = 5; accuracy = 15.384615384615385%
k = 7; accuracy = 15.384615384615385%
k = 9; accuracy = 7.6923076923076925%

```

З даних результатів видно, що для даного набору даних кількість сусідів, при якій алгоритм дає найбільшу точність, дорівнює 3, 5 або 7. Гірша точність при K = 1 або 9.

In [21]:

```

1  knn = KNearest_Neighbours(5)
2  knn.fit(x_train, y_train)
3  y_pred_knn = knn.predict(x_test)

```

In [22]:

```
1 print("К найближчих сусідів: ", y_pred_knn)
2 print("Дерева прийняття рішень: ", y_predicted_test)
```

```
К найближчих сусідів:      [10, 10, 10, 10, 13, 10, 13, 8, 10, 13, 9, 9, 10]
Дерева прийняття рішень:  [11, 10, 11, 11, 11, 10, 11, 11, 10, 11, 13, 14, 1
1]
```

In [23]:

```
1 print("Кількість однаково спрогнозованих значень двома методами: ", sum(np.array(y_pred:
```

```
Кількість однаково спрогнозованих значень двома методами:  3
```

Висновки

1) Було реалізовано два методи класифікації: дерево прийняття рішень та KNN

2) В ході роботи було встановлено, що дерево прийняття рішень дає кращий результат на відкладеній вибірці, ніж KNN