

# 大数据导论lab2 实验报告

220110519 邢瑞龙

## 数据集以及预处理

数据集: [Adult Data Set](#)

训练集数量: 26904

测试集数量: 14130

样本特征: 14种

样本类别: 2种 (个人收入大于50K or 小于等于50K)

预处理:

1. 去除缺失值, 异常值, 重复行:

```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
import collections
import random

columns = ['age', 'workclass', 'fnlwt', 'education', 'educationNum',
           'maritalStatus', 'occupation', 'relationship', 'race', 'sex',
           'capitalGain', 'capitalLoss', 'hoursPerWeek', 'nativeCountry',
           'income']
df_train_set = pd.read_csv('./adult.data', names=columns)
df_test_set = pd.read_csv('./adult.test', names=columns, skiprows=1) #第一行是
非法数据

print(len(df_train_set))
print(len(df_test_set))
df_train_set.to_csv('./train_adult.csv', index=False)
df_test_set.to_csv('./test_adult.csv', index=False)

df_train_set = pd.read_csv('./train_adult.csv')
df_test_set = pd.read_csv('./test_adult.csv')

df_train_set.drop(['fnlwt', 'educationNum'], axis=1, inplace=True) # fnlwt
列用处不大, educationNum与education类似
df_test_set.drop(['fnlwt', 'educationNum'], axis=1, inplace=True)
df_train_set.drop_duplicates(inplace=True) # 去除重复行
df_test_set.drop_duplicates(inplace=True)
df_train_set.dropna(inplace=True) # 去除空行
df_test_set.dropna(inplace=True)

# 去除含有'?'的行
new_columns = ['workclass', 'education', 'maritalStatus', 'occupation',
               'relationship', 'race', 'sex',
               'nativeCountry', 'income']
```

```

for col in new_columns:
    df_train_set = df_train_set[~df_train_set[col].str.contains(r'\?',
regex=True)]
    df_test_set = df_test_set[~df_test_set[col].str.contains(r'\?',
regex=True)]
#save to csv

```

## 2. 连续变量离散化 (为了缩小决策树的最优切分点的搜索范围, 这里将连续变量离散化)

```

#连续变量离散化
continuous_column = ['age', 'capitalGain', 'capitalLoss', 'hoursPerWeek']
allbins = [[0, 20, 40, 60, 80,100],
            [0, 10000, 50000, 100000],
            [0,1,5000],
            [0, 20, 40, 60, 80, 100]]
for col, bins in zip(continuous_column, allbins):
    df_train_set[col] = pd.cut(df_train_set[col], bins, right=False,
labels=False)
    df_test_set[col] = pd.cut(df_test_set[col], bins, right=False,
labels=False)

print(df_train_set.head())
print(df_test_set.head())

```

## 3. 离散变量index化 (为了方便程序处理, 这里将单个属性下的不同类别映射到不同的数字)

```

# #离散变量index化
discrete_column = ['workclass', 'education', 'maritalStatus', 'occupation',
'relationship', 'race', 'sex', 'nativeCountry', 'income']
workclass_mapping = {' Private': 0, ' Self-emp-not-inc': 1, ' Self-emp-inc':
1, ' Local-gov': 2,
                    ' State-gov': 2, ' Federal-gov': 2, ' without-pay': 3, '
Never-worked': 3}
education_mapping = {
    ' 5th-6th': 0,
    ' 7th-8th': 0,
    ' 9th': 0,
    ' 10th': 0,
    ' 11th': 0,
    ' 12th': 0,
    ' HS-grad': 0,
    ' Preschool': 1,
    ' 1st-4th': 1,
    ' Assoc-acdm': 2,
    ' Assoc-voc': 2,
    ' Some-college': 3,
    ' Bachelors': 3,
    ' Doctorate': 4,
    ' Prof-school': 4,
    ' Masters': 4
}

```

```

income_mapping = {' <=50K': 0, ' <=50K.':0, ' >50K': 1, ' >50K.': 1}
special_mapping_name = ['workclass', 'education', 'income']
df_test_set['workclass'] = df_test_set['workclass'].map(workclass_mapping)
df_train_set['workclass'] = df_train_set['workclass'].map(workclass_mapping)
df_test_set['education'] = df_test_set['education'].map(education_mapping)
df_train_set['education'] = df_train_set['education'].map(education_mapping)
df_test_set['income'] = df_test_set['income'].map(income_mapping)
df_train_set['income'] = df_train_set['income'].map(income_mapping)
print(df_train_set.head())
print(df_test_set.head())
for col in discrete_column:
    if(col in special_mapping_name):
        continue
    else:
        res1 = df_test_set[col].value_counts().keys()
        res2 = df_train_set[col].value_counts().keys()
        res = list(set(res1).union(set(res2)))
        mapping = dict(zip(res, range(len(res))))
        print(mapping)
        df_train_set[col] = df_train_set[col].map(mapping)
        df_test_set[col] = df_test_set[col].map(mapping)
print(df_train_set.head())
print(df_test_set.head())
#save to csv
df_train_set.to_csv('../train_adult_pro.csv', index=False)
df_test_set.to_csv('../test_adult_pro.csv', index=False)

```

- 特殊地，对于workclass：我们将其分为四类

```

workclass_mapping = {' Private': 0, ' Self-emp-not-inc': 1, ' Self-emp-
inc': 1, ' Local-gov': 2, ' State-gov': 2, ' Federal-gov': 2, ' Without-
pay': 3, ' Never-worked': 3}

```

- 对于education,这里划分为四类：幼稚园阶段，中小学阶段，大学阶段和更高学历阶段

```

education_mapping = {
    ' 5th-6th': 0,
    ' 7th-8th': 0,
    ' 9th': 0,
    ' 10th': 0,
    ' 11th': 0,
    ' 12th': 0,
    ' HS-grad': 0,
    ' Preschool': 1,
    ' 1st-4th': 1,
    ' Assoc-acdm': 2,
    ' Assoc-voc': 2,
    ' Some-college': 3,
    ' Bachelors': 3,
    ' Doctorate': 4,
    ' Prof-school': 4,
    ' Masters': 4
}

```

- 其余属性不做特殊处理

预处理结果（部分）：

	age	workclass	education	maritalStatus	occupation	...	capitalGain	capitalLoss	hoursPerWeek	nativeCountry	income
0	1	2	3	4	4	...	0	0	2	19	0
1	2	1	3	1	3	...	0	0	0	19	0
2	1	0	0	2	12	...	0	0	2	19	0
3	2	0	0	1	12	...	0	0	2	19	0
4	1	0	3	1	2	...	0	0	2	17	0

  

[5 rows x 13 columns]											
	age	workclass	education	maritalStatus	occupation	...	capitalGain	capitalLoss	hoursPerWeek	nativeCountry	income
0	1	0	0	4	13	...	0	0	2	19	0
1	1	0	0	1	5	...	0	0	2	19	0
2	1	2	2	1	7	...	0	0	2	19	1
3	2	0	3	1	13	...	0	0	2	19	1
5	1	0	0	4	10	...	0	0	1	19	0

  

[5 rows x 13 columns]											
-----------------------	--	--	--	--	--	--	--	--	--	--	--

## 决策树算法

决策树算法选择：CART

### 1. 问题假设：

假设数据集  $D$  中有  $n$  个样本，样本属于  $k$  个属性  $\{C_1, C_2, \dots, C_k\}$ ，每个属性取值个数为  $N_1, N_2, \dots, N_k$ 。对于每一个属性  $C_i$ ， $V_j$ ，我们希望计算其基尼系数  $Gini(C_i)$ ，然后选择最大的属性  $C_i$  和值  $V_j$  作为当前节点的判别依据，其输出类别  $Y \in \{y_1, y_2\}$ ， $y_1, y_2$  分别为个人收入  $\geq 50K$  和  $< 50K$

### 1. 总体信息熵：

当前节点  $father$  的基尼系数  $Gini(father)$  为：

$$Gini(father) = 2p(1 - p)$$

其中， $p$  是类别  $y_1$  在节点  $father$  中的样本比例。

### 2. 选取一个类别 $C_i$ 以及其对应的一个取值 $V$ ：

每次选取一个属性  $C_i$  来计算划分后数据的基尼系数。将当前节点  $father$  按照属于属性  $C_i$  是否等于  $V$  划分为2个子集：

- $father_V$  和  $father_{-V}$

### 3. 划分后基尼系数计算：

根据类别  $C_i$ ， $V$  的划分计算在  $Gini(father|C_i)$ ，即按属性  $C_i$  划分后的基尼系数：

$$Gini(father|C_i) = \frac{|father_V|}{|father|} Gini(father_V) + \frac{|father_{-V}|}{|father|} Gini(father_{-V})$$

### 4. impurity reduction：

衡量划分优劣：

$$\Delta G = G(father) - Gini(father|C_i)$$

### 5. 选取impurity reduction最大的类别以及取值：

计算所有类别  $C_1, C_2, \dots, C_k$  的所有取值  $V$  下对应的划分优劣，选择  $\Delta G$  最大的类别  $C_{best}$ ， $V_{best}$  作为当前节点的判别依据：

$$C_{best}, V_{best} = \arg \max_i \Delta(C_i, V_j)$$

### 6. 继续递归：

根据选中的  $C_{best}$ ， $V_{best}$  将数据集划分为子集，并递归执行以上步骤，直到满足终止条件。

- 为了以防决策树过拟合，这里参数设置最大深度为10，节点样本数量最少为1，节点内基尼系数最小为0.02

```
import numpy as np
import pandas as pd
import random
cart_config = {
    'max_depth': 10,
    'min_samples_split': 1,
    'min_gini': 0.02
}
unused_splits = set() # 记录已经使用过的切分点，避免重复使用

class Node:
    def __init__(self, feature_index, feature_value, left, right, label):
        self.feature_index = feature_index
        self.feature_value = feature_value
        self.left = left
        self.right = right
        self.label = label

    def predict(self, x):
        if self.label is not None:
            return self.label
        if x[self.feature_index] == self.feature_value:
            return self.left.predict(x)
        else:
            return self.right.predict(x)

    def __str__(self):
        return 'feature_index: %d, feature_value: %d, label: %d' %
(self.feature_index, self.feature_value, self.label)

def calc_gini(y):
    """
    计算数据集的基尼指数
    :param x: 当前节点所包含数据
    :return: 基尼指数
    """
    n = len(y)
    if n == 0:
        return 0
    m = y.sum()
    prob = m / n
    gini = 2 * prob * (1 - prob)
    return gini

def split_dataset(X, y, feature_index, feature_value):
    """
    按照给定的列划分数据集
    :param x: 当前节点所包含数据
    :param index: 指定特征的列索引
    :param value: 指定特征的值
    :return: 切分后的数据集
    """
```

```

left = X[X[:, feature_index] == feature_value]
right = X[X[:, feature_index] != feature_value]
left_labels = y[X[:, feature_index] == feature_value]
right_labels = y[X[:, feature_index] != feature_value]
return left, right, left_labels, right_labels

def choose_best_feature_to_split(X, y):
    """
    选择最好的特征进行分裂
    :param X: 当前节点数据
    :return: best_value:(分裂特征的index, 特征的值), best_df:(分裂后的左右子树数据集),
    best_gain:(选择该属性分裂的最大信息增益)
    """
    best_gini = 1
    best_feature = -1
    best_split = None
    best_value = -1
    n = X.shape[0]
    for (i, j) in unused_splits:
        left, right, left_labels, right_labels = split_dataset(X, y, i, j)
        gini = left.shape[0]/n * calc_gini(left_labels) + right.shape[0]/n *
        calc_gini(right_labels)
        if gini < best_gini:
            best_gini = gini
            best_feature = i
            best_value = j
            best_split = (left, right, left_labels, right_labels)
    return best_feature, best_value, best_split, best_gini

def build_decision_tree(X, y, depth, flags):
    """
    构建CART树
    :param X: 数据集
    :param y: 标签集
    :param depth: 当前深度
    :return: CART树
    """
    if(len(np.unique(y)) == 1):
        return Node(None, None, None, None, np.argmax(np.bincount(y)))
    if depth >= cart_config['max_depth']:
        return Node(None, None, None, None, np.argmax(np.bincount(y)))
    if y.shape[0] <= cart_config['min_samples_split']:
        return Node(None, None, None, None, np.argmax(np.bincount(y)))
    gini = calc_gini(y)
    if(gini <= cart_config['min_gini']):
        return Node(None, None, None, None, np.argmax(np.bincount(y)))
    if(len(unused_splits) == 0):
        return Node(None, None, None, None, np.argmax(np.bincount(y)))
    best_feature, best_value, best_split, best_gini =
    choose_best_feature_to_split(X, y)
    if(best_gini >= gini):
        return Node(None, None, None, None, np.argmax(np.bincount(y)))
    node = Node(best_feature, best_value, None, None, None)

```

```

        node.left = build_decision_tree(best_split[0], best_split[2], depth + 1,
flags)
        node.right = build_decision_tree(best_split[1], best_split[3], depth + 1,
flags)
    return node

def save_decision_tree(cart):
    """
    决策树的存储
    :param cart: 训练好的决策树
    :return: void
    """
    np.save('cart.npy', cart)

def load_decision_tree():
    """
    决策树的加载
    :return: 保存的决策树
    """

    cart = np.load('cart.npy', allow_pickle=True)
    return cart.item()

if __name__ == "__main__":
    train_data = np.loadtxt('train_adult_pro.csv', delimiter=',', skiprows=1)
    X,y = train_data[:, :-1], train_data[:, -1]
    y = y.astype(int)
    for i in range(X.shape[0]):
        for j in range(X.shape[1]):
            if((j,X[i][j]) in unused_splits):
                continue
            else:
                unused_splits.add((j,int(X[i][j])))
    cart = build_decision_tree(X, y, 0, unused_splits)
    save_decision_tree(cart)
    cart = load_decision_tree()
    cnt = 0
    test_data = np.loadtxt('test_adult_pro.csv', delimiter=',', skiprows=1)
    for i in range(X.shape[0]):
        if(cart.predict(X[i]) == y[i]):
            cnt += 1
    print(f"test on train data:{cnt/X.shape[0]}")
    cnt = 0
    X_test, y_test = test_data[:, :-1], test_data[:, -1]
    y_test = y_test.astype(int)
    for i in range(X_test.shape[0]):
        if(cart.predict(X_test[i]) == y_test[i]):
            cnt += 1
    print(f"test on test data:{cnt/X_test.shape[0]}")

```

## 测试结果以及准确率

---

- 训练集上准确率: 0.8361953612845674
- 测试集上准确率: 0.8305024769992922

```
PS C:\Users\HiroX\OneDrive\Desktop\bigdata\lab2> python -u "c:\Users\HiroX\OneDrive\Desktop\bigdata\lab2\decision_tree.py"
test on train data:0.8361953612845674
test on test data:0.8305024769992922
PS C:\Users\HiroX\OneDrive\Desktop\bigdata\lab2> █
```

对比sklearn结果:

```
PS C:\Users\HiroX\OneDrive\Desktop\bigdata\lab2> python -u "c:\Users\HiroX\OneDrive\Desktop\bigdata\lab2\by_sklearn.py"
sklearn test on train data: 0.836604222420458
sklearn test on test data: 0.8273885350318472
█
```