

# 决策树机器学习实验——171250574杨逸存

## 实验简介

- 本实验使用python语言及numpy、pandas和matplotlib等python库实现了ID3决策树算法
- 实验数据为给定隐形眼镜的小量数据集，构造决策树预测患者佩戴隐形眼镜的类型，并通过Matplotlib绘制树图形。

## 算法/代码详解

### ID3决策树算法

- ID3算法通过特征的信息增益 $\text{Gain}(D, A) = H(D) - H(D|A)$   
其中 $H(D|A) = \sum p_i * H(D|A=a_i)$ ,  $H(D) = -\sum p_i * \log(p_i)$  来实现最有特征的选择
- 算法输入：训练数据集D，特征集A，阈值 $\epsilon$ ；输出：决策树T
  1. 若当前训练集D中所属实例属于同一类或当前特征集A为空，则返回单节点，节点值为实例所属类或实例数最多的一个类。
  2. 否则，计算A中个特征对D的信息增益，选择信息增益最大的特征 $A_k$ 。
  3. 如果 $A_k$ 的信息增益也小于阈值 $\epsilon$ ，则同样返回值为实例数最多的类的一个单节点。
  4. 否则，对 $A_k$ 的每一种可能值 $a_i$ ，依 $A_k=a_i$ 将D分割为若干非空子集 $D_i$ 。
  5. 对每个 $D_i$ ，以 $A - \{A_k\}$ 作为特征集合，递归调用Step1~4，返回的子树作为原树的一个分支，并返回原树T；

### 核心代码

- 决策树生成主方法解析

```
1 decision_tree = generate_decision_tree(data_set, eps=0.0001)
```

使用dict字典数据结构模拟树，其中子树为嵌套字典对象，字典键为节点值或连接分支值，叶子节点为最终分类对象（具体见实验结果分析）

前面的if分支主要解决递归结束特殊条件，后半部分的for循环中递归调用generate\_decision\_tree，返回的子树作为当前节点特征的一个具体取值分支下的子树。

```
1 def generate_decision_tree(data_set: pd.DataFrame, eps: float):
2     '''
3     生成决策树，迭代调用
4     :param data_set: 当前数据集
5     :param eps: e信息增益阈值
6     :return: 字典（子树）或单个值（叶子）
7     '''
8     # 分类结果列表
```

```

9     class_list = np.array(data_set['class'])
10    # 如果所有实例属于同一个类
11    if sum(class_list == class_list[0]) == len(class_list):
12        return class_list[0]
13    # 如果标签集为空（只剩class列），则为单节点树，把实例数最大的类作为
    此节点标记类
14    if len(data_set.columns.values) == 1:
15        return get_most_common_class(data_set)
16    # 获取信息增益最大的特征及其增益
17    highest_gain_feature, highest_gain =
    get_feature_with_highest_gain(data_set)
18    # 增益小于ε，单一节点，返回实例数最大的类
19    if highest_gain < eps:
20        return get_most_common_class(data_set)
21    # 构建树
22    decision_tree_dict = {highest_gain_feature: {}}
23    # 对每个最高增益特征的取值进行分割数据集，并进行递归调用生成树
24    feature_values = set(data_set[highest_gain_feature])
25    for one_value in feature_values:
26        # 分割Di
27        divided_data_set = data_set[data_set[highest_gain_feature]
    == one_value]
28        # 去除列，A = A - {Ak}
29        divided_data_set =
    divided_data_set.drop(labels=highest_gain_feature, axis=1)
30        # 生成子树
31        decision_tree_dict[highest_gain_feature][one_value] =
    generate_decision_tree(divided_data_set, eps)
32    return decision_tree_dict

```

- 信息熵、条件熵计算

信息熵 $H(D) = -\sum p_i \log(p_i)$

```

1  def H(data_set: pd.DataFrame):
2      '''
3      计算经验熵H(D)
4      :param data_set: 数据集D
5      :return:
6      '''
7      class_list = data_set['class']
8      class_values = set(class_list)
9      H = 0.0
10     for one_value in class_values:
11         this_value_list = data_set[data_set['class'] == one_value]
12         p_i = float(len(this_value_list)) / float(len(data_set))
13         H -= p_i * np.log2(p_i) # H(X) = -Σpi*log(pi), 取2为底的对数
14     return H

```

在条件熵计算中，传入特征取某一具体值的子集合，调用信息熵计算方法，简化了逻辑计算代码

```

1 def H_condition(data_set: pd.DataFrame, feature: str):
2     '''
3     计算经验条件熵H(D|A),  $H(D|A) = \sum p_i * H(D|A=a_i)$ 
4     :param data_set: 数据集D
5     :param feature: 特征A
6     :return:
7     '''
8     class_list = data_set['class']
9     class_values = set(class_list)
10    feature_values = set(data_set[feature])
11    H_con = 0.0
12    for one_feature_value in feature_values:
13        this_value_list = data_set[data_set[feature] ==
14        one_feature_value]
15        pi = float(len(this_value_list)) / float(len(data_set))
16        H_con_this_value = H(this_value_list)
17        H_con += pi * H_con_this_value
18    return H_con

```

## 实验结果分析

- 实验数据为助教给的lenses.txt文件，有24个样本数据量
- 生成的决策树：

```

1  {'tearRate':
2      {'normal':
3          {'astigmatic':
4              {'no':
5                  {'prescript':
6                      {'hyper': 'soft',
7                       'myope':
8                           {'age':
9                               {'pre': 'soft',
10                                'presbyopic': 'noLenses',
11                                 'young': 'soft'}}}},
12                  'yes':
13                      {'prescript':
14                          {'hyper':
15                              {'age':
16                                  {'pre': 'noLenses',
17                                   'presbyopic': 'noLenses',
18                                    'young': 'hard'}}},
19                          'myope': 'hard'}}}},
20                  'reduced': 'noLenses'}
21  }

```

- 调用TreePlotter.py，修改一下figure绘图参数，得到以下树图形：

