

AI and Machine Learning, Homework7

Author: Ying Yiwen
Number: 12210159

Contents

1	Introduction	1
2	Result	3

1 Introduction

We need to compute entropy, conditional entropy and information gain to build decision tree.

Entropy is an important measure of the uncertainty or purity of a data set. The higher the entropy value, the more scattered the samples in the dataset and the less pure the categories. The lower the entropy, the more concentrated the samples in the dataset and the purer the category.

$$H(X) = - \sum_{x \in X} p(x) \log_2 p(x) \quad (1)$$

```
1 count_y = np.bincount(y) # Count the number of each output label
2 prob_y = count_y[np.nonzero(count_y)] / y.size # Compute the probability of each
   output label
3 entropy_y = -np.sum(prob_y * np.log2(prob_y)) # Compute the entropy of output
```

$$H(Y|X) = - \sum_{x \in X} \sum_{y \in Y} p(x, y) \log_2 p(y|x) \quad (2)$$

```
1 for v in feature_values:
2     y_sub = y[feature == v]
3     prob_y_sub = y_sub.size / y.size
4     h += prob_y_sub * self._entropy(y_sub)
```

In the construction of decision trees, we use **information gain** to select the best partition features, that is, to select the features with the maximum information gain to split the data set. We want to choose a feature that minimizes the uncertainty of the dataset and makes the distribution of categories in the dataset consistent.

$$IG(Y|X) = H(Y) - H(Y|X) \quad (3)$$

```
1 ig_feature = self._entropy(y) - self._conditional_entropy(feature, y)
```

Then we need to choose the best attribute, with the highest information gain.

```

1 if features_list:
2     gains = np.apply_along_axis(self._information_gain, 0, X[:, features_list], y)
3     index = np.argmax(gains)
4     if gains[index] > self.gain_threshold:
5         return index

```

The construction of the decision tree usually uses a greedy algorithm, starting from the root node, gradually constructing each subnode according to the information gain of the features, and recursively dividing the dataset into purer and purer subsets by recursively dividing the dataset until the stopping conditions are met (e.g., the maximum depth is reached, the number of samples is less than a certain threshold, all features have been used, the information gain is below the threshold, and the class purity reaches the expectation). Sometimes, after building a complete tree, we prune to avoid overfitting and reduce the complexity of the tree to improve the generalization ability of the model.

```

1 # Divide the training set according to this selected feature
2 # Then use the subset of training examples in each branch to create a sub-tree
3 feature_values = np.unique(X[:, node.feature_index])
4 for v in feature_values:
5     # Obtain the subset of training examples
6     idx = X[:, node.feature_index] == v
7     X_sub, y_sub = X[idx], y[idx]
8     # Build a sub-tree
9     node.children[v] = self._build_tree(X_sub, y_sub, features_list.copy())

```

The training of decision tree is to build the tree:

```

1 self.tree_ = self._build_tree(X_train, y_train, list(range(n)))

```

And the prediction of it is to recursively find the value:

```

1 node = self.tree_
2 while node.children:
3     child = node.children.get(x[node.feature_index])
4     if not child:
5         break
6     node = child
7 return node.value

```

Thus, the main logic of this homework is:

```

1 # Load the dataset
2 data = np.loadtxt('../dataset/lenses/lenses.data', dtype=int)
3
4 # Features and labels
5 X = data[:, 1:-1]
6 y = data[:, -1]
7
8 # Split the data into training and test sets
9 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.1)
10
11 # Create and train the decision tree
12 dt01 = DT.DecisionTree()
13 dt01.train(X_train, y_train)
14
15 # Print the tree

```

```

16 print(dt01.tree_)
17
18 # Use the trained tree to make predictions on the test set
19 y_pred = dt01.predict(X_test)
20
21 # Calculate accuracy
22 accuracy = np.sum(y_pred == y_test) / len(y_test)
23 print(f"Accuracy: {accuracy}")
24
25 # Optionally, you can visualize the decision tree
26 features_dict = {
27     0: { 'name': 'age', 'value_names': {1: 'young', 2: 'pre-presbyopic', 3: 'presbyopic' } },
28     1: { 'name': 'prescript', 'value_names': {1: 'myope', 2: 'hypermetrope' } },
29     2: { 'name': 'astigmatic', 'value_names': {1: 'no', 2: 'yes' } },
30     3: { 'name': 'tear rate', 'value_names': {1: 'reduced', 2: 'normal' } },
31 }
32
33 label_dict = {
34     1: 'hard',
35     2: 'soft',
36     3: 'no_lenses',
37 }
38
39 dtp = DecisionTreePlotter(dt01.tree_, feature_names=features_dict, label_names=
    label_dict)
40 dtp.plot()

```

2 Result

By splitting the dataset into training set and testing set (9:1), we can get the performance of the model.

Accuracy: 1.0

The tree is like:

Internal node <3>:

1 -> Leaf node (3)

2 -> **Internal node <2>:**

1 -> **Internal node <0>:**

1 -> Leaf node (2)

2 -> Leaf node (2)

3 -> **Internal node <1>:**

1 -> Leaf node (3)

2 -> Leaf node (2)

2 -> **Internal node <1>:**

1 -> Leaf node (1)

2 -> **Internal node <0>:**

1 -> Leaf node (1)

2 -> Leaf node (3)

3 -> Leaf node (3)

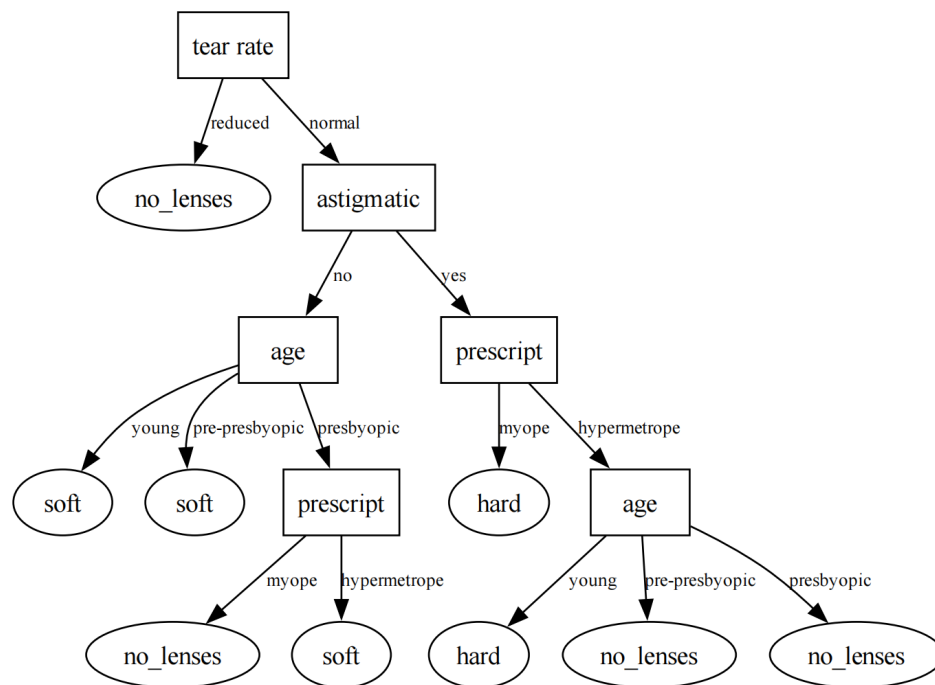


Fig. 1: Tree Structure