

CMSC 5724 Project #1: Decision Tree

1155215526 Ding Luxiao

1155215527 Chen Chen

1155215555 Chen Yubin

1155215807 Zhang Jiarui

Abstract

In this report, we made a predictive analysis using the Adult dataset. We employed Hung's algorithm to construct a decision tree, leveraging the GINI index as the criterion for optimal splits. The dataset was processed through a series of utility functions designed for data reading, frequency counting, label encoding, and model evaluation. Key optimizations reduced the time complexity for calculating ordered features to $O(n \log n)$, enhancing efficiency in identifying the best split points. The experimental setup involved training on a dataset of 30,162 records and evaluating on a separate set of 15,060 records, with specified parameters for minimum samples per split and maximum depth 10. Results indicated an accuracy of 86.72% on the training set and 85.60% on the evaluation set, with precision, recall, and F1 scores illustrating the model's performance. The decision tree's structure was documented, providing insights into the predictive relationships within the dataset.

Dataset

We will use the Adult dataset to predict whether an individual's annual income exceeds \$50K/year based on census data, also known as the "Census Income" dataset. The training set is `adult.data`, and the evaluation set is `adult.test`.

Its description is available <https://archive.ics.uci.edu/dataset/2/adult>.

It can be downloaded in <https://archive.ics.uci.edu/static/public/2/adult.zip>.

Algorithm

Using the optimization logic from ex1 Problem 3, we reduced the time complexity for calculating ordered features to $O(n \log n)$. The algorithm below is the original Hunt's algorithm.

Hunt's Algorithm:

/ S is the training set; the function returns the root of a decision tree */*

1. if all the objects in S belong to the same class
2. return a leaf node with the value of this class

3. if (all the objects in S have the same attribute values) or ($|S|$ is too small)
4. return a leaf node whose class value is the majority one in S
5. find the “best” split attribute A^* and predicate P^* using GINI
6. $S_1 \leftarrow$ the set of objects in R satisfying P^* ; $S_2 \leftarrow S \setminus S_1$
7. $u_1 \leftarrow \text{Hunt}(R_1)$; $u_2 \leftarrow \text{Hunt}(R_2)$
8. create a root u with left child u_1 and right child u_2
9. set $A_u \leftarrow A^*$ and $P_u \leftarrow P^*$
10. return u

We’ve modified the optimized Hunt’s algorithm as below:

1. Define the Dataset: Let S be the set of records in the table. Each record r is in the form (rA, rB) , representing its values on attribute A and class label B , respectively.
2. Sort the Dataset: Sort the dataset S in ascending order by the values of attribute A to efficiently calculate the GINI index for each possible split point.
3. Initialize a Counter: Set up a counter c to track the number of positive (yes) records up to and including the current split point.
4. Traverse the Records: After sorting, traverse the records of S in ascending order of A .
5. Update the Counter:
If the class label B of the record r is positive (yes), increment the counter c by 1.
For each record r , consider its A attribute value $r.A$ as a potential split point and set $cly(a)$ to the current value of the counter c , representing the number of positive classes records up to and including the current record r .
6. Calculate the GINI Index for Each Split: Using the counts obtained in the previous steps, calculate the GINI index for each possible split point. The GINI index measures the impurity of a dataset; the lower the value, the higher the purity of the resulting subsets.
7. Select the Best Split Point: Go through all possible split points and choose the one that results in the smallest GINI index as the best split point.

The key to this algorithm is performing a single sorting operation ($O(n \log n)$ time complexity) and a single traversal ($O(n)$ time complexity), allowing efficient calculation of the GINI index for all possible split points, thus identifying the best split point. This method avoids calculating the GINI index for each split point individually, significantly improving efficiency.

Experimental Process

Utility Functions

Includes functions for data reading, frequency counting, label encoding, and model evaluation.

Can be referred from Table 1.

Table 1: Functions/Classes presented as utilities for data dealing

Function/Class name	Usage
<code>read_data(path, tHead, paraNumeric, paraRemove, isTest)</code>	Reads data files and converts each line into a dictionary format, stored in a list.
<code>count_frequencies(data)</code>	Counts the frequency of each attribute in the data
<code>get_sorted_list_value(data, key, descending)</code>	Sorts attribute frequencies by value
<code>get_sorted_list_key(data, key, descending)</code>	Sorts attribute frequencies by key
<code>label_encode(data, value_dict)</code>	Encodes nonnumeric labels using numeric substitutes
<code>evaluate_model(labels, predictions)</code>	Computes model evaluation metrics, including accuracy, precision, recall, and F1 score
<code>print_model_evaluation_result(data, predictions)</code>	Print the evaluation result of the model(decision tree) based on adult dataset
<code>print_all_results_of_eval(data, predict)</code>	Print all evaluation data's actual income and predicted income
<code>TreePrinter(tree, paraLabel)</code>	Used to print the structure of the decision tree

Decision Tree Construction

Uses the GINI index as the splitting criterion. Below is a detailed explanation of each part:

1. **Data Preparation:** Use the `split_dataset` method to divide the dataset into two subsets based on attributes and split points.

```

1. def split_dataset(self, features, labels, feature_index, split) -> tuple[list]:
2.     left_features, left_labels = [], []
3.     right_features, right_labels = [], []
4.     for i in range(len(features)):
5.         if feature_index in self.numeric_feature_index: # For numeric
6.             if features[i][feature_index] <= split:
7.                 left_features.append(features[i])
8.                 left_labels.append(labels[i])
9.             else:
10.                right_features.append(features[i])
11.                right_labels.append(labels[i])
12.         else: # For non-numeric
13.             if features[i][feature_index] == split:
14.                 left_features.append(features[i])
15.                 left_labels.append(labels[i])
16.             else:
17.                 right_features.append(features[i])
18.                 right_labels.append(labels[i])

```

```
19.     return left_features, left_labels, right_features, right_labels
```

2. **GINI Index Calculation:** Use the gini and gini_split methods to calculate the GINI index for labels, assessing the quality of the splits.

```
1. def gini(self, labels: list) -> float:
```

```
2.     total = len(labels)
```

```
3.     if total == 0:
```

```
4.         return 0
```

```
5.     return 1 - ((labels.count(0) / total) ** 2 + (labels.count(1) / total) ** 2)
```

3. **Best Split Point Selection:** The find_best_split_gini method traverses all features and possible split points to find the feature and split value that minimize the GINI index.

```
1. def find_best_split_gini(self, features, labels) -> tuple:
```

```
2.     best_feature_index = None
```

```
3.     best_split = None
```

```
4.     best_gini = float('inf')
```

```
5.     n_features = len(features[0])
```

```
6.     # optimized version
```

```
7.     for feature_index in range(n_features):
```

```
8.         # find numeric split for numeric data
```

```
9.         if feature_index in self.numeric_feature_index:
```

```
10.            # sort data based on feature
```

```
11.            feature_and_labels = sorted(zip(features, labels), key=lambda x: x[0][feature_index])
```

```
12.            sorted_features = [f for f, _ in feature_and_labels]
```

```
13.            sorted_labels = [l for _, l in feature_and_labels]
```

```
14.            # count yes and no of two subsets
```

```
15.            left_yes_count = 0
```

```
16.            left_no_count = 0
```

```
17.            right_yes_count = sorted_labels.count(1)
```

```
18.            right_no_count = len(sorted_labels) - right_yes_count
```

```
19.            # foreach every possible split and update count number, calculate GINI
```

```
20.            for i in range(1, len(sorted_features)):
```

```
21.                if sorted_labels[i - 1] == 1:
```

```
22.                    left_yes_count += 1
```

```
23.                    right_yes_count -= 1
```

```
24.                else:
```

```
25.                    left_no_count += 1
```

```
26.                    right_no_count -= 1
```

```
27.                # ensure split point available
```

```
28.                if sorted_features[i][feature_index] == sorted_features[i - 1][feature_index]:
```

```
29.                    continue
```

```
30.                # calculate GINI of present split
```

```

31.         total_left = left_yes_count + left_no_count
32.         total_right = right_yes_count + right_no_count
33.         gini_left = 1 - (left_yes_count / total_left) ** 2 - (left_no_count /
total_left) ** 2
34.         gini_right = 1 - (right_yes_count / total_right) ** 2 - (right_no_coun
t / total_right) ** 2
35.         weighted_gini = (total_left / len(labels)) * gini_left + (total_right
/ len(labels)) * gini_right
36.         # update best split
37.         if weighted_gini < best_gini:
38.             best_gini = weighted_gini
39.             best_feature_index = feature_index
40.             best_split = (sorted_features[i - 1][feature_index] + sorted_feature
s[i][feature_index]) / 2
41.         # find non-numeric split for non-numeric feature
42.         else:
43.             unique_values = set([row[feature_index] for row in features])
44.             for value in unique_values:
45.                 _, left_labels, _, right_labels = self.split_dataset(features, labels,
feature_index, value)
46.                 # pass split if size is 0
47.                 if len(left_labels) == 0 or len(right_labels) == 0:
48.                     continue
49.                 # calculate GINI of present split
50.                 current_gini = self.gini_split(left_labels, right_labels)
51.                 # update best split
52.                 if current_gini < best_gini:
53.                     best_gini = current_gini
54.                     best_feature_index = feature_index
55.                     best_split = value
56.         return best_feature_index, best_split
57.
58.         # traditional version
59.         # for feature_index in range(n_features):
60.         #     splits = set([row[feature_index] for row in features])
61.         #     for split in splits:
62.         #         _, left_labels, _, right_labels = self.split_dataset(features, labels,
feature_index, split)
63.         #         if len(left_labels) == 0 or len(right_labels) == 0:
64.         #             continue
65.         #         current_gini = self.gini_split(left_labels, right_labels)
66.         #         if current_gini < best_gini:
67.         #             best_gini = current_gini
68.         #             best_feature_index = feature_index

```

```

69.     #         best_split = split
70.     # return best_feature_index, best_split

```

4. **Tree Construction:** The fit method recursively builds the decision tree until stop conditions are met (e.g., all labels are the same, sample size is less than the minimum split size, or maximum depth is reached).

```

1. def fit(self, features: list, labels: list, depth=0) -> tuple:
2.     n_samples = len(labels)
3.     # return the label if all labels the same
4.     if len(set(labels)) == 1:
5.         return labels[0]
6.     # check for hyper
7.     if len(features) == 0 or \
8.         n_samples < self.min_samples_split or \
9.         (self.max_depth is not None and depth >= self.max_depth):
10.        # ensure labels are not empty or return None
11.        return max(set(labels), key=labels.count) if labels else None
12.    # find the best split of dataset
13.    best_feature_index, best_split = self.find_best_split_gini(features, labels)
14.    # return majority label if best split not found
15.    if best_feature_index is None:
16.        return max(set(labels), key=labels.count) if labels else None
17.    # split based on best split point
18.    left_features, left_labels, right_features, right_labels = self.split_dataset(
19.        features, labels, best_feature_index,
20.        best_split)
21.    # in case for empty subset
22.    if len(left_labels) == 0 or len(right_labels) == 0:
23.        return max(set(labels), key=labels.count) if labels else None
24.    # generate left tree and right tree
25.    left_subtree = self.fit(left_features, left_labels, depth + 1)
26.    right_subtree = self.fit(right_features, right_labels, depth + 1)
27.    return (best_feature_index, best_split, left_subtree, right_subtree)

```

5. **Prediction:** The predict method uses the constructed decision tree to make predictions on new datasets.

```

1. def _predict(self, sample, tree):
2.     if not isinstance(tree, tuple):
3.         return tree
4.     feature_index, split, left_subtree, right_subtree = tree
5.     if feature_index in self.numeric_feature_index:
6.         if sample[feature_index] <= split: # For numeric
7.             return self._predict(sample, left_subtree)
8.         else:
9.             return self._predict(sample, right_subtree)

```

```

10.     else: # For non-numeric
11.         if sample[feature_index] == split:
12.             return self._predict(sample, left_subtree)
13.         else:
14.             return self._predict(sample, right_subtree)
15.
16. def predict(self, features: list[list]) -> list:
17.     return [self._predict(sample, self.tree) for sample in features]

```

Experimental Result

We've conducted evaluations on the training and evaluation sets with training size: 30,162 and evaluation size: 15,060, using different MIN_SAMPLE_SPLIT(MSS) and MAX_DEPTH(MD) value pairs. Results for training sets are experimented and shown in Table 2.

Table 2: Experimental Result of Different Hyperparameter Pairs

MSS \ MD	5	10	15	20	25
5(train)	0.84275	0.84275	0.84275	0.84258	0.84258
5(evaluate)	0.84064	0.84064	0.84064	0.84064	0.84064
10(train)	0.86778	0.86718	0.86652	0.86546	0.86493
10(evaluate)	0.85564	0.85604	0.85578	0.85558	0.85604
15(train)	0.89643	0.89079	0.88631	0.88373	0.88210
15(evaluate)	0.84017	0.84243	0.84420	0.84502	0.84542
20(train)	0.92726	0.91423	0.90624	0.90109	0.89679
20(evaluate)	0.82324	0.82875	0.83353	0.83552	0.83772
25(train)	0.95103	0.93057	0.91874	0.91121	0.90587
25(evaluate)	0.81308	0.82131	0.82709	0.82968	0.83227

From the result mentioned in Table 2, we can easily find out that when MAX_DEPTH is larger than 15, the decision tree is more probable of overfit while the expansion of MIN_SAMPLE_SPLIT can mitigate overfit to some extent. For the result with best evaluation accuracy, having set MIN_SAMPLE_SPLIT=10/25 and MAX_DEPTH=10, the other evaluation scores are shown in Table 3 and Table 4.

Table 3: Evaluation Result of MSS=10 and MD=10

Dataset	Accuracy	Precision	Recall	F1 Score
Training Set	0.86718	0.79779	0.62480	0.70078
Evaluation Set	0.85604	0.76505	0.59757	0.67102

Table 4: Evaluation Result of MSS=25 and MD=10

Dataset	Accuracy	Precision	Recall	F1 Score
Training Set	0.86718	0.79171	0.62067	0.69583
Evaluation Set	0.85604	0.76560	0.59676	0.67072

Besides, we've kept the complete structure of the decision tree and the evaluation set prediction with MSS=10 and MD=10 in the file res.txt.

Appendix:

1. GitHub repo: <https://github.com/Chenchen2001/5724Proj1>

2. Part of Output in res.txt Demo

```
traing data size: 30162    evaluation data set: 15060
min_samples_split=10, max_depth=10
===== TRAINING MODEL on Training SET =====
===== TRAINING MODEL FINISHED =====
===== EVALUATION on Training SET =====
Accuracy: 0.86718
Precision: 0.79779
Recall: 0.62480
F1 Score: 0.70078
===== EVALUATION on Evaluation SET =====
Accuracy: 0.85604
Precision: 0.76505
Recall: 0.59757
F1 Score: 0.67102
===== DECISION TREE RESULT =====
|- martial-status is Married-civ-spouse
L   |- education-num <= 12.5
L   L   |- capital-gain <= 5095.5
L   L   L   |- education-num <= 8.5
L   L   L   L   |- capital-loss <= 1791.5
L   L   L   L   L   |- age <= 36.5
L   L   L   L   L   L   |- occupation is Tech-support
.....
R   R   R   L   R   |- INCOME >50k
R   R   R   R   |- capital-gain <= 30961.5
R   R   R   R   L   |- INCOME >50k
R   R   R   R   R   |- capital-gain <= 67047.0
R   R   R   R   R   L   |- INCOME <=50k
R   R   R   R   R   R   |- INCOME >50k

===== PREDICT RESULT OF EVALUATION SET =====
SAMPLE 1 : actual result - <=50k    predicted result - <=50k    predict True.
SAMPLE 2 : actual result - <=50k    predicted result - <=50k    predict True.
SAMPLE 3 : actual result - >50k     predicted result - <=50k    predict False.
SAMPLE 4 : actual result - >50k     predicted result - >50k    predict True.
SAMPLE 5 : actual result - <=50k    predicted result - <=50k    predict True.
SAMPLE 6 : actual result - >50k     predicted result - >50k    predict True.
...
```