# Artificial Intelligence Project 2
# Decision Tree

### Class 1, Team 6

### December 21, 2014

## 1   Team Members

| Name | Student ID | Job |
|------|-----------|-----|
| Fan Ziyao | 12330081 | Team leader, implementation |
| Chen Yingcong | 12330049 | Implementation |
| Chen Xiongtao | 12330040 | Modeling, implementation |
| Huang Long | 12330132 | Implementation |
| Zhang Qiuyi | 12330402 | Implementation, documentation |

## 2   Problem Description

The dataset contains the records of the United States congressional votes in 1984. There are 435 records, each containing the party (goal attribute) of the congressman/congresswoman and his/her votes on 16 bills. There are only two types of values for the goal attribute: 'democrat' and 'republican', and three types for other attributes: 'n' (can be considered as nay), 'y' (can be considered as yea), and '?' (can be considered as missing).

Given this dataset, the goal is to build a decision tree to classify a record (that is, to find out the party of this congressman/congresswoman) according to its known attributes (that is, his/her votes on these 16 bills). Both the training set and the testing set will be sampled from this dataset.

## 3   Algorithm and Implementation

For this project we implemented the C4.5 algorithm to build the decision tree. Since the dataset only have two effective values ('y' and 'n') for all attributes, there is no need to use the information gain ratio as the measure of attributes. The information gain alone is more suitable for this dataset.

To handle the missing data, when building the tree, we only consider the records with defined value for the attribute being measured. When classifying the record, we use the probability of each possible value to weigh the calculation of each branch.

The data set is split by sampling the training set with a given probability $p$. The rest of the data will be put in to the testing set.

### 3.1   Decision tree building

The algorithm for building the decision tree is described in Algorithm 1.

Since the decision tree built with this data set won't be very deep (after all, there are only 16 attributes, and the data set is fairly skewed), a recursive implementation is sufficient. Records with missing values of the chosen attribute will be ignored when calculating the information gain.

**Algorithm 1** Decision tree bulding
_____
1: **function** BUILD(*data, attributes*)
2:    **if** *data* is empty **then**
3:       **return** an empty node
4:    **end if**
5:    $gain_{best} = 0.0$, $attr_{best} = \text{None}$
6:    $value_{best} = \text{None}$, $sets_{best} = \text{None}$
7:    **for** each attribute $a$ in *attributes* **do**
8:       Divide *data* based on values of $a$ into $set1$, $set2$
9:       Calculate the information gain of spliting *data* into $set1$ and $set2$, ignoring missing values
10:       **if** The new information gain $> gain_{best}$ **then**
11:          Update $gain_{best}$, $attr_{best}$, $sets_{best}$, $value_{best}$
12:       **end if**
13:    **end for**
14:    **if** *attributes* is not empty and $gain_{best} \neq 0.0$ **then**
15:       remove $attr_{best}$ from *attributes*
16:       $left\_branch = \text{BUILD}(set1 \text{ in} sets_{best}, attributes)$
17:       $right\_branch = \text{BUILD}(set2 \text{ in} sets_{best}, attributes)$
18:       Create a node *node*
19:       $node.attr = attr_{best}$, $node.value = value_{best}$
20:       $node.left = left\_branch$, $node.right = right\_branch$, $node.count = \text{size of } data$
21:       **return** *node*
22:    **else**
23:       Create a node *node*
24:       $node.leaves = \text{count of each parties in } data$
25:       $node.count = \text{size of } data$
26:       **return** *node*
27:    **end if**
28: **end function**
_____

Remark 1.
   The *entropy* of a attribute $V$ in the data set $D$ is defined as:

$$H(V, D) = -\sum_k P(v_k) \log_2 \frac{1}{P(v_k)}$$

   where $v_k$ are values of $V$ in $D$, $P(v_k)$ is the probability of $v_k$ in $D$.

Remark 2.
   When a data set $D$ with goal attribute $V$ is split into subsets $D_k$ using attribute $A$, each with $n_k$ records, the *information gain* is defined as:

$$Gain(A) = H(V, D) - Remainder(A)$$

   where

$$Remainder(A) = \sum_k P(n_k) H(V, D_k)$$

   $P(n_k)$ is the proportion of $E_k$ in $D$.

## 3.2   Decision tree pruning

The algorithm for pruning the decision tree is described in Algorithm 2. For the same reason as the tree building, a recursive implementation is enough. Nonetheless, since the data set is quite skewed, there is seldom a need for pruning the tree.

**Algorithm 2** Decision tree pruning

```
 1: function PRUNE(tree, min_gain)
 2:     if tree.left is not a leave then
 3:         PRUNE(tree.left, min_gain)
 4:     end if
 5:     if tree.right is not a leave then
 6:         PRUNE(tree.right, min_gain)
 7:     end if
 8:     if Both tree.left and tree.right are leaves then
 9:         Temprorily combine tree.left and tree.right into new_node
10:         delta = ENTROPY(new_node) − (ENTROPY(tree.left) + ENTROPY(tree.right))/2
11:         if delta < min_gain then
12:             Replace tree with new_node
13:         end if
14:     end if
15: end function
```

## 3.3 Classification

The algorithm for classifying an observation is described in Algorithm 3. It follows all branches at any node for which a value is missing and multiply the weights along each path.

**Algorithm 3** Classification

```
 1: function CLASSIFY(tree, observation)
 2:     if tree is a leave then
 3:         return tree.leaves
 4:     end if
 5:     value = observation[tree.attr]
 6:     if value is missing then
 7:         result_left = CLASSIFY(tree.left, observation)
 8:         result_right = CLASSIFY(tree.right, observation)
 9:         weight_left = tree.left.count/tree.count
10:         weight_right = tree.left.count/tree.count
11:         for each pair of goal, count in result_left do
12:             result[goal]+ = count ∗ weight_left
13:         end for
14:         for each pair of goal, count in result_right do
15:             result[goal]+ = count ∗ weight_right
16:         end for
17:         return result
18:     else
19:         if value = tree.left.value then
20:             return CLASSIFY(tree.left, observation)
21:         else
22:             return CLASSIFY(tree.right, observation)
23:         end if
24:     end if
25: end function
```

# 4 Experiment Result

To evaluate the implementation, we sample the data with probability $p$ and use the sample as training set. The rest of the data set will be the testing set. One of the decision tree built with $p = 0.2$ is visualized in Figure 1. In addition, we take training sets sampled with various $p$ to generate the learning curve for evaluation, as shown in Figure 2. With training sets of different sizes, the precisions of cross validation fluctuate between 0.9 and 1.0, with mean $\mu = 0.9462$ and standard variation $\sigma = 0.0175$.
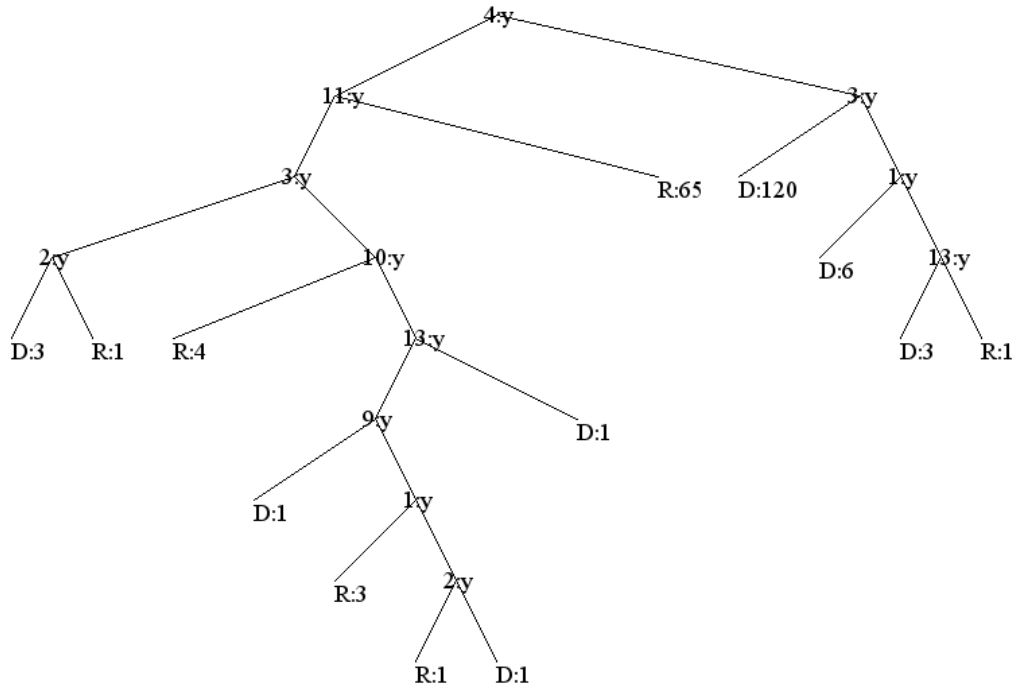
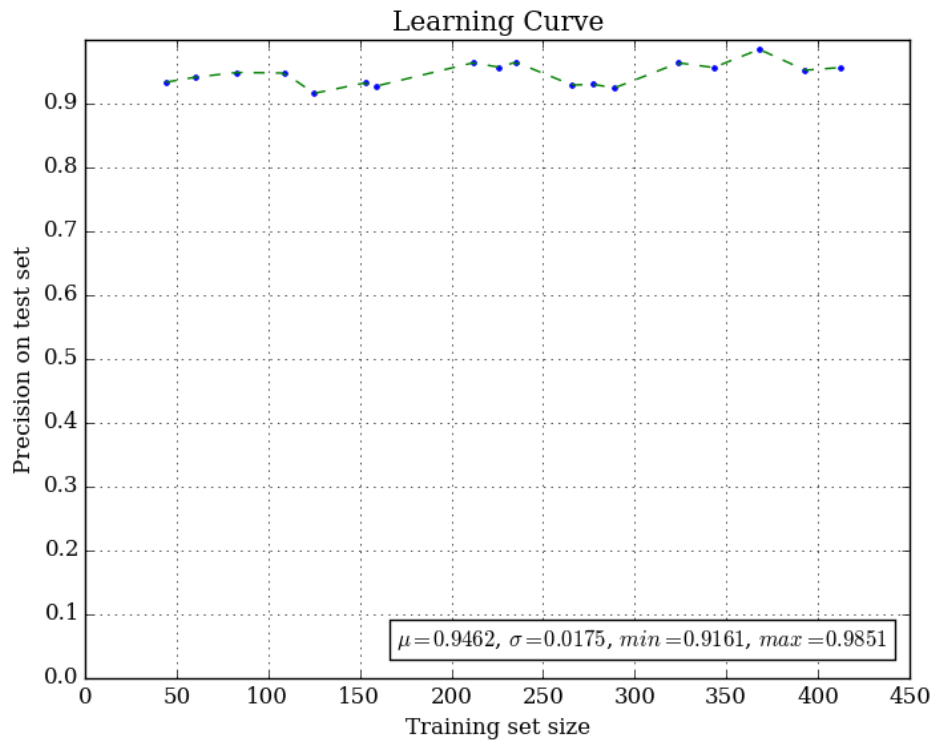Figure 1: One of the decision trees built with training set sampled with $p = 0.5$



Figure 2: Learning Curve