

Decision Trees And Training Strategies

Arnav Gupta

October 18, 2024

Contents

1	Decision Trees	1
1.1	Stopping Criteria	3
1.2	Feature Selection	3
2	Information Theory	3
2.1	Information Gain	4
3	Training Strategies	4
3.1	Final Return Value	4
3.2	Priority Queue	5
3.3	Overfitting	5

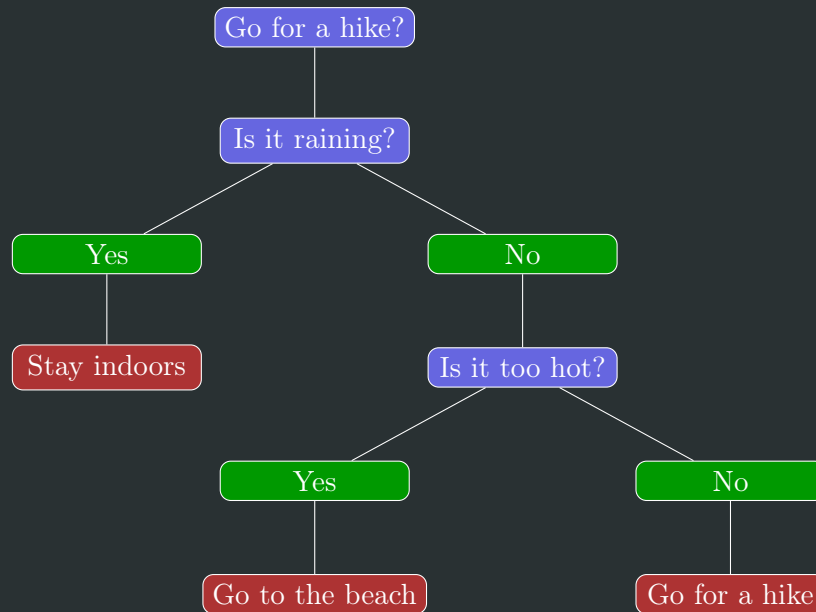
1 Decision Trees

Technique for supervised learning from discrete data:

- representation is a decision tree
- bias towards simple decision tree
- search through space of decision trees, from simple to complex

For each **decision tree**:

- **nodes**: input attributes/features
- **branches**: labeled with input feature values (can have multiple feature values on one branch)
- **leaves**: predictions for target features (point estimates)



To learn a decision tree:

1. split training data based some criteria (bias)
2. recursively solve sub-problems

Criteria for good decision trees: small, good classification (low error on training data), good generalization (low error on test data).

Algorithm 1 Decision Tree Learner

```

1: procedure DECISIONTREELEARNER( $X, Y, E$ )
2:   if stopping criteria is met then
3:     return pointEstimate( $Y, E$ )
4:   else
5:     select input feature  $X_i \in X$ 
6:     for each value  $x_i$  of  $X_i$  do
7:        $E_i \leftarrow$  all examples in  $E$  where  $X_i = x_i$ 
8:        $T_i \leftarrow$  DecisionTreeLearner( $X \setminus \{X_i\}, Y, E_i$ )
9:     end for
10:    return  $\langle X_i, T_1, \dots, T_N \rangle$ 
11:   end if
12: end procedure

```

Algorithm 2 Classify Example

```
1: procedure CLASSIFYEXAMPLE(e, X, Y, DT)
2:    $S \leftarrow DT$ 
3:   while S is an internal node of the form  $\langle X_i, T_1, \dots, T_N \rangle$  do
4:      $j \leftarrow X_i(e)$ 
5:      $S \leftarrow T_j$ 
6:   end while
7:   return S
8: end procedure
```

1.1 Stopping Criteria

Related to the final return value, depends on what must be done.

Possible stopping criteria:

- no more features
- performance on training data good enough

1.2 Feature Selection

Choose sequence of features that result in smallest tree.

In practice, split based on what gives best performance.

Heuristics for best performing feature:

- most even split
- max info gain
- GINI index

2 Information Theory

n bits can distinguish 2^n items, but can do better with probabilities.

In general, need $-\log_2 P(x)$ bits to encode x . Each symbol requires on average

$$-P(x) \log_2 P(x) \text{ bits}$$

To transmit an entire sequence distributed according to $P(x)$,

$$\sum_x -P(x) \log_2 P(x) \text{ bits}$$

bits of info per symbol are needed on average, which is the **info content** or **entropy** of the sequence.

2.1 Information Gain

Given a set E of N training examples, if the number of examples with output feature $Y = y_i$ is N_i , then

$$P(Y = y_i) = P(y_i) = \frac{N_i}{N}$$

is the point estimate.

The total info content for E is

$$I(E) = - \sum_{y_i \in Y} P(y_i) \log(P(y_i))$$

After splitting E into E_1 and E_2 based on input feature X_i , the information content is

$$I(E_{split}) = \frac{N_1}{N} I(E_1) + \frac{N_2}{N} I(E_2)$$

so the desirable X_i is the one that maximizes **info gain**:

$$I(E) - I(E_{split}) = - \sum_{y \in Y} P(y) \log P(y) + \sum_{x \in X, y \in Y} \log \frac{P(x, y)}{P(x)} \geq - \log \left(\sum_{x \in X, y \in Y} P(x, y) \frac{P(x)P(y)}{P(x, y)} \right)$$

Info gain is the reduction in uncertainty about the output feature T given the value of a certain input feature X .

Jensen's inequality: for a convex function $f(x)$, $E[f(x)] \geq f(E[x])$.

3 Training Strategies

3.1 Final Return Value

Point estimate (prediction of target features) of Y over all examples.

A point estimate could be mean, median, mode, or

$$P(Y = y_i) = \frac{N_i}{N}$$

3.2 Priority Queue

Sort leaves using a priority queue ranked by how much info can be gained with the best feature at that leaf. Always expand the leaf at the top of the queue.

Algorithm 3 Decision Tree Learner

```
1: procedure DECISIONTREELEARNER( $X, Y, E$ )
2:    $DT \leftarrow \text{pointEstimate}(Y, E)$   $\triangleright$  initial decision tree
3:    $\{X', \Delta I\} \leftarrow$  best feature and Information Gain value for  $E$ 
4:    $PQ \leftarrow \{(DT, E, X', \Delta I)\}$   $\triangleright$  priority queue of leaves ranked by  $\Delta I$ 
5:   while stopping criteria is not met do
6:      $\{S_\ell, E_\ell, X_\ell, \Delta I_\ell\} \leftarrow$  leaf at the head of  $PQ$ 
7:     for each value  $x_i$  of  $X_\ell$  do
8:        $E_i \leftarrow$  all examples in  $E_\ell$  where  $X_\ell = x_i$ 
9:        $\{X_j, \Delta I_j\} \leftarrow$  best feature and value for  $E_i$ 
10:       $T_i \leftarrow \text{pointEstimate}(Y, E_i)$ 
11:      insert  $\{T_i, E_i, X_j, \Delta I_j\}$  into  $PQ$  according to  $\Delta I_j$ 
12:    end for
13:     $S_\ell \leftarrow \langle X_\ell, T_1, \dots, T_N \rangle$ 
14:  end while
15:  return  $DT$ 
16: end procedure
```

3.3 Overfitting

When there is not enough data, the decision tree does not generalize to test data.

Methods to avoid overfitting:

- **regularization:** prefer small decision trees, so add a complexity penalty to stopping criteria
- **pseudocounts:** add data based on prior knowledge
- **cross validation**



Test set errors are caused by:

- **bias**: error due to the algorithm finding an imperfect model
 - **representation bias**: model too simple
 - **search bias**: not enough search
- **variance**: error due to lack of data
- **noise**: error due to data depending on features not modeled or because process generating data is stochastic
- **bias-variance trade-off**
 - complicated model, not enough data
 - simple model, lots of data

Capacity: ability of a model to fit a wide variety of functions (inverse of bias)

Cross Validation:

- split training data into training and validation set
- use validation set as fake test set
- optimize decision maker to perform well on validation set
- repeat with different validation sets