

Using Decision Trees for Classification

In this chapter we look at a widely-used method of constructing a model from a dataset in the form of a decision tree or (equivalently) a set of decision rules. It is often claimed that this representation of the data has the advantage compared with other approaches of being meaningful and easy to interpret.

4.1 Decision Rules and Decision Trees

In many fields, large collections of examples, possibly collected for other purposes, are readily available. Automatically generating classification rules (often called *decision rules*) for such tasks has proved to be a realistic alternative to the standard Expert System approach of eliciting the rules from experts. The British academic Donald Michie [1] reported two large applications of 2,800 and 30,000+ rules, developed using automatic techniques in only one and 9 man-years, respectively, compared with the estimated 100 and 180 man-years needed to develop the celebrated ‘conventional’ Expert Systems MYCIN and XCON.

In many (but not all) cases decision rules can conveniently be fitted together to form a tree structure of the kind shown in the following example.

4.1.1 Decision Trees: The Golf Example

A fictitious example which has been used for illustration by many authors, notably Quinlan [2], is that of a golfer who decides whether or not to play each day on the basis of the weather.

Figure 4.1 shows the results of two weeks (14 days) of observations of weather conditions and the decision on whether or not to play.

Outlook	Temp (°F)	Humidity (%)	Windy	Class
sunny	75	70	true	play
sunny	80	90	true	don't play
sunny	85	85	false	don't play
sunny	72	95	false	don't play
sunny	69	70	false	play
overcast	72	90	true	play
overcast	83	78	false	play
overcast	64	65	true	play
overcast	81	75	false	play
rain	71	80	true	don't play
rain	65	70	true	don't play
rain	75	80	false	play
rain	68	80	false	play
rain	70	96	false	play

Classes
play, don't play

Outlook
sunny, overcast, rain

Temperature
numerical value

Humidity
numerical value

Windy
true, false

Figure 4.1 Data for the Golf Example

Assuming the golfer is acting consistently, what are the rules that determine the decision whether or not to play each day? If tomorrow the values of *Outlook*, *Temperature*, *Humidity* and *Windy* were sunny, 74°F, 77% and false respectively, what would the decision be?

One way of answering this is to construct a *decision tree* such as the one shown in Figure 4.2. This is a typical example of a decision tree, which will form the topic of several chapters of this book.

In order to determine the decision (classification) for a given set of weather conditions from the decision tree, first look at the value of *Outlook*. There are three possibilities.

1. If the value of *Outlook* is sunny, next consider the value of *Humidity*. If the value is less than or equal to 75 the decision is *play*. Otherwise the decision is *don't play*.

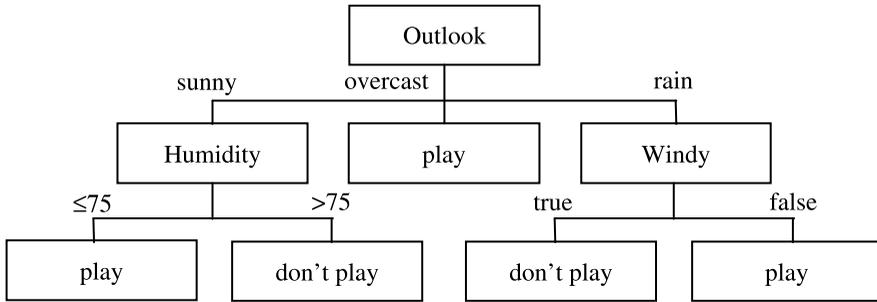


Figure 4.2 Decision Tree for the Golf Example

2. If the value of *Outlook* is overcast, the decision is *play*.
3. If the value of *Outlook* is rain, next consider the value of *Windy*. If the value is true the decision is *don't play*, otherwise the decision is *play*.

Note that the value of *Temperature* is never used.

4.1.2 Terminology

We will assume that the ‘standard formulation’ of the data given in Chapter 2 applies. There is a universe of *objects* (people, houses etc.), each of which can be described by the values of a collection of its *attributes*. Attributes with a finite (and generally fairly small) set of values, such as sunny, overcast and rain, are called *categorical*. Attributes with numerical values, such as *Temperature* and *Humidity*, are generally known as *continuous*. We will distinguish between a specially-designated categorical attribute called the *classification* and the other attribute values and will generally use the term ‘attributes’ to refer only to the latter.

Descriptions of a number of objects are held in tabular form in a *training set*. Each row of the figure comprises an *instance*, i.e. the (non-classifying) attribute values and the classification corresponding to one object.

The aim is to develop *classification rules* from the data in the training set. This is often done in the implicit form of a *decision tree*.

A decision tree is created by a process known as *splitting on the value of attributes* (or just *splitting on attributes*), i.e. testing the value of an attribute such as *Outlook* and then creating a branch for each of its possible values. In the case of continuous attributes the test is normally whether the value is ‘less than or equal to’ or ‘greater than’ a given value known as the *split value*.

The splitting process continues until each branch can be labelled with just one classification.

Decision trees have two different functions: *data compression* and *prediction*. Figure 4.2 can be regarded simply as a more compact way of representing the data in Figure 4.1. The two representations are equivalent in the sense that for each of the 14 instances the given values of the four attributes will lead to identical classifications.

However, the decision tree is more than an equivalent representation to the training set. It can be used to predict the values of other instances not in the training set, for example the one given previously where the values of the four attributes are sunny, 74, 77 and false respectively. It is easy to see from the decision tree that in this case the decision would be *don't play*. It is important to stress that this 'decision' is only a prediction, which may or may not turn out to be correct. There is no infallible way to predict the future!

So the decision tree can be viewed as not merely equivalent to the original training set but as a generalisation of it which can be used to predict the classification of other instances. These are often called *unseen instances* and a collection of them is generally known as a *test set* or an *unseen test set*, by contrast with the original training set.

4.1.3 The *degrees* Dataset

The training set shown in Figure 4.3 (taken from a fictitious university) shows the results of students for five subjects coded as SoftEng, ARIN, HCI, CSA and Project and their corresponding degree classifications, which in this simplified example are either FIRST or SECOND. There are 26 instances. What determines who is classified as FIRST or SECOND?

Figure 4.4 shows a possible decision tree corresponding to this training set. It consists of a number of *branches*, each ending with a *leaf node* labelled with one of the valid classifications, i.e. FIRST or SECOND. Each branch comprises the route from the root node (i.e. the top of the tree) to a leaf node. A node that is neither the root nor a leaf node is called an *internal node*.

We can think of the root node as corresponding to the original training set. All other nodes correspond to a subset of the training set.

At the leaf nodes each instance in the subset has the same classification. There are five leaf nodes and hence five branches.

Each branch corresponds to a *classification rule*. The five classification rules can be written in full as:

IF SoftEng = A AND Project = A THEN Class = FIRST

IF SoftEng = A AND Project = B AND ARIN = A AND CSA = A

SoftEng	ARIN	HCI	CSA	Project	Class
A	B	A	B	B	SECOND
A	B	B	B	A	FIRST
A	A	A	B	B	SECOND
B	A	A	B	B	SECOND
A	A	B	B	A	FIRST
B	A	A	B	B	SECOND
A	B	B	B	B	SECOND
A	B	B	B	B	SECOND
A	A	A	A	A	FIRST
B	A	A	B	B	SECOND
B	A	A	B	B	SECOND
A	B	B	A	B	SECOND
B	B	B	B	A	SECOND
A	A	B	A	B	FIRST
B	B	B	B	A	SECOND
A	A	B	B	B	SECOND
B	B	B	B	B	SECOND
A	A	B	A	A	FIRST
B	B	B	A	A	SECOND
B	B	A	A	B	SECOND
B	B	B	B	A	SECOND
B	A	B	A	B	SECOND
A	B	B	B	A	FIRST
A	B	A	B	B	SECOND
B	A	B	B	B	SECOND
A	B	B	B	B	SECOND

Classes
FIRST, SECOND

SoftEng
A,B

ARIN
A,B

HCI
A,B

CSA
A,B

Project
A,B

Figure 4.3 The *degrees* Dataset

THEN Class = FIRST

IF SoftEng = A AND Project = B AND ARIN = A AND CSA = B

THEN Class = SECOND

IF SoftEng = A AND Project = B AND ARIN = B

THEN Class = SECOND

IF SoftEng = B THEN Class = SECOND

The left-hand side of each rule (known as the *antecedent*) comprises a number of *terms* joined by the logical AND operator. Each term is a simple test on the value of a categorical attribute (e.g. SoftEng = A) or a continuous attribute (e.g. in Figure 4.2, Humidity > 75).

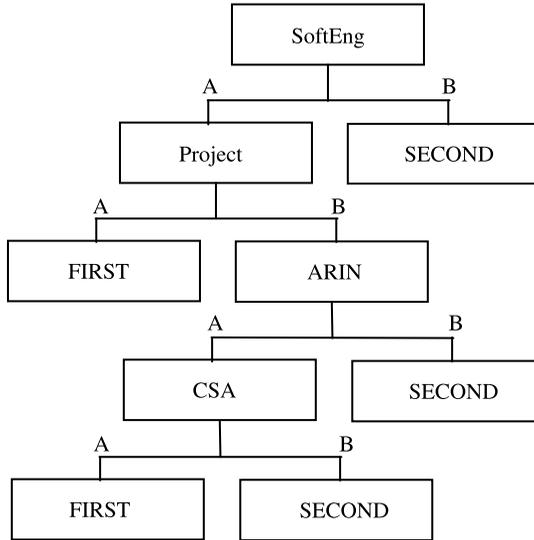


Figure 4.4 Decision Tree for the *degrees* Dataset

A set of rules of this kind is said to be in *Disjunctive Normal Form* (DNF). The individual rules are sometimes known as *disjuncts*.

Looking at this example in terms of data compression, the decision tree can be written as five decision rules with a total of 14 terms, an average of 2.8 terms per rule. Each instance in the original *degrees* training set could also be viewed as a rule, for example

IF SoftEng = A AND ARIN = B AND HCI = A AND CSA = B
AND Project = B THEN Class = SECOND

There are 26 such rules, one per instance, each with five terms, making a total of 130 terms. Even for this very small training set, the reduction in the number of terms requiring to be stored from the training set (130 terms) to the decision tree (14 terms) is almost 90%.

The order in which we write the rules generated from a decision tree is arbitrary, so the five rules given above could be rearranged to (say)

IF SoftEng = A AND Project = B AND ARIN = A AND CSA = B
THEN Class = SECOND
IF SoftEng = B THEN Class = SECOND
IF SoftEng = A AND Project = A THEN Class = FIRST
IF SoftEng = A AND Project = B AND ARIN = B
THEN Class = SECOND

```
IF SoftEng = A AND Project = B AND ARIN = A AND CSA = A
THEN Class = FIRST
```

without any change to the predictions the ruleset will make on unseen instances.

For practical use, the rules can easily be simplified to an equivalent nested set of IF ... THEN ... ELSE rules, with even more compression, e.g. (for the original set of rules)

```
if (SoftEng = A) {
  if (Project = A) Class = FIRST
  else {
    if (ARIN = A) {
      if (CSA = A) Class = FIRST
      else Class = SECOND
    }
    else Class = SECOND
  }
}
else Class = SECOND
```

4.2 The TDIDT Algorithm

Decision trees are widely used as a means of generating classification rules because of the existence of a simple but very powerful algorithm called **TDIDT**, which stands for *Top-Down Induction of Decision Trees*. This has been known since the mid-1960s and has formed the basis for many classification systems, two of the best-known being ID3 [3] and C4.5 [2], as well as being used in many commercial data mining packages.

The method produces decision rules in the implicit form of a decision tree. Decision trees are generated by repeatedly splitting on the values of attributes. This process is known as *recursive partitioning*.

In the standard formulation of the TDIDT algorithm there is a training set of instances. Each instance corresponds to a member of a universe of objects, which is described by the values of a set of categorical attributes. (The algorithm can be adapted to deal with continuous attributes, as will be discussed in Chapter 8.)

The basic algorithm can be given in just a few lines as shown in Figure 4.5.

At each non-leaf node an attribute is chosen for splitting. This can potentially be any attribute, except that the same attribute must not be chosen twice in the same branch. This restriction is entirely innocuous, e.g. in the branch

<p>TDIDT: BASIC ALGORITHM</p> <p>IF all the instances in the training set belong to the same class THEN return the value of the class ELSE (a) Select an attribute A to split on⁺ (b) Sort the instances in the training set into subsets, one for each value of attribute A (c) Return a tree with one branch for each <i>non-empty</i> subset, each branch having a descendant subtree or a class value produced by applying the algorithm recursively</p> <hr/> <p>⁺ Never select an attribute twice in the same branch</p>

Figure 4.5 The TDIDT Algorithm

corresponding to the incomplete rule

IF $SoftEng = A$ AND $Project = B$

it is not permitted to choose *SoftEng* or *Project* as the next attribute to split on, but as their values are already known there would be no point in doing so.

However this harmless restriction has a very valuable effect. Each split on the value of an attribute extends the length of the corresponding branch by one term, but the maximum possible length for a branch is M terms where there are M attributes. Hence the algorithm is guaranteed to terminate.

There is one important condition which must hold before the TDIDT algorithm can be applied. This is the *Adequacy Condition*: no two instances with the same values of all the attributes may belong to different classes. This is simply a way of ensuring that the training set is consistent. Dealing with training sets that are not consistent is the subject of Section 9.1.

A major problem with the TDIDT algorithm, which is not apparent at first sight, is that it is *underspecified*. The algorithm specifies ‘Select an attribute A to split on’ but no method is given for doing this.

Provided the adequacy condition is satisfied the algorithm is guaranteed to terminate and any selection of attributes (even random selection) will produce a decision tree, provided that an attribute is never selected twice in the same branch.

This under-specification may seem desirable, but many of the resulting decision trees (and the corresponding decision rules) will be of little, if any, value for predicting the classification of unseen instances.

Thus some methods of selecting attributes may be *much* more useful than others. Making a good choice of attributes to split on at each stage is crucial to the success of the TDIDT approach. This will be the main topic of Chapters 5 and 6.

4.3 Types of Reasoning

The automatic generation of decision rules from examples is known as *rule induction* or *automatic rule induction*.

Generating decision rules in the implicit form of a decision tree is also often called rule induction, but the terms *tree induction* or *decision tree induction* are sometimes preferred. We will end this chapter with a digression to explain the significance of the word ‘induction’ in these phrases and will return to the topic of attribute selection in Chapter 5.

Logicians distinguish between different types of reasoning. The most familiar is *deduction*, where the conclusion is shown to follow necessarily from the truth of the premises, for example

$$\begin{array}{l} \text{All Men Are Mortal} \\ \text{John is a Man} \\ \hline \text{Therefore John is Mortal} \end{array}$$

If the first two statements (the *premises*) are true, then the conclusion must be true.

This type of reasoning is entirely reliable but in practice rules that are 100% certain (such as ‘all men are mortal’) are often not available.

A second type of reasoning is called *abduction*. An example of this is

$$\begin{array}{l} \text{All Dogs Chase Cats} \\ \text{Fido Chases Cats} \\ \hline \text{Therefore Fido is a Dog} \end{array}$$

Here the conclusion is consistent with the truth of the premises, but it may not necessarily be correct. Fido may be some other type of animal that chases cats, or perhaps not an animal at all. Reasoning of this kind is often very successful in practice but can sometimes lead to incorrect conclusions.

A third type of reasoning is called *induction*. This is a process of generalisation based on repeated observations.

After many observations of x and y occurring together, learn the rule
if x then y

For example, if I see 1,000 dogs with four legs I might reasonably conclude that “if x is a dog then x has 4 legs” (or more simply “all dogs have four legs”). This is *induction*. The decision trees derived from the *golf* and *degrees* datasets are of this kind. They are generalised from repeated observations (the instances in the training sets) and we would expect them to be good enough to use for

predicting the classification of unseen instances in most cases, but they may not be infallible.

4.4 Chapter Summary

This chapter introduces the TDIDT (Top-Down Induction of Decision Trees) algorithm for inducing classification rules via the intermediate representation of a decision tree. The algorithm can always be applied provided the ‘adequacy condition’ holds for the instances in the training set. The chapter ends by distinguishing three types of reasoning: deduction, abduction and induction.

4.5 Self-assessment Exercises for Chapter 4

1. What is the adequacy condition on the instances in a training set?
2. What are the most likely reasons for the condition not to be met for a given dataset?
3. What is the significance of the adequacy condition to automatic rule generation using the TDIDT algorithm?
4. What happens if the basic TDIDT algorithm is applied to a dataset for which the adequacy condition does not apply?

References

- [1] Michie, D. (1990). Machine executable skills from ‘silent’ brains. In *Research and development in expert systems VII*. Cambridge: Cambridge University Press.
- [2] Quinlan, J. R. (1993). *C4.5: programs for machine learning*. San Mateo: Morgan Kaufmann.
- [3] Quinlan, J. R. (1986). Induction of decision trees. *Machine Learning*, 1, 81–106.