

# 1 Parsing and Classification

## 1.1 Choice of Language

With very few exceptions, the code I wrote in support of this thesis was done in Clojure, a dialect of LISP designed to work on top of the Java Virtual Machine (JVM). The choice of a language was easy: a heavy dependence on the Stanford Parser and the WEKA package, both written in Java, necessitated a JVM-based language. The slowness of Java's compile/debug cycle eliminated that language as an option, leaving a handful of possible languages, from which I chose Clojure for its speed, functional style, and elegance.

## 1.2 Parsing

The Stanford Parser software package, version 1.6.7, configured with the probabilistic context-free grammar (PCFG) [Klein and Manning 2003], was used to generate all syntactic parse trees and grammar dependency graphs. In brief, PCFGs have their origins in the work of

## 1.3 The Tests

The crux of this project was the design and creation of a suite of tests, each of which identifies a number of closely related grammatical characteristics of the text samples. These tests operate on the output from the Stanford parser, i.e. parse trees and grammatical dependencies. As output they generate training or testing cases to be used by the Weka classifier. Each of these cases consists of multiple attributes, corresponding to grammatical features, each with continuous values indicating the relative frequency (probability) of that particular feature. For a case with  $n$  attributes where the number of occurrences of the grammatical feature associated with the  $i$ th attribute is  $g_i$ , the value  $f_i$  for that attribute is given by  $g_i / \sum_{i=1}^n g_i$ . For instance, one test measures the relative frequencies of the various

tense/aspect/voice combinations of finite verbs. English has twenty-four such combination, so the case generated by this test has twenty-four attributes.

In addition to the attributes, each case has a class which can be *es* or *en*, indicating that the class is associated with a text sample written by an L1-Spanish speaker or by a native English speaker, respectively. For training cases, the classes are known beforehand and are assigned to the cases manually. For testing cases, the classes have missing values, until such values are determined by a classifier, as discussed in the following section.

## 1.4 Classification

I used the Weka machine learning package, version 3.6 [Hall et al. 2009], to create, train and test classifiers based on the cases discussed above. I primarily used two classifiers: J48, which is Weka's implementation of the C4.5 classifier [Quinlan 1993] and the RandomForest classifier, which is based on the random forest algorithm described by Breiman [2001]. The former is useful for its highly readable decision trees, which clearly indicate which attributes are involved in the classification and their roles. In later sections of this paper are found linguistic explanations for why these particular attributes should be useful in classification.

### 1.4.1 C4.5

This section describes the C4.5 partition as it applies to this project. That is to say, C4.5 can deal with a number of circumstances that do not arise here. What is described here is a version of the C4.5 algorithm that is restricted to continuous attribute value and to exactly two class values, and which does not permit missing attribute values. That having been said, the C4.5 algorithm consists of two phases, *tree construction* and *tree pruning*.

In the tree construction phase a decision tree is built which successively performs binary partitioning of a set of training cases (see fig. 1.1). Let us consider a full binary tree where each edge represents a set of cases and each non-terminal node a partitioning operation.

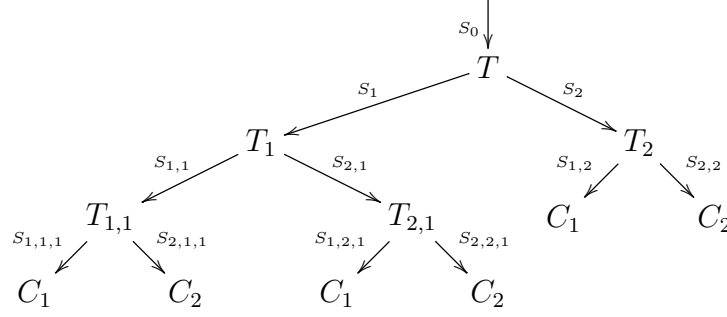


Figure 1.1: A decision tree showing the partitioning of the set of all training cases  $S$  into subsets  $S_{1,2}$ ,  $S_{1,1,1}$ , and  $S_{1,2,1}$  whose elements are of class  $C_1$ , and  $S_{2,2}$ ,  $S_{2,1,1}$ , and  $S_{2,2,1}$  whose elements are of class  $C_2$ . The nodes  $T$ ,  $T_1$ , etc. are partitioning operations.

These partitioning operations take one set, represented by the parent edge, and divide it into two subsets, the daughter edges. The root node operates on an initial set  $S_0$ , and a leaf node simply indicates that its parent edge is a set consisting of cases of a single class. Let the first partitions of  $S_0$  be called  $S_1$  and  $S_2$ , and of  $S_1$  let them be called  $S_{1,1}$  and  $S_{2,1}$  and so forth. Likewise, let the partitioning operation that operates on a particular set be designated by  $T$  with the same subscripts as that set.

The partitioning operations are performed by applying a binary test to each case within  $S$ , the set to be partitioned, and dividing the set based on the results. Each test considers a single attribute  $A$  and compares that value to a threshold value. All cases where the value of attribute  $A$  is less than or equal to that threshold will be put into one subset and all other cases into the other.

into a number of subsets, each containing cases of a single class.

## References

- BREIMAN, L. 2001. Random forests. In *Machine Learning*. 5–32.
- HALL, M., FRANK, E., HOLMES, G., PFAHRINGER, B., REUTEMANN, P., AND WITTEN, I. H. 2009. The WEKA data mining software: An update. *SIGKDD Explor.*

*rations 11, 1.*

KLEIN, D. AND MANNING, C. D. 2003. Accurate unlexicalized parsing. In *Proceedings of the 41st Meeting of the Association for Computational Linguistics*. 423–430.

QUINLAN, J. R. 1993. *C4.5: Programs for Machine Learning*. Morgan Kaufmann Publishers.