

Thesis

Philip White

February 15, 2012

1 Corpora

The data used in this study was drawn from nine different corpora. Of these, three contained only native texts, four only nonnative texts, and two texts of both types. Table 1.1 shows the number of tokens contributed by each corpus. A token is a unit parseable by the Stanford parser, the large majority of which are simply words but which also include punctuation and the genitive suffixes 's and '. As can be seen in the table, the two classes of texts (native and nonnative) are very closely matched in size. Furthermore, the number of samples in each class is identical, 321, giving a total of 642 instances or cases. All classification methods used in this study operated on these same 642 instances.

The following corpora contributed native samples: the Brown University Standard Corpus of Present-Day American English subcorpus of letters-to-the-editor and editorials (BROWN), the International Corpus of English-Hong Kong (ICE-HK), the Michigan Corpus of Upper-level Student Papers (MICUSP), the Open American National Corpus (OANC), and the International Corpus of English-Canada (ICE-CAN). MICUSP and ICE-CAN contributed nonnative samples as well, and the remainder of the nonnative texts came from the International Corpus of Learner English, Spanish Subcorpus (SPICLE), the Santiago University Learner Corpus (SULEC), and the Written Corpus of Learner English (WRICLE). One additional student paper supplied by Missouri State University's English

Language Institute rounded out the nonnative samples. All nonnative samples were written by individuals whose first language was Spanish and who were judged, by the compilers of the corpora, to be advanced English learners. Many of the individuals had a language in addition to English and Spanish. In the cases of the SULEC and WRICLE corpora, both of which were compiled at Spanish universities, a large number of the learners spoke other Romance languages in addition to Spanish, in particular Catalan and Galician. Many of the samples in the ICE-HK corpus were written by individuals whose second language was Cantonese, and a number of the contributors to ICE-CAN had some French as well. Any sample written by an individual who knew a Germanic language (other than English) was not included.

Table 1.1: Corpora Composition

Corpus	Tokens Native	Tokens Nonnative
BROWN	57,809	0
ICE-HK	59,674	0
MICUSP	163,218	29,897
MSUELI	0	538
OANC	84,0522	0
SPICLE	0	216,879
SULEC	0	39,254
WRICLE	0	96,247
ICE-CAN	25,225	2,070
Total	389,978	384,885

2 Parsing and Classification

2.1 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 simple: 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.

2.2 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

2.3 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.

2.4 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.

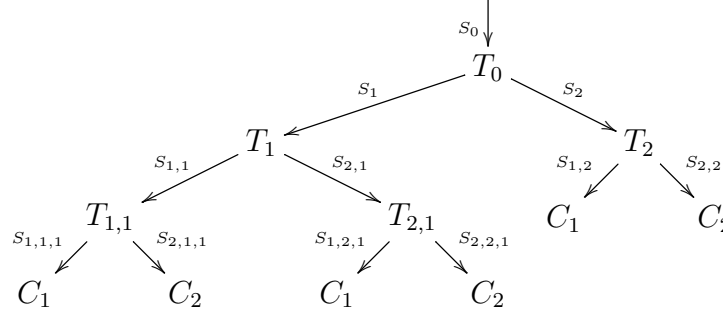
2.4.1 C4.5

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 values 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. Consider a full binary tree where each edge represents a set of cases and each non-terminal node a partitioning operation, as shown in Figure 2.1. 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 where $S_1 \cup S_2 = S_0$ and $S_1 \cap S_2 = \emptyset$, 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

Figure 2.1: A decision tree showing the partitioning of a set of training cases S_0 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_0, T_1 , etc. are partitioning operations such that for any operation T operating on a set S_a the generated sets are $S_{1,a}$ and $S_{2,a}$ where $S_{1,a} \cup S_{2,a} = S_a$ and $S_{1,a} \cap S_{2,a} = \emptyset$.



a single attribute A and compares the value of that attribute, V_A , to a threshold value, V_C . All cases where the $V_A \leq V_C$ will be put into one subset and all other cases into the other.

The decision of the attribute and threshold value for a particular test is determined using what Quinlan calls the “gain ratio criterion” which is calculated as follows. If the probability of randomly drawing a case of class C_1 from a set S is p_1 and of drawing a case of the other class is p_2 where $p_2 = 1 - p_1$, then the average amount of information needed to identify the class of a case in S can be defined in terms of entropy as

$$\text{info}(S) = -p_1 \cdot \log_2(p_1) - p_2 \cdot \log_2(p_2).$$

A similar measure can be applied to the two partitions S_1 and S_2 created by applying the partitioning test T to S . The entropy after partition is given by taking a weighted sum of the entropy of the two sets as

$$\text{info}_T(S) = \frac{|S_1|}{|S|} \cdot \text{info}(S_1) + \frac{|S_2|}{|S|} \cdot \text{info}(S_2)$$

The decrease in entropy, expressed as a positive value (an information gain), due to parti-

tioning S using the test T is then

$$\text{gain}(T) = \text{info}(S) - \text{info}_T(S).$$

Maximizing this gain can be and, in ID3 the predecessor to C4.5, was used as measurement of test fitness. However, in the more general case of C4.5, where one test can partition a set into more than 2 subsets, using this gain criterion to choose tests favors tests that partition sets into numerous subsets. To mitigate this, Quinlan added another factor to the criterion, the split info which for this special case is given by

$$\text{split info}(T) = -\frac{|S_1|}{|S|} \cdot \log_2 \left(\frac{|S_1|}{|S|} \right) - \frac{|S_2|}{|S|} \cdot \log_2 \left(\frac{|S_2|}{|S|} \right).$$

Then the fitness of a test T can be measured using

$$\text{gain ratio}(T) = \frac{\text{gain}(T)}{\text{split info}(T)}$$

It should be noted that in this special case where partitioning operations are always binary, the gain ratio criterion favors tests that split S into disparately sized sets, as split info is at its maximum (unity) when $|S_1| = |S_2|$.

In choosing a test T , the C4.5 algorithm tries each attribute A from the set S of cases to be partitioned. For each, it orders the cases in S on the value of A . If the values of A corresponding to this ordered set are $\{v_1, v_2, \dots, v_m\}$, then any threshold between v_i and v_{i+1} will result in the same partitions. From this it can be seen that the total number of possible partitions is $m - 1$. The algorithm tries all such partitioning schemes, measuring the gain ratio of each. When an optimal attribute and corresponding partitioning scheme has been chosen, the algorithm than chooses a threshold value that will produce this result. Again, to partition S into two sets where the values for A are $\{v_1, v_2, \dots, v_i\}$ and $\{v_{i+1}, v_{i+2}, \dots, v_m\}$, a threshold value v_C must be chosen such that $v_i \leq v_C < v_{i+1}$. For

this, it chooses the largest value for A from the entire training set S_O that does not exceed the midpoint of this range.

3 Grammatical Relations

The simplest classification approach used in this study considered the relative frequency of different grammatical relations. For this approach, the governor and the dependent of the dependencies were ignored, with only the relation itself being used.

Each data set instance contained attributes corresponding to dependency relations. The Stanford parser system in its default configuration does not generate the *punct* or punctuation dependency which connects punctuation symbols to a key element in the associated clause. Since English punctuation is broadly similar to Spanish punctuation, aside from some stark differences such as Spanish’s inverted question and exclamation marks, which should be apparent to even the beginning learner, it did not seem to useful to activate this dependency. Additionally, the *abbrev* or abbreviation dependency was removed. This dependency marks the definition of an abbreviation, as in the example given by de Marneffe and Manning [2008], “Australian Broadcasting Corporation (ABC)”, where the dependency would be *abbrev*(Corporation, ABC). This dependency has little to do with grammar, and thus was ignored for the purposes of this study. Having excluded these two dependencies, each data set instance contained 58 numerical attributes, one for each relation used.

For each attribute A_r corresponding to the relation r , the corresponding value was the floating point number n_r/n_t , where n_r and n_t were the number of occurrences of the relation r and the total number of relations in the text, respectively. A C4.5 decision tree classifier trained on these instances produces the decision tree shown in Figure 3.1, employing 15 different relations. The full names for these relations are shown in Table 3.1. At each terminal node of the tree there is an integer or pair of integers in parentheses. These values indicate the number of the training cases that were categorized (correctly or not)

at that node and the number of cases incorrectly categorized, this latter value only being shown when greater than zero. For any given test node, one can identify one branch as the predominately *en* branch and the other as the *es* branch. For test nodes where one or both branches lead to terminal nodes, this is trivial, as the terminal nodes themselves label the branches. For any other test node, the branches can be identified by summing up the number of test cases at the terminal nodes of that branch. For instance, the root test node, which considers the relation *nn*, divides the training set of 642 cases into a subset of 337 cases, associated with the left branch, and another subset of 305 cases, associated with the right branch. Looking at the left branch, it can be seen that of these 336 cases, 301 of them are nonnative, i.e. of the class *es*, and only 36 are native. This indicates that this is a predominately nonnative branch. Conversely, the right hand branch consists of 205 native cases and only 20 nonnative cases, making it the native branch. This allows one to say, for instance, that fewer occurrences of the *nn* relation are associated with nonnative samples. The following subsections explore the linguistic reasons why these relations should be so useful in making such categorizations.

Table 3.1: Relation abbreviations

<i>auxpass</i>	passive auxiliary
<i>complm</i>	complementizer
<i>cop</i>	copula
<i>det</i>	determiner
<i>expl</i>	expletive
<i>mwe</i>	multi-word expression
<i>nn</i>	noun compound modifier
<i>npadvmod</i>	noun phrase as adverbial modifier
<i>parataxis</i>	parataxis
<i>poss</i>	possession modifier
<i>possessive</i>	possessive modifier
<i>predet</i>	preconjunct
<i>prt</i>	phrasal verb particle
<i>quantmod</i>	quantifier phrase modifier
<i>rel</i>	relative

Figure 3.1: C4.5 decision tree employing relative frequency of dependency relations. Relative frequencies are shown as percentages. Values in parentheses are the number of training case classified at that point and, following the slash when present, the number of those cases which were incorrectly classified.

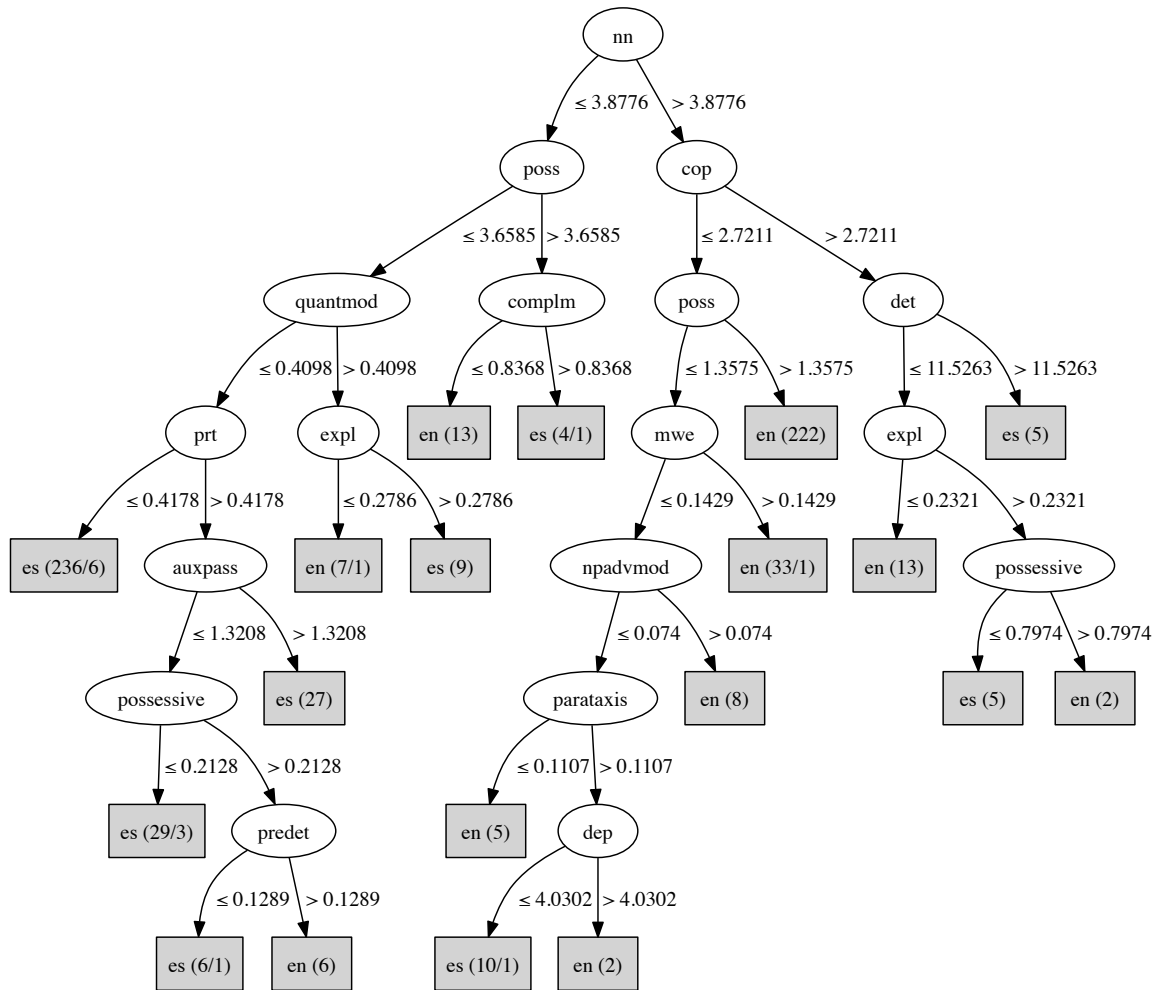
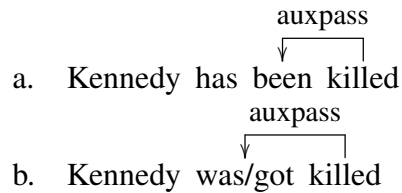


Figure 3.2: The dependencies *auxpass*(killed, been) and *auxpass*(killed, was/got). Taken from de Marneffe and Manning [2008].



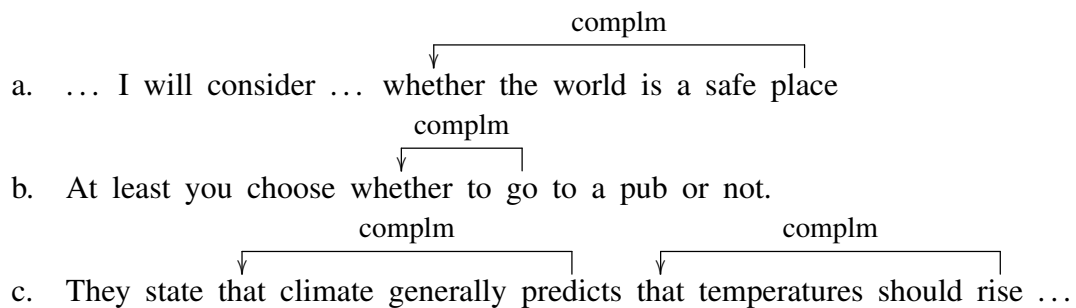
3.1 Passive Auxiliary

The passive auxiliary dependency *auxpass* marks an auxiliary verb which carries the passive information of the clause. In general a parsed sample of text will contain one such dependency for every passive clause and so a high relative frequency of this relation indicates heavy usage of the passive voice. Example 3.2 illustrates this dependency.

3.2 Complementizer

A complementizer is a word that signals the beginning of a clausal complement. The Stanford Parser recognizes the complementizers *that* and *whether* as shown in Example 3.3. The governor of a complementizer dependency is the root of the clause, which is generally a verb or, in the case of copular clauses, the subject complement. The dependent is the complementizer itself.

Figure 3.3: The dependencies *complm*(place, whether), *complm*(go, whether), *complm*(predicts, that), and *complm*(rise, that). Nonnative samples from WRICLE (*a* and *c*) and SULEC (*b*).



Whitley [1986] points out that while English tends to allow complementizers introducing clausal complements in the object position to be deleted, Spanish generally does not (see Example 3.1). Butt and Benjamin [2004, 33.4.6] explain that this rule is occasionally broken, but generally only in two situations, business letters and substandard speech, and when the complementizer *que* appears close to other uses of the word *que*. Since these are restricted cases, it is reasonable to conclude that there would be L1-transfer in the construction of clausal complements, leading to L1-Spanish learners to have some preference for Example 3.1a over 3.1b, particularly considering that they are both perfectly valid constructions.

In a study on differences in complement clause usage between native and nonnative English speakers, Biber and Xeppen [1998] make a number of conclusions relevant to the current study. First, they consider when native speakers omit the complementizer *that* and conclude that it is rarely omitted in academic prose and in opinion and descriptive essays. Since the vast majority of the corpus samples both native and nonnative fall into these categories, this provides encouraging evidence that the differences in complementizer usage identified by the classifier are not due to idiosyncrasies in the samples. Next, while considering four different groups of L1-speakers, French, Spanish, Chinese, and Japanese, Biber and Xeppen find that all groups shows similar levels of *that* omission, and in general these levels of omission are lesser than the levels found in comparable types of native texts. They also find, interesting that L1-Spanish speakers use complement clauses, with and without omission of the complementizer, more often than either native speakers or the other groups of learners.

The decision tree shown in Figure 3.1 uses the *complm* dependency once, and classifies cases with lower occurrences of *complm* as native and larger occurrences as nonnative, without further testing. This this dependency does not part necessarily indicate the presence of a complement clause, but rather the presence of a complementizer, the higher frequency among the learners may be due either to low rates of dropping the complementizer, or

- (3.1) a. *I say that he'll do it.*
b. *I say he'll do it.*
c. *Digo que lo hará.*
d. **Digo lo hará.* (Whitley 1986, p. 278)

high rates of complement clause usage. As shown above, both phenomena have linguistic backing and very likely both are at play.

3.3 Copula

The copula or *cop* dependency marks the copular verb. This dependency takes as its governor the complement of the copular clause and the verb itself as the dependent.

3.4 Determiner

The determiner or *det* dependency connects a determiner to the NP it modifies with the determiner being the dependent and the head of the NP the governor.

3.5 Expletive

An existential *there* and the copular verb associated with it are connected with the expletive or *expl* relation.

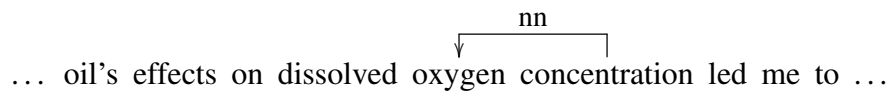
3.6 Multi-Word Expression

The Stanford typed dependency manual [de Marneffe and Manning 2008] defines multi-word expressions as being two or more words that are used together as a single unit such that the relationship between them is difficult to define. In the version of the Stanford parser used here, only the following expressions are considered multi-word expressions: *rather than, as well as, instead of, such as, because of, in addition to, all but, due to.*

3.7 Noun Compound Modifier

Noun-noun compounds (NNCs) are marked with the relation *nn*. The governor of this dependency is the rightmost noun in the compound and the dependent will be one of the nouns to the left. Note that since all dependencies only deal with pairs of words, a compound consisting of more than two nouns would be indicated by multiple dependencies, all sharing a common governor. Example 3.4 demonstrates this dependency.

Figure 3.4: The dependency *nn*(concentration, oxygen). Native sample taken from MICUSP.



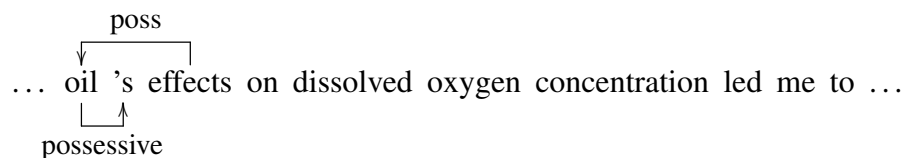
3.8 Noun Phrase as Adverbial Modifier

3.9 Parataxis

3.10 Possession and Possessive Modifiers

Inflected genitive constructions are marked by two dependencies: *poss*, which ties the head of a NP (the governor) to a genitive inflectional suffix ('s or '), indicating that the governor is the possessed element; and *possessive*, which connects a noun to its own genitive inflectional suffix. These two dependencies are illustrated in Figure 3.5. The *poss* dependency can also have as its dependent a possessive determiner such as *its* or *their*. In this type of construction, the *possession* dependency is not used.

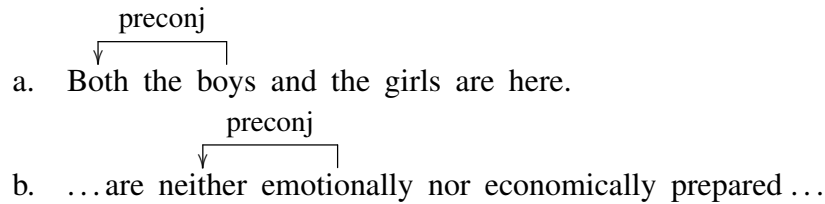
Figure 3.5: The dependencies *poss*(effects, oil) and *possessive*(*textoil*, 's). Native sample taken from MICUSP.



3.11 Preconjunct

The preconjunct (*preconj*) dependency connects the head of a phrase employing a conjunction to a word that emphasizes or brackets that conjunction, such as *either*, *neither*, or *both*. Figure 3.6 demonstrates this dependency.

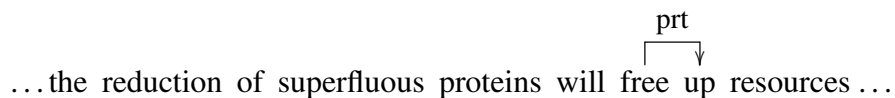
Figure 3.6: The dependencies *preconj*(boys, both) and *preconj*(emotionally, neither). (a) taken from de Marneffe and Manning [2008] and (b) from WRICLE (nonnative).



3.12 Phrasal Verb Particle

The phrasal verb particle relation (*pvt*) ties the head word of a phrasal verb to its particle as shown in Example 3.7. The decision tree in Figure 3.1 contains this relation once. Relative frequencies of less than or equal to 0.4178% lead to the categorization of a text as nonnative, whereas larger values lead to a subtree. It can be seen that a very high percentage, 36.8%, of the training cases terminate at the left, or nonnative, branch of this test node, suggesting that this relation contributes a great deal of useful information to the categorization process.

Figure 3.7: The dependency *pvt*(free, up). Native sample from MICUSP.



Phrasal verbs are multiword verbs consisting of a core word, which can generally stand alone as a distinct verb in other circumstances, and a preposition-like particle appearing after, though in many cases not immediately after, the primary word [Celce-Murcia and Larsen-Freeman 1999]. These verbs appear to be rare in world languages, with few non-Germanic languages containing such constructions [Celce-Murcia and Larsen-Freeman

1999]. Liao and Fukuya [2004] conduct a review of the literature on phrasal verb avoidance in English language learners, starting with [Dagut and Laufer 1985], a study which concluded that L1-Hebrew learners of English do avoid these verbs. They further asserted that the reason for this was syntactic differences between Hebrew and English, though others have questioned their bases for this assertion [Liao and Fukuya 2004]. The review continues with [Hulstijn and Marchena 1989], who investigated the claims of Dagut and Laufer by applying their same data gathering techniques to a group of English learners whose first language was Dutch, a language which also uses phrasal verbs. Contrary to their expectations, they found that the Dutch speakers did not avoid phrase verbs in English, suggesting that L1-interference is, at least in part, the source of phrasal verb avoidance. Finally, the review cites the study of Laufer and Eliasson [1993], which performed a very similar study as Hulstijn and Marchena, but with native Swedish speakers, and made much the same conclusions.

In their own study, Liao and Fukuya investigate L1-Chinese learners of English, and cautiously concluded that the syntactic features of Chinese lead to the avoidance of phrasal verbs in the English of those learners. A later study, Alejo González [2010], uses the Spanish and Swedish subcorpora of ICLE along with the British National Corpus (BNC), a corpus of native written English, to perform a quantitative study of phrasal verb usage. They found that the L1-Swedish learners used phrasal verbs 69% as often as the native speakers and the L1-Spanish learners used phrasal verbs 45% as often. These numbers would seem to indicate that the syntax of the learner's L1 is an important, but not the only, contributing factor to phrasal verb avoidance.

Regardless of the reasons behind L1-Spanish learners avoidance of phrasal verbs, Alejo González [2010] demonstrates that it is a reality of learner English. Considering this, it is not surprising that the C4.5 algorithm uses the *pvt* relation with such success in the categorization process.

Table 3.2: Accuracy results for C4.5 and 100 tree Random Forest classifiers using 20 fold cross-validation on data set of 642 cases.

	C4.5		R. Forest	
Classified as →	es	en	es	en
es	309	12	291	30
en	25	296	36	285
% Correct	89.72		94.24	
MAE	0.1139		0.1707	
κ	0.7944		0.8847	

3.13 Quantifier Phrase Modifiers

3.14 Relative

3.15 Classification Accuracy

Twenty fold cross-validation was used to test the real-world accuracy of the data. There being 642 cases in the data set, thirty-two unique cases were held out at a time and classified using a C4.5 classifier trained on the remaining 610 cases. This produced a correct classification rate of 89.72% with a mean absolute error (MAE) of 0.1139 and a κ value of 0.7944. Using a random forest classifier gave better results; performing 20 fold cross-validation on a 100 tree classifier where each tree was trained on six random features yielded 94.24% accuracy with MAE = 0.1707 and κ = 0.8847. Table 3.2 gives the confusion matrices for these two classifier.

4 Argument Structure

The classification systems discussed in this section considered the argument structure of verbs. In general, every finite verb in English takes one or more arguments, with a subject argument being required in normal speech and writing. The Stanford NLP system marks the arguments of verbs using the dependencies shown in Table 4.1. In the majority of

non-copular sentences, the governors of these dependencies are the core verbs. In copular sentences the governor is generally the subject complement (i.e. the argument generally appearing after the verb which is equated with the subject) though in the case of copular sentences with clausal subjects, the Stanford parser chooses the copula to be the governor. Figure 4.1 show examples of these dependencies.

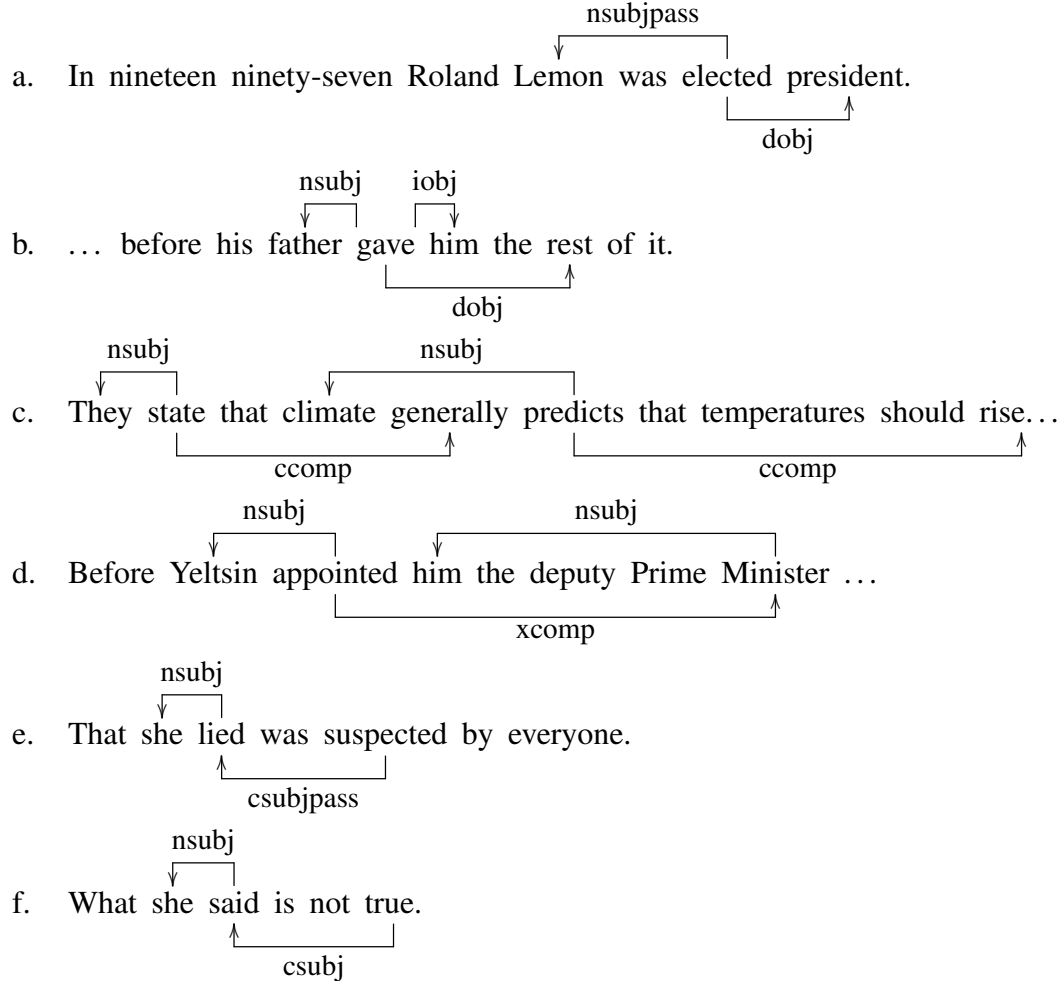
Table 4.1: The dependencies used to identify verbal arguments.

ccomp	Clausal Complement
csubj	Clausal Subject
csubjpass	Passive Clausal Subject
dobj	Direct Object
iobj	Indirect Object
nsubj	Nominal Subject
nsubjpass	Passive Nominal Subject
xcomp	Open Clausal Complement

Du Bois [2003] provides evidence that verbal argument structure may differ between native and nonnative speakers. That study explores when speakers choose to use lexical NPs rather than pronominal NPs as verbal arguments. Considering only native speakers but looking at a number of different languages, Du Bois finds that there is a very strong tendency for speakers to use no more than one lexical argument per finite clause (Du Bois' *One Lexical Argument Constraint*) and to avoid placing an argument in the subject role of transitive sentences (the *Non-Lexical A¹ Constraint*). He makes the case that these rules hold true for a number of world languages, at least in spoken form. However, in presenting data to show that several languages abide by the Non-Lexical A Constraint, he also shows that there are large differences between languages in the likelihood of a lexical argument appearing in the direct object and intransitive subject roles [Du Bois 2003, Table 2.5]. His study shows that, based on his data, in English 21% of lexical arguments are found in intransitive clauses, versus 28% for Spanish, and 79% are found in direct object roles, versus 71% for Spanish. A χ^2 analysis shows that these differences are statistically significant

¹Du Bois [2003] uses the letters **A**, **I**, and **O** to refer to the subject, indirect object, and direct object arguments of a transitive verb, and **S** to refer to the sole argument of an intransitive verb.

Figure 4.1: Examples of the Dependencies Listed in Table 4.1. *a*, native sample from ICE-CAN; *b,c,d*, native samples from MICUSP; *e,f* from de Marneffe and Manning [2008].



considering his sample sizes ($\chi^2 = 2.244$, $df = 2$, $\rho = 0.326$), though admittedly at a relatively low confidence level. As an aside, of the five languages for which data is presented in that table, English and Spanish are actually the most similar in lexical argument distribution. Of the three other languages considered, French, Hebrew, and Sakapultek, all show a greater probability of finding a lexical NP in the intransitive subject role, and a lower probability of finding one in the direct object role, than either Spanish or English. Seeing that there may be differences between Spanish and English in how lexical arguments are distributed among the various argument roles, and noting that other languages show such differences as well, it is worth investigating whether L1-Spanish learners of English differ

in their usage of these lexical arguments as compared to native speakers.

To investigate the feasibility of classifying language based on argument structure, I used Stanford dependency graphs to identify 18 different types of finite clauses: intransitives with and without clausal subjects, copular clauses with and without clausal subjects, simple transitives with and without clausal subjects and objects, ditransitives with and without clausal subjects, complex transitives with and without clausal subjects, passives of simple transitives with or without clausal subjects, and passives of complex transitives with or without clausal subjects and complements. Then for each of these 18 different types of clauses, all possible permutations of lexical and non-lexical arguments were considered, yielding 80 different attributes in total. This system, while providing a generous amount of data for the classifiers, proved to yield decision trees that were difficult to interpret linguistically. Fortunately, I found that from these attributes I could generate three much smaller and more coherent set of attributes which, when combined and used to train a classifier, provided comparable accuracy.

The first of these attribute sets consisted of just three arguments, corresponding to verb valency (i.e. the number of arguments). The values associated with these attributes were the percentage of all finite clauses which used that number of arguments. The next attribute set consisted of four attributes, named *zero*, *one*, *two*, *three* with corresponding values indicating what percentage of finite clauses had that number of lexical arguments. To avoid lengthy periphrasis, it is convenient to call this metric *lexical argument density*. The last attribute set consisted of seven attributes corresponding to types of arguments. The values associates with these attributes were the percentage of lexical arguments were found in that type of argument. The seven types were intransitive, transitive, passive, and copular subjects, indirect and direct objects, and subject complements. This latter category included both the complement of copular clauses and what is usually the third argument in a complex transitive (e.g. *deputy Prime Minister* in Table 4.1). This set of attributes will be referred to as the *lexical argument role* set.

To determine what was a lexical argument and what was not, I compiled a list of pronominal forms, largely taken from Celce-Murcia and Larsen-Freeman [1999, Ch. 16] and shown in Table 4.2, and considered any argument that did not match one of these forms to be lexical. By this standard, clausal arguments were always considered lexical arguments.

Table 4.2: Pronouns Used to Determine Non-Lexical Status of Arguments.

other	another	else	same	one	
this	that	these	those	what	
myself	yourself	herself	himself	itself	
ourselves	yourselves	themselves	oneself		
mine	yours	hers	his	ours	theirs
me	you	her	him	it	us
them	I	she	he	we	they

Figure 4.2 shows a C4.5 decision tree trained on cases consisting of the four verb valency attributes. It is noteworthy that the tree considers only two out of the three attributes, and appears to indicate that native English has a larger proportion of trivalent verbs than does nonnative English, and that the opposite is true of divalent verbs. However, an extremely large number of training cases are misclassified by this tree, particularly by the middle leaf. In fact, Table 4.3 shows that the tree on average only classifies slightly more than half of all test cases correctly, and, worse still, the confidence interval for that accuracy spans 50%, indicating that the tree may be doing no better than a random classifier. It is a bit surprising that the C4.5 algorithm was unable to construct a useful decision tree, as Du Bois presents data showing that there are statistically significant differences in the frequency of usage of transitive and intransitive verbs between English and Spanish [Du Bois 2003, Table 2.3]. He found that in English 58% of finite verbs were intransitive and 42% transitive, whereas in the case of Spanish the numbers were 63% and 37%. Furthermore, calculating from the data he gives, this difference is statistically significant at a high level of confidence: $\chi^2 = 4.693, df = 1, p < 0.05$. Based on Du Bois' data, one might expect a decision tree to classify texts as native based on a low ratio of transitive (di- and triva-

lent) verbs to intransitive (monovalent) verbs. However, there are a number of possible reasons why this is not the case, the most obvious being that perhaps this characteristic is not involved in L1-transfer.

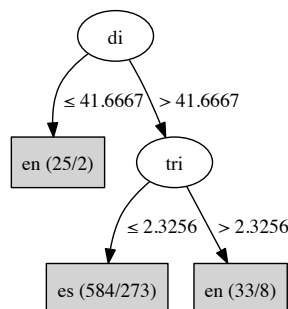


Figure 4.2: Verbal Clause Valency C4.5 Decision Tree.

Table 4.3: Accuracy of Verbal Clause Valency C4.5 Classifier

Nonnative	84.4%
Native	17.4%
Overall	50.9%
Overall C.I. 95%	47.1% — 54.8%

Figure 4.3 shows a decision tree generated by the C4.5 algorithm using cases with attributes indicating lexical argument density. Though Du Bois [2003] did gather this type of data, he only appears to have done so for English and the Mayan language Sakapultek. The differences between those two languages in this regard are pronounced, but suggest little about any such differences between English and Spanish. The tree itself is unusual, as it considers only two of the four attributes, one of which is considered twice. The first node in the tree checks if fewer than 8% of clauses have no lexical arguments and, if so, immediately classifies that case as native. Du Bois' data shows that such clauses are quite common in spoken English, accounting for a startling 87% of all clauses. In written English, they are decidedly less common. In the 321 native training cases finite verbs with no lexical arguments account for only $17 \pm 1\%$ (95% confidence) of all verbs. There does not seem to be much available research on lexical argument density in nonnative language, but it might be conjectured that learners of English have not developed clear

distinctions between the written and spoken registers, and thus tend to show some speech-like elements in their written language, such as a high incidence of finite verbs without lexical arguments. The other two nodes in Figure 4.3 are more opaque to interpretation. These nodes both consider the two lexical argument attribute, essentially performing a ternary test on it, classifying cases with intermediate values as nonnative and cases with high and low values as native. It is hard to imagine that there is a convincing linguistic reason behind this, particularly considering the error associated with these node. Speaking of which, Table ?? shows the accuracy of this classification system. As the confidence interval shows, it is an improvement over the valency system.

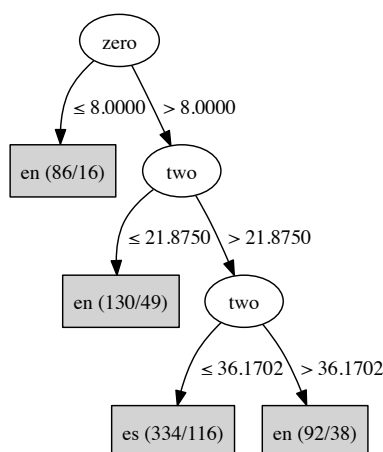


Figure 4.3: Lexical Argument Density C4.5 Decision Tree.

Table 4.4: Accuracy of Lexical Argument Density C4.5 Classifier

Nonnative	70.7%
Native	48.6%
Overall	58.7%
Overall C.I. 95%	55.9% — 63.5%

Training a C4.5 classifier on the lexical argument role attribute set yields the decision tree shown in Figure 4.4. This tree is considerably more complicated than those shown in Figures 4.2 and 4.3, and is more accurate as well. Table 4.5 shows that it classifies testing cases correctly approximately two thirds of the time. Considering this tree in light of Du

Bois's data shows some interesting parallels. Of the three roles that Du Bois considered, he found that there was not a significant difference between English and Spanish in the number of lexical NPs appearing as the transitive subjects, and, in fact, found very few such arguments. It stands to reason then, that the classifier excluded the transitive subject attribute when constructing the decision tree. Du Bois also gathered data for the intransitive subject role, finding that English tends to have fewer new core arguments in this role than does Spanish. As he explains in an earlier paper [?], he lumps the subjects of copular verbs into this role as well. However in the decision tree, intransitive and copular subjects are considered separately, as are the subjects of transitive verbs in the passive voice, which Du Bois presumably considered to be intransitive as well. Considering first the intransitive subject (shown as *intr-subj in the three*), it can be seen that this attribute appears twice in the decision tree. Where it appears closer to the root, the node splits the test cases into those with values less than or equal to 11.6162% and those higher. The lower values, are immediately classified as native, whereas the others undergo an addition test, with the majority ultimately being classified as nonnative. The other usage of this attribute is in another branch, and deeper in the tree. Here, using a higher comparison value, low values lead to a branch where the majority are ultimately classified as nonnative, and the opposite for high values. With these conflicting branches, it is difficult to draw conclusions, but we can see that the latter-mentioned branch, the one that appears to coincide with Du Bois data, does deal with more training cases (43 with 86% accuracy) than does the other (20 with 90% accuracy). This suggests that perhaps while some minority of learners overuse lexical intransitive subjects, a large group of learners underuse them. This certainly holds true for the training data, though whether it applies to learners in general is unclear. With the copular subject a very similar pattern holds. Early in the decision process, one test splits the training cases based on the *cop-sub* attribute and the group with the lower value is immediately classified as native. The other group, which is quite large, undergoes a number of tests more. The second use of the *cop-sub* is much deeper in the tree and is actually

a direct descendant of the first, meaning that the information content of the attribute was much higher during the first test. This is convenient, as the first test is the one that coincides with Du Bois's data.

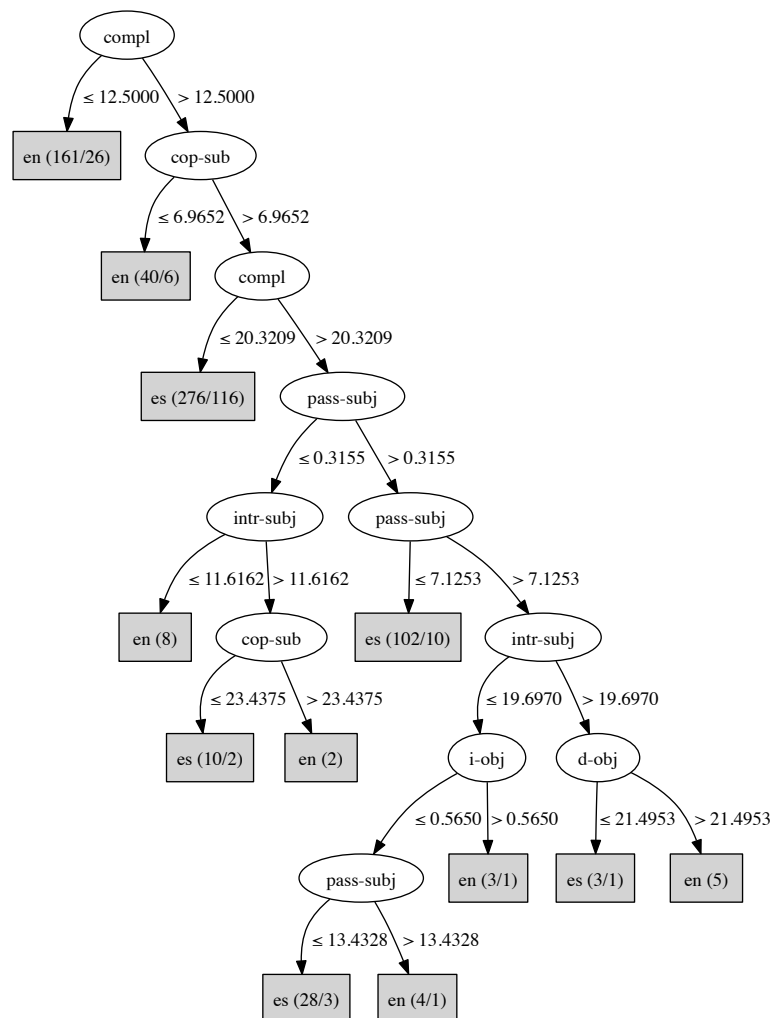


Figure 4.4: Lexical Argument Role C4.5 Decision Tree.

The passive subject attribute (*pass-subj*), is used three times in the classifier. The first test that considers it splits off a small number of test cases with low values and passes them on to further tests. This group consists of approximately equal numbers of native and nonnative cases. The other, much large group is passed onto another test, which also looks at the passive subject attribute. Here, the splitting causes a large number of cases with lower values to be immediately classified as nonnative. Deeper in the graph a pre-terminal

Table 4.5: Accuracy of Lexical Argument Role C4.5 Classifier

Nonnative	78.5%
Native	54.8%
Overall	66.7%
Overall C.I. 95%	63.0% — 70.3%

node again considers this attribute, and once again classifies low values as nonnative. This is not what one would expect based on Du Bois’s data for intransitive subjects. The most likely explanation is that the avoidance of lexical passive subjects by English learners is due to the avoidance of passives altogether. Butt and Benjamin [2004, 28.2.3] notes that the passive in Spanish, particularly the form of the passive which is exactly parallel to the English passive, employing the copular verb plus a non-finite verb form, is quite rarely used in Spanish, and that freedom of word order in Spanish allows simple fronting of an object, which is often why the passive is used in English. It is very likely that this is the source of L1-transfer, resulting in under-usage of the passive, and hence, the lexical passive subject, in learner English.

Du Bois also shows that the direct object role is more likely to be occupied by a lexical NP in English than in Spanish. In the decision tree the direct object attribute (*d-obj*), appears once. Cases with low values are immediately classified as nonnative and those with high values as native. This suggests that L1-interference leads English learners into underusing lexical direct objects.

To gauge the accuracy potential of using these attributes in constructing data models for classification, a 10 tree random forest classifier was trained using the lexical argument role attributes. Another such classifier was trained on all 3 attribute sets combined. The results from these classifiers is shown in Tables 4.6 and 4.7, respectively. Considering first Table 4.6, it does appear that the random forest classifier is able to better take advantage of the data, the average accuracy being boosted from 66.7% to 70.0%. However, the overlapping confidence intervals means that there is not yet sufficient data to conclude that the random forest is better than the C4.5 classifier. Table 4.7 shows that little is gained by

including the other sets of attributes. The valency attributes, as was already shown, are an ineffective basis for classification, and do not seem to improve when combined with other attributes. The failure of the addition of the lexical argument density attributes to improve classification must mean that these attributes do not contribute any information that is not already contained in the lexical role arguments.

Table 4.6: Accuracy of Lexical Argument Role Random Forest Classifier

Nonnative	76.6%
Native	63.2%
Overall	70.0%
Overall C.I. 95%	66.4% — 73.5%

Table 4.7: Accuracy of Combined Random Forest Classifier

Nonnative	76.6%
Native	63.9%
Overall	70.2%
Overall C.I. 95%	66.7% — 73.8%

References

- ALEJO GONZÁLEZ, R. 2010. L2 spanish acquisition of english phrasal verbs. In *Corpus-Based Approaches to English Language Teaching*, M. C. Campoy-Cubillo, B. Bellés-Fortuño, and M. L. Gea-Valor, Eds. Continuum International Publishing, Chapter 11.
- BIBER, D. AND XEPPE, R. 1998. Comparing native and learner perspectives on english grammar: a study of complement clauses. In *Learner English on Computer*, S. Granger, Ed. Addison Wesley Longman, Chapter 11, 148–158.
- BREIMAN, L. 2001. Random forests. In *Machine Learning*. 5–32.
- BUTT, J. AND BENJAMIN, C. 2004. *A New Reference Grammar of Modern Spanish* Fourth Ed. McGraw-Hill.

- CELCE-MURCIA, M. AND LARSEN-FREEMAN, D. 1999. *The Grammar Book, an ESL/EFL Teacher's Course* Second Ed. Heinle and Heinle Publishers.
- DAGUT, M. AND LAUFER, B. 1985. Avoidance of phrasal verbs—a case for contrastive analysis. *Studies in Second Language Acquisition* 7, 01, 73–79.
- DE MARNEFFE, M.-C. AND MANNING, C. D. 2008. Stanford typed dependencies manual.
- DU BOIS, J. W. 2003. Discourse and grammar. In *The New Psychology of Language: Cognitive and Functional Approaches to Language Structure*, M. Tomasello, Ed. Vol. 2. Lawrence Erlbaum Associates, Chapter 2, 47–87.
- HALL, M., FRANK, E., HOLMES, G., PFAHRINGER, B., REUTEMANN, P., AND WITTEN, I. H. 2009. The WEKA data mining software: An update. *SIGKDD Explorations* 11, 1.
- HULSTIJN, J. H. AND MARCHENA, E. 1989. Avoidance. *Studies in Second Language Acquisition* 11, 03, 241–255.
- KLEIN, D. AND MANNING, C. D. 2003. Accurate unlexicalized parsing. In *Proceedings of the 41st Meeting of the Association for Computational Linguistics*. 423–430.
- LAUFER, B. AND ELIASSON, S. 1993. What causes avoidance in L2 learning. *Studies in Second Language Acquisition* 15, 01, 35–48.
- LIAO, Y. AND FUKUYA, Y. J. 2004. Avoidance of phrasal verbs: The case of Chinese learners of English. *Language Learning* 54, 2, 193–226.
- QUINLAN, J. R. 1993. *C4.5: Programs for Machine Learning*. Morgan Kaufmann Publishers.
- WHITLEY, M. S. 1986. *Spanish/English Contrasts: A Course in Spanish Linguistics*. Georgetown University Press.