

A System to Distinguish Native and Nonnative Written English

By Philip White

Goals

- Develop classifiers to distinguish native written English from that written by L1-Spanish speakers.
- Focus on advanced learner English (i.e., don't rely on detecting errors).
- Design the classifier such that it could be used as the basis of a learning tool.

Resources Used

- Corpora of native and nonnative English
- Stanford Parser
- Weka Machine Learning Tools

Corpora

Corpus	Tokens Native	Tokens Nonnative
BROWN	57,809	0
ICE-HK	59,674	0
MICUSP	163,218	29,897
MSUELI	0	538
OANC	84,052	0
SPICLE	0	216,879
SULEC	0	39,254
WRICLE	0	96,247
ICE-CAN	25,225	2,070
Total	389,978	384,885

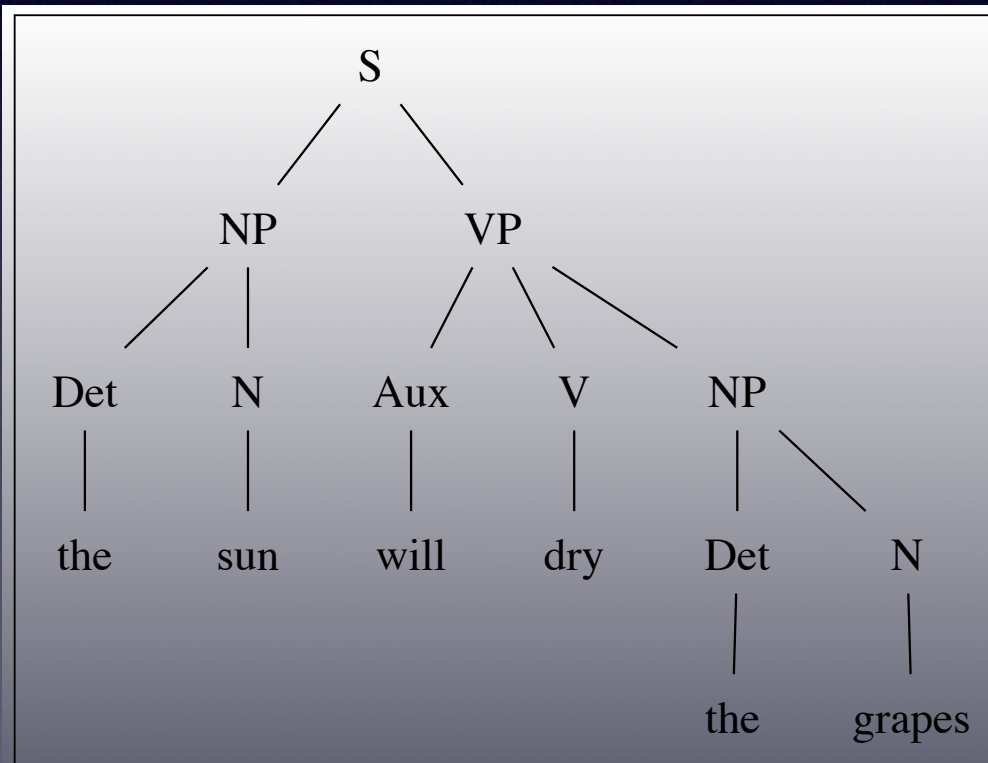
321 Native Samples and 321 Nonnative Samples

Stanford Parser

- Probabilistic parser
- Developed by Stanford
- Generates parse trees and dependency relationship.
- Parse trees can contain up to about 40 different node types.
- Parser generates 58 different dependencies.

Parse Trees

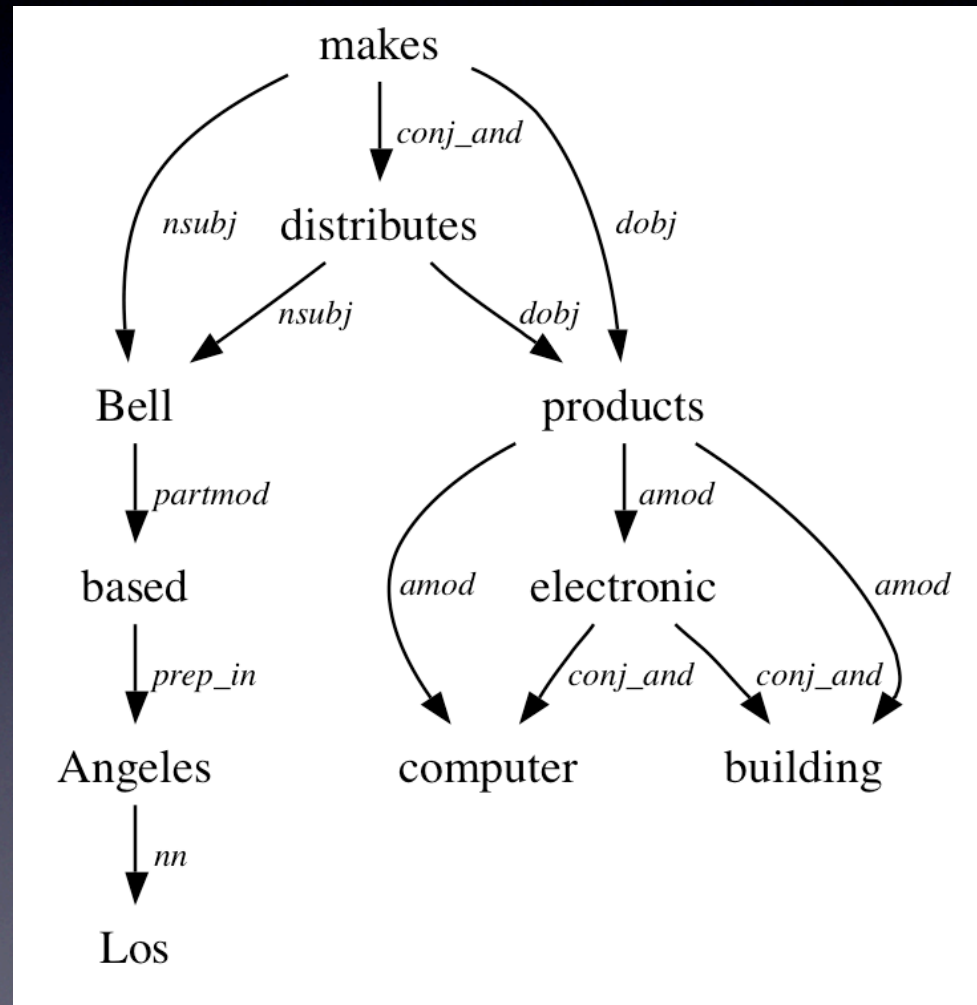
- a. $S \rightarrow NP VP$
- b. $NP \rightarrow (Det) N (PP)$
- c. $VP \rightarrow (Aux) V (NP) (AdvP)^n$
- d. $PP \rightarrow P NP$
- e. $AdvP \rightarrow \begin{Bmatrix} Adv \\ PP \end{Bmatrix}$



Dependencies

- Dependencies are 3-tuples: a relation, a governor, and a dependent.

nn(Angeles,Los)
prep_in(based,Angeles)
partmod(Bell,based)
nsubj(makes,Bell)
nsubj(distributes,Bell)
etc.



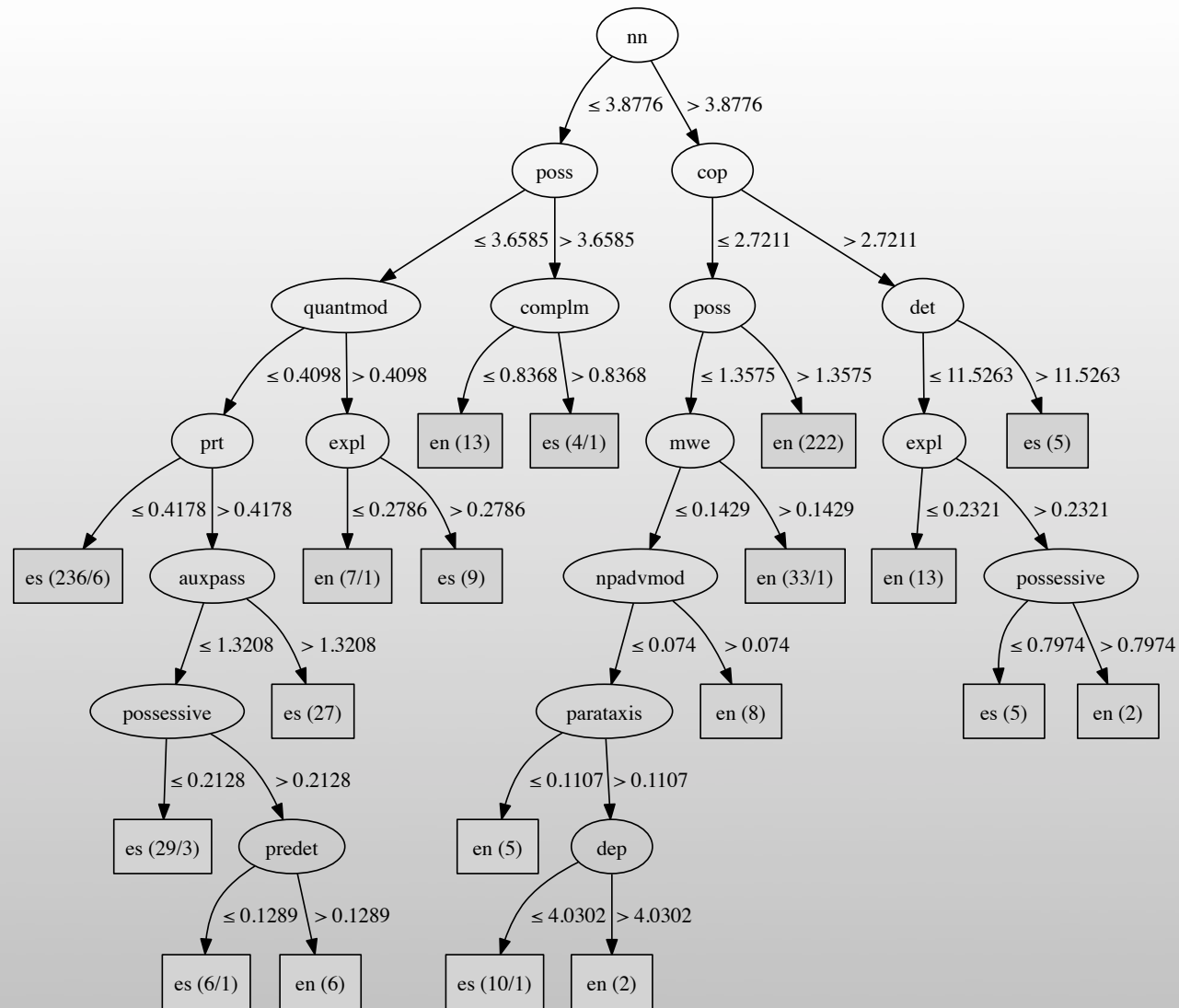
Weka

- Machine Learning
 - Classify instances into classes.
 - Instances are represented as lists of labeled values.
 - Each instance has a class value, which may be unknown.
 - Training and testing instances have known classes.

Classifiers

- C4.5 Decision Tree (J48)
 - Building Stage: constructs tree to divide training sets into subsets of a single class.
 - Pruning Stage: simplifies and generalizes tree.
- Random Forest
 - Uses multiple trees each using a subset of the total number of attributes.
 - Trees vote on best class.

Decision Tree



The Experiments

- Performed three experiments.
- Each experiment consisted of:
 - Extracting grammatical features from parse trees or dependency relations
 - Generating C4.5 classifier trees
 - Analyzing the features used in the trees to determine the linguistic significance of the features used
 - Measuring accuracy of classifiers (C4.5 and R.F.)

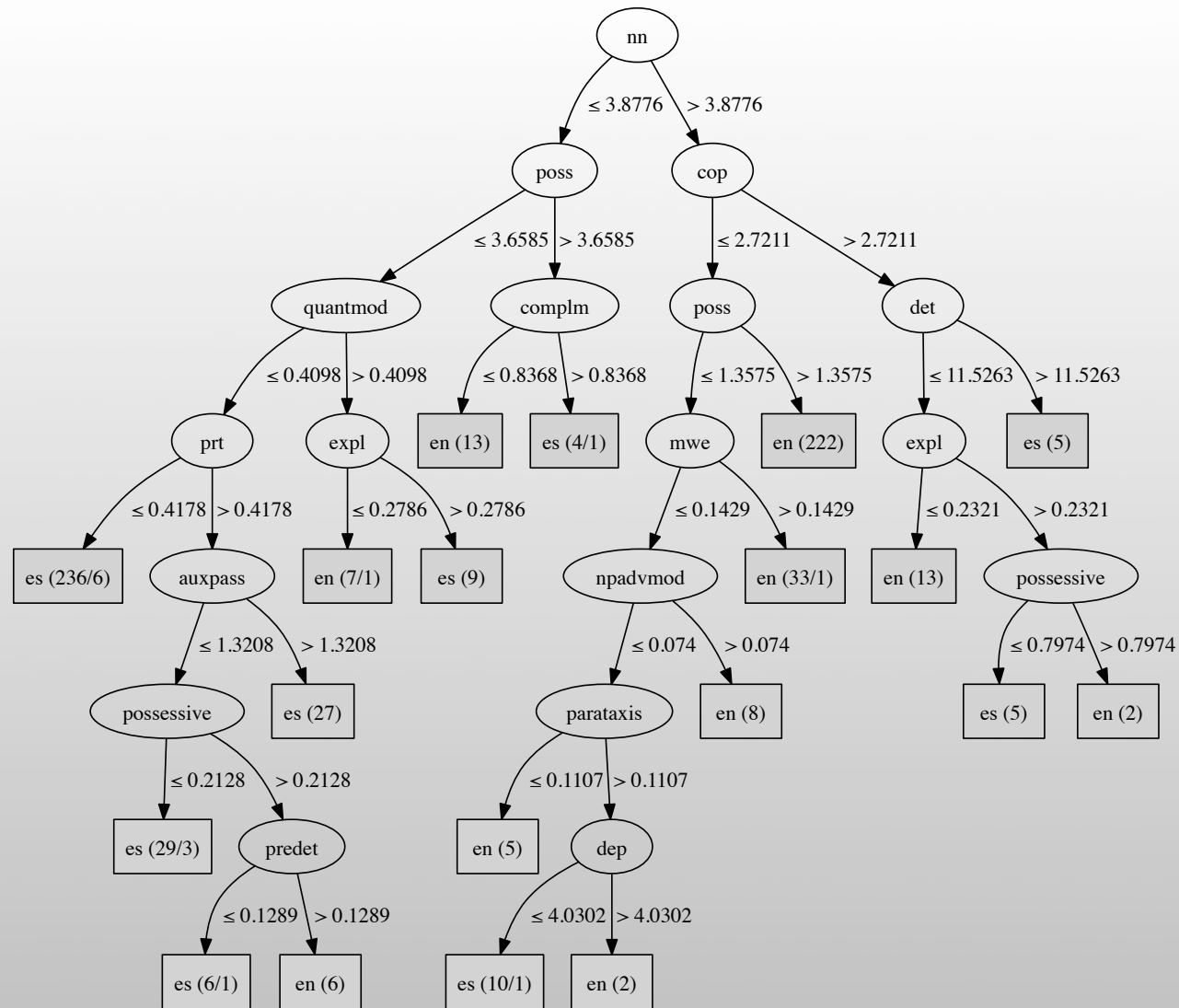
Dependency Experiment

- One attribute per dependency relation
- Value associated with each dependency is the number of occurrences of that relation divided by the number of occurrences of all relation.

Dependency Experiment

Label	Dependency
auxpass	passive auxiliary
complm	complementizer
cop	copula
dep	generic dependency
det	determiner
expl	expletive
mwe	multi-word express.
nn	noun compound
npadvmod	np as adv modifier
parataxis	parataxis
poss	possession modifier
possessive	possessive modifier
predet	predeterminer
prt	phrasal verb particle
quantmod	quantifier

Decision Tree



Dependency Experiment

C4.5 Tree Accuracy	
Nonnative	89.7%
Native	87.9%
Overall	88.8%
95% C.I.	86.3% – 91.2%

Random Forest Accuracy	
Nonnative	96.0%
Native	91.9%
Overall	93.9%
95% C.I.	92.1% – 95.8%

Verbal Argument Experiment

- Verbal argument roles:
 - Subject, object, indirect object, subject complement, patient, etc.
 - Derivable from dependency graphs
- An argument can either be a pronoun or a lexical noun.
- From this a number of features can be derived.

Du Bois, J.W. 2003. *Discourse and Grammar*.

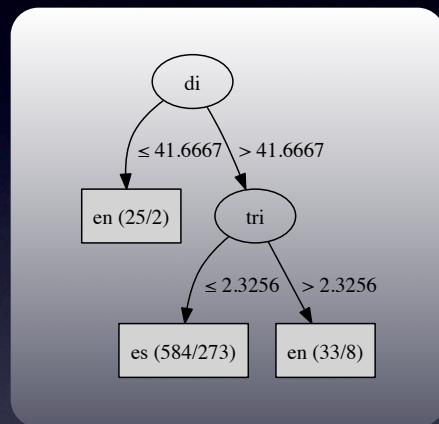
Verbal Arguments

- a. In nineteen ninety-seven Roland Lemon was elected president.
- b. ... before his father gave him the rest of it.
- c. They state that climate generally predicts that temperatures should rise...
- d. Before Yeltsin appointed him the deputy Prime Minister ...
- e. That she lied was suspected by everyone.
- f. What she said is not true.
-
- The diagram illustrates syntactic relations for six sentences (a-f). Arrows indicate the direction of the relation, and labels indicate the type of relation.
- a.** In nineteen ninety-seven Roland Lemon was elected president.
 - nsbjpass: from "was elected" to "Roland Lemon"
 - dobj: from "elected" to "president"
 - b.** ... before his father gave him the rest of it.
 - nsbj: from "gave" to "his father"
 - iobj: from "gave" to "him"
 - dobj: from "gave" to "the rest of it"
 - c.** They state that climate generally predicts that temperatures should rise...
 - nsbj: from "state" to "They"
 - ccomp: from "state" to "that climate generally predicts that temperatures should rise..."
 - nsbj: from "predicts" to "climate"
 - ccomp: from "predicts" to "that temperatures should rise..."
 - d.** Before Yeltsin appointed him the deputy Prime Minister ...
 - nsbj: from "appointed" to "Yeltsin"
 - nsbj: from "appointed" to "him the deputy Prime Minister ..."
 - xcomp: from "appointed" to "Before"
 - e.** That she lied was suspected by everyone.
 - nsbj: from "was suspected" to "That she lied"
 - csbjpass: from "was suspected" to "by everyone"
 - f.** What she said is not true.
 - nsbj: from "is not true" to "What she said"
 - csbj: from "is not true" to "What she said"

Verbal Arguments

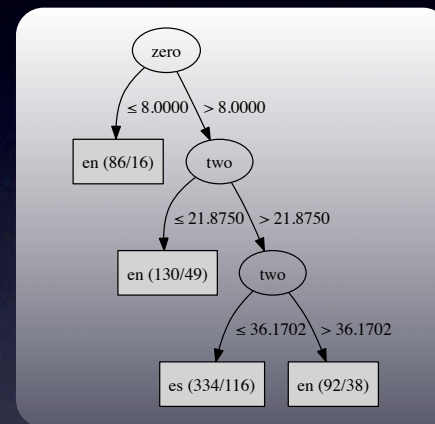
Pronominal Forms					
other	another	else	same	one	
this	that	these	those	what	
myself	yourself	herself	himself	itself	
ourselves	yourselves	themselves	oneself		
mine	yours	hers	his	ours	theirs
me	you	her	him	it	us
them	I	she	he	we	they

Verbal Argument Experiment



C4.5 Tree Accuracy Verb Valency

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



C4.5 Tree Accuracy Lexical Argument Density

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

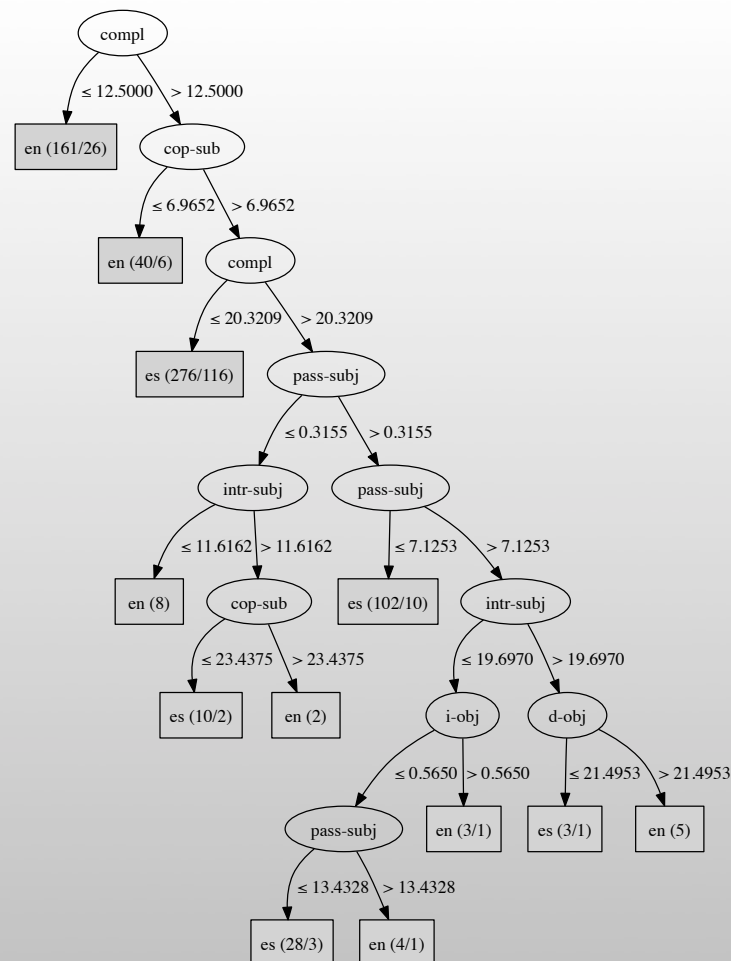
Verbal Argument Experiment

C4.5 Tree Accuracy Lexical Argument Role

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

Random Forest Accuracy Lexical Argument Role

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



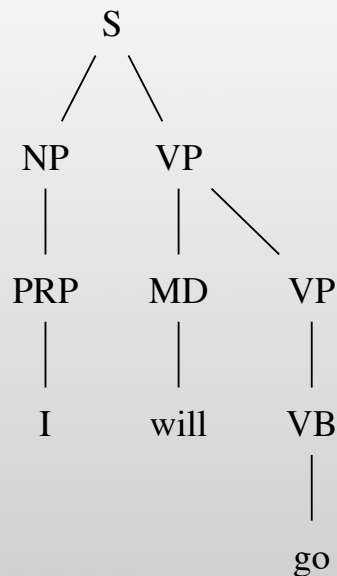
Verbal Argument Experiment

Random Forest Accuracy Combined Attributes	
Nonnative	76.6%
Native	63.9%
Overall	70.2%
95% C.I.	66.7% – 73.8%

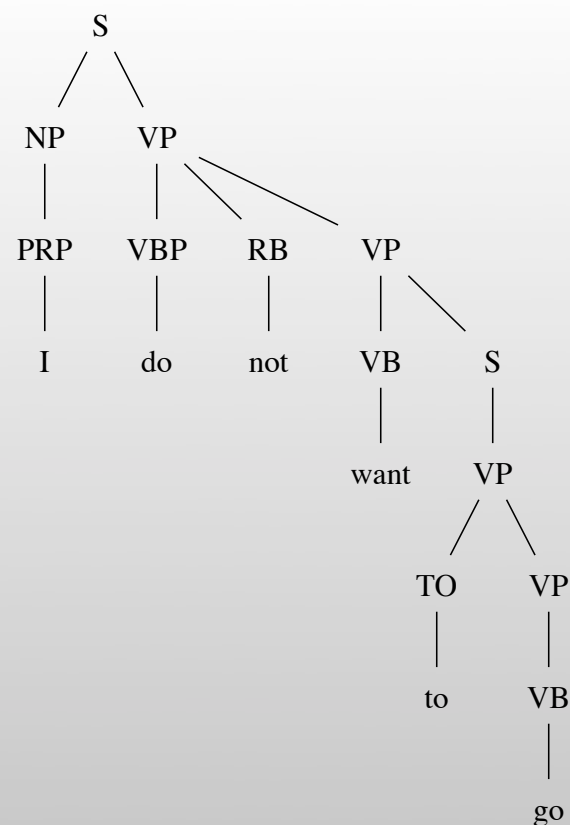
Verb Form Experiment

- Extract verb details from parse trees
- Identifies the following features:
 - Tense, aspect, voice, core modals, phrasal modals, the helping verbs *be*, *get*, *have*, *do*, and the negative particle *not*
- Also considers common verbs

Verb Form Experiment

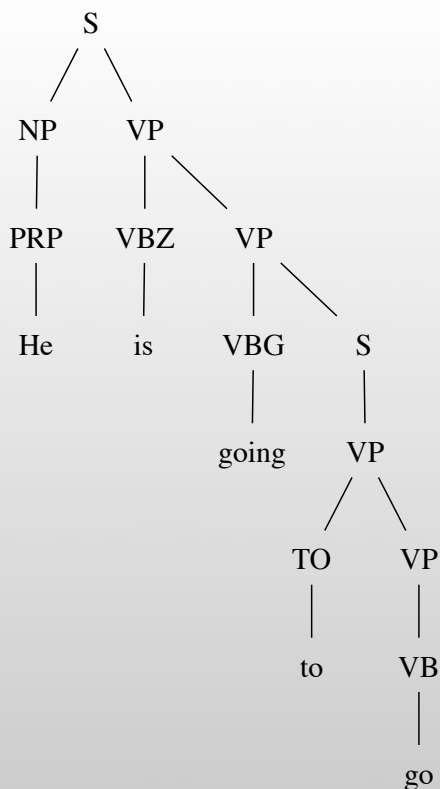


Core Modal

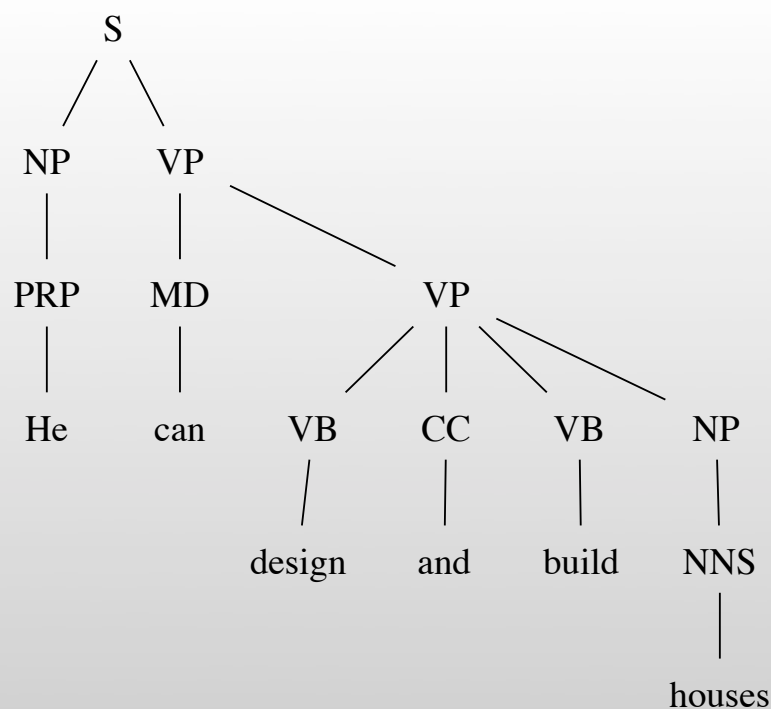


Verb with *not* and *do*

Verb Form Experiment



Phrasal Modal



Core Modal with Two Main Verbs

Verb Form Experiment

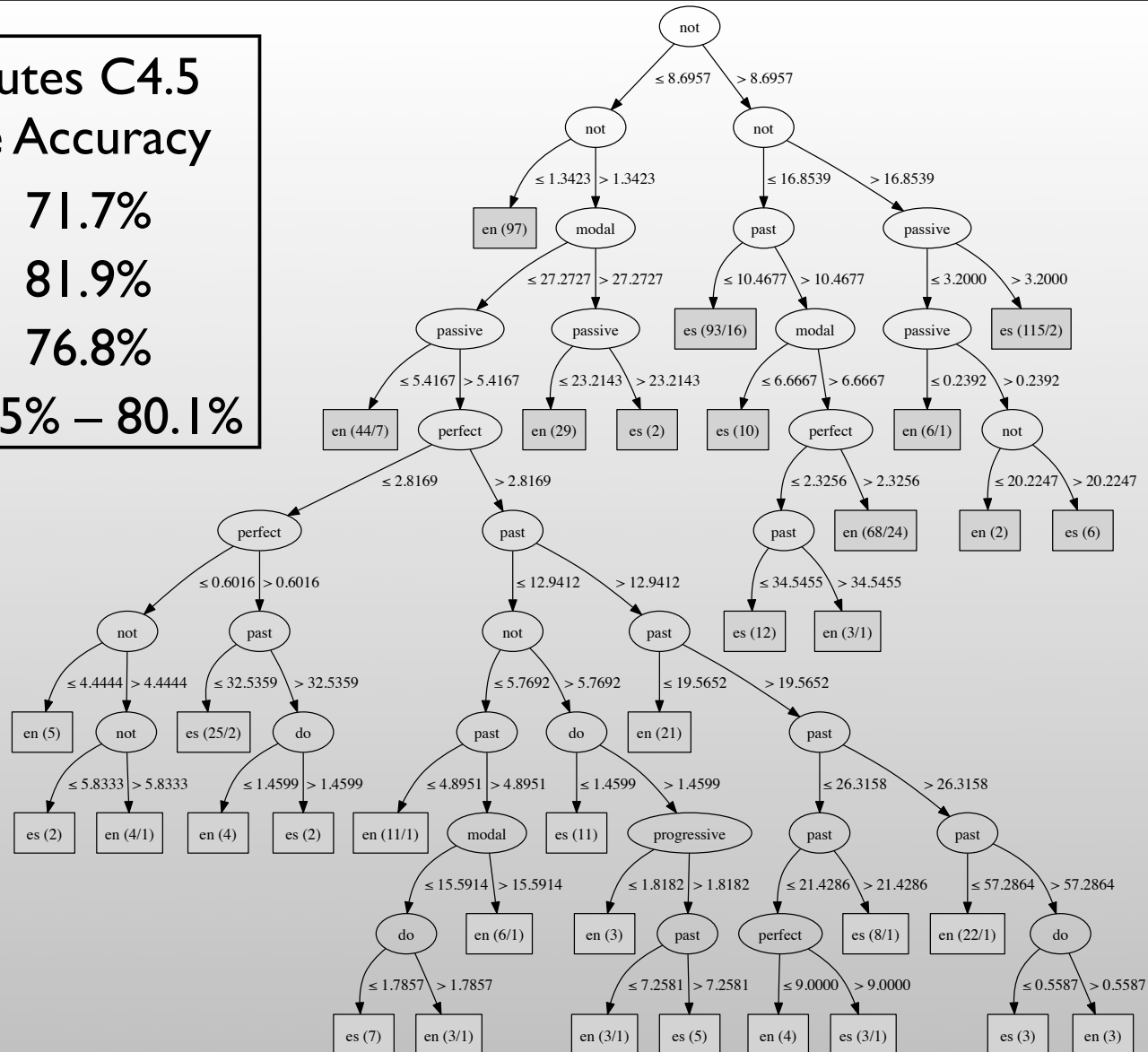
Verbal Attributes C4.5
Decision Tree Accuracy

Nonnative 71.7%

Native 81.9%

Overall 76.8%

95% C.I. 73.5% – 80.1%



Verb Form Experiment

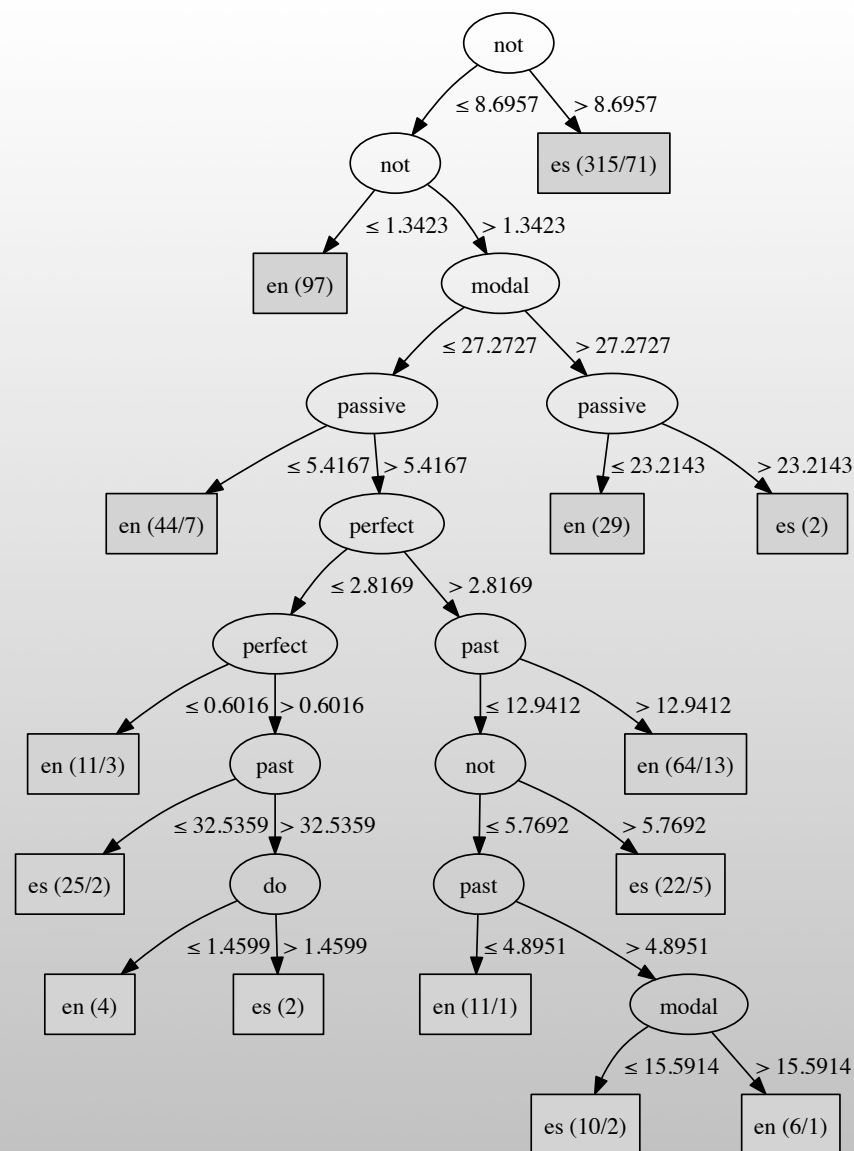
Verbal Attributes Simplified C4.5 Decision Tree Accuracy

Nonnative 71.3%

Native 79.8%

Overall 75.5%

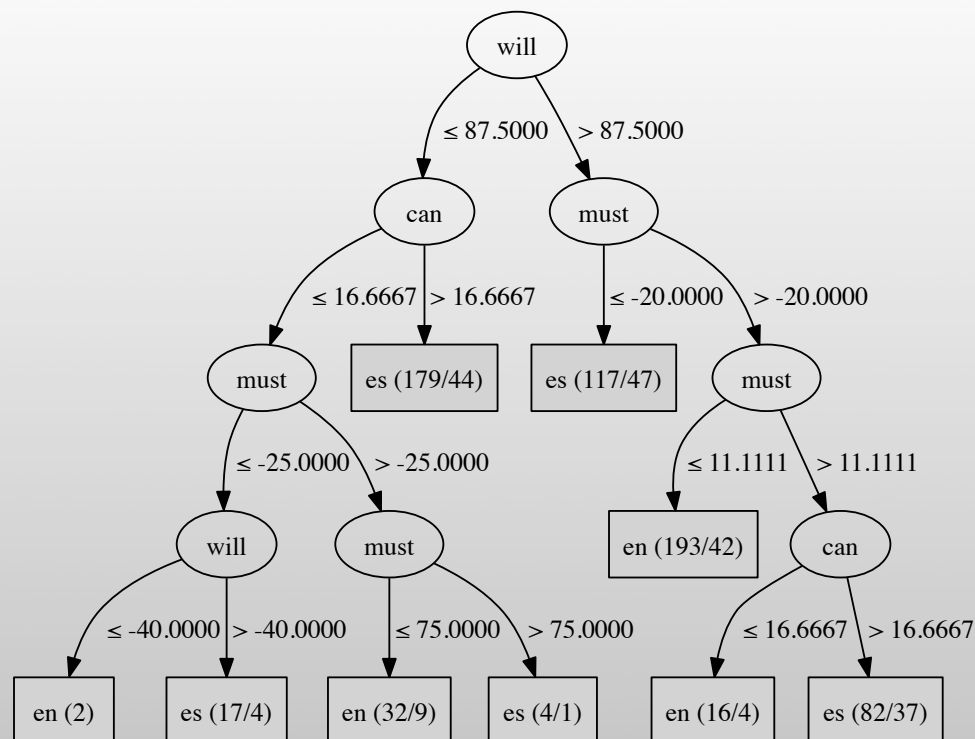
95% C.I. 72.2% – 78.9%



Verb Form Experiment

Core vs Phrasal Modal C4.5 Decision Tree Accuracy

Nonnative	70.7%
Native	61.4%
Overall	66.0%
95% C.I.	62.3% – 69.7%



Verb Form Experiment

This experiment uses the relative frequency of finite usages of certain common verbs as classifier attributes.

make	use	take	see	say
go	become	believe	give	feel
come	find	think	know	look
seem	want	get	live	work

Ringbom, H. 1998. *High-Frequency Verbs in the ICLE Corpus*.

High Frequency Verb C4.5 Decision Tree Accuracy	
Nonnative	70.7%
Native	61.4%
Overall	66.0%
95% C.I.	62.3% – 69.7%

Verb Form Experiment

Combined Verb Attributes Random Forest Accuracy

Nonnative	80.4%
Native	77.3%
Overall	78.8%
95% C.I.	75.7% – 82.0%

Combined Verb Attributes C4.5 Decision Tree Accuracy

Nonnative	90.3%
Native	85.4%
Overall	87.9%
95% C.I.	85.3% – 90.4%

Other Possible Attributes

- Syntactic complexity (parse tree depth, phrase nesting, etc.)
- Types of phrases
- Vocabulary (of structural words in particular)
- etc.

Learning Tool

- Use classifiers to determine if text is identifiably nonnative.
- Inform user which features were responsible for classification as nonnative.
- Show user which parts of the text exhibit these features.
- User edits and reevaluates until text is classified as native.

Learner Tool Mock-up

Keeling et al. 2004), a mesoscale (the spatial scale determined by the aggregation of hosts into communities) and a macroscale (the regional spatial scale defined by set of communities, Keeling et al. 2004, and the connections among them).

In conjunction with the characterization of the spatial scales, the dynamics of the disease also depends on a precise metapopulation description. The parameterization of a metapopulation model consist of estimation of: patch areas, including their spatial location; pairwise distances between them; presence and absence of the species under study; distribution of migrating distances; colonization ability and critical patch area. Each of these parameters may be mapped to epidemiological variables, in particular the critical patch area can be easily linked to the critical community size (Keeling, 1997). The patch areas, distances and distribution of migrating distances, however, are strongly dependent on the transmission of the disease, and the study of the spatial patterns formed during epidemics may provide empirical evidence to determine their realistic values.

In order to find accurate parameter values for spatially explicit model for cholera dynamics, different methodological approach may be used including Point Pattern Analysis, Geostatistical Analysis, and determination of the Critical Community Size, among others.

Previous

Next

Reevaluate

Issues

Information

Core Modal Overuse

Inflected Genitive Underuse

Phrasal Verb Underuse

Pronominal Argument Overuse

Passive Voice Overuse

Consider using phrasal modals in place of certain core modals. Read these resources for information on the subtle semantic differences between core modals and phrasal modals:

<http://www.esl-helper.com/core-modals>

<http://www.esl-helper.com/phrasal-modals>

<http://www.esl-helper.com/future-tense>

Other Possible Applications

- Identifying plagiarism among ESL student.
- Determining first language of the writer of a text sample (forensics?).