

Using Cognitive Biases to Guide Feature Set Selection

Claire Cardie*

Department of Computer Science
University of Massachusetts
Amherst, MA 01003
cardie@cs.umass.edu

Abstract

Although learning is a cognitive task, machine learning algorithms, in general, fail to take advantage of existing psychological limitations. In this paper, we use a learning task from the field of natural language processing and examine three well-known cognitive biases for human information processing: 1) the tendency to rely on the most recent information, 2) the heightened accessibility of the subject of a sentence, and 3) short term memory limitations. In a series of experiments, we modify a baseline instance representation in response to these limitations and show that the overall performance of the learning algorithm improves as increasingly more cognitive biases and limitations are explicitly incorporated into the instance representation.

Introduction

Inductive concept acquisition has always been of primary interest for researchers in the field of machine learning. In this task, a system typically learns one or more concepts by analyzing a set of examples (and possibly counterexamples) of the concepts. In fact, a number of systems for the acquisition of concepts now exist (e.g., ID3 (Quinlan, 1979), ARCH (Winston, 1975), COBWEB (Fisher, 1987), UNIMEM (Lebowitz, 1987)). Independently, psychologists, psycholinguists, and cognitive scientists have examined the effects of numerous psychological limitations on human information processing. However, despite the fact that concept learning is a basic cognitive task, most machine learning systems for concept formation fail to exploit these limitations and make no attempt to model human concept learning.

*This research was supported by the Office of Naval Research, under a University Research Initiative Grant, Contract No. N00014-86-K-0764, NSF Presidential Young Investigators Award NSFIST-8351863 awarded to Wendy Lehnert, and the Advanced Research Projects Agency of the Department of Defense monitored by the Air Force Office of Scientific Research under Contract No. F49620-88-C-0058.

In this paper, we show that the explicit encoding of known cognitive biases into the training instance representation can improve the performance of the learning algorithm for cognitively-based learning tasks. More specifically, we use a well-known concept acquisition system and focus on a single learning task from the field of natural language processing (NLP). After training the system using a baseline instance representation, we modify the representation in response to three cognitive biases: 1) the tendency to rely on the most recent information, 2) the heightened accessibility of the subject of a sentence, and 3) short term memory limitations. Each modification explicitly incorporates one or more cognitive biases into the feature set. In a series of experiments, we compare each of the modified instance representations to the baseline.

Finding the Antecedents of Relative Pronouns

Although the use of cognitive biases to guide feature set selection is a domain-independent technique, we will use a learning task from NLP to illustrate the performance of the technique throughout the paper. Our task for the machine learning system is the following: Given a sentence with the relative pronoun "who," learn to recognize the phrase or phrases that represent the relative pronoun's antecedent. (Note: This paper focuses only on the technique with respect to machine learning issues. For a detailed discussion of the viability of this approach for the disambiguation of relative pronouns from the NLP perspective, see (Cardie, 1992a) and (Cardie, 1992b)). Finding the antecedents of relative pronouns is a crucial task for natural language systems because the antecedent must be made available to the subsequent clause where it implicitly fills the *actor* or *object* roles.¹ Consider the following example:

Igor shook hands with *the skater who*
beat him in the race.

¹In practice, the antecedent of "who" sometimes fills semantic roles other than the actor or object.

A correct semantic interpretation of this sentence should include the fact that “the skater” is the actor of “beat” even though the phrase does not appear in the embedded clause. Only after the natural language system associates “the skater” with “who” can it make this inference. Locating the antecedent of “who” may initially appear to be an easy problem because the antecedent often immediately precedes the word “who.” Unfortunately, this is not always the case as shown in S1 and S2 of Figure 1. Even when the antecedent does immediately precede the relative pronoun, it does not appear in a consistent syntactic constituent. In S3, for example, the antecedent is the subject of the preceding clause; in S4, it is the direct object; in S5, it is the object of a preposition. Furthermore, the antecedent of “who” may contain more than one phrase. In S6, for example, the antecedent is a conjunction of three phrases and in S7, either “our sponsors” or its appositive “Gatorade and GE” is a semantically valid antecedent. Occasionally, there is no apparent antecedent at all (e.g., S8).

- | |
|--|
| <p>S1. <i>The woman</i> from Philadelphia who played soccer was my sister.</p> <p>S2. I spoke to <i>the man</i> in the black shirt and green hat over in the far corner of the room who demanded to meet the skiers.</p> <p>S3. <i>The skater</i> who won the medal was from Japan.</p> <p>S4. I saw <i>the skater</i> who won the medal.</p> <p>S5. Igor ate dinner with <i>the skater</i> who won the medal.</p> <p>S6. I'd like to thank <i>Nike, Reebok, and Adidas</i>, who provided the uniforms.</p> <p>S7. I'd like to thank <i>our sponsors, Gatorade and GE</i>, who provide financial support.</p> <p>S8. We wondered who would win the race.</p> |
|--|

Figure 1: Antecedents of “who”

Despite these ambiguities, we will describe how a machine learning system can learn to locate the antecedent of “who” given a description of the clause that precedes it. In effect, we are teaching the system to recognize the “relative pronoun antecedent” concept. More importantly, we will show that performance of the learning system improves as the instance description explicitly encodes increasingly more cognitive limitations and cognitive biases.

COBWEB and the Representation of Training Instances

For our experiments we chose COBWEB (Fisher, 1987) – a well-known concept formation system that is one of a relatively small number of concept acquisition systems designed to model some aspects

of human concept learning.² Given a set of training instances, COBWEB discovers a classification scheme that covers the instances. Instead of forming concepts at a single level of abstraction, however, COBWEB organizes instances into a classification hierarchy where leaves represent instances and internal nodes represent concepts that increase in generality as they approach the root of the tree. In addition, COBWEB’s construction of the hierarchy is cognitively economical in that new objects are incrementally added to the hierarchy as they arrive. To evaluate the concepts it creates, COBWEB employs the *category utility* metric (Gluck, & Corter, 1985) – a measure developed in psychological studies of basic level categories.

COBWEB takes as input a set of training instances described as a list of attribute-value pairs. Because the antecedent of a relative pronoun usually appears as one or more phrases in the clause preceding “who,” the attribute-value pairs in each training case represent the constituents that precede “who.” At first glance, it may seem that only syntactic information needs to be encoded. However, finding the antecedent of a relative pronoun actually requires the assimilation of syntactic and semantic knowledge. For this reason, each *constituent attribute-value pair* takes the following form:

- The *attribute* describes the syntactic class and position of the phrase.
- The *value* provides its semantic classification.

Consider, for example, the sentences in Figure 2. In the training instance for S1, we represent “the man” with the attribute-value pair (*s human*) because it is the subject of the sentence and the noun “man” is human. We represent “from Oklahoma” with the pair (*s-pp1 location*) because it is the first prepositional phrase that follows the subject and “Oklahoma” is a location. All noun phrases are described by one of seven general semantic features: human, proper-name, location, entity, physical-target, organization, and weapon.³ When clauses contain conjunctions and appositives, each phrase in the construct is labelled separately. In S2, for example, the real direct object of “thank” is the conjunction “Nike and Reebok.” However, in our instance representation, “Nike” is tagged as the direct object (*do*) and “Reebok” as the

²The COBWEB/3 system was provided by Kevin Thompson, NASA Ames Research Center.

³These features are specific to the domain from which the training instances were extracted. A different set would most likely be required for nouns in a different domain.

experiments, 2 sets were used for training and the third reserved for testing. The results are shown in Figure 4 and indicate that COBWEB finds the correct antecedent of “who” an average of 59% of the time when using the baseline instance representation. In the next two sections, we modify this baseline representation in response to three cognitive biases and show the results of these modifications on COBWEB’s performance.

Exp #	Training Sets (# instances)	Test Set (# instances)	Baseline Rep
1	set1 + set2 (170)	set3 (71)	63%
2	set2 + set3 (159)	set1 (82)	47%
3	set1 + set3 (153)	set2 (88)	66%

Figure 4: Baseline Results (% correct)

Incorporating the Recency Bias

In processing language, people consistently show a bias towards the use of the most recent information (e.g., (Kimball, 1973), (Frazier, & Fodor, 1978), (Gibson, 1990)). In particular, the mechanisms people use for finding the antecedents of pronouns and missing subjects have been investigated in a series of recent experiments (see (Nicol, 1988)). The results show that in locating antecedents during language processing, people consider all noun phrases preceding the pronoun starting with the most recent noun phrase and working backwards to the most distant noun phrase.

We translate this recency bias into representational changes for the training instances in two ways. First, we label the constituent attribute-value pairs with respect to the relative pronoun. This establishes a right-to-left labelling rather than the left-to-right labelling of the baseline. In Figure 5, for example, “in Congress” receives the attribute *pp1* because it is a prepositional phrase one position to the left of “who.” Similarly, “the hardliners” receives the attribute *np2* because it is a noun phrase two positions to the left of “who.” Notice, however, that the subject of the sentence retains its original *s* attribute. We based this decision on studies that indicate that the subject of a sentence remains highly accessible even at the end of a sentence (e.g., (Gernsbacher, Hargreaves, & Beeman, 1989)). Consider the following sentences: 1) “it was a message from *the hardliners* in Congress, who...” and

Message Understanding System Evaluation and Message Understanding Conference (Sundheim,1991).

<p>Sentence: [It] [was] [the hardliners] [in Congress] who...</p> <p>Baseline Representation: (<i>s entity</i>) (<i>v t</i>) (<i>do human</i>) (<i>do-pp1 entity</i>) (<i>antecedent ((do))</i>)</p> <p>Right-to-Left Labelling: (<i>s entity</i>) (<i>v t</i>) (<i>np2 human</i>) (<i>pp1 entity</i>) (<i>antecedent ((np2))</i>)</p> <p>Duplicate Information: (<i>s entity</i>) (<i>v t</i>) (<i>do human</i>) (<i>do-pp1 entity</i>) (<i>most-recent entity</i>) (<i>part-of-speech prep-phrase</i>) (<i>antecedent ((do))</i>)</p>

Figure 5: Incorporating the Recency Bias

2) “it was from *the hardliners* in Congress who ...”. The right-to-left labelling tags the antecedents in each sentence with the same attribute (i.e., *pp2*), indicating the similarity of the examples with respect to the location of the relative pronoun antecedent. In the baseline representation, however, the antecedents retain distinct attributes – *do-pp1* and *v-pp1*, respectively.

Alternatively, given the baseline instance representation, we can incorporate the recency bias by including more than one attribute-value pair for the most recent information. Figure 5 also shows this second representational change. The most recent constituent (“in Congress”) is represented three times⁷: 1) as a constituent attribute-value pair – (*do-pp1 entity*), 2) as the most recent constituent – (*most-recent entity*), and 3) via its part of speech – (*part-of-speech prep-phrase*). In this representation, we also allow the antecedent attribute-value pair to refer to the more general *most-recent* constituent rather than the equivalent, but more specific, constituent attribute-value pair. If, for example, the antecedent in Figure 5 had been *do-pp1*, it would become *most-recent* in the new representation.

The results of experiments that use each of these representations separately and in a combined form are shown in Figure 6. In this table, the MR1 representation used the right-to-left labelling, the MR2 representation included extra information about the most recent constituent, and the MR1+MR2 representation combined both the right-to-left labelling and the duplicate information formats. In general, it is clear that incorporating the recency bias into the instance representation improves performance. On average, the right-to-left labelling

⁷We used all information about the most recent constituent readily available from the parser.

Exp #	Training Sets (# of instances)	Test Set (# of instances)	Baseline Rep	MR1: R-to-L Labelling	MR2: Duplicate Info	MR1 + MR2
1	set1 + set2	set3	63%	75%	83%	84%
2	set2 + set3	set1	47%	62%	65%	73%
3	set1 + set3	set2	66%	66%	71%	74%

Figure 6: Experiments Using the Recency Bias (% correct)

increased the percentage of correctly identified antecedents from 59% to 68% while including extra information for the most recent constituent increased the percentage correct to 73%. The best results, however, occurred using the combined representation, where the percentage correct increased to an average of 77%.

Incorporating the Short Term Memory Bias

Psychological studies have determined that people can keep at most seven plus or minus two facts in short term memory (Miller, 1956). More recently, Daneman and Carpenter ((Daneman, & Carpenter, 1980), (Daneman, & Carpenter, 1983)) show that working memory capacity affects a subject's ability to find the referents of pronouns over varying distances. Also, King and Just (King, & Just, 1991) show that differences in working memory capacity can cause differences in the reading time and the comprehension of certain classes of relative clauses. Moreover, it has been hypothesized that language learning in humans is successful precisely because limits on information processing capacities allow children to ignore much of the linguistic data they receive (see (Newport, 1990)).

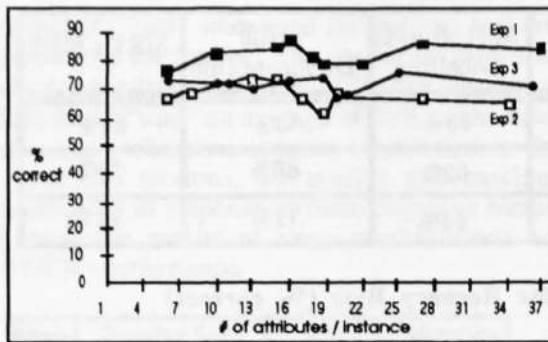
COBWEB, however, clearly does not make use of short term memory (STM) limitations either in its learning algorithm or in its attribute-value instance representation. Each training and test instance has to be normalized with respect to all attributes across the training instances.⁸ In the baseline representation, this normalization resulted in instances of 35 attribute-value pairs as compared to an average of 5 attribute-value pairs in the original, unnormalized instances. The short term memory bias implies that not all of the 35 features should be retained for the task of finding relative pronoun antecedents. In an attempt to incorporate this limitation, we ran a series

⁸Our fixed feature set includes every attribute that appears in the training set. To create a training instance, we generate a unique value for any missing attribute, i.e., for any attribute that is irrelevant for the instance.

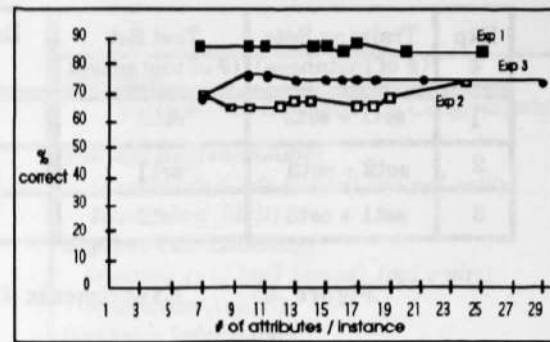
of experiments using instances with successively fewer features. We let $n = 1, 2, 3, 4, 5, 7, 9, 15, 20$, and 50, and included in the training and test instances only those features that occurred at least n times in the original, unnormalized instances of the training set. As n increases, the number of attributes per instance decreases.

When this STM cutoff was applied to the baseline representation, the performance of the learning algorithm gradually declined as n increased. The percentage correct declined from 63% to 37%, from 47% to 19%, and from 66% to 41% for experiments 1, 2, and 3, respectively. Although the STM cutoff did not improve performance when applied to the MR1 training sets that use the right-to-left labelling, the decline in percentage correct was not nearly so drastic. For experiment 1 (originally 75% hit rate), the percentage correct never dropped below 69%. For experiment 2, results ranged from 62% (with no cutoff) to 49%; and in experiment 3 (originally 66% correct), results ranged from 51% to 67% correct.

Figure 7 shows the results of the STM cutoff for the instance representations of MR2 (extra information for most recent phrase) and MR1+2 (right-to-left labelling and extra information for most recent phrase). In these experiments, the STM bias actually improved COBWEB's performance. In the MR2 experiments, the original hit rate for experiment 1 increased from 83% (37 attributes / instance) to 87% at $n = 7$ (16 attributes / instance). In experiment 2, the percentage correct moved from 65% (34 attributes / instance) in the original representation to 74% at $n = 9$ (13 attributes / instance). In experiment 3, the percentage correct increased from 71% (36 attributes / instance) in the original representation to 76% at $n = 2$ (25 attributes / instance). Similar results occurred for the MR1+2 instance representation. There were increases from 84% (25 attributes / instance) to 87% (17 attributes / instance, $n = 4$) and from 74% (29 attributes / instance) to 76% (10 attributes / instance) for experiments 1 and 3, respectively. Performance for experiment 2, however, declined.



(a) Applying STM Bias to MR2



(b) Applying STM Bias to MR1+2

Figure 7: Experiments Using the STM Bias

Conclusions

Based on the preliminary experiments presented in the last three sections, we conclude that explicit incorporation of cognitive biases into the instance representation can greatly improve learning algorithm performance. In addition, although the technique was tested on only one task from NLP, the use of cognitive biases to guide feature set selection is a domain-independent technique that can be applied to any cognitively-based learning task. It is clear, however, that further experimentation is required to explore the effects of additional cognitive limitations, to determine the biases that work well together, and to find the correct parameters for those biases. Finally, further research is required before we can use cognitive biases to automate, rather than guide, feature set selection.

References

- Cardie, C. (1992a). Corpus-Based Acquisition of Relative Pronoun Disambiguation Heuristics. *Proceedings, 30th Annual Meeting of the Association for Computational Linguistics*. University of Delaware, Newark.
- Cardie, C. (1992b). Learning to Disambiguate Relative Pronouns. *Proceedings, Ninth National Conference on Artificial Intelligence*. San Jose, CA.
- Daneman, M., & Carpenter, P. A. (1980). Individual differences in working memory and reading. *Journal of Verbal Learning and Verbal Behavior*, 19, 450-466.
- Daneman, M., & Carpenter, P. A. (1983). Individual differences in integrating information between and within sentences. *Journal of Experimental Psychology: Learning, Memory, and Cognition*, 9, 561-584.
- Fisher, D. H. (1987). Knowledge Acquisition Via Incremental Conceptual Clustering. *Machine Learning*, 2, 139-172.
- Frazier, L., & Fodor, J. D. (1978). The sausage machine: A new two-stage parsing model. *Cognition*, 6, 291-325.
- Gernsbacher, M. A., Hargreaves, D. J., & Beeman, M. (1989). Building and accessing clausal

representations: The advantage of first mention versus the advantage of clause recency. *Journal of Memory and Language*, 28, 735-755.

Gibson, E. (1990). Recency preferences and garden-path effects. *Proceedings, Twelfth Annual Conference of the Cognitive Science Society*. Cambridge, MA.

Gluck, M. A., & Corter, J. E. (1985). Information, uncertainty, and the utility of categories. *Proceedings, Seventh Annual Conference of the Cognitive Science Society*. Lawrence Erlbaum Associates.

Kimball, J. (1973). Seven principles of surface structure parsing in natural language. *Cognition*, 2, 15-47.

King, J., & Just, M. A. (1991). Individual differences in syntactic processing: the role of working memory. *Journal of Memory and Language*, 30, 580-602.

Lebowitz, M. (1987). Experiments with Incremental Concept Formation: UNIMEM. *Machine Learning*, 2, 103-138.

Lehnert, W., Cardie, C., Fisher, D., Riloff, E., & Williams, R. (1991). University of Massachusetts: Description of the CIRCUS System as Used for MUC-3. *Proceedings, Third Message Understanding Conference (MUC-3)*. San Diego, CA. Morgan Kaufmann Publishers.

Miller, G. A. (1956). The magical number seven, plus or minus two: Some limits on our capacity for processing information. *Psychol. Review*, 63(1).

Newport, E. (1990). Maturation Constraints on Language Learning. *Cognitive Science*, 14, 11-28.

Nicol, J. (1988). *Coreference processing during sentence comprehension*. Ph. D. Thesis. Massachusetts Institute of Technology.

Quinlan, J. R. (1979). Discovering Rules from Large Collections of Examples: A Case Study. In D. Michie (Ed.), *Expert Systems in the Microelectronics Age*. Edinburgh: Edinburgh University Press.

Sundheim, B. M. (1991). Overview of the Third Message Understanding Evaluation and Conference. *Proceedings, Third Message Understanding Conference (MUC-3)*. San Diego, CA. Morgan Kaufmann Publishers.

Winston, P. H. (1975). Learning Structural Descriptions from Examples. In P. H. Winston (Ed.), *The Psychology of Computer Vision*. New York: Mc Graw-Hill.