

What Is This, Anyway: Automatic Hypernym Discovery

Alan Ritter and Stephen Soderland and Oren Etzioni

Turing Center

Department of Computer Science and Engineering

University of Washington

Box 352350

Seattle, WA 98195, USA

{aritter,soderlan,etzioni}@cs.washington.edu

Abstract

Can a system that “learns from reading” figure out on its own the semantic classes of arbitrary noun phrases? This is essential for text understanding, given the limited coverage of proper nouns in lexical resources such as WordNet. Previous methods that use lexical patterns to discover hypernyms suffer from limited precision and recall.

We present methods based on lexical patterns that find hypernyms of arbitrary noun phrases with high precision. This more than doubles the recall of proper noun hypernyms provided by WordNet at a modest cost to precision. We also present a novel method using a Hidden Markov Model (HMM) to extend recall further.

Introduction and Motivation

A goal of “Machine Reading” is an automatic system that extracts information from text and supports a wide range of inferencing capabilities (Etzioni, Banko, and Cafarella 2006). To enable such inferencing, an Information Extraction (IE) system must go beyond representing extracted entities as text strings and *ontologize* the text strings (Pennacchiotti and Pantel 2006; Soderland and Mandhani 2007).

This involves semantically typing the text strings, grounding the string in real world entities where possible, and mapping the string to a concept in a formal taxonomy. In many cases, particularly for proper nouns, the text string is not found in an existing taxonomy and must be added on the fly as a child of an existing concept.

This paper focuses on just one step in this ontologizing process: finding hypernyms for an arbitrary noun phrase (NP). This is a necessary first step in semantically typing the NP and mapping it to a node in a taxonomy. The problem that we address here is as follows. Given an NP e , find a set of NPs c_i such that each c_i is a hypernym of e in some sense as judged by a human.

$$\text{hypernymOf}(e) = \{c_1, c_2, \dots, c_k\}$$

Note that we define hypernymy as a relation between surface forms and not “synsets”, as is done in WordNet.

Consider the example of an NP extracted from the Web, “Borland Delphi”. This term is not found in WordNet, but lexical patterns in a large corpus suggest that Borland Delphi may have the following hypernyms:

hypernymOf(Borland Delphi)
= *development environment*
= *software development system*
= *software tool*
= *language*
= *technology*

The rest of this paper is organized as follows. We discuss previous attempts at solving this problem. The manually engineered WordNet thesaurus (Miller et al. 1990) has good coverage of common nouns, but contains relatively few entries for Proper Nouns. Several researchers have used methods based on lexical patterns to discover hypernyms (Hearst 1992; Roark and Charniak 1998; Caraballo 1999; Snow, Jurafsky, and Ng 2005), although each of these methods suffers from limited precision and recall. We then present a series of methods for finding hypernyms of arbitrary NPs: *HYPERNYMFINDER_{freq}* which is based on the frequency of Hearst pattern matches and *HYPERNYMFINDER_{svm}* that uses a Support Vector Machine (SVM) classifier to incorporate additional features. Each of these systems learns to find hypernyms from lexical patterns and corpus statistics. Since this is an inherently error-prone process, each version of *HYPERNYMFINDER* assigns a probability to each c_i it finds. Finally, we present *HYPERNYMFINDER_{hmm}*, which uses an HMM language model to extend the recall of *HYPERNYMFINDER* to cover NPs that do not match any of the lexical patterns used by the previous versions.

Previous Work on Hypernym Discovery

The WordNet thesaurus (Miller et al. 1990) is a high-precision lexical resource that has served as a baseline for hypernym discovery since the 1990’s. On a test set of 953 Noun phrases randomly selected from our Web corpus, we found that only 17% of the proper nouns and 64% of the common nouns were covered by WordNet. Information is often spotty even for NPs that are found in WordNet. For example, WordNet has an entry for “Mel Gibson” as an instance of the class “actor”, a subclass of “human”, but does not indicate that he is a “man”. We would need to use other

knowledge sources or extractions from a corpus to find that Mr. Gibson is also an instance of “movie star”, “Australian”, or “Catholic”.

Some of the earliest research in automatically discovering hypernyms from text was by Hearst (1992). She presented the manually discovered lexical patterns listed in Figure 1 that detect hypernyms. These patterns are fairly good evidence that the entity E is a subclass or member of C . For the remainder of this paper, we will refer to these patterns as the *Hearst patterns*.

E is a C
E [,] and other C
E [,] or other C
C [,] such as E
C [,] including E
C [,] especially E
such C as E

Figure 1: The “Hearst patterns”: Lexical patterns that predict that entity E belongs to class C .

The Hearst patterns suffer from limited recall – not all pairs $\langle E, C \rangle$ will be found in such patterns. Also, due to the local nature of the patterns, they inevitably make errors. A sentence with

“... urban birds in cities such as pigeons ...”

matches the pattern “C such as E” with C bound to city and E bound to pigeon, leading to city as a hypernym of pigeon.

We made use of the Hearst patterns in our IE systems, KNOWITALL (Etzioni et al. 2004a; 2004b; 2005) and TEXTRUNNER (Banko et al. 2007). We found that we needed to refine the rules somewhat. We increased recall by replacing the NP E with a list of NPs for most patterns. For example, the pattern “C [,] such as NList” would identify k NPs with hypernym C when $NList$ is a list of k simple NPs.

We also found that some patterns only have high precision when the class name is identified as a plural noun by a statistical part of speech tagger.

Correct match on plural class name: “Developing countries such as China and India have emerged ...” ⇒ $hypernymOf(\text{China}) = \text{country}$ ⇒ $hypernymOf(\text{India}) = \text{country}$
Error from singular class name: “I like to listen to country such as Garth Brooks.” ⇒ $hypernymOf(\text{Garth Brooks}) = \text{country}$

The accuracy of proposed hypernyms is correlated with extraction frequency, the number of distinct sentences that match a Hearst pattern. We estimated that, on a corpus of 117 million Web pages, hypernyms with frequency 1 had accuracy of 30% for common nouns and 49% for proper nouns. With a frequency ≥ 40 , the accuracy increased to 86% (at recall less than 0.03) for common nouns and accuracy of 87% (at recall 0.01) for proper nouns.

Snow, Jurafsky, and Ng (2005) induced a large number of weak patterns to detect hypernyms, using WordNet to automatically tag pairs of NPs from a corpus of 6 million newswire sentences. They then built a classifier using nearly 70,000 of these patterns and the frequency that each pattern occurs for a given pair of noun phrases in their corpus. Most of the patterns gave very weak evidence of a hypernym relation. The Hearst patterns, which were among those induced, were among the highest precision patterns.

Snow *et al.* found that their classifier outperformed a baseline of simply looking for a single occurrence of a Hearst pattern and also outperformed a WordNet classifier. They evaluated their classifier on a set of 5,300 hand-tagged pairs of NPs taken from random paragraphs in their newswire corpus. Their recall-precision graph indicates precision 0.85 at recall 0.10 and precision 0.25 at recall of 0.30 for their hypernym classifier. A variant of their classifier that included evidence from coordinate terms (terms with common ancestor classes) also had precision 0.85 at recall 0.10 but increased precision to 0.35 at recall 0.30.

More recently, (McNamee et al. 2008) applied the techniques described in (Snow, Jurafsky, and Ng 2005) to detecting hypernyms of Named Entities (i.e. Proper Nouns) in order to improve performance of Question Answering systems. Their corpus consisted of about 16 million sentences which was larger than that used by Snow *et al.*, but still an order of magnitude smaller than the corpus used in this work. They obtained 53% mean average precision on automatically detected hyponyms of 75 categories. Their automatically discovered hyponyms resulted in a 9% boost in performance on a TREC Question Answering data set.

Exploiting Lexico-syntactic Patterns

Our goal is an efficient, practical HYPERNYMFINDER that takes an arbitrary NP e and returns a ranked list of up to k terms that are hypernyms of e for any sense of e . We first present HYPERNYMFINDER_{freq}, a classifier based on the frequency of Hearst pattern matches.

Then we present HYPERNYMFINDER_{svm} that combines a variety of features using an SVM classifier to distinguish the correct hypernyms from among matches to the Hearst patterns. Later we explore methods that extend recall beyond hypernyms found by the Hearst patterns.

Frequency-based Classifiers

As a first step in building each of the following versions of HYPERNYMFINDER, we ran TEXTRUNNER on corpus of 117 million Web pages and recorded candidate hypernym pairs that matched any of the Hearst patterns. For each of these pairs of noun phrases $\langle e, c \rangle$ we counted the frequency of matches from distinct sentences for each Hearst pattern.

Although the precision of a candidate hypernym generally correlates with the frequency of pattern matches, we found that even very high frequency candidates are often errors. Simply setting a threshold on frequency is not sufficient to obtain high precision.

The reason for this is the local nature of the Hearst patterns, which means that they can be fooled by certain sentence structures that occur with high frequency. Here are some examples of high frequency matches on Hearst patterns that produce incorrect hypernyms.

- | | |
|----|------------------------------------------------------------------------------------------|
| A. | “... all around the world including Australia, ...”
⇒ $hypernymOf(Australia) = world$ |
| B. | “... information about this hotel such as rates, ...”
⇒ $hypernymOf(rates) = hotel$ |
| C. | “... first prize is a trip to ...”
⇒ $hypernymOf(prize) = trip$ |

Figure 2: Certain sentence structures that cause errors match Hearst patterns with high frequency. Fortunately, most of these sentence structures match only “right” patterns (as in A and B) or only “left” patterns (as in C). High precision can be obtained by filtering out hypernyms that do not have at least one right and one left pattern.

Most of these high frequency errors occur only with “right” patterns (where e comes to the right of c) or only with “left” patterns (where e comes to the left of c). Precision is improved by requiring at least one match on both right and left pattern.

Figure 3 compares the recall-precision curve for two methods of using Hearst pattern frequencies to find hypernyms. Both methods rank the proposed hypernyms by the number of distinct sentences that match a Hearst pattern. The baseline frequency method simply ranks by total pattern matches. The filtered frequency method (labeled as “1 left and 1 right”) assigns low probability to any proposed hypernyms that are not found in at least one right pattern and at least one left pattern.

For this and each of the other experiments, we used a test set of 953 well-formed simple NPs e that were randomly selected from TEXTRUNNER extractions from a corpus of 117 million Web pages. We had each method output up to $k = 5$ proposed hypernyms for each e . We hand-tagged the proposed hypernyms as correct, incorrect, or vague.

We compute recall in these experiments as the percentage of NPs in the test set that are “covered” by each method – those NPs for which at least one correct hypernym is found. Precision is the total number of correct hypernyms divided by the total number of guessed hypernyms (up to a maximum of 5 for each e). Note that this differs from the method used to compute Precision and recall used by (Snow, Jurafsky, and Ng 2005). Their method computed precision and recall over candidate hypernym/hyponym word pairs. In contrast we are more concerned with whether or not we can find a good hypernym for a given noun phrase.

The precision-recall curve for the baseline frequency method has recall 0.11 at precision 0.90 after an initial dip below 0.90 at recall 0.02. The filtered frequency method avoids this initial dip in precision and reaches recall 0.23 at precision 0.90, more than twice the recall of the baseline method. Each method reaches recall of 0.66 at precision

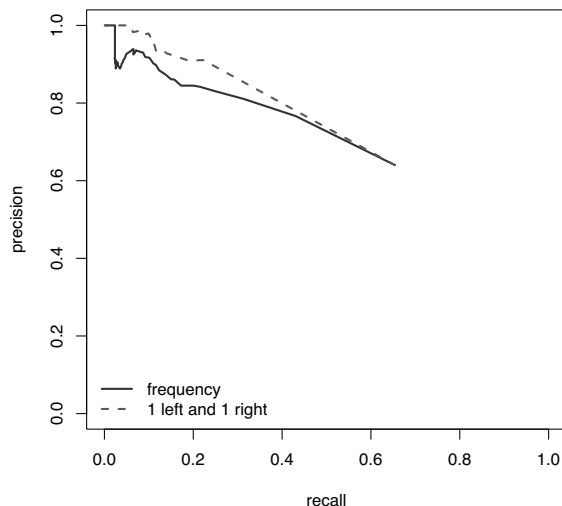


Figure 3: Ranking hypernyms by the frequency of Hearst pattern matches covers 66% of the NPs in our test set. Filtering these hypernyms to require at least one “left” and one “right” pattern gives a boost in precision.

0.64.

Note that there is a limit in recall for any method built on the Hearst patterns – about 1/3 of the noun phrases in our test set do not occur in any Hearst patterns in our large web corpus.

Combining Evidence with SVM

Although simple rules based on extraction frequency and the existence of both right and left patterns proved surprisingly effective, there are other types of evidence that can help distinguish correct hypernyms. We created $HYPERNYMFINDER_{svm}$ that uses an SVM classifier to assign a probability of correctness (Platt 1999) to a proposed hypernym pair $\langle e, c \rangle$ that has matched at least one Hearst pattern on a large text corpus.

We used the following features to classify the proposed hypernym $\langle e, c \rangle$:

- Total number of distinct pattern matches.
- Total number of *left* pattern matches.
- Total number of *right* pattern matches.
- Total number of matches for “ e is a c ”.
- Fraction of matches e which follow articles, determiners, or quantifiers.
- The rank by frequency among hypernyms of e .

While frequency is the best indicator of hypernym correctness, we wanted to combine this flexibly with other weaker evidence. Our rule-based classifier filtered out proposed hypernyms where either left or right pattern matches had zero frequency, but this is a bit extreme. Many hypernyms with unbalanced left-right patterns are actually correct.

We included the frequency of matches on “ e is a c ” because this is somewhat less reliable than the other patterns. The ratio of e which follow articles, determiners and quantifiers helps avoid matches on sentences such as “Our dog is a bloodhound ...”. This does not imply that all dogs are bloodhounds.

We trained a Support Vector Machine (SVM) classifier on hand-labeled data using 10 fold cross validation. For these experiments, we split the test referenced earlier into a set of 583 proper nouns and a set of 370 common nouns, and evaluated each set separately¹.

We consider a pair of NPs $\langle e, c \rangle$ in our test sets as correct if e belongs to the class c for some sense of e . Thus the entity “cedar” has hypernyms “wood”, “softwood”, and “material” for one sense of cedar, and also has hypernyms “tree” and “conifer” for another sense. We let each method output up to 5 hypernyms for each NP in the test set.

WordNet is a baseline for high-precision hypernyms – if a word is found in WordNet, it is guaranteed to have at least one hypernym. By our definition of correct hypernyms, WordNet has nearly perfect precision. So a practical hypernym finder would first look up a word in WordNet, before augmenting this with other methods.

In the remainder of the experiments, we begin by looking up each e in WordNet and consider the words found in WordNet to be “covered” with at least one correct hypernym. Thus, each of our recall-precision curves starts at recall 0.17 and precision 1.00 for proper nouns, and starts at recall 0.64 and precision 1.00 for common nouns. These represent the words found in WordNet. We found 43% of the proper nouns and 9% of the common nouns which were covered by WordNet to be multi-word phrases. Our interest is in extending the recall to words not in WordNet.

As in earlier experiments, we count recall as the number of e for which at least one correct hypernym c is found. Precision is computed as the fraction of c which are correct.

Figures 4 and 5 present results for proper nouns and common nouns, respectively. The lines labeled “Rule-based Classifier” are for the filtered frequency implementation of $\text{HYPERNYMFINDER}_{freq}$. The lines labeled “SVM” are from $\text{HYPERNYMFINDER}_{svm}$ using the features enumerated in this section.

A simple WordNet lookup covers 17% of proper nouns. $\text{HYPERNYMFINDER}_{freq}$ extends the recall to 0.32 at precision 0.90 and recall 0.65 at precision 0.64. The curve for $\text{HYPERNYMFINDER}_{svm}$ is consistently higher than for $\text{HYPERNYMFINDER}_{freq}$ and has recall 0.37 at precision 0.90 and a higher precision of 0.73 at the tail of the curve.

There was less room for improvement over WordNet for common nouns, of which WordNet covers 64%. Here $\text{HYPERNYMFINDER}_{freq}$ extends recall to 0.67 at precision 0.90 and gives recall of 0.80 at precision 0.61. Again, the curve for $\text{HYPERNYMFINDER}_{svm}$ is consistently higher than for $\text{HYPERNYMFINDER}_{freq}$ with recall 0.70 at precision 0.90.

These methods are straightforward ways to find hyper-

¹56% of the proper nouns and 27% of the common nouns were multi-word phrases

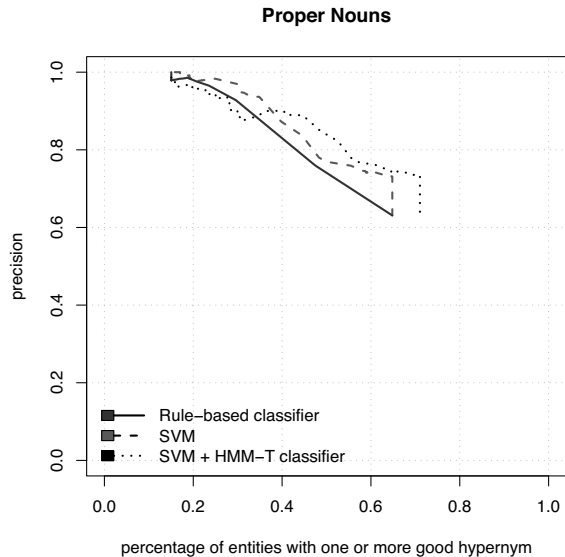


Figure 4: Each of our methods extends recall for proper nouns far beyond the recall of 0.17 from WordNet lookup. An SVM classifier uniformly raises precision over simple rules using Hearst pattern frequency. Adding an HMM language model extends recall even further.

nyms for arbitrary proper nouns that extend recall beyond that of WordNet, particularly the SVM classifier, with a graceful decline in precision at higher recall. Yet these results are limited by the coverage of Hearst patterns. For 20% of the common nouns and 36% of the proper nouns in our test set, there were no matches of Hearst patterns that give a correct hypernym. In the following section, we describe experiments to extend the coverage of our HYPERNYMFINDER to find hypernyms for some of these NPs that are beyond the reach of Hearst patterns.

Extending Recall Beyond Patterns

Here is the intuition of our approach to extending recall beyond the Hearst patterns. If we cannot find hypernyms for entity e using the methods presented thus far, perhaps we can find a sufficiently similar entity e' whose hypernyms we know. Given a good similarity metric, there should be a high probability that e and e' share the same hypernyms.

HMM Language Models

Inspired by REALM (Downey, Schoenmackers, and Etzioni 2007) we investigated using Hidden Markov Models (HMMs) (Rabiner 1989) as the similarity metric for another version of HYPERNYMFINDER . REALM uses HMMs to perform type checking of arguments in an open-domain IE system.

Given an HMM with n hidden states, which has been trained to model the generation of our corpus, we can compute the probability that each NP is generated by each of the

