

1. Ontology Learning Part One — On Discovering Taxonomic Relations from the Web

Alexander Maedche¹, Viktor Pekar², and Steffen Staab³

¹ FZI
University of Karlsruhe
Germany
<http://www.fzi.de/wim>

² Bashkir State University
Okt.Revolutsii 3a
Ufa, Russia, 450000
vpekar@ufanet.ru

³ Institute AIFB	&	Learning Lab	&	Ontoprise GmbH
University of		Lower Saxony		Karlsruhe, Germany
Karlsruhe		Hannover,		http://www.ontoprise.de
Germany		Germany		
http://www.aifb.uni-karlsruhe.de/WBS		http://www.learninglab.de		

1.1 Introduction

The Web in its' current form is an impressive success with a growing number of users and information sources. However, the growing complexity of the Web is not reflected in the current state of Web technology. The heavy burden of accessing, extracting, interpreting and maintaining information is left to the human user. Tim Berners-Lee, the inventor of the WWW, coined the vision of a Semantic Web in which background knowledge on the meaning of Web resources is stored through the use of machine-processable (meta-)data. The Semantic Web should bring structure to the content of Web pages, being an extension of the current Web, in which information is given a well-defined meaning. Thus, the Semantic Web will be able to support automated services based on these descriptions of semantics. These descriptions are seen as a key factor to finding a way out of the growing problems of traversing the expanding Web space, where most Web resources can currently only be found through syntactic matches (e.g., keyword search), providing a new level of Web Intelligence.

Ontologies have shown to be the right answer to these problems by providing a formal conceptualization of a particular domain that is shared by a group of people. Thus, in the context of the Semantic Web, ontologies describe domain theories for the explicit representation of the semantics of the data. The Semantic Web relies heavily on these formal ontologies that structure underlying data enabling comprehensive and transportable machine understanding. Though ontology engineering tools have matured over the last decade, the manual building of ontologies still remains a tedious, cumbersome task which can easily result in a knowledge acquisition bottleneck. The success of the Semantic Web strongly depends on the

proliferation of ontologies, which requires that the engineering of ontologies be completed quickly and easily. When using ontologies as a basis for Semantic Web applications, one has to face exactly this issue and in particular questions about development time, difficulty, confidence and the maintenance of ontologies. Thus, what one ends up with is similar to what knowledge engineers have dealt with over the last two decades when elaborating methodologies for knowledge acquisition or workbenches for defining knowledge bases [1.23, 1.24]. A method which has proven to be extremely beneficial for the knowledge acquisition task is the integration of knowledge acquisition with machine learning techniques

In this section we focus on an essential part of ontology engineering, namely the development of the taxonomic backbone of the ontology. The purpose of the section is to give a survey of existing work on learning taxonomic relations from texts and an example of how such learning may be performed and evaluated.

The article is organized as following. We start with two survey on existing work (also cf. [1.16]). We consider symbolic (Section 1.2) and statistics-based approaches (Section 1.3). The trade-off between the two is that statistics-based approaches allow for better scaling, but symbolic approaches might eventually turn up being more precise. Therefore, we aim at a reconciliation between the two paradigms and propose new algorithms for taxonomy learning including existing taxonomic relations as background knowledge (Section 1.4). The Sections 1.5 through 1.7 will focus on the evaluation of algorithms for ontology learning and elucidate the evaluation with a typical Web scenario. In particular, Section 1.5 introduces the overall setting that we used as an example for evaluating the learning of taxonomies, presenting the input data for our learning method. Section 1.6 shows how to perform evaluation and finally Section 1.7 outlays the results we have obtained so far.

1.2 Survey of Symbolic Approaches

1.2.1 Extraction of Taxonomic Relations

The idea of using lexico-syntactic patterns in the form of regular expressions for the extraction of semantic relations, in particular taxonomic relations has been introduced by [1.6]. Pattern-based approaches in general are heuristic methods using regular expressions that originally have been successfully applied in the area of information extraction (see [1.8]). In this lexico-syntactic ontology learning approach the text is scanned for instances of distinguished lexico-syntactic patterns that indicate a relation of interest, e.g. the taxonomic relation. Thus, the underlying idea is very simple: Define a regular expression that captures re-occurring expressions and map the results of the matching expression to a semantic structure, such as taxonomic relations between concepts.

Example. This examples provides a sample pattern-based ontology extraction scenario. For example, in [1.6] the following lexico-syntactic pattern is considered

... $NP\{, NP\} * \{, \}$ or other NP ...

When we apply this pattern to a sentence it can be inferred that the NP's referring to concepts on the left of *or other* are sub concepts of the NP referring to a concept on the right. For example from the sentence

Bruises, wounds, broken bones or other injuries are common.

we extract the taxonomic relations (BRUISE,INJURY), (WOUND,INJURY), and, (BROKEN-BONE,INJURY).

In [1.6] the patterns have been defined manually, which is a time-consuming and error-prone task. In [1.18] the work proposed by [1.6] is extended by using a symbolic machine learning tool to refine lexico-syntactic patterns. In this context the PROMETHEE system has been presented that supports the semi-automatic acquisition of semantic relations and the refinement of lexico-syntactic patterns.

The work of [1.1] reports a practical experiment of construction of a regional ontology in the field of electric network planning. He describes a clustering approach that combines linguistic and conceptual criteria. As an example he gives the pattern $\langle NP, line \rangle$ which results in two categorizations by modifiers. The first categorization is motivated by the `function_of_structure` modifiers, resulting in a clustering of `connection line`, `dispatching line` and `transport line` (see Table 1.1). For the other concepts the background knowledge lacks adequate specifications such that further categorizations could have been proposed.

Table 1.1. Example Categorization

A proposal categorization	The other candidate terms
connection line	mountain line
dispatching line	telecommunication line
transport line	input line

[1.3] have presented a cooperative machine learning system called ASIUM which is able to acquire taxonomic relations from syntactic parsing. The ASIUM system is based on a conceptual clustering algorithm. Basic clusters are formed on head words that occur with the same verb after the same preposition. ASIUM successively aggregates clusters to form new concepts and the hierarchies of concepts form the ontology. The ASIUM approach differs from the approach in this work because the relation learning is restricted to taxonomic relations.

An ontology learning system where the different techniques have been applied on dictionary definitions in the context of the insurance and telecommunication domains is described in [1.11, 1.15]. An important aspect in this system and approach is that existing concepts are included in the overall process. Thus, in contrast to [1.6, 1.18] the extraction operations have been performed on the concept level, thus, patterns have been directly matched onto concepts. Thus, the system is, beside ex-

tracting taxonomic relations from scratch, able to refine existing relations and refer to existing concepts.

1.2.2 Refinement of Taxonomic Relations

[1.4] introduced a methodology for the maintenance and refinement of domain-specific taxonomies. An ontology is incrementally updated as new concepts are acquired from real-world texts. The acquisition process is centered around linguistic and conceptual “quality” of various forms of evidence underlying the generation and refinement of concept hypotheses. In particular they consider semantic conflicts and analogous semantic structures from the knowledge base into the ontology in order to determine the quality of a particular proposal. Thus, they extend an existing ontology with concepts and taxonomic relations between concepts.

The system Camille¹ was developed as a natural language understanding system, e.g. when the parser comes across words that it does not know, Camille tries to infer whatever it can about the meaning of the unknown word [1.5]. If the unknown word is a noun, semantic constraints on slot-fillers provided by verbs give useful limitations about what the noun could mean. The meaning of a noun can be derived, because constraints are associated with verbs. Learning unknown verbs is more difficult, thus, verb acquisition has been the main focus of the research on Camille. Camille was tested on several real-world domains within information extraction tasks (MUC), where the well-known scoring methods precision and recall, taken from the information retrieval community, have been calculated. For the lexical acquisition task *recall* is defined as the percentage of correct hypobooks. A hypobook was counted as correct if one of the concepts in the hypobook matched the target concept. Precision is the total number of correct concepts divided by the number of concepts generated in all the hypobooks. Camille has achieved a recall of 42% and a precision of 19% on a set of 50 randomly-selected sentences containing 17 different verbs.

1.3 Survey of Statistics-based Approaches

1.3.1 Statistics-based Extraction of Taxonomic Relations

Clustering can be defined as the process of organizing objects into groups whose members are similar in some way (see [1.10]). In general there are two major styles of clustering: **non-hierarchical clustering** in which every object is assigned to exactly one group and **hierarchical clustering**, in which each group of size greater than one is in turn composed of smaller groups. Hierarchical clustering algorithms are preferable for detailed data analysis. They produce hierarchies of clusters, and therefore contain more information than non-hierarchical algorithms. However, they

¹ Contextual Acquisition Mechanism for Incremental Lexeme Learning

are less efficient with respect to time and space than non-hierarchical clustering². [1.17] identify two main uses for clustering in natural language processing³: The first is the use of clustering for exploratory data analysis, the second is for generalization. Seminal work in this area of so-called *distributional clustering* of English words has been described in [1.19]. Their work focuses on constructing class-based word co-occurrence models with substantial predictive power. In the following the existing and seminal work of applying statistical hierarchical clustering in NLP (see [1.19]) is adopted and embedded into the framework.

Baseline Hierarchical Clustering. The tree of hierarchical clusters can be produced either bottom-up, by starting with individual objects and grouping the most similar ones, or top-down, whereby one starts with all the objects and divides them into groups. Algorithm 1 (adopted from [1.17]) describes the bottom-up algorithm. It starts with a separate cluster for each object. In each step, the two most similar clusters are determined, and merged into a new cluster. The algorithm terminates when one large cluster containing all objects has been formed. The most important aspect in clustering is the selection of an appropriate computation strategy and a similarity measure. We will introduce a number of computation strategies (e.g. single-link, complete link or group-average) and similarity measures (e.g. cosine, kullback-leibler) later in this subsection.

Algorithm 1 Hierarchical Clustering Algorithm — Bottom-Up

Require: a set $X = \{x_1, \dots, x_n\}$ of objects, n as the overall number of objects,
a function $\text{sim}: 2^X \times 2^X \rightarrow R$

Ensure: the set of clusters K (or cluster hypothesis)

```

for  $i:=1$  to  $n$  do
     $k_i := x_i$ .
end for
 $K := \{k_1, \dots, k_n\}$ .
 $j := n + 1$ .
while  $|K| > 1$  do
     $(k_{n1}, k_{n2}) := \arg \max_{(k_u, k_v) \in K \times K} \text{sim}(k_u, k_v)$ .
     $k_j = k_{n1} \cup k_{n2}$ .
     $K := K \setminus \{k_{n1}, k_{n2}\} \cup \{k_j\}$ .
     $j := j + 1$ 
end while

```

Algorithm 2 (adopted from [1.17]) roughly describes the top-down algorithm. It starts out with one cluster that contains all objects. The algorithm then selects the least coherent cluster in each iteration and splits it. Clusters with similar objects are more coherent than clusters with dissimilar objects. Thus, the strategies single-link,

² Hierarchical clustering has in the average quadratic time and space complexity.

³ A comprehensive survey on applying clustering in NLP is also available in the EAGLES report, see <http://www.ilc.pi.cnr.it/EAGLES96/rep2/node37.htm>

complete link and group-average can also serve as measures of cluster coherence (function coh) in top-down clustering.

The reader may note that splitting a cluster (function split) is also a clustering task (namely the task of finding two sub-clusters of a cluster). Thus, there is a recursive need for a second clustering algorithm. Any clustering algorithm may be used for the splitting operation, including bottom-up algorithms.

Algorithm 2 Hierarchical Clustering Algorithm — Top-Down

Require: a set $X = \{x_1, \dots, x_n\}$ of objects, n as the overall number of objects,

a function coh: $2^X \rightarrow \mathbf{R}$

a function split: $2^X \times 2^X \rightarrow 2^X$

$K := \{X\} (= k_1)$

$j := 1$

while $\exists k_i \in K$ s.t. $|k_i| > 1$ **do**

$k_u := \arg \min_{k_v \in K} \text{coh}(k_v)$

$(k_{j+1}, k_{j+2}) = \text{split}(k_u)$

$K := K \setminus \{k_u\} \cup \{k_{j+1}, k_{j+2}\}$

$j := j + 2.$

end while

As mentioned earlier an important aspect is the selection of an appropriate computation strategy and a similarity measure. In the following the most important ones are presented.

Computation strategies used in hierarchical clustering. In this work it is focused on the three functions single link, complete link and group-average that have shown to perform good in statistical hierarchical clustering. Their advantages and disadvantages a shortly introduced. The interested reader is referred to a more detailed introduction given in [1.10]. Measuring similarity based on **single linkage** means that the similarity between two clusters is the similarity of the two closest objects in the clusters. Thus, one has to search over all pairs of objects that are from the two different clusters and select the pair with the greatest similarity. Single-link clustering have clusters with local coherence. If similarity is based on **complete linkage** the similarity between two clusters is computed based on the similarity of the two least similar members. Thus, the similarity of two clusters is the similarity of their two most dissimilar members. Complete-link clustering has a similarity function that focuses on global cluster quality. The last similarity function considered is **group-average**. Group average may be considered as a bit of both, single linkage and complete linkage. The criterion for merges is the average similarity between members.

Similarity Measures. As mentioned earlier clustering requires some kind of similarity measure that is computed between objects using the functions described above. Different similarity measures (e.g a good overview is given in [1.12]) and their evaluation [1.2] are available from the statistical natural language processing community. The two most important measures within our work, namely the cosine measure

(see Definition 1.3.1) and the kullback leibler divergence (see Definition 1.3.2) are briefly introduced. The cosine measure and the kullback leibler divergence proved to be the most important ones in the area of statistical NLP.

Definition 1.3.1. *The cosine measure or normalized correlation coefficient between two vectors x and y is given by*

$$\cos(x, y) = \frac{\sum_{x \in X, y \in Y} xy}{\sqrt{\sum_{x \in X} x^2 \sum_{y \in Y} y^2}} \quad (1.1)$$

Using the cosine measure it is computed how well the occurrence of a specific lexical entry correlates in x and y and then divided by the Euclidean length of the two vectors to scale for the magnitude of the individual length of x and y .

Though, the following measure is not a metric in the strong sense, it has been quite successfully applied in statistical NLP. The kullback leibler divergence has its roots in information theory and is defined as follows:

Definition 1.3.2. *For two probability mass functions $p(x)$, $q(x)$ their relative entropy is computed by*

$$D(p||q) = \sum_{x \in X} p(x) \log \frac{p(x)}{q(x)} \quad (1.2)$$

The kullback leibler divergence is a measure of how different two probability distributions (over the same event space) are. The kullback leibler divergence between p and q is the average number of bits that are wasted by encoding events from a distribution p with a code based on a not-quite-right distribution q . The quantity is always non-negative, and $D(p||q) = 0$ iff $p = q$. An important aspect is that kullback leibler divergence is not defined for $p(X) > 0$ and $q(x) = 0$. In cases where propability distributions of objects have many zeros, the usage of bottom-up clusering becomes nearly impossible. Thus, for using kullback leibler divergence top-down clustering is the more natural choice.

Example. To explain the similarity measures a small example is given in the following. Imagine a simple concept-concept matrix as given by Table 1.2 consisting of 5 concepts.

Using the cosine measure one may compute the similarity between the concepts HOTEL and ACCOMODATION as follows. The vector of the concept HOTEL is given by $x^T = (0, 14, 7, 4, 6)$, the vector of the concept ACCOMODATION is given by $y^T = (14, 0, 11, 2, 5)$.

$$\cos(x, y) = \frac{7 \cdot 11 + 4 \cdot 2 + 6 \cdot 5}{101 \cdot 150} \approx 0.93 \quad (1.3)$$

For computing the kullback leibler divergence one has first calculate the probability mass functions for each concept and its corresponding frequencies. The probability mass functions for HOTEL are given as $(0, 0.45, 0.22, 0.13, 0.19)$ the probability mass functions for the concept ACCOMODATION are given as $(0.44, 0, 0.34, 0.06, 0.16)$

Table 1.2. Example Similarity Matrix

ID	HOTEL	ACCOMODATION	ADDRESS	WEEKEND	TENNIS
HOTEL	-	14	7	4	6
ACCOMODATION	14	-	11	2	5
ADDRESS	7	11	-	10	3
WEEKEND	4	2	10	-	5
TENNIS	6	5	3	5	-

Based on these values one can compute the kullback leibler divergence as follows

$$D(\text{HOTEL}||\text{ACCOMODATION}) = 0.22 \cdot \frac{0.22}{0.34} + \dots + 0.19 \cdot \frac{0.19}{0.16} \approx 0.65(1.4)$$

We refer the reader to [1.17] where a detailed introduction into further similarity measures between two sets X and Y such as the matching coefficient $X \cap Y$, the dice coefficient $\frac{2|X \cap Y|}{|X| + |Y|}$, the Jaccard or Tanimoto coefficient $\frac{|X \cap Y|}{|X \cup Y|}$ or the overlap coefficient $\frac{|X \cap Y|}{\min(|X|, |Y|)}$ is given.

In the area of information retrieval some work of automatically deriving a hierarchical organization of concepts from a set of documents without use of training data or standard clustering techniques has been presented by [1.21]. They use a subsumption criterion to organize the salient words and phrases extracted from documents hierarchically. In [1.9] a novel statistical latent class model is used for text mining and interactive information access. In his work the author introduces a Cluster-Abstraction Model (CAM) that is purely data-driven and utilizes context-specific word occurrence statistics. CAM extracts hierarchical relations between groups of documents as well as an abstract organization of keywords.

1.3.2 Statistics-based Refinement of Taxonomic Relations

Classification techniques previously applied to distributional data can be summarized according to the following methods: the k nearest neighbor (kNN) method, the category-based method and the centroid-based method. They all operate on vector-based semantic representations, which describe the meaning of a word of interest (target word) in terms of counts of its cooccurrence with context words, i.e., words appearing within some delineation around the target word. The key differences between the methods stem from different underlying ideas about how a semantic class of words is represented, i.e. how it is derived from the original cooccurrence counts, and, correspondingly, what defines membership in a class. The kNN method is based on the assumption that membership in a class is defined by the new instance's similarity to one or more individual members of the class. Thereby, similarity is defined by a similarity score as, for instance, by the cosine between cooccurrence vectors. To classify a new instance, one determines the set of k training instances that are most

similar to the new instance. The new instance is assigned to the class that has the biggest number of its members in the set of nearest neighbors. In addition, the classification decision can be based on the similarity measure between the new instance and its neighbors: each neighbor may vote for its class with a weight proportional to its closeness to the new instance. When the method is applied to augment a thesaurus, a class of training in-stances is typically taken to be constituted by words belonging to the same synonym set, i.e. lexicalizing the same concept (e.g., [1.7]). A new word is assigned to that synonym set that has the biggest number of its members among nearest neighbors. The major disadvantage of the kNN method that is often pointed out is that it involves significant computational expenses to calculate similarity between the new instance and every instance of the training set. A less expensive alternative is the category-based method (e.g., [1.20]). Here the assumption is that membership in a class is defined by the closeness of the new item to a generalized representation of the class. The generalized representation is built by adding up all the vectors constituting a class and normalising the resulting vector to unit length, thus computing a probabilistic vector representing the class. To determine the class of a new word, its unit vector is compared to each vector representing a class. Thus the number of calculations is reduced to the number of classes. Thereby, a class representation may be derived from a set of vectors corresponding to a synonym set or a set of vectors corresponding to a synonym set and some or all subordinate synonym sets. In the straightforward case of standard kNN a class is represented by its synonym set. The other extreme is to represent the class in the same manner as is used, e.g., in [1.20] to represent a concept of a thesaurus. Here a class vector would be built from data on all hyponyms of the corresponding concept; new words found similar to this vector would be assigned to that class that lexicalizes the corresponding concept. Another way to prepare a representation of a word class is what may be called the centroid-based approach (e.g., [1.19]). It is almost exactly like the category-based method, the only difference being that the vector representing a class is computed slightly differently. All n vectors corresponding to class members are added up and the resulting vector is divided by n to compute the centroid between the n vectors.

1.4 Making use of the structure of the thesaurus

The classification methods described above presuppose that semantic classes being augmented exist independently of each other. For most existing thesauri this is not the case: they typically encode taxonomic relations between word classes. It seems worthwhile to employ this information to enhance the performance of the classifiers.

1.4.1 Tree descending algorithm

One way to factor the taxonomic information into the classification decision is to employ the tree-descending classification algorithm, which is a familiar technique

in text categorization. The principle behind this approach is that the semantics of every concept in the thesaurus tree retains some of the semantics of all its hyponyms in such a way that the upper the concept, the more relevant semantic characteristics of its hyponyms it reflects. It is thus feasible to determine the class of a new word by descending the tree from the root down to a leaf. The semantics of concepts in the thesaurus tree can be represented by means of one of the three methods to represent a class described in Section 1.3. At every tree node, the decision which path to follow is made by choosing the child concept that has the biggest distributional similarity to the new word. After the search has reached a leaf, the new word is assigned to that synonym set, which lexicalizes the concept that is most similar to the new word. This manner of search offers two advantages. First, it allows to gradually narrow down the search space and thus save on computational expenses. Second, it ensures that, in a classification decision, more relevant semantic distinctions of potential classes are given more preference than less relevant ones. As in the case with the category-based and the centroid-based representations, the performance of the method may be greatly dependent on the number of subordinate synonyms sets included to represent a concept.

1.4.2 Tree ascending algorithm

Another way to use information about inter-class relations contained in a thesaurus is to base the classification decision on the combined measures of distributional similarity and taxonomic similarity (i.e., semantic similarity induced from the hierarchical organization of the thesaurus) between nearest neighbors. Suppose words in the nearest neighbors set for a given new word, e.g., *trailer*, all belong to different classes as in the following classification scenario: *box* (similarity score to *trailer*: 0.8), *house* (0.7), *barn* (0.6), *villa* (0.5) (Figure 1.1). In this case, the kNN method will classify *trailer* into the class CONTAINER, since it appears to have biggest similarity to *box*. However, it is obvious that the most likely class of *trailer* is in a different part of the thesaurus: in the nearest neighbors set there are three words which, though not belonging to one class, are semantically close to each other. It would thus be safer to assign the new word to a concept that subsumes one or all of the three semantically similar neighbors. For example, the concepts DWELLING or BUILDING could be feasible candidates in this situation.

The crucial question here is how to calculate the voting weight for these two concepts to be able to decide which of them to choose or whether to prefer the concept of *box*. Clearly, one cannot sum or average the distributional similarity measures of neighbors below a candidate concept. In the first case the root will always be the best-scoring concept. In the second case the score of the candidate concept will always be smaller than the score of its biggest-scoring hyponym. We propose to estimate the voting weight for such candidate concepts based on taxonomic similarity between relevant nodes. The taxonomic similarity between two concepts is measured according to the procedure elaborated in [1.14]. Assuming that a taxonomy is given as a tree with a set of nodes N , a set of edges $E \subset N \times N$, a unique

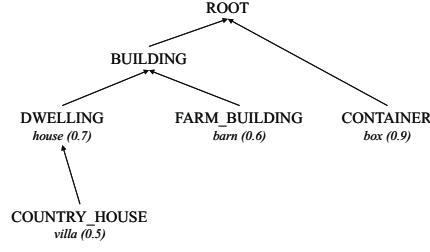


Fig. 1.1. A semantic classification scenario

root $ROOT \in N$, one first determines the least common superconcept of a pair of concepts a, b being compared.

It is defined by

$$lcs(a, b) := \iota c \in N : \delta(a, c) + \delta(b, c) + \delta(root, c) \text{ is minimal} \quad (1.5)$$

where $\delta(a, b)$ describes the number of edges on the shortest path between a and b . The taxonomic similarity between a and b is then given by

$$\Omega(a, b) := \frac{\delta(ROOT, c)}{\delta(a, c) + \delta(b, c) + \delta(ROOT, c)} \quad (1.6)$$

where $c = lcs(a, b)$. Ω is such that $0 \leq \Omega \leq 1$, with 1 standing for the maximum taxonomic similarity. Ω is directly proportional to the number of edges from the least common superconcept to the root, which agrees with the intuition that a given number of edges between two concrete concepts signifies greater similarity than the same number of edges between two abstract concepts.

The first method to calculate the voting weight for a candidate concept is to sum the distributional similarity measures of its hyponyms to the target word t each weighted by the taxonomic similarity measure between the hyponym and the candidate node:

$$W(n) := \sum_{h \in I_n} \text{sim}(t, h) \cdot \Omega(n, h) \quad (1.7)$$

where I_n is the set of hyponyms below the candidate concept n , $\text{sim}(t, h)$ is the distributional similarity between a hyponym h and the word to be classified t , and $\Omega(n, h)$ is the taxonomic similarity between the candidate concept and the hyponym h .

1.5 Data and Settings of the Experiments

The machine-readable thesaurus we used in this study was derived from GET-ESS [1.22], an ontology for the tourism domain. Each concept in the ontology is associated with one lexical item, which expresses this concept. From this ontology, word classes were derived in the following manner. A class was formed by

words lexicalizing all child concepts of a given concept. For example, the concept CULTURAL_EVENT in the ontology has successor concepts PERFORMANCE, OPERA, FESTIVAL, associated with words *performance*, *opera*, *festival* correspondingly. Though these words are not synonyms in the traditional sense, they are taken to constitute one semantic class, since out of all words of the ontology's lexicon their meanings are closest. The thesaurus thus derived contained 1052 words and phrases (out of these, 756 cropped up in the corpus at least once) grouped into 182 classes. The corpus from which distributional data about the words were obtained was extracted from a web site advertising hotels around the world. It contained about 6 megabytes of text (988 000 words).

Collection of distributional data was carried out in the following settings. The preprocessing of corpus included a very simple stemming (most common inflections were chopped off; irregular forms of verbs, adjectives and nouns were changed to their first forms). The context of usage was delineated by a window of 3 words on either side of the target word, without transgressing sentence boundaries. In case a stop word other than a proper noun appeared inside the window, the window was accordingly expanded. The stoplist included 50 most frequent words of the British National Corpus, words listed as function words in the BNC, and proper nouns not appearing in the sentence-initial position. The obtained frequencies of cooccurrence were weighted by the $1+\log$ weight function. The distributional similarity was measured by means of three different similarity metrics: the Jaccard's coefficient, L1 distance, and the skew divergence, a weighted version of the Kullback-Leibler divergence (cf., [1.13]).

1.6 Evaluation method

The performance of the algorithms was assessed in the following manner. For each algorithm, we held out a single word of the thesaurus as the test case, and trained the system on the remaining 755 words. We then tested the algorithm on the held-out vector, observing if the assigned class for that word coincided with its original class in the thesaurus, and counting the number of correct classifications (direct hits). This was repeated for each of the words of the thesaurus.

However, given the intuition that a semantic classification may not be simply either right or wrong, but rather of varying degrees of appropriateness, we believe that a clearer idea about the quality of the classifiers would be given by an evaluation method that takes into account near misses as well. We therefore evaluated the performance of the algorithms also in terms of how close the proposed class for a test word was to the correct class. For this purpose we measured the taxonomic similarity between the assigned and the correct classes of words so that the appropriateness of a particular classification was estimated on a scale between 0 and 1, with 1 signifying assignment to the correct class. Thus this measure of accuracy of classifications was compatible with the counting of direct hits, which, as will be shown later, may be useful for evaluating the methods. In the following, the evaluation of the classification algorithms is reported both in terms of the average of direct

hits and the average of the taxonomy similarity between the assigned and the correct classes (direct+near hits) over all words in the thesaurus.

To have a benchmark for evaluation of the algorithms, a baseline was calculated, which was the average hit value a given word gets, when its class label is chosen at random. The baseline for direct hits was estimated at 0.012; for direct+near hits, it was 0.15741.

1.7 Results

We first conducted experiments evaluating performance of the three standard classifiers. To determine the best version for each particular classifier, only those parameters were varied that, as described above, we deemed to be critical for a specific algorithm in the setting of thesaurus augmentation. Other parameters can have a serious impact on the absolute quality of their performance, but we, however, were interested in comparing the classifiers relative to each other.

In order to get a view on how the accuracy of the algorithms was related to the amount of available distributional data on the target word, all words of the thesaurus were divided into three groups depending on the amount corpus data available on them (see Table 1.3). The amount of distributional data for a word (the frequency in the left column) is the total of frequencies of its context words.

Table 1.3. Distribution of words of the thesaurus into frequency ranges

Frequency range	# words in the range
0-40	274
40-500	190
>500	292

The results of the evaluation of the methods are summarized in the tables below. Rows specify the metrics used to measure distributional similarity and columns specify frequency ranges. Each cell describes the average of direct+near hits / the average of direct hits over words of a particular frequency range and over all words of the thesaurus. The statistical significance of the results was measured in terms of the one-tailed chi-square test.

kNN. The evaluation of the method was conducted with $k = 1, 3, 5, 7, 10, 15, 20, 25$, and 30. The accuracy of classifications increased with the increase of k . However, starting with $k = 20$ the increase of k yielded only insignificant improvement. Table 1.4 describes results of evaluation of kNN using 30 nearest neighbors, which was found to be the best version of kNN.

Category-based method. To determine the best version of this method, we experimented with the number of levels of hyponyms below a concept that were used to build a class vector). The best results were achieved when a class was represented by data from its hyponyms at most three levels below it (Table 1.5).

Table 1.4. kNN, $k=30$

	0-40	40-500	>500	Overall
JC	.33773 /.17142	.33924 /.15384	.40181 /.12457	.37044 /.15211
L1	.33503 /.16428	.38424 /.21025	.38987 /.14471	.37636 /.17195
SD	.31505 /.14285	.36316 /.18461	.45234 /.17845	.38806 /.17063

Table 1.5. Category-based method, 3 levels

	0-40	40-500	>500	Overall
JC	.26918 /.12142	.34743 /.17948	.47404 /.28282	.37554 /.2023
L1	.27533 /.125	.41736 /.25128	.56711 /.38383	.43242 /.26190
SD	.28589 /.12857	.34932 /.18461	.51306 /.31649	.39755 /.21957

Centroid-based method.. As in the case with the category-based method, we varied the number of levels of hyponyms below the candidate concept. Table 1.6 details results of evaluation of the best version of this method (a class is represented by 3 levels of its hyponyms).

Table 1.6. Centroid-based method, 3 levels

	0-40	40-500	>500	Overall
JC	.17362 /.07831	.18063 /.08119	.30246 /.14434	.22973 /.10714
L1	.21711 /.09793	.30955 /.13938	.37411 /.1687	.30723 /.12698
SD	.22108 /.09972	.23814 /.11374	.36486 /.16147	.28665 /.10714

Comparing the three algorithms we see that overall, kNN and the category-based method exhibit comparable performance (with the exception of measuring similarity by L1 distance, when the category-based method outperforms kNN by a margin of about 5 points; statistical significance $p < 0.001$). However, their performance is different in different frequency ranges: for lower frequencies kNN is more accurate (e.g., for L1 distance, $p < 0.001$). For higher frequencies, the category-based method improves on kNN (L1, $p < 0.001$). The centroid-based method exhibited performance, inferior to both those of kNN and the category-based method.

Tree descending algorithm.. In experiments with the algorithm, candidate classes were represented in terms of the category-based method, 3 levels of hyponyms, which proved to be the best generalized representation of a class in previous experiments. Table 1.7 specifies the results of its evaluation.

Its performance turns out to be much worse than that of the standard methods. Both direct+near and direct hits scores are surprisingly low, for 0-40 and 40-500 much lower than chance. This can be explained by the fact that some of top concepts in the tree are represented by much less distributional data than other ones. For example, there are less than 10 words that lexicalize the top concepts MASS_CONCEPT and MATHEMATICAL_CONCEPT and all of their

Table 1.7. Tree descending algorithm

	0-40	40-500	>500	Overall
JC	.00726 /0	.01213 /.00512	.02312 /.0101	.014904 /.005291
L1	.08221 /.03214	.05697 /.02051	.21305 /.11111	.128844 /.060846
SD	.08712 /.03214	.07739 /.03589	.16731 /.06734	.011796 /.047619

hyponyms (compare to more than 150 words lexicalizing THING and its hyponyms up to 3 levels below it). As a result, at the very beginning of the search down the tree, a very large portion of test words was found to be similar to such concepts.

Tree ascending algorithm. The experiments were conducted with the same number of nearest neighbors as with kNN. Table 1.8 describes the results of evaluation of the best version (equation 1.7, $k = 15$).

Table 1.8. Tree ascending algorithm, voting weight according to equation 1.7, $k=15$.

	0-40	40-500	>500	Overall
JC	.32112 /.075	.33553 /.0923	.40968 /.08754	.36643 /.08597
L1	.33369 /.07142	.34504 /.0923	.42627 /.09764	.38005 /.08862
SD	.31809 /.06785	.32489 /.05128	.45529 /.11111	.38048 /.08201

There is no statistically significant improvement on kNN overall, or in any of the frequency ranges. The algorithm favored more upper concepts and thus produced about twice as few direct hits than kNN. At the same time, its direct+near hits score was on par with that of kNN! This algorithm thus produced much more near hits than kNN, what can be interpreted as its better ability to choose a superconcept of the correct class. Based on this observation, we combined the best version of the tree ascending algorithm with kNN in one algorithm in the following manner. First the former was used to determine a superconcept of the class for the new word and thus to narrow down the search space. Then the kNN method was applied to pick a likely class from the hyponyms of the concept determined by the tree ascending method. Table 1.9 specifies the results of evaluation of the proposed algorithm.

Table 1.9. Tree ascending algorithm combined with kNN, $k=30$

	0-40	40-500	>500	Overall
JC	.34444 /.16428	.35858 /.14358	.41260 /.10774	.38215 /.14021
L1	.35147 /.16428	.36545 /.15384	.41086 /.11784	.38584 /.14682
SD	.32613 /.13571	.36485 /.1641	.45732 /.16498	.39456 /.1574

The combined algorithm demonstrated improvement both on kNN and the tree ascending method of 1 to 3 points in every frequency range and overall for direct+near hits (except for the 40-500 range, L1). The improvement was statistically

significant only for L1, > 500 ($p = 0.05$) and for L1, overall ($p = 0.011$). For other similarity measures and frequency ranges it was insignificant (e.g., for JC, overall, $p = 0.374$; for SD, overall, $p = 0.441$). The algorithm did not improve on kNN in terms of direct hits. The hits scores set in bold in Table 1.9 are those which are higher than those for kNN in corresponding frequency ranges and similarity measures.

1.8 Conclusion

In this section we have surveyed symbolic and statistical means for extracting taxonomic relations from text. We have proposed new algorithms that combine the advantages of scalable statistical approaches with symbolic approaches that consider background knowledge.

For this proposal we have shown how to evaluate learning approaches in a typical web setting. Though these first step have not yet resulted in a significant improvement we may conjecture that with improved and further evaluation of combined symbolic/statistic approaches we may arrive at a good quality classification of unknown words. In particular, the aggregation of multiple such approaches will yield a basis for multi-strategy learning as well as for estimating the reliability of learned taxonomies — which is necessary in order to embed these algorithms into actual ontology development environments like [1.24, ?].

In particular, our study demonstrated that taxonomic similarity between nearest neighbors, in addition to their distributional similarity to the new word, may be a useful evidence on which classification decision can be based. We have proposed a “tree ascending” classification algorithm which extends the kNN method by making use of the taxonomic similarity between nearest neighbors. This algorithm was found to have a very good ability to choose a superconcept of the correct class for a new word. On the basis of this finding, another algorithm was developed that combines the tree ascending algorithm and kNN in order to optimize the search for the correct class. Although only limited statistical significance of its improvement on kNN was found, the results of the study indicate that this algorithm is a promising possibility to incorporate the structure of a thesaurus into the decision as to the class of the new word. In particular, we conjecture that the tree ascending algorithm leaves a lot of room for improvements and combinations with other algorithms like kNN. The tree descending algorithm, a technique widely used for text categorization, proved to be much less efficient than standard classifiers when applied to the task of augmenting a domain-specific thesaurus. Its poor performance is due to the fact that in such a thesaurus there are great differences between top concepts in the amount of distributional data used to represent them.

In order to have a better understanding of the role of different parameters on the performance of the classifiers, they can be further studied on the material of a general thesaurus, where richer information about its structural organization is available. Also study on further domains like organizational memories or texts on genomics let us expect fruitful results with regard to methods and applications.

10

- 1.1 H Assadi. Construction of a regional ontology from text and its use within a documentary system. In *N. Guarino (ed.), Formal Ontology in Information Systems, Proceedings of FOIS-98, Trento, Italy, 1999*, pages 236–249, 1999.
- 1.2 I. Dagan, L. Lee, and F. Pereira. Similarity-based models of word cooccurrence probabilities. *Machine Learning*, 34(1):43–69, 1999.
- 1.3 D. Faure and C. Nedellec. A corpus-based conceptual clustering method for verb frames and ontology acquisition. In *In LREC workshop on Adapting lexical and corpus resources to sublanguages and applications, Granada, Spain, May 1998.*, 1998.
- 1.4 U. Hahn and Schnattinger. Ontology engineering via text understanding. In *IFIP'98 — Proceedings of the 15th World Computer Congress*, Vienna and Budapest, 1998.
- 1.5 Peter M. Hastings. *Automatic Acquisition of Word Meaning from Context*. PhD thesis, University of Michigan, 1994.
- 1.6 M.A. Hearst. Automatic acquisition of hyponyms from large text corpora. In *Proceedings of the 14th International Conference on Computational Linguistics*. Nantes, France, 1992.
- 1.7 Marti A. Hearst and Hinrich Schütze. Customizing a lexicon to better suit a computational task. In *Proc. of the ACL-SIGLEX Workshop on Acquisition of Lexical Knowledge from Text, Columbus, Ohio, USA, 1993*.
- 1.8 J. Hobbs. The generic information extraction system. In *Proceedings of the Fifth Message Understanding Conference (MUC-5)*, Morgan Kaufmann, 1993., 1993.
- 1.9 T. Hofmann. The Cluster-Abstraction Model: Unsupervised Learning of Topic Hierarchies from Text Data. In *Proceedings of 16th International Conference on Artificial Intelligence (IJCAI-99)*, Stockholm, Sweden, 1999, pages 682–587, 1999.
- 1.10 L. Kaufman and P.J. Rousseeuw. *Finding Groups in Data: An Introduction to Cluster Analysis*. John Wiley, 1990.
- 1.11 J.-U. Kietz, R. Volz, and A. Maedche. Semi-automatic ontology acquisition from a corporate intranet. In *International Conference on Grammar Inference (ICGI-2000)*, to appear: Lecture Notes in Artificial Intelligence, LNAI, 2000.
- 1.12 L. Lee. Measures of distributional similarity. In *Proceedings of the ACL'99*, pages 25–32, 1999.
- 1.13 Lilian Lee. Measures of distributional similarity. In *Proc. of the 37th Annual Meeting of the Association for Computational Linguistics*, pages 25–32, 1999.
- 1.14 A. Maedche and S. Staab. Discovering conceptual relations from text. In *ECAI-2000 - European Conference on Artificial Intelligence. Proceedings of the 13th European Conference on Artificial Intelligence*. IOS Press, Amsterdam, 2000.
- 1.15 A. Maedche and S. Staab. Mining ontologies from text. In *Proceedings of EKAW-2000, Springer Lecture Notes in Artificial Intelligence (LNAI-1937)*, Juan-Les-Pins, France, 2000. Springer, 2000.
- 1.16 A. Maedche and S. Staab. Ontology learning for the semantic web. *IEEE Intelligent Systems*, 16(2), 2001.
- 1.17 C.D. Manning and H. Schuetze. *Foundations of Statistical Natural Language Processing*. MIT Press, Cambridge, Massachusetts, 1999.
- 1.18 E. Morin. Automatic acquisition of semantic relations between terms from technical corpora. In *Proc. of the Fifth International Congress on Terminology and Knowledge Engineering - TKE'99*, 1999.
- 1.19 F. Pereira, N. Tishby, and L. Lee. Distribution Clustering of English Words. In *Proceedings of the ACL-93*, 1993, pages 183–199, 1993.
- 1.20 Philip S. Resnik. *Selection and Information: A Class-based Approach to Lexical Relationships*. PhD thesis, University of Pennsylvania, 1993.

- 1.21 M. Sanderson and B. Croft. Deriving Concept Hierarchies from Text. In *Proceedings of the International Conference on Information Retrieval — SIGIR'99, August 1999, Berkley CA, USA*, 1999.
- 1.22 S. Staab, C. Braun, A. Düsterhöft, A. Heuer, M. Klettke, S. Melzig, G. Neumann, B. Prager, J. Pretzel, H.-P. Schnurr, R. Studer, H. Uszkoreit, and B. Wrenger. GET-ESS — searching the web exploiting german texts. In *Proceedings of the 3rd Workshop on Cooperative Information Agents*, LNCS-1652, Berlin, 1999. Springer.
- 1.23 S. Staab, H.-P. Schnurr, R. Studer, and Y. Sure. Knowledge processes and ontologies. *IEEE Intelligent Systems*, 16(1), 2001.
- 1.24 York Sure, Michael Erdmann, Juergen Angele, Steffen Staab, Rudi Studer, and Dirk Wenke. Ontoedit: Collaborative ontology development for the semantic web. In *Proceedings of the 1st International Semantic Web Conference (ISWC2002), June 9-12th, 2002, Sardinia, Italia*, LNCS. Springer, 2002.