

Measuring the Similarity between Implicit Semantic Relations from the Web

Danushka Bollegala*
The University of Tokyo
Hongo 7-3-1, Tokyo
113-8656, Japan
danushka@mji.ci.i.u-
tokyo.ac.jp

Yutaka Matsuo
The University of Tokyo
Hongo 7-3-1, Tokyo
113-8656, Japan
matsuo@biz-model.t.u-
tokyo.ac.jp

Mitsuru Ishizuka
The University of Tokyo
Hongo 7-3-1, Tokyo
113-8656, Japan
ishizuka@i.u-tokyo.ac.jp

ABSTRACT

Measuring the similarity between semantic relations that hold among entities is an important and necessary step in various Web related tasks such as relation extraction, information retrieval and analogy detection. For example, consider the case in which a person knows a pair of entities (e.g. *Google, YouTube*), between which a particular relation holds (e.g. acquisition). The person is interested in retrieving other such pairs with similar relations (e.g. *Microsoft, Powerset*). Existing keyword-based search engines cannot be applied directly in this case because, in keyword-based search, the goal is to retrieve documents that are relevant to the words used in a query – not necessarily to the relations implied by a pair of words. We propose a relational similarity measure, using a Web search engine, to compute the similarity between semantic relations implied by two pairs of words. Our method has three components: representing the various semantic relations that exist between a pair of words using automatically extracted lexical patterns, clustering the extracted lexical patterns to identify the different patterns that express a particular semantic relation, and measuring the similarity between semantic relations using a metric learning approach. We evaluate the proposed method in two tasks: classifying semantic relations between named entities, and solving word-analogy questions. The proposed method outperforms all baselines in a relation classification task with a statistically significant average precision score of 0.74. Moreover, it reduces the time taken by Latent Relational Analysis to process 374 word-analogy questions from 9 days to less than 6 hours, with an SAT score of 51%.

Categories and Subject Descriptors

H.3.3 [Information Systems]: Information Search and Retrieval

General Terms

Algorithms

Keywords

Relational Similarity, Web Mining, Natural Language Processing

*Research Fellow of the Japan Society for the Promotion of Science (JSPS)

1. INTRODUCTION

Similarity measures can be categorized broadly into two types: *attributional* similarity measures and *relational* similarity measures. For attributional similarity measures, the objective is to compute the similarity between two given words by comparing the attributes of each word. For example, the two words *car* and *automobile* share many attributes (e.g. *has wheels, is used for transportation*). Consequently, they are considered as synonyms. On the other hand, relational similarity is the correspondence between semantic relations that exist between two *word pairs*. Word pairs that show a high degree of relational similarity are considered as *analogies*. For example, the two word pairs (*ostrich, bird*) and (*lion, cat*). Ostrich is a large bird and lion is a large cat are illustrative of high relational similarity. The semantic relation, *is a large*, pertains between the two words in each word pair.

The information available on the Web can be considered as a vast, hidden network of classes of objects (e.g. named entities) that is interconnected by various semantic relations applying to those objects. Measuring the similarity between semantic relations is an important intermediate step in various tasks in information retrieval and natural language processing such as relation extraction [7, 8, 40], in which the goal is to retrieve instances of a given relation. For example, given the relation, ACQUIRER-ACQUIREE, a relation extraction system must extract the instance (*Google, YouTube*) from the sentence *Google completed the acquisition of YouTube*. Bootstrapping methods [25, 6, 14], which require a few seeds (ca. 10 pairs of instances per relation) have extracted numerous candidate instance pairs from a text corpus. Given a set of candidate instance pairs, a relational similarity measure can be used to compute the similarity between the relations in the seeds and in the candidates. Candidate instance pairs with high relational similarity with the seed pairs can then be selected as the correct instances of a relation.

Relational similarity measures have been used to find word analogies [10, 24, 31, 33, 38]. Word analogy questions have been used from the Scholastic Aptitude Test (SAT; Educational Testing Service) to benchmark relational similarity measures. An SAT word analogy question consists of a stem word pair that acts as the question and five choice word pairs, out of which only one is analogous to the stem. A relational similarity measure is used to compare the stem word pair with each choice word pair and to select the choice word pair with the highest relational similarity as the answer.

An interesting application of relational similarity in information retrieval is to search using implicitly stated analogies [21, 37]. For example, the query “Muslim Church” is expected to return “mosque”, and the query “Hindu bible” is expected to return “the Vedas”. These queries can be formalized as word pairs: (Christian,

Church) vs. (Muslim,X), and (Christian, Bible) vs. (Hindu,Y). We can then find the words X and Y that maximize the relational similarity in each case.

Despite the wide applications of relational similarity measures, accurately measuring the similarity between implicitly stated relations remains a challenging task for several reasons. First, relational similarity is a dynamic phenomenon: it varies with time. For example, two companies can be competitors initially; subsequently one company might acquire the other. Second, there can be more than one relation between a given word pair. For example, between the two words *ostrich* and *bird*, aside from the relation *is a large*, there is also the relation *is a flightless*. A relational similarity measure must first extract all relations between the two words in each word pair before it can compute the similarity between the word pairs. Third, there can be more than one way to express a particular semantic relation in a text. For example, the three patterns – *X was acquired by Y*, *Y completed the acquisition of X*, and *Y buys X* – all indicate an acquisition relation between X and Y. In addition to the problems described above, measuring relational similarity between pairs in which one or both words are named entities (e.g., company names, personal names, locations, etc.) is even more difficult because such words are not well covered by manually created dictionaries such as WordNet¹ [23].

As described herein, we propose a relational similarity measure that uses a Web search engine to measure the similarity between implicitly stated semantic relations in two word pairs. Formally, given two word pairs, (a,b) and (c,d) , we design a function, $relsim((a,b), (c,d))$, that returns a similarity score in the range $[0, 1]$. The proposed relational similarity measure first extracts implicitly stated relations that exist between the two words in each word pair. The measure then compares the extracted relations between word pairs.

Our contributions are summarized as follows:

- We propose a shallow, lexical-patterns-based approach to represent the various semantic relations that pertain between the two words in a given word pair. The proposed pattern extraction algorithm requires no language dependent preprocessing steps such as part-of-speech tagging or dependency parsing, which can be time consuming or even infeasible at the Web scale. We extract numerous lexical patterns that describe various semantic relations.
- We present an efficient sequential clustering algorithm to cluster lexical patterns, to identify the different patterns that describe a particular semantic relation. The proposed clustering algorithm requires only one pass through the set of extracted patterns. For that reason, it scales linearly with the number of patterns. We then use the clusters to define features for a supervised metric learning algorithm.
- We evaluate the proposed method in two tasks: classifying semantic relations between named entities, and solving SAT word-analogy questions. In the relation classification task, the proposed method significantly outperforms all baselines, including the state-of-the art Latent Relational Analysis (LRA) [33]. Moreover, the proposed method achieves an SAT score of 51.1 and reduces the time taken to answer 374 questions by LRA from 9 days to less than 6 hours.

2. RELATED WORK

The Structure Mapping Theory (SMT) [15] is based on the premise that an analogy is a mapping of knowledge from one domain (base) into another (target), which conveys that a system of relations known to hold in the base also holds in the target. The target objects need not resemble their corresponding base objects. This structural view of analogy is based on the intuition that analogies are about relations, rather than simple features. Although this approach works best when the base and the target are rich in higher-order causal structures, it can fail when structures are missing or flat [39].

Turney et al. [35] combined 13 independent modules by considering the weighted sum of the outputs of each module to solve SAT analogy questions. The best performing individual module was based on the Vector Space Model (VSM). In the VSM approach [34], a vector is first created for a word pair (X,Y) by counting the frequencies of various lexical patterns containing X and Y. In their experiments, they used 128 manually created patterns such as “X of Y”, “Y of X”, “X to Y”, and “Y to X”. These patterns are then used as queries to a search engine. The numbers of hits for respective queries are used as elements in a vector to represent the word pair. Finally, the relational similarity is computed as the cosine of the angle between the two vectors that represent the two word pairs. Turney et al. [35] introduced a dataset containing 374 SAT analogy questions to evaluate relational similarity measures. An SAT analogy question consists of a stem word pair that acts as the question, and five choice word pairs. The choice word pair that has the highest relational similarity with the stem word pair is selected by the system as the correct answer. The average SAT score reported by high school students for word-analogy questions is 57%. The VSM approach achieves a score of 47% on this dataset.

Turney [31, 33] proposed Latent Relational Analysis (LRA) by extending the VSM approach in three ways: a) lexical patterns are automatically extracted from a corpus, b) the Singular Value Decomposition (SVD) is used to smooth the frequency data, and c) synonyms are used to explore variants of the word pairs. Similarly, in the VSM approach, LRA represents a word pair as a vector of lexical pattern frequencies. First, using a thesaurus, he finds related words for the two words in a word pair and create additional word pairs that are related to the original word pairs in the dataset. Second, n -grams of words are extracted from the contexts in which the two words in a word pair cooccur. The most frequent n -grams are selected as lexical patterns to represent a word pair. Then a matrix of word pairs vs. lexical patterns is created for all the word pairs in the original dataset and the additional word pairs. Elements of this matrix correspond to the frequency of a word pair in a lexical pattern. Singular value decomposition is performed on this matrix to reduce the number of columns (i.e. patterns). Finally, the relational similarity between two word pairs is computed as the average cosine similarity over the original word pairs and the additional word pairs derived from them. In fact, LRA achieves a score of 56.4% on SAT analogy questions.

Both VSM and LRA require numerous search engine queries to create a vector to represent a word pair. For example, with 128 patterns, the VSM approach requires at least 256 queries to create two pattern-frequency vectors for two word pairs before it can compute the relational similarity. In fact, LRA considers synonymous variants of the given word pairs. For that reason, it requires even more search engine queries. Methods that require numerous queries impose a heavy load on search engines. Despite efficient implementations, singular value decomposition of large matrices is time consuming. In fact, LRA takes over 9 days to process the 374 SAT analogy questions [33]. This is problematic when computing

¹<http://wordnet.princeton.edu/>

relational similarity on the scale of the Web. Moreover, in the case of named entities, thesauri of related words are not usually available or are not complete, which becomes a problem when creating the additional word pairs required by LRA.

Veale [38] proposed a relational similarity measure based on the taxonomic similarity in WordNet. The quality of a candidate analogy $A:B::C:D$ (i.e. A to B as C to D) is evaluated through comparison of the paths in the WordNet, joining A to B and C to D . Relational similarity is defined as the similarity between the $A:B$ paths and $C:D$ paths. However, WordNet does not fully cover named entities such as personal names, organizations and locations, which becomes problematic when using this method to measure relational similarity between named entities.

Using a relational similarity measure, Turney [32] proposed an unsupervised learning algorithm to extract patterns that express implicit semantic relations from a corpus. His method produces a ranked set of lexical patterns that unambiguously describes the relation between the two words in a given word pair. Patterns are ranked according to their expected relational similarity (i.e. pertinence); they are computed using an algorithm similar to LRA. To answer an SAT analogy question, first, ranked lists of patterns are generated for each of the six word pairs (one stem word pair and five choice word pairs). Then each choice is evaluated by taking the intersection of its patterns with the stem's patterns. The shared patterns are scored by the average of their rank in the stem's list and the choice's lists. The algorithm picks the choice with the lowest scoring shared pattern as the correct answer. This method reports an SAT score of 54.6%.

Relational similarity measures have been applied in natural language processing tasks such as generating word analogies [10], and classifying noun-modifier compounds based on the relation between the head and the modifier [33, 24, 9, 24]. Davidov and Rappoport [10] proposed an unsupervised algorithm to discover general semantic relations that pertain between lexical items. They represent a semantic relation with a cluster of patterns. They use the pattern clusters to generate SAT-like word analogy questions for English and Russian languages. The generated questions are then solved by human subjects. They do not evaluate their method for relational similarity between named entities.

Relational similarity measures have been used to classify the relationships between the head and the modifier in noun-compounds [33, 24, 9]. For example, in the compound *viral flu*, the *flu* (head) is *caused by* a *virus* (modifier). The *Diverse* dataset of Barker and Szpakowicz [1], which consists of 600 head-modifier pairs (noun-noun, adjective-noun and adverb-noun) is used as a benchmark dataset to evaluate relation classification of noun-compounds. Each noun-modifier pair in this dataset is annotated with one of the following five relations: *causal*, *temporal*, *spatial*, *participant*, and *quality*. Nakov and Hearst [24] proposed a linguistically motivated method that utilizes verbs, prepositions, and coordinate conjunctions that can help make explicit the hidden relations between the target nouns. They report a classification accuracy of 40.5% on the *Diverse* dataset using a single nearest neighbor classifier.

3. METHOD

3.1 Outline

Given two pairs of words (or named entities), (a,b) and (c,d) , the problem of measuring the similarity of implicit semantic relations between the two pairs can be viewed as a two-stage process.

First, we must extract the semantic relations that pertain in each word pair. We use a web search engine to retrieve the various contexts in which the two words in a word-pair cooccur. We then ex-

Google to acquire YouTube for \$1.65 billion in stock. Combination will create new opportunities for users and content owners everywhere...

Figure 1: A snippet returned for the query “*Google * * * YouTube*”.

tract lexical patterns from the retrieved contexts to represent the various semantic relations that hold between two words. However, not all patterns represent different semantic relations. A single semantic relation can be expressed using more than one lexical pattern. For example, both lexical patterns X *acquired* Y and Y *was bought by* X indicate an ACQUISITION relation between entities X and Y . We present an efficient clustering algorithm to identify the various lexical patterns that denote a particular semantic relation.

Second, we must compare the extracted semantic relations between the two word pairs to compute their relational similarity. We model this problem as one of learning a distance metric between relationally similar and dissimilar word pairs. Unlike previously proposed relational similarity measures, we do not assume semantic relations to be independent, and learn a non-Euclidean Mahalanobis distance metric.

3.2 Retrieving Contexts

We must first identify the implicitly stated relations that hold between the two words in each word pair to compute the relational similarity between two given word pairs. The context in which two words cooccur provides useful clues about the semantic relations that pertain between those words. We propose the use of text snippets retrieved using a Web search engine as an approximation of the context of two words. Snippets (also known as *dynamic teasers*) are brief summaries provided by most Web search engines along with the search results. Typically, a snippet contains a window of text selected from a document that includes the queried words. Snippets are useful for search because, most of the time, a user can read the snippet and decide whether a particular search result is relevant, without even opening the url. Using snippets as contexts is also computationally efficient because it obviates the need to download the source documents from the Web, which can be time consuming if a document is large.

A snippet for a query containing two words captures the local context in which they cooccur. For example, consider the snippet shown in Figure 1, returned by *Yahoo*² for the query “*Google * * * YouTube*”. Here, the wildcard operator “*” matches one word or none in a document. The snippet in Figure 1 is extracted from an online newspaper article about the acquisition of YouTube by Google.

To retrieve snippets for a word pair (A,B) , we use the following seven types of queries: “ $A * B$ ”, “ $B * A$ ”, “ $A * * B$ ”, “ $B * * A$ ”, “ $A * * * B$ ”, “ $B * * * A$ ”, and $A B$. The queries containing the wildcard operator “*” returns snippets in which the two words, A and B appear within a window of specified length. We designate such queries a *wildcard* queries. We search for snippets in which the query words cooccur within a maximum window of three words (tokens). This process is intended to approximate the local context of two words in a document. The quotation marks around a query will ensure that the two words appear in the specified order (e.g. A before B in snippets retrieved for the query “ $A * B$ ”). As a fallback in the case that all wildcard queries fail to return any snippets, we use the query $A B$ (without wildcards or quotations) to retrieve snippets where A and B appear in any order.

²<http://developer.yahoo.com/search/boss/>

Once we collect snippets for a word pair using the procedure described above, we remove duplicate search results. We consider two snippets to be duplicates if they contain the exact sequence of all words. Duplicate snippets exist mainly for two reasons. First, a web page can be mirrored in more than one location, and the default de-duplication mechanism of the search engine might fail to filter out the duplicates. Second, the queries we construct for a word pair are not independent. For example, a query with two wildcards might return a snippet that can also be retrieved using a query with one wildcard. However, we observed that the ranking of search results vary with the number of wildcards used. A search engine usually returns only the top ranking results (in the case of Yahoo, only the top 1000 snippets can be downloaded). We use multiple queries per word pair that induce different rankings, and aggregate search results to circumvent this limitation.

3.3 Extracting Lexical Patterns

Lexical syntactic patterns have been used in various natural language processing tasks such as extracting hypernyms [17, 30], or meronyms [2], question answering [28], and paraphrase extraction [3]. Following these previous works, we present a shallow lexical pattern extraction algorithm to represent the semantic relations between two words. The proposed method requires no language-dependent preprocessing such as part-of-speech tagging or dependency parsing, which can be both time consuming at Web scale, and likely to produce incorrect results because of the fragmented and ill-formed snippets. The pattern extraction algorithm consists of the following three steps.

Step 1: Given a context S , retrieved for a word pair (A, B) according to the procedure described in section 3.2, we replace the two words A and B , respectively, with two variables X and Y . Legal abbreviations such as *Inc.*, *Ltd.*, *Corp.*, and titles such as *Mr.*, *Ms.*, *Prof.*, *Dr.*, *Rev.* are considered as occurrences of the query terms. For example, *Google Inc.* is considered as an occurrence of the entity *Google*. We replace all numeric values by D , a marker for digits. Punctuation marks are not removed.

Step 2: We generate all subsequences of the context S that satisfy all of the following conditions.

- (i). A subsequence must contain exactly one occurrence of each X and Y (i.e., exactly one X and one Y must exist in a subsequence).
- (ii). The maximum length of a subsequence is L words.
- (iii). A subsequence is allowed to have gaps. However, we do not allow gaps of more than g number of words. Moreover, the total length of all gaps in a subsequence should not exceed G words.
- (iv). We expand all negation contractions in a context. For example, *didn't* is expanded to *did not*. We do not skip the word *not* when generating subsequences. For example, this condition ensures that from the snippet *X is not a Y*, we do not produce the subsequence *X is a Y*.

Step 3: We count the frequency of all generated subsequences for all word pairs in the dataset. We select subsequences with frequency greater than N as lexical patterns to represent the semantic relations between words.

Our pattern extraction algorithm has four parameters (ca. L , g , G and N). We set the values of those parameters experimentally, as

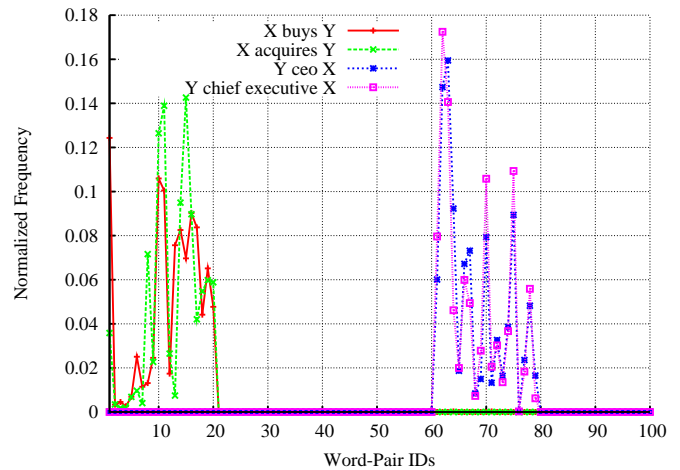


Figure 2: Distribution of four lexical patterns in word pairs.

explained later in section 4. It is noteworthy that the proposed pattern extraction algorithm considers all the words in a snippet, and is *not* limited to extracting patterns only from the mid-fix (i.e., the portion of text in a snippet that appears between the queried words). Moreover, the consideration of gaps enables us to capture relations between distant words in a snippet. We use a modified version of the *prefixspan* algorithm [26] to generate subsequences. The conditions in Step 2 are used to prune the search space, thereby reducing the number of generated subsequences in *prefixspan*. For example, some patterns extracted from the snippet shown in Figure 1 are: X to acquire Y , X acquire Y , and X to acquire Y for.

3.4 Identifying Semantic Relations

A semantic relation can be expressed using more than one pattern. For example, consider the two distinct patterns, X acquired Y , and X completed the acquisition of Y . Both these patterns indicate that there exists an acquisition relation between X and Y . It is important to know whether any correspondence pertains between the sets of patterns extracted for each word pair when we compute the relational similarity between two word pairs. We can expect a high relational similarity if there are many related patterns between two word pairs.

We use the distributional hypothesis [16] to find semantically related lexical patterns. The distributional hypothesis states that words that occur in the same context have similar meanings. The distributional hypothesis has been used in various related tasks, such as identifying related words [18], discovering inference rules [19], and extracting paraphrases [3]. If two lexical patterns are similarly distributed over a set of word pairs (i.e. occurs with the same set of word pairs), then from the distributional hypothesis it follows that the two patterns must be similar. For example, consider the distributions shown in Figure 2 for four lexical patterns: X buys Y , X acquires Y , Y CEO X , and Y chief executive X , over a set of 100 word pairs. Each distribution is normalized such that the sum of frequencies over all word pairs equals one. Figure 2 shows that the distributions of patterns Y CEO X , and Y chief executive X have a high overlap (i.e., cosine similarity of 0.969). Similarly, the distributions of patterns X buys Y , and X acquires Y show a high overlap (i.e. cosine similarity of 0.853). However, almost no overlap is apparent between other combinations of distributions. Consequently, to recognize semantically related patterns, we cluster lexical patterns using the similarity of their distributions over word pairs.

We represent a pattern p by a vector \mathbf{p} of word-pair frequencies. We designate \mathbf{p} , the *word-pair frequency* vector of pattern p . It is analogous to the *document frequency* vector of a word, as used in information retrieval. The value of the element corresponding to a word pair (a_i, b_i) in \mathbf{p} , is the frequency, $f(a_i, b_i, p)$, that the pattern p occurs with the word pair (a_i, b_i) . As demonstrated later in the experiments of this study, the proposed pattern extraction algorithm typically extracts numerous lexical patterns (more than 140,000). Clustering algorithms based on pairwise comparisons among all patterns are not feasible when the patterns are numerous. Next, we present a sequential clustering algorithm to efficiently cluster the extracted patterns.

Given a set P of patterns and a clustering similarity threshold θ , Algorithm 1 returns clusters (of patterns) that express similar semantic relations. First, in Algorithm 1, the function *SORT* sorts the patterns into descending order of their total occurrences in all word pairs. The total occurrence of a pattern p is the sum of frequencies over all word pairs (i.e., $\sum_i f(a_i, b_i, p)$). After sorting, the most common patterns appear at the beginning in P , whereas rare patterns (i.e., patterns that occur with only few word pairs) get shifted to the end. Next, in line 2, we initialize the set of clusters, C , to the empty set. The outer for-loop (starting at line 3), repeatedly takes a pattern \mathbf{p}_i from the ordered set P , and in the inner for-loop (starting at line 6), finds the cluster, $c^* \in C$ that is most similar to \mathbf{p}_i . First, we represent a cluster by the centroid of all word pair frequency vectors corresponding to the patterns in that cluster to compute the similarity between a pattern and a cluster. Next, we compute the cosine similarity between the cluster centroid (\mathbf{c}_j), and the word pair frequency vector of the pattern (\mathbf{p}_i). If the similarity between a pattern \mathbf{p}_i , and its most similar cluster, c^* , is greater than the threshold θ , we append \mathbf{p}_i to c^* (line 14). We use the operator \oplus to denote the vector addition between c^* and \mathbf{p}_i . Then we form a new cluster $\{\mathbf{p}_i\}$ and append it to the set of clusters, C , if \mathbf{p}_i is not similar to any of the existing clusters beyond the threshold θ .

The only parameter in Algorithm 1, the similarity threshold, θ , ranges in $[0, 1]$. It decides the *purity* of the formed clusters. Setting θ to a high value ensures that the patterns in each cluster are highly similar. However, high θ values also yield numerous clusters (increased model complexity). In section 4, we investigate, experimentally, the effect of θ on the overall performance of the proposed relational similarity measure.

The computational time complexity of Algorithm 1 is $O(n|C|)$, where n is the number of patterns to be clustered and $|C|$ is the number of clusters. Usually, n is much larger than $|C|$ (i.e. $n \gg |C|$). Therefore, the overall time complexity of Algorithm 1 linearly scales with the number of patterns. The sequential nature of the algorithm avoids pairwise comparisons among all patterns. Moreover, sorting the patterns by their total word-pair frequency prior to clustering ensures that the final set of clusters contains the most common relations in the dataset.

3.5 Measuring Relational Similarity

Evidence from psychological experiments suggest that similarity can be context-dependent and even asymmetric [36, 22]. Human subjects have reportedly assigned different similarity ratings to word pairs when the two words were presented in reverse order. However, experimental results investigating the effects of asymmetry, report that the average difference in ratings for a word pair is less than 5 percent [22]. Consequently, we assume relational similarity to be symmetric and limit ourselves to symmetric similarity measures. This assumption is in line with previous work on relational similarity described in section 2.

We model the problem of measuring relational similarity be-

Algorithm 1 Sequential pattern clustering algorithm.

Input: patterns $P = \{\mathbf{p}_1, \dots, \mathbf{p}_n\}$, threshold θ
Output: clusters C

```

1: SORT( $P$ )
2:  $C \leftarrow \{\}$ 
3: for pattern  $\mathbf{p}_i \in P$  do
4:    $max \leftarrow -\infty$ 
5:    $\mathbf{c}^* \leftarrow null$ 
6:   for cluster  $\mathbf{c}_j \in C$  do
7:      $sim \leftarrow \text{cosine}(\mathbf{p}_i, \mathbf{c}_j)$ 
8:     if  $sim > max$  then
9:        $max \leftarrow sim$ 
10:       $\mathbf{c}^* \leftarrow \mathbf{c}_j$ 
11:    end if
12:  end for
13:  if  $max > \theta$  then
14:     $\mathbf{c}^* \leftarrow \mathbf{c}^* \oplus \mathbf{p}_i$ 
15:  else
16:     $C \leftarrow C \cup \{\mathbf{p}_i\}$ 
17:  end if
18: end for
19: return  $C$ 

```

tween word pairs as one of learning a Mahalanobis distance metric from a given set of relationally similar and dissimilar word pairs. Given two points $\mathbf{x}_i, \mathbf{x}_j$, the (squared) Mahalanobis distance between them, $d_A(\mathbf{x}_i, \mathbf{x}_j)$, is parametrized using a positive definite matrix \mathbf{A} as follows,

$$d_A(\mathbf{x}_i, \mathbf{x}_j) = (\mathbf{x}_i - \mathbf{x}_j)^T \mathbf{A} (\mathbf{x}_i - \mathbf{x}_j). \quad (1)$$

The Mahalanobis distance is a straightforward extension of the standard Euclidean distance. In fact, if we let \mathbf{A} be the identity matrix, then the Mahalanobis distance reduces to the Euclidean distance.

The motivation behind using Mahalanobis distance to measure relational similarity is two-fold. First, Mahalanobis distance can be learned from a few data points, and efficient algorithms that can scale well to high-dimensional feature spaces are known [13, 12]. Second, unlike Euclidean distance, Mahalanobis distance does not assume that features are independent. This is particularly important for relational similarity measures because semantic relations are not always independent. A posterior analysis of the Mahalanobis matrix (\mathbf{A}) can provide useful information related to the correlation between semantic relations.

To learn a Mahalanobis distance metric, we first represent each word pair (a_i, b_i) as a feature vector \mathbf{x}_i . The j -th element of \mathbf{x}_i is the total frequency of the word pair (a_i, b_i) in the j -th cluster; it is given as $\sum_{p \in c_j} f(a_i, b_i, p)$. Here, p is a pattern in the cluster c_j , and $f(a_i, b_i, p)$ is the number of times that the word pair (a_i, b_i) appears with the pattern p . We L_2 normalize all feature vectors.

Given a set of relationally similar pairs S and dissimilar pairs D , the problem of learning a relational similarity measure becomes one of finding a positive definite matrix \mathbf{A} , such that $d_A(\mathbf{x}_i, \mathbf{x}_j) \leq u$ for all $(i, j) \in S$, and $d_A(\mathbf{x}_i, \mathbf{x}_j) \geq l$ for all $(i, j) \in D$. Here u and l respectively signify upper and lower bounds of the decision threshold, and are set experimentally as described later in section 4. Intuitively, word pairs that share identical semantic relations must have a higher relational similarity. We set an additional constraint that the learned Mahalanobis matrix \mathbf{A} must be “close” to the identity matrix \mathbf{I} to incorporate this prior knowledge in the learning problem at hand. This keeps the Mahalanobis distance similar to the Euclidean distance; it also helps to prevent overfitting of the

data. We follow the information theoretic metric learning (ITML) approach proposed by Davis et al. [13] to optimize the matrix \mathbf{A} .

We observe the fact that there exists a simple bijection (up to a scaling function) between the set of Mahalanobis distances and the set of equal mean multivariate Gaussian distributions to quantify the “closeness” between \mathbf{A} and \mathbf{I} . Assuming the equal mean to be μ , for a Mahalanobis distance parameterized by \mathbf{A} , the corresponding Gaussian is given by $p(\mathbf{x}; A) = \frac{1}{Z} \exp(-\frac{1}{2}d_A(\mathbf{x}, \mu))$, where Z is a normalizing constant and A^{-1} is the covariance of the distribution. Then, the closeness between \mathbf{A} and \mathbf{I} is measurable using the Kullback-Liebler (KL) divergence between their corresponding multivariate Gaussians:

$$KL(p(\mathbf{x}; I) \parallel p(\mathbf{x}; A)) = \int p(\mathbf{x}; I) \log \frac{p(\mathbf{x}; I)}{p(\mathbf{x}; A)} d\mathbf{x}. \quad (2)$$

Using Formula 2, the learning problem can be stated as

$$\begin{aligned} & \min_{\mathbf{A}} KL(p(\mathbf{x}; I) \parallel p(\mathbf{x}; A)) & (3) \\ \text{s.t.} & \quad d_A(\mathbf{x}_i, \mathbf{x}_j) \leq u \quad (i, j) \in S \\ & \quad d_A(\mathbf{x}_i, \mathbf{x}_j) \geq l \quad (i, j) \in D. \end{aligned}$$

The integral form of the KL divergence presented in Formula 2 is difficult to numerically optimize directly. However, it has been shown that the KL divergence between two multivariate Gaussians can be expressed as the convex combination of a Mahalanobis distance between mean vectors and the LogDet divergence between the covariance matrices [11]. Therefore, assuming that the means of the Gaussians are equal, we have

$$KL(p(\mathbf{x}; I) \parallel p(\mathbf{x}; A)) = \frac{1}{2} D_{ld}(A, I). \quad (4)$$

Here, $D_{ld}(A, B)$ is the LogDet divergence of $n \times n$ positive-definite matrices \mathbf{A} , \mathbf{B} . It is given as

$$D_{ld}(A, B) = \text{tr}(AB^{-1}) - \log \det(AB^{-1}) - n. \quad (5)$$

Finally, we incorporate slack variables into the formulation 3 to guarantee the existence of a feasible solution for \mathbf{A} , and pose the following optimization problem:

$$\begin{aligned} & \min_{\mathbf{A} \geq 0, \xi} D_{ld}(A, I) + \gamma D_{ld}(\text{diag}(\xi), \text{diag}(\xi_0)) & (6) \\ \text{s.t.} & \quad \text{tr}(A(\mathbf{x}_i - \mathbf{x}_j)(\mathbf{x}_i - \mathbf{x}_j)^T) \leq \xi_{c(i,j)} \quad (i, j) \in S \\ & \quad \text{tr}(A(\mathbf{x}_i - \mathbf{x}_j)(\mathbf{x}_i - \mathbf{x}_j)^T) \geq \xi_{c(i,j)} \quad (i, j) \in D, \end{aligned}$$

where $c(i, j)$ is the index of the (i, j) -th constraint, ξ is a vector of slack variables, initialized to ξ_0 (components of ξ_0 are initialized to u and v , respectively, for similar and dissimilar constraints), and γ is the parameter that controls the tradeoff between satisfying the constraints and minimizing $D_{ld}(A, I)$. Algorithm 2 solves the optimization problem 6 by repeatedly projecting the current solution onto a single constraint. Unlike Latent Relational Analysis [33], Algorithm 2 requires no eigen-decomposition, which is time consuming for large matrices. In Algorithm 2, a single iteration of looping through all constraints costs $O(cd^2)$, where c signifies the number of constraints, and d represents the dimensionality of feature vectors.

Once we obtain a Mahalanobis matrix \mathbf{A} from Algorithm 2, we can use Formula 1 to compute relational distances. Distance and similarity are inversely related. Therefore, it is possible to use Formula 1 directly to compare word pairs. However, if one wants to convert distance values ranging in $[0, +\infty)$ to similarity scores ranging in $[0, 1]$, it can be done using sigmoid functions [27].

Algorithm 2 Information-theoretic metric learning.

Input: \mathbf{X} , ($d \times n$ matrix); S , set of similar pairs; D , set of dissimilar pairs; u, l : distance thresholds; \mathbf{I} , identity matrix; γ , slack parameter; c , constraint index function

Output: \mathbf{A} : Mahalanobis matrix

```

1:  $\mathbf{A} \leftarrow \mathbf{I}$ ,  $\lambda_{ij} \leftarrow 0 \forall i, j$ 
2:  $\xi_{c(i,j)} \leftarrow u$  for  $(i, j) \in S$ ; otherwise  $\xi_{c(i,j)} \leftarrow l$ 
3: repeat
4:   Pick a constraint  $(i, j) \in S$  or  $(i, j) \in D$ 
5:    $p \leftarrow (\mathbf{x}_i - \mathbf{x}_j)^T \mathbf{A} (\mathbf{x}_i - \mathbf{x}_j)$ 
6:    $\delta \leftarrow 1$  if  $(i, j) \in S$ ,  $-1$  otherwise
7:    $\alpha \leftarrow \min \left( \lambda_{ij}, \frac{\delta}{2} \left( \frac{1}{p} - \frac{\gamma}{\xi_{c(i,j)}} \right) \right)$ 
8:    $\beta \leftarrow \delta \alpha / (1 - \delta \alpha \xi_{c(i,j)})$ 
9:    $\xi_{c(i,j)} \leftarrow \gamma \xi_{c(i,j)} / (\gamma + \delta \alpha \xi_{c(i,j)})$ 
10:   $\lambda_{ij} \leftarrow \lambda_{ij} - \alpha$ 
11:   $\mathbf{A} \leftarrow \mathbf{A} + \beta \mathbf{A} (\mathbf{x}_i - \mathbf{x}_j) (\mathbf{x}_i - \mathbf{x}_j)^T \mathbf{A}$ 
12: until convergence
13: return  $\mathbf{A}$ 

```

4. EXPERIMENTS

We use two different datasets to evaluate the proposed relational similarity measure in two tasks: classifying semantic relations between named entities, and solving SAT word-analogy questions. Solving SAT word analogy questions was first proposed by Turney et al. [35] as a benchmark to evaluate relational similarity measures. An SAT analogy question consists of a stem word pair that acts as the question, and five choice word pairs. A relational similarity measure under evaluation will compare the stem word pair with each choice word pair separately, and select the choice word pair with the highest relational similarity as the correct answer. The dataset contains 374 questions.

A limitation frequently associated with the SAT dataset is that it contains no named entities or relations that Web users are typically interested in, such as relations pertaining to companies or people. Consequently, in addition to the SAT dataset, we created a dataset³ containing only entity pairs to evaluate the proposed relational similarity measure. Hereinafter, we designate this as the **ENT** dataset. The ENT dataset contains 100 instances (i.e. named-entity pairs) of the following five relation types.

ACQUIRER-ACQUIREE This relation holds between pairs of company names (A, B) , where the company B (acquiree) is acquired by the company A (acquirer). We only consider acquisitions that have already completed.

PERSON-BIRTHPLACE This relation holds between pairs (A, B) , where A is the name of a person, and B is the location (place) where A was born. We consider city names and countries as locations.

CEO-COMPANY This relation holds between pairs (A, B) , where A is the chief executive officer (CEO) of a company B . We consider both current as well as past CEOs of companies.

COMPANY-HEADQUARTERS This relation holds between pairs A, B , where company A 's headquarters is located in a place B . We select names of cities as B .

PERSON-FIELD This relation holds between pairs (A, B) , where a person A is an expert or is known for his or her abilities in a

³<http://www.miv.t.u-tokyo.ac.jp/danushka/reldata.zip>

