

An Evidential Reasoning Approach to Weighted Combination of Classifiers for Word Sense Disambiguation

Cuong Anh Le¹, Van-Nam Huynh², and Akira Shimazu¹

¹ School of Information Science

² School of Knowledge Science

Japan Advanced Institute of Science and Technology

1-1, Asahidai, Nomi, Ishikawa 923-1292 Japan

Email: {cuonganh,huynh,shimazu}@jaist.ac.jp

Abstract. Arguing that various ways of using context in word sense disambiguation (WSD) can be considered as distinct representations of a polysemous word, a theoretical framework for the weighted combination of soft decisions generated by experts employing these distinct representations is proposed in this paper. Essentially, this approach is based on the Dempster-Shafer theory of evidence. By taking the confidence of individual classifiers into account, a general rule of weighted combination for classifiers is formulated, and then two particular combination schemes are derived. These proposed strategies are experimentally tested on the datasets for four polysemous words, namely *interest*, *line*, *serve*, and *hard*.

Keywords: Computational linguistics, Weighted combination of classifiers, Word sense disambiguation, Dempster-Shafer theory of evidence.

1 Introduction

Word sense disambiguation is a computational linguistics task recognized since the 1950s. Roughly speaking, word sense disambiguation involves the association of a given word in a text or discourse with a particular sense among numerous potential senses of that word. As mentioned in [5], this is an “intermediate task” necessarily to accomplish most natural language processing tasks. It is obviously essential for language understanding applications, while also at least helpful for other applications whose aim is not language understanding such as machine translation, information retrieval, among others. Since its inception, many methods involving WSD have been developed in the literature (see, e.g., [5] for a survey). During the last decade, many supervised machine learning algorithms have been used for this task, including Naïve Bayesian (NB) model, decision trees, exemplar-based model, support vector machine, maximum entropy, etc. As observed in studies of machine learning systems, although one could choose one of learning systems available to achieve the best performance

for a given pattern recognition problem, the set of patterns misclassified by the different classification systems would not necessarily overlap. This means that different classifiers may potentially offer complementary information about the patterns to be classified. This observation highly motivated the interest in combining classifiers during the recent years. Especially, classifier combination for WSD has unsurprisingly received much attention recently from the community as well, e.g., [6, 4, 12, 8, 3, 15].

As is well-known, there are basically two classifier combination scenarios. In the first scenario, all classifiers use the same representation of the input pattern. In the context of WSD, the work by Kilgarriff and Rosenzweig [6], Klein et al. [8], and Florian and Yarowsky [3] could be grouped into this first scenario. In the second scenario, each classifier uses its own representation of the input pattern. An important application of combining classifiers in this scenario is the possibility to integrate physically different types of features. In this sense, the work by Pedersen [12], Wang and Matsumoto [15] can be considered as belonging to this scenario. In this paper, we focus on the weighted combination of classifiers for WSD in the second scenario of combination strategies. Particularly, we first consider various ways of using context in WSD as distinct representations of a polysemous word under consideration, then all these representations are used as providing individual information sources to identify the meaning of the target word. We then develop a general framework for the weighted combination of individual classifiers corresponding to distinct representations. Essentially, this approach is based on Dempster-Shafer (DS) theory of evidence [13], which has been recently increasingly applied to classification problems, e.g. [2, 16]. Moreover, two combination strategies are developed and experimentally tested on the datasets for four polysemous words, namely *interest*, *line*, *serve*, and *hard*, and compared with previous studies.

The paper is organized as follows. In the next section, basic notions of DS theory will be briefly recalled. Section 3 reformulate the WSD problem so that the general framework for the weighted combination of classifiers can be formulated, and the two combination strategies can be developed. Section 4 discuss about context representation of a target word and presents our selection. Next, section 5 presents experimented results and the comparison with previous known results on the same test datasets. Finally, some conclusions are presented in Section 6.

2 Dempster-Shafer Theory of Evidence

In DS theory, a problem domain is represented by a finite set Θ of mutually exclusive and exhaustive hypotheses, called *frame of discernment* [13]. In the standard probability framework, all elements in Θ are assigned a probability. And when the degree of support for an event is known, the remainder of the support is automatically assigned to the negation of the event. On the other hand, in DS theory mass assignments are carried out for events as they know, and committing support for an event does not necessarily imply that the remaining

support is committed to its negation. Formally, a basic probability assignment (BPA, for short) is a function $m : 2^\Theta \rightarrow [0, 1]$ verifying

$$m(\emptyset) = 0, \text{ and } \sum_{A \in 2^\Theta} m(A) = 1$$

The quantity $m(A)$ can be interpreted as a measure of the belief that is committed exactly to A , given the available evidence. A subset $A \in 2^\Theta$ with $m(A) > 0$ is called a *focal element* of m . A BPA m is called to be *vacuous* if $m(\Theta) = 1$ and $m(A) = 0$ for all $A \neq \Theta$.

Two evidential functions derived from the basic probability assignment m are the belief function Bel_m and the plausibility function Pl_m , defined as

$$Bel_m(A) = \sum_{\emptyset \neq B \subseteq A} m(B), \text{ and } Pl_m(A) = \sum_{B \cap A \neq \emptyset} m(B)$$

Two useful operations that play a central role in the manipulation of belief functions are *discounting* and *Dempster's rule of combination* [13]. The discounting operation is used when a source of information provides a BPA m , but one knows that this source has probability α of reliable. Then one may adopt $(1 - \alpha)$ as one's *discount rate*, which results in a new BPA m^α defined by

$$m^\alpha(A) = \alpha m(A), \text{ for any } A \subset \Theta \quad (1)$$

$$m^\alpha(\Theta) = (1 - \alpha) + \alpha m(\Theta) \quad (2)$$

Consider now two pieces of evidence on the same frame Θ represented by two BPAs m_1 and m_2 . Dempster's rule of combination is then used to generate a new BPA, denoted by $(m_1 \oplus m_2)$ (also called the orthogonal sum of m_1 and m_2), defined as follows.

$$\begin{aligned} (m_1 \oplus m_2)(\emptyset) &= 0, \\ (m_1 \oplus m_2)(A) &= \frac{1}{1 - \sum_{B \cap C = \emptyset} m_1(B)m_2(C)} \sum_{B \cap C = A} m_1(B)m_2(C) \end{aligned} \quad (3)$$

It is worth noting that Dempster rule of combination has some attractive features such as: it is commutative and associative; given two BPAs m_1 and m_2 , if m_1 is vacuous then $m_1 \oplus m_2 = m_2$.

3 Weighted Combination of Classifiers for WSD

In this section, after reformulating the WSD problem in terms of a pattern recognition problem with multi-representation of patterns. The general framework for weighted combination of classifiers is developed for WSD problem and then, two particular combination schemes are explored.

3.1 WSD with Multi-Representation of Context

Given a polysemous word w , which may have M possible senses (classes): c_1, c_2, \dots, c_M , in a context C , the task is to determine the most appropriate sense of w . Generally, context C can be used in two ways [5]: in the *bag-of-words approach*, the context is considered as words in some window surrounding the target word w ; in the *relational information based approach*, the context is considered in terms of some relation to the target such as distance from the target, syntactic relations, selectional preferences, phrasal collocation, semantic categories, etc. As such, for a target word w , we may have different representations of context C corresponding to different views of context. Assume we have such R representations of C , say $\mathbf{f}_1, \dots, \mathbf{f}_R$, serving for the aim of identifying the right sense of the target w .

Now let us assume that we have R classifiers, each representing the context by a distinct set of features. The set of features \mathbf{f}_i , which is considered as a representation of context C of the target w , is used by the i -th classifier. Furthermore, assume that each i -th classifier (expert) is associated with a weight α_i , $0 \leq \alpha_i \leq 1$, reflecting the relative confidence in it, which may be interpreted as reliable probability of the i -th classifier in its prediction. As such representations \mathbf{f}_i 's ($i = 1, \dots, R$) are considered as distinct information sources associated with corresponding weights serving for identifying the sense of the target w . The problem now is how to combine these information sources to reach a consensus decision for identifying the sense of w .

3.2 A General Framework

Given a target word w in a context C and $\mathcal{S} = \{c_1, c_2, \dots, c_M\}$ is the set of its possible senses. Using the vocabulary of DS theory, \mathcal{S} can be called the *frame of discernment* of the problem. As mentioned above, various ways of using the context could be considered as providing different information sources to identify the meaning of the target word. Each of these information sources does not by itself provide 100% certainty as a whole piece of evidence for identifying the sense of the target. Formally, we have the available information for making the final decision on the sense of w given as follows

- R probability distributions $P(\cdot|\mathbf{f}_i)$ ($i = 1, \dots, R$) on \mathcal{S} ,
- the weights α_i of the individual information sources ($i = 1, \dots, R$)³.

From the probabilistic point of view, we may straightforwardly think of the combiner as a weighted mixture of individual classifiers defined as

$$P(c_k|\mathbf{f}_1, \dots, \mathbf{f}_R) = \frac{1}{\sum_i \alpha_i} \sum_{i=1}^R \alpha_i P(c_k|\mathbf{f}_i), \text{ for } k = 1, \dots, M \quad (4)$$

³ Note that the constraint $\sum_i \alpha_i = 1$ does not need to be imposed.

Then the target word w should be naturally assigned to the sense c_j according to the following decision rule

$$j = \arg \max_k P(c_k | \mathbf{f}_1, \dots, \mathbf{f}_R) \quad (5)$$

However, by considering the problem as that of weighted combination of evidence for decision making, in the following we will formulate a general rule of combination based on DS theory. To this end, we first adopt a probabilistic interpretation of weights. That is, the weight α_i ($i = 1, \dots, R$) is interpreted as reliable probability of the i -th classifier. This interpretation of weights seems to be especially appropriate when defining weights in terms of the accuracy of individual classifiers.

Under such an interpretation of weights, the piece of evidence represented by $P(\cdot | \mathbf{f}_i)$ should be discounted at a discount rate of $(1 - \alpha_i)$. This results in a BPA m_i verifying

$$m_i(\{c_k\}) = \alpha_i P(c_k | \mathbf{f}_i) \triangleq p_{i,k}, \text{ for } k = 1, \dots, M \quad (6)$$

$$m_i(\mathcal{S}) = 1 - \alpha_i \triangleq p_{i,\mathcal{S}} \quad (7)$$

$$m_i(A) = 0, \forall A \in 2^{\mathcal{S}} \setminus \{\mathcal{S}, \{c_1\}, \dots, \{c_M\}\} \quad (8)$$

That is, the discount rate of $(1 - \alpha_i)$ can not be distributed to anything else than \mathcal{S} , the whole frame of discernment. We are now ready to formulate our belief on the decision problem by aggregating all pieces of evidence represented by m_i 's in the general form of the following

$$m = \bigoplus_{i=1}^R m_i \quad (9)$$

where m is a BPA and \oplus is a combination operator in general.

3.3 The Discounting-and-Orthogonal Sum Combination Strategy

As discussed above, we consider each $P(\cdot | \mathbf{f}_i)$ as the belief quantified from the information source \mathbf{f}_i and the weight α_i as a “degree of trust” of \mathbf{f}_i supporting the identification for the sense of w as a whole. As mentioned in [13], an obvious way to use discounting with Dempster’s rule of combination is to discount all BPAs $P(\cdot | \mathbf{f}_i)$ ($i = 1, \dots, R$) at corresponding rates $(1 - \alpha_i)$ ($i = 1, \dots, R$) before combining them. Thus, Dempster’s rule of combination now allows us to combine BPAs m_i ($i = 1, \dots, R$) under the independent assumption of information sources for generating the BPA m , i.e. \oplus in (9) is the orthogonal sum operation.

Note that, by definition, focal elements of each m_i are either singleton sets or the whole set \mathcal{S} . It is easy to see that m also verifies this property if applicable. Interestingly, the commutative and associative properties of the orthogonal sum operation with respect to a combinable collection of BPAs m_i ($i = 1, \dots, M$) and the mentioned property essentially form the basis for developing a recursive algorithm for calculation of the BPA m . This can be done as follows.

Let $I(i) = \{1, \dots, i\}$ be the subset consisting of first i indexes of the set $\{1, \dots, R\}$. Assume that $m_{I(i)}$ is the result of combining the first i BPAs m_j , for $j = 1, \dots, i$. Let us denote

$$p_{I(i),k} \triangleq m_{I(i)}(\{c_k\}), \text{ for } k = 1, \dots, M \quad (10)$$

$$p_{I(i),\mathcal{S}} \triangleq m_{I(i)}(\mathcal{S}) \quad (11)$$

With these notations and (6)–(7), the key step in the combination algorithm is to inductively calculate $p_{I(i+1),k}$ ($k = 1, \dots, M$) and $p_{I(i+1),\mathcal{S}}$ as follows

$$p_{I(i+1),k} = \frac{1}{\kappa_{I(i+1)}} [p_{I(i),k}p_{i+1,k} + p_{I(i),k}p_{i+1,\mathcal{S}} + p_{I(i),\mathcal{S}}p_{i+1,k}] \quad (12)$$

$$p_{I(i+1),\mathcal{S}} = \frac{1}{\kappa_{I(i+1)}} (p_{I(i),\mathcal{S}}p_{i+1,\mathcal{S}}) \quad (13)$$

for $k = 1, \dots, M$, $i = 1, \dots, R-1$, and $\kappa_{I(i+1)}$ is a normalizing factor defined by

$$\kappa_{I(i+1)} = \left[1 - \sum_{j=1}^M \sum_{\substack{k=1 \\ k \neq j}}^M p_{I(i),j}p_{i+1,k} \right] \quad (14)$$

Finally, we obtain m as $m_{I(R)}$. For the purpose of decision making, we now define a probability function P_m on \mathcal{S} derived from m via the *pignistic transformation* as follows

$$P_m(c_k) = m(\{c_k\}) + \frac{1}{M}m(\mathcal{S}) \text{ for } k = 1, \dots, M \quad (15)$$

and we have the following decision rule:

$$j = \arg \max_k P_m(c_k) \quad (16)$$

It would be interesting to note that an issue may arise with the orthogonal sum operation, that is the use of the total probability mass κ associated with conflict as defined in the normalization factor. Consequently, applying it in an aggregation process may yield counterintuitive results in the face of significant conflict in certain situations as pointed out in [17]. Fortunately, in the context of the weighted combination of classifiers, by discounting all $P(\cdot|\mathbf{f}_i)$ ($i = 1, \dots, R$) at corresponding rates $(1 - \alpha_i)$ ($i = 1, \dots, R$), we actually reduce conflict between the individual classifiers before combining them.

3.4 The Discounting-and-Averaging Combination Strategy

In this strategy, instead of using Dempster's rule of combination after discounting $P(\cdot|\mathbf{f}_i)$ at the discount rate of $(1 - \alpha_i)$, we apply the averaging operation over BPAs m_i ($i = 1, \dots, R$) to obtain the BPA m defined by

$$m(A) = \frac{1}{R} \sum_{i=1}^R m_i(A) \quad (17)$$

for any $A \in 2^{\mathcal{S}}$. By definition, we get

$$m(\{c_k\}) = \frac{1}{R} \sum_{i=1}^R \alpha_i P(c_k | \mathbf{f}_i), \text{ for } k = 1, \dots, M \quad (18)$$

$$m(\mathcal{S}) = 1 - \frac{\sum_{i=1}^R \alpha_i}{R} \triangleq 1 - \bar{\alpha} \quad (19)$$

$$m(A) = 0, \forall A \in 2^{\mathcal{S}} \setminus \{\mathcal{S}, \{c_1\}, \dots, \{c_M\}\} \quad (20)$$

Note that the probability mass is unassigned to individual classes but the whole frame of discernment \mathcal{S} , $m(\mathcal{S})$, is the average of discount rates. Therefore, instead of allocating the average discount rate $(1 - \bar{\alpha})$ to $m(\mathcal{S})$ as above, we use it as a normalization factor and easily obtain

$$m(\{c_k\}) = \frac{1}{\sum_i \alpha_i} \sum_{i=1}^R \alpha_i P(c_k | \mathbf{f}_i), \text{ for } k = 1, \dots, M \quad (21)$$

$$m(A) = 0, \forall A \in 2^{\mathcal{S}} \setminus \{\{c_1\}, \dots, \{c_M\}\} \quad (22)$$

which interestingly turns out to be the weighted mixture of individual classifiers as defined in (4). Then we have the decision rule (5).

4 Representations of Context for WSD

Context plays an essentially important role in WSD and the representation choice of context is a factor which may be more important than the algorithm used for the task itself on the aspect of affecting the obtained result. For predicting senses of a word, information usually used in all studies is the topic context which is represented by bag of words. Ng and Lee [11] proposed the use of more linguistic knowledge resources including topic context, collocation of words, and a syntactic relationship verb-object, which then became popular resources for determining word sense in many papers. In [9], the authors use another information type, which is words or part-of-speech and each is assigned with its position in relation with the target word. However, in the second scenario of classifier combination strategies, according to our knowledge, only topic context with different sizes of context windows is used for creating different representations of a polysemous word, such as in Pedersen [12] and Wang and Matsumoto [15].

On the other hand, we observe that two of the most important information sources for determining the sense of a polysemous word are the topic of context and relational information representing the structural relations between the target word and the surrounding words in a local context. Under such an observation, we have experimentally designed five kinds of representation defined as follows: \mathbf{f}_1 is a set of unordered words in the large context; \mathbf{f}_2 is a set of words assigned with their positions in the local context; \mathbf{f}_3 is a set of part-of-speech tags assigned with their positions in the local context; \mathbf{f}_4 is a set of collocations of words; \mathbf{f}_5 is a set of collocations of part-of-speech tags. Symbolically, we have

- $\mathbf{f}_1 = \{w_{-n_1}, \dots, w_{-2}, w_{-1}, w_1, w_2, \dots, w_{n_1}\}$
- $\mathbf{f}_2 = \{(w_{-n_2}, -n_2), \dots, (w_{-1}, -1), (w_1, 1), \dots, (w_{n_2}, n_2)\}$
- $\mathbf{f}_3 = \{(p_{-n_3}, -n_3), \dots, (p_{-2}, -2), (p_{-1}, -1), (p_1, 1), (p_2, 2), \dots, (p_{n_3}, n_3)\}$
- $\mathbf{f}_4 = \{w_{-l} \cdots w_{-1} w w_1 \cdots w_r \mid l + r \leq n_4\}$
- $\mathbf{f}_5 = \{p_{-l} \cdots p_{-1} p p_1 \cdots p_r \mid l + r \leq n_5\}$

where w_i is the word at position i in the context of the ambiguous word w and p_i be the part-of-speech tag of w_i , with the convention that the target word w appears precisely at position 0 and i will be negative (positive) if w_i appears on the left (right) of w . In the experiment, we design the window size of topic context (for both left and right windows) as 50 for the representation \mathbf{f}_1 , i.e. $n_1 = 50$, while the window size n_i of local context as 3 for remaining representations.

5 Experiments

5.1 Computing the probabilities and determining weights

In the experiment, each individual classifier is a naive Bayesian classifier built on a context representation. We have five individual classifiers corresponding to five context representations as mentioned above. As we have seen above, in the weighted combination of classifiers we need to compute the a posteriori probabilities $P(c_k | \mathbf{f}_i)$. For the context C , suppose that the representation \mathbf{f}_i of C is represented by a set of features $\mathbf{f}_i = (f_{i,1}, f_{i,2}, \dots, f_{i,n_i})$ with the assumption that the features $f_{i,j}$ are conditionally independent, and then $P(c_k | \mathbf{f}_i)$ is computed using the following formula based on Bayes theorem.

$$P(c_k | \mathbf{f}_i) = \frac{P(\mathbf{f}_i | c_k) P(c_k)}{P(\mathbf{f}_i)} = \frac{P(c_k) \prod_{j=1}^{n_i} P(f_{i,j} | c_k)}{P(\mathbf{f}_i)} \quad (23)$$

In the experiment, we used 10-fold cross validation on the training data and then the obtained accuracies of the individual classifiers are used for weights α_i . Although we determine the weights based on the accuracies of individual classifiers, other methods of identifying the weights α_i such as using linear regression and least-squares-fit could be used. However, this is left for the long version of this paper.

5.2 Data and Result

We tested on the datasets of four words, namely *interest*, *line*, *serve*, and *hard*, which are used in numerous comparative studies of word sense disambiguation methodologies such as Pedersen [12], Ng and Lee [11], Bruce & Wiebe [1], and Leacock and Chodorow [9]. There are 2369 instances of *interest* with 6 senses, 4143 instances of *line* with 6 senses, 4378 instances of *serve* with 4 senses, and 4342 instances of *hard* with 3 senses.

In the experiment, we obtained the results using 10-fold cross validation. Table 1 shows the results obtained by using two strategies of weighted combination

of classifiers and the best results obtained by individual classifiers respectively. It is shown that both combination strategies give better results than the best individual classifier in all cases. Interestingly also, the results showed that in all cases the orthogonal sum based combination strategy is better than that based on weighted sum. This can be experimentally interpreted as follows. In our multi-representation of context, each individual classifier corresponds to a type of features so that the conditional independence assumption seems to be realistic and, consequently, the orthogonal sum based combination strategy is a suitable choice for this scheme of multi-representation of context. In addition, Table 1 also shows that both combination strategies also give better results than previous work in all cases, with the exception of *line* which corresponds to Pedersen’s method as the best.

Table 1. Results using the proposed methods and some results from previous studies. In the table, BW, M, NL, LC, and P respectively abbreviate for Bruce & Wiebe [1], Mooney [10], Ng & Lee [11], Leacock & Chodorow [9], and Pedersen [12].

(%)	BW	M	NL	LC	P	The proposed method		
						best individual classifier	based on weighted sum	based on orthogonal sum
<i>interest</i>	78	–	87	–	89	86.8	90.7	90.9
<i>line</i>	–	72	–	84	88	82.8	85.6	87.2
<i>hard</i>	–	–	–	83	–	90.2	91	91.5
<i>serve</i>	–	–	–	83	–	84.4	89	89.7

6 Conclusion

In this paper we first argued that various ways of using context in WSD can be considered as distinct representations of a polysemous word under consideration, then these representations assigned with weights are jointed into an account to identify the meaning of the target word. Based on DS theory of evidence, we developed a general framework for the weighted combination of individual classifiers corresponding to distinct representations. Moreover, two combination strategies have been developed and experimentally tested on the datasets for four polysemous words, namely *interest*, *line*, *serve*, and *hard*, and compared with previous studies. It has been shown that considering multi-representation of context significantly improves the accuracy of WSD by combining classifiers, as individual classifiers corresponding to different types of representation suitably offer complementary information about the target to be assigned to a sense. The experiment also shows that the combination strategy based on orthogonal sum is a suitable choice for this scheme of multi-representation of context.

Acknowledgement

This research is partly conducted as a program for the of “Fostering Talent in Emergent Research Fields” in Special Coordination Funds for Promoting Science and Technology by the Japanese Ministry of Education, Culture, Sports, Science and Technology.

References

1. Bruce, R. and Wiebe, J. 1994. Word-Sense Disambiguation using Decomposable Models. *Proceedings of the 32nd Annual Meeting of the Association for Computational Linguistics (ACL)*, pp. 139–145.
2. Denoeux, T., A k -nearest neighbor classification rule based on Dempster-Shafer theory, *IEEE Transactions on Systems, Man and Cybernetics* **25** (1995) 804–813.
3. Florian, R., and D. Yarowsky, Modeling consensus: Classifier combination for Word Sense Disambiguation, *Proceedings of EMNLP 2002*, pp. 25–32.
4. Hoste, V., I. Hendrickx, W. Daelemans, and A. van den Bosch, Parameter optimization for machine-learning of word sense disambiguation, *Natural Language Engineering* **8** (3) (2002) 311–325.
5. Ide, N., J. Véronis, Introduction to the Special Issue on Word Sense Disambiguation: The State of the Art, *Computational Linguistics* **24** (1998) 1–40.
6. Kilgarriff, A., and J. Rosenzweig, Framework and results for English SENSEVAL, *Computers and the Humanities* **36** (2000) 15–48.
7. Kittler, J., M. Hatef, R. P. W. Duin, and J. Matas, On combining classifiers, *IEEE Transactions on Pattern Analysis and Machine Intelligence* **20** (3) (1998) 226–239.
8. Klein, D., K. Toutanova, H. Tolga Ilhan, S. D. Kamvar, and C. D. Manning, Combining heterogeneous classifiers for Word-Sense Disambiguation, *ACL WSD Workshop*, 2002, pp. 74–80.
9. Leacock, C., M. Chodorow, and G. Miller, Using corpus statistics and WordNet relations for Sense Identification, *Computational Linguistics* **24** (1998) 147–165.
10. Mooney, R. J., Comparative experiments on Disambiguating Word Senses: An illustration of the role of bias in machine learning, *Proceedings of the Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 1996, pp. 82–91.
11. Ng, H. T., and H. B. Lee, Integrating multiple knowledge sources to Disambiguate Word Sense: An exemplar-based approach, *Proceedings of the 34th Annual Meeting of the Society for Computational Linguistics (ACL)*, 1996, pp. 40–47.
12. Pedersen, T., A simple approach to building ensembles of Naive Bayesian classifiers for Word Sense Disambiguation, *Proceedings of the North American Chapter of the Association for Computational Linguistics (NAACL)*, 2000, pp. 63–69.
13. Shafer, G., *A Mathematical Theory of Evidence* (Princeton University Press, Princeton, 1976).
14. Smets, P. and R. Kennes, The transferable belief model, *Artificial Intelligence* **66** (1994) 191–234.
15. Wang, X. J., and Y. Matsumoto, Trajectory based word sense disambiguation, *Proceedings of the 20th International Conference on Computational Linguistics*, Geneva, August 2004, pp. 903–909.
16. Wang, H., and D. Bell, Extended k -nearest neighbours based on evidence theory, *The Computer Journal* **47** (6) (2004) 662–672.
17. Zadeh, L. A., Reviews of Books: A Mathematical Theory of Evidence, *The AI Magazine* **5** (1984) 81–83.