

AN ANALYSIS OF THE RELATEDNESS BETWEEN SIMILARITY MODELS FOR WORDS

HUI LIU

School of Business Information
Shanghai University of International Business and Economics
No. 1900, Wenxiang Rd., Shanghai 201620, P. R. China
liuh@suiibe.edu.cn

Received October 2015; accepted January 2016

ABSTRACT. *Similarity computing is a popular task in natural language processing. Therefore, there are many similarity computing methods based on different models with different resources. This paper analyzes the relation between models. We first calculated the Kendall's τ correlation between the results of different models. Then, we clustered the models based on the τ . We compared 12 models on the Miller and Charles data set, and 6 models on a named entity data set. The results show that there exist some clusters of models for the similarity measurement of common nouns. The models in the clusters are all based on WordNet. The implication is that models will contain complementary information if they are not based on the same resource. Further analysis shows that we can obtain a combination model using weighted average but with lowered weights for models in a cluster. The combination model has a correlation of 0.95 with human labeling.*

Keywords: Similarity computing, Relatedness analysis

1. Introduction. Word similarity computing is the foundation for nearly all semantic applications. It is the core component of document clustering, semantic search, Q&A, etc. As a separate task, similarity computing aims at emulating human's similarity measurement by computer algorithms. Since there are many semantic models around, a practical problem for researchers is to choose one or several models to use in a specific application.

The answer for the above problem seems to be explicit. Naturally, researchers would compare their similarity models with previous works. Though most comparisons were brief, some were quite extensive at the time. For example, Li and Bandar [10] suggested an approach using multiple information sources, with detailed comparison of previous works. Iosif and Potamianos [8] suggested a method using semantic networks created from Web. The paper also had a whole section comparing 11 methods.

However, there are not many independent works which are devoted to the comparison of different models. Lin [11] analyzed several methods under the information theory view. Budanitsky and Hirst [3] evaluated the performance of several WordNet-based metrics. Pedersen et al. [20] compared different methods of semantic similarity and relatedness in the biomedical domain. Agirre et al. [1] studied the strengths and weaknesses of WordNet-based and distributional methods. Meng et al. [18] also gave a review of WordNet based methods, though their comparison is purely descriptive without evaluations. Similarly, Harispe et al. [7] reviewed different similarity models in a textbook way.

Nearly all comparisons in these works are constructed on the "performance view". The performance of a model is quantified as the Pearson correlation between model outputs and human labeling on a data set. The superiority of a model is assumed as granted if it gets better performance. We do not think it is a good assumption as we will show later in this paper. Different from previous works, this paper focuses on the relatedness of different models. We believe that one similarity model actually maps some facets of

a concept. Therefore, less-related models contain complementary information. We can know which models are complementary only by studying the relatedness between models, not just their performance.

The contributions of this paper reside mainly in two aspects. First, we suggest a novel view to evaluate a group of similarity models by exploiting their inter-relation through clustering, while previous works did not consider the relatedness of models in an empirical way. Second, we show that a better-performed model cannot simply replace an inferior one if they are from different clusters. While we are combining different models, we shall choose less-related models, since they contain complementary information. Following this criteria, we reach a correlation of 0.95 on a shortened version of the Miller and Charles set [19], by combining the outlier models with the models within one cluster, and lowering the weights for the models in that cluster. The correlation result is among the highest in literature [1, 8].

This paper is structured as the following. In the next section we give a very brief overview of similarity models discussed in this paper. In Section 3, we discuss the relatedness between models, for common words and for named entities. The last section is the conclusion.

2. Overview of Models in Analysis. The majority of researchers focus on the similarity of common words, mainly nouns. A popular method is to use WordNet. Resnik [21] introduced a method to compute similarity based on the distance of two concepts on WordNet. Wu and Palmer [22] suggested a scaled metric to calculate concept similarity, using “global depth” of a concept as the scaling factor. Lin [11] introduced an information-theoretic definition of similarity, and implemented the idea on the WordNet. Liu et al. [12] suggested using concept tree and graph derived from the structure of WordNet to calculate similarity.

Some researchers use WordNet and other resources together. Jiang and Conrath [9] combined edge and node based techniques, using corpus statistics as correction. Li and Bandar [10] combined the information from WordNet and corpus. Agirre et al. [1] introduced an SVM model to combine WordNet-based and distributional models. [5] suggested another model combining edge-counting and information content.

There are other resources in previous works, such as the Web or a dictionary. Chen et al. [4] used a double checking method to compute the similarity of two words. It is a variation of “Web co-occurrence”. If the search results of word A contains word B, and vice versa, they are somewhat similar. Bollegala et al. [2] exploited the similarity by checking lexico-syntactic patterns of the snippets returned from a search engine for measuring similarity. Liu et al. [16] proposed a method that computes the similarity of two words using expanded dictionary definition. Han et al. [6] improved the traditional PMI-IR method by augmenting the number of polysemy of each word.

This paper also contains a preliminary analysis on our previous models for named entities. [17] calculated the similarity of two (complex) named entities using the URLs of the their web hits from a search engine. [13] introduced a method similar to [16], but using Wikipedia instead of a traditional dictionary. [14] suggested several methods: simple co-occurrence on the Web by search, PMI-IR, a Web generated rough hierarchy and the hybrid model of some methods.

Table 1 shows a summary of the models compared in this paper.

3. Empirical Analysis. In this paper, we compare models for two kinds of data: common words in the Miller&Charles (M&C) data set, and a named entity set. For both sets, the analysis was performed in three steps. First, we collected the similarity values calculated by different models. Second, we calculated correlations between each pair of models. Third, the models were clustered using these correlation values.

TABLE 1. Overview of similarity models analyzed in this paper

Model	Method/Resource description	Origin
1	Definitions of machine readable dictionaries.	[16]
2	Information content from a taxonomy for the task.	[21]
3	Web-based double checking statistics.	[4]
4	Multiple information resources including taxonomy and corpus.	[10]
5	Scaled metric on the WordNet.	[22]
6	An information-theoretic similarity measure from a set of assumptions.	[11]
7	Page counts and lexico-syntactic patterns in a Web mining scheme.	[2]
8	Combination of edge counting and corpus statistics.	[9]
9	Combination of WordNet and distributional methods by SVM.	[1]
10	WordNet based hierarchy concept tree and hierarchy concept graph.	[12]
11	Combining edge-counting and information content by WordNet.	[5]
12	PMI-IR with estimates of word's number of senses.	[6]
13	A hybrid method using hierarchy, PMI and URL.	[14]
14	Web generated hierarchy, part of 13.	[14]
15	PMI-IR, part of 13.	[14]
16	Simple co-occurrence on the Web, part of 13.	[14]
17	URL-based similarity.	[17]
18	Wikipedia based similarity.	[13]

3.1. Analysis of models for the M&C data set. We used a shortened version of the M&C data set containing 26 word pairs. Model 1 to 12 in Table 1 were used in the analysis. The similarity values were collected from the previous papers listed in Table 1. We performed a simple normalization to map the values onto $[0, 1]$, if they were not in that range. To simplify discussion, we denote the outputs by the i th model as \mathbf{s}_i . The manual labeled similarity values are denoted as \mathbf{s}_m . Table 2 shows \mathbf{s}_m and each \mathbf{s}_i .

The Pearson correlations (denoted as Co. in the tables hereafter) between each \mathbf{s}_i and \mathbf{s}_m are shown in Table 2. Please be aware that since our data set is shorter than the original data set, the correlations are a little different from the previous reported values. We calculated the correlations between different models using Kendall's τ coefficient, which is a non-parametric rank correlation. The symmetric correlation matrix is shown in Table 3. All correlations are significant at the 0.01 level (1-tailed). Another point to note here is that the original data of model 12 lacks one word pair in the original paper. So the correlations between model 12 and other models are calculated from 25 pairs instead of 26 pairs.

We use Kendall's τ , instead of Pearson correlation, because unlike Pearson correlation, Kendall's τ does not assume that the population has a normal distribution, which is questionable.¹

We designed a distance metric $d_{i,j}$ of the models based on the τ values simply as $d_{i,j} = 1 - \tau_{i,j}$. Using this metric, we performed a DBScan clustering on the 12 models, since we only had distance value for clustering. Table 4 shows the resulting clusters with various ϵ and k settings.

The clustering result on the M&C data set is quite straightforward. We have a very clear set of models $C_1 = \{2, 5, 6, 8\}$. In a closer look, the four models are from Resnik [21], Wu and Palmer [22], Lin [11] and Jiang and Conrath [9]. They are all early methods based on WordNet. A further expansion of the cluster by setting ϵ to 0.3 will result in a larger cluster $C_2 = \{2, 4, 5, 6, 8, 9, 10, 11\}$, which contains all models employing WordNet as the

¹Consider a simple situation: similarity values are calculated based on the edges between nodes on the hierarchy which is a binary tree. We designed such a simulation, and the results showed that the distribution of such similarity is not normal.

TABLE 2. Comparison of models on the M&C data set

words	M&C	1	2	3	4	5	6	7	8	9	10	11	12
rooster, voyage	0.02	0	0	0	0	0	0	0.02	0.03	0	0	0	0
noon, string	0.02	0.17	0	0	0	0	0	0.02	0.05	0.07	0	0	0.14
glass, magician	0.03	0.05	0.07	0	0	0.11	0.06	0.18	0.28	0.15	0.01	0	0.45
chord, smile	0.03	0.14	0.16	0	0	0.41	0.2	0	0.31	0.07	0.06	0	0.35
lad, wizard	0.11	0.03	0.2	0.34	0.36	0.55	0.2	0.22	0.48	0.2	0.01	0.28	0.27
coast, forest	0.11	0.08	0	0	0.17	0.33	0.16	0.41	0.2	0.28	0.08	0.1	0.57
monk, slave	0.14	0	0.2	0	0.36	0.55	0.18	0.38	0.49	0.26	0.01	0.28	0.19
forest, graveyard	0.21	0.02	0	0	0.13	0	0	0.55	0.14	0.43	0.01	0.03	0.4
coast, hill	0.22	0.32	0.42	0	0.37	0.63	0.58	0.87	0.75	0.54	0.06	0.46	0.6
food, rooster	0.22	0.04	0.07	0	0	0.7	0.04	0.06	0.31	0.03	0.14	0	0.35
journey, car	0.29	0.23	0	0.41	0	0	0	0.29	0.31	0.59	0	0	0.54
lad, brother	0.42	0.32	0.2	0.36	0.36	0.55	0.2	0.34	0.46	0.62	0.01	0.28	0.45
crane, implement	0.42	0.04	0.2	0	0.37	0.63	0.39	0.13	0.42	0.41	0.2	0.39	n/a
brother, monk	0.71	0.61	0.2	0.39	0.78	0.5	0.16	0.38	0.44	0.67	1	0.84	0.51
tool, implement	0.74	0.1	0.4	0.5	0.78	0.9	0.8	0.68	0.96	0.89	0.75	0.82	0.86
bird, crane	0.74	0.32	0.62	0	0.47	0.78	0.67	0.88	0.69	0.61	0.65	0.56	0.81
bird, cock	0.76	0.69	0.62	0.46	0.78	0.91	0.83	0.59	0.79	0.74	0.66	0.83	0.77
food, fruit	0.77	0.34	0.33	0.47	0.17	0.33	0.24	1	0.65	0.75	0.62	0.05	0.74
furnace, stove	0.78	0.22	0.11	0.4	0.59	0.41	0.18	0.89	0.32	0.8	0.39	0.58	0.77
midday, noon	0.86	0.81	0.83	0.6	1	1	1	0.82	1	0.98	1	1	0.68
magician, wizard	0.88	1	0.91	0.42	1	1	1	1	1	0.79	1	1	0.88
coast, shore	0.93	0.26	0.72	0.58	0.78	0.9	0.93	0.95	0.93	0.85	0.99	0.83	0.85
boy, lad	0.94	0.67	0.56	0.57	0.78	0.9	0.85	0.97	0.77	0.92	0.99	0.83	0.7
journey, voyage	0.96	0.64	0.45	0.53	0.78	0.9	0.89	1	0.86	0.87	0.89	0.83	0.6
gem, jewel	0.96	0.86	0.99	0.71	1	1	1	0.69	1	0.98	1	1	0.97
car, automobile	0.98	0.73	0.54	0.85	1	1	1	0.98	1	1	1	1	0.92
Co.	1	0.78	0.79	0.83	0.88	0.76	0.82	0.83	0.83	0.94	0.94	0.87	0.85

TABLE 3. The correlation matrix of models on the M&C data set

	1	2	3	4	5	6	7	8	9	10	11	12
1	1.000	.582	.574	.578	.540	.575	.483	.550	.590	.591	.566	.553
2	.582	1.000	.542	.720	.771	.834	.475	.810	.571	.693	.709	.569
3	.574	.542	1.000	.627	.484	.590	.500	.613	.782	.603	.593	.568
4	.578	.720	.627	1.000	.765	.766	.546	.730	.709	.794	.958	.607
5	.540	.771	.484	.765	1.000	.824	.400	.797	.539	.711	.727	.552
6	.575	.834	.590	.766	.824	1.000	.506	.829	.625	.738	.761	.607
7	.483	.475	.500	.546	.400	.506	1.000	.531	.645	.509	.534	.607
8	.550	.810	.613	.730	.797	.829	.531	1.000	.648	.697	.697	.613
9	.590	.571	.782	.709	.539	.625	.645	.648	1.000	.633	.696	.649
10	.591	.693	.603	.794	.711	.738	.509	.697	.633	1.000	.794	.61
11	.566	.709	.593	.958	.727	.761	.534	.697	.696	.794	1.000	.583
12	.553	.569	.568	.607	.552	.607	.607	.613	.649	.61	.583	1

TABLE 4. Model clustering using DBScan on the M&C set

ϵ	$k = 1$	$k = 2$	$k = 3$
0.1	{4, 11}	\emptyset	\emptyset
0.2	{2, 5, 6, 8}, {4, 11}	{2, 5, 6, 8}	{2, 5, 6, 8}
0.3	{2, 3, 4, 5, 6, 8, 9, 10, 11}	{2, 4, 5, 6, 8, 9, 10, 11}	{2, 4, 5, 6, 8, 9, 10, 11}

main resource. The outlier models that are not in C_2 all employ other resources, regardless of their correlations with human labeling. Therefore, models using taxonomies are similar in essence, though they perform differently against human labeling. Our finding is in accordance with Lin’s [11], in which Lin said that Wu and Palmer’s work is a special situation of their work. In other words, we can say that models using different resources may contribute differently to the similarity measurement, so they are not interchangeable.

3.2. Combination of models. We can further prove our findings and justify the meaningfulness of the clusters by designing combination models. Our above discussion in Section 3.1 actually implies:

- The combination of models will lead to better evaluation results;
- Combining models in a single cluster is not as good as combining models of different resources;
- One model in a cluster can more or less represent the cluster.

Our combination method is simple. Let B be a set of the subscripts of component models. The output s_c of the combination model is the arithmetic mean of the outputs from the component models.

$$s_c = \frac{\sum_{i \in B} s_i}{|B|}$$

In Table 5 we show that if we only combine the models within the clusters C_1 or C_2 , the results are not the best. However, the best result (0.95) can be obtained by the combination of *all* the outlier models $\{1, 3, 7, 12\}$ and the combination model 2’. In other words, we lowered the weights for models in C_2 to $1/8$.

Interestingly, if we replace model 2’ with a member of C_2 , such as model 8, the performance of the new model is not affected much (-0.1), though the performances of model 2’ and 8 have big difference (0.91 vs 0.83). It shows that model 8 can represent C_2 in some ways, so our cluster is meaningful. The results also correspond to our earlier findings [15].

TABLE 5. The performance of combination of models

No.	Model Included	Co.
1’	C_1	0.83
2’	C_2	0.91
3’	2’, 1, 3, 7, 12	0.95
4’	8, 1, 3, 7, 12	0.94
5’	1-12	0.93

3.3. Analysis on the models for named entities. For the similarity computing of named entities, there is not a common data set. Therefore, we used our previous data set in [14] and tested other models on the data set. The data set is shown in Table 6. The named entities are in Chinese, but English translations are given. The data set consists of four groups. Each group has 7 to 8 word pairs. We calculated the similarity between the first word of each group with the others in that group.

We analyzed models 13 to 18 in Table 1. Table 7 shows the Pearson correlations between s_i and the human labeling results s_m . Table 8 shows the Kendall’s τ correlations between the six methods. Most correlations are significant at the 0.01 level (1 tailed). Only one is significant at the 0.05 level (noted by *). We also used DBScan to cluster the methods. The cluster is shown in Table 9.

For the clustering on the name entity data, we only start to have a cluster $\{13, 14, 17\}$ when ϵ is 0.3. Model 13 is a hybrid model, so no wonder it is related with model 14 and 17. However, even if we exclude model 13, we will find a cluster $\{14, 17, 18\}$ at $\epsilon = 0.5$. The difference between model 15/16 and the others is that these two are co-occurrence

TABLE 6. The Chinese named entity data set and translations

Group 1		Group 2	
三个火枪手	The Three Musketeers	酷睿	Core Duo
基督山伯爵	The Count of Monte Cristo	奔腾	Pentium
巴黎圣母院	The Hunchback of Notre Dame	AMD Athlon	/
雾都孤儿	Oliver Twist	Intel GMA	/
罗密欧与朱丽叶	Romeo and Juliet	AMD 700芯片组系列	AMD 700 Chipset
暮光之城	Twilight	NVIDIA GeForce 9	/
万历十五年	1587, a Year of No Significance	成都	Chengdu
大仲马	Alexandre Dumas, pere	宾得K-x	Pentax K-x
eyes on me	/		
Group 3		Group 4	
新民晚报	Xinmin Evening	windows xp	/
东方早报	Oriental Morning Post	windows 7	/
扬子晚报	Yangtze Evening	Microsoft office	/
京华时报	Beijing Times	dos	/
环球时报	Global Times	encarta	/
南风窗	South Reviews	7-zip	/
上海	Shanghai	魔兽世界	World of Warcraft
windows xp	/	微软中国研发集团	MSRA
似水流年(歌曲)	Sishui Liunian (Song)		

TABLE 7. The Pearson correlation of 6 models with human labeling on the named entity set

	13	14	15	16	17	18
Co.	0.91	0.81	0.70	0.76	0.81	0.70

TABLE 8. The correlation matrix of models on the named entity data set

	13	14	15	16	17	18
13	1.000	.776	.367	.457	.789	.530
14	.776	1.000	.385	.321*	.580	.515
15	.367	.385	1.000	.325	.370	.339
16	.457	.321*	.325	1.000	.405	.461
17	.789	.580	.370	.405	1.000	.501
18	.530	.515	.339	.461	.501	1.000

TABLE 9. Model clustering using DBScan on the named entity set

ϵ	$k = 1$	$k = 2$	$k = 3$
0.2	\emptyset	\emptyset	\emptyset
0.3	{13, 14, 17}	{13, 14, 17}	\emptyset
0.4	{13, 14, 17}	{13, 14, 17}	\emptyset
0.5	{13, 14, 17, 18}	{13, 14, 17, 18}	{13, 14, 17, 18}

based models, while the other models rely on a kind of description of the CNE. From this we can see that co-occurrence information is quite different from other information. So we agree with [3] that distributional similarity is different from similarity in concept.

4. Conclusion. In this paper we analyzed the relatedness among similarity models for common words and named entities. We calculated the Kendall's correlation between the outputs of each model and clustered the models according to a distance metric based

on these correlation values. The empirical results show that there exist some clusters for similarity models of common words. The clusters contain methods utilizing WordNet. The results imply that different resources contribute differently to the similarity measurement. Using our findings, we combined the models by average while lowering the weights for models within the cluster. The result of the combination has 0.95 correlation with human labeling on a shortened version of the Miller&Charles set.

For future work, we want to study the inter-relation between similarity models on Chinese words. It is a tough task because unlike English, there is not an agreed data set for evaluation, which means we have to implement the compared methods.

Acknowledgement. This work is supported by the Ministry of Education of China (Project of Humanities and Social Sciences, Grant. 13YJC740055).

REFERENCES

- [1] E. Agirre, E. Alfonseca, K. Hall, J. Kravalova, M. Paşca and A. Soroa, A study on similarity and relatedness using distributional and wordnet-based approaches, *Proc. of Human Language Technologies: The 2009 Annual Conference of the North American Chapter of the Association for Computational Linguistics*, pp.19-27, 2009.
- [2] D. Bollegala, Y. Matsuo and M. Ishizuka, Measuring semantic similarity between words using web search engines, *Proc. of NAACL-HLT*, 2007.
- [3] A. Budanitsky and G. Hirst, Evaluating wordnet-based measures of lexical semantic relatedness, *Computational Linguistics*, vol.32, no.1, pp.13-47, 2006.
- [4] H. Chen, M. Lin and Y. Wei, Novel association measures using web search with double checking, *Proc. of the 21st International Conference on Computational Linguistics and the 44th Annual Meeting of the ACL*, pp.1009-1016, 2006.
- [5] J.-B. Gao, B.-W. Zhang and X.-H. Chen, A wordnet-based semantic similarity measurement combining edge-counting and information content theory, *Engineering Applications of Artificial Intelligence*, vol.39, pp.80-88, 2015.
- [6] L. Han, T. Finin, P. McNamee, A. Joshi and Y. Yesha, Improving word similarity by augmenting pmi with estimates of word polysemy, *IEEE Trans. Knowledge and Data Engineering*, vol.25, no.6, pp.1307-1322, 2013.
- [7] S. Harispe, S. Ranwez, S. Janaqi and J. Montmain, Semantic similarity from natural language and ontology analysis, *Synthesis Lectures on Human Language Technologies*, vol.8, no.1, 2015.
- [8] E. Iosif and A. Potamianos, Similarity computation using semantic networks created from web-harvested data, *Natural Language Engineering*, vol.21, no.01, pp.49-79, 2015.
- [9] J. Jiang and D. Conrath, Semantic similarity based on corpus statistics and lexical taxonomy, *Proc. of International Conference on Research in Computational Linguistics*, vol.33, 1997.
- [10] Y. Li and Z. Bandar, An approach for measuring semantic similarity between words using multiple information sources, *IEEE Trans. Knowledge and Data Engineering*, vol.15, no.4, pp.871-882, 2003.
- [11] D. Lin, An information-theoretic definition of similarity, *Proc. of the 15th International Conference on Machine Learning*, pp.296-304, 1998.
- [12] H. Liu, H. Bao and D. Xu, Concept vector for semantic similarity and relatedness based on wordnet structure, *Journal of Systems and Software*, vol.85, no.2, pp.370-381, 2012.
- [13] H. Liu and Y. Chen, Computing semantic relatedness between named entities using wikipedia, *International Conference on Artificial Intelligence and Computational Intelligence*, vol.1, pp.388-392, 2010.
- [14] H. Liu and Y. Chen, Semantic similarity between complex named entities: An approach using multiple web resources, *ICIC Express Letters*, vol.5, no.1, pp.71-76, 2011.
- [15] H. Liu and R. Lu, Word similarity based on an ensemble model using ranking SVMs, *Proc. of the 2008 IEEE/WIC/ACM International Conference on Web Intelligence and Intelligent Agent Technology*, vol.3, pp.283-286, 2008.
- [16] H. Liu, J. Zhao and R. Lu, Computing semantic similarities based on machine-readable dictionaries, *IEEE International Workshop on Semantic Computing and Systems*, pp.8-14, 2008.
- [17] H. Liu, J. Zhao and R. Lu, Mining the URLs: An approach to measure the similarities between named-entities, *The 3rd International Conference on Innovative Computing Information and Control*, Dalian, China, pp.299-302, 2008.
- [18] L. Meng, R. Huang and J. Gu, A review of semantic similarity measures in wordnet, *International Journal of Hybrid Information Technology*, vol.6, no.1, pp.1-12, 2013.

- [19] G. Miller and W. Charles, Contextual correlates of semantic similarity, *Language and Cognitive Processes*, vol.6, no.1, pp.1-28, 1991.
- [20] T. Pedersen, S. V. Pakhomov, S. Patwardhan and C. G. Chute, Measures of semantic similarity and relatedness in the biomedical domain, *Journal of Biomedical Informatics*, vol.40, no.3, pp.288-299, 2007.
- [21] P. Resnik, Using information content to evaluate semantic similarity in a taxonomy, *Proc. of the 14th International Joint Conference on Artificial Intelligence*, vol.1, pp.448-453, 1995.
- [22] Z. Wu and M. Palmer, Verbs semantics and lexical selection, *Proc. of the 32nd Annual Meeting on Association for Computational Linguistics*, pp.133-138, 1994.