

*

Improving the Performance of Document Clustering with Distributional Similarities

(Jae Yun Lee)**

가 . KL 가 Jansen-Shannon
 , , 가
 . 가
 . 가
 1 2
 2 가 .
 2

ABSTRACT

In this study, measures of distributional similarity such as KL-divergence are applied to cluster documents instead of traditional cosine measure, which is the most prevalent vector similarity measure for document clustering. Three variations of KL-divergence are investigated; Jansen-Shannon divergence, symmetric skew divergence, and minimum skew divergence. In order to verify the contribution of distributional similarities to document clustering, two experiments are designed and carried out on three test collections. In the first experiment the clustering performances of the three divergence measures are compared to that of cosine measure. The result showed that minimum skew divergence outperformed the other divergence measures as well as cosine measure. In the second experiment second-order distributional similarities are calculated with Pearson correlation coefficient from the first-order similarity matrixes. From the result of the second experiment, second-order distributional similarities were found to improve the overall performance of document clustering. These results suggest that minimum skew divergence must be selected as document vector similarity measure when considering both time and accuracy, and second-order similarity is a good choice for considering clustering accuracy only.

: , , 2 , , , 2
 distributional similarity, divergence, second-order similarity, document clustering,
 automatic classification

* 2006

()
(memexlee@kgu.ac.kr)

** :2007 11 30
:2007 12 10

1.

(2005).

가

(, 2001).

(Salton & McGill 1983)

(divergence)

가

(Dagan & Lee 1999; Lee 1999; Lin 1998; Pereira, Tishby, Lee 1993; Weeds 2003).

Griffith(1981)가

White

1

2

2

A B

A B

가

2

가

$$D(q|r) = \sum_y q(y) (\log q(y) - \log r(y))$$

$$= \sum_y (q(y) \log \frac{q(y)}{r(y)})$$

KL 가 0

2

가 0

KL

3.1

가

KL

3.2

2

J

(Kullback 1968).

$$J(q,r) = D(q|r) + D(r|q)$$

4

KL

2.

r(y) 0

2.1

0

Jensen - Shannon

(skew) 가

Jensen - Shannon

Kullback - Leibler

(KL - Divergence; Kullback 1968;

Kullback & Leibler 1951)가

KL

Kullback - Leibler

Kullback -

(Lin

Leibler

(Theodoridis &

1991).

Koutroumbas 2003)

KL

0

q(y)

0

r(y)

KL

avg(q, r)

2.2

KL

$$JS(q,r) = \frac{1}{2} [D(q||aug(q,r)) + D(r||aug(q,r))]$$

가가
가

Jensen - Shannon

KL

가

Jensen - Shannon

가 KL

KL

(Lee 1999).

$$s_{\alpha}(q,r) = D(r||\alpha q + (1-\alpha)r)$$

0 1 가
=1 KL

가 . Lee(2001)
=0.99

가가
=0.99 KL

0
가 Lee(1999)

=0.99

가

가

d_i d_j t_k

가 가 $w(t_k, d_i)$, $w(t_k, d_j)$

d_i (가)가

$|d_i|$, d_i d_j KL

$D(d_i || d_j)$

$$\begin{aligned} D(d_i || d_j) &= \sum_x \frac{w(t_x, d_i)}{|d_i|} \left(\log \frac{w(t_x, d_i)}{|d_i|} - \log \frac{w(t_x, d_j)}{|d_j|} \right) \\ &= \sum_x \frac{w(t_x, d_i)}{|d_i|} \log \frac{w(t_x, d_i) |d_j|}{w(t_x, d_j) |d_i|} \\ &= \frac{1}{|d_i|} \sum_x w(t_x, d_i) \log \frac{w(t_x, d_i) |d_j|}{w(t_x, d_j) |d_i|} \end{aligned}$$

KL

d_i

d_j

가

가

Jenson - Shannon

)

(

Jenson -

Shannon

가

가 3.

SSD(d_i, d_j)

3.1

가

< 1 >

Kullback(1968)

KL

J

SSD

0

Kullback

J

가

$$SSD(d_i, d_j) = s_{0.99}(d_i, d_j) + s_{0.99}(d_j, d_i)$$

MQ Medline

MSD(d_i, d_j)

. Medline

1,033

696

30

. 30

가

8

$$MSD(d_i, d_j) = \min(s_{0.90}(d_i, d_j), s_{0.90}(d_j, d_i))$$

276

MQ

MSD

CQ

가

CACM

< 1 >

가

CQ	8	263	32.9	0.331
MQ	8	276	34.5	0.150
HQ	5	351	70.2	0.367

가 8
 263
 CACM
 CQ
 HQ
 HANTEC
 HANTEC 4
 30 12
 가
 5 351
 HQ

< 2> < 1>
 가
 가
 MQ가가 CQ HQ
 가 가 CQ
 가 가 MQ가 HQ
 MQ가가 HQ가 CQ가
 가 CQ가 가 MQ가가
 가
 , Ward 가

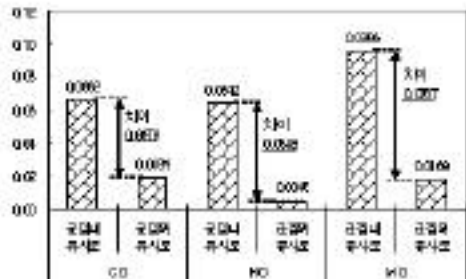
< 2> /

	CQ	HQ	MQ
	0.0663	0.0642	0.0956
	0.0184	0.0045	0.0169
	0.0479	0.0596	0.0787

(Griffith et al. 1984; Griffith et al. 1986).

가
 COS :

COS(ltf) :



TF

TF

1
 가

< 1> /

JSD : Jenson - Shannon

SSD : 가
 가
 가

MDS :
 3.2

COS(Itf) COS 가
 가

가 IDF IDF 18가
 < 3> < 4>
 IDF 가
 가 MSD
 18가 12가 1
 < 2> < 3> 가
 IDF 가

가 WACS (, 2001)
 가 WACS
 CSIM, MSD
 F 가 < 3>
 IDF 가 MSD
 COS 24.2%, COS(Itf)
 22.3%
 IDF 가 JSD SSD
 COS
 12.4% 10.9%
 MSD
 가
 가 가
 < 4>

$$WACS = \frac{1}{D} \sum_{i=1}^n \sum_{j=1}^n \frac{2|M \cap C_j|^2}{|M_j| + |C_j|}$$

< 5>, < 6>
 가 MSD 가

< 3> IDF 가 ()

		COS	COS(ltf)	JSD	SSD	MSD
CQ		0.2852 (2)	0.2537 (3)	0.2357 (4)	0.2332 (5)	0.3533 (1)
		0.2941 (5)	0.3001 (4)	0.3211 (2)	0.3486 (1)	0.3208 (3)
	Ward	0.4115 (4)	0.4197 (3)	0.4403 (2)	0.3905 (5)	0.4755 (1)
MQ		0.6859 (3)	0.6836 (5)	0.6841 (4)	0.6995 (2)	0.7154 (1)
		0.2764 (5)	0.3316 (2)	0.2867 (4)	0.3030 (3)	0.3396 (1)
	Ward	0.6383 (5)	0.7828 (4)	0.8196 (1)	0.8067 (3)	0.8171 (2)
HQ		0.7480 (4)	0.7792 (1)	0.7728 (3)	0.7783 (2)	0.6799 (5)
		0.3342 (4)	0.3395 (3)	0.3492 (1)	0.3399 (2)	0.3328 (5)
	Ward	0.5607 (5)	0.7730 (3)	0.7158 (4)	0.8057 (1)	0.7751 (2)
		0.4705 (5)	0.5182 (3)	0.5139 (4)	0.5228 (2)	0.5344 (1)

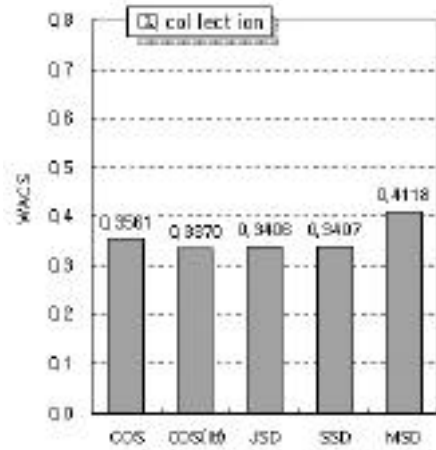
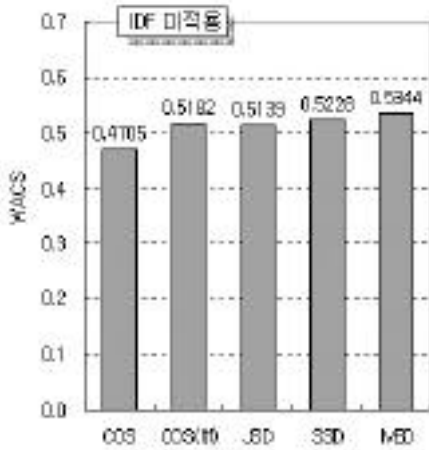
< 4> IDF 가 ()

		COS	COS(ltf)	JSD	SSD	MSD
CQ		0.4366 (1)	0.3945 (3)	0.3253 (4)	0.3174 (5)	0.3960 (2)
		0.3002 (2)	0.2959 (3)	0.2797 (4)	0.2725 (5)	0.3576 (1)
	Ward	0.4091 (4)	0.3584 (5)	0.4416 (3)	0.4818 (2)	0.5679 (1)
MQ		0.7905 (2)	0.7501 (5)	0.7552 (4)	0.7890 (3)	0.8357 (1)
		0.2923 (4)	0.2674 (5)	0.5147 (2)	0.4602 (3)	0.6019 (1)
	Ward	0.6686 (5)	0.7198 (4)	0.7518 (2)	0.7377 (3)	0.7618 (1)
HQ		0.7256 (5)	0.7557 (4)	0.7944 (1)	0.7863 (3)	0.7913 (2)
		0.3293 (5)	0.3531 (4)	0.5131 (2)	0.4829 (3)	0.6296 (1)
	Ward	0.6263 (5)	0.7558 (2)	0.7701 (1)	0.7518 (3)	0.7456 (4)
		0.5087 (5)	0.5168 (4)	0.5718 (2)	0.5644 (3)	0.6319 (1)

JSD SSD 가 , < 4> CQ
 MQ HQ 가 가
 (COS) <

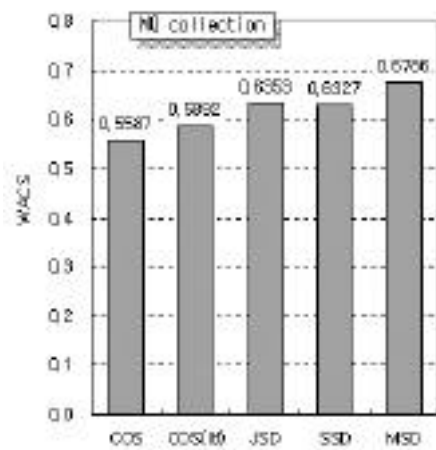
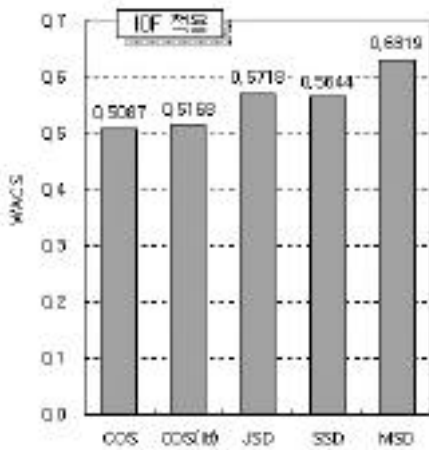
7>, < 8>, < 9> .
 MSD 가 . JSD
 가
 SSD 가
 Ward COS COS(tf)

< 7> .
 가
 가
 가 18가



< 2> IDF

< 4> CQ



< 3> IDF

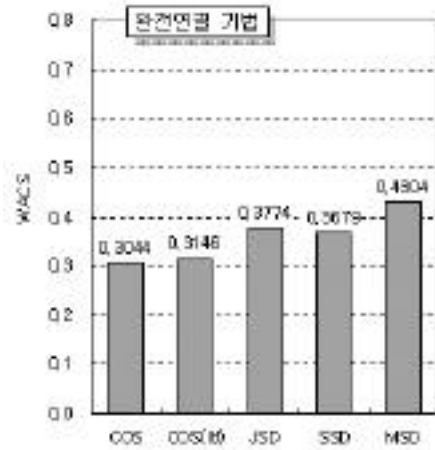
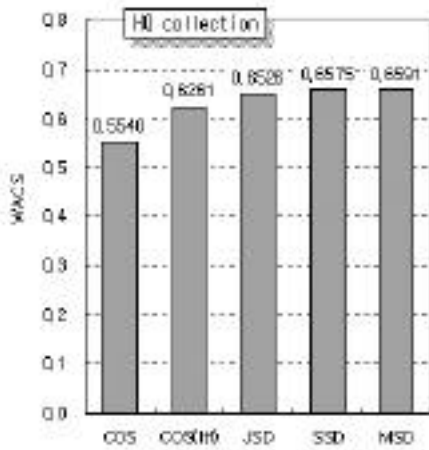
< 5> MQ

Wilcoxon

< 5> MSD COS COS(tf)
99% ,
JSD SSD 95%

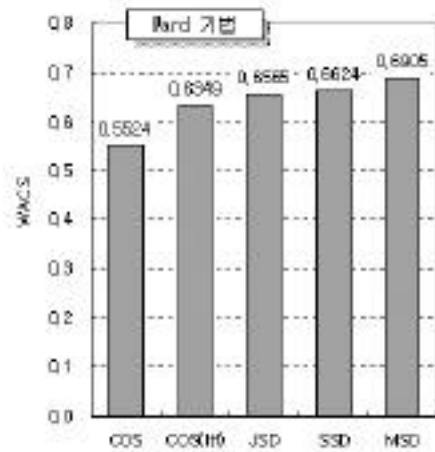
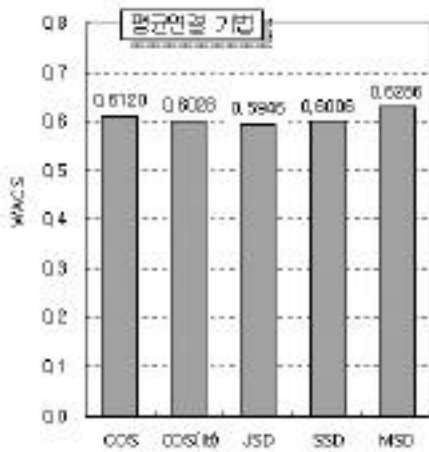
MSD

95%



< 6> HQ

< 8>



< 7>

< 9> Ward

JSD SSD COS
 95% , TF 가
 COS(ltf)

, MSD 가
 가 .
 3.3 2

< 5>
 Wilcoxon
 (99% 가 , 95% 가)

가
 1
 2
 가

	COS	COS (ltf)	JSD	SSD	MSD
COS		-	<<	<<	<<<
COS(ltf)	-		-	-	<<<
JSD	>>	-		-	<<
SSD	>>	-	-		<<
MSD	>>>	>>>	>>	>>	

1
 IDF 가
 , ,
 18가
 IDF 가
 < 6> < 7>
 가 2

가
 MSD
 가

가
 가1
 IDF 가
 < 7> , 가

가
 MSD
 MSD가 KL

MSD
 SSD
 JSD

SSD
 가

COS
 < 7> 가 MSD
 1
 24.2%(COS)
 22.3%(COS(ltf))
 2.6%(COS) 2.5%(COS(tf))

2 가 1
 < 8>, < 9>, < 10>, < 11>
 , 2 IDF 가
 < 6> 2 (IDF 가 ,)

		COS	COS(ltf)	JSD	SSD	MSD
CQ		0.3248 (5)	0.3495 (2)	0.3436 (3)	0.3379 (4)	0.3614 (1)
		0.3763 (3)	0.4006 (2)	0.3331 (5)	0.4031 (1)	0.3489 (4)
	Ward	0.3872 (4)	0.4053 (3)	0.4205 (2)	0.4263 (1)	0.3867 (5)
MQ		0.6719 (5)	0.6989 (3)	0.6936 (4)	0.7242 (1)	0.7231 (2)
		0.6567 (5)	0.7913 (1)	0.6938 (4)	0.7256 (3)	0.7800 (2)
	Ward	0.7691 (5)	0.7921 (2)	0.7906 (3)	0.8568 (1)	0.7712 (4)
HQ		0.7431 (5)	0.7858 (2)	0.7750 (4)	0.7867 (1)	0.7819 (3)
		0.6480 (5)	0.7857 (4)	0.8045 (2)	0.8004 (3)	0.8050 (1)
	Ward	0.7769 (5)	0.8520 (3)	0.7855 (4)	0.8599 (1)	0.8573 (2)
		0.5949 (5)	0.6513 (2)	0.6267 (4)	0.6579 (1)	0.6462 (3)

< 7> 2 (IDF 가 ,)

		COS	COS(ltf)	JSD	SSD	MSD
CQ		0.5862 (1)	0.4858 (3)	0.3870 (5)	0.4612 (4)	0.5160 (2)
		0.4861 (3)	0.4485 (4)	0.3656 (5)	0.5008 (1)	0.4914 (2)
	Ward	0.5081 (5)	0.5334 (4)	0.5605 (2)	0.5913 (1)	0.5602 (3)
MQ		0.8145 (3)	0.8156 (2)	0.7732 (5)	0.8162 (1)	0.8084 (4)
		0.7404 (4)	0.7358 (5)	0.7999 (1)	0.7958 (2)	0.7898 (3)
	Ward	0.6841 (5)	0.7912 (4)	0.8135 (2)	0.8132 (3)	0.8192 (1)
HQ		0.7702 (5)	0.7968 (3)	0.7793 (4)	0.8007 (1)	0.7989 (2)
		0.7711 (2)	0.7598 (3)	0.8159 (1)	0.6857 (5)	0.7438 (4)
	Ward	0.8073 (1)	0.8032 (2)	0.7463 (5)	0.7947 (4)	0.7983 (3)
		0.6853 (4)	0.6856 (3)	0.6713 (5)	0.6955 (2)	0.7029 (1)

< 9>

1

가 , 가 41가
45가
2 IDF 가 < 8>

< 8> 1

2

(IDF 가)

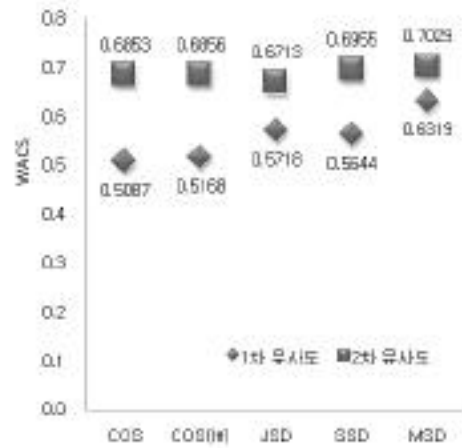
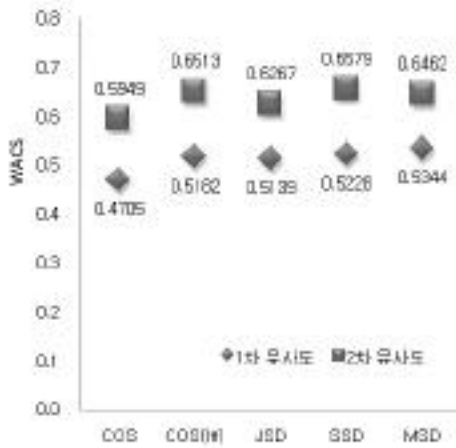
		COS	COS(ltf)	JSD	SSD	MSD
CQ		13.9%	37.7%	45.8%	44.9%	2.3%
		27.9%	33.5%	3.7%	15.6%	8.8%
	Ward	-5.9%	-3.4%	-4.5%	9.2%	-18.7%
MQ		-2.1%	2.2%	1.4%	3.5%	1.1%
		137.6%	138.6%	142.0%	139.5%	129.7%
	Ward	20.5%	1.2%	-3.5%	6.2%	-5.6%
HQ		-0.7%	0.8%	0.3%	1.1%	15.0%
		93.9%	131.4%	130.4%	135.5%	141.9%
	Ward	38.6%	10.2%	9.7%	6.7%	10.6%
		26.4%	25.7%	21.9%	25.8%	20.9%

< 9> 1

2

(IDF 가)

		COS	COS(ltf)	JSD	SSD	MSD
CQ		34.3%	23.1%	19.0%	45.3%	30.3%
		62.0%	51.6%	30.7%	83.7%	37.4%
	Ward	24.2%	48.8%	26.9%	22.7%	-1.4%
MQ		3.0%	8.7%	2.4%	3.4%	-3.3%
		153.3%	175.1%	55.4%	72.9%	31.2%
	Ward	2.3%	9.9%	8.2%	10.2%	7.5%
HQ		6.1%	5.4%	-1.9%	1.8%	1.0%
		134.2%	115.1%	59.0%	42.0%	18.1%
	Ward	28.9%	6.3%	-3.1%	5.7%	7.1%
		34.7%	32.7%	17.4%	23.2%	11.2%



< 10> 1 2
(IDF 가)

< 11> 1 2
(IDF 가)

1 2
MSD
20.9% COS 26.4%

1
IDF 가

IDF 가 < 9>
MSD 11.2%
COS 34.7%

< 10> < 11>
1 2

2
< 10> < 11>
IDF 가
10> 1
2
가 가

IDF 가 가
가 1
0.1232 , 2
0.0176
가

IDF 가 < 11>
1
2
가

2
IDF 가

IDF 가 , 1
 . 가 가 1 . 2 1
 가 가 1
 MSD . 1
 2 가
 4. 2 가가 . 1
 가가
 .
 2
 가 , 1
 (MSD 0.6319)
 . 2
 가 , 2
 .
 (MSD) 가
 24.2% . 가
 .
 1 2 2
 2 .

가

가

. 2005. 『』 :
() .
, . 2001. 『』
』 18(2): 203-230.

Dagan, Ido, Lillian Lee, and Fernando Pereira. 1999. "Similarity-based models of cooccurrence probabilities." *Machine Learning* 34(1-3): 43-69.

Griffith, A., L. A. Robinson, and P. Willett. 1984. "Hierarchic agglomerative clustering methods for automatic document classification." *Journal of Documentation*, 40(3): 175-205.

Griffiths, A., H. C. Luckhurst, and P. Willett. 1986. "Using inter document similarity information in document retrieval systems." *Journal of the American Society for Information Science*, 37(1): 3-11.

Kullback, S., and R. A. Leibler. 1951. "On information and sufficiency." *Annals of Mathematical Statistics* 22(1): 79-86.

Kullback, Solomon. 1968. *Information Theory and Statistics*, 2nd ed. New York: Dover Books.

Lee, Lillian. 1999. "Measures of distributional similarity." *Proceedings of the 37th Annual Meeting of the Association for Computational Linguistics* 25-32.

Lee, Lillian. 2001. "On the effectiveness of the skew divergence for statistical language analysis." *Proceedings of the Eighth International Workshop on Artificial Intelligence and Statistics(AISTATS-2001)*, 65-72.

Lee, Lillian, and Fernando Pereira. 1999. "Distributional similarity models: Clustering vs. nearest neighbors." *Proceedings of the 37th Annual Meeting of the*

- Association for Computational Linguistics 33-40.
- Lin, Dekang. 1998. "Automatic retrieval and clustering of similar words." Proceedings of the COLING-ACL '98, 768-773.
- Lin, Jianhua. 1991. "Divergence measures based on the Shannon entropy." IEEE Transactions on Information Theory, 37(1): 145-151.
- Pereira, Fernando, Naftali Tishby, and Lillian Lee. 1993. "Distributional clustering of English words." Proceedings of the 31st Annual Meeting of the ACL, 183-190.
- Salton, Gerard, and Michael J. McGill. 1983. Introduction to Modern Information Retrieval. New York: McGraw Hill.
- Theodoridis, S., and K. Koutroumbas. 2003. Pattern Recognition. 2nd ed. Oxford, UK: Elsevier.
- Weeds, J. E. 2003. Measures and Applications of Lexical Distributional Similarity. Ph. D. diss., University of Sussex.
- White, H. D., and B. C. Griffith. 1981. "Author cocitation: a literature measure of intellectual structure." Journal of the American Society for Information Science, 32: 163-171.

