

Latent Semantic Analysis

A Study on Content-Based Information
Retrieval System using LSA

2000 2

Latent Semantic Analysis

A study on Content-Based Information Retrieval
System using LSA

1999 12

1999 12

印

印

印



1.	1
1.1	1
1.2	3
1.3	6
2. Singular Value Decomposition (SVD)		7
2.1	7
2.2 SVD	8
2.3 SVD	11
3. Latent Semantic Analysis (LSA)		12
3.1	12
3.2 LSA	17
4.	27
4.1	27
4.2 LSA	29
4.3	29
4.4 LSA	-	33
4.5 LSA	-	36
5.	38
	39
Abstract	41

가

가

가 가

가

가

(co- occurrence)

LSA(Latent Semantic Analysis)

SVD(Singular Value Decomposition)

LSA

가

- , -

LSA

가

LSA

20,000

LSA

: LSA, SVD,

, , , ,

: 98132- 504

1.

1.1

가

가

(thesaurus)

가

Miller가 WordNet [Miller, 1990]. WordNet

(semantic network)

. WordNet

WordNet

WordNet

가

가

. WordNet 8

가

가

가

가

가 가

, WordNet

(co- occurrence)

가 .

, LSA

가

가

(co- occurrence)

. LSA

?

가

LSA가

1.2

LSA

. LSA

[Landauer, 1997].

(clustering)

(mental lexicon)

. (Warrington,

1984) LSA

가 LSA

가

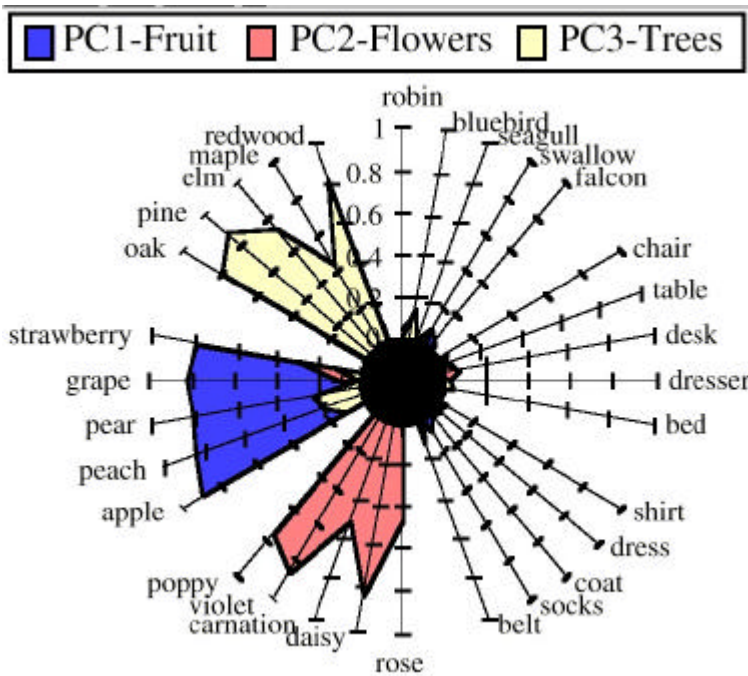
가

1 LSA

1,2,3

loading

chair, table, desk



1 Principal Components Factor Loadings

가 [Laham, 1997].

LSA

. Landauer Dumais

5

TOEFL(Test of English as a Foreign

Language)

.

가

. LSA

LSA 80

64%

[Landauer and Dumais, 1997].

LSA

100- 350

90- 200

가

10

가

.83, .65, .82

LSA

.80, .64, .84

LSA

가

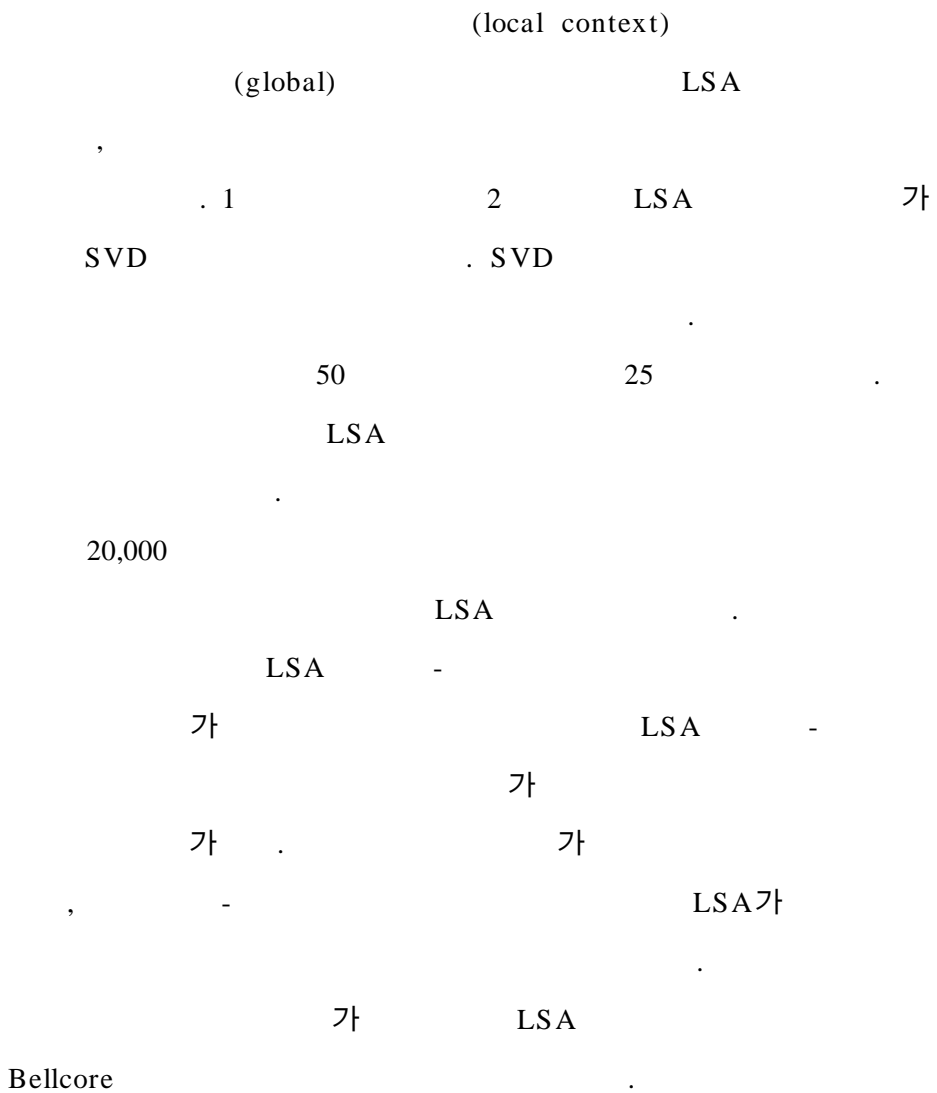
가 가

가

LSA가

[Landauer and Laham, 1997].

1.3



2. Singular Value Decomposition (SVD)

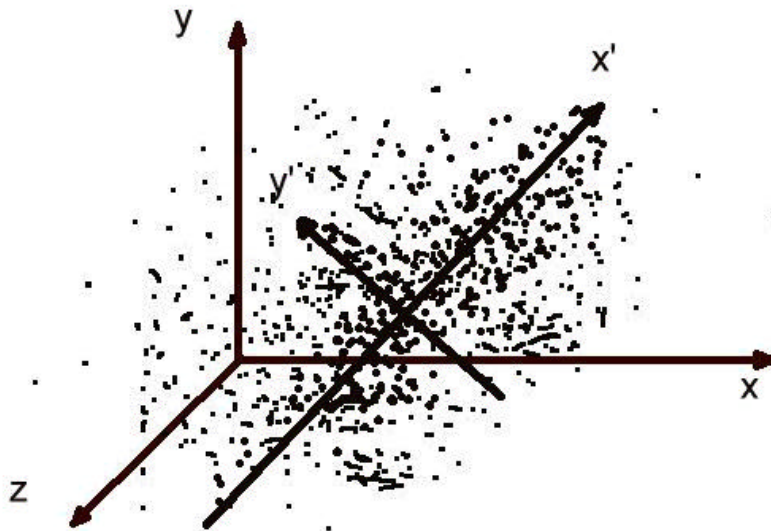
2.1

LSA

가

Singular

Value Decomposition (SVD)



2 LSA

SVD

가

2

LSA

x,y,z 3 LSA

x',y' (projection) , 3

2

가

가

2.2 SVD

Singular Value Decomposition(SVD)

(unconstrained linear least squares problem), rank

, (canonical correlation analysis)

SVD

sparse SVD 가
 Lanczos, block-Lanczos,
 (subspace iteration), (trace minimization
 method)

A m-by-n ($m \gg n$), $\text{rank}(A)=r$,
 A SVD

$$A = U \Sigma V^T \quad (1)$$

U, V orthonormal $U^T U = V^T V = I_n$ 가
 $\Sigma = \text{diag}(\sigma_1, \dots, \sigma_n)$ singular value
 $1 \leq i \leq r$ ($\sigma_i > 0$) $i \geq r+1$ ($\sigma_i = 0$)
 U, V orthogonal U, V r
 $A A^T, A^T A$ r eigenvalue orthonormal
 eigenvector A singular value $A A^T$ n
 eigenvalue Σ A
 i singular triplet(vectors) $\{u_i, \sigma_i, v_i\}$ (
) singular singular () singular
 U r left singular vector, V
 r right singular vector
 LSA m-by-n (sparse)
 (1) p singular
 가

SVD가

theorem 1. Let the SVD of A be given by (1) and

$$\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_r > \sigma_{r+1} = \dots = \sigma_n = 0$$

and let $R(A)$ and $N(A)$ denote the *range* and *null space* of A , respectively, then

1. *Rank property:*

$$\text{rank}(A) = r, \quad N(A) \equiv \text{span} \{v_{r+1}, \dots, v_n\},$$

$$R(A) \equiv \text{span} \{u_1, \dots, u_r\} \text{ where}$$

$$U = [u_1 u_2 \dots u_m] \text{ and } V = [v_1 v_2 \dots v_n].$$

2. *Dyadic decomposition:*

$$A = \sum_{i=1}^r u_i \cdot \sigma_i \cdot v_i^T$$

3. *Norms:*

$$\|A\|_F^2 = \sigma_1^2 + \dots + \sigma_r^2, \text{ and } \|A\|_2 = \sigma_1$$

Rank property A singular (qualitative)
rank (quantitative)

. *Dyadic decomposition* ,

singular

theorem 2. Let the SVD of A be given by (1) with

$$r = \text{rank}(A) \leq p = \min(m, n) \text{ and define:}$$

$$A_k = \sum_{i=1}^k u_i \cdot \sigma_i \cdot v_i^T \text{ with } k < r,$$

then

$$\min_{r(B)=k} \|A - B\|_F^2 = \|A - A_k\|_F^2 = \sigma_{k+1}^2 + \dots + \sigma_p^2.$$

(least squares)

A_k is A

rank- k

2.3 SVD

SVD

1. A bidiagonal reduce .

$$A = U_1 B V_1^T \quad U_1 \quad V_1 \text{ orthogonal}, \quad B \quad m \geq n$$

- bidiagonal, $m < n$

- bidiagonal .

2. bidiagonal B SVD $B = U_2 \Sigma V_2^T$

, $U_2 \quad V_2$ orthogonal, Σ diagonal

.

3. , A singular vector $U = U_1 U_2$,

$$V = V_1 V_2 \quad .1)$$

1) <http://nurapt.kaist.ac.kr/lapack/lug/node55.html>

, SVD (factorization) 가
 [Landauer, 1997].

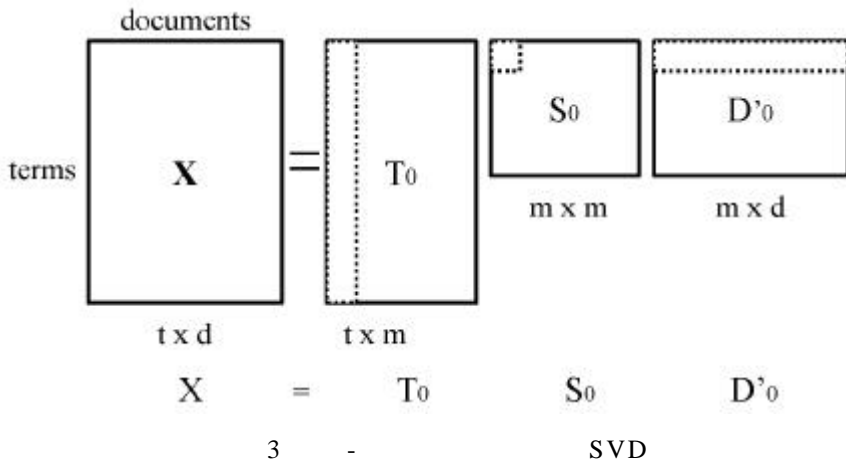
3. Latent Semantic Analysis (LSA)

3.1.

SVD
 LSA .
 가 . LSA
 d_i
 가

$$d_i = \langle t_1, \dots, t_n \rangle \quad (2)$$

t_j 가 1 0
 가
 $tf \cdot idf$
 가 [Salton and McGill, 1983].



terms)

sparse

SVD

가 d $t \times d$ X (1)-1 가 t

$$X = T_0 S_0 D_0 \quad (1)-1$$

T_0, D_0 orthonormal left singular vector right singular vector, S_0 singular value diagonal singular value

3 - SVD

T_0 : orthogonal, unit-length columns

D_0 : orthogonal, unit-length columns

S_0 : diagonal

t : X ()

d : X ()

m : X rank ($\leq \min(t,d)$)

X $T_0 S_0 D_0$
full rank . full rank
 k ($k < m$)

X 가 .
가

$$X \approx \hat{X} = TSD' \quad (1)-2$$

T : $t \times k$

S : $k \times k$

D' : $k \times d$

(1)-2 X \hat{X} k (reduce) .
 k 가
. SVD 3가 :
 i j 가, i j
가, i j 가 ?
,
.

3.2. LSA

SVD

LSA

50

c01 =

c02 = 가

c03 =

c04 =

c05 =

c06 =

c07 = 가

c08 = , ,

c09 = 가

c10 =

c11 =

c12 =

c13 =

c14 = :

c15 =

c16 =

c17 =

c18 = :

c19 =

c20 = :

c21 =

c22 =

c23 = -

c24 =

c25 =

c26 = 가가

c27 = 가

c28 =

c29 =

c30 = 가

c31 =

c32 = 가

c33 =

c34 =

c35 =
 c36 =
 c37 =
 c38 = : 가
 c39 =

 c40 =
 c41 =
 c42 = :

 c43 = :
 c44 =
 c45 =
 c46 =
 c47 =
 c48 =
 c49 =
 c50 =

t01= , t02= , t03= , t04= , t05= ,
 t06= , t07= , t08= , t09= , t10= ,
 t11= , t12= , t13= , t14= , t15= ,
 t16= , t17= , t18= , t19= , t20= ,
 t21= , t22= , t23= , t24= , t25=

	c01	c02	c03	c04	c05	c06	c07	c08	c09	c10	...
t01	0	0	0	0	0	0	0	0	0	0	...
t02	0	0	0	0	0	1	0	0	0	0	...
t03	0	0	0	0	0	0	1	0	0	0	...
t04	0	0	0	0	0	0	0	0	1	1	...
t05	0	0	0	0	1	0	0	0	0	0	...
t06	0	0	0	0	0	0	0	1	0	0	...
t07	0	0	0	0	0	0	0	0	0	0	...
t08	1	0	0	0	0	0	0	0	0	2	...
t09	0	0	0	0	0	0	0	0	0	0	...
t10	0	0	0	1	0	0	0	0	0	0	...
...

2 -

	c01	c02	c03	c04	c05	c06	c07	c08	c09	c10	...
t01	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	...
t02	0.00	0.00	0.00	0.00	0.00	2.77	0.00	0.00	0.00	0.00	...
t03	0.00	0.00	0.00	0.00	0.00	0.00	1.79	0.00	0.00	0.00	...
t04	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	2.48	2.48	...
t05	0.00	0.00	0.00	0.00	2.08	0.00	0.00	0.00	0.00	0.00	...
t06	0.00	0.00	0.00	0.00	0.00	0.00	0.00	3.91	0.00	0.00	...
t07	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	...
t08	2.08	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	4.16	...
t09	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	...
t10	0.00	0.00	0.00	2.30	0.00	0.00	0.00	0.00	0.00	0.00	...
...

3 tf*idf 가 -

1 25 ()
 . 2
 - . 가

가 (3)

$tf*idf$

3 (1)-1 SVD 3 가 .

	v01	v02	v03	v04	v05	v06	v07	v08	v09	v10	...
v01	9.41	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	...
v02	0.00	8.75	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	...
v03	0.00	0.00	8.67	0.00	0.00	0.00	0.00	0.00	0.00	0.00	...
v04	0.00	0.00	0.00	7.68	0.00	0.00	0.00	0.00	0.00	0.00	...
v05	0.00	0.00	0.00	0.00	7.43	0.00	0.00	0.00	0.00	0.00	...
v06	0.00	0.00	0.00	0.00	0.00	6.90	0.00	0.00	0.00	0.00	...
v07	0.00	0.00	0.00	0.00	0.00	0.00	6.67	0.00	0.00	0.00	...
v08	0.00	0.00	0.00	0.00	0.00	0.00	0.00	6.53	0.00	0.00	...
v09	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	6.08	0.00	...
v10	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	5.97	...
...

4 diagonal eigen value S

	t01	t02	t03	t04	t05	t06	t07	t08	t09	t10	...
v01	0.03	-0.00	0.02	-0.06	-0.09	0.00	0.02	0.41	-0.14	-0.17	...
v02	0.11	0.00	0.05	-0.01	-0.23	-0.24	-0.01	-0.08	0.23	-0.21	...
v03	0.21	-0.00	-0.04	-0.03	-0.41	-0.41	-0.19	-0.07	0.18	-0.2	...
v04	0.26	-0.00	-0.08	-0.37	0.49	-0.15	0.13	-0.09	-0.04	-0.21	...
v05	0.28	-0.00	0.43	0.34	0.13	0.06	0.16	-0.13	-0.18	-0.17	...
v06	-0.00	-0.45	0.00	0.00	0.00	-0.00	-0.00	0.00	-0.00	-0.00	...
v07	0.08	-0.00	0.05	-0.12	-0.07	0.09	0.04	0.48	0.12	-0.21	...
v08	0.21	-0.00	0.09	-0.41	0.15	-0.37	0.14	0.10	-0.01	0.47	...
v09	0.22	0.00	-0.02	-0.12	0.13	0.23	-0.06	0.18	0.24	-0.01	...
v10	0.19	0.00	0.03	-0.24	0.04	0.25	-0.30	-0.27	-0.05	-0.1	...
...

5 orthonormal eigen vector T

	c01	c02	c03	c04	c05	c06	c07	c08	c09	c10	...
v01	0.06	-0.00	0.03	-0.13	0.02	-0.11	-0.01	0.05	-0.13	0.22	...
v02	0.11	-0.00	0.00	-0.13	-0.21	0.19	-0.06	0.05	-0.23	0.09	...
v03	0.04	0.00	0.02	-0.04	-0.08	-0.02	-0.05	0.13	-0.21	0.04	...
v04	0.09	0.00	0.04	-0.12	0.07	0.04	-0.18	-0.16	-0.05	-0.04	...
v05	0.19	-0.00	0.28	0.23	0.09	0.04	-0.07	0.10	-0.03	-0.03	...
v06	0.07	0.00	0.08	0.04	-0.07	-0.09	-0.02	0.03	0.09	-0.12	...
v07	0.05	-0.00	-0.00	-0.03	-0.13	-0.11	-0.11	-0.00	-0.07	-0.00	...
v08	-0.00	-1.00	-0.00	0.00	0.00	0.00	-0.00	0.00	-0.00	-0.00	...
v09	0.11	-0.00	-0.01	-0.20	0.13	-0.03	0.08	0.00	0.08	-0.10	...
v10	0.20	-0.00	0.05	-0.40	0.31	-0.32	0.06	-0.04	-0.05	0.26	...
...

6 orthonormal eigen vector D

- 가

가 .3)

$$(a, b) = \cos \theta = \frac{a \cdot b}{\|a\| \|b\|} \quad (4)$$

(4) a b

1 0

1 8 가

‘ , ‘ , 가

‘c08 = , , ,

가

‘ , c08

‘ , ‘ , ‘ , ‘ ,

‘ (ism)’ ‘ (attention)’가

· LSA (同音異議語)

가 , ‘

2) 가 가 가

3) 가

가

	LSA		가			
	k=25		k=5		k=3	
	(0.50)	(0.41)	(0.97)	(0.96)	(1.00)	(1.00)
	(0.35)	(0.26)	(0.95)	(0.95)	(1.00)	(1.00)
	(0.40)	(0.35)	(0.98)	(0.95)	(1.00)	(1.00)
	(0.27)	(0.25)	(0.87)	(0.87)	(1.00)	(1.00)
	(0.47)	(0.46)	(0.99)	(0.98)	(1.00)	(0.99)
	(1.00)		(1.00)		(1.00)	
	(0.44)	(0.41)	(0.98)	(0.95)	(1.00)	(1.00)
	(0.33)	(0.25)	(0.96)	(0.90)	(1.00)	(1.00)
	(0.32)	(0.26)	(0.88)	(0.88)	(1.00)	(0.99)
	(0.45)	(0.34)	(0.96)	(0.89)	(1.00)	(0.99)
	(0.50)	(0.22)	(0.98)	(0.95)	(1.00)	(1.00)
	(0.87)	(0.20)	(0.99)	(0.82)	(1.00)	(0.93)
	(0.35)	(0.20)	(0.86)	(0.81)	(1.00)	(1.00)
	(0.50)	(0.24)	(0.98)	(0.96)	(1.00)	(1.00)
	(0.48)	(0.26)	(0.81)	(0.74)	(1.00)	(0.98)
	(0.40)	(0.26)	(0.83)	(0.82)	(0.97)	(0.96)
	(0.50)	(0.24)	(0.98)	(0.97)	(1.00)	(1.00)
	(0.87)	(0.25)	(0.99)	(0.75)	(1.00)	(0.94)
	(0.45)	(0.33)	(0.90)	(0.88)	(1.00)	(1.00)
	(1.00)		(1.00)		(1.00)	
	(0.44)	(0.29)	(0.98)	(0.91)	(0.99)	(0.99)
	(0.46)	(0.38)	(0.98)	(0.98)	(1.00)	(1.00)
	(0.34)	(0.29)	(0.94)	(0.93)	(1.00)	(1.00)
	(0.47)	(0.32)	(0.87)	(0.77)	(1.00)	(0.99)
	(0.48)	(0.38)	(0.99)	(0.98)	(1.00)	(1.00)

8 singular vector

-

, ' (ism)' , ' (attention)'
 , ' (transitive) ' ,
 ' ,
 ' ,
 (asymmetric) ' , ' (ism)'
 , ' ' (attention)'

eigen vector	1	2	3	4	5	6
k=25						
k=5						
k=3						

9 LSA Self Organizing Map

가

4).

9 8

Self Organizing Map

singular value 가 ,

9 가

4)

5) <http://krcogsci.snu.ac.kr/~krishna>

4.

5)

4.1

(query)

(information need)

가

10

(DDC)

가

(term- based)

(posting file)

(inverted

file)

6).

(indexing)

(retrieval)

5) <http://krcogsci.snu.ac.kr/~krishna>

6)

(presentation)

(Boolean)

SQL(Structured Query Language)

가

가

가

'1998

SQL

**SELECT BOOK FROM LIBRARY WHERE
YEAR=1998 AND AUTHOR=**

. SQL

WHERE

가

4.2 LSA

가 (key fact)
LSA, LSA
가
LSA

4.3

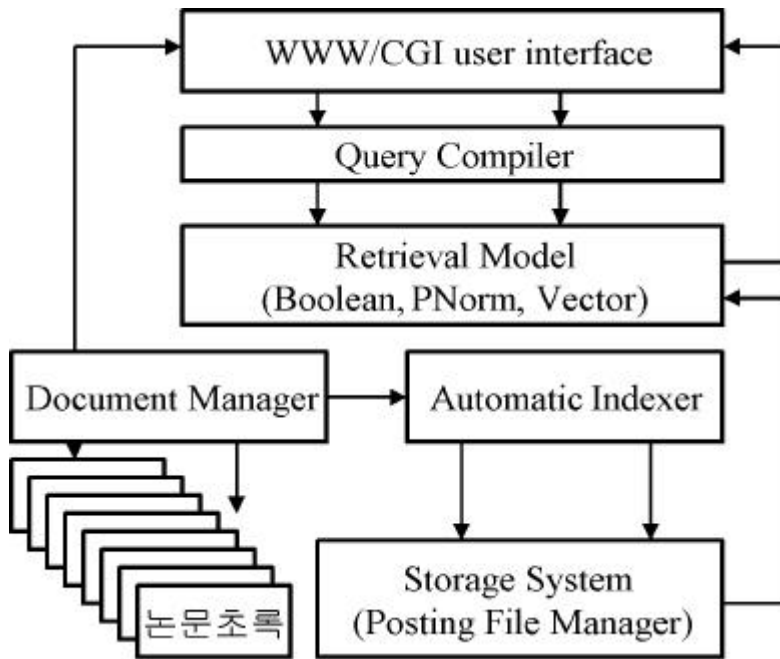
LSA가
가
1988 1998 21,556
(100MB) 7)
60 가 (document
frequency)가 4 1,500
7) 가

51,738

51,738-by- 21,556

4

Query Compiler가



4

8)

B

가

8)

1

expression) (boolean

가 (strictly)
, ' AND '
'가
' AND (OR)'

가 (fuzzy) PNorm

(vector space model)

가 [Shin
and Zhang, 1998]. LSA

(Document Manager)

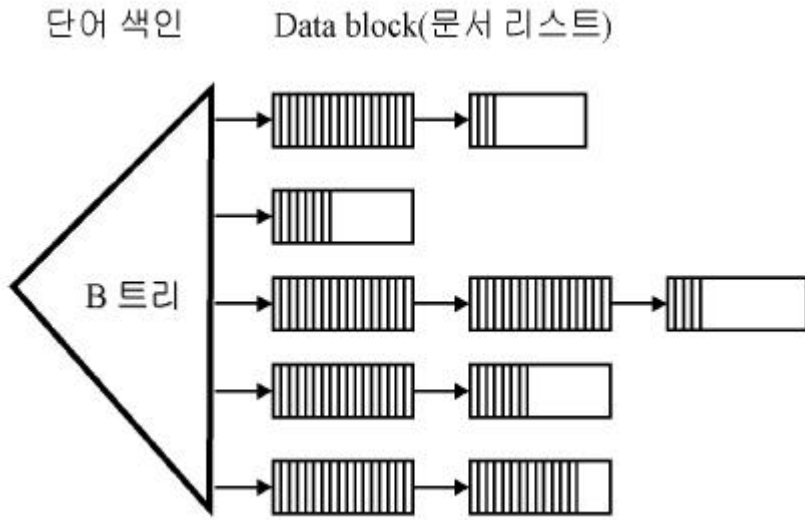
(Automatic Indexer)

가

9).

(Posting File Manager) 5 B

가가



5 Posting File Manager

9)

()

4.4 LSA

(1)-2

LSA

\hat{X}

D

orthonormal

$$DD' = I$$

$$\begin{aligned} \hat{X}\hat{X}' &= (TSD')(TSD')' \\ &= TSD'DST' \\ &= TS^2T' \end{aligned} \quad (5)$$

(5)

(4)

10

10

LSA

LSA

가

(cross-linguistic)

10)

10) (cross linguistic retrieval):

	backpropagation
가	가 가 가 households
	emotion
	가
	Forms
	存在
	實在

10 LSA

-

LSA가

가

LSA

Bellcore

Bellcore

가

Bellocore

11

.11)

11)

: <http://krcogsci.snu.ac.kr/~krishna/termsim.html>

Bellcore

: <http://superbook.bellcore.com/lsi-bin/lqiQuery>

question		question, essence, notion, communicative, conception, sciences, usages, arrive, understanding, vision
	Bellcore	question, ma, adresse, mr, my, the, le, directed, of, que
law		law, legal, Court, prosecutor, criminal, legality, reserved, judicial, Crime, provisions
	Bellcore	law, grief, real, employees, employee, grievance, problemes, dise, precise, based
faith		faith, faithful, debated, principle, originally, ruler, bringing, backdrop, grounded, Confucian
	Bellcore	faith, confiance, document, coding, respond, vast, classification, histoire, biggest, analyser
economy		economy, domestic, Economies, economic, underdeveloped, protectionism, reinforce, initiative, Privatization, internationalization
	Bellcore	economy, economie, derniers, compared, heavily, cereales, subvention, americans, subsidized, speaking
politics		politics, political, republic, faction, politically, ruling, democracy, corrupt, regime, regionalism
	Bellcore	politics, assainissement, risque, progressiste, heart, donnent, defend, ecouter, laissez, themselves

11 Bellcore

-

4.5 LSA

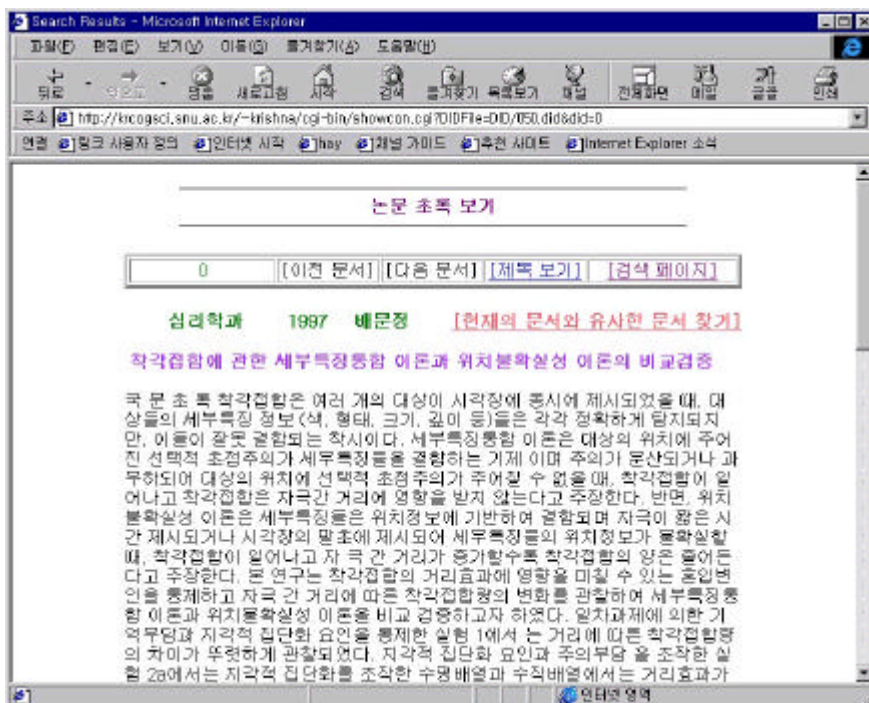
-

가

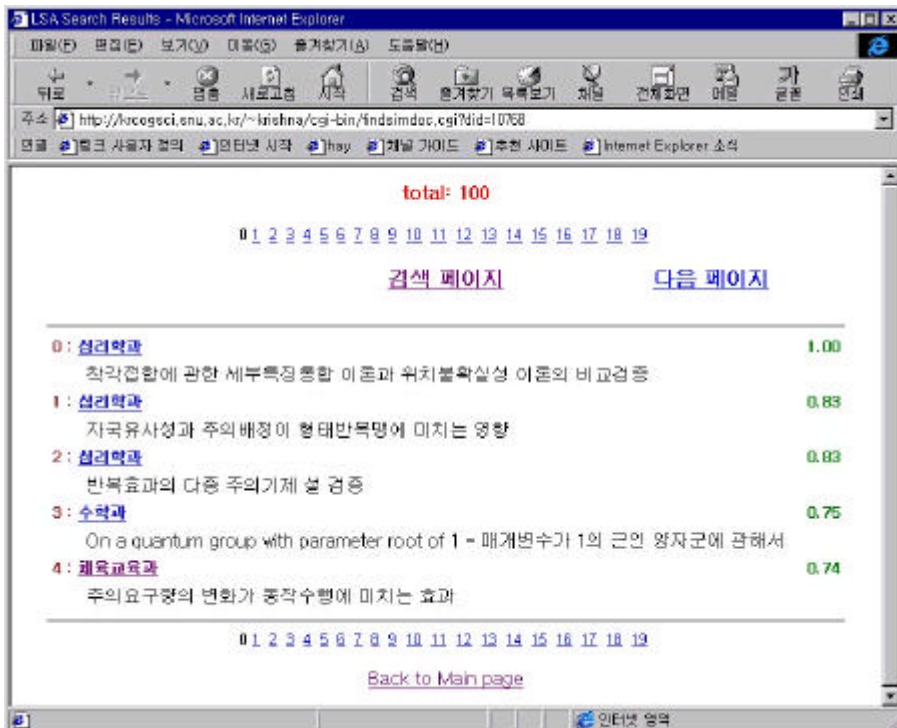
가

3

[



6



5.

가

LSA

. LSA가

LSA

가

가

가

가

가

가

[

, 1999].

- [, 1999] , ,
JCEANF-99, PP. 293-296. 1999.
- [Berry et al, 1993] Berry, M.W., Do, T., O'Brien, G.W., Krishna, V., and Varadhan, S, SVDPACKC (Version 1.0) User's Guide. 1993.
- [Berry et al, 1995a] Berry, M.W., Dumais, S.T., and Letsche., T.A., Computational Methods for Intelligent Information Access, *Proceedings of Supercomputing'95*, San Diego, CA. 1995.
- [Berry et al, 1995b] Berry, M.W., Dumais, S.T., and Shippy, A.T., A Case Study of Latent Semantic Indexing. 1995.
- [Berry et al, 1995c] Berry, M.W., Dumais, S.T., and O'Brien, G.W., Using Linear Algebra for Intelligent Information Retrieval. *SIAM Review* 37:4, pp. 573-595, 1995.
- [Berry et al, 1996] Berry, M.W. and Ricardo D.F., Low-Rank Orthogonal Decompositions for Information Retrieval Applications., *Numerical Linear Algebra with Applications* 3:4, pp. 301-328. 1996.
- [Jiang, 1997] Jiang, J., Using Latent Semantic Indexing for Data Mining., MS Thesis, Department of Computer Science, University of Tennessee. 1997.
- [Landauer and Dumais, 1997] Landauer, T. K. and Dumais, S. T. A solution to Plato's problem: The Latent Semantic Analysis theory of the acquisition, induction, and representation of knowledge. *Psychological Review*, 104, 211-240. 1997.
- [Landauer and Laham, 1997] Landauer, T. K. and Laham, D. Learning Human-like Knowledge by Singular Value Decomposition: A

- Progress Report. *NIPS*, 10, 45- 51. 1997.
- [Laham, 1997] Laham, D., Latent Semantic Analysis Approaches to Categorization. *Proceedings of the Cognitive Science Society*. 1997.
- [Letsche, 1996] Letsche, T.A., Toward Large-Scale Information Retrieval Using Latent Semantic Indexing. MS Thesis, Department of Computer Science, University of Tennessee. 1996.
- [Letsche and Berry, 1997] Letsche, T.A. and Berry, M.W., Large-Scale Information Retrieval with Latent Semantic Indexing, *Information Sciences - Applications* 100, pp. 105- 137. 1997.
- [Miller, 1990] Miller G., Five Papers on WordNet, *Special Issue of International Journal of Lexicography* 3(4). 1990.
- [O'Brien, 1994] O'Brien, G.W., Information Management Tools for Updating an SVD- Encoded Indexing Scheme. 1994.
- [Salton and McGill, 1983] Salton, G. and McGill, M.J., Introduction to Modern Information Retrieval. McGraw- Hill. 1983.
- [Shin and Zhang, 1988] Shin, D.H. and Zhang, B.T., Automatic Query Generation for the TREC-7 Ad-Hoc Task, *Proceedings of the 7th Text REtrieval Conference (TREC-7)*, NIST Special Publication, 1998.
- [Warrington and Shallice, 1984] Warrington, E. K. and Shallice, T. Category- specific semantic impairments. *Brain*, 107, 829- 853. 1984.

Abstract

Natural language processing is one of the most difficult research areas because of its complexity and diversity. So, syntactic rules based on human intuition have been mainly used to analyze human language until now. Whereas a rule-based approach is effective for the analysis of a small corpus, its performance degrades rapidly as the corpus size increases (over tens of thousands of documents) due to the many exceptional cases. Contrary to this, a statistical approach can attack large text corpus easily using simple information such as co-occurrence. A statistical approach can be viewed as processing lower level data to produce more higher and more general level representation. LSA (Latent Semantic Analysis) can compute higher level representation directly by statistical method of SVD(Singular Value Decomposition).

In the previous learning method, only some pre-selected features are considered in learning concept for its computational burden. But LSA produce high dimensional vector space called 'concept space' through the entire large text corpus statistically. These 'concept space', made by purely statistical method, showed similar results with human mental lexicon in many experiments. In this paper, The human-like property (semantic property) of LSA is applied to information retrieval system. Information

retrieval system using LSA can really be called 'content-based information retrieval system'. Experiments were done under the 20,000 paper abstracts provided by Seoul National University Library (<http://solarsnet.snu.ac.kr>). The results were excellent, even better than Bellcore LSI demo system.

Key words: LSA, SVD, IR, mental lexicon, concept space, similarity measure, learning

student number : 98132-504

2

2

가

가

가

1

가

가

가

가

가

1

가

2

98

1

가

PCA

ETRI

(知音)

1

, 가