

*

-

-

()

bagging

()

()

/

Rocchio

/

-

-

- SVM -

- Rocchio -

:

bagging

[]

Rocchio

"

"

SVM

[]

SVM

Rocchio

[] bagging

[]
()
Joachims SVM []
[] [] " "
[] -
[]
KNN Rocchio)
) - ()
() []
[] ()
KNN KNN SVM
Rocchio (KNNM)
20-NewsGroup
()
(KNNM SVM) []
[]
[]
(SVM) F1
[] []
Rocchio :
[] [] KNN
[] []
[] SVM
[] bagging [] boosting []

.....

boosting

d

(ACC)

d d

WLC bagging

[] DCS

()

[]

[]

(WLC) (MV)

(DCS) (ACC)

[]

(MV)

K

$\frac{k+1}{2}$

d

K-means HAC

K- HAC

[] means

(WLC)

d

W_j

Φ_j

[] d

[] (DCS)

Φ_i $\{\Phi_1, \dots, \Phi_i\}$

k

Leave-one- d

[] out

d

(b) []

(c) TFIDF

LTC TFC

() () ()

TFIDF

$$P = a/(a + c)$$

()

n_i

N

)

$$R = a/(a + b)$$

()

a_{ik}

k

i

(

k

i

)

f_{ik}

i

k

i

$$F1 = \frac{2 * P * R}{P + R}$$

()

()

($a_{ik} = f_{ik}$)

k

i

TFIDF

()

()

)

($a_{ik} = f_{ik} \times \log(\frac{N}{n_i})$)

n

c_i b_i a_i

[]

$$b = \sum_{i=1}^n b_i \quad a = \sum_{i=1}^n a_i$$

$$c = \sum_{i=1}^n c_i$$

[]

(DF)

:

(MI)

(IG)

[]

SCHI

CHI

(DF)

(DF)

F1

[] F

(R)

(P)

[]

(T_i)

[]

(T)

(a :

d

c_j

$p(c_j | d) = p(c_j) \prod_{k=1}^M p(w_k | c_j)$

(T_i) True False

(R) (P)

Rocchio SVM

F1

$C_{max} = \arg \max_{c_j} p(c_j) \prod_{k=1}^M p(w_k | c_j)$

\ln

$\ln p(c_j | d) = \ln(p(c_j)) + \sum_{k=1}^M \ln p(w_k | c_j)$

SVM Rocchio

Rocchio

[] SVM

()

(DF)

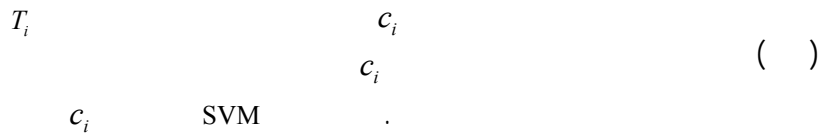
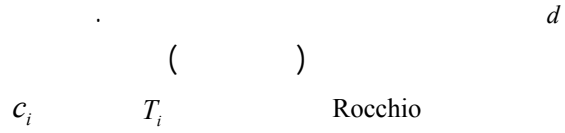
Rocchio

c_j

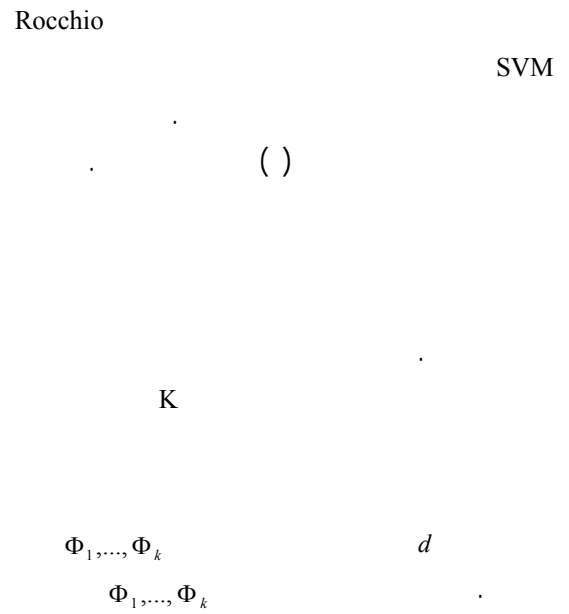
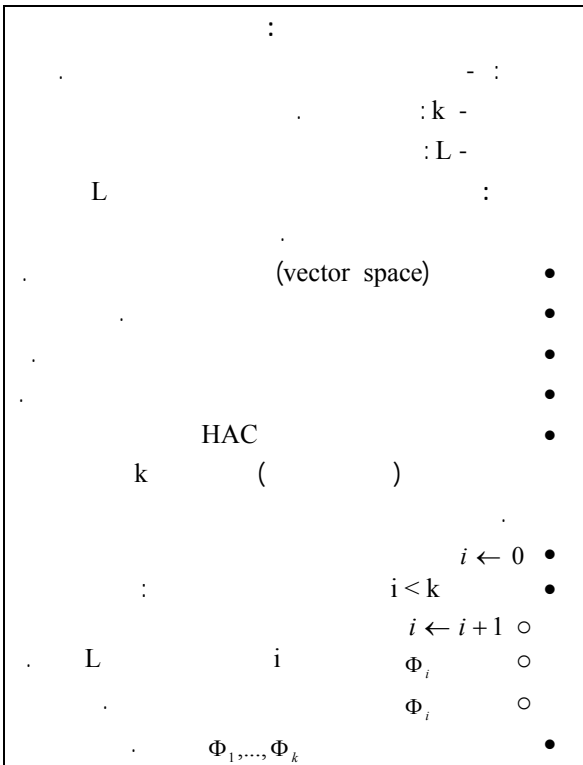
c_j S

		$C_j = \frac{1}{ S } \sum_{i=1}^{ S } doc_i, doc_i \in c_j$	$: [\quad]$	()
SVM				()
SVM	(Rocchio)	c_j	doc_i	
		d		
			$[\quad]$	
		d	d	
		$: ()$	C_{max}	
	$[\quad]$	$C_{max} = \arg \max_{c_j} Sim(C_j, d)$		()
		d	$Sim(C_j, d)$	
		c_j		
	SVM		TFIDF	
		(DF)		
	$[\quad]$ Rocchio			SVM
				SVM
HAC		n		
)	$[\quad]$		
		(
)			
HAC		(
				SVM
		()	$[\quad]$	
			SVM	
			$[\quad]$ SVM_Light	
	$sim(cluster_m, cluster_n) = \min(sim(c_i, c_j))$	()	(
	$c_i \in cluster_m \ \& \ c_j \in cluster_n$)
		()		
				()

$$dissim(c_i, c_j) = 1 - sim(c_i, c_j) \quad ()$$



[]



(f)

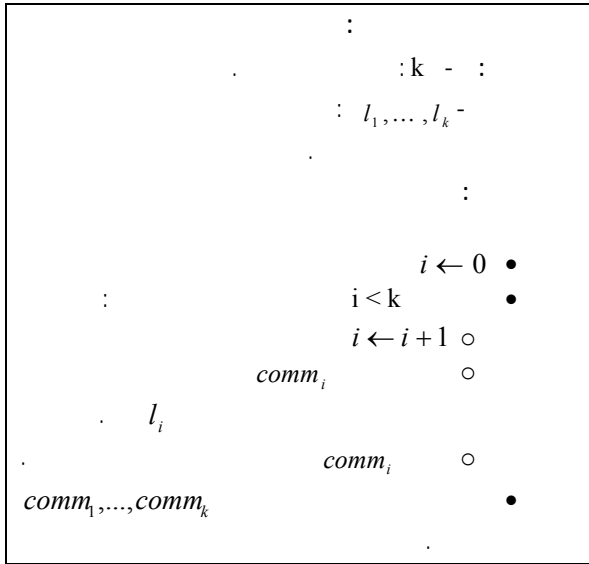
d

[]

SVM

Rocchio

[]



()

... SVM Rocchio

)

(

()

()

)

SVM
(

Rocchio

SVM

)

(

Rocchio

Rocchio

SVM

SVM

-
ModApte

()

[]

SVM

%

%

)

SVM (

()

()

()

Rocchio

() ()

()

Multi Class Approach Max Approach

(Max Approach)

SVM :

SVM

	With Max Approach		With Multi Class Approach	
	R	P	R	P
SVM	0.85	0.95	0.93	0.85

()

SVM :

SVM

	With Max Approach		With Multi Class Approach	
	R	P	R	P
SVM	0.68	0.90	0.81	0.64

Rocchio :

	Use Max Approach		Use Break Even Point Per Class	
	R	P	R	P
Rocchio	0.86	0.95	0.8	0.8
Naïve Bayesian	0.85	0.93	0.87	0.87

ModApte

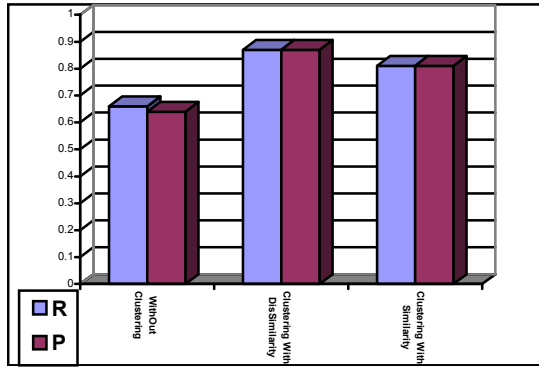
Rocchio :

	Use Max Approach		Use Breakeven Point Per Classes	
	R	P	R	P
Rocchio	0.7	0.86	0.7	0.7
Naïve Bayesian	0.65	0.80	0.66	0.64

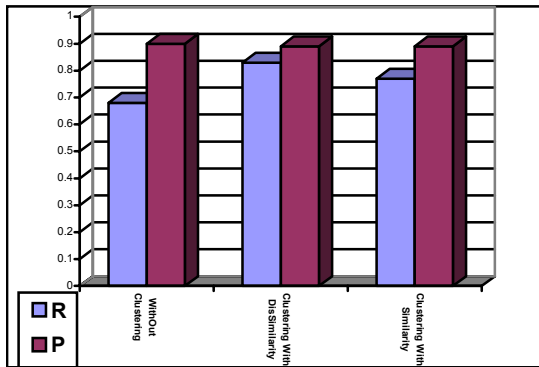
SVM :

() ()

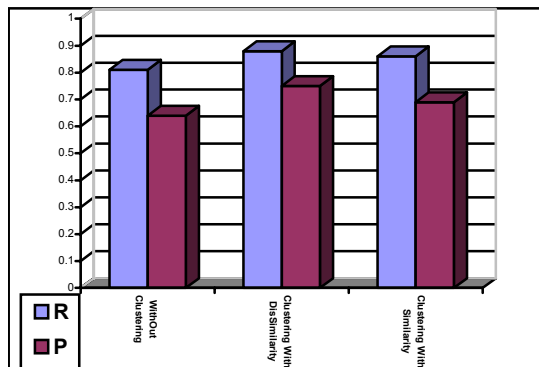
	For first 10 Classes		For All Classes	
	R	P	R	P
SVM	0.96	0.97	0.92	0.97



.) (Breakeven per Classes

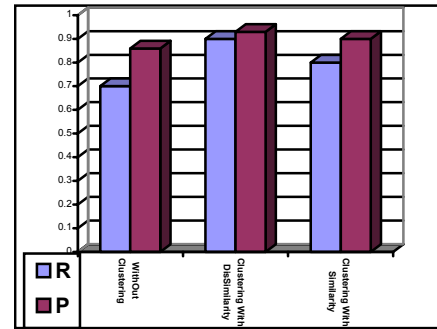


SVM
.) (Max Approach

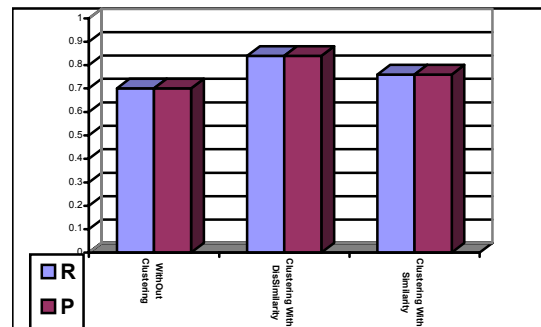


Multi Class SVM
.) (Approach

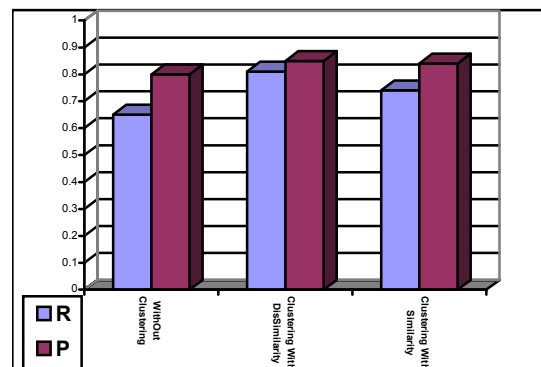
ModApte



Rocchio
.) (Max Approach



Rocchio
.) (Breakeven per Classes



.) (Max Approach

(F1)

SVM ()
SVM

() SVM Rocchio (P R)

() ()
Y

X (P R)
()

() Rocchio ()
SVM () ()
()

SVM Rocchio

() () ()
/

() ()

Rocchio Rocchio Rocchio
()
/

() SVM

SVM F1

SVM SVM)

(
(F1)

)
(Rocchio

SVM

() ()

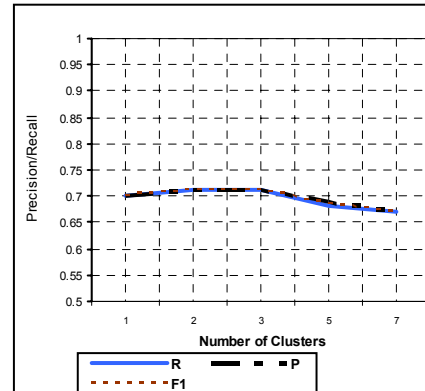
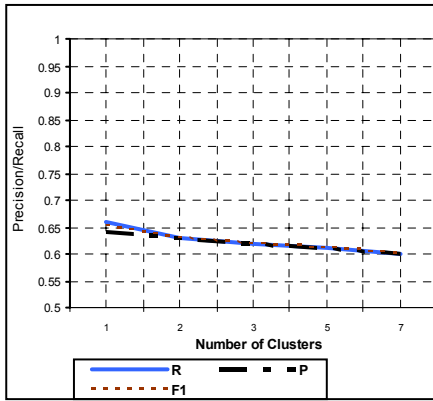
(SVM) (F1)

F1 ModApte

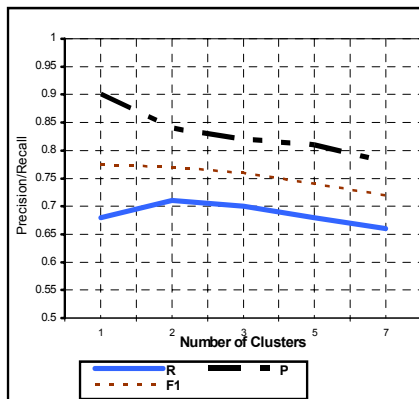
(SVM)

20- []

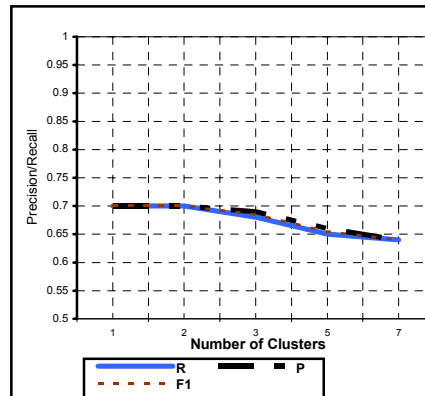
F1



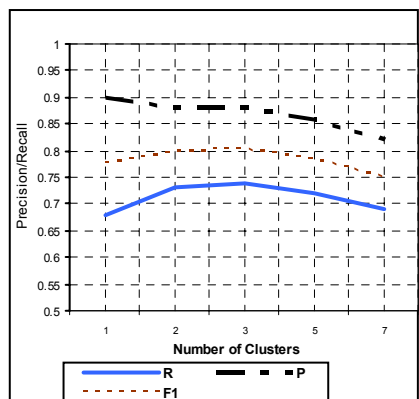
Rocchio



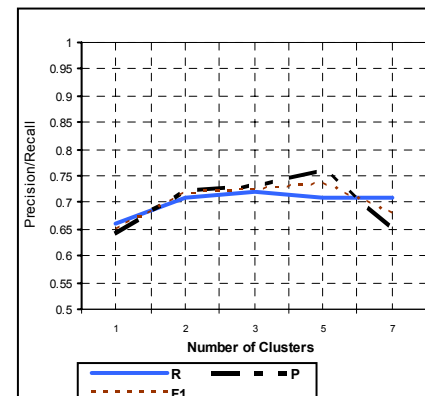
SVM



Rocchio



SVM



()

()

SVM Rocchio

SVM

SVM

() ()

()

/

()

/

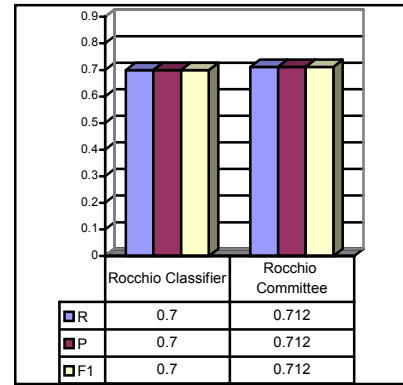
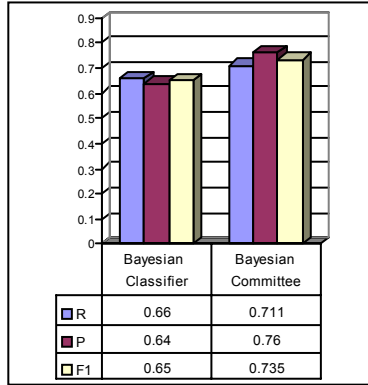
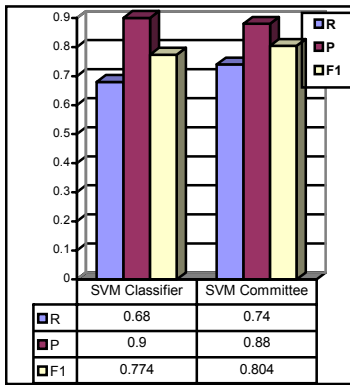
)

()

Rocchio

(

/



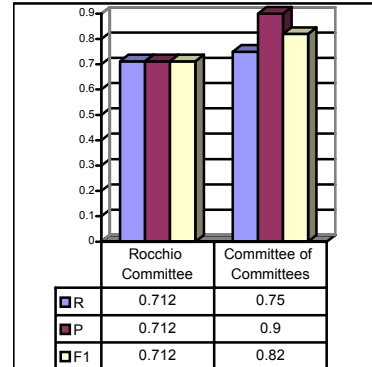
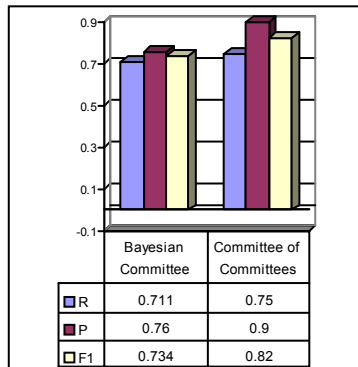
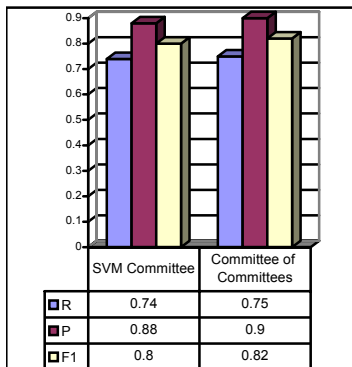
SVM

SVM

Rocchio

Rocchio

()



SVM

Rocchio

"Committee 2"

F1

(Committee 3-1)

ModApte

F1

SVM

[]

(SVM)

20-NewsGroup

F1

F1

SVM

)

(

:

		10 First Classes			All Classes		
		R	P	F1	R	P	F1
Use Max Approach	Rocchio	0.86	0.95	0.90	0.7	0.86	0.77
	Naïve Bayesian	0.85	0.93	0.89	0.65	0.80	0.72
	SVM	0.85	0.95	0.90	0.70	0.90	0.79
Committee Use MV Method		0.87	0.97	0.92	0.70	0.93	0.80

SVM Rocchio

(MV)

() ()

F1

breakeven

Rocchio

SVM

SVM

		First 10 Classes			All Classes		
		R	P	F1	R	P	F1
Use Breakeven Point Per Classes	Rocchio	0.80	0.80	0.80	0.70	0.70	0.70
	Naïve Bayesian	0.87	0.87	0.87	0.66	0.64	0.65
SVM Use Multi Class Approach		0.93	0.85	0.89	0.81	0.65	0.72
Committee Use MV Method		0.90	0.89	0.90	0.70	0.82	0.76

(Max Approach)

F1

()

()

()

ModApte

()

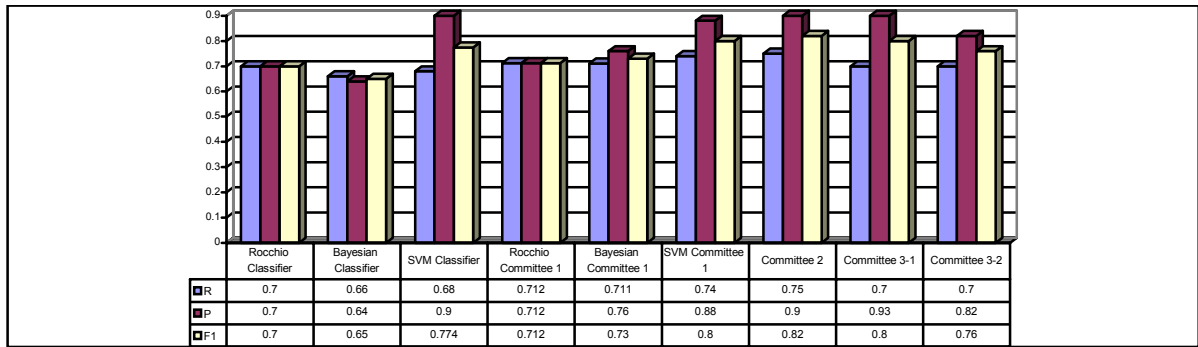
()

()

()

"Committee 3-2"

"Committee 3-1"



.ModeApte

(

Rocchio

SVM / ()

SVM

SVM (

F1

SVM

) [] SVM

(

Rocchio

[]

[] SVM

20-NewsGroup

- SVM

) ModApte

) (/)

-
- (
- []
- () d -)
- (ModApte
- []
- 1 - Giorgetti, D. and Sebastiani, F. (2003). "Multiclass text categorization for automated survey coding." *Proceedings of SAC-03, 18th ACM Symposium on Applied Computing*, Melbourne, US, PP. 798-802.
 - 2 - Breiman, L. (1996). "Bagging predictors." *Machine Learning*, Vol. 24, PP. 123-140.
 - 3 - Debole, F. and Sebastiani, F. (2004). "An analysis of the relative difficulty of Reuters-21578 subsets." *In Proceedings of LREC-04, 4th Int. Conf. on Language Resources and Evaluation*, Lisbon, PT, PP. 971-974.
 - 4 - Lewis, D. D., Yang, Y., Rose, T. and Li, F. (2004). "RCV1: A New Benchmark Collection for Text Categorization Research." *Journal of Machine Learning Research*, Vol. 5, PP. 361-397.
 - 5 - Li, Y. H. and Jain, A. K. (1998). "Classification of text documents." *The Computer Journal*, Vol. 41, No. 8, PP.537-546.
 - 6 - Joachims, T., Freitag, T. and Mitchell, D. (1997). "WebWatcher: A tour guide for the world wide web." *In proceedings of 15th Int. joint Conf. on artificial intelligence*.
 - 7 - Aas, K. and Eikvil, L. (1999). "Text categorization: A survey." *Int. Conf. on machine learning*, PP.128-156.
 - 8 - Sebastiani, F. (2002). "Machine learning in automated text categorization." *Journal of ACM Computing Surveys*, Vol. 34, No. 1, PP. 1-47.
 - 9 - Yang, Y. and Chute, C. G. (1994). "An example-based mapping method for text categorization and retrieval." *ACM Transactions on Information System*, Vol. 12, No. 3, PP.252-277.
 - 10 - Yang, Y. and Liu, X. (1999). "A re-examination of text categorization methods." *In proceedings of SIGIR-99, 22nd ACM Int. Coference on Research and Development in Information Retrieval* (Berkeley, US), PP. 42-49.
 - 11 - Dagan, I., Karov, Y. and Roth, D. (1997). "Mistakedriven learning in text categorization." *In Proceedings of 2nd Conf. on Empirical Methods in Natural Language Processing* (Providence, RI), PP.55-63.
 - 12 - Li, H. and Yamanishi, K. (1999). "Text classification using ESC-based stochastic decision lists." *In Proceedings of 8th ACM Int. Conf. on Information and Knowledge Management*, PP.122-130.
 - 13 - Cohen, W. and Singer, Y. (1999). "Context sensitive learning methods for text categorization." *ACM Trans. Inform. Syst.*, Vol. 17, No. 2, PP.141-173.
 - 14 - Apte, C., Damerau, F. J. and Weiss, S. M. (1999), "Automated learning of decision rules for text categorization." *ACM Trans. On inform. Syst.*, Vol. 12, No. 3, PP.233- 251.
 - 15 - Lewis, D. D. (1998). "Naive(Bayes) at Forty: The independence assumption in information retrival." *ECML*.
-

-
- 16 - Koller, D. and Sahami, M. (1997). "Hierarchically classifying documents using very few words." *Proceedings of the 14th Int. Conf. on Machine Learning (ML)*, Nashville, Tennessee, PP. 170-178, July.
 - 17 - Hull, D. A. (1994). "Improving text retrieval for routing problem using latent semantic indexing." *In Proceedings of 17th ACM Int. Conf. on Research and Development in Information Retrieval*, PP. 282-289.
 - 18 - Jalili, S. and Aghaee A. (2003). *Text Classification by class centroid method*. 11th conference on electrical engineering (ICEE 2003).
 - 19 - Joachims, T. (1998). "Text categorization with support vector machines: Learning with many relevant features." *ECML*, PP.137-142.
 - 20 - Klinkenberg, R. and Joachims, T. (2000). "Detecting concept drift With support vector machines." *In Proceedings of 17th Int. Conf. on Machine Learning*, PP.487-494.
 - 21 - Dumais, S. T., Platt, J., Heckerman, D. and Sahami, M. (1998). "Inductive learning algorithms and representations for text categorization." *In Proceedings of CIKM-98, 7th ACM Int. Conf. on Information and Knowledge Management (Bethesda, MD)*, PP.148-155.
 - 22 - Larkey, L. S. and Croft, W. B. (1996). "Combining classifiers in text categorization." *In Proceedings of SIGIR-96, 19th ACM Int. Conf. on Research and Development in Information Retrieval (Switzerland)*, PP.289-297.
 - 23 - Bi, Y. and Be, D. (2004). *Classification Decision Combination for Text Categorization: An Experimental Study*, DEXA 2004, LNCS 3180, PP. 222 - 231.
 - 24 - Bell, D. A. and Guan, J. W. and Bi, Y. X. (2005). "An evidential approach to classification combination for text categorisation." *Studies in Fuzziness and Soft Computing*, Vol. 185, PP. 13-22.
 - 25 - Bi, Y. and Be, D. (2004). *Combining Multiple Classifiers Using Dempster's Rule of Combination for Text Categorization*, LNAI 3131, PP. 127-138.
 - 26 - Zheng, Z. and Zhou, S. and Zhou, A. (2004). *Sequential Classifiers Combination for Text Categorization: An Experimental Study*, WAIM 2004, LNCS 3129, PP. 509–518.
 - 27 - Nardiello, P., Sebastiani, F. and Sperduti, A. (2003). "Discretizing continuous attributes in AdaBoost for text categorization." *Proceedings of ECIR-03, 25th European Conf. on Information Retrieval*, PP. 320—334.
 - 28 - Steinbach, M., Karypis, G. and Kumar, V. (2000). "A comparison of document clustering techniques." *Proc. TextMining Workshop, KDD*.
 - 29 - Tsirikia, T. and Lalmas, M. (2001). "Merging techniques for performing data fusion on the web." *Proceedings of the 10th Int. Conf. on information and Knowledge management*, Atlanta, Georgia, USA, October.
 - 30 - Bordogna, G. (2002). "Soft fusion of information accesses." *Proceedings of IEEE Int. Conf. on Fuzzy Systems*, Vol 2, PP.1466-1471.
 - 31 - Beil, F. Ester, M. and Xu, X. (2002). "Frequent term-based text clustering." *Proc. 8th Int. Conf. on Knowledge Discovery and Data Mining (KDD'02)*.
 - 32 - Zamir, O., Etzioni, O., Madani, O., Karp, R. M. (1997). "Fast and intuitive clustering of web documents." *KDD '97*, PP. 287-290.
 - 33 - Aggarwal, C. C., Gates, S. C. and Yu, P. S. (1999). "On the merits of building categorization systems by supervised clustering." *Proc.of the 5th ACM Int. onf. on Knowledge Discovery and Data Mining*, PP. 352 – 356.
 - 34 - Cutting, D. R., Karger, D. R., Pedersen, J. O. and Tukey, J. W. (1992). "Scatter/Gather: A Cluster-based Approach to Browsing Large Document Collections." *SIGIR '92*, PP. 318 – 329.
 - 35 - Debole, F. and Sebastiani, F. (2003). "Supervised term weighting for automated text categorization." *Proceedings of SAC-03, 18th ACM Symposium on Applied Computing*, Melbourne, US, PP. 784-788.
-

-
- 36 - Galavotti, L., Sebastiani, F. and Simi, M. (2000). "Experiments on the use of feature selection and negative evidence in automated text categorization." *In Proceedings of 4th European Conf. on Research and Advanced Technology for Digital Libraries*, PP. 59-68.
- 37 - Yang, Y. and Pedersen, J. O. (1997). "A comparative study on feature selection in text categorization." *In Proce. of 14th Int. Conf. on Machine Learning*, PP.412-420.
- 38 - Jalili, S. and Bitarafan, M. (2006). "Increasing performace of text classification based on improving feature selection." *Journal of faculty of engineering (Tehran university), Special Issue:Electrical Engineering*, Vol 40, No.3.
- 39 - Sadri, A.A. (2005). *Text Classification Based on Committee and Feature Selection*. Thesis of Master of Science (M.S.) in Computer Engineering (Software)" , Tarbiat Modares University, Faculty of Engineering.
- 40 - Sadri, A.A. and Jalili, S. (2006). *Quality Enhancement in Text Classification with classification committees*. 14th conference on electrical engineering(ICEE 2006) , AmirKabir University.
- 41 - Moens, M. and Dumortier, J. (2000). "Text categorization: the assignment of subject descriptors to magazine articles." *Information Processing and Management*, Vol. 36, Elsevier Science, PP.841-861.
- 42 - El-Hamdouchi, A. and Willet, P. (1989). "Comparison of Hierarchic Agglomerative Clustering Methods for Document Retrieval." *The Computer Journal*, Vol. 32, No.3.
- 43 - Cauwenberghs, G. and Poggio, T. (2000). "Incremental and Decremental support vector machine learning." *Aadvances In Neural Information Processing*, Vol. 13.

- | | |
|----------------------------------|--|
| 1 - Information Retrieval | 2 - k Nearest Neighbor |
| 3 - Support Vector Machines | 4 - Pair Wise Coupling |
| 5 - Majority voting | 6 - Weighted linear Combination |
| 7 - Dynamic classifier selection | 8 - Adaptive classifier combination |
| 9 - Local accuracy | 10 - Decision Fusion |
| 11 - Text Mining | 12 - Hierarchic Agglomerative Clustering |
| 13 - Stop words | 14 - Document Frequency |
| 15 - Information Gain | 16 - Mutual Information |
| 17 - Correlation Coefficient | 18 - Precision |
| 19 - Recall | 20 - Breakeven |
| 21 - Micro Average | 22 - Macro Average |
| 23 - Monotone | 24 - Discriminative |
| 25 - Supprot Vectors | 26 - Incremental |