

{hjshin, btzhang}@scai.snu.ac.kr, ytkim@comp.snu.ac.kr

Feature Selection with Non-linear PCA in Text Categorization

Hyung Joo Shin, Byoung-Tak Zhang, Yung Taek Kim

School of Computer Engineering and Science, Seoul National University

PCA (Document Frequency) (PCA) (nonlinear PCA) locally linear PCA kernel

1.

(category) (regression model), (Naive Bayesian probabilistic model, NB), (Decision tree model), (Inductive rule learning model), (Neural Networks, NNets), Support Vector Machines(SVM), k-Nearest Neighbor(kNN), Linear Least Square Fit(LLSF) [1]

Hebbian networks, associated multi-layer perceptrons, principal curves, locally linear PCA, kernel PCA DF locally linear PCA [5] kernel PCA [4]

feature selection, (Automatic feature selection, Lewis & Ringuette (Information Gain, IG), Wiener (Mutual Information, MI) χ -square(CHI), Yang Schutze LLSF (linear PCA), Yang & Wilbur kNN, Lang Minimum Description Length(MDL) DF(documet frequency) IG, CHI [2]

n-ary kNN [10] (PCA) (eigenvectors) (correlation matrix) (iteratively) (principal components) (noise) X

Reuters 21578 SVM, LLSF, kNN NB NNets [1], LLSF[3] (multivariate) 가

$$X \cdot V^k \quad (1)$$

(eigenvalue) , k k

(linear PCA) 가

(nonlinear PCA)
 , principal curves [7]
 [8] (discrete) principal curves SOM
 , [9] 4-layer MLP(Multi Layer Perceptron)
 principal curves [4]
 kernel PCA 가, [5] locally linear PCA 가
 Kernel PCA [4]
 (feature space) (mapping)

$$\Phi(\mathbf{x}_i) = \sum_{j=1}^M K(\mathbf{x}_i, \mathbf{x}_j) \mathbf{v}_j$$

(2) feature space
 kernel
 kernel PCA

$$\mathbf{V}^k \cdot \Phi(\mathbf{x}_i) = \sum_{j=1}^M K(\mathbf{x}_i, \mathbf{x}_j) \mathbf{v}_j$$

(2) feature space
 , α kernel
 kernel polynomial kernel, radial-basis kernel, sigmoid kernel
 가 M kernel

$M \times M$
 Locally linear PCA [5] clustering
 cluster [5] clustering
 Vector Quantization(VQ)
 Organizing Maps(SOM) Self-

3.
 3.1.
 Reuters-21578 Reuters-21578
 5 가
 (TOPICS) 135 가 10
 (positive data)
 (negative data)가 n-ary
 negative
 negative 9
 가

Corn	133	36
Ship	180	80
Wheat	185	54
Interest	275	97
Trade	304	103
Crude	334	156
Grain	370	125
Money-fx	428	130
Acq	1483	640
Earn	2706	1043
	6397	2464

Porter algorithm
 stemming , 524 stop word ,
 8754 feature() , tf/idf(term

frequency/inverse term frequency)
 8754 6397¹⁾
 weight 가 0 7218
 6397×7218
 2464×7218

3.2. Feature Selection

5000, 4000, 3000,
 2000, 1000, 500 가

3.2.1. Document Frequency (DF)
 weight 가 0 feature
 feature

3.2.2.
 (1)

3.2.3. Locally linear PCA
 Self-Organizing Map Program Package Ver. 3.1²⁾ 가 feature
 100 cluster clustering cluster feature
 7218 feature
 100 clustering cluster feature
 가 2 cluster feature 5000 cluster
 5000 cluster 가 feature cluster
 100 feature 가 SOM 가 cluster
 50 cluster cluster 가
 feature 가 cluster

2. Feature		clustering	SOM	map					
47	17	11	9	98	3	3	77	6	158
1	9	67	3	11	97	2	18	8	15
97	15	12	4	50	29	27	26	18	118
11	26	13	51	7	24	84	33	5	20
70	5	64	3	7	105	18	16	49	58
44	54	9	86	7	8	109	37	12	86
236	76	130	25	12	99	22	14	146	88
124	1065	76	93	35	15	104	58	98	257
465	109	109	76	12	63	40	14	158	88
201	288	165	64	83	44	93	66	22	109

3.2.4. Kernel PCA
 polynomial kernel radial-basis kernel
 polynomial kernel d=4 , radial-basis kernel
 $\sigma=1$ kernel
 가
 kernel

1) McCallum Bowlibrary
<http://www.cs.cmu.edu/~mccallum/bow>

2) Teuvo Kohonen 1995
 anonymous ftp site cochlea.hut.fi /pub/som_pak

3.3. n-ary 500

kNN [6] Similarity measure cosine feature
 measure neighborhood size k 100
 2464 가 , 3 DF
 가 가 feature 가 2000

4. 가
 3 4

4.1. kernel PCA
 kNN randomness 가
 81.73% kernel 가
 1 DF
 가 3 PCA locally linear PCA DF, kernel
 Linear PCA (:%)
 3. DF Linear PCA (:%)

DF	Linear PCA (%)
5000	80.03
4000	79.98
3000	79.12
2000	78.84
1000	76.19
500	66.24

4. Kernel PCA Locally linear PCA (:%)

	Kernel PCA (polynomial)	Kernel PCA (radial basis)	Locally linear PCA
5000	80.98	81.25	82.28
4000	81.73	82.71	82.65
3000	82.22	84.13	83.18
2000	82.50	82.08	84.33
1000	80.18	80.22	81.17
500	76.97	78.45	76.84

5. (:%)

		(:%)
Corn	133	76.04
Wheat	180	77.05
Ship	185	78.98
Interest	275	81.29
Trade	304	80.57
Grain	334	82.55
Money-fx	370	81.59
Crude	428	82.47
Acq	1483	85.15
Earn	2706	86.39
	6397	84.33

4 linear PCA kernel PCA locally
 5 locally linear PCA
 2000

4. 3, 4 weight 0

PCA locally linear PCA DF, kernel
 Locally linear PCA
 가
 4 가
 가
 가
 kernel PCA locally linear
 PCA 가 가

(c1-98-006800)

[1] Yiming Yang and Xin Liu, "A re-examination of text categorization methods," in *Proceedings of the 22th Ann Int ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR '99)*, 1999.

[2] Yiming Yang and Jan O. Pedersen, "A comparative study on feature selection in text categorization," in *Proceedings of the 14th Int Conference on Machine Learning (ICML '97)*, pp. 412-420, 1997.

[3] Yiming, Yang, "Noise reduction in a statistical approach to text categorization," in *Proceedings of the 18th Annual Int ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR'95)*, pp. 256-263, 1995.

[4] Bernhard Schölkopf and Alexander Smola, "Nonlinear component analysis as a kernel eigenvalue problem," *Neural Computation*, vol. 10, no. 5, pp. 1299-1319, 1998.

[5] N. Kambhatla and Todd K. Leen, "Dimension reduction by local principal component analysis," *Neural Computation*, vol.9, no.7, pp. 1493-1516, 1997.

[6] Erkki Oja, "Principal components, minor components, and linear neural networks," *Neural Networks*, vol. 5, pp. 927-935, 1992.

[7] T. Hastie and W. Stuetzle, "Principal curves," *Journal of the American Statistical Association*, vol. 84, no. 406, pp. 502-516, 1989.

[8] Helge Ritter, Thomas Martinetz, and Klaus Schuler, *Neural Computation and Self-Organizing Maps, An Introduction.*, Addison-Wesley, Reading, Massachusetts, 1992.

[9] Mark A. Kramer, "Nonlinear principal component analysis using autoassociative neural networks," *American Institute of Chemical Engineers (AIChE)*, vol. 37, no. 2, pp. 233-243, 1991.

[10] B. Masand, G. Linoff, and D. Walts, "Classifying news stories using memory based reasoning," in *Proceedings of the 15th Annual Int ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR '92)*, pp. 59-64, 1992.