

應用非參數型聚類法於聚類分析

范仕遠、陳木松

E-mail: 9304702@mail.dyu.edu.tw

摘要

聚類分析不僅是擷取取樣資料結構與特性的重要工具，也可應用於萃取複雜系統的輸入與輸出的關係。聚類分析已廣泛應用於許多工程與科學的領域，如圖像辨識與分類、機器學習、電腦視覺、與系統建模等。聚類分析試圖將取樣資料分割成若干個群組或聚類，使得同聚類的資料距離總和愈小(或相似度愈大)，而且不同聚類的距離愈大。然而在僅有取樣資料而無任何先前知識的情況下，估測正確的聚類數目是非常困難的工作。本研究應用資料的距離資訊、密度、熵值、格子聚類法、及共變數矩陣等相關資訊，提出橢圓密度熵值聚類演算法，以僅有取樣資料而無任何先前知識(prior knowledge)的情況下，減少參數的敏感度即可以自動決定各個聚類群組，以解決對於取樣資料結構與各種類型的聚類演算法的認知不足，導致於設定無意義或錯誤的參數，以致於對最後分析結果產生誤判的情況。本論文以高斯分佈的資料結構分別對於(1).聚類間隔較大、(2).聚類資料部份重疊，與(3).含有雜訊等情況，以驗證本研究方法的可行性。藉由與FCM(Fuzzy C-means)的比較，實驗結果展示本法的可靠性與穩定性。

關鍵詞：資料探勘、知識發掘、聚類分析、熵值、橢圓密度熵值聚類演算法

目錄

第一章 緒論	1	1.1 研究背景與動機	1	1.2 研究方法	3	1.3 本文架構	4
第二章 聚類分析相關理論	6	2.1 參數型聚類法	7	2.1.1 高斯混合模式	7	2.1.2 切割式聚類法	8
	7	2.2 非參數型聚類法	11	2.2.1 階層式聚類法	11	2.2.2 密度基礎聚類法	14
	11	2.2.3 格子基礎聚類法	14	2.3 相關研究	16	2.3.1 以熵值為主的聚類個數預測法則	18
第三章 橢圓密度熵值聚類演算法	23	3.1 橢圓密度熵值聚類演算法概述	23	3.2 橢圓密度熵值聚類演算法的架構	24	3.2.1 有效聚類範圍	24
	23	3.2.2 聚類中心的選擇	25	3.2.3 以密度為主的聚類門檻	31	第四章 實驗模擬結果與分析	38
	24	4.1 聚類難易度與結果分析之評估	38	4.2 不同高斯分佈的資料結構	41	4.3 高斯分佈聚類資料部份重疊	47
	31	4.4 以含有雜訊的高斯分佈資料	52	第五章 結論與展望	57	參考文獻	59
	38		52		57		59

參考文獻

- [1] 范仕遠、陳木松, “以熵值為主的模糊聚類分析”, 2003 人工智慧、模糊系統及灰色系統聯合研討會.
- [2] N.Zahid, M. Limouri and A. Essaid, “A new cluster-validity for fuzzy clustering,” Pattern Recognition, vol.32 pp.1089-1097, 1999.
- [3] C.E. Shannon, “A mathematical theory of communication,” Bell Syst. Tech. J., vol. 27, pp. 379-423, 1948.
- [4] I. Csiszk and J. Korner. Information Theory: Coding Theorems for Discrete Memoryless System. Academic Press, 1981.
- [5] 范仕遠、陳木松, “以熵值為主的聚類個數預測法則”, 彭雲嘉 2003 年研發成果聯合發表會 [6] Glenn Fung “A Comprehensive Overview of Basic Clustering Algorithms,” May. 2001.
- [7] C. J. C. Burges. A tutorial on support vector machines for pattern recognition. Data Mining and Knowledge Discovery, 2(2):121-167, [8] McLachlan G. J., Basford K. E. Mixture Models: Inference and Applications to Clustering. Marcel Dekker, 1988.
- [9] Pan, W, Lin, J, & Le, C: A mixture model approach to detecting differentially expressed genes with microarray data. Technical Report 2001-011, Division of Biostatistics.
- [10] S. M. Zabin and G. A. Wright, “Non-parametric density estimation and detection in impulsive interference channels-Part II: Detectors,” IEEE Trans. Commun., vol. 42, pp. 1698 – 1711, 1994.
- [11] Ferguson, T.S., “A Bayesian Analysis of Some Nonparametric Problems,” The Annals of Statistics, 1, 209-230, 1973.

- [12] C. Fraley and E. Raftery, "How many clusters? which clustering method?, answers via model-based cluster analysis." Technical Report 329, Dept. of Statistics. University of Washington, Seattle, 1998.
- [13] C. Bishop, Neural networks for pattern recognition. Oxford, U.K. Oxford Univ. Press, 1995.
- [14] A. P. Dempster, N. M. Laird and D. B. Rubin, "Maximum likelihood from incomplete data via the EM algorithm." Journal of the Royal Statistical Society Series B vol.39:1-38, 1977.
- [15] Geoffrey J. McLachlan and Kaye E. Basford. Mixture Models: Inference and Applications to Clustering. Marcel Dekker, 1988.
- [16] S.J. Roberts, Dirk Husmeier, lead Rezek, and William Penny, "Bayesian Approaches to Gaussian Mixture Modeling," IEEE Trans. on PAMI, vol. 20, no. 11, Nov. 1998.
- [17] J. C. Bezdek. Pattern Recognition with Fuzzy Objective Function Algorithms. New York, 1981.
- [18] J.S. Jang, C.T. Sun, and E. Mizutani, Neuro-Fuzzy and Soft Computing, Prentice Hall, 1997.
- [19] D.E. Gustafson and W. Kessel, "Fuzzy clustering with a fuzzy covariance matrix," Proc. IEEE-CDC 2, 761-766, 1979.
- [20] Tran, D.; Wagner, M "Fuzzy entropy clustering," Fuzzy Systems, The Ninth IEEE International Conference on, Vol.1, Pages:152 - 157 May 2000.
- [21] L. Kaufman and P. J. Rousseeuw, Finding Groups in Data: an Introduction to Cluster Analysis. John Wiley & Sons, 1990.
- [22] Han J. and Kamber M., Data Mining: Concepts and Techniques, Morgan Kaufmann, 2000.
- [23] M. Ester, H. Kriegel, J. Sander, and X. Xu, "A density-based algorithm for discovering clusters in large spatial databases with noise," in Proc. of the 2nd Int'l Conf. on Knowledge Discovery in Databases, Menlo Park, CA., pp. 226-231, 1996.
- [24] W. Wang, J. Yang, and R. Muntz, "STING: A Statistical Information Grid Approach to Spatial Data Mining," Technical Report CSD-97006, Computer Science Department, University of California, Los Angeles, Feb. 1997.
- [25] G. Sheikholeslami, S. Chatterjee, and A. Zhang, "WaveCluster: A multi-resolution clustering approach for very large spatial databases," In Proc. 1998 Int. Conf. Very Large Databases (VLDB '98), pp. 428-439, New York, Aug. 1998.
- [26] R. Agrawal, J. Gehrke, D. Gunopulos, and P. Raghavan, "Automatic subspace clustering of high dimensional data for data mining applications," SIGMOD 1998.
- [27] 陳奕學,空間資料叢集演算法之設計, 義守大學資訊工程學系, 碩士論文, 2002.