

國立交通大學

資訊科學系

碩士論文

一個針對二維資料的機率式重新分配模擬退火分群法的介面

A GUI of Simulated-Annealing K-means Clustering with
Probabilistic Redistribution for 2D Data

研究生：卓明宏

指導教授：李嘉晃 教授

中華民國九十三年七月

一個針對二維資料的機率式重新分配模擬退火分群法的介面
A GUI of Simulated-Annealing K-means Clustering with
Probabilistic Redistribution for 2D Data

研究生：卓明宏

Student : Ming-Hung Cho

指導教授：李嘉晃

Advisor : Chia-Hoang Lee

國立交通大學
資訊科學系
碩士論文



Submitted to Institute of Computer and Information Science
College of Electrical Engineering and Computer Science

National Chiao Tung University

in partial Fulfillment of the Requirements

for the Degree of

Master

in

Computer and Information Science

July 2004

Hsinchu, Taiwan, Republic of China

中華民國九十三年七月

一個針對二維資料的機率式重新分配模擬退火分群法的介面

研究生：卓明宏

指導教授：李嘉晃 教授

國立交通大學電機資訊學院 資訊科學系碩士班

中文摘要

分群法是一個被廣泛使用的自動化資料分類技術，常用於機器學習或資料探勘的先置處理步驟前置。其中模擬退火法是一種模擬多粒子物理系統由較高能量熱平衡狀態降溫到較低能量熱平衡狀態過程以取得最少 TSSE 的分群法。SAKM-clustering 整合了模擬退火法取得系統最小能量及 K-means 快速搜尋的能力。它適用於搜尋多維特徵空間中適當分群並使得分群結果在相似度測度上最佳化。在本論文中我們針對二維資料實作了一個介面以觀察 SAKM-clustering 及 K-means 的分群過程效能比較。

A GUI of Simulated-Annealing K-means Clustering with Probabilistic Redistribution for 2D Data

Student : Ming-Hung Cho

Advisor : Prof. Chia-Hoang Lee

Department of Computer and Information Science

National Chiao Tung University

ABSTRACT

Clustering is an extensively used technique for automatic data classification, such as in the preprocessing of machine learning and data mining. Simulated-annealing is a clustering technique, which obtains the minimum TSSE of a group of data by simulating the cooling down process of a many-particle physical system from a state in thermal equilibrium with higher energy to another state in thermal equilibrium with lower energy. The SAKM-Clustering integrates the power of simulated-annealing for obtaining minimum energy configuration and the searching capability of K-means algorithm. It is used to search proper clusters in multidimensional feature space such that a similarity metric of the resulting clusters is optimized. In this thesis we implement a GUI to observe the clustering processes and to compare the performances of the SAKM-clustering as well as the K-means clustering.

誌謝

本論文得以順利完成，首先要感謝的是指導教授李嘉晃老師，他給予我不斷的機會去嘗試，給予我充分的空間發展，並且給予我最大的耐心與包容。要特別謝謝老師讓我能有機會輔修電子物理系的課程，完成多年來追尋基礎科學課程的願望。也感謝口試委員蘇豐文教授與楊武教授的指正建議，讓我獲益良多。

要謝謝國立交通大學的師長們及許多單位的長官與工作人員，他們提供了一個高度專業與安全舒適的環境讓我可以安心於課業，並且可以追尋自己所喜愛的知識。要謝謝資訊科學系及電子物理系的助理小姐們熱心的幫助與關心。要謝謝電子物理系授課老師們的教導與幫助，讓我能有機會一窺這個神秘殿堂的奧妙與神奇之處。

還要感謝實驗室所有學長，同學與學弟們的幫助，充實了我的實驗室生活。尤其要謝謝湯隆宜，林駿豪，王明德，高鳴遠，黃明超，林建良幾位，謝謝大家幫我度過許多忙亂的階段。一路走來曾幫助我的老師，助教，學長姊及交大清大的同學們，我也在此表達對各位誠摯的謝意。

最後我要感謝我的家人，因為他們多年來的支持與鼓勵是我努力的最大動力。謹以此論文獻給我最親愛的家人。

目錄

中文摘要	i
英文摘要	ii
誌謝	iii
目錄	iv
圖目錄	v
第一章 緒論	1
1.1 簡介	1
1.2 分割式分群法	2
第二章 背景知識	4
2.1 McQueen' s K-means Algorithm	4
2.2 Forgy' s K-means Algorithm	6
2.3 Hierarchical Methods	8
2.4 Peak-Climbing Clustering	11
2.5 Fuzzy K-means Clustering	12
2.6 Simulated Annealing	14
2.7 Deterministic Annealing	15
第三章 機率式重分配的模擬退火 K-means 演算法	17
3.1 簡介	17
3.2 SAKM 所使用的 K-means 演算法	17
3.3 Metropolis 演算法	18
3.4 SAKM-clustering 中的重分配機制	18
3.5 SAKM-clustering 中的 temperature schedule	19
3.6 SAKM-clustering 演算法	19
第四章 系統實作	20
4.1 變數與參數設定	20
4.2 操作介面與流程	20
第五章 SAKM 與 K-means 的效能比較	24
5.1 點數固定時的比較	24
5.2 分群數固定時的比較	26
5.3 SAKM 中點數固定時分群數與溫度的關係	28
5.4 SAKM 中分群數固定時點數與溫度的關係	29
5.5 SAKM 中的退化現象	29
參考文獻	31

圖目錄

圖 2.1	Centroid-Linkage Method	9
圖 2.2	Complete-Linkage Method	10
圖 2.3	Single-Linkage Method	10
圖 2.4	Peak-climbing	12
圖 4.1	操作介面	21
圖 4.2	資料與系統參數設定	21
圖 4.3	資料初始化	22
圖 4.4	執行結果	22
圖 5.1	資料點為 300 點時 SAKM 及 K-means 的 TSSE 比較	24
圖 5.2	資料點為 3000 點時 SAKM 及 K-means 的 TSSE 比較	25
圖 5.3	資料點為 300 點時 SAKM 及 K-means 的執行次數比較	25
圖 5.4	資料點為 3000 點時 SAKM 及 K-means 的執行次數比較	26
圖 5.5	分群數為 5 群時 SAKM 及 K-means 的 TSSE 比較	26
圖 5.6	分群數為 20 群時 SAKM 及 K-means 的 TSSE 比較	27
圖 5.7	分群數為 5 群時 SAKM 及 K-means 的執行次數比較	27
圖 5.8	分群數為 20 群時 SAKM 及 K-means 的執行次數比較	28
圖 5.9	SAKM 中點數固定時分群數與溫度的關係	28
圖 5.10	SAKM 中分群數固定時點數與溫度的關係	29
圖 5.11	SAKM 中的退化現象	30