

行政院國家科學委員會專題研究計畫 成果報告

以階層式叢集研究多重外傷

計畫類別：個別型計畫

計畫編號：NSC91-2213-E-038-004-

執行期間：91年08月01日至92年07月31日

執行單位：臺北醫學大學醫學資訊研究所

計畫主持人：蔣以仁

報告類型：精簡報告

處理方式：本計畫涉及專利或其他智慧財產權，1年後可公開查詢

中 華 民 國 92 年 12 月 26 日

以階層式叢集分析研究醫學資料庫之多重外傷

Hierarchical Clustering of Multiple Traumas in Medical Database

蔣以仁, 劉致和, 陳瑞杰, 李友專, 翁昭旻

Chiang, I-Jen^{ab}; Liu, Charles C.H.^{abc}; Chen, Ray-Jade^{ad}; Li, Yu-Chuan^a; Wong, Chau-Min^b

^a 台北醫學大學醫學資訊研究所 ^b 台灣大學醫學工程學研究所

^c 國泰醫學中心外科 ^d 長庚醫學中心外科及外傷科

主持人: 蔣以仁

計畫編號: NSC91-2213-E-038-004

一、中文摘要

醫學資料探勘是 2001 年台灣地區健保局釋出健保資料庫 (NHIRD) 供研究使用後的熱門題目。本團隊以多重外傷為例, 探討階層式叢集分析在此種大型資料庫上協助醫學研究及醫院管理之重要性與適用性。

階層式叢集分析具有顯現叢集過程的優點, 提供專家在互動環境下依研究目的來探討資料結構。我們以 Matlab 實作與資料庫連結, 產生資料倉儲的線上分析 (OLAP) 所需之維度資料表, 供進一步資料分析。對於健保資料庫來源龐雜的主副 ICD 診斷資訊的整理, 我們引入臨床分類系統 (CCS) 作資料的初步整理。

我們在兩個主題上介紹在醫學研究及發現的細節 (1) 以叢集排序後之樹狀圖展現台灣地區多重外傷之分布情況, 可提供政策制定與監測之用; (2) 以叢集分析各級醫院在燙傷住院時的不同治療類型, 而分群可協助發掘同質性的病人族群, 做更精細的統計分析, 或用以發現異常的病例或醫院。

關鍵詞: 資料探勘, 叢集分析, 醫學資料庫, 申報資料庫, 全民健保研究資料庫, 外傷, 燙傷

ABSTRACT

Medical data mining becomes more mandatory after the release of national healthcare insurance research database (NHIRD) in 2001. In multiple trauma injuries, our research demonstrated the feasibility and benefits of the hierarchical clustering methodology in medical research and in hospital administration.

The process of clustering could be

explicit shown in hierarchical cluster analysis. The domain experts could elucidate the structure of data in interactive setting. We implement in Matlab the realtime linkage with the data sources, and propose the methodology of dimensional table generation for further online analytical processing (OLAP) in a data warehouse. For the widely varied coding policy in the International Classification of Disease (ICD) diagnosis fields, we introduced Clinical Classification System (CCS) for more medically relevant data preparation.

The medical findings were illustrated in details in two topics – (1) Visualization of the distribution and correlation of multiple trauma in Taiwan, which could provide abundant information for research needs from clinical and administrative fields; (2) To find the major treatment patterns of hospitalized burn patients in various hospitals, to facilitate sophisticated statistical analysis of more homogeneous patient populations after clustering, or as the surveillance of normal pattern and to find patient or hospital/doctor outliers.

Keywords: Data mining, Cluster analysis, Medical database, Claim database, National Healthcare Insurance Research Database (NHIRD), Trauma, Burn injury.

BACKGROUND:

The release of National Health Insurance Research Database (NHIRD) in 2001 brought a new challenge to the medical informatics society in Taiwan [11]. In addition to its huge size and the need of data warehouse for management, the issue of data quality and complexity, as mixture of various

clinical severity and indications from heterogeneous sources, rendered more difficulties on the researchers from clinical medicine or hospital administration.

To find the risk factors contributing to the clinical outcome is always one of the major aims of clinical research. The interplay of co-morbidities and complications is especially important in the retrospective insurance claim database [7]. In our preliminary experience in NHRID, to add more clinical details, and to focus on more homogeneous patient subpopulations was the key factor of yielding data mining in the heterogeneous medical datasets [10][13].

Hierarchical clustering (HC) is recently revisited vigorously for phylogenomics study and the genetic expression pattern of microarray [3][5][8]. We use the agglomerative algorithm (HAC) for the same reason as in bioinformatics, in order to allow intervention by domain knowledge in data rescaling and in the intermediate stages of clustering, to explore the unknown pattern of diagnosis and treatment pattern in the dataset of 96% coverage of Taiwan population.

MATERIALS AND METHODS:

Data preparation and warehousing

The trauma subsets of NHIRD from 1997 to 2001 were pooled into a Oracle 8i data warehouse. Another copy was managed in a Microsoft SQL 2000 for comparison. About 20G of disk space were required for each copy, and the OLAP manipulation might need additional 10 to 40 G.

The hospitalized trauma patients with more than two diagnoses were extracted, totally 1,319,176 cases.



Fig.1. OLAP display of hierarchical ICD categories of disease and statistics of subgroup

Table 2-1. Serial table of cluster information, at 200-cluster-size/level, only the clusters with case number larger than 3 were shown. (Data for 3D clustering in DCSLET section 4.2.6).

c00	c01	c02	c03	c04	c05	c06	c07	c08	c09	c10	c11	c12	c13	c14	c15	c16	c17	c18	c19	CaseNo	ICD	DE	STSG	NT
5	3	1	11	30	7	5	11	51	19	60	23	9,097	1,293	0	0									
5	3	1	11	30	7	5	11	51	19	60	27	6,672	932	0	0									
5	3	1	11	30	7	5	12	52	38	76	9	4,990	4,580	0	0									
5	3	1	11	30	7	5	12	52	42	15	6	4,737	8,830	0	0									
5	3	1	11	30	7	5	12	52	42	16	9	5,009	12,500	0	0									
5	3	1	11	30	7	6	45	92	1	47	12	11,638	469	0	0									
5	3	1	11	30	7	29	61	40	16	59	209	305	1,850	0	0									
5	3	1	11	30	7	29	61	40	17	19	47	1,677	4,373	0	0									
8	3	1	11	30	7	29	61	40	17	20	16	3,001	6,740	0	0									
5	3	1	11	30	7	29	61	40	17	36	34	274	6,990	0	0									
5	3	1	11	30	7	29	61	40	18	21	47	4,273	674	0	0									
8	3	1	11	30	7	29	61	40	18	22	173	2,429	762	0	0									
5	3	1	11	30	7	33	64	38	76	17	19	49	32,247	0	0									
5	3	1	11	30	13	42	71	15	4	50	6	7,239	36,299	0	0									
5	3	1	11	30	13	42	71	15	5	51	7	9,340	38,407	0	0									
8	3	1	11	30	30	9	34	72	28	8	17	1,672	1,888	5,699	0									
0	0	1	11	30	30	9	34	72	28	8	32	81	2,189	3,113	0									
0	0	1	11	30	30	9	34	72	28	0	74	583	5,382	3,104	0									
0	0	1	11	30	30	0	0	0	0	89	118	0	0	0	0	0	0	0	0	0	0	0	0	0

Table.1. Dimensional table for OLAP service, generated directly from the clustering results

To smooth out the various coding policies of different hospitals of similar clinical conditions, we introduced AHRQ clinical classification system (CCS) ICD remapping [2].

To study the treatment pattern of burn trauma, 1519 cases with information about area of injury were extracted from the 2860 cases in the 1/20 sampled datasets with clinical details.

Clustering programs

The TreeView and Hierarchical Cluster

Explorer (HCE) available in the public domain during the bioinformatics hype were surveyed [4][6].

Mathwork Matlab 6.1, with the Database and Statistical toolboxes, was used for data connection and clustering. The cases were sorted in clusters, and the dimensional table was then generated, in regular user-specified intervals (based on the number of clusters or criteria of similarity), then delivered to the OLAP service. An example were given in Table 1 and Figure 1.

Exploratory data analysis was performed in the java-based IndexMiner data mining package, developed by the senior author [1].

Various combinations of data rescaling, normalization, linkage and similarity function were tried [3][8].

RESULT & DISCUSSION:

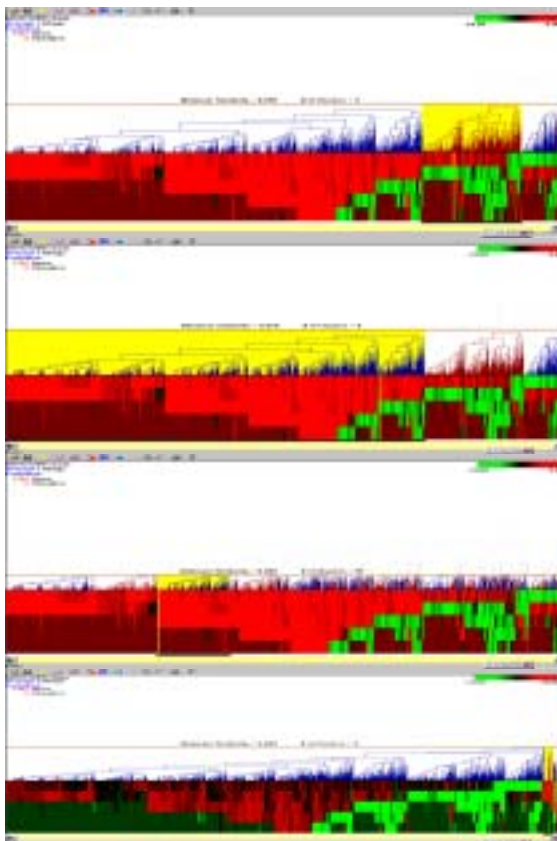


Fig.2a,b,c.. HAC of multiple trauma in Taiwan. Red designated for various trauma injuries; green, for associated nontraumatic diseases. Fig.2d, burn cases highlighted after data rescaling. (by HCE)

Distribution of multiple trauma

The CCS numerical codes were clustered after rescaling the trauma codes (from 16.1 to 16.9) to between 1 and 9 (red), and recoding of the other non-traumatic disease to their negative values (green) with single linkage (Fig.2). If the criterion of similarity was set too high, only the clusters with much variable disease combination were shown (Fig.2a). The lower, and more homogeneous, clusters need change of the threshold for separation, comparing Fig.2c with Fig.2b. The interactive nature supported the need of customizable views in a data warehouse.

To show the burn injury of particular interest, the burn diagnosis was recoded to 18, and then the burn clusters with or without associated injuries or diseases were highlighted clearly.

Treatment patterns of burn injury

According to the “vector models” in text mining[9], the most frequently present medical orders or the orders with statistical significance according to users’ interest yielded small clusters (Fig.3). Some of the purified patient groups could arouse clinical interest, but the overall pattern were not clear.

After data normalization and “windowing” specific to the clinical interest, the importance of management of extreme values and missing values was stressed.

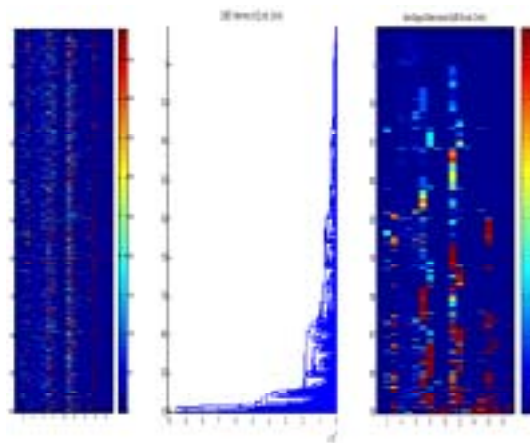


Fig.3a. Unsorted burn data set; Fig.3b,c, the dendrogram and respective clusters after sorting. (by Matlab)

Dressing-change, debridement operations, and skin grafting were finally used as criteria of clustering (Table 2). Six groups of treatment patterns were noted (Fig.4). More heavy-grafting type and heavy-dressing-change type occurred in hospitals with burn specialties. No-grafting-and-low-fee and more-dressing-change-than-debridement styles were favored by hospitals without burn specialties.

Clustering of treatment patterns improved the fit and decreased the errors of regression models for prediction of total medical expenses, and helped the judgment of outliers (Fig.5). After clustering, R square of multiple linear regression was improved from 0.650 to 0.702 and 0.876, and R square of regression trees were improved from 0.697 to 0.790 and 0.924, by demographic factors and cluster information, or by addition of length of stay as an independent variable. (Table 3).

The patterns of clustering could be used

CaseID	Case No	BURN_CD	DEBR	SENG	VENI_LAYD	SI_LAYD	SI_DEBR	SI_SENG	SI_VENI	SI_LAYD	SI_DEBR	SI_SENG	SI_VENI
Outlier	27	16,821	46,476	4,724	17,333	14,362	44,480	8,563	17,332				
MGCD	99	15,281	34,811	30,686	3,408	14,240	28,411	6,982	11,778				
HGCD	26	10,829	30,853	4,987	208	11,923	6,240	902	1,859				
KGCD	82	1,165	3,179	3,586	8	1,844	2,597	1,346	0				
MG	118	11,311	1,873	8	8	7,964	2,859	0	0				
MGD	73	5,731	17,575	8	168	4,259	7,074	0	1,439				
MGCD	373	1,293	1,787	13	44	1,364	2,300	386	423				

Table 2. Clustering center of the 4 attributes.

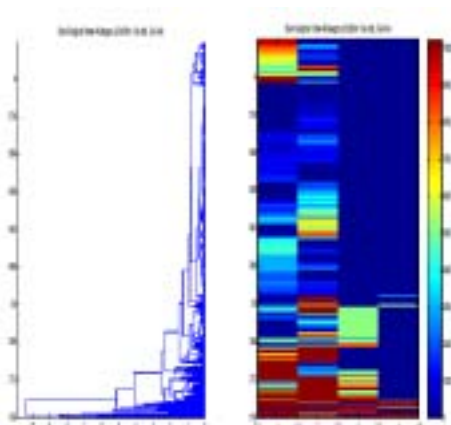


Fig.4. Six major treatment pattern after regrouping of attributes to 1) Dressing-change, 2) debridement operatDressing-change, debridement operations, and 3), 4) two kinds of skin grafting.

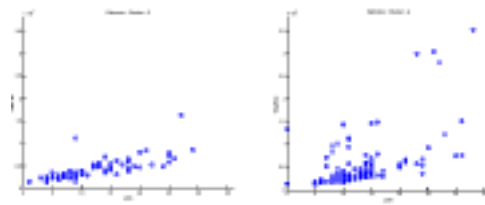


Fig. 5. Hospitalization fee versus days of hospitalization, cleaner data in more homogeneous patient population after clustering.

Table 3-1. Prediction of TotalFee by demographic factors without or with the help of clustering.

		Correlation coefficient	R square	Mean absolute error	Root mean absolute error
Prediction by demographic factors only	Multiple linear regression	0.606	0.603	56,193	99,890
	ML Regression tree	0.835	0.697	44,313	97,219
By demographic and cluster information	Multiple linear regression	0.838	0.782	46,489	88,438
	ML Regression tree	0.889	0.780	31,256	74,356
With days of burn ward admission	Multiple linear regression	0.856	0.876	27,599	56,962
	ML Regression tree	0.965	0.924	16,495	44,861

Table 3-2. Prediction of TotalFee in each cluster.

Case No	Multiple linear regression				ML Regression tree				
	Correlation coefficient	R square	Mean absolute error	Root mean absolute error	Correlation coefficient	R square	Mean absolute error	Root mean absolute error	
MG	0.9	0.842	0.887	78,497	115,874	0.949	0.906	69,830	108,297
MGCD	26	0.739	0.576	18,566	53,481	0.786	0.634	31,747	46,793
HGCD	82	0.795	0.629	10,002	17,452	0.880	0.640	9,823	16,400
MG	118	0.802	0.814	15,040	25,286	0.947	0.997	11,733	18,739
MGCD	73	0.547	0.299	17,785	87,281	0.942	0.887	20,446	28,688
MGCD	373	0.768	0.590	6,262	12,887	0.81	0.631	5,416	11,132
All Cases	973	0.836	0.876	27,598	56,962	0.961	0.924	16,495	44,861

Table.3. Benefits of clustering on the regression tree model for prediction of burn outcome.

as reference for hospital administration, for comparison of high-risk and low-risk patient components. More delicate studies could be arose by the homogeneous patient partitions after cluster analysis.

CONCLUSION & FUTURE WORKS:

We demonstrated the feasibility of hierarchical clustering methodology for research multiple trauma in a huge claim database. The intervention of domain knowledge in many processes of analysis and the interactive nature were addressed, which is especially important for participation of domain experts to handle complicated problems in medical informatics.

In the first example, the overall trauma pattern of the five year NHIRD were explored by the numerical values of CCS

codes, the numerical values were for proximity, and their order was only a byproduct. In the future, we should treat the CCS or ICD codes as simply nominal categories, and define newer distance function.

In the beginning of the second example, the vector model of medical orders could be further explored. The “meaning” of the smaller clusters might be understood with the help of other external utility scores. We think the advancements in text mining and information retrieval could contribute to part of the solution.

REFERENCES:

- [1] Chiang I.J., Lin T.Y., Index Miner: a data mining system, In Proceedings of the Annual International Computer Software and Applications Conference (COMPSAC'2001), Chicago, IL, 2001, 613-614.
- [2] Clinical Classifications Software (CCS), developed by and available at the Agency for Healthcare Research and Quality (AHRQ), <http://www.ahrq.gov/data/hcup/ccsfact.htm>.
- [3] Eisen JA. Phylogenomics: Improving Functional Predictions for Uncharacterized Genes by Evolutionary Analysis. *Genome Research* 1998; 8:163-167.
- [4] Eisen MB, Cluster, and TreeView, current version April 17, 2000. Available at <http://rana.lbl.gov/EisenSoftware.htm>
- [5] Everitt B, Landau S. Cluster analysis. 4th ed ed. London, New York: E. Arnold. Halsted Press, 2001.
- [6] Seo J, Bshneiderman B, "Interactively Exploring Hierarchical Clustering Results," *IEEE Computer*, Volume 35, Number 7, pp. 80-86, July 2002. Software available at <http://www.cs.umd.edu/hcil/hce/>
- [7] Iezzoni LI, Foley SM, Daley J, Hughes J, Fisher ES, Heeren T. Comorbidities, complications, and coding bias. Does the number of diagnosis codes matter in predicting in-hospital mortality? *JAMA* 267: 2197-203, 1992.
- [8] Kaufman L. Finding groups in data an introduction to cluster analysis. New York: Wiley, 1990.
- [9] Koller, Daphne and Sahami, Mehran. Hierarchically classifying documents using very few words . Proc. of the 14th International Conference on Machine Learning ICML97. pp.170---178.
- [10] Lee, Chien-Chih; Liu, Charles C. H.; Li, Yu-Chuan; Chiang, I-Jen, and Lu, Shiu-Yien. Epidemiology of limb replantation from 1996 to 2000 in Taiwan in National Health Insurance database. Proceedings of Annual Meetings of Taiwan Society of Plastic Surgery, 2002; Taipei. 2002.
- [11] National Healthcare Insurance Research Database, <http://www.nhri.org.tw/nhird/>
- [12] Scully KW, Schubart JR, Einbinder J. S. Improving Search Results for Diagnoses or Procedures Using an ICD-9-CM Clustering Standard. Proc AMIA Symp 11332000.
- [13] Yang I-Chih, Liu, Charles CH, Li Yu-Chuan, Chiang I-Jen, and Lu Shiu-Yien. Epidemiology of burn hospitalization from 1996 to 2000 in Taiwan in National Health Insurance database. Proceedings of Annual Meetings of Taiwan Society of Plastic Surgery 2002, 2002.