

論文名稱：應用詞頻以改良多元貝氏定理

總頁數：88

於文件分類之研究

校(院)所組別：中國文化大學商學院資訊管理研究所

畢業時間及提要別：100 學年度第 1 學期碩士學位論文提要

研究生：羅仁君

指導教授：吳宏文

李中彥

論文提要內容：

多元貝氏定理(multimembership Bayesian，簡稱 MMB)近幾年曾運用於醫療、網站、郵件和文件的自動分類與推論，到目前 MMB 的在知識分類推論領域上的相關研究一直都持續進行著。而本研究是針對 MMB 文件自動分類提出以基因演算法為概念的改良方法，使用動態產生篩選門檻來達到真正完全的自動分類。

以基因演算法為概念的 MMB 改良方式「自動調適篩選門檻」來提取重要字詞進行 MMB 分類，最後也以「自動評估」取得最佳結果。經實驗發現，當類別彼此差異度大時，最佳分類準確率為 83.93%；類別彼此差異度小時，其分類準確率為 70.60%。

關鍵字：多元貝氏定理(multimembership Bayesian)、文件分類(document classification)、文件自動分類(automatic document classification)、基因演算法(genetic algorithm)。

A research on applying term frequency to improve
multimembership Bayesian theorem
on document classification

Student: Jen-Chun Lo

Advisor: Prof . Homer Wu

Prof . Chong-Yen Lee

Chinese Culture University

ABSTRACT

In recent years, multimembership Bayesian (MMB) has had a wide application for medical, website, E-mail and other document processing use – the practices list above utilize the automatic classification and knowledge inference function of MMB to improve efficiency. Given MMB's practicality and popularity across all walks of lives, the research around MMB remains a constant focus academically. To further improve the strength of MMB's core automatic document classification function, our study proposes the additional application of genetic algorithm before traditional MMB.

Based on the law of probability, the extra step of genetic algorithm helps develop the "automatic adaptable screening threshold" mathematically, thus with more accuracy. Such calculation pin-points the significant, frequently-used words to form the threshold for further MMB classification. Since the application of genetic algorithm acts as the initial screen, consequently, the extracted leftovers are more precise for any document classification. Further, based on the mathematic results, the threshold leads to automatic assessment, which selects the most desirable word choice automatically.

The research presents significantly improved results. When class differences are relatively in large degree, its classification accuracy achieves 83.93%. Even when class differences are to a lesser extent, its classification accuracy rate of is able to reach 70.60%.

Key Words: multimembership Bayesian, document classification, automatic document classification, genetic algorithm

誌 謝 辭

原本以為研究所的日子很漫長，繁多的報告、可怕的英文、數不清且令人緊張的師生討論時間，最後還是得咬牙撐著。回頭一望才驚覺時間飛逝，終於可以結束這個煎熬。

一路走來要感謝的人很多，首先我要感謝的是我的兩位指導教授吳宏文老師與李中彥老師，有您們的支持幫助與寶貴意見，讓我在進行論文時，遇到任何問題都能迎刃而解，在此向兩位指導教授致上最高的感謝與敬意。

接著感謝的是論文口試期間陳武倚老師與王福星老師的細心指導與建議，不僅讓我受益良多，最重要的是讓我的論文能更趨完整。另外，還有百忙之中抽空前來指導論文的校外口試委員辜輝超老師，感謝您的寶貴建議與鼓勵。

再來是研究所期間，感謝那群一同跟著老師作研究、相互扶持鼓勵的同伴們。最後，要感謝我最親愛的家人，你們的支持讓我可以沒有任何煩惱的從大學念到研究所。

承蒙諸位恩師的悉心指導，在您們的諄諄教誨下，研究所的課程與論文才能順利完成，師恩浩翰，無以回報，在此謹向諸位恩師致上最高的敬意！

內容目錄

中文摘要	iii
英文摘要	iv
誌謝辭	v
內容目錄	vi
表目錄	viii
圖目錄	x
第一章 緒論	1
第一節 研究背景	1
第二節 研究目的	3
第三節 研究限制	4
第四節 研究流程	5
第二章 文件分類之文獻探討	7
第一節 文件分類	7
第二節 文件分類方法	8
第三節 特徵選擇	23
第三章 多元貝氏定理之文獻探討	27
第一節 MMB 於醫療診斷之應用	28
第二節 MMB 於網站分類之應用	29
第三節 MMB 於郵件分類之應用	36
第四節 MMB 於文件分類之應用	39
第四章 研究方法與系統架構	45
第五章 系統實作與實驗結果	54
第一節 系統實作	54
第二節 實驗結果	65
第六章 結論、貢獻與未來建議	79

第一節	結論與貢獻	79
第二節	未來建議	83
參考文獻		84



表 目 錄

表 3-1	MMB 醫療診斷之變數名稱說明	28
表 3-2	網站字詞出現記錄	30
表 3-3	MMB 網站自動分類之變數名稱說明	30
表 3-4	網站字詞 W_1 、 W_2 、 W_3 、 W_4 、 W_5 與 W_6 的 P_{ij} 與 \bar{P}_{ij} 值	31
表 3-5	MMB 郵件自動分類之變數名稱說明	37
表 3-6	MMB 文件自動分類之變數名稱說明	40
表 4-1	P_{ij} 值計算表	48
表 4-2	\bar{P}_{ij} 值計算表	49
表 4-3	類別 C_l 的 MMB 推論知識表	49
表 5-1	複數還原單數	55
表 5-2	文件比例 7:3 之樣本分布表	56
表 5-3	文件比例 8:2 之樣本分布表	57
表 5-4	名詞資料庫	58
表 5-5	文件與名詞之關係	58
表 5-6	類別之字詞出現文章數	59
表 5-7	P_{ij} 值與 \bar{P}_{ij} 值計算結果	60
表 5-8	大類別分類之自動調適出現次數門檻	65
表 5-9	大類別分類之自動調適字詞貢獻差異度門檻	66
表 5-10	大類別分類之自動調適出現次數與字詞貢獻差異度 門檻	66
表 5-11	大類別分類之第三四分位數門檻	67
表 5-12	大類別分類之字詞貢獻差異度門檻	68
表 5-13	大類別分類之跨類別分析	70
表 5-14	小類別分類之自動調適出現次數門檻	71
表 5-15	小類別分類之自動調適字詞貢獻差異度門檻	72

表 5-16 小類別分類之自動調適出現次數與字詞貢獻差異度 門檻	73
表 5-17 小類別分類之第三四分位數門檻	74
表 5-18 小類別分類之字詞貢獻差異度門檻	76
表 5-19 小類別分類之跨類別分析	78
表 6-1 大類別綜合比較	79
表 6-2 小類別綜合比較	81



圖 目 錄

圖 1-1	研究架構圖	6
圖 2-1	支持向量機示意圖	8
圖 2-2	貝氏文件分類示意圖	12
圖 2-3	基因演算法示意圖	15
圖 2-4	k-means 演算法示意圖	18
圖 2-5	決策樹示意圖	21
圖 3-1	多元貝氏定理之網站分類知識產生示意圖	29
圖 3-2	多元貝氏定理之網站知識建構示意圖	32
圖 3-3	多元貝氏定理之網站知識推論示意圖	34
圖 3-4	多元貝氏定理之郵件知識建構示意圖	36
圖 3-5	多元貝氏定理之郵件知識推論示意圖	38
圖 3-6	多元貝氏定理之文件知識建構示意圖	39
圖 3-7	多元貝氏定理之文件知識推論示意圖	41
圖 3-8	多元貝氏定理之文件知識推論(字詞篩選)示意圖	42
圖 3-9	多元貝氏定理之文件知識推論(差異度篩選)示意圖	43
圖 4-1	系統架構圖	45
圖 4-2	前置處理流程圖	46
圖 4-3	知識建構模組流程圖	47
圖 4-4	MMB 公式的推論流程圖	50
圖 4-5	自動調適出現次數門檻	51
圖 4-6	自動調適字詞貢獻差異度門檻	52
圖 4-7	自動調適出現次數與字詞貢獻差異度門檻	53
圖 5-1	詞性分析後的文件(字詞/詞性)	54
圖 5-2	程式撰寫畫面之示意圖	61
圖 5-3	程式執行畫面之示意圖	62

圖 5-4	程式執行結果之示意圖	63
圖 5-5	字詞貢獻差異度執行結果之示意圖	64
圖 5-6	大類別分類之四分位數門檻	67
圖 5-7	大類別分類之字詞貢獻差異度門檻	69
圖 5-8	小類別分類之四分位數門檻	75
圖 5-9	小類別分類之字詞貢獻差異度門檻	77

