

論文名稱：以多元貝氏定理建構文件自動分類
系統

總頁數:86

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：97 學年度第 2 學期碩士學位論文提要

研究生：周政淇

指導教授：李中彥

論文提要內容：

由於網際網路日益發達，使得數位化文件增加的速度越來越快，所以文件自動分類的重要性也慢慢增加，如何在更短的時間內將文件做最快與最正確的分類是文件分類領域中永遠的課題。

本系統是利用多元貝式定理，建構文件自動分類機制，可以有效的對文件進行分類。本系統的建置分為兩個階段，第一階段為訓練階段，透過建立字詞庫，再利用字詞庫針對訓練文件全文進行字詞的比對與過濾，最後在計算所有比對出字詞的 P_{ij} 值與 \bar{P}_{ij} 值。第二階段為測試推論，將 P_{ij} 值與 \bar{P}_{ij} 值代入多元貝式定理後便可計算出文件分別所屬類別的機率。多元貝式定理的優點在於機率的計算是可隨著已知文件的增加而逐次調整與修正，當有新樣本加入時只需局部調整某些機率值，即可得到新的分類模型，其分類模型的機動性相當高，在資料不斷的增加的情況下，可以得到較好的分類效能。實驗結果顯示本研究中使用的方法在文件分類上的正確率可以達到 95%。

關鍵字：多元貝氏(multimembership Bayesian)，文件分類(document classification)，文件自動分類(automatic document categorization)。

A Document Classification System Based on Multimem-
bership Bayesian theorem

Student: Zeng Chi Chou

Advisor: Prof. Chong-Yen Lee

Chinese Culture University

ABSTRACT

As a result of the development of Internet, it makes the increasing speed of digital documents faster. So the important of document automatic classification increases too. How to classify the documents quickly and correctly in shorter time is a very important question at the domain of document automatic classification. In this paper we establish the document automatic classification by Multimembership Bayesian. It could classify and manage documents usefully. In the training phrase, we establish the database of information word and to compare with the training documents. Finally, we compute the P_{ij} and \bar{P}_{ij} to the Multi-membership Bayesian formula then we will get the probability of document belong the class. The merit of Multi-membership Bayesian is the value of probability will be modifying by the mount of document increase. We only need to modify some value of opportunity and we will get a new classification module when the new samples get in on. For this reason, the classification module has high mobility. The more documents increase, the more effective module we have.

Key Words: multimembership Bayesian, document classification, automatic document categorization

誌 謝 辭

總算將論文告了個段落，這一路走來所遇到的酸苦至今已化為煙雲，取而代之的是滿滿的感動、喜悅，以及對支持陪伴我的家人與朋友難以言喻的感激，在向未來勇敢邁步之前，我想對我的指導教授李中彥博士致上最高的感謝與敬意，李老師在本論文的寫作上給了我相當多的啟示與提點，讓我在遇到挫折時不至困於迷霧中，而其對完美的堅持及給予我們的肯定與信任，更讓我對人生有另一番的體認與想法，希望自己在未來的旅途上也能帶給身邊的人如此的溫暖與智慧。論文口試期間王美慈老師、郭乃文老師細心的指導與建議不僅使本論文更臻完滿，其對學生的提攜愛護之情亦讓人感動萬分。謹以本文獻給我最可愛也最親愛的家人，不善於情感表達的我想對你們說我愛你們，感謝你們所付出的一切。



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 研究背景與動機	1
第二節 研究目的	2
第三節 研究範圍與限制	3
第四節 研究流程	4
第二章 文件分類相關文獻探討	6
第一節 文件分類的定義	6
第二節 文件自動分類	8
第三節 文件分類的方法	14
第四節 相關研究	24
第三章 研究方法與系統架構	32
第一節 系統架構	32
第二節 知識建構與推論	39
第四章 實作與評估	44
第一節 實作流程	44
第二節 實作結果	52
第五章 結論和未來建議	77
第一節 結論	77
第二節 未來建議	81
參考文獻	82

表 目 錄

表 2-1	國內相關碩士研究論文一	24
表 2-2	國內相關碩士研究論文二	25
表 2-3	國內相關碩士研究論文三	25
表 2-4	國內相關碩士研究論文四	26
表 2-5	國內相關碩士研究論文五	27
表 2-6	國內相關碩士研究論文六	27
表 2-7	國外相關研究一	28
表 2-8	國外相關研究二	29
表 2-9	國外相關研究三	30
表 2-10	國外相關研究四	31
表 3-1	P_{ij} 值計算表	40
表 3-2	\bar{P}_{ij} 值計算表	41
表 3-3	C_1 類別的MMB推論知識表	42
表 4-1	訓練樣本	44
表 4-2	作者處理規則	48
表 4-3	測試文件類別分佈	52
表 4-4	以題目字詞進行分類結果	53
表 4-5	以作者人名進行分類結果	55
表 4-6	以摘要字詞進行分類結果	56
表 4-7	以本文字詞進行分類結果	58
表 4-8	以摘要字詞進行25%門檻值分類結果	60
表 4-9	以本文字詞進行25%門檻值分類結果	62
表 4-10	以摘要字詞進行50%門檻值分類結果	63
表 4-11	以本文字詞進行50%門檻值分類結果	65
表 4-12	以摘要字詞進行75%門檻值分類結果	67
表 4-13	以本文字詞進行75%門檻值分類結果	68

表 4-14	以 P_{ij} 與 \bar{P}_{ij} 相差值為0.5分類結果	70
表 4-15	以 P_{ij} 與 \bar{P}_{ij} 相差值為0.4分類結果	71
表 4-16	以 P_{ij} 與 \bar{P}_{ij} 相差值為0.35分類結果	73



圖 目 錄

圖 1-1	研究流程圖	5
圖 2-1	向量表示圖	15
圖 2-2	類神經網路示意圖	20
圖 3-1	系統架構圖	32
圖 3-2	知識建構模組流程圖	33
圖 3-3	題目、摘要、本文推論流程圖	36
圖 3-4	加入門檻值後的推論流程	37
圖 3-5	加入反頻率後的推論流程	38
圖 3-6	作者推論流程圖	39
圖 4-1	資訊字詞庫	46
圖 4-2	字詞比對結果	47
圖 4-3	字詞知識表示法	47
圖 4-4	作者知識表示法	48
圖 4-5	類別文件出現次數	49
圖 4-6	P_{ij} 值計算結果	50
圖 4-7	\bar{P}_{ij} 值計算結果	50
圖 4-8	本文門檻值表示法	51