

維基百科瀏覽輔助介面—— 整合連結探勘與語意關聯分析

Browsing and Visualizing Wikipedia through Link
Mining and Semantic Relatedness Analysis

吳怡瑾

I-Chin Wu

天主教輔仁大學資訊管理學系 助理教授
Assistant Professor
Department of Information Management
Fu Jen Catholic University

張鈞甯

Chun-Ning Chang

天主教輔仁大學資訊管理學系 研究生
Graduate Student
Department of Information Management
Fu Jen Catholic University

【摘要 Abstract】

隨著網際網路與Web 2.0技術的推陳出新，以使用者貢獻為本之新型態的社會媒體服務（social media service）網站紛紛崛起。由於網站易於開發與網頁易於存取的特性，造成網路資訊快速的成長，網路世界逐漸成為使用者獲取資訊的來源，其中維基百科（Wikipedia）更為使用者快速獲取定義、解釋等資訊的重要網路服務。由於網路資訊不斷倍增，故其延伸之主要問題為資訊超載，因此使用者經常花費許多時間尋找與過濾所需資訊。本研究即以Wikipedia為研究對象，以連結探勘與語意關聯分析技術為理論基礎，試圖建構特定主題之知識網路圖。本研究首先提出藉由Wikipedia頁面連結型態（type）與連結頻率（frequency）之連結

投稿日期：2010.12.11；接受日期：2011.5.28
email: 吳怡瑾icwu.fju@gmail.com

關聯強度法 (link strength measure) 以建構初始網路，再進一步採用以搜尋結果為依據之Normalized Google Distance (NGD) 演算法計算節點間的語意關係以建構主題網路。本研究最後採用社會網路分析指標來分析主題間之關係，並以視覺化的方式呈現研究結果。本研究透過不同使用者搜尋任務設計以評估所提出方法與建構之主題導向維基百科地圖介面之有效性，研究結果顯示該發展介面有助於協助使用者快速瀏覽Wikipedia資訊，且能協助使用者完成較複雜的任務搜尋。

With the ubiquity of the Internet and the emergence of Web 2.0 technologies, social web sites (i.e., social networking websites and, micro-blogging services) are providing unprecedented opportunities for creating user-generated content, as well as for promoting communication, collaboration and information-sharing among users. Wikipedia, one of the most famous collaborative projects on the Web, has become an extremely popular reference database for people seeking information or knowledge. However, since the number of articles and the wide variety of topics in Wikipedia is constantly expanding, it is difficult for users to find information efficiently via the hypertext links, i.e., the network of linked documents. To address the problem, we propose a hybrid approach that is based on the theories and techniques of link-based analysis and semantic relatedness analysis. Specifically, we employ a link strength measure to establish a preliminary topic network by analyzing the relationships between articles. We also refine the “Normalized Google Distance” to quantify the strength of the relationship between two articles via key terms. Then, we apply social network analysis indicators to determine the relationships between topics and visualize the analysis results in order to help users browse Wikipedia efficiently. Finally, a topic-based WikiMap is generated based on the proposed hybrid approach. We conducted a user-task oriented evaluation study to confirm that the derived topic-based WikiMap can help users browse topics and execute complicated tasks easily and efficiently.

[關鍵字 Keywords]

連結探勘；正規化Google距離；語意關聯分析；主題導向維基百科地圖
Link mining; Normalized Google distance; Semantic relatedness analysis;
Topic-based WikiMap

壹、緒論

隨著網路資訊量與日俱增，改變了以往獲取知識的途徑，書本或圖書館不再是知識獲取的唯一來源，越來越多人透過網路世界的資源來取得知識。近幾年，以使用者為本而提供新型態的社會媒體服務（social media service）網站紛紛崛起，如：MySpace、Facebook、Flicker、Wikipedia、Del.icio.us、Twitter 等，其中Blog與Wiki等更成為網路促進知識交流與傳遞的重要媒體與平臺（Hill, Gaudiot, Hall, Marks, Prinetto, & Baglio, 2006; Rosenblom, 2004），哈佛商學院教授McAfee（2006）由Web 2.0的技術中萃取六項元素，透過這些元件集成的技術，提出Enterprise 2.0的概念，這六個元件分別是：搜尋（search）、連結（links）、書寫（authoring）、標籤（tags）、延伸（extensions）與信號（signals），企圖為企業的知識工作者架構出方便的技術與工具，例如Blog、Wiki 等技術，達成企業協同運作以應用群體智慧。維基百科是Web 2.0概念下的代表性產物，也是最廣為人知的線上百科全書，由網路上的用戶們共同編輯並且茁壯，是目前使用率極高的知識獲取來源之一。在資訊不斷倍增時，延伸出了另外一個問題：「資訊超載」（information overload），當有大量的歷史性資訊需要被探勘的同時，又有高比例的新資訊出現，因此造成使用者很難分辨哪些是真正有用的資訊，大量且隨時在變化的環境，造成使用者在找尋所需資訊時花費許多精神與時間。網頁中的資料大多由非結構化與半結構化的資料所組成，要從中有效的找到使用者所需資訊，經常運用「網路探勘」（web mining）技術以提高獲取所需資訊的效率。

現今網路知識已經成為非常重要的角色，其中維基百科是使用率極高的網路知識獲取來源之一。Wikipedia是基於wiki技術的全球性多語言百科全書協作計畫，使用者可透過瀏覽器進行閱覽與修改大部分的頁面。但Wikipedia的條目超連結非常多，使用者在進行資訊瀏覽或搜尋時，容易耗費過多時間在網頁連結之間找尋真正所需的資訊，部分使用者甚至迷失在超連結（hyperlink）之間。近期，搜尋引擎與網路探勘技術之相關研究針對網頁中的超連結架構，提出排名的演算法以便進行搜尋結果的排序。相關研究如：Carrière與Kazman（1997）提出WebQuery；Kleinberg（1998）提出Hypertext-Induced Topic Search（HITS）演算法，透過分析網路拓樸來找到中心（hub）與權威（authority）頁面；各演算法中最廣為人知的PageRank演算法為Google搜尋引擎的基礎，PageRank演算法定義搜尋結果的等級，並排序搜尋結

果 (Brin & Page, 1998; Page, Brin, Motwani, & Winograd, 1998)。以上研究主要使用網路連結探勘技術解決網頁排名問題，Google更是採用了全民投票的精神進行網頁排名且有十分好的結果，但仍無法有效協助使用者進行有效率的網頁瀏覽。Chin與Chignell (2006)的研究運用部落格連結關係與社會網路技術以分析部落格社群關係，其研究提出的社會超文本模型 (social hypertext model)，主要為應用社會網路分析指標找出對應之社群意識特徵，以便更準確的判別社群關係。江憲坤、陳鴻文與楊境榮 (2008)提出自動鏈結分析演算法並輔以社會網路圖形分析技術 (social network analysis, SNA)，以判別安隆公司 (Enron) 員工電子郵件寄送關係。

近期相關研究探討透過維基百科的內部連結分析，可進行語意知識的建構，並進而透過Wikipedia之語意關係，進行搜尋工具改良或文件探勘方法之改善 (Hu, Zhang, Lu, Park, & Zhou, 2009; Milne, Witten, & Nichols, 2007; Völkel, Kröttsch, Vrandecic, Haller, & Studer, 2006)。Milne等人提出並建構Koru搜尋工具與介面，該搜尋工具以Wikipedia為知識來源以建構Koru的詞典 (thesaurus) 並進而達成字詞擴展。該研究指出透過Wikipedia之連結關係可分析主題語意關聯，研究結果顯示以主題為基礎的Koru介面較傳統關鍵字搜尋更能找到相關與有幫助的文章。此外，Hu等人 (2009) 研究以Wikipedia為知識來源協助文件分群，作者指出Wikipedia的每一篇文章即代表一個主題，而透過內部連結可得到等價關係的主題概念 (equivalent concepts)，研究結果顯示透過Wikipedia的類別確實能強化原本文件的特徵字詞並進而有較佳的分群結果。基於以上相關研究，本研究探討以連結為基礎的網站中，使用者藉由點擊網頁中的連結以瀏覽與搜尋資訊，透過這種方式使用者往往無法輕易得知網頁間的關係，且需要花許多時間逐一閱讀連結的文章，使用者因此易於迷失在網頁文章之中。維基百科中包含了廣泛的主題，這讓使用者在網頁與網頁間找尋資訊的過程中變得困難；一般使用者使用維基百科都是使用最基本的方式，點擊網頁中的連結往返 (click forwards or backwards) 來瀏覽與搜尋資訊，在傳統的介面中，使用者無法得知網頁間的關係，所以必須逐一閱讀文章內容以決定資訊需求，但結果卻往往是無意間迷失在文章之中。本研究基於過去的相關研究，探討透過Wikipedia網站的內部連結以分析特定主題之相關主題，並進而協助使用者瀏覽相關資訊。技術上，本研究嘗試提出整合網路探勘、語意關聯技術與社會網路分析 (SNA) 之資訊搜尋與瀏覽介面研究，研究以Wikipedia為研究對象，試圖透過網路探勘之連結分析與Normalized Google Distance

(NGD) 演算法 (Cilibrasi & Vitányi, 2007) 建立網路導覽地圖，並輔以社會網路分析進行網路節點分析，以協助使用者快速瀏覽，進而搜尋所需資訊。因此，主要的研究問題為以下三點：(1) 探討基於Wikipedia網頁內部連結，所提出之網頁關聯強度法 (link strength measure) 與結合NGD演算法過濾語意關係薄弱節點，所建構維基主題網路圖之差異；(2) 提出運用社會網路分析指標 (SNA) 分析網頁主題與主題關聯性，該方法主要為在介面呈現節點關係，並標註於網路圖扮演重要角色的節點；(3) 設計不同使用者搜尋任務情境，以評估所提出模式與介面相較於傳統介面對於資訊瀏覽與搜尋之效益。實作上，研究使用Java Universal Network/Graph Framework (JUNG) 開發工具，透過相關函式庫使得我們可以視覺化呈現社會網路地圖與相關分析指標；在介面呈現部分，我們評估兩種介面，分別為維基百科傳統瀏覽介面與研究提出之主題導向之維基百科知識地圖介面 (topic-based WikiMap interface)。計畫評估方法主要設計不同型態的瀏覽與搜尋任務，分別評估所提出模式與介面對於資訊瀏覽與搜尋的幫助。

貳、文獻探討

一、連結探勘與分析 (Link Mining and Analysis)

連結探勘 (link mining) 為近幾年廣為探討的研究議題，連結探勘主要應用連結分析 (link analysis)、超文本 (hypertext) 與網路探勘之網路結構探勘 (web structure mining) 技術，試圖在網頁本文之外，透過網頁本文之間的連結資訊發掘出潛藏的關係 (Liu, 2007; Srinivas, Kiran Kumar Reddy, & Govardhan, 2010)。網路上普遍常見之連結型態為out-link、in-link、reciprocal links，如圖1與圖2所示。假設以A節點為基準點，連結至B節點的連結稱為out-link，反之從B節點連結到A節點的連結稱為in-link。當兩個節點同時擁有in-link跟out-link時，可將此情形稱為擁有reciprocal links (Björneborn & Ingwersen, 2004)。Katz (2004) 對原本的連結型態定義進行整理，圖3包含了上述三種連結類型外的三種連結的類型，表1解釋了此六種類型連結的連結方式。連結分析與資料探勘分析資料不同之處在於其對象為分析連結本身；近年來連結探勘已經逐漸受到矚目，透過連結的分析可以得到更豐富準確的訊息，而連結探勘目前包含了主要七種不同的任務類型，整理如表2。

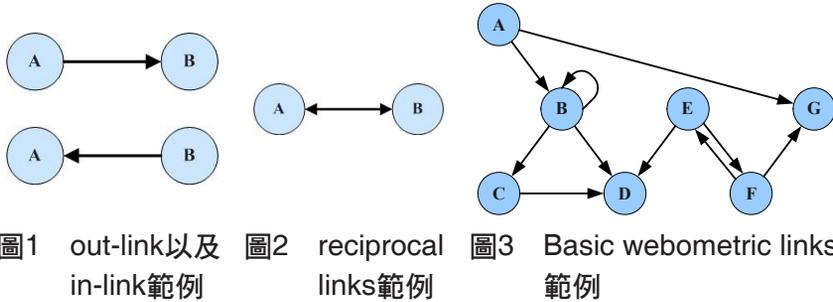


表1 連結類型的定義

連結類型	節點與節點間連結
Single links	
inlink	從A連結到B (以B為基準)
outlink	從B連結到C (以B為基準)
self-link	從B連結到B (以B為基準)
Co-links	
reciprocal links	從E連結到F，且從F連結到E
co-inlinks	從B連結到C，且從B連結到D (C與D擁有co-inlinks)
co-outlinks	從B連結到D，且從E連結到D (B與E擁有co-outlinks)

資料來源：“Co-Link web indicators of the European research area,” by J. S. Katz, 2004, from <http://www.sussex.ac.uk/Users/sylvank/pubs/Co-Link.pdf>

表2 常見的連結分析任務及分類

任務類型	分類
節點相關任務	基於連結的節點排序 (Link-Based Object Ranking)
	基於連結的節點分類 (Link-Based Object Classification)
	節點分群 (Object Clustering)
連結相關任務	連結預測 (Link Prediction)
圖形相關任務	子圖發現 (Subgraph Discovery)
	圖形分類 (Graph Classification)
	圖形的產生式模型 (Generative Model for Graphs)

資料來源：“Link mining: A survey,” by L. Getoor, & C. P. Diehl, 2005, *ACM SIGKDD Explorations Newsletter*, 7(2), pp. 3-12.

目前最具代表性的兩個連結分析的演算法為HITS演算法（Hyperlink-Induced Topic Search）與PageRank演算法。HITS演算法與PageRank演算法為著名的連結分析演算法，HITS演算法為Cornell大學的Kleinberg（1998）提出，為了解決冗長而無序的查詢目錄所開發的。HITS演算法是根據節點之間的連接關係建構關係矩陣，運用關係矩陣通過迭代的方法計算每個節點的權威分數（authority）和匯集分數（hub），根據兩種分數的高低為節點進行排名，而PageRank演算法為Google創始人Larry Page與Sergey Brin於1998年在Stanford大學發明的技術，此亦為Google搜尋引擎的基礎（Brin & Page, 1998; Page, et al., 1998）。PageRank演算法和HITS演算法一樣是透過分析WWW上的連結結構計算搜尋結果排名。PageRank的基本想法是將頁面之間的連結關係看作頁面之間權重的投票（weighted votes），一個頁面的重要性由指向它的頁面和指向它的頁面所發出的連結多寡決定。其方法是將根據節點之間的連結關係所構成的關係矩陣，運用關係矩陣通過迭代的方法為每個節點評分，再根據這個分數為節點進行排名（Page, et al.）。本研究即參考PageRank演算法與HITS演算法，提出考量連結型態（type）與連結頻率（frequency）之演算法。由於研究主要目的為由大量連結為主的網站中，過濾與分析主題並進而協助使用者瀏覽與搜尋所需資訊，因此與HITS與PageRank演算法進行網頁排序的目的略有不同。

二、NGD 語意關聯性分析（NGD Semantic Relation Analysis）

Normalized Google Distance（NGD）演算法（Cilibrasi & Vitányi, 2007）應用於判別字詞概念的問題上。NGD演算法的發展為Cilibrasi與Vitányi發現Google搜尋引擎能用來檢測兩個字詞之間的語意關聯性，例如電腦想要知道「hat」的意義，首先得建立一個字詞樹狀結構，亦即一個足以顯現字詞之間關聯性的資料庫，此種樹狀結構資料庫可以從任兩個字詞開始，以瞭解字詞間的關聯性，例如：於Google以hat與cowboy為關鍵字查詢，可以搜尋到約6,730,000筆資料，不過若鍵入cowboy與banana卻僅有1,320,000筆資料，因此可用來判定hat與cowboy之間的關聯性較cowboy與banana為強。為了要量測其關聯性，Cilibrasi與Vitányi基於搜尋到的資料筆數發展了一種統計指標，能顯示一對字詞之間的邏輯距離，並稱之為Normalized Google Distance，NGD越低代表兩個字詞之間的關聯性越緊密，其公式如下：

$$MGD(x, y) = \frac{\max\{\log f(x), \log f(y)\} - \log f(x, y)}{\log M - \min\{\log f(x), \log f(y)\}} \quad (1)$$

$f(x)$ 與 $f(y)$ 代表在Google分別以 x 與 y 為關鍵字搜尋時所回傳的結果數量，而 $f(x, y)$ 則為同時以 x 跟 y 搜尋所回傳的結果數， M 則為Google有給予索引的結果總數。若 $f(x), f(y) > 0$ 且 $f(x, y) = 0$ ，則 $NGD(x, y) = 1$ 。由公式(1)得知NGD值域如下：

1. $NGD(x, y)$ 未定義，當 $f(x) = f(y) = 0$ ；
2. $NGD(x, y) = 1$ ，當 $f(x, y) = 0$ 或 $f(x) > 0 \mid f(y) > 0$ ；
3. $NGD(x, y) \geq 0$ ，其餘狀況

Milne與Witten (2008) 之Relatednestane演算法參考Normalized Google Distance演算法將公式改良為連結導向並應用於消歧語義模糊，以兩網頁間彼此之連結數量取代Google搜尋傳回的數量，計算出兩頁面之間的關聯強度，如公式(2)所示：

$$relatedness(a, b) = \frac{\log(\max(|A|, |B|)) - \log(|A \cap B|)}{\log(|W|) - \log(\min(|A|, |B|))} \quad (2)$$

其中 a 跟 b 代表兩則有興趣的文章， A 與 B 表示分別連結至 a 跟 b 的文章集合， $A \mid B$ 代表同時連接到文章 a 與 b 的文章集合， $|W|$ 代表維基百科的總文章數目。

三、社會網路 (Social Network)

社會網路是社會科學的分支，目的在於檢視人們於社會、經濟，文化等框架 (framework) 中所扮演的角色。以人為節點 (vertices)，框架中的關係作為連結 (edges) 以建構社會關係的模型，並藉此分析個體與個體、個體與群體、群體與群體之間的互動關係及影響。社會網路分析 (SNA) 認為社會結構是客觀存在的，但並非僅是簡單的分層關係，而是複雜的網路結構。SNA便是針對這種結構，利用量化的方式進行描述分析。社會網路的架構是由actors與relations所組成，actors (角色) 可為一個個體或是組織 等，又可稱為節點；relations (關聯) 為連接actor的邊 (edges)，表示某種社會關係 (Scott, 2000; Wasserman & Faust, 1994)。

透過社會網路分析的工具可以發現actors之間的關係。社會網路分析的發展中，包含了許多指標可以衡量一個網路中的特性，而其中指標可分為兩個類型，第一種為centrality類型指標，可用來鑑別出在一個

社會網路中最為重要的actor，其中包含，如：betweenness centrality、closeness centrality、degree centrality 等，而第二種類型為cohesive可判別社會網路中凝聚力高的子集合，有n-cliques、k-cores 等指標，下面介紹研究所採用指標的基本定義。

Betweenness Centrality：此指標是為了測量出網路中，某個節點扮演著中介者的角色，以接通最多原本彼此無法直接連結的節點（Anthonisse, 1971; Freeman, 1977）。公式(3)中betweenness centrality定義為 $BC(n_i)$ ， g_{jk} 表示為節點 j 與節點 k 之連結總數， $g_{jk}(n_i)$ 表示為所有連結節點 j 與節點 k 並包含節點 i 的連結總數。

$$BC(n_i) = \sum_{j < k} \frac{g_{jk}(n_i)}{g_{jk}} \quad (3)$$

圖4中節點H之betweenness centrality最高，若將該節點移除則節點I、J將與大量節點（節點A至G）失去聯繫，相較之下節點A到G之間相互連結性高，若移除其中任一節點，剩餘節點仍能透過其他路徑接通彼此，因此相對於節點H而言其重要性則較低。

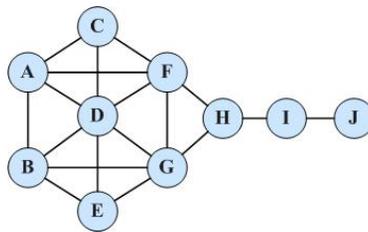


圖4 風箏結構的社會網路節點圖

資料來源：“Assessing the political landscape: Structure, cognition, and power in organizations,” by D. Krackhardt, 1990, *Administrative Science Quarterly*, 35, pp. 342-369.

Closeness Centrality：Closeness centrality指標可找出群集中距離所有其它節點平均距離最短的節點（Freeman, 1979）。公式(4)中closeness centrality定義為 $CC(n_i)$ ，其中 k 值為群集中節點的數目， $d(n_i, n_j)$ 為節點 i 與節點 j 之距離函式，公式(4)為計算所有節點至 i 節點的距離總和之倒數，故 $CC(n_i)$ 值愈小表示該節點為距離其它節點的最短節點。

$$CC(n_i) = \left[\sum_{j=1}^k d(n_i, n_j) \right] \text{ where } j | i \quad (4)$$

Degree Centrality：Degree centrality概念為某節點到其它節點的連結總數，可用來估算該節點的受歡迎程度，此項指標較大的節點被稱為中心（hub），即中心節點為擁有最大degree centrality數值的節點（Freeman, 1979; Niemincn, 1974）。公式(5)中定義degree centrality為 $DC(n_i)$ ，其中 k 值為群集中節點的數目；而degree of the node, $d(n_i)$ ，為節點 n_i 所擁有的連結數。

$$DC(n_i) = d(n_i) \quad (5)$$

參、研究目的與介面發展程序

一、研究目的

本研究將參考Kleinberg（1998）提出的HITS演算法與Brin與Page（1998）提出之PageRank演算法（Page, et al., 1998），針對Wikipedia提出網頁關聯強度法（link strength measure, LS measure）尋找出關係強烈之主題，並進而建立初始主題網路（preliminary topic network）；此外，研究將運用NGD演算法過濾掉語意關係薄弱節點，以完成維基主題網路的建立。NGD主要精神為使用Google搜尋引擎的搜尋結果，以進行字詞間語意關聯判別，如果兩組字詞間有較多的搜尋結果，表示此兩組字詞有較為強烈的關係（Cilibrasi & Vitányi, 2007）。相關研究亦應用本體論與NGD演算法進行網頁之解模糊語意（disambiguous）的問題（Gracia, et al., 2006）。該研究顯示透過NGD演算法可不用另外建立詞庫，並可有效分析語意關聯以避免搜尋到過多但不一定相關的網頁資訊。研究進一步將運用社會網路指標分析文章之角色，主要以中心性指標定義中介文章、關鍵文章、延伸性主題文章與旁支節點文章之文章角色，並以不同符號標註，以幫助使用者對特定主題之相關文章有初步的瞭解，並進而搜尋相關文章。

研究將進行連結分析之前測（pre-test）工作，首先由Wikipedia選擇特定具有代表性的種子網頁（seed web pages），如：History of personal computer，進行連結型態與數量（frequency）分析，並將提出基於連結型態與數量之網頁連結強度演算法。由於研究所提出之關聯強度演算法

只考慮網頁間的連結關係，並未考慮網頁內容相似度，在我們的研究中發現，只透過連結分析會容易出現內容相似度較低之關聯網頁（Wu & Wu, 2009），例如：使用者想了解有關積架汽車的（Jaguar）知識，但是因為使用者沒有選用適合的關鍵字，因此搜尋結果中可能會有美洲虎（Jaguar）的知識參雜其中。研究透過使用NGD演算法算出兩網頁之間連結關聯程度，可排除上述誤將連結指向不相關文章的情況。以下簡述本研究之目的：

1. 提出考量網頁連結型態與頻率之連結強度法（link strength measure）與應用語意關聯性（NGD）演算法以建立維基百科主題網路圖，研究將分析應用NGD演算法是否可建構語意關聯較高的網路圖形，相關公式與實驗將於本文第肆章詳細說明。
2. 應用社會網路分析指標，分析與標註文章關聯性，以幫助使用者透過不同文章於網路中所扮演的角色以有效瀏覽與搜尋維基文章，將於本文第肆章說明。
3. 透過提出之方法建立一個知識領域主題地圖並設計不同任務情境，以準確率（precision）與正確率（accuracy）指標評估所提出模式與介面於資訊瀏覽與搜尋之效益，相關評估程序與指標將於本文第伍章說明。

二、研究界面發展程序

基於研究問題的闡述，此節介紹本研究所提出之主題導向維基百科知識地圖介面發展流程，如圖5所示。

1. 維基百科文章前處理：維基百科的文章是由維基語法編寫而成，而維基百科中每個文章都有相對應的XML頁面顯示其基本語法。圖6中所示維基百科中的超連結可分為外部連結及內部連結（連結維基百科內其它文章），外部連結是使用[與]符號，而內部連結是用[[以及]]符號表現，而我們將藉由分析相對應XML頁面中的內容，擷取出所需頁面中內部連結資訊，並儲存於文章資料庫中。從維基百科下載文章並且轉換成資料庫儲存的格式，抓取下來的維基百科之文章皆會以XML的格式儲存到文章資料庫，而該格式檔案中會包含文章編號（article ID）、文章主題名稱（article name）、編輯日期（editing date）以及該文章連結到外部的位址（article out-link），藉以分析各個文章的連結強度。
2. 連結強度法建立網路：以使用者所下的關鍵字所找到之頁面為初始節點，運用連結強度法，考量網頁連結頻率與三種連結型態：

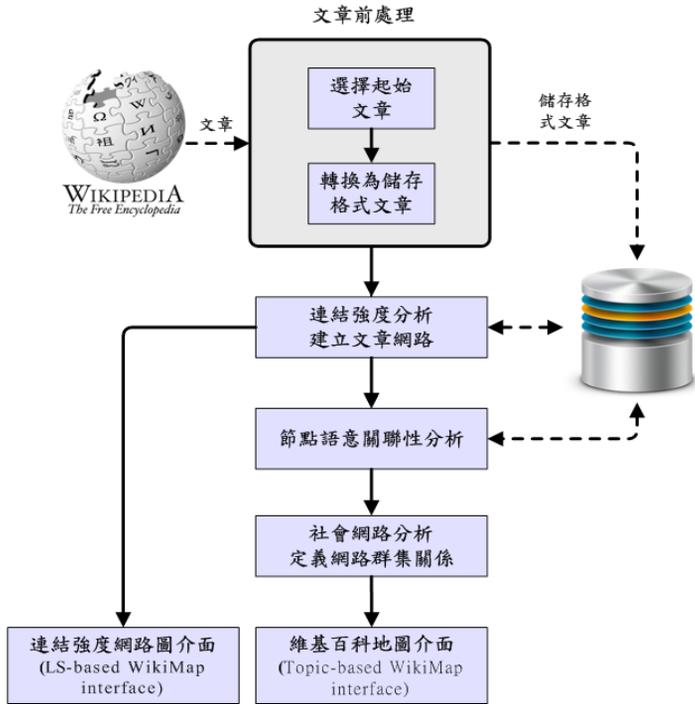


圖5 維基百科地圖產生程序

維基百科，自由的百科全書

TF-IDF

TF-IDF (term frequency-inverse document frequency) 是一種用於資訊檢索與文本挖掘的常用加權技術。TF-IDF 是一種統計方法，用以評估一詞對於一個文件集或一個語料庫中的其中一份文件的重要程度。字詞的重要性隨著它在文件中出現的次數成正比增加，但同時會隨著它在語料庫中出現的頻率成反比下降。TF-IDF 加權的各種形式常被搜索引擎應用，作為文件與用戶查詢之間相關程度的度量或評級。除了TF-IDF以外，網際網路上的搜尋引擎還會使用基於連結分析的評級方法，以確定文件在搜尋結果中出現的順序。

```

- <mediawiki xsis:schemaLocation="http://www.mediawiki.org/xml/export-0.4/ http://www.mediawiki.org/xml/export-0.4.xsd"
version="0.4" xml:lang="zh">
+ <siteinfo></siteinfo>
- <page>
<title>TF-IDF</title>
<id>406339</id>
- <revision>
<id>13186169</id>
<timestamp>2010-05-20T06:13:16Z</timestamp>
+ <contributor></contributor>
<minor/>
<comment>機器人 修改: [[فا:وزني تي:فا:فا:فا]]</comment>
- <text xml:space="preserve">

```

“TF-IDF” (term frequency-inverse document frequency) 是一種用於[[資訊檢索]]與[[文本挖掘]]的常用加權技術。TF-IDF 是一種統計方法，用以評估一詞對於一個文件集或一個[[語料庫]]中的一份文件的重要程度。字詞的重要性隨著它在文件中出現的次數或[[正比]]增加，但同時會隨著它在語料庫中出現的頻率成反比下降。TF-IDF 加權的各種形式常被[[搜索引擎]]應用，作為文件與用戶查詢之間相關程度的度量或評級。除了TF-IDF以外，互聯網上的搜尋引擎還會使用基於連結分析的評級方法，以確定文件在搜尋結果中出現的順序。 == 原理 == 在一份給定的文件裡，“詞頻” (term frequency, TF) 指

圖6 維基百科文章及XML文件之維基語法內部連結範例

in-link、out-link與reciprocal link，計算出連結強度並篩選後建立初始連結強度網路圖。節點語意關聯性分析：使用NGD演算法對文章進行語意關聯性分析，分析前項步驟留下來的網路節點之間的語意關聯性，移除關聯性低的節點以消除雜訊進而建立維基知識。

3. 社會網路分析定義網路群集關係：此處定義網路群集關係代表的含意，為找出網路中各個節點在群集中扮演的角色，以及文章群中的子群組（subgroup），例如聯繫其它群集的節點可以定義為「中介者」，群集的中心節點可定義為「關鍵文章」。研究採用社會網路分析指標，可以有效找出屬於特定性質節點在網路圖中的角色。

肆、維基百科主題地圖生成方法

一、連結強度法建立維基網路

連結強度分析（LS measure）分析兩個網頁之間不同型態連結的頻率，藉此判斷兩網頁相關程度高低，研究主要預期保留維基百科中相互連結數目高的網頁。過去相關研究以引用為基礎分析兩網頁之間的關聯程度，如：SimRank（Jeh & Widom, 2002）考慮兩網頁共同的in-link；P-Rank進一步考慮兩網頁共同的in-link與out-link以計算網頁相似度（Zhao, Han, & Sun, 2009）。本研究主要著眼於：

1. 以引用為基礎的分析需要大量的計算時間，勢必降低系統的即時性；
2. 研究同時考慮普遍常見之三種連結型態為out-link、in-link、reciprocal links，並以out-link與reciprocal links為主要分析的連結型態，主要著眼於瞭解由特定主題為種子節點（seed node）之相關維基網頁並考慮相互引用之網頁必定有較高的關聯性。

本研究所提出連結強度分析如公式(6)所示：

$$LinkStrengt|(a_i, a_j) = f_{b(a_i, a_j)}^{w_1} + f_{in(a_i, a_j)} \times w_2 + f_{out(a_i, a_j)} \times w_3 \quad (6)$$

a_i 與 a_j 分別代表 i 與 j 兩篇文章， f_{bi} 、 f_{in} 與 f_{out} 分別為reciprocal link、in-link與out-link的頻率， w_1 、 w_2 與 w_3 為根據不同連結類型的重要性差異所給予之權重， w_2 與 w_3 之總和為1；權重的選定方面，本研究預先測試了不同權重值以及門檻值設定，再對分別產生的圖形評定結果優劣，

以選定權重數值。圖7的(a)到(c)為 $w_2 = 0.9$ 與 $w_3 = 0.1$ 權重設定下所產生之圖形，在圖7(c)發現所分析出的中心節點，通常為在編輯維基百科時非常容易被指向的頁面（in-link比例高），如：Microprocessor，而我們設置之起始頁面History of personal computer反而變成Microprocessor節點的旁枝；最後圖7(d)為本研究認為權重與門檻設置之理想圖形。相關權重分析結果將於第陸章第一節說明。

二、以NGD進行網路語意關聯分析

節點語意關聯分析即為使用NGD計算網頁（節點）之間關聯性，其演算法可用來判斷兩個概念或是字詞的關聯程度，在本研究的實驗中，我們將概念定義為維基百科文章中的標題，透過兩網頁標題的分析計算，可得到兩網頁之間關聯性強度，也就是求得網路中兩節點之間的距離，透過此演算法在節點之間皆有一關聯距離的數值，而其數值為越低越好，而再將高於門檻值的文章過濾，去除掉關聯性低的節點以過濾雜訊，其演算法計算式如第貳章公式(1)。NGD門檻之相關分析結果將於第陸章第二節說明。

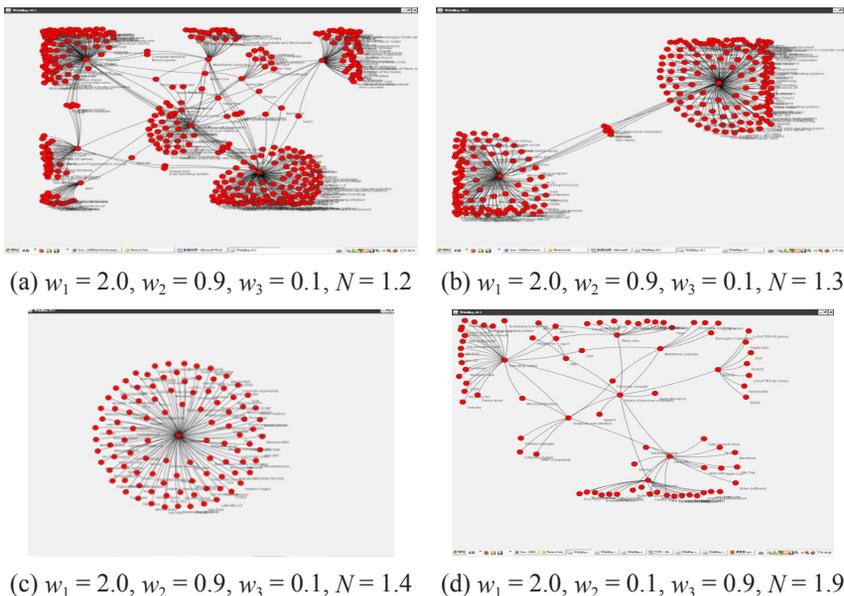


圖7 連結強度演算法在不同權重及門檻值設定所產生之圖形

三、以社會網路分析指標分析維基主題地圖

(一)社會網路分析指標

研究運用社會網路指標分析網路中文章相關性，以有效幫助使用者瀏覽文章。本研究中使用的社會網路分析指標，並且對於各指標名詞及視覺上的定義，表3為其各節點計算之社會網路指標數值，以及在網路圖中標記之外觀，以下為各名詞定義說明。

中介文章 (Betweenness Articles)：本研究使用社會網路指標中的 betweenness centrality 指標來找出中介文章節點，並標註為延伸主題的入口。使用者點選該節點會有功能選項顯示該節點所在的另一主題的文章網路，幫助使用者在瀏覽時可以延伸閱讀至相近的主題，在搜尋任務或瀏覽時可以獲得幫助。

關鍵文章 (Key Articles)：本研究運用 closeness centrality 找出文章網路中距離所有節點距離總和最小的節點，使用者在此頁面瀏覽時，可以透過平均最少的點擊超連結動作，就可以到達文章網路中的任一節點，本研究將此節點的文章定義為關鍵文章。

延伸主題文章 (Extension Topic Articles)：此指標 degree centrality 單純計算與目標節點相鄰的節點數量，本研究認為當一個文章擁有大量的外部連結（即大量的相鄰節點），則可判別為主題性的文章，例如：英文維基百科中的文章「Mouse (Computing)」，介紹電腦滑鼠且文章中許多連結延伸之其相關技術與類似之發明，可視為滑鼠這個發明的主題文章。

表3 本研究使用指標、社會網路指標與使用標記

名詞定義 (角色)	視覺定義	社會網路指標	標記
中介文章 (Betweenness Article)	群組間的中介點	Betweenness Centrality	
關鍵文章 (Key Article)	文章網路的中心	Closeness Centrality	
延伸主題文章 (Extension Topic Article)	星狀網路	Degree Centrality > 5	
旁枝節點文章	其它節點	Degree Centrality = 1	

(二) 建構維基百科主題地圖

本研究系統產生維基百科地圖的流程如下：

1. 選定起始頁面：在此我們以人工判斷找尋一個可作為主題文章的網頁，如上述提到的「History of personal computer」。
2. 設定運算的最高階層與計算相鄰階層之連結強度：研究設定層級為2，因此程式將以「History of personal computer」為起點抓取兩階層內之節點。程式計算出所有第一階層節點與起始節點的連結強度，之後將所有連結強度平均後乘以一參數值當作門檻，過濾掉第一階層中強度較弱節點。剩餘的第一階層節點的群集，從該群集中與起始節點強度最高的節點開始，與第二階層的節點做連結強度計算，與上述步驟相同依序過濾掉第二階層節點後，就生成了連結強度網路。
3. 語意關聯演算法計算連結強度網路：連結強度網路在我們的觀察下，會有實質上與主題不相關的節點留下，所以我們再運用判斷語意概念距離值的Normalized Google Distance演算法，計算兩兩節點之間的關聯距離，計算兩個節點的語意關聯性，消除網路圖中之雜訊。
4. 社會網路分析：上個步驟生成的網路地圖，再分別運用本研究使用的四個社會網路指標運算後，標示出節點在網路地圖中所扮演的角色，幫助使用者搜尋及瀏覽。

經過上述四個步驟之後，即生成維基百科主題地圖。圖8為本研究產生「History of personal computer」維基百科地圖之部分圖形範例。

伍、介面說明與實驗設計

本章節主要目的為針對維基百科傳統瀏覽介面（Wikipedia, 介面A），與所建構之主題導向之維基百科地圖瀏覽介面（topic-based WikiMap, 介面B）進行實驗評估。評估流程、任務設計與指標將於各小節詳細說明。

一、實驗流程與發展介面簡介

使用者進入實驗室後會先填寫基本資料，在實驗開始前會由實驗控制者先介紹兩種使用者介面上之注意事項，隨後開始讓使用者執行任務，任務結束後讓使用者根據其觀感填寫問卷，實驗流程如下圖9。研究提供使用者兩種瀏覽介面：維基百科傳統瀏覽介面與主題導向之維基百

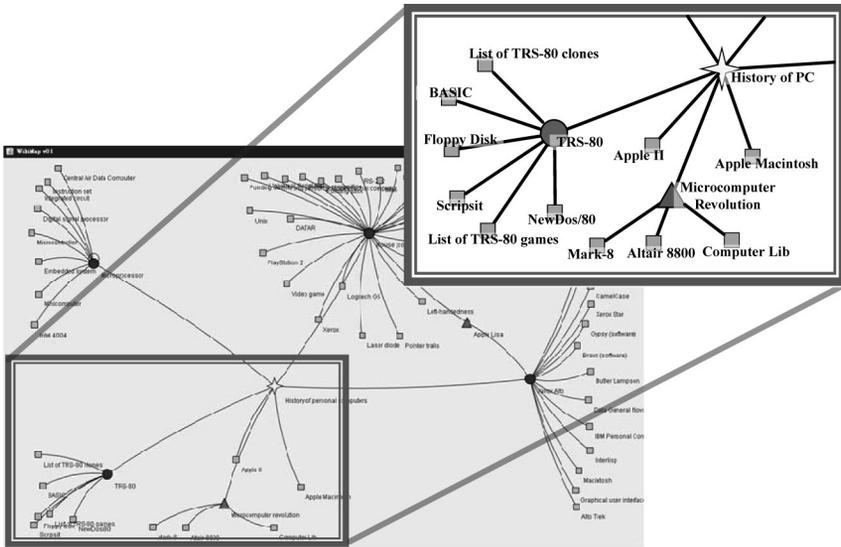


圖8 維基百科地圖與角色標記範例
(起始頁面：History of personal computer)

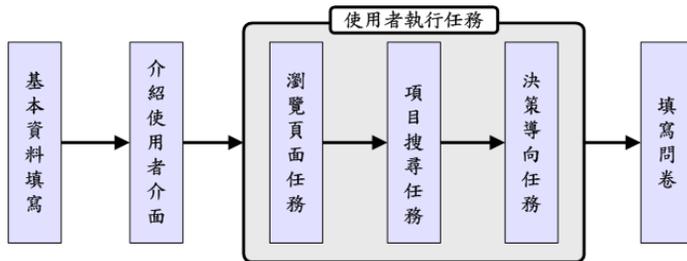


圖9 實驗流程圖

科地圖瀏覽介面以分別執行任務，圖10與圖11為本研究改良先前研究設計之使用者瀏覽介面。

二、評估對象與任務設計

評估對象根據使用維基百科的經驗，分為經驗豐富與鮮少使用的使用者，主要分類標準為使用者使用維基百科的年限與頻率。而本研究實驗使用英文版維基百科，因此為了避免語言不熟悉所造成的誤差，我們挑選英文程度較佳的受測者，如：TOEIC成績均在750分以上的學生。評估對象總數為30人，經驗豐富者與鮮少使用者各15人。而實驗任務之設

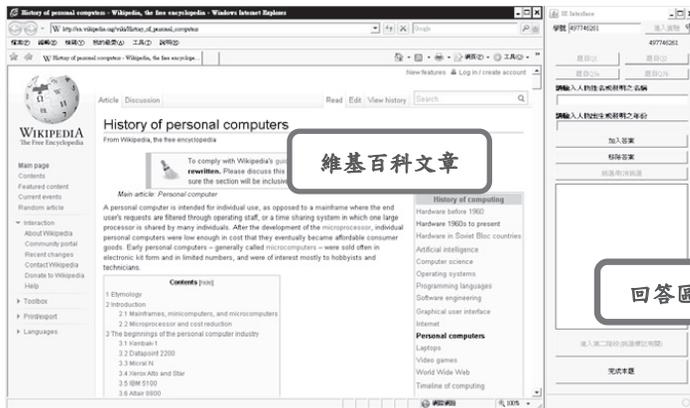


圖10 維基百科傳統瀏覽器介面 (Wikipedia, 介面A)



圖11 主題導向之維基百科地圖瀏覽介面 (topic-based WikiMap, 介面B)

計如下表4，實驗題目設計共有三個類型，每個類型有兩小題，每名受測者皆會使用介面A與B執行類型一的2個小題，類型二與三的4個題目，避免因為題目的不同對各介面表現造成差異。

三、評估指標

(一)實驗法

研究對於使用者回答的答案評估，在Q2、Q3任務的評估方式採用準確率，其公式如下：

表4 維基百科地圖使用者搜尋任務

Q1. 瀏覽頁面任務
Q1a 從History of PC這個頁面開始，藉由點擊超連結到達Macintosh這個頁面。
Q1b 從History of PC這個頁面開始，藉由地圖的瀏覽找到iPod這個頁面。
Q2. 短時間（5min），項目搜尋任務
Q2a 從History of PC這個頁面開始，找尋電腦相關的人物，並且附上出生年分。
Q2b 從History of PC這個頁面開始，找尋電腦相關的發明以及其發表日期。
Q3. 長時間（15min），決策導向任務
Q3a-1 從America Civil War這個頁面開始，找尋美國獨立戰爭中相關的事件。
Q3a-2 並標記該事件是否與Abraham Lincoln有關。
Q3b-1 從Star Trek這個頁面開始，找尋影集中出現的角色。
Q3b-2 標記其種族是否為人類。

$$Precision = \frac{|使用者回覆的答案中，為正確的數量|}{|使用者回覆的答案總數|} \quad (7)$$

基於Q3任務的特性（決策導向任務），本研究使用正確率加以評估效能，其公式如下：

$$Accuracy = \frac{|由正確答案中，正確判斷正、負相關的數量|}{|正確答案數量|} \quad (8)$$

透過以上之定量分析找出使用者在使用不同介面時，對答案選擇的影響與差異，藉此評估本研究系統之優劣。

(二)問卷調查

本研究問卷分為兩個部分，第一個部分會詢問使用者基本資訊，如：使用維基百科的經驗、過去參與過的英文檢定分數等，以分類使用者類型，可輔以定量分析得知有經驗者與無經驗者在各介面使用上是否有所差異；另外必須過濾可能會因英文能力的優劣造成的實驗誤差。第二部分針對使用者使用兩介面完成搜尋任務後的觀感進行後測，測量尺度採用李克特尺度量表，使用者根據設計的問題在五個尺度中擇一回答。

陸、實驗結果與討論

一、建立維基網路分析

(一)連結強度法權重分析

本研究連結強度分析法門檻值的設定為平均連結強度乘以 N 值。平均連結強度 ($AvgLS$) 為一個頁面若有 m 個超連結，則該頁面與 m 個頁面之間的連結強度加總後平均即為平均連結強度。而研究中取門檻的方式為先以 $AvgLS$ 設定門檻，但是由於各個Wikipedia頁面的平均連結數差異極大，故 N 值為調節網路節點數目的常數， c ，門檻 ($Threshold$) 公式如下。

$$Threshold = c \times AvgLS,$$
$$where, AvgLS = \left(\sum_{k=1}^m LS \right) / m \quad (9)$$

表5中可以發現當 w_2 與 w_3 分別設定為0.9跟0.1時，節點數的變化十分不穩定；圖12中我們考慮了雙向連結的權重設置，若是降低了該權重值，同樣僅有非常少數的節點能通過門檻。第肆章所示之圖7的(a)到(c)為 $w_2 = 0.9$ 與 $w_3 = 0.1$ 權重設定下所產生之圖形，即考慮in-link的重要度高於out-link的狀態；我們在此設定值之下發現(1) 門檻值逐漸提升後會導致餘下節點快速減少，如表5的分析；(2) 在圖7(c)發現所分析出的中心節點，通常為在編輯維基百科時非常容易被指向的頁面 (in-link比例高)，如：Microprocessor，而我們設置之起始頁面History of personal computer反而變成Microprocessor節點的旁枝。圖7(d)為設定out-link的重要度高於in-link的狀態，在不同門檻設定下，除了節點數目的差異，均呈現如圖7(d)之結果，即能夠以種子節點為出發點分析出多個相關主題並不會有第(2)點之狀況，本研究僅列出其中一分析結果，圖7(d)為本研究認為權重與門檻設置之理想圖形。

(二)以NGD進行網路語意關聯分析結果

NGD門檻值根據Evangelista與Kjos-Hanssen (2006) 之研究測定NGD期望值會在0.7左右，本研究請了7名資訊管理系之碩士生判斷在連結分析後餘下節點與主題「History of personal computer」是否相關，表6為判別節點與主題無關之節點名稱，我們計算了同時有兩人以上判別與電腦歷史無關的節點之加權平均結果為0.695，與前面學者研究所

表5 不同連結型態權重與門檻值之節點變化

	節點數目 ($w_1 = 2, w_2 = 0.1, w_3 = 0.9$)	節點數目 ($w_1 = 2, w_2 = 0.9, w_3 = 0.1$)
主題一：History of personal computer (Not a main article)		
$N = 2.1$	51	0
$N = 1.9$	93	0
$N = 1.7$	147	0
$N = 1.5$	438	0
主題二：Star Trek (A main article)		
$N = 2.1$	43	106
$N = 1.9$	46	221
$N = 1.7$	102	1312
$N = 1.5$	357	1590
主題三：American Civil War (A main article)		
$N = 2.1$	183	368
$N = 1.9$	204	521
$N = 1.7$	356	721
$N = 1.5$	550	1339

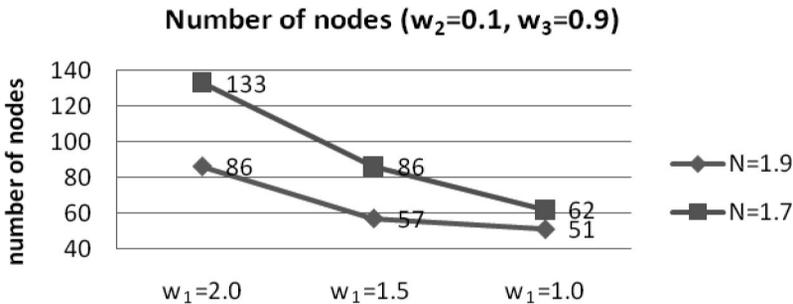


圖12 權重值 w_1 調整之節點數量之變化

述之期望值約為0.7極為相近，因此我們將門檻值設置為0.7。我們將此門檻值進行另一主題「Star Trek」測試，在連結強度分析後留下之節點「Cult Following」（引領另類潮流作品的追隨者）與其連結之節點，如：影集The X-file、網路血腥動畫Happy Tree Friend或美國樂團Devo等其它節點，這些與Star Trek無關的節點都將會被過濾。此外，研究進一步採用cosine measure比較連結強度分析法與透過NGD分析所產

表6 「History of personal computer」之NGD門檻值測試結果

節點名稱	判別無關次數	NGD值
74181 (ALU)	2	0.464
BASIC	2	0.876
Bluetooth	2	0.971
Burroughs_Corporation	3	0.444
Laser_diode	2	1.000
Left-handedness	3	1.000
Pixar	2	0.896
PlayStation_2	3	0.736
The_Walt_Disney_Company	3	0.421
Video_game	3	0.836
Wii_Remote	5	0.421
加權平均 = 0.695		

生之主題網路圖的整體相似度之差異。研究結果顯示主題「History of personal computer」在兩種網路圖有相似的結果，但主題「Star Trek」與「Abraham Lincoln」在結合NGD過濾之連結強度分析法均有較佳的結果，研究結果如表7所示。

表7 NGD門檻值測試結果

主題名稱	連結強度分析法 (LS measure)	結合NGD過濾之 連結強度分析法
History of PC	0.248	0.240
Star Trek	0.149	0.275
Abraham Lincoln	0.103	0.148

二、評估結果分析

(一) 瀏覽頁面任務

此一任務主要比較使用者在使用不同介面找尋特定頁面時，所需耗費的時間以及經過的路徑長度（閱讀之頁面數量，不計算起始頁面），例如：找尋Touchpad頁面時，受測者的路徑為History of Personal

Computer→Mouse (computing) →Touchpad，並且耗時兩分鐘，則紀錄為2 (120)，若耗時超過600秒還無法找到指定頁面，則屬於「迷失」之使用者。表8所示為鮮少使用及經驗豐富的使用評估結果。

觀察1：在表8可發現，使用介面B (維基百科地圖瀏覽介面) 進行瀏覽頁面任務時，不論是對於鮮少使用維基百科或是經驗豐富之受測者，其所耗費之時間與到達目的頁面所走的路徑，都僅有使用介面A (維基百科傳統瀏覽介面) 時之一半左右；而介面A在全部30名受測對象中，則出現了9人在執行任務時迷失了方向，耗費了十分鐘仍然無法完成任務設計之問題，找尋到指定之頁面。

觀察2：介面B在有網路圖的輔助下，受測者能花較少時間完成指定任務，我們亦可從表8看出經驗豐富之受測者在完成此任務的耗時及路徑，都比鮮少使用的對象還來得少。

表8 瀏覽頁面任務結果

受測對象	介面A (baseline)			介面B (topic-based WikiMap)		
	迷失總計	平均路徑	平均耗時 (秒)	迷失總計	平均路徑	平均耗時 (秒)
鮮少使用	5	3.20	408.40	0	1.47	202.67
經驗豐富	4	3.08	354.53	0	1.25	110.75

(二)項目搜尋任務

1. 準確率 (precision)

觀察1：分析受測對象使用介面A以及介面B在執行項目搜尋任務之結果，表9為鮮少使用維基百科的受測對象回覆答案之準確率結果，其中介面B的15名受測對象回覆的答案完全正確，主要差異出現於Q2a之結果，使用介面B找尋相關人物能有效提升準確率。

觀察2：在經驗豐富使用者的分析其結果與觀察1相似，使用介面B執行Q2a任務 (找尋人物) 能有效提升準確率。在Q2b部分，平均而言為介面A (基礎介面) 略優於介面B。

觀察3：根據實驗現場之觀察，有經驗之使用者在使用介面B執行任務時，由於題目為受測者所熟知之領域，大多都是憑藉個人經驗使用傳統瀏覽維基百科方式找尋答案，而鮮少使用維基百科的受測者，則對維基百科地圖的使用上較為依賴。

表9 項目搜尋任務結果——準確率

受測對象	介面	介面A (baseline)				介面B (topic-based WikiMap)			
		各題平均		正確總數	總平均	各題平均		正確總數	總平均
		Q2a	Q2b			Q2a	Q2b		
鮮少使用		0.87	1.00	11	0.93	1.00	1.00	15	1.00
經驗豐富		0.89	1.00	13	0.94	0.98	0.96	13	0.97

(三) 決策導向任務

1. 準確率

決策導向任務主要為決策導向問題，且主題內容對於本研究的受測對象而言較為陌生，主題為歷史「美國獨立戰爭」與影集「星艦迷航」。

觀察1：經驗豐富之維基百科使用者在執行決策導向任務時，對於維基百科地圖之依賴程度比執行項目搜尋任務時來得高。表10顯示，經驗豐富的使用者使用介面B答題完全正確總人數多於介面A，分別為11人與8人；介面B之平均準確率亦皆優於介面A。

觀察2：鮮少使用的受測者使用介面B執行決策導向任務，完全正確與準確率的總數及平均皆高於介面A。

觀察3：若我們進一步分析兩個決策導向任務，可以明顯觀察到題目Q3a不論是對於鮮少使用或經驗豐富受測者而言，介面B都有顯著的幫助，該題為找出與美國獨立戰爭相關「事件」。反之Q3a與Q3b找出與星迷航的「劇中角色」任務相比，對於受測者來說，前者（Q3a）較需要費時判斷與耗費心思閱讀文章內容的題目，後者（Q3b）則比較偏向直覺性；判斷文章是否為描述劇中角色的過程，可依賴維基百科文章旁的角色劇照即能快速判別，並且在文章中要快速瀏覽找尋角色名稱之超連結也較找尋事件來得容易，介面B在此題型下對於使用者之幫助則較為有限。因此，任務之難易度為後續研究需要考慮的重要因素。

表10 決策導向任務結果——準確率

受測對象	介面	介面A (baseline)				介面B (topic-based WikiMap)			
		各題平均		正確總數	總平均	各題平均		正確總數	總平均
		Q3a	Q3b			Q3a	Q3b		
鮮少使用		0.81	1.00	8	0.87	0.95	0.92	11	0.93
經驗豐富		0.82	0.96	11	0.90	0.97	0.96	12	0.96

2. 正確度 (accuracy)

觀察1：表11顯示使用者在決策任務中的正確度及答題完全正確的總人數，介面B的表現皆優於介面A。對於鮮少使用之受測對象，題目Q3a的表現差異較為明顯；而對經驗豐富者則兩題目都有明顯差異。

觀察2：若我們進一步觀察任務三的答題總數，表12顯示，平均而言經驗豐富的使用者答題數量均遠多於鮮少經驗的使用者。其結果顯示，經驗豐富的使用者能夠選擇較多的答案且不會降低正確率；而鮮少經驗的使用者在作答行為上則較為保守。

表11 決策任務結果——正確度

受測對象	介面	介面A (baseline)				介面B (topic-based WikiMap)			
		各題平均		正確總數	總平均	各題平均		正確總數	總平均
		Q3a	Q3b			Q3a	Q3b		
鮮少使用		0.94	0.96	11	0.95	1.00	0.97	13	0.98
經驗豐富		0.88	0.92	11	0.90	0.97	1.00	13	0.98

表12 不同介面、不同使用者各任務答題與正確數量

受測對象	介面	介面A (baseline)				介面B (topic-based WikiMap)			
		鮮少使用		經驗豐富		鮮少使用		經驗豐富	
		Q3a	Q3b	Q3a	Q3b	Q3a	Q3b	Q3a	Q3b
答題數量		92		120		85		113	
		55	37	53	67	37	48	68	45
正確數量		63		116		83		104	
		35	28	50	66	36	47	66	38

三、問卷調查

研究在實驗任務結束後對於受測者進行兩介面觀感之後測問卷，如表13所示，問卷結果顯示不論是對於鮮少使用維基百科的使用者或是面對不熟悉的主題，維基百科地圖可幫助使用者過濾節點並呈現主題之知識網路。兩介面最大差異為受測者要到達特定頁面時，若不知道其正確名詞時，在舊有的瀏覽方式下必須閱讀大量的資訊，而藉由瀏覽維基百科地圖來找尋，能降低在頁面與頁面間迷失的情形發生。此外，使用者觀感上認為當瀏覽維基百科時，有維基百科地圖之輔助能幫助他們取

表13 使用者觀感問卷結果

問項	受測對象/介面		鮮少使用		經驗豐富		Significant		
	介面A		介面B		介面A			介面B	
	Avg.	Sdv.	Avg.	Sdv.	Avg.	Sdv.		Avg.	Sdv.
Q1 此介面能幫助我瀏覽時不迷失方向	2.73	1.16	3.93	0.59	2.53	0.92	4.27	0.59	yes, $p = 0.000$
Q2 此介面的使用並不需花費過多的時間精力	3.00	1.46	4.00	0.76	2.53	1.19	3.87	0.64	yes, $p = 0.000$
Q3 此介面的使用是清楚且容易理解的	2.93	1.33	4.13	0.52	2.67	0.72	3.73	0.80	yes, $p = 0.000$
Q4 此介面的使用可幫助我容易取得所需資訊及應用	2.67	1.18	4.00	0.65	2.47	1.13	4.00	0.65	yes, $p = 0.000$
Q5 此介面整體而言是容易使用的	3.00	1.41	3.80	0.94	2.80	1.08	3.47	0.83	yes, $p = 0.004$
Q6 此介面的使用可提升我的任務績效	2.67	1.11	3.67	0.82	2.27	1.10	4.20	0.56	yes, $p = 0.000$
Q7 此介面的使用可增加我的任務效率	2.67	1.29	3.73	0.96	2.13	0.92	4.47	0.52	yes, $p = 0.000$
Q8 此介面的使用可增加我在任務上的生產力	2.40	1.18	4.00	1.00	2.40	1.06	4.07	0.59	yes, $p = 0.000$
Q9 此介面整體而言對於我的任務是有用的	2.80	1.15	4.27	0.80	2.53	0.92	4.27	0.59	yes, $p = 0.000$

註：Significant at $p < .01$.

得更多資訊、提升任務執行上的效率。問卷調查結果顯示，每一問項在99%信心水準進行成對樣本 t 檢定皆達顯著差異，可發現受測者對兩介面之使用觀感確實有明顯之差異，受測者對於維基百科地圖之輔助持正面之觀感。

四、評估結果討論

以下我們根據評估結果提出發現及討論：

1. 由圖13可以觀察整體平均而言，研究提出的介面B其準確率均優於介面A；其中決策導向任務（Q3a與Q3b）因問題難度較高，故平均準確率均低於項目搜尋任務（Q2a與Q2b）。

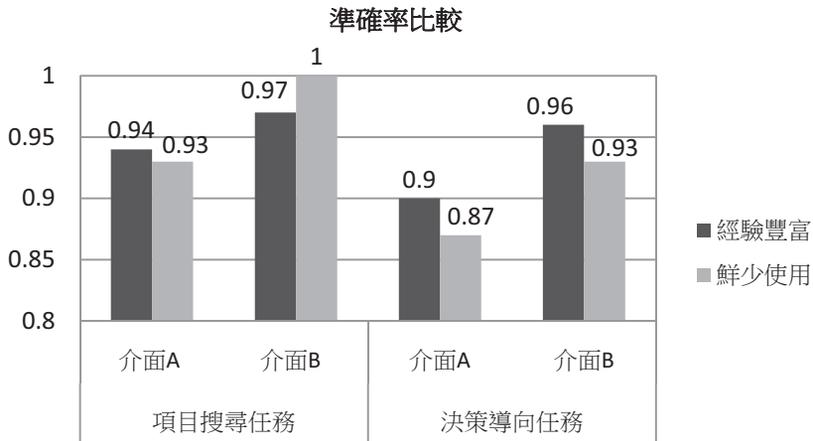


圖13 不同任務型態與不同介面之準確率比較

2. 由圖14可以觀察整體平均而言，維基百科經驗豐富的使用者其答題總數，均較維基百科使用經驗少的使用者明顯高出許多，但卻可以保持一定的準確率與正確率。表13滿意度可以發現，透過研究設計的介面B瀏覽維基百科文章，普遍對於維基百科使用經驗較為不足的使用者會有較大的幫助，維基百科地圖可以幫助其更快找到適合文章並瀏覽其內容；而對於維基百科使用經驗較為豐富的使用者而言，由於他們已經充分了解維基百科文章之編輯架構及其中規律，因此大部分的經驗豐富者在沒有維基百科地圖輔助下執行任務，亦能維持一定的作答品質。

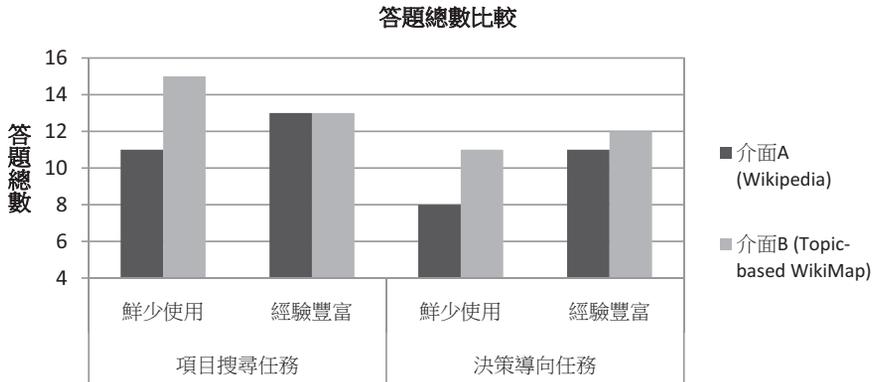


圖14 不同任務型態與不同介面之答題數量與正確總數比較

3. 使用者的觀感普遍認為研究提出之維基百科主題地圖（介面B），對於他們的任務執行有所幫助，而根據實驗過程的觀察及實際數據的結果可發現，在較為複雜且不熟悉的主題任務中，有維基百科主題地圖的幫助，能提升使用者作答的品質，並能更快速理解主題，例如：不了解星艦迷航角色名稱的受測者，能先縱觀維基百科地圖，幫助其了解主題並連結至其它角色之文章，而不需花費過多精力在閱讀主題頁面「Star Trek」，此點對於鮮少使用維基百科之受測者十分明顯。

柒、結論與未來展望

本研究提出以連結強度法、語意關聯性演算法與社會網路分析技術，建立維基主題網路地圖，除提出網路地圖建立模式並完成介面開發。本研究結果顯示，所提出之維基百科主題網路介面（topic-based WikiMap），可有效協助使用者瀏覽文章，且不易迷失方向。此外研究亦發現，搜尋結果之優劣與使用者對任務的難易度認知有關聯，未來將針對任務型態進行深入的設計與分析。本研究建構之主題導向維基百科地圖，首先經由連結分析方法過濾後，其後使用語意關聯性來過濾雜訊節點，因此本研究所產生之圖形與頁面超連結的多寡有非常密切之關係。然而維基百科為許多人共同編輯而成，我們無法控制其個人編輯習慣，例如：有的人編輯文章，某單詞每次出現在文章之中，皆會將其設置為超連結，但有些人僅會在文章中第一次出現某單詞時，將其設定為超連結，而這樣的情形也間接影響了以連結分析為主產生的維基百科主

題地圖，該限制為研究選擇任務需要考慮的要點之一。本研究設計之系統所面臨的另一限制為語意關聯性分析演算法（NGD）的計算，由於研究中產生之節點數目眾多，再加上演算法之性質，必須大量的向搜尋引擎Google送出查詢，並取得搜尋結果數目，而導致之結果為遭Google發現程式自動的送出查詢而被阻擋，因此必須耗費時間排程送出查詢，無法即時性動態產生出地圖，此為未來技術要突破之處。未來研究將朝以下三點發展：

1. 目前本研究之過濾流程為先經由連結強度法過濾節點後，再以語意關聯性演算法去除雜訊，未來可調整兩演算法使用之方法，再進一步考量圖形結果。
2. 影響實驗變數之一為介面之設計是否切合使用者的使用習慣，研究將進一步改善使用者瀏覽介面。
3. 在主題式網路圖型部分，另一進行研究結果顯示透過群集分群並輔以概念摘要呈現可以顯著提升使用者的使用滿意度與對主題的理解度（Wu & Lin, in press）；此外，維基百科的文章內容繁瑣，使用者在有時間限制下的搜尋與閱讀效率不佳，因此輔以文件摘要技術預期可提升完成搜尋任務的效率。

整體而言，本研究結果顯示該發展介面有助於協助使用者快速瀏覽Wikipedia資訊，且能協助完成較複雜任務的搜尋，該研究成果可應用在教育中之專業學科的學習輔助系統，或者研發單位中有目標主題的工作者。

致謝

本研究經費承國科會專題研究計畫補助（編號：NSC 99-2416-H-030-047-MY3）與輔仁大學學術研究計畫（編號：409931074078），僅以此致謝。

參考文獻

- 江憲坤、陳鴻文、楊境榮（2008）。自動鏈結分析演算法在社會網絡之開發與應用。*中華民國資訊管理學報*，15(3)，157-180。
- Anthonisse, J. M. (1971). *The rush in a directed graph* (Mathematische Besliskunde, No. BN 9/71). Amsterdam: Stichting Mathematisch Centrum.

- Björneborn, L., & Ingwersen, P. (2004). Towards a basic framework of webometrics. *Journal of the American Society for Information Science and Technology*, 55, 1216-1227.
- Brin, S., & Page, L. (1998). The anatomy of a large-scale hypertextual web search engine. *Computer Networks and ISDN Systems*, 30, 107-117.
- Carrière, J., & Kazman, P. (1997). WebQuery: Searching and visualizing the web through connectivity. *Computer Networks and ISDN Systems*, 29, 1257-1267.
- Chin, A., & Chignell, M. (2006). A social hypertext model for finding community in blogs. *Proceedings of the 17th ACM Conference on Hypertext and Hypermedia, Denmark*, 11-22.
- Cilibrasi, R. L., & Vitányi, P. M. B. (2007). The Google similarity distance. *IEEE Transactions on Knowledge and Data Engineering*, 19, 370-383.
- Evangelista, A. J., & Kjos-Hanssen, B. (2006). *Google distance between words*. Storrs, CT: University of Connecticut, Frontiers in Undergraduate Research.
- Freeman, L. C. (1977). A set of measures of centrality based on betweenness. *Sociometry*, 40, 35-41. Retrieved December 5, 2010, from <http://www.jstor.org/stable/3033543?seq=1>
- Freeman, L. C. (1979). Centrality in social networks: Conceptual clarification. *Social Networks*, 1, 215-239.
- Getoor, L., & Diehl, C. P. (2005). Link mining: A survey. *ACM SIGKDD Explorations Newsletter*, 7(2), 3-12.
- Gracia, J., Trillo R., Espinoza, M., & Mena, E. (2006). Querying the web: A multontology disambiguation method. *Proceedings of International Conference on Web Engineering (ICWE'06), USA*, 241-248.
- Hill, M. D., Gaudiot, J.-L., Hall, M., Marks, J., Prinetto, P., & Baglio, D. (2006). A Wiki for discussing and promoting best practices in research. *Communications of the ACM*, 49(9), 63-64.
- Hu, X., Zhang, X., Lu, C., Park, E. K., & Zhou, X. (2009). Exploiting Wikipedia as external knowledge for document clustering. *Proceedings of the 15th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, France*, 389-396.
- Jeh, G., & Widom, J. (2002) SimRank: A measure of structural-context similarity. *Proceedings of the 8th ACM SIGKDD International*

- Conference on Knowledge Discovery and Data Mining, Canada*, 538-543.
- Katz, J. S. (2004). Co-link web indicators of the European research area. Retrieved December 5, 2010, from <http://www.sussex.ac.uk/Users/sylvank/pubs/Co-Link.pdf>
- Kleinberg, J. M. (1998). Authoritative sources in a hyperlinked environment. *Proceedings of the 9th Annual ACM-SIAM Symposium on Discrete Algorithms (SODA'98), USA*, 668-677.
- Krackhardt, D. (1990). Assessing the political landscape: Structure, cognition, and power in organizations. *Administrative Science Quarterly*, 35, 342-369.
- Liu, B. (2007). *Web data mining: Exploring hyperlinks, contents and usage data*. New York: Springer.
- McAfee, A. P. (2006). Enterprise 2.0: The dawn of emergent collaboration. *Management of Technology and Innovation*, 47(3), 21-28.
- Milne, D., Witten, I. H., & Nichols, D. M. (2007). A knowledge-based search engine powered by Wikipedia. *Proceedings of the 6th ACM Conference on Information and Knowledge Management, Portugal*, 445-454.
- Milne, D., & Witten, I. H. (2008). Learning to link with Wikipedia. *Proceedings of the 17th ACM Conference on Information and Knowledge Management, USA*, 509-518.
- Niemincn, J. (1974). On the centrality in a graph. *Scandinavian Journal of Psychology*, 15(1), 322-336. Retrieved December 5, 2010, from <http://onlinelibrary.wiley.com/doi/10.1111/j.1467-9450.1974.tb00598.x/abstract>
- Page, L., Brin, S., Motwani, R., & Winograd, T. (1998). The pagerank citation ranking: Bringing order to the web. Retrieved December 5, 2010, from <http://dbpubs.stanford.edu:8090/pub/1999-66>
- Rosenblom, A. (2004). The blogosphere. *Communications of the ACM*, 47(12), 31-33.
- Scott, J. P. (2000). *Social network analysis: A handbook*. London: Sage.
- Srinivas, K., Kiran Kumar Reddy, L., & Govardhan, A. (2010). A theoretical approach to link mining for personalization. *International Journal of Computer Science Issues*, 7(3), 41-44.
- Völkel, M., Krötzsch, M., Vrandečić, D., Haller, H., & Studer, R. (2006). Semantic Wikipedia. *Proceedings of the 15th International Conference on World Wide Web, UK*, 585-594.

- Wasserman, S., & Faust, K. (1994). *Social network analysis: Methods and applications*. Cambridge, UK: Cambridge University Press.
- Wu, I.-C., & Lin, Y. H. (in press). WNav: Constructing a SNA-based navigation interface for Wikipedia. *Proceedings of the Conference on Education and Education Management, China*.
- Wu, I.-C., & Wu, C.-Y. (2009). A user-oriented topic discovery approach for effective browsing of Wikipedia. *Proceedings of the 13th International Conference on Human-Computer Interaction, USA*, 574-579.
- Zhao, P., Han, J., & Sun, Y. (2009). P-rank: A comprehensive structural similarity measure over information networks. *Proceeding of the 18th ACM Conference on Information and Knowledge Management (CIKM), Hong Kong*, 553-562.

Browsing and Visualizing Wikipedia through Link Mining and Semantic Relatedness Analysis

I-Chin Wu

Assistant Professor
Department of Information Management
Fu Jen Catholic University

Chun-Ning Chang

Graduate Student
Department of Information Management
Fu Jen Catholic University

1. Introduction

With the growing popularity of Web 2.0 technologies, social web sites (i.e., social networking websites and, micro-blogging services) are providing unprecedented opportunities for creating user-generated content, as well as for promoting communication, collaboration and information-sharing among users. Wikipedia, one of the most successful collaborative projects on the Web, has become an extremely popular reference database for people seeking information or knowledge. However, since the number of articles and the variety of topics in Wikipedia is constantly expanding, it might be difficult for users to find information efficiently by traversing through hyperlinks. To facilitate users' exploration of Wikipedia, a visualization tool was created that represented a map-like structure of topics associated with users' initial query. A task-based user evaluation was conducted to test the performance of the proposed tool. Besides demonstrating the overall performance of our system, we also discussed its effectiveness in different types of tasks and for users with difference search expertise.

2. System framework

One distinctive feature of our system is utilizing both the hypertextual and semantic structure of the Wikipedia topics. The linking structure of topics associated with users' initial query was mined. An algorithm for link strength

was proposed that taking into consideration of well-known PageRank and HITS algorithms (Formula 1), which is calculated as follow.

$$LinkStrength(a_i, a_j) = f_{bi}(a_i, a_j)^{w_1} + f_{in}(a_i, a_j) \times w_2 + f_{out}(a_i, a_j) \times w_3 \quad (1)$$

Where a_i and a_j represent i and j Wikipedia entries; f_{bi} , f_{in} and f_{out} represent the frequency of reciprocal link, in-link and out-link. w_1 , w_2 , and w_3 were weighting assigned to different types of linkages, with the sum of w_2 and w_3 to be 1. Figure 1 shows the results of varying different weighting at different linkage strength thresholds N .

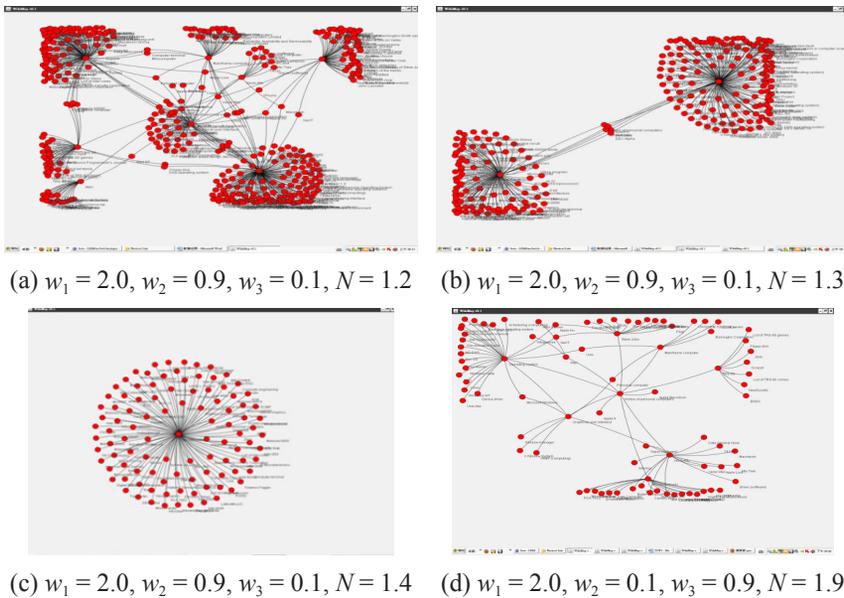


Figure 1. Graphs of Varying Different Weighting at Different Linkage Strength Thresholds N

We then apply social network analysis indicators to determine the relationships between topics and visualize the results in order to help users browse Wikipedia efficiently. Finally, a topic-based WikiMap is generated based on the proposed hybrid approach. To present the “landmarks” of the topic map, betweenness, closeness, and degree centrality measures in Social

Network Analysis (SNA) were used to help identify significant nodes in the network. We define the typology of the significant nodes as followed (see Table 1).

Table 1. Definition of Landmarks

Role	Definition	SNA metrics	Icon
Betweenness Article	Bridge nodes	Betweenness Centrality	▲
Key Article	Center	Closeness Centrality	★
Extension Topic Article	Star graph	Degree Centrality > 5	●
Other Article	Other nodes	Degree Centrality = 1	■

Excerpt from the map of topics associated with “Personal Computer History” can be seen below (see Figure 2).

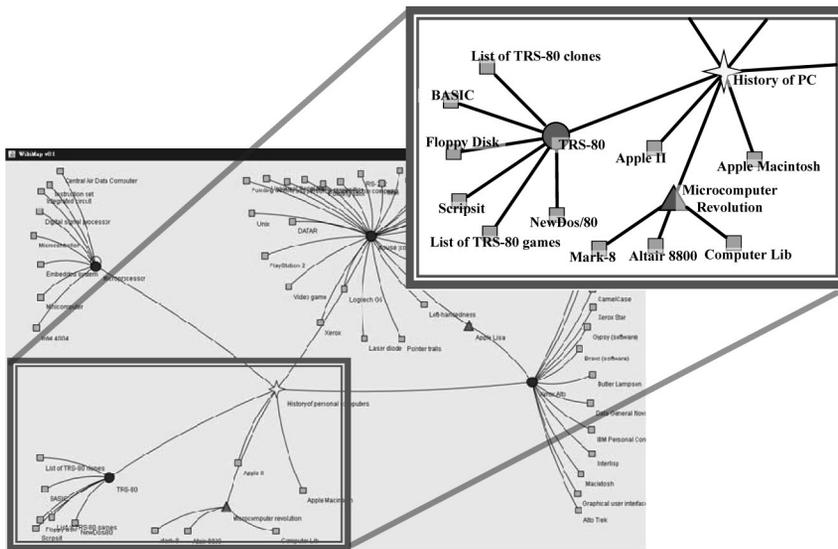


Figure 2. An Example of Using SNA Indicators to Label Topics in a Network

Besides linking structure, Normalized Google Distance (NGD) was used to quantify the strength of the relationship between two articles via key terms, which help disambiguate synonyms.

3. User study

We conducted a user-task oriented evaluation study to assess the effectiveness the topic-based WikiMap in helping users in different task situations. Three types of tasks were created (see Table 2).

Table 2. Subjects' Search Tasks in WikiMap

Q1. Task of browsing webpages
Q1a: Started from webpage of "History of PC" and clicked hyperlinks to access "Macintosh."
Q1b: Started from "History of PC" and browsed WikiMap to access "iPod."
Q2. Short time (5min), item-based search task
Q2a: Started from "History of PC" and found <u>related persons</u> with computers. Please noted their <u>year of birth</u> .
Q2b: Started from "History of PC" and found related <u>inventions</u> and their <u>published dates</u> .
Q3. Longer time (15min), decision-making oriented task
Q3a-1: Started from "America civil war" and found related <u>events</u> in American War of Independence.
Q3a-2: Noted that if the event was related with Abraham Lincoln.
Q3b-1: Started from "Star Trek" and found roles appearing in films.
Q3b-2: Noted that if the role was humanity.

A total of 30 subjects were recruited for the study, half of whom had substantial search experience, the other were novices. A within-subject design was adopted where all subject search alternately the proposed WikiMap and regular Wikipedia interface. No matter experienced users or novices, the time and number of links traversed when using the experimental interface was the half of that with baseline; besides, less subjects were reported lost with WikiMap. Not surprisingly, experienced users could use less time finishing the task than novice (see Table 3).

Table 3. Results of Page Browsing

Interface	Baseline			Topic-based WikiMap		
	Times of lost	Average routs	Average time (second)	Times of lost	Average routs	Average time (second)
Subjects						
Novices	5	3.20	408.40	0	1.47	202.67
Experienced users	4	3.08	354.53	0	1.25	110.75

Table 4 and 5 showed that the WikiMap performed slightly better in terms of precision rate in both item-based searching task and judgment tasks.

Table 4. Results in Item-Based Search Task -- Precision

Interface	Baseline				Topic-based WikiMap			
	Average		Correct number of people	Total average	Average		Correct number of people	Total average
	Q2a	Q2b			Q2a	Q2b		
Subjects								
Novices	0.87	1.00	11	0.93	1.00	1.00	15	1.00
Experienced users	0.89	1.00	13	0.94	0.98	0.96	13	0.97

Table 5. Results in Judgment Tasks -- Precision

Interface	Interface A (Baseline)				Interface B (Topic-based WikiMap)			
	Average		Correct number of people	Total average	Average		Correct number of people	Total average
	Q3a	Q3b			Q3a	Q3b		
Subjects								
Novices	0.81	1.00	8	0.87	0.95	0.92	11	0.93
Experienced users	0.82	0.96	11	0.90	0.97	0.96	12	0.96

Table 6 and 7 compared the accuracy of the two interface in judgment task. The results showed a general pattern of the experimental interface performing better than the baseline.

Table 6. Results in Judgment Tasks -- Accuracy

Interface	Baseline				Topic-based WikiMap			
	Average		Correct number of people	Total average	Average		Correct number of people	Total average
	Q3a	Q3b			Q3a	Q3b		
Novices	0.94	0.96	11	0.95	1.00	0.97	13	0.98
Experienced users	0.88	0.92	11	0.90	0.97	1.00	13	0.98

Experienced users were shown to do come up with more correct answers (see Table 7).

Table 7. Number of Answering and Correct Questions in Different Interfaces and Users

Interfaces	Baseline				Topic-based WikiMap			
	Novice		Experienced users		Novice		Experienced users	
	Q3a	Q3b	Q3a	Q3b	Q3a	Q3b	Q3a	Q3b
Answering questions	92		120		85		113	
	55	37	53	67	37	48	68	45
Correct questions	63		116		83		104	
	35	28	50	66	36	47	66	38

Table 8 showed the results of post-search question aimed at assessing user experience. There were significant differences in all question items regarding subjects' perception toward two interfaces, both in experienced and expert users groups (see Table 8).

4. Discussion

In the article we reported a framework for mining Web linkage and semantic structure for the creating of a visualization tool for browsing the complicated network of Wikipedia. An experiment was conducted to validate the effectiveness of the proposed interface, the Wikimap. Different task types were create to simulate user tasks with Wikipedia. The results showed that the Wikimap performed better both in efficiency (time spent, links traversed),

Table 8. Users' Satisfaction with the Two Interfaces

Question	Subjects/Interface		Novice				Experienced users				Significant
			Interface A		Interface B		Interface A		Interface B		
	Avg.	Sdv.	Avg.	Sdv.	Avg.	Sdv.	Avg.	Sdv.	Avg.	Sdv.	
Q1: The interface helped me not get lost in browsing	2.73	1.16	3.93	0.59	2.53	0.92	4.27	0.59	yes, $p = 0.000$		
Q2: It costs little effort to use interface	3.00	1.46	4.00	0.76	2.53	1.19	3.87	0.64	yes, $p = 0.000$		
Q3: The interface was easy to use	2.93	1.33	4.13	0.52	2.67	0.72	3.73	0.80	yes, $p = 0.000$		
Q4: The interface could help me to collect useful information	2.67	1.18	4.00	0.65	2.47	1.13	4.00	0.65	yes, $p = 0.000$		
Q5: Generally speaking, it was easy to use	3.00	1.41	3.80	0.94	2.80	1.08	3.47	0.83	yes, $p = 0.004$		
Q6: The interface helped improve my effectiveness when performing the tasks	2.67	1.11	3.67	0.82	2.27	1.10	4.20	0.56	yes, $p = 0.000$		
Q7: The interface helped improve my efficiency	2.67	1.29	3.73	0.96	2.13	0.92	4.47	0.52	yes, $p = 0.000$		
Q8: The interface helped improve my productivity	2.40	1.18	4.00	1.00	2.40	1.06	4.07	0.59	yes, $p = 0.000$		
Q9: Generally speaking, the interface was useful	2.80	1.15	4.27	0.80	2.53	0.92	4.27	0.59	yes, $p = 0.000$		

Note: Significant at $p < .01$.

effectiveness (precision and accuracy), as well as user experience. It was also found that task difficulties as perceived by the users might influence the effectiveness of the tool. Further study can explore the interaction between task types and interfaces. To better assist browsing, we also plan to include automatic summarization and clustering techniques to help users make better sense of the WikiMap, further studies can be done to explore how the proposed interface can be further improved.

References

- Anthonisse, J. M. (1971). *The rush in a directed graph* (Mathematische Besliskunde, No. BN 9/71). Amsterdam: Stichting Mathematisch Centrum.
- Björneborn, L., & Ingwersen, P. (2004). Towards a basic framework of webometrics. *Journal of the American Society for Information Science and Technology*, 55, 1216-1227.
- Brin, S., & Page, L. (1998). The anatomy of a large-scale hypertextual web search engine. *Computer Networks and ISDN Systems*, 30, 107-117.
- Carrière, J., & Kazman, P. (1997). WebQuery: Searching and visualizing the web through connectivity. *Computer Networks and ISDN Systems*, 29, 1257-1267.
- Chiang, Heien-Kun, Chen, Hown-Wen, & Yang, Jing-Rong. (2008). The development and application of an automatic link analysis algorithm for social networks. *Journal of Information Management*, 15(3), 157-180. [Text in Chinese].
- Chin, A., & Chignell, M. (2006). A social hypertext model for finding community in blogs. *Proceedings of the 17th ACM Conference on Hypertext and Hypermedia, Denmark*, 11-22.
- Cilibrasi, R. L., & Vitányi, P. M. B. (2007). The Google similarity distance. *IEEE Transactions on Knowledge and Data Engineering*, 19, 370-383.
- Evangelista, A. J., & Kjos-Hanssen, B. (2006). *Google distance between words*. Storrs, CT: University of Connecticut, Frontiers in Undergraduate Research.
- Freeman, L. C. (1977). A set of measures of centrality based on betweenness. *Sociometry*, 40, 35-41. Retrieved December 5, 2010, from <http://www.jstor.org/stable/3033543?seq=1>

- Freeman, L. C. (1979). Centrality in social networks: Conceptual clarification. *Social Networks*, 1, 215-239.
- Getoor, L., & Diehl, C. P. (2005). Link mining: A survey. *ACM SIGKDD Explorations Newsletter*, 7(2), 3-12.
- Gracia, J., Trillo R., Espinoza, M., & Mena, E. (2006). Querying the web: A multontology disambiguation method. *Proceedings of International Conference on Web Engineering (ICWE'06), USA*, 241-248.
- Hill, M. D., Gaudiot, J.-L., Hall, M., Marks, J., Prinetto, P., & Baglio, D. (2006). A Wiki for discussing and promoting best practices in research. *Communications of the ACM*, 49(9), 63-64.
- Hu, X., Zhang, X., Lu, C., Park, E. K., & Zhou, X. (2009). Exploiting Wikipedia as external knowledge for document clustering. *Proceedings of the 15th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, France*, 389-396.
- Jeh, G., & Widom, J. (2002) SimRank: A measure of structural-context similarity. *Proceedings of the 8th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, Canada*, 538-543.
- Katz, J. S. (2004). Co-link web indicators of the European research area. Retrieved December 5, 2010, from <http://www.sussex.ac.uk/Users/sylvank/pubs/Co-Link.pdf>
- Kleinberg, J. M. (1998). Authoritative sources in a hyperlinked environment. *Proceedings of the 9th Annual ACM-SIAM Symposium on Discrete Algorithms (SODA'98), USA*, 668-677.
- Krackhardt, D. (1990). Assessing the political landscape: Structure, cognition, and power in organizations. *Administrative Science Quarterly*, 35, 342-369.
- Liu, B. (2007). *Web data mining: Exploring hyperlinks, contents and usage data*. New York: Springer.
- McAfee, A. P. (2006). Enterprise 2.0: The dawn of emergent collaboration. *Management of Technology and Innovation*, 47(3), 21-28.
- Milne, D., Witten, I. H., & Nichols, D. M. (2007). A knowledge-based search engine powered by Wikipedia. *Proceedings of the 6th ACM Conference on Information and Knowledge Management, Portugal*, 445-454.
- Milne, D., & Witten, I. H. (2008). Learning to link with Wikipedia. *Proceedings of the 17th ACM Conference on Information and Knowledge Management, USA*, 509-518.

- Nieminen, J. (1974). On the centrality in a graph. *Scandinavian Journal of Psychology*, 15(1), 322-336. Retrieved December 5, 2010, from <http://onlinelibrary.wiley.com/doi/10.1111/j.1467-9450.1974.tb00598.x/abstract>
- Page, L., Brin, S., Motwani, R., & Winograd, T. (1998). The pagerank citation ranking: Bringing order to the web. Retrieved December 5, 2010, from <http://dbpubs.stanford.edu:8090/pub/1999-66>
- Rosenbloom, A. (2004). The blogosphere. *Communications of the ACM*, 47(12), 31-33.
- Scott, J. P. (2000). *Social network analysis: A handbook*. London: Sage.
- Srinivas, K., Kiran Kumar Reddy, L., & Govardhan, A. (2010). A theoretical approach to link mining for personalization. *International Journal of Computer Science Issues*, 7(3), 41-44.
- Völkel, M., Krötzsch, M., Vrandečić, D., Haller, H., & Studer, R. (2006). Semantic Wikipedia. *Proceedings of the 15th International Conference on World Wide Web, UK*, 585-594.
- Wasserman, S., & Faust, K. (1994). *Social network analysis: Methods and applications*. Cambridge, UK: Cambridge University Press.
- Wu, I.-C., & Lin, Y. H. (in press). WNAvi: Constructing a SNA-based navigation interface for Wikipedia. *Proceedings of the Conference on Education and Education Management, China*.
- Wu, I.-C., & Wu, C.-Y. (2009). A user-oriented topic discovery approach for effective browsing of Wikipedia. *Proceedings of the 13th International Conference on Human-Computer Interaction, USA*, 574-579.
- Zhao, P., Han, J., & Sun, Y. (2009). P-rank: A comprehensive structural similarity measure over information networks. *Proceeding of the 18th ACM Conference on Information and Knowledge Management (CIKM), Hong Kong*, 553-562.