

# 聲符部件排序與形聲字發音規則探勘

## Pronunciation Rules Discovery for Picto-Phonetic Chinese Characters

張嘉惠 Chia-Hui Chang

國立中央大學資訊工程學系

Department of Computer Science and Information Engineering

National Central University

[chia@csie.ncu.edu.tw](mailto:chia@csie.ncu.edu.tw)

林書彥 Shu-Yen Lin

國立中央大學資訊工程學系

Department of Computer Science and Information Engineering

National Central University

[985202041@cc.ncu.edu.tw](mailto:985202041@cc.ncu.edu.tw)

### 摘要

近年來台灣有相當多的新移民的加入，這些新移民在口語的學習上雖然有地利之變，但是在漢字的認識上則是相當弱勢。由於漢字乃是圖形文字，學習單一字的成本相對的高。如果可以讓漢字教一個字，可以學到十個字，對於漢字教學的成效應有相當的助益。本文從部件教學的概念出發，考慮聲符的發音強度、出現頻率、及筆劃數，做為聲符部件教學順序的準則。我們利用部件發音強度 [8]，以線性加總、幾合乘積、及調和平均三種方法對部件排序。根據此部件排序學習，前五個部件便可延伸學習多達 140 個相似發音的漢字。進一步，我們應用中研院文獻處理實驗室所建立的「漢字構形資料庫」，以及標記所得之形聲字，拆解形聲字組成的部件，挖掘串連漢字之間關係的形音關聯規則。我們從 600 萬條發音規則中篩選與分群出 8 條高信賴度與兩組各約 10 條高支持度的規則，並藉由這些規則來輔助漢語發音的學習效率。

關鍵詞：形聲字、部件教學、聲符強度、機率分佈、學習曲線、關聯規則

### 一、簡介

漢字是世界上最古老的文字之一，也是至今仍廣為使用一種形系文字。近年來由於中國市場的興起，以華語做為第二外語的學習也連帶地愈來愈受到重視，華語學習者的人數也倍數成長，據 *China Daily 2010* 的文章指出，目前全世界超過四千萬的非華裔人士正在學習華語文。由此可見未來華語文學習市場的龐大需求；再者，台灣近年來外籍與大陸配偶的人數從 2002 年的二十三萬人成長至今四十四萬人，其中外籍配偶約十四萬六千多人，已取得國籍者約九萬人，在在顯示了漢語學習的重要性。

過去學習漢語只能靠資深的中文老師的教導或是學習者慢慢累積經驗，不僅對於海

外華語師資的培育緩不濟急，對於學習者而言更是一條漫長的路。然而，漢語字形讀音繁複，初學者並不易掌握學習要訣，尤其漢語的發音更是複雜多變。事實上華語作為第二語言的學習，比起英文作為第二語言的學習更是難上許多，因為漢語的字形與音調相較拼音文字複雜，學習者要同時進行形、音、義三者的連結，如果沒有適當的聯想，將需要很大的記憶力，比起傳統的拼音拉丁文字，即使會說華語的海外華人對於漢字的認識也可能相當有限。其最主要的原因在於漢字是圖形文字(pictograph system)，無法像英文等拼音文字(alphabet system)一樣，一旦學會拼音方法(phonetic representation)，即有基本的閱讀能力。相較之下，一般漢字學習者讀寫的學習進展則會比較緩慢，而且必須搭配注音符號(Chinese phonetic symbols)或是其他拼音方法，才可知道每個漢字的發音。這樣的限制，對於漢字的學習相當不利，這也是為什麼二十世紀初期許多專家欲將漢字拉丁化的主要原因。

漢字的構成包含象形、指事、會意、形聲、轉注、假借(總稱六書[1])。據統計資料，7000 個現代漢語通用字中，屬於「形聲」結構的有 5631 個，約佔總字數的 80.5%，這麼多的形聲字在整字的組合上，多數採用「1+1」的方式，也就是一個意符加上一個聲符。基於這樣一個語言事實，我們可以借助部件教學，充分發揮部件的組合關係強化學習者對於漢字的識記。

本篇論文中，我們應用[8]，以部件發音分佈的集中性計算聲符強度，加以部件延伸字數及筆劃數的考量，提出線性加總、幾合乘積、及調和平均三種結合方法，對部件加以排序。利用此排序做為漢字部件教學的順序，可以幫助學習者在短時間內提高閱讀效率。我們以累計延伸字個數做為學習成效的比較，發現有效的排序，可以在學習完前五個部件，便可藉此延伸學習多達 140 個具有高度相似發音的漢字，同時累計筆劃數也是可以接受的範圍，顯示適當排序的重要性。

除了考量聲符部件學習順序之外，我們也試圖分析漢字發音規則，做為學習發音的參考。為了要產出易懂的發音規則，讓中文的學習者可以應用形聲字的特性來推測漢字的發音，在本文中我們應用關聯規則探勘挖掘形聲字發音所存在的規則。我們應用中研院文獻處理實驗室所建立的「漢字構形資料庫」，拆解其組成的部件，挖掘串連漢字發音關係的形音關聯規則，來輔助學習者學習，讓漢字不是教一個字才學到一個字，而能搭配關聯規則「一舉數字」，發揮數位學習的優點。我們從 600 萬條發音規則中篩選與分群出 8 條高信賴度與兩組各約 10 條高支持度的規則，並藉由這些規則來輔助漢語發音的學習效率。

## 二、相關研究

最早有關漢字構造的研究，應屬中央研究院資訊科學研究所文獻處理實驗室，從 1993 年開始，陸續建構古今文字的源流演變、字形結構及異體字表，做為記錄漢字形體知識的資料庫，也就是漢字構形資料庫[10]。漢字構形資料庫不僅銜接古今文字以反映字形源流演，也記錄了不同歷史時期的文字結構。另外也由於開發漢字部件檢字系統，得以解決缺字問題。然而漢字構形資料庫過去的研究著重在字形知識的整理，尚未涉及字音與字義的處理；因此文獻處理實驗室近年來開始文字學入口網站建置計畫[2,3]。一如其文所述：“漢字構形資料庫目前只著重在字形知識的整理，尚未涉及字音與字義；建立一個形、音、義俱備的漢字知識庫，仍是我們長遠的目標”。因此本論文的目的即是以挑戰漢字的發音規則知識庫為出發，除了了解漢字發音規則外，也希望藉

由此項研究找出一套形聲字發音轉換規則，讓華語學習者可以在聲符與規則的輔助下，順利讀出字的發音出來。

與本研究最為相關的研究計畫是淡江大學中文系高柏園、郭經華、胡映雪教授所主持之“字詞教學模式與學習歷程研究”。其概念是藉由即時回饋的寫字練習（學文 Easy Go!），比較部件拆解做為漢字教學策略成效（洪文斌 2010），輔以線上教學平台「IWILL Campus」（郭經華 2010），進行「以字帶詞」之詞彙學習策略（高柏園 2010）。此計畫在美國加州地區 Saratoga High School 針對 26 名修習 AP 中文課程之學生，實施四週約八堂之主題課程，用以評估漢字部件教學之學習策略對於海外華語文學習者之成效。從國科會期中報告顯示，採用多媒體自習一組的學生在認字、書寫、及字的結構上，比傳統標示筆劃順序的習字方法呈現較佳的成果，顯示以部件拆解做為漢字教學策略的可行性。

張等人於 2010 年提出了兩種自動化判定形聲字聲符的方法[8]。第一種方法為發音相似度比較法，由於聲符構件通常與原字的發音相似度高於非聲符構件與原字的發音相似度，因此經由語言學專家的協助，分別制訂聲母、韻母之間發音相似度。進一步，為了提升經由發音相似度比較法判斷聲符之準確率，採用限制性最佳化技術，求得發音相似度分數。另一種為構件發聲分佈比較法，通常做為聲符構件的漢字，其衍生字的發聲分佈比非聲符構件的漢字發聲分佈更為集中。因此作者利用一個可以計算兩個機率分佈差距的公式 KL divergence，來計算每個構件的發聲分佈與所有漢字的發聲分佈 KL 值做為構件做為聲符的強度。實驗結果顯示，發音相似度比較法在 7340 個形聲字中的判定聲符準確率為 93.35%，而構件發聲分佈比較法則可達到 98.66% 的準確率，顯示兩種方法做為聲符判斷問題的可行性。

### 三、部件重要性排序

首先我們從部件教學的概念出發，希望對於聲符的教學順序，提出一個考慮聲符發音強度、出現頻率、及筆劃數的排序方法，做為聲符部件教學順序的準則。由於構件發聲分佈比較法對於判定形聲字聲符有高達九成八的準確率，因此我們此處即採用做為聲符發音強度。根據[8]的定義，每一個部件的聲母發音強度、韻母的發音強度、及調號的發音強度可由下列三式計算而得：

$$I(w) = KL(P_I(W) \parallel P_I(A)) \quad (1)$$

$$F(w) = KL(P_F(W) \parallel P_F(A)) \quad (2)$$

$$T(w) = KL(P_T(W) \parallel P_T(A)) \quad (3)$$

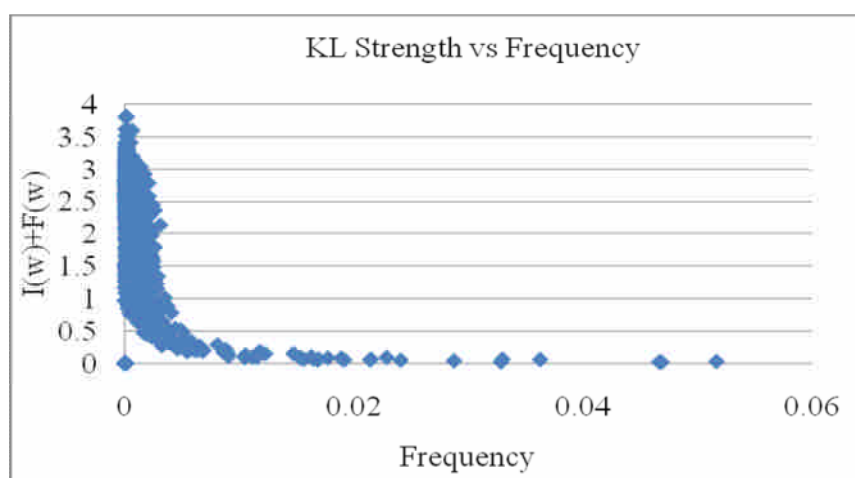
其中A表示所有漢字所成的集合，W則表示部件w所延伸的字所成的集合。函數 $P_I(A)$ 、 $P_F(A)$ 、 $P_T(A)$ 分別表示A集合中漢字的聲母、韻母及調的分佈機率。 $KL(P \parallel Q)$ 則代表兩個機率分佈的KL-divergence:

$$KL(P \parallel Q) = \sum_i P(i) \log \frac{P(i)}{Q(i)} \quad (4)$$

對於聲符而言，由於發音集中度較高，因此 w 的聲母分佈  $P_I(W)$  與所有漢字的聲母分佈  $P_I(A)$  會有較大的差異。同理韻母分佈  $P_F(W)$  與  $P_F(A)$  差異，以及聲調分佈  $P_T(W)$  與  $P_T(A)$  差異也會較大。因此我們即可以 KL-divergence 公式對此差異值計算出其程度，

換句話說我們利用公式 1, 2, 3 分別計算一個部件的聲母、韻母、及調號的 KL 值，這三種數值分別反應出此部件的聲母、韻母、及調號的發音強度。

除了部件的發音強度，在部件學習排序上，我們也必須考慮部件的頻率。因為對於漢字學習者來說，發音強的部件，也要有一定的出現頻率，才能發揮其做為聲符的功能。因此若單純以發音強度來決定教學順序，並不是非常適當的選擇。再者，對於學習者來說，漢字的筆畫數多寡也會影響學習的效率。因此如何將三者同時考慮於部件教學的順序，是此處最主要的挑戰。常見的結合方式是以線性加總，然而在此處並非最佳的結合方法，如圖一部件發音強度與頻率散佈圖顯示，若以線性加總發音強度與部件的頻率（部件頻率定義為包含部件  $w$  的形聲字字數 $|W|$ 除以全部字數），可能先找到的是頻率高但發音強度較弱的部件，或是發音強的部件但是頻率較低的部件，而非同時據有高頻及高發音強度的部件。



圖一、部件發音強度與頻率散佈圖

為了找出頻率高且發音強度強的部件，且同時也希望能將筆劃數較少的部件優先排序。我們提出三種排序部件的依據：

1. 線性加總： $ScoreA(w)=a*Freq(w)+ I(w)+ F(w) + b*Strokes(w)$
2. 幾何乘積： $ScoreG(w)= Freq(w) * (I(w) + F(w)) / \sqrt{Strokes(w)}$
3. 調和平均： $ScoreH(w)=ScoreG(w)/ScoreA(w)$

其中  $Freq(w)$ 代表部件  $w$  的頻率， $Strokes(w)$ 為部件  $w$  的筆畫數； $a$  與  $b$  則是線性加總的權重。由圖一可知發音強度約為頻率的  $a=90$  倍，同理，我們求得筆畫數的權重  $b=0.01$ ，可使線性加總的三個因素間取得平衡。第二種結合方法則是三個因素的幾何乘積，最後調和平均則是取線性加總與幾何乘積的調和平均做為部件排序的評估。

### 3.1 實驗評估

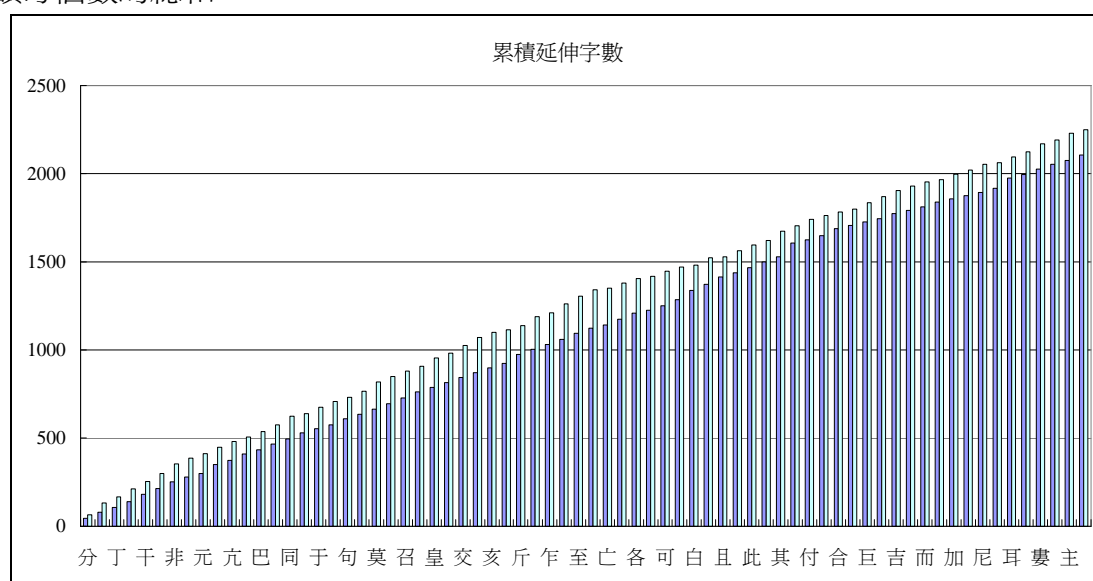
為了評估三個部件排序是否能有效率地提昇學習效率，我們繪製出以幾何乘積做為

部件排序，與其累積延伸字數的關係<sup>1</sup>。如圖二所示，橫軸表示排序過的部件，從左而右依序是：分令丁方干包等字，縱軸淺色代表累積延伸字的個數  $Y_1$ ，縱軸深色則代表聲符能正確預測聲母個數與韻母個數的總和  $Y_2$ ，兩者分別定義如下：

$$Y_1 = \sum_i |W_i|, \quad (5)$$

$$Y_2 = \sum_i (\text{Imatch}(w_i, W_i) + \text{Fmatch}(w_i, W_i)) \quad (6)$$

其中  $\text{Imatch}(w_i, W_i)$ 代表部件  $w_i$  延伸字集合  $W_i$  中與部件  $w_i$  具有相同聲母的字數，同理， $\text{Fmatch}(w_i, W_i)$ 代表部件  $w_i$  延伸字集合  $W_i$  中與部件  $w_i$  具有相同韻母的字數。舉例來說若  $w = \text{包}(\text{ㄅㄠ})$ ， $W = \{\text{炮}(\text{ㄅㄠ}), \text{胞}(\text{ㄅㄠ}), \text{苞}(\text{ㄅㄠ})\}$ ，那麼  $w$  與  $W$  中相同聲母的字數為 2 {胞、苞}，相同韻母數為 3 {炮、胞、苞}。因此兩者相加後可得正確預測聲母個數與韻母個數的總和=5。



圖二、幾何乘積排序與累積延伸字關係（縱軸取對數以減少高度）

正確預測聲母個數與韻母個數的總和 ( $Y_2$ ) 愈接近兩倍累積延伸字的個數 ( $2Y_1$ )，表示預測正確的準確率愈高，將上述兩值相除，可得準確發音比例。從圖二可以看出排序在前面的字即有相當多的延伸字，同時準確發音的比例也相當的高。表一列出排序前五個部件及其可延伸學習的形聲字，如表一所示，這些部件都具有延伸字發音高度相似、出現頻率高、筆數少的特性，益於先行學習。

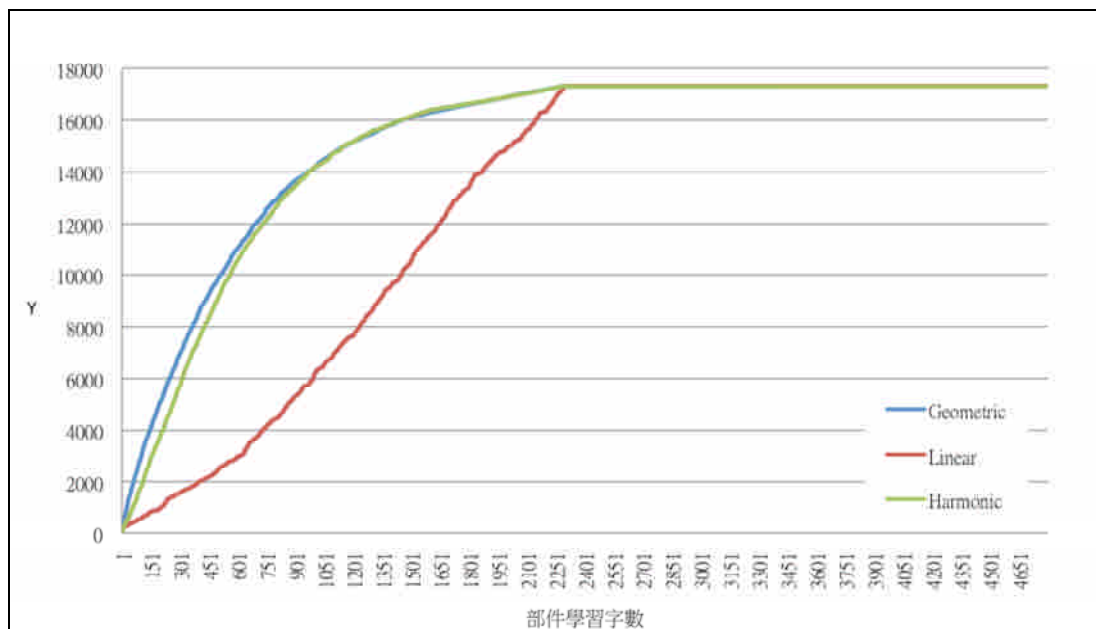
接著我們比較三種排序公序的學習曲線如圖三，同樣地橫軸為部件排序，縱軸為正確預測聲母個數與韻母個數的總和。從圖三中可看出幾何乘積排序較線性加總法來的有效，在學到 1000 字以前幾何乘積排序呈現大幅度的成長，也就是說若我們依照乘積排序的部件順序來學習，一開始便能達到快速學習到大量的延伸字。調和平均排序採用幾何乘積與線性加總算數平均法的調和，不過其走勢幾乎與幾何乘積排序相同，這點也顯示出幾何乘積排序明顯優於線性加總。

<sup>1</sup> 所有漢字相關資料來源則是使用中研院所開發的漢字構形資料庫。

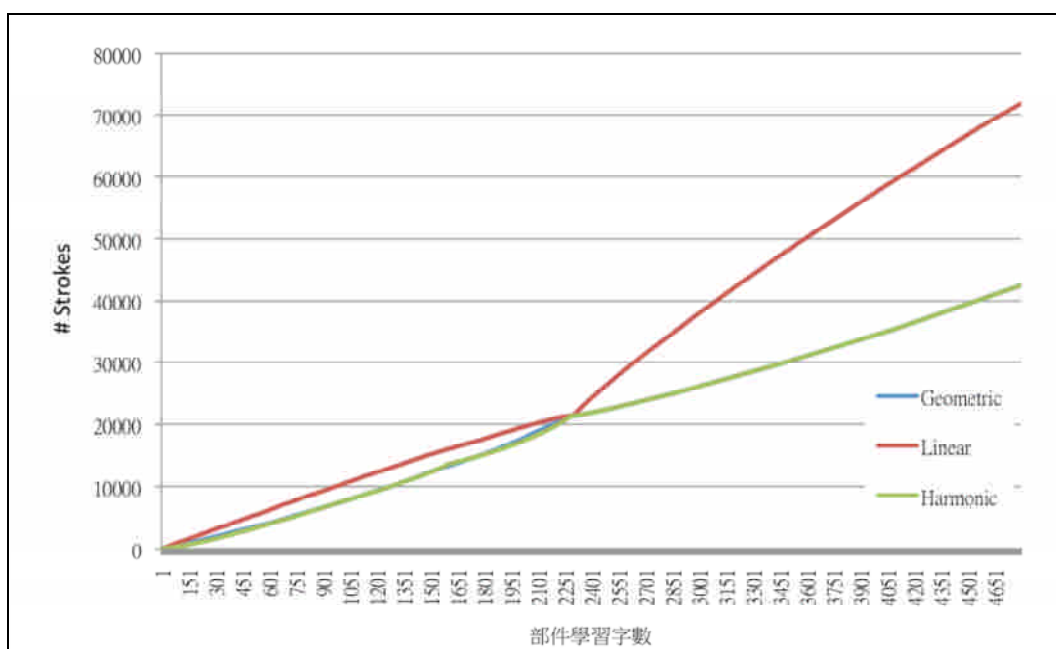
最後我們以累積筆畫數的學習曲線來看(圖四),幾何乘積排序的累積筆畫數學習曲線也較線性加總排序所得來的優異。圖三的收斂點與圖四的筆劃數大增的轉折點也顯示了在學習了 2200 個部件後,累積延伸字數已呈飽和狀態,顯示接續其後的部件已是複合部件,同時筆畫數增加速度較快。也因此可判斷排序大於 2200 後的部件並不是迫切的學習對象。

表一、幾何乘積排序之部件

部件 $w_i$	延伸字 $ W_i $	$Y_1$	$Y_2$	準確發 音比例	筆劃 數	累積筆 劃數	延伸字
分	45	45	64	0.71	4	4	份份盆奔吩吩粉榜芬...
令	35	80	132	0.83	5	9	伶冷玲呤囹嶺伶伶冷...
丁	27	107	167	0.78	2	11	汀亭叮可叮叮宁寧玎...
方	33	140	211	0.75	4	15	仿坊彷彿枋枋瓶放昉防...
干	42	182	253	0.70	3	18	刊平幹杆犴盱旱汗扞...
包	32	214	298	0.70	5	23	刨匏咆庖抱胞炮炮砲...
非	38	252	353	0.70	8	31	菲啡扉緋斐腓翡徘排...
屯	26	278	386	0.69	4	35	沌沌囤鈍屯屯炖饨饨...
元	20	298	412	0.69	4	39	兀阮完阮玩阮沅阮莞...
工	51	349	448	0.64	3	42	巨仞功左巧巫差式攻...



圖三、部件排序學習曲線比較圖



圖四、部件排序與筆畫數學習曲線比較圖

#### 四、發音規則探勘

本文第二個重點在於形聲字發音規則的探勘，藉由已標記的形聲字聲符，找出聲符與延伸的形聲字之間是否有常見的發音規則。爲了要產出易懂的發音規則，讓中文的學習可以應用形聲字的特性來推測漢字的發音，在本文中我們將應用關聯規則探勘 Apriori 演算法做爲探勘形聲字發音規則的方法。每一條關聯規則必須符合最小支持度(support)及最小信賴度(confidence)，對於學習者才算有用。以下我們首先介紹如何準備形聲字成爲關聯規則探勘所需要的交易資料，以及規則的篩選與分群，以及最終所得的發音規則。

##### 4.1 形聲字交易資料

關聯規則探勘原本的目的是從超市購買交易記錄的資料庫中，找出產品之間被購買的關聯程度，其主要依據爲支持度(support)及信賴度(confidence)。其中支持度代表一個規則的涵蓋率（全部交易資料中有多少百分比讓規則爲真），而信賴度則代表一個規則的準確率（前提爲真的情況下，有多少百分比資料讓結果也同時爲真）。爲了推測發音規則，我們以常用字中的 3000 個形聲字準備成 3000 筆交易資料。

形聲字的發音分成三個部份：聲母、韻母、以及調號，分別將其記爲 INITIAL、FINAL、TONE。另將形聲字的聲符(Phonetic component)，以及聲符的發音以 PC\_INITIAL、PC\_FINAL、PC\_TONE 三個屬性標記。其次漢字的部首(Radical component)、形聲字排列方式(單體字、左右連接、上下連接、包圍式、其他)、形聲字筆劃(Stroke)、聲符筆劃(PC\_Stroke)、兩者差值(diff\_STROKE)等特徵都列爲表達發音規則的探勘項目之一。最後，形聲字的發音若與其聲符的發音相同，則標記成聲母發音不

變(IU)、韻母不變(FU)、音調不變(TU)等項目，做為交易資料的一部份。值得一提的部份是，由於筆劃數及數值性屬性，考慮到記憶的方便性，我們統計了漢字構形資料庫中所有的漢字的筆劃數將其平分為三類。如此一來，若筆劃數在規則條件及以筆數是否大於或是小於某個範圍表示，降低規則的複雜性。每筆形聲字交易資料所包含的項目屬性如表二所示。

表二、漢字特徵對照表及”炮”的交易範例

符號	意義	數值範圍	範例:炮
INITIAL	聲母	{ø, ㄅ, ㄆ, ㄇ, …, ㄌ}	ㄆ
FINAL	韻母	{ø, 一, ㄨ, …, 儿}	么
TONE	調號	{1,2,3,4,5}	4
CONNECT	形聲字的連接方法	{單體字,左右連接,上下連接,包圍式,其他}	左右
PC	聲符	形聲字	包
PC_LOCATION	聲符所在形聲字之位置	{左,右,上,下,內,其他}	右
PC_INITIAL	聲符的聲母	{ø, ㄅ, ㄆ, ㄇ, …, ㄌ}	ㄅ
PC_FINAL	聲符的韻母	{ø, 一, ㄨ, …, 儿}	么
PC_TONE	聲符的調號	{1,2,3,4,5}	1
STROKE	形聲字筆劃數 L16 表示 >=16 14-15 表示 14 與 15 s11 表 <=11	{L16,14-15,s11}	s11
PC_STROKE	聲符筆劃數	{L16,14-15,s11}	s11
Diff_STROKE	形聲字與其聲符筆劃差值	{s3, 4-5, L6}	4-5
INITIAL_UNCHANGED(IU)	形聲字與其聲符之聲母不變	{false, true}	IU=false
FINAL_UNCHANGED(FU)	形聲字與其聲符之韻母不變	{false, true}	FU=true
TONE_UNCHANGED(TU)	形聲字與其聲符之聲調不變	{false, true}	TU=false



我們使用 weka<sup>2</sup>來進行形聲字發音規則探勘。針對最小支持度取 0.3%、0.5%與 1% 對應各種不同的最小信賴度 60%~100%，進行 Apriori 運算後，得到不同數量的規則數如表三。雖然在最小支持度 1%及最小信賴度 100%時，即可探勘出 50,054 條發音規則，然許多高支持度的規則不符合信賴度，為避免錯失重要的發音規則，以上各項參數設定中，我們取最多規則數的參數組合(最小支持度 0.3%，最小信賴度 60%情形下)，共 6,625,518 條規則存入資料庫中，做為進一步的篩選過濾。

表三、關聯規則探勘後規則數

sup \ conf	60%	70%	80%	90%	100%
0.3%	6,625,518	5,144,742	3,879,619	2,809,951	1,810,585
0.5%	1,573,613	1,149,779	802,029	500,708	314,523
1%	304,330	217,346	143,301	87,324	50,054

## 4.2 規則篩選

每條關聯規則皆是由“左邊條件[左支持度] → 右邊結果[右支持度,信賴度]”組成。雖然關聯規則探勘可以取得為數不少的發音規則，但其中有許多是不符合我們預期的規則。舉例來說：

PC\_LOCATION=右 (sup=2054) → CONNECT =左右 (sup=2054, conf=1)

上述這條規則表示“若聲符位置在右，則形聲字連接方式為左右連接”。像這樣的規則對發音的推測其實並沒有幫助。又如

INITIAL=ㄅ (sup=20) → PC\_INITIAL=ㄅ (sup=20, conf=1)

表四、篩選後規則數

sup \ conf	60%	70%	80%	90%	100%
0.3%	368,810	272,957	195,735	152,152	106,740
0.5%	61,171	32,089	15,243	7,561	5,190
1%	13,470	6,340	1,889	505	42

上述規則描述“若形聲字聲母發音為ㄅ，則其聲符聲母發音為ㄅ”。像這樣的規則也無助於推測發音，原因在於我們的本意是讓學習者在具備基礎聲符的閱讀能力下，利用對聲符的相關認知，來推測出更多尚未認識的形聲字發音。因此合法的規則應該具備：“聲

<sup>2</sup> <http://www.cs.waikato.ac.nz/ml/weka/>

符條件 或 形聲字筆劃數” → “形聲字發音 或 形聲字發音與聲符發音的關係”。根據此一篩選原則，我們統計出最小支持度與最小信賴度不同參數下合法的規則數如表四。

### 4.3 規則分群

雖然在最小支持度 0.3%，最小信賴度 60%情形下，規則篩選已將的規則數減少至 368,810 筆規則，但由於規則中有許多同質性的規則散佈在資料庫中，我們需要有系統地將它們分群。以圖五條件集為例，可以發現 1、2、3 具有相同條件「聲符的聲母=ㄉ」，且這些規則均具有相近的支持度。仔細深入查看符合這些條件的字後發現，支持這些規則的字組也相當程度的重疊（如「老」、「呂」、「里」等聲符的延伸字），所以聲符的聲母條件可以是分群的重要參考因素。

- |   |
|---|
| <ol style="list-style-type: none"><li>1. 聲符的聲母=ㄉ，聲符的調=2，聲符所在位置=右，形聲字筆劃數=12-15 (sup=17)</li><li>2. 聲符的聲母=ㄉ，聲符的筆劃數=L16，漢字與其聲符筆劃差值=4-5 (sup=16)</li><li>3. 聲符的聲母=ㄉ，聲符的調=3，漢字與其聲符筆劃差值=s3 (sup=16)</li></ol> |
|---|

圖五、發音規則條件範例

同理聲符的韻母也多涉即相同性質的規則，因此規則中若有指定相同的聲符韻母，也是我們分群的依據之一。緊接著我們繼續觀察其他規則的左方條件：

1. 部首=艸，形聲字的連接方法=上下連接，support=22
2. 部首=女，形聲字的連接方法=左右連接，support=15
3. 部首=艸，形聲字的連接方法=上下連接，support=22
4. 部首=艸，聲符所在位置=下，support=22
5. 部首=金，形聲字筆劃數=L16，support=31
6. 部首=女，形聲字的連接方法=左右連接，support=15
7. 部首=言，聲符所在位置=右，形聲字筆劃數=L16，support=28

我們發現相同部首的規則具有相近的 support 值、因此可分群成{部首=艸|1、3、4}；{部首=女|2、6}分爲同群。而形聲字的連接方法在具有相同部首的狀況下通常也會有特定的連接方法如上規則{1、3、4}；{2、6}，因此形聲字的連接方法也在我們的分群條件之內。然而，有時候會出現規則中具有相同聲符與部首等條件同時出現的規則，這時我們便要設定一個判斷分群條件優先權如下：

1. 聲符
2. 聲符聲母
3. 聲符韻母
4. 部首
5. 形聲字的連接方法

根據這些分群優先條件，便可將相同性質規則分爲同群。表五爲篩選後合法規則數及對應上述分群後的結果。

表五、篩選後合法規則數/分群後規則數

sup \ conf	60%	70%	80%	90%	100%
IU, FU	3454/332	2004/225	1097/139	597/73	264/39
IU	9002/486	5383/383	3067/262	1758/161	809/91
FU	12171/690	8373/608	4855/470	2673/325	1392/189

#### 4.4 結果

最後我們將規則分為兩類，高支持度(Support)與高信賴度(Confidence)。其中**高支持度**的規則可涵蓋形聲字較廣泛，且需要的條件較少，但具有較多不符合規則的例外字。舉其中的規則 R1 來說，當一形聲字的聲符發音為ㄉ時，藉由它的結果“聲母發音=不變”可預測這個聲符加上其他部首或是部件時有 9 成(178/197)的比例也會發音ㄉ。像是聲符「蘆(ㄉ)」+「艹」=「蘆(ㄉ)」。但由於此類規則所需條件較少，涵蓋範圍大，因此例外字也有約  $197*(1-0.1)=20$  字，像是聲符發音ㄉ的「立(ㄉ)」+水部=「泣(ㄑ)」。像這樣易於記憶(條件少)且含蓋範圍廣的規則就適合初次使用本系統者。由於高信賴度的規則往往需要搭配較多條件，但可以更準確的推測發音，例外字較少，如下面**高信賴度**規則中 R1 “符的聲母=ㄉ，聲符的調=3, w 與聲符筆劃數差值=s3”，可以看出除了上述的聲符發音ㄉ以外，加上另外兩個條件“聲符的調=3 聲符筆劃數差值=s3”便可以達到幾乎百分之百預測的效果。值得一提的是，筆劃數差值 s3(小於等於 3)也透露出聲符加上比劃數很小的部首如口、人、水...等等這樣的部首是不太影響聲符本身的發音，如里+口=哩、呂+人=侶、老+人=佬。觀察這樣的發音規則似乎也能透露出部首本身的特性。因此進階學習者較適合高信賴度的規則類別學習。<sup>3</sup>

**高支持度規則 [Query: supp>=3% and conf>=80% and IU (聲母不變)]。共 15 規則，3 群。**

1. (R1)聲母的聲母=ㄉ (supp:197) → 聲母發音=不變 (supp:178, conf:0.9)
2. (R2)聲母的聲母=ㄇ (supp:128) → 聲母發音=不變 (supp:105, conf:0.82)
3. (R3)w 的筆劃數=L16，聲符的筆劃數=L16 (supp:123) → 聲母發音=不變 (supp:98, conf: 0.8)

**高支持度規則 [Query: supp>=3% and conf>=80% and FU (韻母不變)]。24 規則，8 群。**

1. (R1)聲母的聲母=ㄉ (supp:197) → 韻母發音=不變 (supp:158, conf:0.8)

<sup>3</sup> 更多規則查詢請連結本系統 <http://hanzi.ncu.edu.tw/picpho/pronrule.php>。

2. (R2)聲符的聲母=ㄉ (supp:124)→韻母發音=不變 (supp:100, conf:0.81)
3. (R3)聲符的聲母=ㄘ, 聲符的筆劃數= $\leq 11$  (supp:121)→韻母發音=不變 (supp:99, conf:0.82)
4. (R4)聲符的韻母=ㄨ (supp:111)→韻母發音=不變 (supp:104, conf:0.94)
5. (R5)聲符的韻母=ㄨㄥ (supp:106)→韻母發音=不變 (supp:90, conf:0.85)
6. (R6)聲符的韻母=ㄨㄥ (supp:114)→韻母發音=不變 (supp:93, conf:0.82)
7. (R7)聲符的調=2, w 的筆劃數= $\geq 16$ (supp:221)→韻母發音=不變 (supp:176, conf:0.8)
8. (R8) 部首=艸, w 的連接符號=上下連接, 聲符的筆劃數= $\leq 11$  (supp:113)→韻母發音=不變 (supp:91, conf:0.81)

高信賴度規則 [Query: conf $\geq$ 100% and supp $\geq$ 0.4% an IU (聲母不變) and FU (韻母不變)]。共 34 規則，5 群。

1. (R1)聲符的聲母=ㄉ, 聲符的調=3, w 與聲符筆劃數差值=s3 (supp:16)→聲母發音=不變, 韻母發音=不變 (supp:16, conf:1)
2. (R2)聲符的聲母=ㄉ, 聲符的調=2, 聲符所在位置=右, w 的筆劃數=12-15 (supp:17) @聲母發音=不變, 韻母發音=不變 (supp:17, conf:1)
3. (R3)聲符的聲母=ㄉ, 聲符的筆劃數=L16, w 與聲符筆劃數差值=4-5 (supp:16)→聲母發音=不變, 韻母發音=不變 (supp:16, conf:1)
4. (R4)聲符的聲母=ㄚ, 聲符的調=1, w 的筆劃數=12-15, 聲符的筆劃數=s11 (supp:16) →聲母發音=不變, 韻母發音=不變 (supp:16, conf:1)
5. (R5)聲符的韻母=ㄨ, 聲符的調=1, w 的連接符號=左右連接, w 的筆劃數=s11 (supp:17)→聲母發音=不變, 韻母發音=不變(supp:17, conf:1)
6. (R6)聲符的韻母=ㄨ, w 的連接符號=左右連接, 聲符的筆劃數= $\leq 11$ , w 與聲符筆劃數差值= $\geq 6$  (supp:13)→聲母發音=不變, 韻母發音=不變 (supp:13, conf:1)
7. (R7)聲符的韻母=ㄨㄥ, w 的筆劃數=12-15, 聲符的筆劃數= $\leq 11$  (supp:13)→聲母發音=不變, 韻母發音=不變 (supp:13, conf:1)
8. (R8) 聲符的韻母=ㄨ, 聲符的調=3, 聲符的筆劃數=12-15 (supp:13)→聲母發音=不變, 韻母發音=不變 (supp:13, conf:1)

## 五、結論及未來研究

本論結論可分為二個主要方向說明。第一部份，延續機率分佈比較法，考慮到部件筆劃數少、發音強度高、出現頻率高三種因素，我們提出三種部件排序方法及兩種評量法，其中幾何平均法在筆劃數與學習字數曲線圖的表現上較為出色。第二部份，我

們藉由形聲字的特徵，運用關聯探勘法則挖掘出許多發音規則。而發音規則經由我們歸納後可分為，高支持度與高信賴度兩大類。藉由這兩大類的規則能幫助不同程度的初學者更易於推測未知漢字的發音。

然而，仍有許多地方尚待我們改進。目前形聲字規則的演進過程不夠明朗且稍嫌不夠深入，除此之外，還欠缺有效的發音規則排序及評量方法。或者，搭配部件排序，讓重要部件的規則先行教學。另一方面，還可加強形聲字查詢介面的效率以及加入破音字作為發音規則考據等等。最終目的，希望能充分發揮數位學習的優點，讓漢字的學習更為生動簡易。

## 六、致謝

本論文的完成感謝李淑萍、廖湘美及孫致文教授，以及陳怡如、葉博榮、鍾哲宇、趙婕妤等人的幫助。

## 參考文獻 [References]

- [1] 許慎撰，段玉裁注，《說文解字注》，台北藝文印書館，1988年。
- [2] 莊德明、謝清俊，[漢字構形資料庫的建置與應用](#)，漢字與全球化國際學術研討會，台北，2005年。
- [3] 莊德明、鄧賢瑛，[文字學入口網站的規畫](#)，第四屆中國文字學國際學術研討會，山東煙台，2008年。
- [4] 董鵬程，台灣華語文教學的過去、現在與未來展望. 2007多元文化與族群和諧國際研討會，台北教育大學。[http://r9.ntue.edu.tw/activity/multiculture\\_conference/memoirs.html](http://r9.ntue.edu.tw/activity/multiculture_conference/memoirs.html)。
- [5] 許聞廉、呂明秦、胡志偉、柯華蕙、辜玉旻、呂菁菁、張智凱、莊宗嚴，構建一個新移民者有機成長的多元認同平台的整合研究（期中進度報告），2009- 2011。
- [6] 高柏園、郭經華、胡映雪，華語文作為第二語言之字詞教學模式與學習歷程研究，2009-2010。
- [7] 洪文斌，華語文作為第二語言之字詞教學模式與學習歷程研究——子計畫一：中文字部件拆解教學模式與電腦輔助學習系統之研發（期中進度報告），2010。
- [8] 張嘉惠，李淑瑩，林書彥，黃嘉毅，陳志銘，《以最佳化及機率分佈判斷漢字聲符之研究》，ROCLING XXI, 2010。
- [9] 萬雲英，《兒童學習漢字的心理特徵與教學》，載於楊中芳、高尚仁主編，中國人、中國心—發展與教學篇，403-448。台北：遠流。
- [10] 盛繼豔，《華文教學中漢語的部件教學》。
- [11] 梁彥民《漢字部件區別特徵與對外漢字教學》，《語言教學與研究》2004。
- [12] 李思維、王昌茂編著，《漢字形音學》，武漢：華中師範大學出版社，2000年版。
- [13] 中研院文獻處理實驗室，「漢字構形資料庫」網站。