Assessing Geo-Location and Gender Information in Han Chinese Personal Names

Bruce Brown and Deryle Lonsdale Brigham Young University

Sixth Annual Family History Technology Workshop Provo, Utah March 9, 2006

List #	Country	Total	List #	Country	Total
1	Albania	20	51	Kenya	8
2	Angola	1	52	Korea	194
3	Argentina	32	53	Kuwait	1
4	Armenia	7	54	Kyrgyztan	1
5	Australia	19	55	Latvia	1
6	Austria	4	56	Lithuania	5
7	Bangladesh	3	57	Macao	1
8	Barbados	1	58	Madagascar	2
a	Belgium	1	50	Malaysia	2
10	Bolivia	15	60	Mali	6
11	Brozil	02	61	Mouritiuo	2
11	Didzii Diitiala Minaia Jalaa	93	01	Maulius	2
12	British Virgin Isles	1	62	Niexico	235
13	Bulgaria	13	63	Moldova	3
14	Canada	247	64	Mongolia	37
15	Chile	38	65	Morocco	4
16	China P.R.	171	66	Namibia	3
17	Colombia	40	67	Nepal	35
18	Costa Rica	1	68	Netherlands	5
19	Croatia	5	69	New Zealand	14
20	Czech Republic	4	70	Nicaragua	1
21	Denmark	2	71	Niger	2
22	Dominican Republic	3	72	Nigeria	10
23	Ecuador	40	73	Norway	17
24	Favot	3	74	Pakistan	11
25	El Salvador	11	75	Panama	2
26	Estonia	4	76	Paraquay	3
27	Fiii	5	77	Poru	65
28	Finland	8	78	Philippines	6
20	France	17	70	Poland	7
20	French Polynesia	3	80	Portugal	5
21	Coordia	5	00	Pomonio	15
20	Georgia	5	01	Rumania	15
32	Germany	34	82	Russia	31
33	Gnana	9	83	Sierra Leone	1
34	Guatemala	21	84	Singapore	24
35	Haiti	1	85	Slovak Republic	4
36	Honduras	5	86	Slovenia	1
37	Hong Kong	32	87	South Africa	11
38	Hungary	5	88	Spain	24
39	Iceland	4	89	Sri Lanka	1
40	India	34	90	Sudan	1
41	Indonesia	6	91	Sweden	16
42	Iran	2	92	Switzerland	8
43	Ireland	1	93	Syria	2
44	Israel	5	94	Taiwan	50
45	Italy	24	95	Tajikistan	1
46	Ivory Coast	2	96	Tanzania	1
47	Jamaica	6	97	Thailand	7
48	Japan	96	98	Tonga	2
49	Jordan	24	99	Turkev	2
50	Kazakhstan	2	100	Uganda	7
				U	

.ist #	Country	Total
101	Ukraine	29
102	United Kingdom	62
103	Uruguay	11
104	Uzbekistan	2
105	Venezuela	22
106	Vietnam	20
107	West Bank	3
108	West Samoa	1
109	Yugoslavia	5
110	Zimbabwe	4
	TOTAL =	2198
1	Vinter Semester 200	5

L

Count of Brigham Young University Students from Each of 110 Nations (Winter Semester, 2005)

List #	Country	Tota
1	Albania	34
2	Argentina	713
3	Armenia	9
4	Asia North	1
5	Australia	186
6	Austria	11
7	Baltic	62
8	Baltic States	10
9	Belgium	137
10	Bolivia	120
11	Brazil	137
12	Bulgaria	62
13	Cambodia	30
14	Canada	357
15	Cape Verde Praia	4
16	Chile	619
17	China Hong Kong	107
18	Colombia	66
19	Costa Rica	56
20	Croatia	37
21	Czech Republic	51
22	Denmark	46
23	Dominican Republic	180
24	Ecuador	212
25	El Salvador	71
26	England	237
27	Fiji	28
28	Finland	41
29	France	199
30	Germany	377
31	Ghana	8
32	Greece	25
33	Guatemala	223
34	Haiti	15
35	Honduras	133
36	Hungary	74
37	India	15
38	Ireland	30
39	Italy	325
40	Ivory Coast	13
41	Jamaica	25
42	Japan	468
43	Kenya	24

st #	Country	Total	
14	Korea	319	
45	Madagascar	28	
16	Mexico	587	
17	Micronesia Guam	19	
18	Mongolia Ulaanbaata	33	
19	Netherlands	22	
50	New Zealand	45	
51	Nicaragua	45	
52	Nigeria	3	
53	Norway	52	
54	Panama	41	
55	Paraguay	106	
56	Peru	222	
57	Philippines	383	
58	Poland	80	
59	Portugal	169	
50	Puerto Rico	68	
51	Romania	80	
52	Russia	436	
53	Samoa	12	
64	Scotland	28	
65	Singapore	26	
66	South Africa	60	
67	Spain	408	
68	Sweden	74	
59	Switzerland	125	
70	Tahiti	12	
71	Taiwan	318	
72	Thailand	86	
73	Tonga	5	
74	Ukraine	170	
75	Uruguay	129	
76	Venezuela	246	
77	West Indies	45	
78	Zimbabwe	10	
	тота	40050	1
	IUIAL =	10252]
		5387	
		15639	

Count of Brigham Young University Students Who Have Served Missions in Various Foreign Nations (Fall Semester, 2004)

ten additional nations

Study 1. Pilot Study of Subjective Judgments

- Purpose: To identify subjective collateral information in Han Chinese personal names.
- Six native Chinese informants provided judgments of (1) gender, (2) location, (3) ethnicity, (4) language/dialect, and (5) religion.
- There were four parts to the electronic questionnaire process.

Part A. Categorization and confidence rating of 269 names.

Part B. Textual explanation of the reasons for the categorizations.

Part C. Ratings of 269 names on scales reflecting basis of judgment.

P



Study 1. Pilot Study Part A. Categorization & Confidence

Gender Identification Accuracy

Six Native Chinese Informants:	Female vs Male
NCI 1	71%
NCI 2	69%
NCI 3	68%
NCI 4	65%
NCI 5	76%
NCI 6	71%

Signal Detection Theory Applied to Gender Identifications

Judged:				F					М						Test of known Female and "Other"									er"							
		1	2	3	4	5	5	4	3	2	1				F	M															
Actual	F	14	1	1	0	0	7	2	2	2	33	62		F	16	46	62	0.26	0.74	1											
	M&b	2	0	0	1	0	8	2	2	3	189	207		М	3	204	207	0.01	0.99	1 82% = percent correct											
		16	1	1	1	0	15	4	4	5	222	269			19	250	269	0.27	1.73	2	2										
	conservative ideal receiver liberal																														
	conservative														receiv	er										libera	1				
frequent	cies:	F	M		F	M		F	M		F	M		F	M		F	M		F	M		F	M		F	M				
Actual	F	14	48	62	15	47	62	16	46	62	16	46	62	16	46	62	23	39	62	25	37	62	27	35	62	29	33	62			
	M&b	2	205	207	2	205	207	2	205	207	3	204	207	3	204	207	11	196	207	13	194	207	15	192	207	18	189	207			
		16	253	269	17	252	269	18	251	269	19	250	269	19	250	269	34	235	269	38	231	269	42	227	269	47	222	269			
probabi	li <u>ties:</u>	F	M		F	M		F	M		F	M		F	M		F	M		F	M		F	M		F	M				
Actual	F	0.23	0.77	1	0.24	0.76	1	0.26	0.74	1	0.26	0.74	1	0.26	0.74	1	0.37	0.63	1	0.4	0.6	1	0.44	0.56	1	0.47	0.53	1			
	M&b	0.01	0.99	1	0.01	0.99	1	0.01	0.99	1	0.01	0.99	1	0.01	0.99	1	0.05	0.95	1	0.06	0.94	1	0.07	0.93	1	0.09	0.91	1			
		0.24	1.76	2	0.25	1.75	2	0.27	1.73	2	0.27	1.73	2	0.27	1.73	2	0.42	1.58	2	0.47	1.53	2	0.51	1.49	2	0.55	1.45	2			
																										,					
	0.75 d'= 0.7 d'= 0.65 d'= 0.65 d'=													0.65	d'=		0.33	d'=		0.25	d'=		0.16	d'=		0.08	d'=				
	2.34 1.59 2.34 1.64 2.34 1.69 2.18 1.53													2.18	1.53		1.62	1.29		1.53	1.29		1.46	1.3		1.36	1.28				

A Signal Detection Theory (TSD) paradigm was used to evaluate the accuracy and the confidence level of native Chinese informants in identifying gender from the 269 names.

The d-prime statistics are stable across confidence boundaries and also similar across the six native Chinese informants.

Study 1. Pilot Study Part A. Categorization & Confidence

Location Identifications

Surprisingly, native Chinese informants were able to identify location

from the names 20% or better, well beyond the chance level.

Six Native Chinese Informants:	Number of Judgments	Percent Correct	Normal Deviate	Probability	Chance of Guess
NCI 1	235	19.6%	z = 4.89	.0000005	five in ten million
NCI 2	229	21.8%	z = 5.97	.000000001	one in a billion
NCI 3	11	54.5%	z = 4.92	.0000004	four in ten million
NCI 4	15	26.7%	z = 2.15	.0157122	1.6 in a hundred
NCI 5	0				
NCI 6	0				

Example Accuracy Matrix for Identification of Location

	t.		st		Judg	jed	st		bu	e	
Actual:	A.Northeas	B.North	C.Northwe	D.East	E.Central	F.South	G.Southwe	ll.Taiwan	III.Hong Ko	IV.Singapo	
A.Northeast	0	9	3	5	1	1	0	1	0	0	0.0%
B.North	6	26	4	13	6	6	3	8	0	0	36.1%
C.Northwest	0	11	2	6	3	2	3	2	0	0	6.9%
D.East	6	5	1	10	4	2	4	8	0	0	25.0%
E.Central	2	9	1	3	2	2	0	1	0	0	10.0%
F.South	0	2	1	1	6	7	2	2	0	0	33.3%
G.Southwest	0	7	0	1	1	2	2	1	0	0	14.3%
II.Taiwan	0	4	0	0	0	1	2	1	1	0	11.1%
III.Hong Kong	0	1	0	1	0	1	0	0	0	0	0.0%
IV.Singapore	0	0	0	0	0	0	0	1	0	0	0.0%
	14	74	12	40	23	24	16	25	1	0	229
	0.0%	35.1%	16.7%	25.0%	8.7%	29.2%	12.5%	4.0%	0.0%	0.0%	21.8%
		north	northwe	east		south					



one in a billion probability by chance

🔎 Microsoft /	Access				
Eile Edit	Study 1 Pilo	t Study Par	t B Toxtual	Evolana	ion for help
U Wenpeil	Study 1. Filo	i Sludy Fai	t D. Textual I	схріана	
Open 🕍	0	Exit			
Tables	Chinese Name Mainland Traditional	Meaning of Given Name/ Comments on Given Name or Surname	Location	F D	Ethnic Langua
Forms	268 张 川田 張		E 2	汉 1	
Pages	269 张 五岳 張		B 2	汉 1	Zh 1
🖉 Macros			A 2	汉 3	Zh 1
Groups	白寿彝		G 2	2	OT 3
	3		B 1	汉 1	Zh 1
	4 蔡 子民		II 2	汉 1	MN 1
	5 岁 国荣		III 2	汉 1	Yu 1
			D 2	₹₩ 1	Zh 1
	Record: II I III * 1	69			>
	A. Name Categorization	US B. Textual Reasons	Example of for sed to obtain to commentary v respect to na properties that	orm extual with me help	
<		V	vith categoriza	ation.	>
Form View					NUM
🛃 start	🔁 🧿 🙆 🥥 🛛 🐻 Microsoft PowerP	🔄 🤷 Part B Rating Ap 🛛 📠 WenpeiPartB	: D 😑 Main Form : Form	🔚 Textual Reason F	🦯 🧨 💯 🔏 🛛 🏷 🏹 N 10:36 PM

Part B of the electronic questionnaire obtained textual commentary on the basis by which names were categorized.

The results of this qualitative aspect of the study were used to create rating scales for quantitative classification of the names in Part C.

FB N Study 1. Pilot Study Part C. Analysis of Rating Scales Chinese Name Familiar Military Heroic Patriotic An <Modern Traditional> <LowHigh> <NoneMuch> <NoneMuch> <NoneMuch> ≤ 1 Mainland Traditional 45 李 慈君 46 李 | 夫| 克 $\bullet \bullet \bullet \bullet \bullet \bullet \bullet \bullet \bullet$ 47 李 荒 $\bullet \bullet \bullet \bullet \bullet \bullet \bullet \bullet \bullet$ $\bullet \bullet \bullet \bullet \bullet \bullet \bullet \bullet$ $\bullet \bullet$ 48 李 际泰 $\bullet \bullet \bullet \bullet \bullet \bullet \bullet \bullet \bullet$ ••• 49 李 梦夫 $\bullet \bullet \bullet \bullet \bullet \bullet \bullet \bullet \bullet$ $\bullet \bullet \bullet \bullet \bullet \bullet \bullet \bullet$ 50 李 人林 $\bullet \bullet \bullet \bullet \bullet \bullet \bullet \bullet \bullet$ ••• 51 李 树元 Example of form Record: > used to obtain ratings of names, to quantify C. Name Rating Scale **D.** Combined Decision the basis on which they are categorized.

-8

🛎 Data Galaxy



Input Files vectors.xml FactorA.xm

FactorB.xml FactorAsp.xml FactorBsp.xml Study 1. Pilot Study Part

Part C. Analysis of Rating Scales

Structure discovery tool displaying possible vector space corresponding to eleven dimensions of onomastic variance.

Editor refresh

🛎 💽 健 🖸

🛃 start

🚳 Microsoft PowerPoint ...

🎑 Part B Rating Applicat... 🛛 🚺 JPad

👹 Data Galaxy

_ 7 🗙

🛎 Data Galaxy



Input Files vectors.xm FactorA.xm FactorB.xm FactorAsp.xml FactorBsp.xml



🛃 start

Plotting of 269 names In hypothetical space of eleven dimensions colored according to category of interest..

Part C. Analysis of Rating Scales

of onomastic variance,

J JPad

A 🖸 🙆 🔾 🔄 Part B Rating Applicat...

Study 1. Pilot Study

😹 Data Galaxy

1 12 ? 🗞 🔁 N 11:05 PM

P

Initial work: develop analytical tools to provide precise comparisons of the accuracy of onomastic categorization.

This kind of precise analysis lends itself well to making cross-language and cross-cultural onomastic comparisons.

One particularly useful analytical tool for these purposes is the Brunswick Lens Model.

The Lens Model of Proximal Cues as the mediators of accurate subjective judgments:



The Lens Model Equation:

$$r_a = GR_e R_s + C\sqrt{1 - R_e^2}\sqrt{1 - R_s^2}$$

• r_a = Correlation of the subjects' judgments with the distal variable •G = Correlation of predicted scores from the two models • R_e = Multiple correlation of the distal variable and cues • R_s = Multiple correlation of subjects' judgments and cues •C = Correlation of the residuals from the two models

Figure 1. Cross Tabulation of the Thirty-Four Most Common Han Names in the 2180 Dataset Crossed with Geo-Location, the Thirty-Two Provinces of China

	northeast													east							S	ou	th	_													
	a1 Heilongjiang	a2 Jilin	a3 Liaoning	b1m Beijing	b2m Tianjin	b3 Hebei	b4 Shandong	b5 Henan	b6 Shanxi	c1 Shaanxi	c2 Gansu	c3ar NeiMonggol	c4ar Ningxia	c5 Qinghai	d1 Zhejiang	d2 Jiangsu	d3m Shanghai	d4 Anhui	e1 Hubei	e2 Hunan	e3 Jiangxi	f1 Guangdong	f2 arGuangxi	f3 Fujian	g1 Sichuan	g2m Changqing	g3 Guizhou	g4 Yunnan	h1 Hainan	i1 Taiwan							
FN3 FN3 FN4 FN11 FN21 FN21 FN21 FN23 FN24 FN25 FN35 FN40 FN40 FN42 FN40 FN42 FN41 FN42 FN42 FN44 FN55 FN61 FN66 FN65 FN74 FN76 FN78 FN78 FN78 FN78 FN80 FN80 FN80 FN80	$\begin{smallmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 $	$ \begin{array}{c} 1 \\ 0 \\ 2 \\ 0 \\ 1 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0$	$\begin{array}{c} 2\\ 2\\ 2\\ 0\\ 0\\ 0\\ 0\\ 0\\ 0\\ 0\\ 0\\ 0\\ 0\\ 0\\ 0\\ 0\\$	$3 \\ 2 \\ 1 \\ 0 \\ 1 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0$	$ 2 \\ 0 \\ 1 \\ 1 \\ 0 \\ $	8502100310002410129334196 1800111	33011161203030001532403223731070	$\begin{array}{c} 2\\ 3\\ 1\\ 0\\ 0\\ 0\\ 1\\ 0\\ 0\\ 0\\ 0\\ 0\\ 0\\ 0\\ 0\\ 0\\ 0\\ 0\\ 0\\ 0\\$	$\begin{array}{c}1\\4\\0\\1\\1\\0\\0\\0\\1\\0\\0\\2\\2\\7\\0\\1\\1\\1\\5\\2\\2\\1\\7\\0\end{array}$			$\begin{smallmatrix} 0 \\ 1 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\$	$ 2 \\ 0 \\ $		$\begin{array}{c} 5 \\ 0 \\ 1 \\ 1 \\ 0 \\ 2 \\ 1 \\ 1 \\ 0 \\ 0 \\ 2 \\ 3 \\ 0 \\ 1 \\ 0 \\ 0 \\ 1 \\ 1 \\ 1 \\ 2 \\ 0 \\ 0 \\ 5 \\ 9 \\ 0 \\ 0 \\ 2 \\ 6 \\ 0 \\ 0 \\ 2 \\ 6 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0$	3308506020082103111312100161100516	$\begin{array}{c}1\\1\\0\\0\\1\\0\\0\\0\\1\\0\\0\\0\\0\\0\\0\\0\\0\\0\\0\\0$	$\begin{array}{c} 2 \\ 1 \\ 1 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0$	$\begin{array}{c} 10\\ 1\\ 0\\ 1\\ 0\\ 0\\ 0\\ 0\\ 0\\ 0\\ 0\\ 0\\ 0\\ 0\\ 0\\ 0\\ 0\\$	$\begin{smallmatrix} 6 \\ 1 \\ 0 \\ 0 \\ 1 \\ 1 \\ 0 \\ 2 \\ 0 \\ 1 \\ 1 \\ 0 \\ 0 \\ 1 \\ 2 \\ 0 \\ 0 \\ 1 \\ 2 \\ 0 \\ 0 \\ 1 \\ 2 \\ 0 \\ 0 \\ 1 \\ 2 \\ 0 \\ 0 \\ 1 \\ 2 \\ 0 \\ 0 \\ 1 \\ 2 \\ 0 \\ 0 \\ 1 \\ 2 \\ 0 \\ 0 \\ 1 \\ 2 \\ 0 \\ 0 \\ 1 \\ 2 \\ 0 \\ 0 \\ 1 \\ 2 \\ 0 \\ 0 \\ 1 \\ 2 \\ 0 \\ 0 \\ 1 \\ 2 \\ 0 \\ 0 \\ 1 \\ 2 \\ 0 \\ 0 \\ 1 \\ 2 \\ 0 \\ 0 \\ 1 \\ 0 \\ 0 \\ 1 \\ 0 \\ 0 \\ 1 \\ 0 \\ 0$	$\begin{smallmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 $	$\begin{array}{c} 12\\ 3\\ 0\\ 0\\ 1\\ 0\\ 0\\ 0\\ 0\\ 0\\ 0\\ 0\\ 0\\ 0\\ 0\\ 0\\ 0\\ 0\\$	$ \begin{array}{c} 0 \\ 2 \\ 1 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0$	51010101013400002010100210010	$\begin{array}{c} 11 \\ 2 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0$		$ \begin{array}{c} 0\\ 0\\ 0\\ 0\\ 0\\ 0\\ 0\\ 0\\ 0\\ 0\\ 0\\ 0\\ 0\\ $		$\begin{array}{c} 1 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\$		$\begin{array}{c} 86\\ 34\\ 13\\ 18\\ 28\\ 15\\ 10\\ 24\\ 6\\ 5\\ 10\\ 35\\ 20\\ 10\\ 17\\ 5\\ 12\\ 11\\ 10\\ 13\\ 47\\ 18\\ 22\\ 22\\ 101\\ 66\\ 10\\ 63\\ 63\\ 6\end{array}$						
FN84 高 FN85 黄 FN87 龚	0 1 1	3 1 0	13 2 0	3 1 0	1 3 0	9 6 1	16 3 0	3 2 0	0 4 0	13 3 1	0 1 0	0 0 0	0 1 0	0 0 0	6 9 1	6 10 3	1 5 3	2 6 3	0 16 1	2 8 2	0 21 0	1 24 1	0 11 0	2 22 2	0 9 1	0 1 0	0 1 1	1 3 0	0 1 0	0 2 0	82 177 21						

Figure 2. Geo-Location of the Han Names: The Thirty-two Provinces of China Grouped into Nine Regions



Figure 3. Metrika Vector Plot of Thirty Chinese Provinces in the Anthroponomastic Space of the Thirty-Two Most Common of the 2180 Names, Colored According to Region



Figure 4. Metrika Datapoint Plot of the Thirty-Four Most Common Names, Plotted in the Same Anthroponomastic Space as Figure 3, with Three Names in Extreme Positions Labeled



Conclusions

- Gender: Chinese given names less informative than American given names.
- Geo-Location: Chinese family names well beyond chance level (contra experts' view).
- Similar studies are being conducted for six national groups: India, Korea, Jordan, Russia, Singapore, and Nepal.