



National Center for Health Statistics

Data Linkage

Evaluating Hispanic and Asian Name Algorithms as a Tool for Improving Linkage to the National Death Index

Dean H. Judson,
Eric Miller,
Jennifer D. Parker,
Hannah Day
National Center for Health Statistics (NCHS)

2013 Federal Committee on Statistical Methodology

Disclaimer: These analyses were performed on an “almost final” file; a small handful of cases still need final resolution.

Outline

- Describe goals and process
- Describe innovations
- Focus on name algorithms/SSN4's
- Describe operational results and evaluate effects
- Think about future evaluations



NCHS Linked Mortality Files

- Last released in 2009 (up to NDI data year 2006), linkage of survey respondents from NCHS surveys (now ~2.5M records) to the National Death Index (~70M records) to ascertain mortality status & causes of death
- Auxiliary data from the Social Security Death Master File (DMF) also used to ascertain status
- Used for passive follow-up epidemiological studies



Process

- Construct submission files
 - Up to ~50 alternative records per person accounting for:
 - Alternative names (nicknames, middle/first flip)
 - Alternative dates of birth (esp. with imputations)
 - Alternative identifying numbers
 - Unique Hispanic and Asian naming conventions
- Submit to NDI linkage
- Potential links returned



Process (continued)

- Reweight and rescore
 - Construct detailed weights based on name, dates, digits, and other characteristics' frequencies
 - $Wt_i = \log_2(1/\text{Proportion})$, prorate for inexact matches
 - Sum weights and select “best” (highest scoring) pair
 - Tentatively flag records above threshold
- Rule-based post-processing
 - Alternative to clerical review
- Clerical review



Innovations this time

- Hispanic and Asian name algorithms
- Hispanic nickname tables
- New use of phonetic algorithms:
 - NYSIIS, Soundex, Jaro-Winkler score, Edit distance, “within”
- Use of enhanced matching macro (citation: Dean Resnick) to search the DMF
- Enhanced weights
- Rule-based post-processing
- Expedited/Computerized clerical review
- 4-digit SSN NDI search (completely new)



Example of potential for linkage error - Hispanic Paradox

- Hispanics have lower mortality rates compared to non-Hispanic whites
 - Health selective immigration
 - Salmon bias (return migration)
 - Advantageous health behaviors and social support
 - **Data quality / Insufficient linkage**
 - Naming conventions for Hispanics differ from other US populations
 - Use of mother's and father's surname
 - May not have single middle name
 - Less likely to have social security number
 - Especially among older adults and foreign born



Methods – Name Clean-up

- Accounting for Hispanic and Asian naming conventions
 - Hispanic
 - Hispanic nickname lookup table
 - switch middle and last
 - Asian
 - switch first and last



Hispanic Nickname Table

Sex	Formal Name	Nicknames						
F	<i>Adelina</i>	<i>Deli</i>	<i>Lina</i>					
F	<i>Adelaida</i>	<i>Ade</i>	<i>Adela</i>					
M	<i>Adrián</i>	<i>Adri</i>						
F	<i>Adriana</i>	<i>Adri</i>						
M	Alberto	<i>Alber</i>	<i>Albertito</i>	Beto	Berto	Tico	Tuco	Tito
M	<i>Alejandro</i>	<i>Ale</i>	<i>Álex</i>	Alejo	Jandro	Jano	Sandro	
F	<i>Alejandra</i>	<i>Sandra</i>	<i>Ale</i>	Álex	Aleja	Jandra	Jana	
M	<i>Alfonso</i>	<i>Alfon</i>	<i>Fon</i>	Fonso	Fonsi	Poncho		
F	<i>Alicia</i>	<i>Ali</i>	<i>Licha</i>					



NDI Selection Algorithm

NDI Selection Criteria

	Criterion one	Criterion two	Criterion three
Social Security Number	X		X
First name	X	X	
Middle initial			
Last name	X	X	
Day of birth		X	X
Month of birth		X	X
Year of birth		X	X
Sex			X
Father's surname			
Race			
State of residence			
State of birth			
Marital Status			



Basic Tabulations

-> eligible = NO

Rejected by	NDI	Freq.	Percent	Cum.
	NO	169	0.12	0.12
	YES	143,835	99.88	100.00
	Total	144,004	100.00	

-> eligible = YES

Rejected by	NDI	Freq.	Percent	Cum.
	NO	2,478,065	99.09	99.09
	YES	22,862	0.91	100.00
	Total	2,500,927	100.00	



Eligible Record Dispositions

Final disposition class, among eligibles	Final NDI Status		Total
	ALIVE	DEAD	
Assumed alive	2,087,704	0	2,087,704
	100.00	0.00	84.25
Dead by score/class	0	382,334	382,334
	0.00	97.94	15.43
Dead by rule	0	3,327	3,327
	0.00	0.85	0.13
Dead by manual review	0	1,848	1,848
	0.00	0.47	0.07
Dead by mult. sources	0	2,852	2,852
	0.00	0.73	0.12
Total	2,087,704	390,361	2,478,065
	100.00	100.00	100.00



Rescore

- Five “classes”
 - One: Full hit on 8 or 9 fields, incl. SSN
 - Two: SSN present almost exact
 - Three: SSN unknown, 8 or 9 other fields
 - Four: SSN unknown, <8 other fields
 - Five: SSN < 8 digits, not in classes one-four
- Scores = sum of fields weights with proration for partial matches
- Thresholds by class (NHEFS)



Eligible Record Dispositions by Class

Final disposition class, among eligibles	Rescore Class						Total
	1	2	3	4	5	.	
Assumed alive	4 0.00	1,228 1.32	23,083 24.00	802,162 98.00	474,065 99.73	787,162 100.00	2,087,704 84.25
Dead by score/class	207,690 100.00	90,392 97.01	71,032 73.87	13,220 1.62	0 0.00	0 0.00	382,334 15.43
Dead by rule	0 0.00	966 1.04	1,265 1.32	760 0.09	336 0.07	0 0.00	3,327 0.13
Dead by manual review	0 0.00	137 0.15	543 0.56	1,159 0.14	9 0.00	0 0.00	1,848 0.07
Dead by mult. sources	0 0.00	453 0.49	238 0.25	1,225 0.15	936 0.20	0 0.00	2,852 0.12
Total	207,694 100.00	93,176 100.00	96,161 100.00	818,526 100.00	475,346 100.00	787,162 100.00	2,478,065 100.00



Innovations

- Rules
 - E.g. Rule 2: if social security matches on 8 or 9, last name matches, sex matches, exact DOB, and birth state matches or is missing, then match=yes (dead).
- Collectively referred to as “autodead”, “autoalive”, “autoreview”
- Clerical review screens to follow



Dispositions by “Auto” Rules

Final disposition class, among eligibles	Auto-dead rule fired			Total
	N	Y	.	
Assumed alive	1,299,062	1,480	787,162	2,087,704
	99.56	0.38	100.00	84.25
Dead by score/class	979	381,355	0	382,334
	0.08	98.76	0.00	15.43
Dead by rule	0	3,327	0	3,327
	0.00	0.86	0.00	0.13
Dead by manual review	1,848	0	0	1,848
	0.14	0.00	0.00	0.07
Dead by mult. sources	2,852	0	0	2,852
	0.22	0.00	0.00	0.12
Total	1,304,741	386,162	787,162	2,478,065
	100.00	100.00	100.00	100.00



Dispositions by “Auto” Rules

Final disposition class, among eligibles	Auto-alive rule fired			Total
	0	1	.	
Assumed alive	797,632 67.15	502,910 99.95	787,162 100.00	2,087,704 84.25
Dead by score/class	382,331 32.19	3 0.00	0 0.00	382,334 15.43
Dead by rule	3,327 0.28	0 0.00	0 0.00	3,327 0.13
Dead by manual review	1,828 0.15	20 0.00	0 0.00	1,848 0.07
Dead by mult. sources	2,631 0.22	221 0.04	0 0.00	2,852 0.12
Total	1,187,749 100.00	503,154 100.00	787,162 100.00	2,478,065 100.00



Submission Record DMF Record String Comparators

Submission



NDI

	Name	Intl.	Match	Name	Intl.	
First:	D	D	X	Dean	D	Exact Match: <input type="checkbox"/>
Middle:	Harold	H	I	H	H	
Last:	Judson		X	Judson	FSN-LN	<input type="text"/>
Surname:				Judson	LN-FSN	<input type="text"/>

Submission	Matching Digits:	NDI	Sub.	Match	NDI
9-Digit: XXXXX2654	-----XXXX	4	XXXXXX2654	Sex: M	X M
4-Digit: 2654		2654	ST Birth: CT		MI
			ST Res: MD	X	MD
			Race: White	X	White
			Marital: Married	?	Unknown
			Dth Age:	-	48
			P_Dead_2011: 0.131400		ST Dth: VA

	Birth Date		NDI Death Date	
	Sub.	Match	NDI	Date
Month:	07	X	07	12
Day:	20	X	20	24
Year:	1962	-1	1963	2011

Last Alive Year: 2008	Linkage ID: 10000000000001	Reviewer: VVV6	Date: <input type="text"/>	Review: <input type="button" value="Find"/>
Rec. # 1156	SCORE: 45	Review Notes: <input type="text"/>		<input type="button" value="Save"/>
CLASS: 2				<input type="button" value="Back"/> <input type="button" value="Next"/>
Reason: DEAN'S SPECIAL PII				<input checked="" type="checkbox"/> DMF Match <input type="button" value="Quit"/>

Submission Record DMF Record String Comparators

	DMF Name	Match	Submission Name	Match	NDI Name
First:	Harry	<input type="checkbox"/>	D	<input checked="" type="checkbox"/>	Dean
Middle:	D	<input type="checkbox"/>	Harold	<input checked="" type="checkbox"/>	H
Last:	Judson	<input checked="" type="checkbox"/>	Judson	<input checked="" type="checkbox"/>	Judson
Suffix:					

Verify/Proof:

DMF SSN	Matching Digits	Submission	Matching Digits	NDI SSN		
000012654	-----XXXX	4	XXXXX2654	-----XXXX	4	XXXXXX2654

Birth Date			Death Date			
DMF	Match	Sub.	DMF	Match	NDI	
Month:	07	<input checked="" type="checkbox"/>	07	12	<input checked="" type="checkbox"/>	12
Day:	28	<input type="checkbox"/>	20	01	<input type="checkbox"/>	24
Year:	1988	<input checked="" type="checkbox"/>	+26	2011	<input checked="" type="checkbox"/>	2011

State Birth:	CT	MI
State Res:	MD	MD
DMF State Country Code:	<input type="text"/>	
Zip Last Res:	MD	
Zip Lump:	<input type="text"/>	

Last Alive Year: Linkage ID:

Rec. #: SCORE: CLASS:

Reason:

Reviewer: Date:

Review Notes:

Review:

- No Match
- Unsure
- Match
- DMF Match

Buttons: Find, Save, Back, Next, Quit

Some Operational Results

- Interrater reliability (clerical review)
- Overall link rates by race/ethnicity
- 4-digit SSN vs. 9-digit SSN match rate

Interrater Reliability Table

	NO MATCH	MATCH	TOTAL
NO MATCH	108	11	119
	49.1	5.0	54.09
	90.8	9.2	-
	83.7	12.1	-
MATCH	18	79	97
	8.2	35.9	44.09
	18.6	81.4	-
	14.0	86.8	-
TOTAL	126	90	216
	58.64	41.36	100

Cohen's Kappa: 0.727



National Center for Health Statistics

Data Linkage

Disposition by Race/Ethnicity

Final disposition class, among eligibles	Race/Ethnicity						Other	Total
	White-Non	Black-Non	API-Non	AIAN-Non	Hispanic			
Assumed alive	1,306,474 81.11	306,850 85.30	66,917 93.24	13,488 85.46	313,848 93.86	80,127 93.48	2,087,704 84.25	
Dead by score/class	298,531 18.53	51,760 14.39	4,755 6.63	2,250 14.26	19,713 5.90	5,325 6.21	382,334 15.43	
Dead by rule	2,690 0.17	352 0.10	27 0.04	25 0.16	170 0.05	63 0.07	3,327 0.13	
Dead by manual review	1,209 0.08	305 0.08	26 0.04	9 0.06	229 0.07	70 0.08	1,848 0.07	
Dead by mult. sources	1,789 0.11	458 0.13	42 0.06	11 0.07	424 0.13	128 0.15	2,852 0.12	
Total	1,610,693 100.00	359,725 100.00	71,767 100.00	15,783 100.00	334,384 100.00	85,713 100.00	2,478,065 100.00	



Nine-digit SSNs

Final disposition class, among eligibles	Full 9-digit SSN present		Total
	NO	YES	
Assumed alive	1,488,698 94.74	599,006 66.07	2,087,704 84.25
Dead by score/class	77,896 4.96	304,438 33.58	382,334 15.43
Dead by rule	1,923 0.12	1,404 0.15	3,327 0.13
Dead by manual review	1,569 0.10	279 0.03	1,848 0.07
Dead by mult. sources	1,315 0.08	1,537 0.17	2,852 0.12
Total	1,571,401 100.00	906,664 100.00	2,478,065 100.00



Four-digit SSNs, 2007+

Final disposition class, among eligibles	last four SSN digits only		Total
	NO	YES	
Assumed alive	2,064,166	23,538	2,087,704
	84.15	93.86	84.25
Dead by score/class	380,806	1,528	382,334
	15.52	6.09	15.43
Dead by rule	3,319	8	3,327
	0.14	0.03	0.13
Dead by manual review	1,843	5	1,848
	0.08	0.02	0.07
Dead by mult. sources	2,852	0	2,852
	0.12	0.00	0.12
Total	2,452,986	25,079	2,478,065
	100.00	100.00	100.00

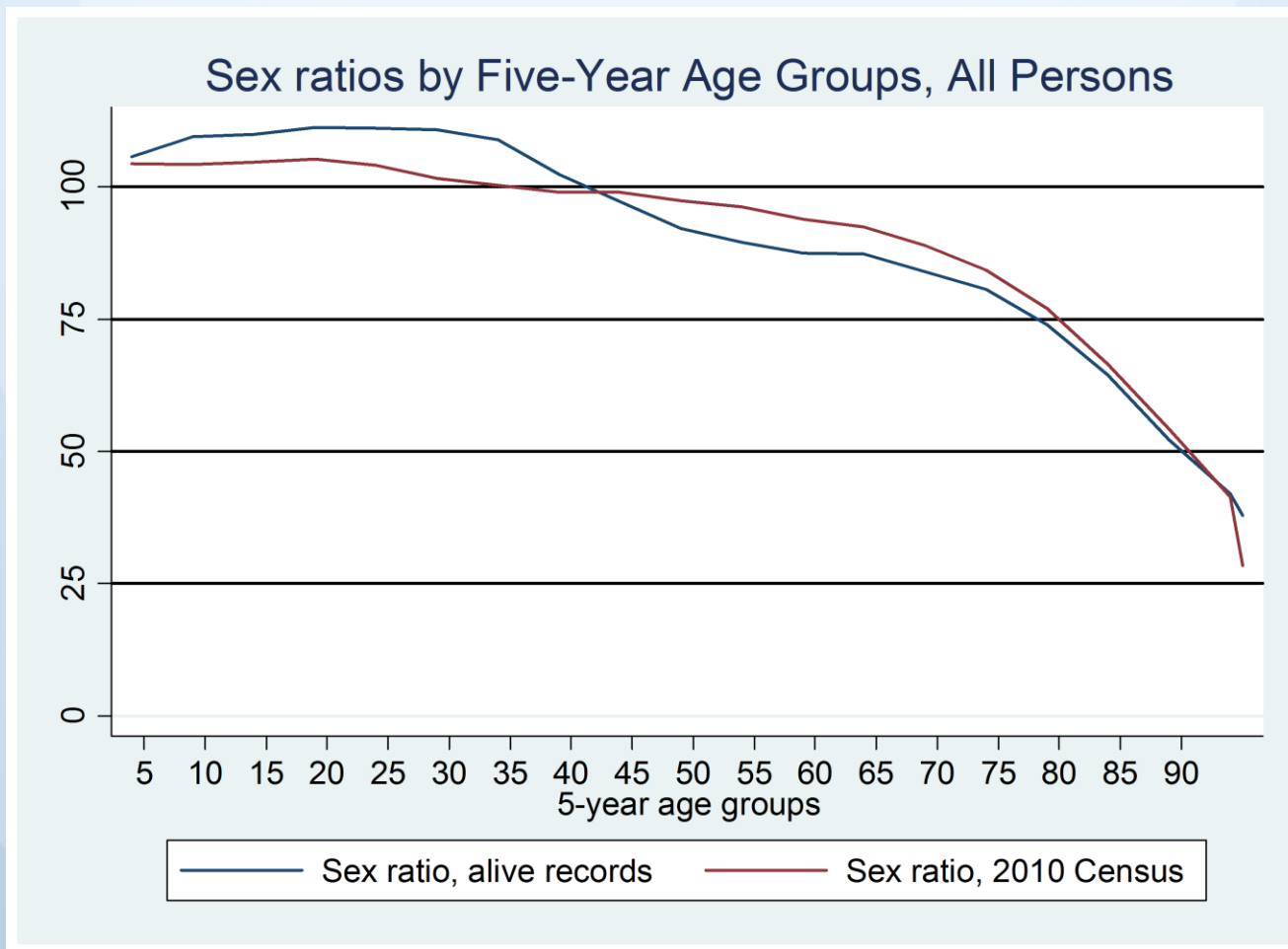


9-digit SSNs, 2007+, NHANES

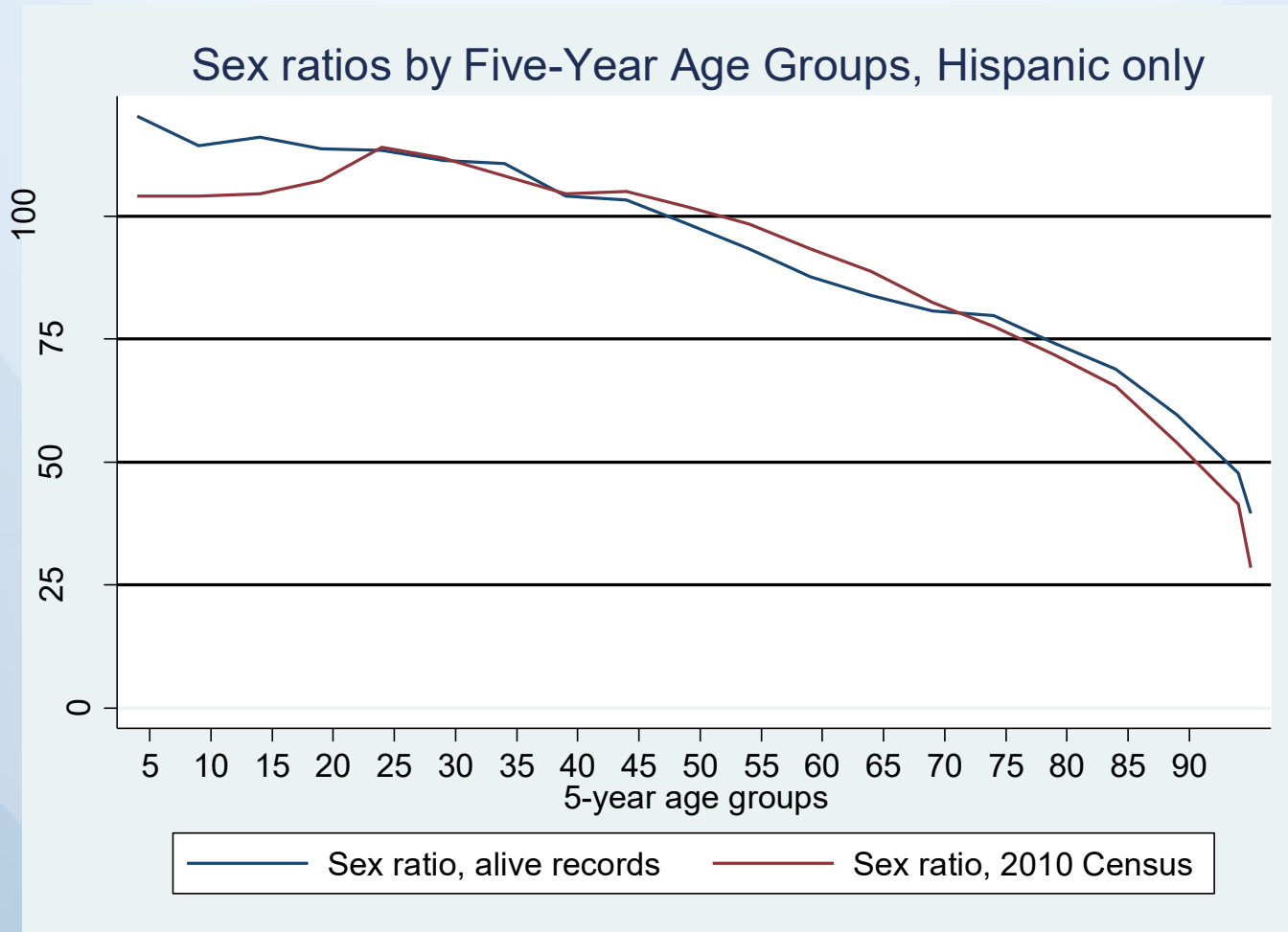
Final disposition class, among eligibles	9-digit SSN, 2007+		Total
	NO	YES	
Assumed alive	2,080,957 84.22	6,747 94.59	2,087,704 84.25
Dead by score/class	381,949 15.46	385 5.40	382,334 15.43
Dead by rule	3,326 0.13	1 0.01	3,327 0.13
Dead by manual review	1,848 0.07	0 0.00	1,848 0.07
Dead by mult. sources	2,852 0.12	0 0.00	2,852 0.12
Total	2,470,932 100.00	7,133 100.00	2,478,065 100.00



Demographic Analysis



Demographic Analysis, Hispanics



Conclusions and Future Research

- Limited success with Hispanic and Asian name algorithms
- Some success with four-digit SSN algorithm
- Good success with demographics
- More mysteries, more evaluation!
- Future: Latent class modeling, error bars, etc.

