



Linking Maternal and Child Health Data to Create a Comprehensive Longitudinal Dataset: The Florida Experience

Hamisu Salihu, Ph.D., Principal Investigator
University of South Florida College of Public Health
R01 HS1997-01
9/30/2010 – 9/29/2013 (3 years)

April 6, 2011

Specific Aims

1. To create an **expanded clinically enhanced maternal-infant dataset** for the State of Florida by augmenting the current statewide hospitalization data files through linkages to other data sources.
2. To **validate the created dataset** in Specific Aim 1 through a rigorous process that will establish confidence in the use of the dataset.
3. To demonstrate the utility of the newly created, enriched dataset in conducting comparative effectiveness analysis using early term elective delivery as a case study.



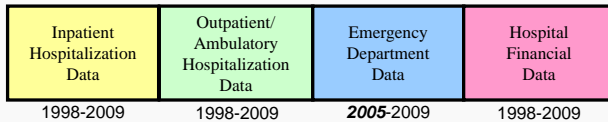
2

Data Sources

FDOH (Office of Vital Statistics) Data

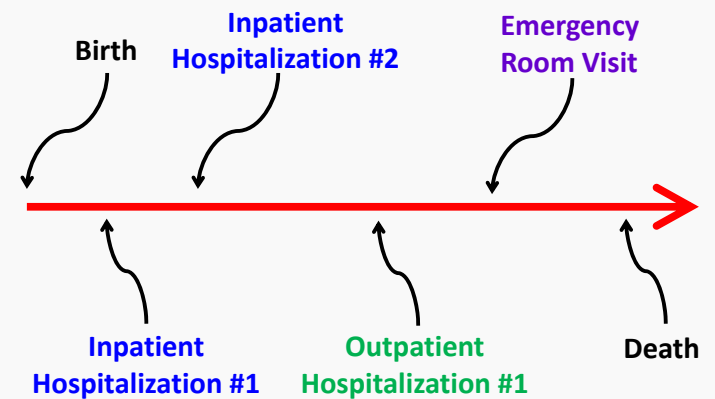


Agency for Health Care Administration Data



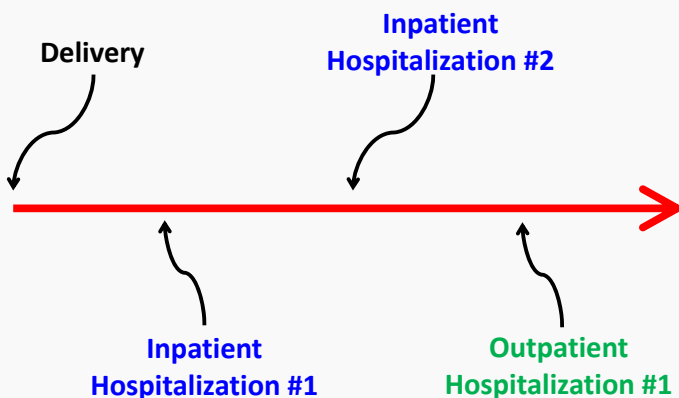
3

"Follow" Infants Over Time Through Linkage



4

"Follow" Moms Over Time Through Linkage



5

Special Challenges to Our Data Linkage

- **Birth vital records** contain a significant amount of identifying information
- **Hospital records** (inpatient, ambulatory, ED) for the infant contain limited identifying information
 - ✓ No infant SSN, name, address
 - ✓ Primary identifier is mother's SSN (INFANTLINK), but it is missing >10% and disproportionately among certain subgroups
 - ✓ Previous investigation reveal that maternal SSN has a typo or transposition in over 1,000 instances (**ASSUME** identifiers have errors)
 - ✓ Missing mother's date of birth, a key linking and/or confirmatory variable



6

Our Approach to Linking AHCA to VS

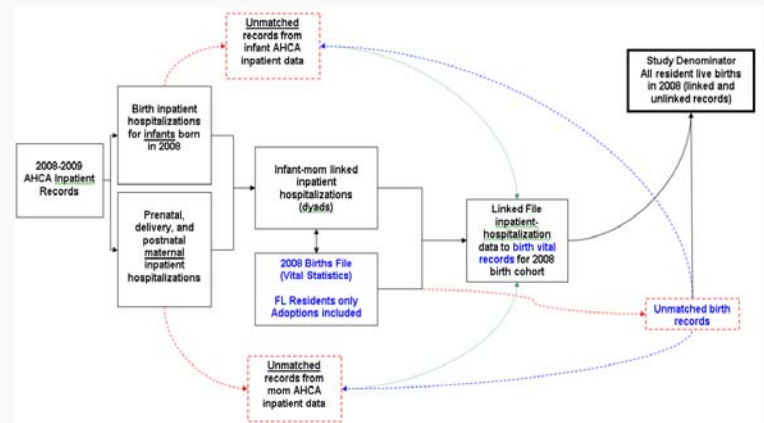
- Stage I**
 - Within the inpatient hospital discharge data, we first attempt to link infants to their mothers (so called **dyad** links) with the primary goal of obtaining **maternal DOB**, an important linking variable (**FIND** other identifying information)
- Stage II**
 - Link these dyad pairs to birth vital records, now incorporating infant's and mom's DOB, mom's SSN, and facility of birth as the primary linking variables
- Stage III**
 - Attempt to link infant and mom hospitalizations that did not link to a maternal record from Stage #1 directly to the birth record



7

Example of Overarching Linkage Approach

Live birth records to inpatient birth hospitalizations in 2008



8

Software To Facilitate Data Linkage

- LinkSolv
- AutoMatch
- LinkageWiz
- FRIL
- LinkPlus
- Link King
- SQL Match
- FEBRL
- SQL Server (SSIS)
- SAS**, SPSS, Stata, S-Plus, R
- ...many more!



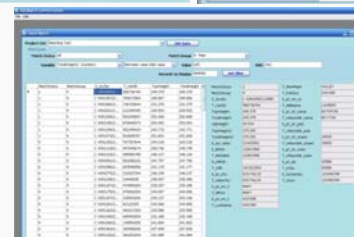
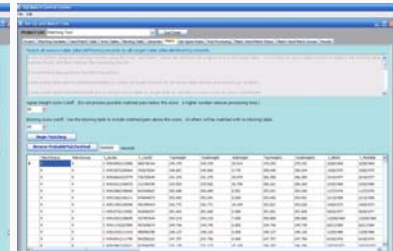
9

SQL Match

Set up data linkage



Linkage Results



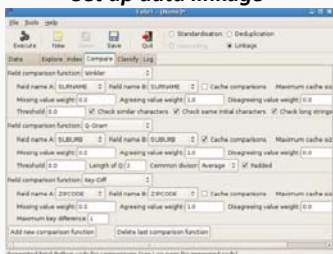
Manual Review



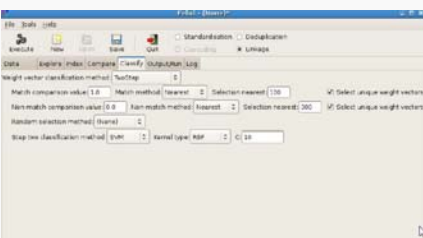
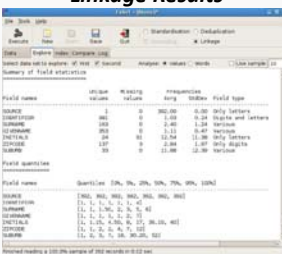
10

Freely Extensible Biomedical Record Linkage

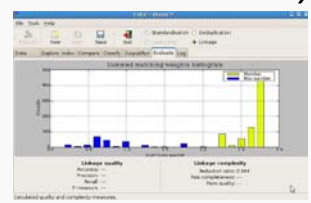
Set up data linkage



Linkage Results



Manual Review and Summary



11

Link Plus

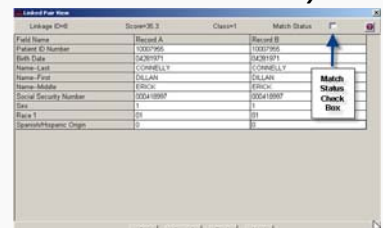
Set up data linkage



Linkage Results



Manual Review and Summary



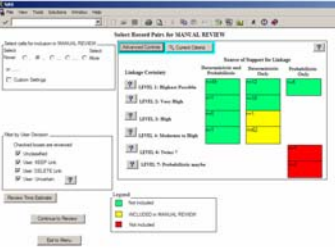
12

Link King

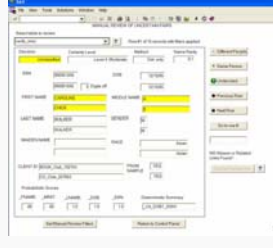
Set up data linkage



Select Variables

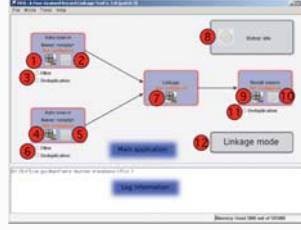


Manual Review

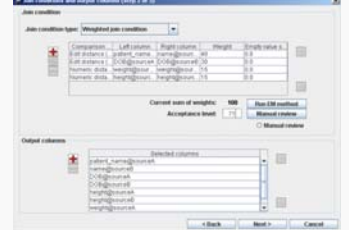


Fine-Grained Record Linkage (FRIL)

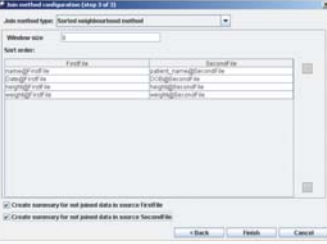
Set up data linkage



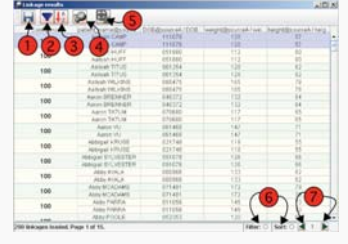
Select Variables and Weights



Join Method (Blocking, SNM)



Manual Review



But our choice...SAS

```

432,434,435,436,437,438,439 (abortions)
proc; 4901,4901,7491,750 (abortions)

flagM_female= upcase(SEXSEX) in: ("F", "M");
flagM_v27=0; flagM_450_459=0; flagM_dry=0; flagM_proc=0;

*Algorithm for identifying live birth hospitalizations for infants is much more simpler:
Flag 1: Must be born in 2008 (establishing a 2008 cohort)
Flag 2: V3 code (includes singletons and multiples)
Flag 3: Does the date of birth equal the date of admissions (may not work for home births)
Flag 4: Type of admission = 4, indicative of newborn admission or initial admission w/ 24 hours of birth

flagB1_yob = (year(BIRTHDATE)+2008);
flagB2_V3 = 0;
flagB3_dobAdmit = (ADMITDATE=BIRTHDATE);
flagB4_admType = (ADM_P8100="4");

*Also creating an exclusion flag for mom records (see EXCLUSIONS: above);
flagMexc1=0;

*Using diagnosis, procedure, and DRG codes to apply select "flags";

do i=1 to 31;
  if _dcl(i) = "427" then flagM1_v27=1;
  if "450" <= _dcl(i) < "4509" then flagM1_450_459=1;
  if _px(i) in ("720","721","722","723","724","725","726","727","728","729","730","731","732","733","734","735","736","737","738","739","740","741","742","743","744","745") then flagM1_proc=1;
  if "430" <= _dcl(i) < "440" then flagMexc1=1;
  if _px(i) in ("4901","4902","7491","750") then flagMexc1=1;
  if _dcl(i) = "493" then flagM1_V3=1;
end;
drop i;

if DRG in ("745","746","747","748","774","775") then flagM1_dry=1;
    
```

Linking Mechanics

- Developed a SAS macro
- Hierarchical, stepwise** series of linking stages, using various combinations of variables, proceeding from highest to lowest confidence
 - Exact and partial matching, **linking with replacement**
 - Primarily **deterministic**, includes **probabilistic** elements
 - CREATES** potential matches
- Coding algorithm to calculate a "linking confidence" score to **GRADE** matches
 - Also incorporate a "delivery confidence" score
- Records above a certain score are **SELECTED** as links, borderline scores require **manual validation**
 - May find false + we need to **CORRECT**
 - We try to minimize manual review



Linking Mechanics

- We do not use **blocking**
 - Too concerned about flawed data
 - Linking approximately 230,000 birth hospitalization records to approximately 1.4 million "women" records using the merging macro takes approximately 1 hour
 - Will sacrifice extra time for greater sensitivity
- SAS
 - Not as automated or "point-and-click" as other software
 - Extremely **customizable** through coding
 - Easy to incorporate a large number of variables (Link King)
 - Easy to allow "**crossover**" links
 - Mom's SSN in AHCA links to father's SSN in vital stats
 - Can process extremely large datasets quickly given powerful computers

```

***** LINK STAGE 57 *****
***** GLOBAL CONDITION *****
GLOBAL CONDITION ----> |
*****
Number of records linked in this stage (prior to unduplication) ----> 202,574
Number of NEW records linked in this stage ----> 3,204
Number of records linked OVERALL (unduplicated) ----> 207,516
*****

WARNING: Data too long for column 'linkFlagM1'; truncated to 100 characters to fit.

***** COLUMN 6 C ***** COLUMN 10 ***** COLUMN 14 ***** COLUMN 23 *****
DATASET WITH ONE RECORD PER SUBJECT ----> MORN2.LINK
Variable serving as unique identifier ----> sys_recid
Suffix to be added to all variables ----> n
Total # of records ----> 129141
Total # of variables ----> 49

DATASET WITH MULTIPLE RECORDS PER SUBJECT ----> INF2.LINK
Variable serving as unique identifier ----> sys_recid
Suffix to be added to all variables ----> i
Total # of records ----> 991409
Total # of variables ----> 57

TOTAL NUMBER OF LINKING 'PASSES' ----> 57
FROM STAGE 1 TO 57

DATASET WITH LINKED RECORDS ----> LINKED
Total # of LINKED records ----> 297016
(unduplicated linkage pairs, false positive links may exist)

UNLINKED RECORDS FROM THE ONE DATASET ----> UNLINKED_MOMS
Total # of UNLINKED "mom" records ----> 193769

UNLINKED RECORDS FROM THE MANY DATASET ----> UNLINKED_INFANTS
Total # of UNLINKED "many" records ----> 3333

NOTE: Global condition applied to all passes ----> |
    
```

Additional Challenges

- **Disentangling multiples** (twins, triplets, etc)
 - ✓ No infant SSN, no names in hospital data
 - ✓ Multiples will share all mom characteristics
 - ✓ Ordering of variables in AHCA does not match birth order
 - ✓ Can use **sex** to differentiate between opposite-sex dizygotic twins
 - ✓ Can use diagnosis codes that reflect 500 gram **birth weight categories** to disentangle same sex multiples that may differ in birth weight
 - ✓ For multiples that have the same sex, similar birth weights, it may be impossible to determine, given the available data, which hospital record goes with which birth record
 - Investigating other options
 - Random assignment
 - Allocation to “family” as unit



19

Next Steps

- Finalize enhancements to linkage of birth record to birth hospitalizations
- Link in **post-birth hospitalizations**, ambulatory records, and ED data
 - ✓ Challenging for those with missing/incorrect maternal SSN
- Develop an identifier crosswalk to link in **cost-to-charge ratios (CCRs)** from CMS to convert hospital charges to costs



20