

IPUMS

Multigenerational Longitudinal Panel

Life histories for the U.S. population, 1850-1940

- Censuses
- Social Security
- Military records (draft, enlistment)
- Vital records (birth, death, marriage, divorce)

IPUMS

Multigenerational Longitudinal Panel

Life histories for the U.S. population, 1850-1940

- Impact of early life conditions on later health and well-being
- Social, Economic, Geographic Mobility
- Life course transitions

IPUMS

Multigenerational Longitudinal Panel

Life histories for the U.S. population, 1850-1940

Link across 5+ generations

- Impact of forebears on health and well-being
- Socioeconomic mobility across generations:
Do we have dynasties?

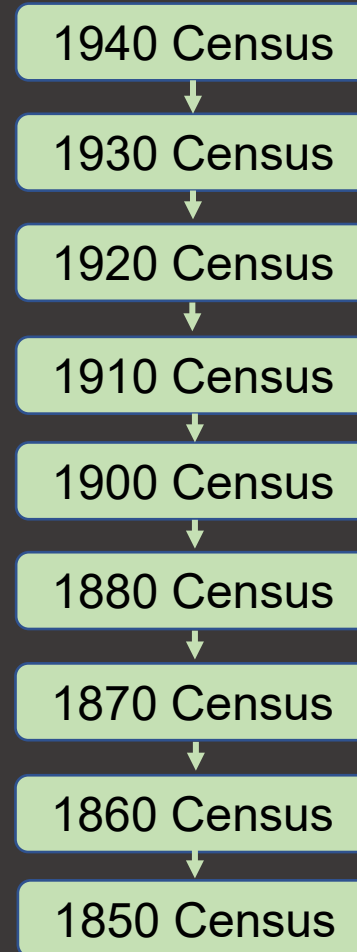
IPUMS

Multigenerational Longitudinal Panel

Life histories for the U.S. population, 1850-1940

Understanding the great transformations:
demographic transition, family transition,
urbanization, immigration, industrialization

IPUMS Multigenerational Longitudinal Panel



IPUMS MLP

- Availability of complete-count data makes it feasible to link most people with extremely low Type I errors, provided that we use all information available, including characteristics of others in the household, neighbors, and location
- This can introduce selection bias, but this bias turns out to no more severe than other linkage methods.



Census Linking Project

The Census Linking Project offers researchers the ability to create longitudinal datasets using historical US Census data (1850-1940). We provide links between each pair of complete-count Censuses using a wide variety of linking algorithms.

church street	450	Smith, Jane	mother
		— Janet	daughter
		— Jack	son
		— Alfred	son
		— Cindy	daughter
		— Michael	son
church street	420	Smith, Sharon	mother
		— Jackie	daughter
		— Tom	son
		— Alfred	son
		— Jackie	daughter
		— Michael	son

Get the Data

[Download the crosswalk files](#) →

Learn more

[Watch our video series](#) →

MLP vs. CLP

- As expected, MLP has far lower Type I error rate (false positives) than CLP
- MLP also has lower Type II error rates than CLP (higher linkage rate)
- More surprisingly, MLP and CLP have similar overall selection bias (representativeness)
- MLP will soon be much easier to use than CLP as it is integrated into the IPUMS data access system

Links and legibility: Making sense of historical U.S. Census automated linking methods*

Arkadev Ghosh, Sam Il Myoung Hwang, Munir Squires

University of British Columbia

July 1, 2021

Abstract

This paper explores the effect of handwriting legibility on the performance of algorithms that link individuals across census rounds. We propose a measure of legibility which we implement at scale for the 1940 US Census, and find strikingly wide variation in enumerator-level legibility. Using boundary discontinuities in enumeration districts, we estimate the causal effect of low legibility on the performance of a set of popular automated linking algorithms. We show that one algorithm out-performs the rest across the spectrum of high to low legibility, and find that it provides a better measure of 10-year interstate migration.

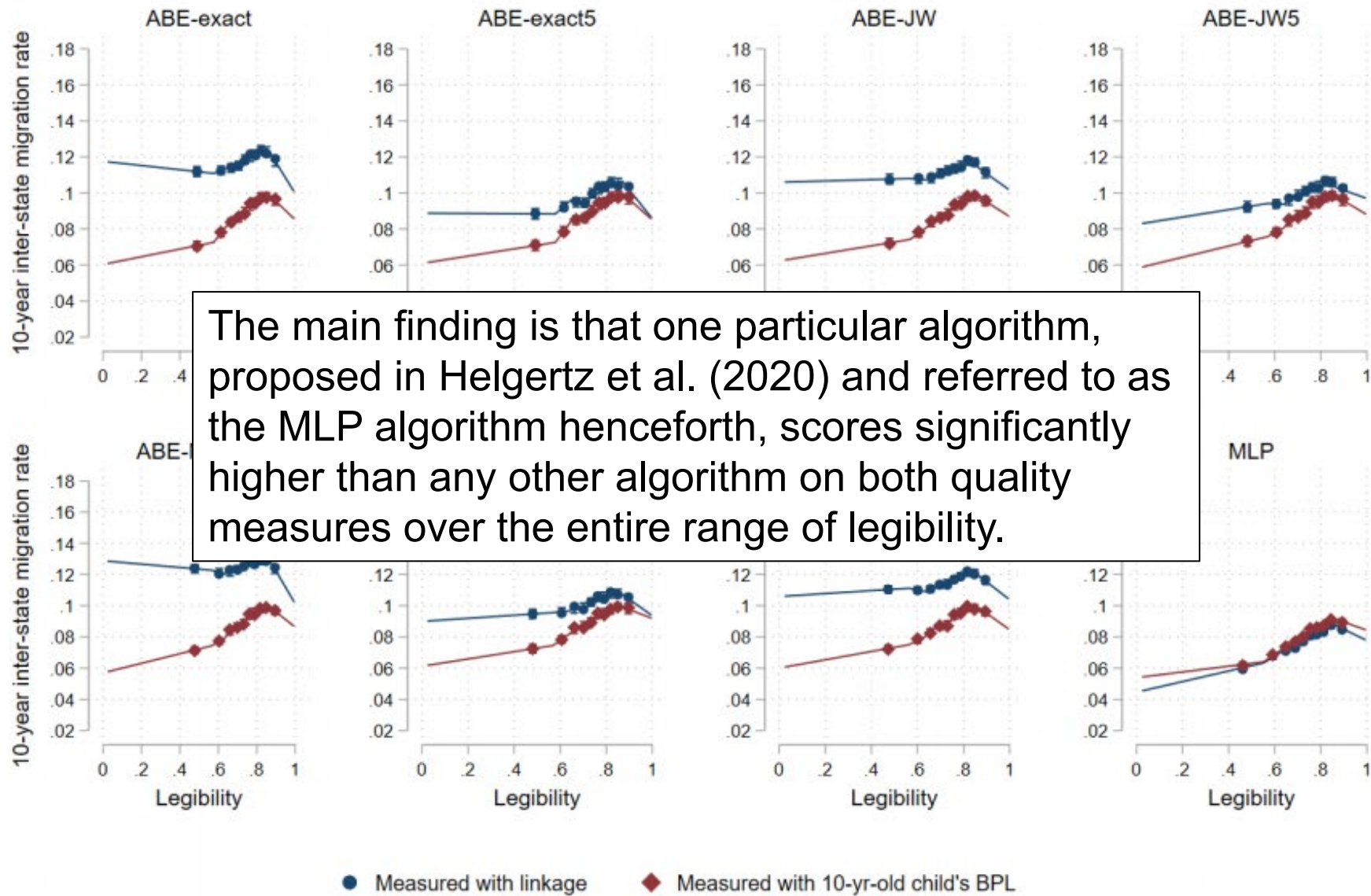
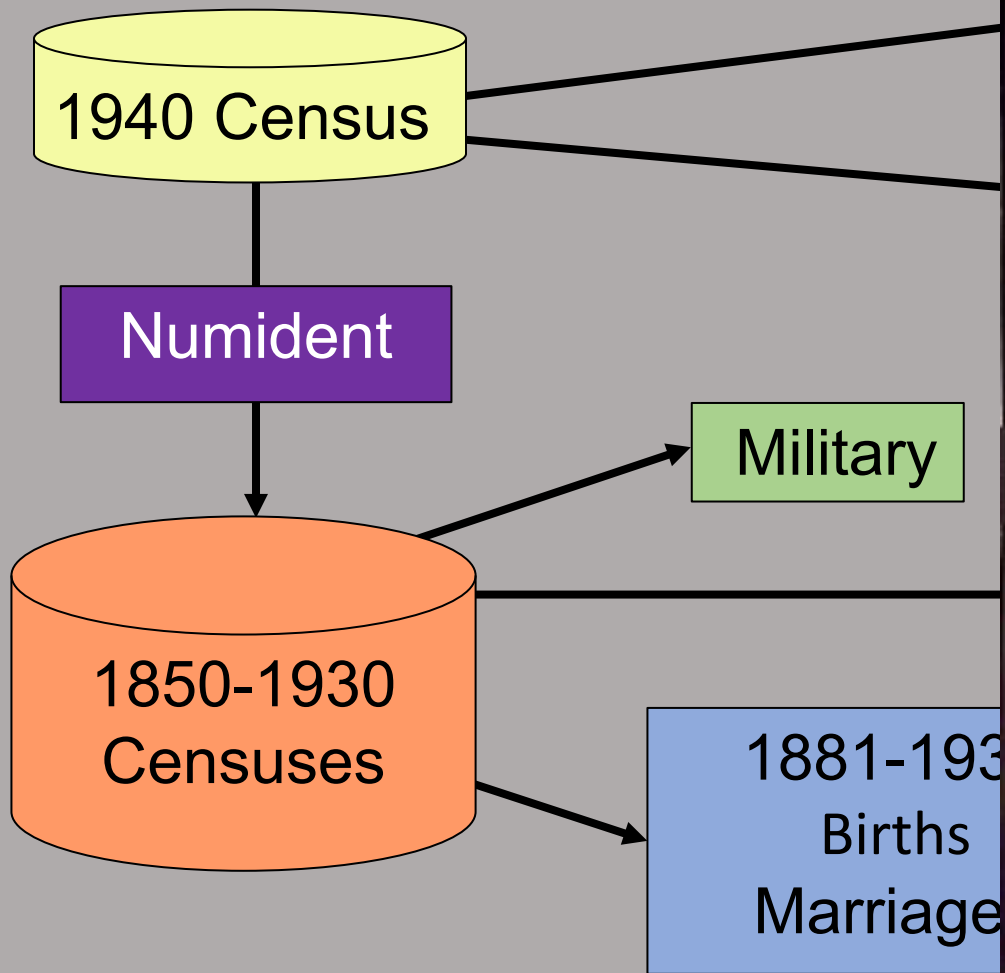


Figure 7: Comparison of 10-year inter-stage migration rates: the BPL10 vs. the linkage migration rates



5 Aging Surveys

MLP Linking Strategy

IPUMS-MLP Goals

- Use all available information to minimize errors and maximize linkage rates
- Weights to adjust for selection bias
- Links of individuals and families across nine censuses and other sources
- Large
- Easy-to-use
- Maintainable
- Expandable
- Interoperable with CLIP

NUMIDENT (Social Security Claims Database)

- Includes persons who had Social Security and had died or reached the age of 110 (not clear by when)
- Includes: name, up to three previous names used with date of use, date of birth, date of death, place of birth, sex, race, citizenship, mother's maiden name, father's name.
- Released by NARA in October 2018

HLink

- Hlink provides an end-to-end linking solution, replacing our earlier process using multiple programs (FEBRL, LIBSVM, C, Stats Packages) and data formats (ASCII, MySQL, binary files).
- Written in Python and Spark SQL, some Scala
- Leverages Apache Spark
- Enables parallel processing throughout
- Uses Parquet column-store data structure
- Two orders of magnitude faster than the old system



Thank You.