## Complete microdata for the U.S., 1790-1940
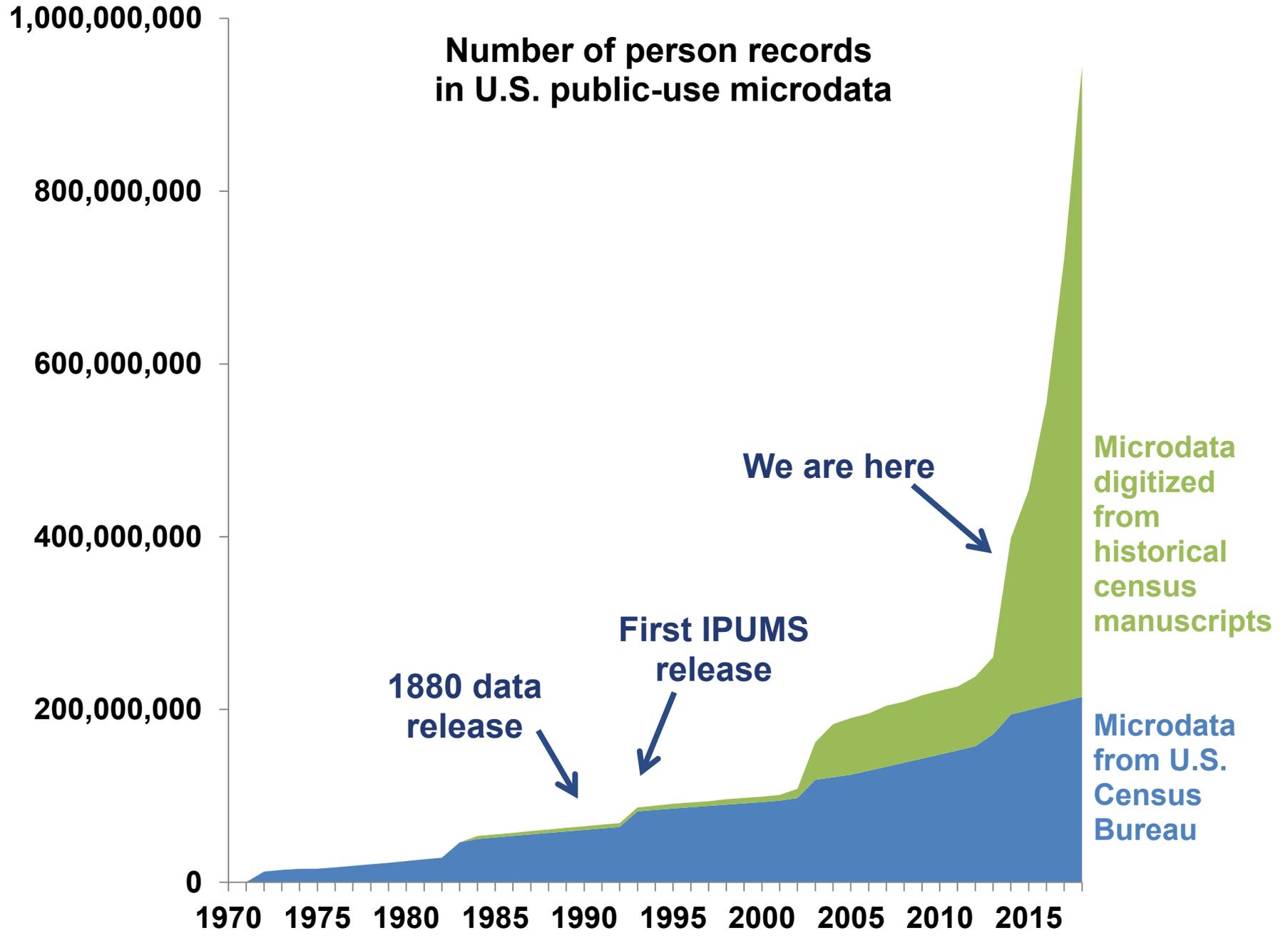
- 750 million records
- Data-entry required over 10,000 years of effort
- Donated data would cost $520 million to replicate

Number of person records
in U.S. public-use microdata

We are here

First IPUMS
release

1880 data
release

Microdata
digitized
from
historical
census
manuscripts

Microdata
from U.S.
Census
Bureau

1,000,000,000

800,000,000

600,000,000

400,000,000

200,000,000

0

1970  1975  1980  1985  1990  1995  2000  2005  2010  2015

Note A.—The Census Year begins June 1, 1879, and ends May 31, 1880.

Note B.—All persons will be included in the Enumeration who were living on the 1st day of June, 1880. No others will. Children BORN SINCE June 1, 1880, will be OMITTED. Members of Families who have DIED SINCE June 1, 1880, will be **Received July 22, 1880.**

Note C.—Questions Nos. 13, 14, 22 and 23 are not to be asked in respect to persons under 10 years of age.

1st District Princeton Township

E 1.—Inhabitants in _Borough of Princeton_, in the County of _Mercer_, State of _New Jersey_

enumerated by me on the _____ day of June, 1880.

_Chas O Hudnut_

| | The Name of each Person whose place of abode, on 1st day of June, 1880, was in this family. | Color | Sex | Age | If born within the Census year | Relationship | Single | Married | Widowed | Married during Census year | Profession, Occupation or Trade | Number of months unemployed | Sickness or disability | Blind | Deaf and Dumb | Idiotic | Insane | Maimed, Crippled | Attended school | Cannot read | Cannot write | Place of Birth of this person | Place of Birth of Father | Place of Birth of Mother |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 88 | Schenck Daniel | | M | 40 | | | | 1 | | | Laborer | ✓ | | | | | | | | | | New Jersey | New Jersey | New Jersey |
| | Phoebe | | F | 32 | | Wife | | 1 | | | Keeping House | | | | | | | | | | 1 1 | New Jersey | N.J. | N.J. |
| 89 | Johnson Clara | B | F | 51 | | | | | 1 | | Keeping House | | | | | | | | | | 1 1 | New Jersey | N.J. | N.J. |
| | Jane | B | F | 29 | | Daughter | 1 | | | | Servant | | | | | | | | | | | New Jersey | Rodney town | N.J. |
| | Mary E | B | F | 21 | | Daughter | 1 | | | | At Home | | | | | | | | | | | New Jersey | Penna | N.J. |
| | Georgiana | | F | 14 | | Daughter | 1 | | | | at Home | | | | | | | | | 1 | | New Jersey | Penna | N.J. |
| | Viola | | F | 2 | | G. Daughter | | | | | | | | | | | | | | | | New Jersey | Penna | N.J. |
| | Harry E | | M | 10/12 | | G. Daughter | | | | | | | | | | | | | | | | New Jersey | Penna | N.J. |
| 90 | Dickson James | B | M | 48 | | | | 1 | | | Laborer | 3 ✓ | | | | | | | | | | New Jersey | | |
| | Lucinda | B | F | 48 | | Wife | | 1 | | | Keeping House | | | | | | | | | | | New York | | |
| 91 | Golden Ellen | W | F | 55 | | | | | 1 | | Keeping House | ✓ | | | | | | | | | 1 1 | Ireland | Ireland | Ireland |
| | Mary | W | F | 17 | | Daughter | 1 | | | | Servant | ✓ | | | | | | | | | | New Jersey | Ireland | Ireland |
| | Philip | W | M | 20 | | Son | 1 | | | | Laborer | ✓ | | | | | | | | | 1 | New Jersey | Ireland | Ireland |
| 92 | O'Connor Ellen | W | F | 51 | | | | 1 | | | Servant | ✓ | | | | | | | | | 1 1 | Ireland | Ireland | Ireland |

# 1999: Collaboration with the Latter-Day Saints

FAMILY HISTORY RESOURCE FILE — CD-ROM LIBRARY

FamilySearch

1880 UNITED STATES CENSUS AND NATIONAL INDEX

55 CD-ROM Set

Presented to the

MINNESOTA POPULATION CENTER UNIVERSITY OF MINNESOTA

For their exceptional contribution to the development of the

1880 UNITED STATES CENSUS ON COMPACT DISC

Crowdsourcing Project
2006-2011

Form 15-6

DEPARTMENT OF COMMERCE—BUREAU OF THE CENSUS

FIFTEENTH CENSUS OF THE UNITED STATES: 1930

POPULATION SCHEDULE

State Illinois
County Cook
Township or other division of county Precinct 6

Incorporated place Chicago City
Ward of city 41   Block No. 58
Unincorporated place

Institution

Enumeration District No. 16-2792
Supervisor's District No. 3

Enumerated by me on April 7, 1930, Irene L. Burroughs, Enumerator.

| Line | House number | Dwelling number | Family number | Name | Relation | Home Data | | | Personal Description | | | | | Education | | Place of Birth — Person | Place of Birth — Father | Place of Birth — Mother | Mother Tongue | Code | | | Citizenship | | | Occupation | Industry | Code | Class | Emp. | | Vet |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 51 | | 1217 | | Stearns, Allen N. | Grandson | O | 5000 | | No | M | W | 5 | S | | Yes | No | Illinois | Missouri | Illinois | | 61 | | | | | | None | | | | | | |
| 52 | 708 | 128 | 128 | Retter, William | Head | O | 10,000 | R | No | M | W | 65 | M | 25 | No | Yes | Germany | Germany | Germany | German | 13 | 13 | V | 1870 | Na | Yes | Trunkmaker | Leather | 7740 | W | Yes | No |
| 70 | 693 | 134 | 134 | Schwiesow, Theodore | Head | O | 5000 | | No | M | W | 76 | M | 24 | No | Yes | Germany | Germany | Germany | German | 13 | 13 | V | 1873 | Na | Yes | None | | | | | |
| 71 | | | | — Minnie | Wife-H | | V | | | | F | W | 67 | M | 17 | No | Yes | Germany | Germany | Germany | German | 13 | 13 | V | 1874 | Na | Yes | None | | | | | |
| 72 | 947 | 135 | 135 | Schwiesow, William | Head | O | 10,000 | R | No | M | W | 45 | M | 33 | No | Yes | Illinois | Germany | Germany | | 61 | 13 | O | | | Yes | Plasterer | Building | 0HX1 | E | No | 33 | No |
| 73 | | | | — Rose | Wife-H | | V | | | | F | W | 38 | M | 26 | No | Yes | Illinois | Bohemia | Bohemia | | 61 | 15 | O | | | Yes | None | | | | | |
| 74 | 945 | 136 | 136 | Schwiesow, David | Head | O | 12,000 | R | No | M | W | 42 | M | 21 | No | Yes | Illinois | Germany | Germany | | 61 | 13 | O | | | Yes | Trimmer | Roofer | 37X1 | W | No | 57 | No |
| 75 | | | | — Amanda | Wife-H | | V | | | | F | W | 40 | M | 21 | No | Yes | Illinois | Illinois | Illinois | | 61 | | | | | Yes | None | | | | | |
| 76 | 935 | 137 | 137 | Olson, Octavius | Head | O | 10,000 | R | No | M | W | 46 | M | 22 | No | Yes | Illinois | Norway | Norway | | 61 | 05 | O | | | Yes | Decorator | China | 8394 | E | Yes | No |
| 77 | | | | — Augusta | Wife-H | | V | | | | F | W | 46 | M | 23 | No | Yes | Illinois | Germany | Germany | | 61 | 13 | O | | | Yes | None | | | | | |
| 78 | | | | — Perry | Son | | | | | M | W | 23 | S | | No | Yes | Illinois | Illinois | Illinois | | 61 | | | | | Yes | Decorator | China | 8394 | W | Yes | No |
| 79 | | | | — Evelyn | Daughter | | V | | | | F | W | 20 | S | | No | Yes | Illinois | Illinois | Illinois | | 61 | | | | | Yes | Stenos. | China | 7190 | W | Yes | |
| 80 | 923 | 138 | 138 | Reiling, Charles | Head | O | 7,500 | R | No | M | W | 64 | M | 24 | No | Yes | Illinois | Illinois | Illinois | | 61 | | | | | Yes | None | | | | | |
| 81 | | | | — Julia | Wife-H | | | V | | | F | W | 60 | M | 20 | No | Yes | Germany | Germany | Germany | German | 13 | 13 | V | 1872 | Na | Yes | None | | | | | |
| 82 | 923 | | 139 | Bredfield, Henry C. | Head | R | 70. | R | No | M | W | 31 | M | 24 | No | Yes | Illinois | Illinois | Germany | | 61 | 05 | O | | | Yes | Engineer | Heating | 5394 | W | Yes | Yes | WW |
| 83 | | | | — Muriel | Wife-H | | V | | | | F | W | 25 | M | 19 | No | Yes | Illinois | Missouri | California | | 61 | | | | | Yes | None | | | | | |
| 84 | | | | — Marilyn | Daughter | | V | | | | F | W | 3/12 | | | | | Illinois | Illinois | Illinois | | 61 | | | | | | None | | | | | |
| 85 | 919 | 139 | 140 | Blankenhagen, Henry | Head | O | 18,000 | R | No | M | W | 66 | M | 29 | No | Yes | Illinois | Germany | Germany | | 61 | 13 | O | | | Yes | Letter carrier | Mail | 8876 | W | Yes | No |
| 86 | | | | — Elizabeth | Wife-H | | V | | | | F | W | 59 | M | 19 | No | Yes | Illinois | Germany | New York | | 61 | 13 | | | | Yes | None | | | | | |
| 87 | 910 | 140 | 141 | Blankenhagen, Walter | Head | R | 12,000 | R | No | M | W | 38 | M | 21 | No | Yes | Illinois | Illinois | Illinois | | 61 | | | | | Yes | Steamfitter | Heating | 3028 | W | Yes | No |
| 88 | | | | — J. Florella | Wife-H | | V | | | | F | W | 38 | M | 21 | No | Yes | Illinois | Norway | Norway | | 61 | 05 | O | | | Yes | None | | | | | |
| 89 | | | | — Walter | Son | | | | | M | W | 10 | S | | | Yes | Illinois | Illinois | Illinois | | 61 | | | | | Yes | None | | | | | |
| 90 | | | | — Elaine | Daughter | | V | | | | F | W | 7 | S | | | Yes | Illinois | Illinois | Illinois | | 61 | | | | | Yes | None | | | | | |
| 91 | 907 | 141 | 142 | Heubach, Henry | Head | O | 10,000 | R | No | M | W | 49 | M | 32 | No | Yes | Illinois | Germany | Germany | | 61 | 13 | O | | | Yes | Clerk | Post office | 7X76 | W | Yes | Yes | WW |
| 92 | | | | — Lillian | Wife-H | | V | | | | F | W | 47 | M | 20 | No | Yes | Illinois | Illinois | Illinois | | 61 | | | | | Yes | None | | | | | |
| 93 | | | | — Henry G. | Son | | | | | M | W | 17 | S | | Yes | Yes | Illinois | Illinois | Illinois | | 61 | | | | | Yes | None | | | | | |
| 94 | | | | — Robert W. | Son | | | | | M | W | 15 | S | | Yes | Yes | Illinois | Illinois | Illinois | | 61 | | | | | Yes | None | | | | | |
| 95 | | | | — Edward W. | Son | | | | | M | W | 12 | S | | Yes | Yes | Illinois | Illinois | Illinois | | 61 | | | | | Yes | None | | | | | |
| 96 | 904 | 142 | 143 | Bretscher, John M. | Head | O | 12,000 | R | No | M | W | 28 | S | | No | Yes | Wisconsin | Missouri | Michigan | | 63 | | | | | Yes | Realtor | Real Estate | 8286 | E | Yes | No |
| 97 | | | | — Marie | Mother-H | | V | | | | F | W | 57 | Wd | | No | Yes | Michigan | Missouri | Michigan | | 62 | | | | | Yes | None | | | | | |
| 98 | | | | Hutchinson, Henry | Bro-in-law | | | | | M | W | 38 | M | 31 | No | Yes | Ohio | Ohio | Virginia | | 59 | | | | | Yes | Salesman | Insurance | 8880 | W | Yes | Yes | WW |
| 99 | | | | — Margaret | Sister | | V | | | | F | W | 34 | M | 27 | No | Yes | Wisconsin | Missouri | Illinois | | 63 | | | | | Yes | Stenos. | Exports | 7190 | W | Yes | No |
| 100 | 910 | 143 | 144 | Koolehan, Paul | Head | R | 90. | R | No | M | W | 40 | M | 27 | No | Yes | Illinois | Illinois | Illinois | | 61 | | | | | Yes | Electrician | Wiring | 1V41 | W | Yes | No |

# 1940 Census Project

Collaboration of MPC and Ancestry.com

- 7.8 billion keystrokes
- 134 million persons

State New York
County Bronx
Bronx Burough
Incorporated place New York
Township or other division of county A.S.1
Ward of city 7
Block Nos. A'B'C'F Institution
DEPARTMENT OF COMMERCE—BUREAU O
SIXTEENTH CENSUS OF THE UNITE
POPULATION SCHED
I.D.No. 63-142
Sheet No.
April X, 1940. 1 A
Enumerator

| LOCATION | HOUSEHOLD DATA | NAME | RELATION | PERSONAL DESCRIPTION | EDUCATION | PLACE OF BIRTH | CITIZENSHIP | RESIDENCE, APRIL 1, 1935 | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

**EDUCATION**

## PERSONS 14 YEARS OLD AND OVER—EMPLOYMENT STATUS

| Was this person AT WORK for pay or profit in private or nonemergency Govt. work during week of March 24–30? (Yes or No) | If not, was he at work on, or assigned to, public EMERGENCY WORK (WPA, NYA, CCC, etc.) during week of March 24–30? (Yes or No) | If neither at work nor assigned to public emergency work. ("No" in Cols. 21 and 22) | | | | OCCUPATION, INDUSTRY, AND CLASS OF WORKER | | | | | |
| | | Was this person SEEKING WORK? (Yes or No) | If not seeking work, did he HAVE A JOB, business, etc.? (Yes or No) | Indicate whether engaged in home housework (H), in school (S), unable to work (U), or other (Ot) | CODE | Number of hours worked during week of March 24–30, 1940 | Duration of unemployment up to March 30, 1940—in weeks | OCCUPATION Trade, profession, or particular kind of work | INDUSTRY Industry or business | Class of worker | CODE (Leave blank) | Number of weeks worked in 1939 |
| 21 | 22 | 23 | 24 | 25 | E | 26 | 27 | 28 | 29 | 30 | F | 31 |
| YES | — | — | — | — | 1 | 40 | | TIMEKEEPER | CUT STONE YARD | P.W | 266 24 | 52 |
| YES | — | — | — | — | 1 | 40 | 28 | HELPER | CUT STONE YARD | P.W | 496 24 | 40 |
| No | No | No | No | H | 5 | | | | | | | 0 |
| B | | | | | | | | | | | | |
| YES | — | — | — | — | 1 | 48 | | SUPERINTENDANT LOFT BUILDING | PHILLIP JONE SHIRT COMPANY | P.W | 129 81 | 52 |
| No | No | No | No | H | 5 | | | | | | | |

For Persons R SAME HOUSE

NO H4 30

NO H-4 36

| | | FATHER | MOTHER | CODE | | | | | | | | USUAL OCCUPATION | USUAL INDUSTRY | Usual class of worker | CODE (Leave blank) | | | | | |
| 14 | IMOR. AUGUSTUS | NEW YORK | NEW YORK | ENGLISH | | No No No | No | — | 0 | | | | | | | | | | | |
| 29 | TENKE VILMA | YUGOSLAVIA | GERMANY | 17 ENGLISH | | No No No | No | — | 0 | | NONE | STUDENT | S. | | 0 7 2 1 6 8 7 | | 0 00 0 L |

# Terms of Ancestry License

- **Public data**
  - Numerically-coded data
  - For education and scholarly research

- **Restricted data**
  - Alphabetic strings available through institutional license with data security plan

Number of person records in historical American census data

# Data Improvements

- New variables: 1860 - 1930

- Variable coding

- Data cleaning

- Missing data allocation

- IPUMS constructed variables

# New Variables – Housing

Housing Characteristics

- Farm residence

- Ward

- House number and street

- Homeownership

- Mortgage status

- Value of home/cost of rent

# Demographic and Health Variables

Demographic

- Birth month

- Married in year

- Age at first marriage

- Number of times married and duration of marriage

- Children ever-born and children surviving

Disability

- Deaf, blind

- Sickness

# Socioeconomic Variables

**Wealth variables**

- Value of real estate and personal property

**Work**

- Occupation

- Industry

- Class of worker

**Education**

- School attendance

- Literacy

# Data Availability

| Public | Restricted |
|---|---|
| 1790-1840 - soon | 1790-1840 |
| 1850 – soon | 1850 |
| 1860 - TBA | 1860 |
| 1870 - TBA | 1870 |
| 1880 - now | 1880 |
| 1900 - TBA | 1900 |
| 1910 - TBA | 1910 |
| 1920 – mid 2016 | 1920 |
| 1930 - early 2016 | 1930 |
| 1940 - now | 1940 |

- Census Longitudinal Infrastructure Project (CLIP)
- National Historical Census Files Project

# Census Longitudinal Infrastructure Project (CLIP)

1. PIK the 1940 Census (Protected Identification Key)

2. Create a longitudinal resource, linking 1940 census to later censuses, surveys (e.g. CPS, SIPP, NHIS), and administrative records (e.g. Social Security, Medicare, Medicaid)

3. Disseminate CLIP through the RDC network

# National Historical Census Files Project

1. Recover and verify of 1960, 1970, 1980, and 1990 long-form and short-form microdata

2. Create an IPUMS-compatible version of internal census microdata files from 1960-present (including ACS)

3. Create new PUMS and summary files for 1960

# Integrated U.S. microdata available for research, 1970-2018
(number of person records)



2,000,000,000

1,500,000,000

**IPUMS-Format Microdata in the Census Research Data Centers**

1,000,000,000

**IPUMS Microdata digitized from historical manuscripts**

**We are here**

500,000,000

**Public-use IPUMS data from Census**

0

1970   1980   1990   2000   2010

# Creating linked samples

# 2003 onwards – linking samples to complete count datasets

| Canada | Great Britain | Iceland | Norway | Sweden | United States |
|--------|--------------|---------|--------|--------|---------------|
| 1852 | 1851 | 1703 | 1801 | 1890 | 1850 |
| 1871 | 1881 | A very old person | Some very old people | 1900 | 1860 |
| | | 1835 | 1865 | | 1870 |
| 1881 | | 1845 | 1875 | | 1880 |
| 1891 | | 1870 | 1900 | | 1900 |
| 1901 | | 1880 | | | 1910 |
| 1911 | | 1900 | | | 1920 |
| 1921 | | | | | 1930 |
| 1931 | | | | | |
| 1941 | | | | | |
| 1951 | | | | | |

Number of person records in historical American census data

Number of males surviving decade 1850-1860 and potentially linkable

Total white males surviving 1850-1860 is approximately 10,006,097

# Linking processes

# Blocking strategy

- Sex of individual or dyad

- Race (U.S.)

- Birthplaces
  - US / Canadian data
    - State or province for native born
    - Foreign country for migrants
  - Scandinavian / British
    - Parish of birth [human memory is poor] for natives
    - Foreign country for migrants
  - For all countries: some random error with more or less precision than required by enumeration

MPC

# Name cleaning

- Remove stray characters (non-alpha), titles, middle initials

- Parse into separate fields for first and last names

- Some standardization of names for common names with variants and abbreviations [practice has varied over time …]
  - Tom / Thomas
  - Wm. / William
  - Edward / Theodore

- Original name retained as separate variable

# Age windows

- Integer age at census transformation into year of birth OR expected age at next census has inherent error
    - Month of enumeration varies
    - Age heaping from numeracy / enumeration practices
- Potential links restricted to age window of +/- 4 to 5 years
    - Have experimented with up to 7

**MPC**

**Distribution of Birth Year Differences**

IPUMS Matched Samples between 1880 and 1850–1930

# Comparison of names

- Jaro-Winkler scores calculated for
    - First name pairs
    - Standardized first name pairs
    - Last name pairs

# Jaro-Winkler string comparator

$$\text{sim}_{\text{jaro}}(s_1, s_2) = \frac{1}{3|} \left( \frac{c}{|s_1|} + \frac{c}{|s_2|} + \frac{c-t}{c} \right).$$

Where s1 and s2 are the two strings

    e.g. Knutsen Haugen and Knudtson

c : Number of characters that are the same in half the length of the longer string

t: Number of transpositions of adjacent characters

**MPC**

# Jaro-Winkler string comparator

$$\mathrm{sim_{winkler}}(s_1, s_2) = \mathrm{sim_{jaro}}(s_1, s_2) + (1.0 - \mathrm{sim_{jaro}}(s_1, s_2))\frac{p}{10},$$

Winkler modification increases the score (similarity) if beginning is similar and differences appear later.

p: Number of agreeing characters at start,

$$0 \le p \le 4$$

**MPC**

# Composite score

- Jaro-Winkler scores summed to create composite score

- Pairs of records ≥ threshold value written out
  - Decision: What is threshold

# Training data

- Clerical review of pairs of links to mark records as true or false links

- Use household information to confirm which links are true

# Additional indicator variables

- NYSIIS match

- Last name matches on double metaphone

- Is first name match a pair of initials?

- Middle initial present in both records

# Age scores

- age70pct
  $1 - |age_1 - age_2| / [\max\{age_1, age_2\}/0.71]$
  Intuition: Deviations at younger ages more important

- mrecage_plus_norm
  Classifies age from earlier year
  $Floor\{(age/5) + 1\} / 100$

- agediff_norm_abs
  $1 / (|age_1 - age_2| + 1)$

**MPC**

# Support vector machine

- Training data provides set of true links with associated characteristics on variables previously described

- SVM selects best combination of data from potential links

- Select positive confidence records (we think they're a match) that don't conflict with matches to another record

# Steps in process

| Step | Format | Software |
|------|--------|----------|
| Edit and format input data<br>    Blocks<br>    Expected ages | Starts as IPUMS files<br>Fixed width<br>Coded | Statistics package<br>Perl tools on ASCII files<br>(on MPC servers) |
| Name score algorithm | Delimited files | Python (on MSI computer) |
| Re-attach additional variables | Delimited files | Python (on MSI computer) |
| Run data through classifier | Delimited files | Libsvm (MSI) |
| Evaluate links | | Statistics package |
| Post-linking filtering | | Statistics package |
| Format dataset | | Statistics package |
| Make weights | | Statistics package |
| Make documentation | | Text editor / Word |

# Acknowledgements

National Institutes of Health

National Science Foundation

Ancestry.com

Family Search

U.S. Census Bureau

**MPC**

**Minnesota Population Center**