

SynLBD inputs

Synthetic Longitudinal Business Data

International User Seminar

May 9, 2017

What is the basis for the synthetic data?

- Should the cleaned data be the input to synthetic data, or data with missing values, data failing plausibility edits, etc.?
- Current approach:
 - cleaned data is used
 - (SSB) fully completed (through imputation)
 - Neither data synthesizes missing data patterns

Issues

- Changing concepts over time
 - Employment measures (Germany)
 - Industry coding (all countries)
 - Geography coding (counties in the US?)
 - Coverage (geography: Germany)

Versioning and updating

- Versions of the input data
 - Should they matter?
 - Need to be tracked
- Updated methodology (e.g. longitudinal linkage)
- Updates of the Synthetic Data

Example: German SynLBD

Replicating the synthetic LBD with German establishment data

59th World Statistics Congress
28. August 2013, Hong Kong

▪
Excerpt
▪

Jörg Drechsler
Institute for
Employment Research

&

Lars Vilhuber
Cornell University

The Longitudinal Business Database

- created from the U.S. Census Bureau's Business Register
- data available from 1976 to 2011
- contains information on:
 - birth
 - death
 - industry
 - location
 - payroll
 - firm affiliation

- in the synthetic version location is not available

The German Employment History Panel (BHP)

- no business register available at the IAB
- all establishment level information is derived by aggregating the German Social Security Data via the establishment id
- BHP is one of the data products derived from the GSSD
- BHP will be the main data source to build the German LBD
- data available from 1976/1992 (Western Germany/Eastern Germany)
- contains detailed information on the personnel structure
- not all the variables available in the LBD are also available in the BHP

Differences between the LBD and BHP

- information whether establishment belongs to multi-unit firm not available
- until 1999 the BHP only contains establishments that had at least one employee covered by social security
- payroll information in the BHP is for the reference date June 30 for each year
- LBD contains yearly payroll

Building the German LBD (GLBD)

Table 1: Variables from the BHP that were used for generating the GLBD

Name	Description
ID	Unique Random Number for Establishment
County	Geographic Information on the County Level
State	Geographic Information on the State Level
WZ73	Industry code according to 1973 classification
WZ93	Industry code according to 1993 classification
WZ03	Industry code according to 2003 classification
WZ08	Industry code according to 2008 classification
Firstyear	First Year Establishment is Observed
Lastyear	Last Year Establishment is Observed
Employment _{tot}	Total Number of Employees on June 30
Employment _{ss}	Number of Employees covered by Social Security on June 30
Employment _{me}	Number of Employees with Marginal Employment on June 30

Ensuring Consistent Establishment Size

- until 1999 employers only had to report all their employees covered by social security
- since 1999 all employees need to be reported
- significant changes in the data between 1999/2000
 - many establishments report more employees although they didn't grow
 - increase in the number of establishments since establishments with only marginally employed are also included
- to ensure consistency, we
 - subtract the number of marginally employed from total number of employees
 - set all establishment sizes = 0 to missing
 - drop all establishments that never report their establishment size after the adjustments
- final dataset contains 6,916,183 establishments

Generating a unique geographic location and industry code

- geographic location and industry code are constant in the LBD
- this is not true for the BHP
- select the mode of both variables over the lifespan of the establishment
- if two modes exist, the first one is selected
- might be improved
 - select mode randomly
 - weight the years by establishment size

Generating a unique geographic location and industry code

Table 2: number of establishments with a change in location or industry

Variable	number of status changes (% changes based on entire dataset)	years in which informa- tion is available	# of records with at least one reported value
County	214,354 (2.72)	1975–2008	7,851,109
State	45,638 (0.58)	1975–2008	7,851,109
WZ73	229,759 (2.91)	1975-2002	6,037,241
WZ93	21,866 (0.28)	1999-2003	3,502,881
WZ03	49,773 (0.63)	2003-2008	4,081,497

- only a small number of establishments report a change
- even fewer have two or more modes
- we stick with the simple approach

Updating the information on establishment births and deaths

- information on the first/last year establishment is observed is not necessarily equivalent with the birth/death of the establishment
 - data are left and right censored
 - new establishment appears whenever a new establishment id is generated
- new ids are not necessarily equivalent to a new establishment
- several other reasons possible, e.g.
 - change of ownership
 - change in declared industry classification
- these new ids should not be treated as establishment births
- similarly disappearance of ids should not always be treated as establishment deaths

Updating the information on establishment births and deaths

- use the employee flows to identify real births and deaths based on a similar approach by Benedetto et al (2007)
- flow-files generated by Hethey and Schmieder (2010)
- basic idea
 - if (almost) all employees of an exiting establishment work in the same new establishment in the following year this is most likely an id change
 - if (almost) all employees of a new establishment worked in the same establishment in the year before but this establishment still exists in the current year, the new establishment is most likely a spin-off.
- the files also contain suggestions how the observed births and deaths should be categorized
- we use the suggested classification

Updating the information on establishment births and deaths

- all birth and death categories are treated as births and deaths
- establishments with unknown status are treated according to the information in the BHP
- spin-offs are left unchanged
- id changers are merged and employment and payroll information is aggregated
- industry and geographic information are based on the mode of the observed variables in the linked record
- Some establishments identified as dead reappear in the data later
- Since number is small (3,941 establishments) we ignored this

Adding payroll information

- payroll information only available at a reference date in the BHP
- possible to derive yearly payroll by aggregating the information from each employee that was ever employed in a specific establishment for a given year
- aggregated yearly payroll information also available from another project at the IAB
- only includes the payroll of all full time employees
- for almost 230,000 (3.3%) records no payroll information is available
- payroll information for all establishments in the BHP based on all employees from the underlying administrative data could be incorporated in the future