

IRS County-to-County Migration Data *

Mathew E. Hauer ^{1,2*} *Florida State University*

¹ Department of Sociology, Florida State University. 605 Bellamy Building, 113 Collegiate Loop, Tallahassee, FL USA 30622.

² Center for Demography and Population Health, Florida State University

BACKGROUND: The Internal Revenue Service’s (IRS) county-to-county migration data are an incredible resource for understanding migration in the United States. Produced annually since 1990 in conjunction with the US Census Bureau, the IRS migration data represent 95 to 98 percent of the tax filing universe and their dependents, making the IRS migration data one of the largest sources of migration data. However, any analysis using the IRS migration data must wrangle at least seven legacy formats of these public data across more than 2000 data files – a serious burden for migration research.

OBJECTIVE: To produce a single, flat data file containing complete county-to-county IRS migration flow data.

METHODS: This paper uses R to wrangle more than 2,000 IRS migration files into a single, flat data file for use in migration research.

CONTRIBUTION: To encourage and facilitate the use of this data, I provide the full R script to download and flatten the IRS migration data as counts and a finalized data set covering the period 1990-2010.

Introduction

Migration flow data (ie, the number of migrants from location i to location j) is typically difficult to obtain (Willekens et al., 2016). In the United States, migration flow data is available from two primary sources: the decennial census and the Internal Revenue Service’s (IRS) county-to-county migration data¹. The IRS annual series of county-to-county migration data cover 95 to 98 percent of the tax filing universe and their dependents making these data the largest migration data source for count flows between counties in the United States.

The IRS began using tax data to estimate migration in the 1970s and 1980s (Engels and Healy, 1981; Franklin and Plane, 2006). Beginning with tax year 1991 (migration year 1990), the IRS produces these data in conjunction with the US Census Bureau using the IRS Individual Master File which contains every Form 1040, 1040A, and 1040EZ. A migrant is identified when a current year’s tax return contains an address that is different from the matched preceding years’ return. For the 2002 tax year, the IRS migration data contained approximately 130 million returns (Gross, 2005). If an address is identical over two years, a filer is considered a non-migrant.

Because of the annual availability, relatively long time series (the series starts in 1990), and large universe, the IRS data are an attractive data source for conducting migration research (e.g. (Curtis, Fussell and DeWaard, 2015; Molloy, Smith and Wozniak, 2011; Shumway

*The data and code that supports this analysis are available in the supplementary materials.

¹The IRS migration data are available at <https://www.irs.gov/statistics/soi-tax-stats-migration-data>

and Otterstrom, 2001; Frey, 2009)). Unfortunately, these data exist in seven legacy formats, split over 2,000+ files making analysis with this data rather burdensome and the prevented the widespread use of this valuable resource.

To encourage and facilitate the use of this tremendous migration resource, I flattened the IRS migration data into a single, standardized file containing all migration flows for the period 1990-2010. By publishing both the R code to collate, wrangle, and flatten the IRS migration data as well as the migration data itself, my hope is to save time for other migration scholars and facilitate their use. Scholars who wish to use these data should still familiarize themselves with the strengths and weaknesses, idiosyncracies, and design of these data (see (Gross, 2005; Engels and Healy, 1981; Franklin and Plane, 2006; Pierce, 2015) and with the procedures outlined in this document and in the corresponding R code.

While the IRS migration data presently continues to 2015, I limit the data to the period 1990-2010. Starting in tax year 2012 (migration year 2011), the IRS introduced a change in their processing (Pierce, 2015) creating a break in the time series.

Data Preparation

The IRS migration data for the period 1990-2010 are available in seven legacy formats. **Table 1** summarises some of the similarities and differences in these formats. For every year, the IRS publishes at least 104 data files (52 state entities by in/out-migration). The underlying file organization, file format, naming schema, and coding can differ between these legacy formats. Migration years 1990 and 1991 are available as fixed-width text file, while 1992-2010 are available as excel files. For years 1990-2003, the IRS separated in/out migration into separate folders while 2004-2010 are published in a single file. Each legacy format utilizes a different file naming scheme as well, making grepping difficult. Importantly, the IRS treats non-migrants and total migrants differently in the seven legacy formats. For 1990 and 1991, the IRS simply has a field that reads “County Non-Migrants” for non-migrants; for 1992-1994, the IRS introduced a State code 63 but two difference County codes (010 for 1992 and 1994 and 050 for 1993) creating a 5-digits FIPS code of 63010 or 63050. After 1995, the IRS smartly set the origin FIPS equal to the destination FIPS for non-migrants. Lastly, Total Migrants are treated differently too. For 1990 and 1991, the destination field simply reads “Total Migration.” For 1992-1994, the IRS introduced a State Code 00 and county code 001 for total migrants. After 1995, the IRS used State Code 96 and county Code 000 for a combined 5-digit FIPS code of 96000.

The differences described above and in **Table 1** are only some of the differences that are of interest to the data I produce here. Total Migrants, ie FIPS 96000 for migration data after 1995, is also broken down into Total Mig - US (FIPS 97000), Total Mig - US Same State (FIPS 97001), Total Mig - US Diff St (FIPS 97003), and Total Mig - Foreign (FIPS 98000). The IRS did not code these migration flows in this manner for all years, and in some cases (such as Total Mig - Foreign) migration flows are not reported. For simplicity and data continuity purposes, I simply create a new origin/destination (FIPS 99999) that contains all unspecified migration flows. I do this by subtracting the number of migrants from migration flows with greater than 10 migrants from the total number of migrants. This way, the sum of all migrants will equal the total number of migrants.

Besides the aggregation to FIPS 99999, I performed no other post-processing of the IRS data.

Table 1: Select differences in the file formats, file organizations, naming, and treatment of various migration statistics.

Years	Data Format	File Organization	Sample File naming	Treatment of non-migrants	Treatment of Total Migrants
1990-1991	txt	Separate in/out migration	C9091pao.tx	Destination field reads 'County Non-Migrants'	Destination field reads 'Total Migration'
1992, 1994	xls	Separate in/out migration	C9293Pao.xls	State code = 63, County code = 010	State code = 00, County code = 001
1993	xls	Separate in/out migration	co934pao.xls	State code = 63, County code = 050	State code = 00, County code = 001
1995-2003	xls	Separate in/out migration	co956paor.xls	Origin FIPS = Destination FIPS	State code = 96, County code = 000
2004-2006	xls	Single folder	co0405PAo.xls	Origin FIPS = Destination FIPS	State code = 96, County code = 000
2007-2008	xls	Single folder	co0708oPa.xls	Origin FIPS = Destination FIPS	State code = 96, County code = 000
2009-2010	xls	Single folder	co0910oPA.xls	Origin FIPS = Destination FIPS	State code = 96, County code = 000

R Code

The R code used to produce these data is available in the **Supplementary Materials** and can also be found on github. The code makes use of multi-core processing to speed up computation time. There are three main sections in the code: A setup section; a data download section; and a data processing section. The final flat file, `county_migration_data.txt`, contains the # of exemptions and can be either downloaded at github or produced by running the R code.

Setup

The first part of the code simply sets up the R workspace to facilitate the data wrangling. The appropriate R packages are downloaded and installed if the user does not already have these

packages installed. The parallel computing environment is also set up as `DetectCores() - 1` to ensure one user core is left for other tasks that computer might need. A single reference tab separated (tsv) file is required in this section. `ref_state.tsv` contains FIPS code information for US states. I simply add an additional FIPS code for ‘unknown’ and assign it FIPS state 99.

Data Download

The IRS data that can be downloaded is in two primary formats. 1990-2003 and 2004-onwards. This section of code will download and unzip the migration data into a folder standardized format into subdirectory, `MigData/`. There are eight files that the IRS includes in their zip archives that contain no data. After being downloaded and unzipped, these files are then deleted. If they are not deleted, they will cause the subsequent `for loops` to fail in the next section. These files do not contain any state specific migration information.

Data Processing

The third and final section contains several `foreach` parallel processing loops to wrangle the seven legacy formats into a common data format. These files are then row-bound using `rbindlist` and transformed into a ‘short’ data frame. **Table 2** demonstrates the general file layout.

Usage Notes

Any origin-destination pair with fewer than 10 tax filers over the entire period is excluded from the final datafile since no data would be recorded in the IRS datafile.

The count data come from the `exemptions` field of the IRS migration data. The original IRS migration data contains two consistent fields across all years of data: a `returns` field and an `exemptions` field. Returns are the actual number of tax returns filed while exemptions are a proxy for the members of the household.

US Counties are fairly stable geographic units but some changes in counties do exist. To try and keep as close to the original data fidelity as possible, I did not recode any geographic changes. For instance, Broomfield County CO (FIPS 07014) was created out of parts of Adams, Boulder, Jefferson, and Weld counties and thus contains data only after 2002. Users should be aware of any changes in geography that could substantially alter any analyses².

Users should be aware that I prepared these data to stay as close to the original data fidelity as possible, documented their creation, and with open-source computer code. These data should be used only with full awareness of the inherent limitations of the IRS migration data and with the knowledge of the procedures outlined in this document and in the corresponding R code.

²More detailed information about county boundary changes can be found at the following locations
<https://www.census.gov/geo/reference/county-changes.html> http://www.nber.org/asg/ASG_release/County_City/FIPS/FIPS_Changes.pdf https://www.cdc.gov/nchs/data/nvss/bridged_race/County_Geography_Changes.pdf https://www.ddorn.net/data/FIPS_County_Code_Changes.pdf

Table 2: Selected file format for the final flat file. Origin and Destinations are the five-digit FIPS codes with 99999 representing all flows with fewer than 10 filers.

origin	destination	1990	1991	1992	...	2010
01001	01001	10785	11010	11196	...	17278
01001	01003	0	0	12	...	11
...
01001	99999	524	403	479	...	732

References

- Curtis, Katherine J, Elizabeth Fussell and Jack DeWaard. 2015. “Recovery migration after Hurricanes Katrina and Rita: Spatial concentration and intensification in the migration system.” *Demography* 52(4):1269–1293.
- Engels, Richard A and Mary K Healy. 1981. “Measuring interstate migration flows: an origin—destination network based on internal revenue service records.” *Environment and Planning A* 13(11):1345–1360.
- Franklin, Rachel S and David A Plane. 2006. “Pandora’s box: The potential and peril of migration data from the American Community Survey.” *International Regional Science Review* 29(3):231–246.
- Frey, William. 2009. “The great American migration slowdown.” *Brookings Institution, Washington, DC*.
- Gross, Emily. 2005. Internal revenue service area-to-area migration data: Strengths, limitations, and current trends. In *Proceedings of the Section on Government Statistics*. p. 2005.
- Molloy, Raven, Christopher L Smith and Abigail Wozniak. 2011. “Internal migration in the United States.” *Journal of Economic perspectives* 25(3):173–96.
- Pierce, K. 2015. “SOI migration data. A new approach: Methodological improvements for SOIC’s United States population migration data, calendar years 2011–2012.” *Statistics of Income, Internal Revenue Service*.
- Shumway, J Matthew and Samuel M Otterstrom. 2001. “Spatial patterns of migration and income change in the Mountain West: the dominance of service-based, amenity-rich counties.” *The Professional Geographer* 53(4):492–502.
- Willekens, Frans, Douglas Massey, James Raymer and Cris Beauchemin. 2016. “International migration under the microscope.” *Science* 352(6288):897–899.