

PHP 2410E - Assignment 2

Kevin W. McConeghy

Compiled: 2019-09-27

Contents

Introduction	2
Statement of work	2
Note on R markdown	2
Assignment Overview	2
Assessment of Minimum Dataset Assessments	3
Dataset	3
Reformat	5
Consistency of data across MDS assessments	5
Assessment of Medicare Enrollment File	6
Dataset	6
Consistency of data across enrollment years	7
Matching MDS and Medicare file	8
Prepare for matching	8
Comparison of MDS and denominator data	8
Comparison of Date of Birth	10
Summary discussion	10
Session Info	11

Introduction

This is the completed assignment 2 for the ‘Medicare data’ course at Brown University. All code is stored in a Github repository, <https://github.com/kmcconeghy/PHP2410E>

Statement of work

This document was created solely by the author, guidance in the homework solutions was driven by class instruction, materials or prior experience. The solutions were not shared with anyone else.

Note on R markdown

This report was generated using R markdown, LaTeX, and several non-base R packages (e.g. tidyverse).

```
## Non-base packages loaded:  Scotty tidyverse rJava kableExtra
```

Assignment as written:

- Data Assignment #2
- Working with Medicare Public Use Files
- Due October 2nd, 2019

Assignment Overview

There are two data sets each containing identifying information on Medicare beneficiaries. The information includes date of birth, gender and race. The two sources of data are:

- 1) The Minimum Data Set, a clinical patient record assessment that is completed each time a patient is admitted to a nursing home in the US that is certified by Medicare/Medicaid. The MDS data file (in either STATA or SAS format) is likely to include multiple records per ID (unique individual) for many individuals;
- 2) The Medicare Enrollment Record is the individual identifier of every Medicare beneficiary. The overall file includes detailed data from Social Security as well as whether and when the beneficiary had joined a Medicare Advantage Plan and when. The current file includes only identifying information such as date of birth, gender and race. There is only one record per person per year. The ID number on the Medicare Enrollment Record is the same as that on the MDS data file.

There are two parts to this data assignment. Each is described below. There are various ways to complete these components. It is up to the student to choose how s/he goes about completing the assignment.

- Using the MDS data determine how consistent the identifying data are across records for the same persons who have multiple records. This means, among those with more than one record, what is the rate of inconsistency for date of birth, gender and race. DoB, gender and race have multiple ways in which they could disagree, in addition to calculating the rate of inconsistency, characterize the different ways (e.g. day, month or year of birth; missing categories of information in which this inconsistency manifests itself). (Note:

The Codebook for the MDS2.0 data can be found here: <https://www.cms.gov/Medicare/Quality-Initiatives-Patient-Assessment-Instruments/NursingHomeQualityInits/downloads/MDS20MDSAllForms.pdf>)

- Match the MDS and Medicare enrollment records based upon patient ID, using the first MDS record to match for those cases that have multiple MDS records. Next, compare the DoB, gender and Race from the two different data sources. Note that DoB is necessarily measured the same way across the two data sources as is gender (although there may be differing levels of missing data). The race variables across the two data sets are coded differently so it may be important to separately estimate the degree of agreement across the categories that are comparably labeled (e.g. white vs. white; black vs. black) since the Medicare Enrollment record is known to underestimate the number of Hispanic and Asian Americans.

Assessment of Minimum Dataset Assessments

The overall objective is to better understand the validity of MDS assessments by examining variation in patient characteristics which would typically be fixed across time (e.g. age, gender, race/ethnicity).

Dataset

```
## Dataframe: MDS dataset
## Memory Size: 45 Mb Rows: 1,481,145 Columns: 5
## id: 150,757 Missing: 0
```

Data structure

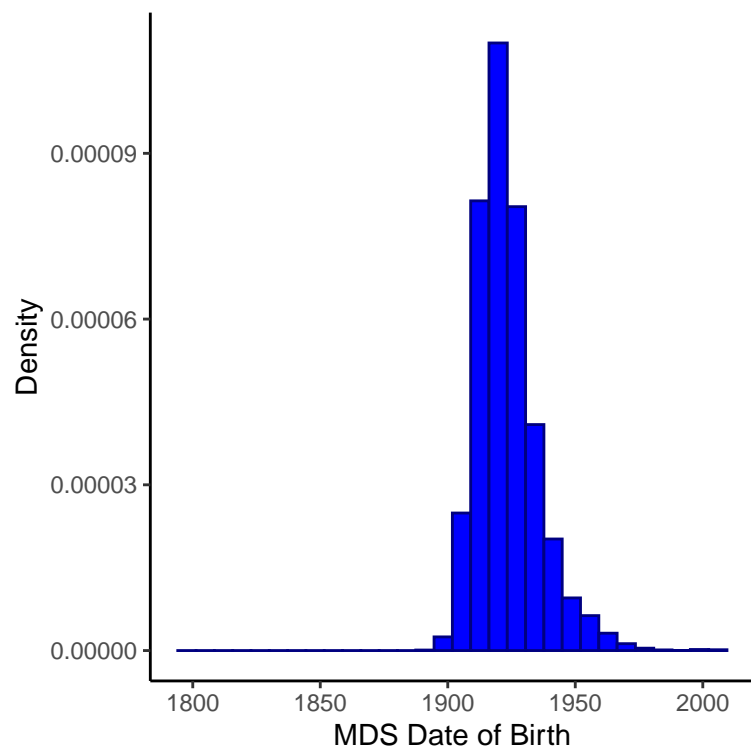
```
## Classes 'tbl_df', 'tbl' and 'data.frame': 1481145 obs. of 5 variables:
## $ idgendr: Factor w/ 2 levels "Male","Female": 1 1 1 1 1 1 1 1 1 1 ...
## $ idbirdt: Date, format: "1943-11-02" "1943-11-02" ...
## $ idrace : Factor w/ 5 levels "Native-American",...: 5 5 5 5 5 5 5 5 5 5 ...
## $ dmdate : Date, format: "2003-09-25" "2003-12-15" ...
## $ id : num 1 1 1 1 1 1 1 1 2 2 ...
## ..- attr(*, "format.stata")= chr "%12.0g"
```

Summary statistics

Date of birth

```
##           Min.      1st Qu.      Median      Mean      3rd Qu.
## "1800-09-17" "1915-03-16" "1921-07-08" "1923-05-02" "1928-12-29"
##           Max.      NA's
## "2008-11-19"      "302"
```

Figure 1. Histogram of MDS date of birth



Year

Year	No. assessments	No. persons	%
1998	55,503	21,920	3.7
1999	131,571	30,733	8.9
2000	131,169	30,109	8.9
2001	135,622	30,423	9.2
2002	140,206	30,606	9.5
2003	142,857	30,780	9.6
2004	144,912	30,872	9.8
2005	148,893	31,598	10.1
2006	151,363	31,539	10.2
2007	151,430	31,645	10.2
2008	147,619	30,992	10.0

Sex

Gender	No. assessments	percent
Male	485,132	32.8
Female	995,466	67.2
	547	0.0

Where 1 is male and 2 is female.

MDS Race	No. assessments	%
Native-American	5,290	0.4
Asian	17,130	1.2
Black	164,435	11.1
Hispanic	49,427	3.3
White	1,238,941	83.6
	5,922	0.4

Summary

Median age date of birth is in 1921, most recent is in 2008 (highly suspicious for error), 302 missing DOB. The number of residents per year peaks in 2007. In the raw MDS file, 67% of assessments are female, 84% white, 11% Afr. American.

- Using the MDS data determine how consistent the identifying data are across records for the same persons who have multiple records. This means, among those with more than one record, what is the rate of inconsistency for date of birth, gender and race. DoB, gender and race have multiple ways in which they could disagree, in addition to calculating the rate of inconsistency, characterize the different ways (e.g. day, month or year of birth; missing categories of information in which this inconsistency manifests itself).

Basic strategy: We recode race according to the denominator definitions. The MDS dataset will be grouped by ID, a series of flags will be generated to identify disparate records among those variables expected to be fixed. Then the findings will be summarized in a set of tables.

Reformat

Consistency of data across MDS assessments

	Gender	Race cat.	Race recode	DOB	Birth Year-Month
	0: 32	0: 477	1:145706	0: 10	1:146002
	1:148822	1:146607	2: 4841	1:143600	2: 4537
	2: 1903	2: 3550	3: 194	2: 6706	3: 210
		3: 115	4: 15	3: 416	4: 7
		4: 8	5: 1	4: 22	7: 1
				5: 2	
				8: 1	

Discussion

By individual a value of 0 denotes no entries, a value of 1 or more means 1 or more distinct values for that variable (within person). There are a few missing entries but it's not the primary concern. While most only have one value, many have multiple values for the same variable (which is incoherent) and would need to be adjudicated in an analysis. In the extreme, one individual has 8 different birth dates across assessments. Birth dates may vary by a few days, or only be entered as the first of the month etc. If you recode the birth date to the month-year you may be able to collapse conflicting birthdates into the same category of month-year (**birmonyr**). This resolves about 3000 conflicted birthdates but still leaves many conflicts. Note: The summary statistics presented are by record, and multiple records per person so somewhat misleading for population inference.

Assessment of Medicare Enrollment File

Dataset

```
## Dataframe: Denominator dataset
## Memory Size: 146 Mb Rows: 4,768,016 Columns: 5
## id: 671,591 Missing: 0
```

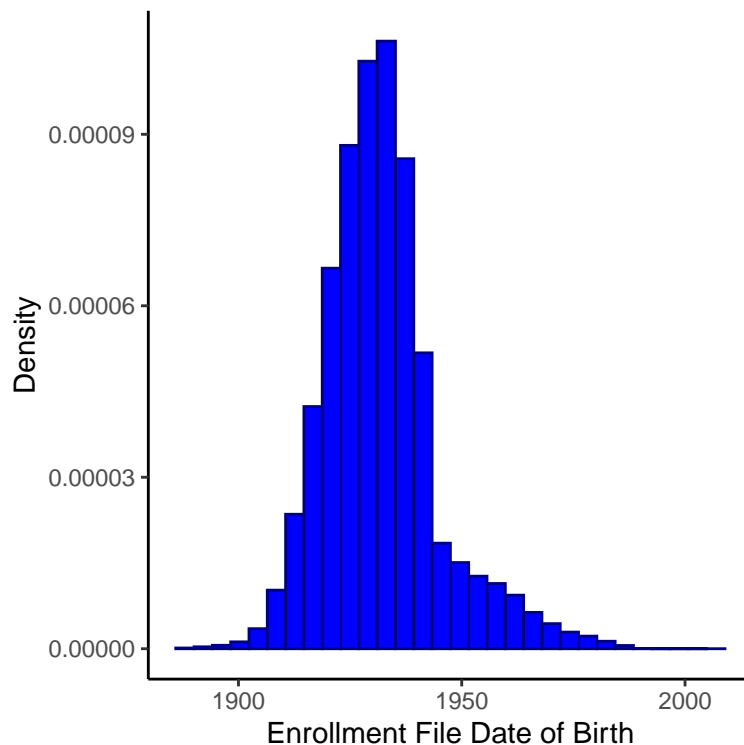
Data structure

```
## Classes 'tbl_df', 'tbl' and 'data.frame': 4768016 obs. of 4 variables:
## $ hkyear: num 2008 2008 1998 1999 2000 ...
## ..- attr(*, "label")= chr "Denom: Reference year (year of data)"
## ..- attr(*, "format.stata")= chr "%8.0g"
## $ hkdob : Date, format: "1943-11-02" "1947-02-15" ...
## $ hksex : Factor w/ 2 levels "Male","Female": 1 1 2 2 2 2 2 2 2 ...
## $ hkrace: Factor w/ 6 levels "Unknown","White",...: 2 2 2 2 2 2 2 2 2 ...
```

Summary statistics

Figure 2. Denominator file date of birth

```
##           Min.          1st Qu.          Median          Mean          3rd Qu.
## "1888-01-10" "1923-10-13" "1930-11-06" "1931-11-16" "1937-08-22"
##           Max.
## "2006-11-28"
```



Year

Year	No. records	No. persons	%
1998	403,999	403,999	8.47
1999	408,188	408,188	8.56
2000	412,847	412,847	8.66
2001	417,321	417,321	8.75
2002	422,779	422,779	8.87
2003	429,255	429,255	9.00
2004	436,078	436,078	9.15
2005	444,430	444,430	9.32
2006	453,160	453,160	9.50
2007	464,059	464,059	9.73
2008	475,900	475,900	9.98

Importantly, there are no duplicates within year. Each enrollee should have one record per year.

Year cross-over

No. years / person	No. persons	%
1	52,666	1.10
2	53,245	1.12
3	49,758	1.04
4	47,558	1.00
5	45,446	0.95
6	43,802	0.92
7	40,704	0.85
8	39,152	0.82
9	37,761	0.79
10	35,170	0.74
11	226,329	4.75

15% of the Medicare enrollees have been enrolled across all 11 years.

Sex

Denom. Sex	No. records	%
Male	2,093,428	43.91
Female	2,674,588	56.09

Denom. Race	No. records	%
Unknown	10,160	0.21
White	4,019,074	84.29
Black	464,395	9.74
Asian	76,473	1.60
Hispanic	112,571	2.36
North American Native	18,404	0.39
	66,939	1.40

Consistency of data across enrollment years

	Gender	Race cat.	DOB
	1:671591	0: 11238	1:671591
		1:660353	

Discussion. Sex, race, and date of birth do not vary within ID number (across enrollment years).

Summary description

Median date of birth is in the 1930s (i.e. the denominator file represents younger persons than the MDS assessments). The years span 1998 - 2008, with no duplicates within year. Most recent date of birth is in 2006 (highly suspicious for error). The number of enrollees increases every year. The entire Medicare population is younger, more male, and about the same % white than the nursing home population. Note: The summary statistics presented are by record, and multiple records per person so somewhat misleading for population inference.

Matching MDS and Medicare file

Match the MDS and Medicare enrollment records based upon patient ID, using the first MDS record to match for those cases that have multiple MDS records. Next, compare the DoB, gender and Race from the two different data sources. Note that DoB is necessarily measured the same way across the two data sources as is gender (although there may be differing levels of missing data). The race variables across the two data sets are coded differently so it may be important to separately estimate the degree of agreement across the categories that are comparably labeled (e.g. white vs. white; black vs. black) since the Medicare Enrollment record is known to underestimate the number of Hispanic and Asian-Americans.

Although not explicitly said above, each beneficiary may have many records in the enrollment file. Particular instruction for joining wasn't given on this dimension, however as shown above the person characteristics do not vary by year so not a practical issue.

Prepare for matching

MDS record, first row per person

```
## Dataframe: limited to first row by person, year
## Memory Size: 9 Mb  Rows: 150,757  Columns: 11
## id: 150,757 Missing: 0
```

Match 1:many join of first MDS assessment with denom. file

```
## N rows, MDS file: 150757 N rows, Denom file: 4768016
```

```
## N IDs, MDS file: 150757 N IDs, Denom file: 671591
```

```
## Dataframe: Inner join MDS / Denom.
## Memory Size: 12 Mb  Rows: 150,757  Columns: 15
## id: 150,757 Missing: 0
```

All ID's in the MDS file had at least one matching record in the enrollment file.

Comparison of MDS and denominator data

Comparison of Gender

```
##          hksex  Male Female
```



```
## idgendr
## Male      53892    829
## Female    1626  94295
## NA        42     73
```

Row percentages

```
##          hksex
## idgendr  Male Female
##   Male   0.985 0.015
##   Female 0.017 0.983
##   <NA>   0.365 0.635
```

The MDS variable 'idgendr' is missing in a few cases, and is discordant in about 1.5% of cases. 1.5% where the MDS says male and enrollment says female, and 1.7% where MDS says female and enrollment says male. Note: 'NA' means missing.

Comparison of Race

```
##          hkrace Unknown  White  Black  Asian Hispanic North American Native  NA
## rc_recode
## Unknown          1    745    106    35      23          5    11
## White          395 126272   1039   279    317        182   378
## Black           44    499  12762    12     33         18    68
## Asian           14    160    36   1110    12         2   514
## Hispanic        18   2862   120    23   1873        15   222
## North American Native  1   273    32    15     9        192   30
```

Note: The MDS race/ethnicity variable has been recoded to match categories with the enrollment file race variable. MDS missing was recoded to 0, unknown. The denominator file has a category, 'other' with no comparable category in the MDS assessment. 0 - Unknown, 1 - White, 2 - Black, 3 - Other, 4 - Asian, 5 - Hispanic, 6 - Native American.

Row percentages

```
##          hkrace
## rc_recode  Unknown White Black Asian Hispanic
##   Unknown   0.001 0.805 0.114 0.038   0.025
##   White     0.003 0.980 0.008 0.002   0.002
##   Black     0.003 0.037 0.950 0.001   0.002
##   Asian     0.008 0.087 0.019 0.601   0.006
##   Hispanic  0.004 0.558 0.023 0.004   0.365
##   North American Native 0.002 0.495 0.058 0.027   0.016

##          hkrace
## rc_recode  North American Native  <NA>
##   Unknown          0.005 0.012
##   White            0.001 0.003
##   Black            0.001 0.005
##   Asian            0.001 0.278
##   Hispanic         0.003 0.043
##   North American Native 0.348 0.054
```

There is 98% agreement between the MDS and denominator on white race (category 1), 95% on black race (cat. 2), 60% agreement on asian race (cat. 4), and poor agreement on hispanic or native american race (cat. 5, 6). Many of the individuals coded as hispanic or native american in the MDS, are listed white in the denominator file.

Comparison of Date of Birth

Match by exact date of birth

DOB match	No. records	%
0	14,284	9.47
1	136,413	90.49
	60	0.04

The date of birth exactly matches between MDS and enrollment in 90% of cases.

Match by same month-year of birth

Month of birth match	No. records	%
0	8,345	5.54
1	142,352	94.42
	60	0.04

The month and year of birth exactly matches in 94% of cases, increasing match rate over exact date.

Evaluate a systematic vs. random error in date-match

Understanding the dates do not match in a small subset of cases, one question would be are the unmatched dates randomly different, or does one date consistently run over- under- the other.

MDS DOB > Denom. DOB	No. records	%
0	143,113	94.93
1	7,584	5.03
	60	0.04

MDS DOB < Denom. DOB	No. records	%
0	143,997	95.52
1	6,700	4.44
	60	0.04

4.4% of MDS DOBs are less than the denominator file DOB, and 5% are more than the denominator DOB. So if the date of birth from CMS enrollment records is considered more valid, the MDS may be biased towards moderately younger (i.e. later) dates of birth than the actual.

Summary discussion

In this assignment, two limited data files, 1) Nursing home minimum dataset assessments and 2) The Centers for Medicaid and Medicare Services enrollment file ('denominator') were summarized and joined. The primary objective was to compare how consistent reported gender/sex, race and date of birth are between the two files. To simplify the exercise the first assesment by nursing home resident was used. The results reported above show general agreement between the two files, but some notable discrepancies for minority race categories, particularly hispanics, asians and native-americans. Dates of birth are consistent for 90%

of cases outright, and the match can be improved by allowing matches by month-year instead. Overall the MDS assessments are imperfect but in reasonable agreement with the denominator file.

Session Info

Thank you for taking the time to review my work!

```
## setting value
## version R version 3.6.1 (2019-07-05)
## os Windows 10 x64
## system x86_64, mingw32
## ui RTerm
## language (EN)
## collate English_United States.1252
## ctype English_United States.1252
## tz America/New_York
## date 2019-09-27
```