

PHPM672: Data Science for Health Services Research PHPM677: Data Science in Public Health

Hye-Chung Kum (kum@tamu.edu)

Associate Professor

Population Informatics Lab (<https://pinformatics.org/>)

Course URL: <http://pinformatics.org/phpm672>

License:
Data Science for Health by Hye-Chung Kum is licensed under a
Creative Commons Attribution-NonCommercial-ShareAlike 4.0 International License



1

Assignment 3: Arrays & Loops



2

Reshaping & Combining Tables

Unit of analysis

Combining

set: concatenate tables (stack rows)

merge: link tables (attach columns)

Reshaping

proc summary: consolidate rows

proc transpose: reshape table



3

Few things...

4

subset obs

- P1.2 subset obs;
 - How? Why?

```
/**P1.2 subset obs**/
data diabetes;
set diabetes;
if age= "[0-10]" or age="[10-20]" then delete;
if race= "?" then delete;
if gender="Unknown/Invalid" then delete;
```

File I/O

- <https://pinformatics.org/resources/sas-sample-code.php>

```
output
infile/input/datalines
proc import / proc export
libname name xport 'folder location' ;
```

5

Assignment 4

- Concatenate multiple tables (more rows)
 - stack tables on top of each other to increase the number of rows
 - using **set**
 - Be sure to understand the different behavior given different situations (i.e., what happens to shared variables? What happens to not shared variables?)
- Link up multiple tables using a shared key (more columns)
 - align the rows using the shared key, and link multiple tables to increase the number of variables in the tables
 - using **merge**
 - Be sure to understand the different behavior given different situations (i.e., what happens to shared vars? What happens to not shared vars?)
 - What is a 1-to-1 link
 - What is a 1-to-N link
 - What is a N-to-N link (you will not be doing this, but need to understand what this is. This must be done with proc sql in SAS)
- New keyword **in=**

Data Science Knowledge Discovery & Data mining (KDD)

Table Operations: multiple table → 1 table

- set (Append)**

Table A

Table B

→

Table A
Table B
- merge (link)**

Table A

Table B

→

Table A	Table B
---------	---------

Table Operations: 1 table → 1 table (reshaping)

- Proc Transpose**

1	2
a	d
b	e
c	f

→

1	a	b	c
2	d	e	f
- Proc Summary**

A
B
C

→

D

Where D= function(A,B,C)
Examples of function are
Sum(A,B,C) Mean(A,B,C) Max(A,B,C) Min(A,B,C)

Assignment 4 continued

- Combine multiple rows into one row
 - by group processing **proc summary**
- Reshape table to flip rows & columns
 - using **proc transpose**
 - Also transpose (flip rows & columns) by groups or row



UNIT OF ANALYSIS

12

Basic Regression



- $y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \epsilon$
- y : dependent variable
- x_i : independent variables
 - β_i : coefficient
- ϵ : error term

14

Unit of analysis

- $y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \epsilon$
- Table
 - column: y, x_1, x_2
 - row: ? (unit of analysis)
- What is unit of y/x ?
 - DV: capacity of hospital (unit: ?)
 - DV: service use (unit: ?)

15

Unit of analysis

- $y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \epsilon$
- Table
 - column: y, x_1, x_2
 - row: ? (unit of analysis)
- What is unit of y/x ?
 - DV: capacity of hospital (unit: hospital)
 - DV: service use (unit: patient=person)

16

Reshaping to correct unit



- What do you have?
- What do you want? (unit of analysis)

17

Example

- Flu data
 - Weekly estimates
- NSDUH
 - Person
- Tx Discharge Data
 - Per hospital

18

Converting to the desired unit

- Consolidating multiple rows
 - Flu: Weekly estimates to monthly estimates
 - NSDUH: Per person to per race
 - Tx Discharge: Per hospital to per region
- Transposing: changing row/column
 - Flu: Weekly estimates to estimates per state
 - Tx Discharge: Per hospital to per hospital year

20

Consolidating multiple rows

- Must first determine how to consolidate
 - Sum, max, min, count (of nonmissing) etc
 - Think about each variable and decide on the correct method [per variable](#)
- **MUST be sorted first by the by varlist**
- Example
 - Flu: SUM - Weekly estimates to monthly estimates
 - NSDUH: MEAN - Per person to per race
 - Tx Discharge: SUM- Per hospital to per region

21

proc summary (try it)

```
proc sort data= srofn [out= fn nodupkey];
  by byvar1 byvar2 ...;
run;

proc summary data= fn;
  [by byvar1 byvar2 ...];
  var var1 var2 ...;
  output out= outfn (drop=_type_ _name);
run;

proc summary data= fn;
  [by byvar1 byvar2 ...];
  var var1 var2 ...;
  output out= outfn (drop=_type_);
  sum var1 = outvar1;
  mean var2 = outvar2;
run;
```

22

Transposing: changing row/column

- Must first determine unit of transpose
 - Per time period
- **MUST be**
 - sorted first by the by varlist (unit of transpose)
 - one row per unit
- Example
 - Flu: Weekly estimates to estimates per state
 - Full table
 - Tx Discharge: Per hospital to per hospital year
 - Group transpose

22

proc transpose (try it)

```
proc sort data= srofn [out= fn] nodupkey;
  by byvar1 byvar2 ...;
run;

proc transpose data= fn out= outfn [prefix=prefix];
  [by byvar1 byvar2 ...];
  var var1 var2 ...;
  id idvar;
run;
```

23

Lab 4

- Lab 4 (2 pts): Due in 1 week
 - Learn how each command behaves
 - Submit excel file with answers
 - Will post answer one week from now
 - Will be on midterm
- Midpoint email (1 pt): Due in 1 week
 - Separate from lab
 - Must have started the assignment to answer
 - Review together

24

Lab 4: midpoint email (answer questions) SEPARATE from Lab



- Describe in one sentence, what each of the tables are (there is a total of 8).
- What is the unit (row) of each table?
- For each table that does not have the required unit of analysis as “county year”, explain how you will convert the given table into the required “county year” table. If not applicable write NA.
- When linking up all the tables to have all the variables in one table,
 - Which tables link up as 1-to-1 matching? What are the matching variables?
 - Which tables link up as 1-to-N matching? What are the matching variables?

20

Assignment 4 (9 pts)



- Most difficult
 - Covers ALL topics we have done so far. (final grade: 12)
 - Assignment 5: extension to assignment 4 (4 pt)
 - You have to think about what task is required, and then which commands to use
 - 4 weeks (2/20-3/20): spring break in the middle
- Look at the assignment together

26

Reminder



- Read the required readings
- Do the lab this week to learn the behavior of each command
 - Set
 - Merge
 - Proc summary
 - Proc transpose

27