Managing and Manipulating Data Using R Lecture 5, Survey data and exploratory data analysis (for data quality)

Ozan Jaquette

- 1. Introduction
- 2. haven and labelled package
- 3. Exploratory data analysis (EDA)
- 4. Brainstorm next assignment: create GPA from course-level data



Libraries we will use today [install if you don't have them]

Data we will use today

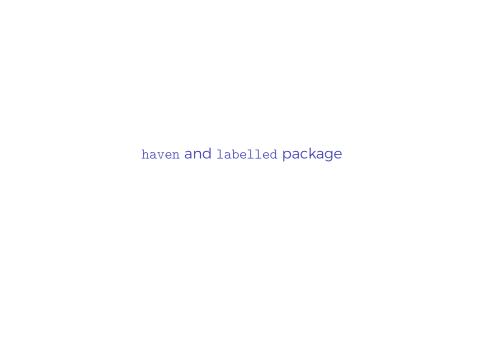
High school longitudinal surveys from National Center for Education Statistics (NCES)

o Follow U.S. students from high school through college, labor market

We will be working with High School Longitudinal Study of 2009 (HSLS:09)

- o Follows 9th graders from 2009
- Data collection waves
 - ▶ Base Year (2009)
 - ▶ First Follow-up (2012)

 - ▶ High School Transcripts (2013-2014)
 - Second Follow-up (2016)



haven package

haven, which is part of **tidyverse**, "enables R to read and write various data formats" from the following statistical packages:

- SAS
- SPSS
- Stata

When using haven to read data, resulting R objects have these characteristics:

- o Are **tibbles**, a particular type of data frame we discuss future weeks
- Transform variables with "value labels" into the labelled() class [our focus today]
 - Helpful description HERE
- Dates and times converted to R date/time classes
- Character vectors not converted to factors

haven package

Use read_dta() function from haven to import Stata dataset into R

```
hsls <- read_dta(file="../../data/hsls/hsls_stu_small.dta")
```

Let's examine the data

```
names(hsls)
names(hsls) <- tolower(names(hsls)) # convert names to lowercase
str(hsls)
str(hsls$s3classes)</pre>
```

labelled package

labelled package is to work with data imported from SPSS/Stata/SAS using the haven package.

 Specifically, "The purpose of the labelled package is to provide functions to manipulate metadata as variable labels, value labels and defined missing values using the labelled class and the label attribute introduced in haven package. LINK

Functions in labelled package

- Full list
- A couple relevant functions
 - ▶ val labels: get or set variable value labels
 - var_label: get or set a variable label

```
hsls %>% select(s3classes) %>% var_label hsls %>% select(s3classes) %>% val_labels
```

Understanding labelled data

First, let's review core concepts:

atomic vectors (and lists) the underlying data

- o data structures: vector or list
- o data type: numeric (integer or double); character; logical

```
typeof(hsls$s3classes)
#> [1] "double"
```

augmented vectors are atomic vectors with attributes attached

attributes are "metadata" attached to an object. Examples

- names: names of elements of a vector or list (e.g., variable names)
- o **levels**: display output associated with values of a factor variable
- o class: e.g., factor, labelled

```
attributes(hsls$s3classes)
```

class is an object oriented programming concept. The class of an object determines which functions can be applied to the object and what those functions do

o e.g., can't apply sum() to an object where class=character

Understanding labelled data

Let's investigate the attributes of hsls\$s3classes

```
typeof(hsls$s3classes)
class(hsls$s3classes)
str(hsls$s3classes)
attributes(hsls$s3classes)

#use attr(object_name, "attribute_name") to refer to each attribute
attr(hsls$s3classes, "label")
attr(hsls$s3classes, "labels")
attr(hsls$s3classes, "class")
attr(hsls$s3classes, "format.stata")
```

Understanding labelled data

What is class==labelled?

- An object class created by the haven package for importing variables from SAS/SPSS/Stata that have value labels
- o value labels [in Stata] are labels attached to specific values of a variable:
 - ▶ e.g., variable value 1 attached to value label "married", 2="single", 3="divorced"
- Variables in an R data frame with class==labelled:
 - ▶ data type can be numeric(double) or character
 - The value labels associated with each value: attr(data frame name\$variable name."labels")
 - ▶ In filter() refer to variable value, not the value label

Working with class==labelled variables

```
#show variable label
hsls %>% select(s3classes,s3clglvl) %>% var_label
#show value labels
hsls %>% select(s3classes,s3clglvl) %>% val_labels
#Frequency table
hsls %>% select(s3classes) %>% count(s3classes)
#Frequency table, value labels instead of variable values
hsls %>% select(s3classes) %>% count(s3classes) %>% as_factor()
#in filters, refer to variable value, not value label
hsls %>% select(s3classes) %>% filter(s3classes==-8) %>% count(s3classes)
```

Converting class==labelled to class==factor

The $as_factor()$ function from haven package converts variables with class=labelled to class=factor

Can be used for descriptive statistics

```
hsls %>% select(s3classes) %>% count(s3classes) %>% as_factor()
```

Can create object with some or all labelled vars converted to factor

```
hsls_f <- as_factor(hsls,only_labelled = TRUE)
```

Let's examine this object

```
glimpse(hsls_f)
hsls_f %>% select(s3classes,s3clglvl) %>% str()
typeof(hsls_f$s3classes)
class(hsls_f$s3classes)
attributes(hsls_f$s3classes)
hsls_f %>% select(s3classes) %>% var_label()
hsls_f %>% select(s3classes) %>% val_labels()
```

Working with class==factor data

Showing values associated with factor levels

In code, refer level attribute not variable value

Comparing class==labelled to class==factor

	class==labelled	class==factor
data type	numeric or character	integer
name of value label attribute	labels	levels
refer to data using	variable values	levels attribute





What is exploratory data analysis (EDA)?

The Towards Data Science website has a nice definition of EDA:

"Exploratory Data Analysis refers to the critical process of performing initial investigations on data so as to discover patterns, to spot anomalies, to test hypothesis and to check assumptions with the help of summary statistics and graphical representations."

This course focuses on "data management":

- investigating and cleaning data for the purpose of creating analysis variables
- o Basically, everything that happens before you conduct analyses

I think about "Exploratory data analysis for data quality"

- o Investigating values and patterns of variables from "input data"
- o Identifying and cleaning errors or values that need to be changed
- Creating analysis variables
- Checking values of analysis variables agains values of input variables

Tools of EDA

To do EDA for data quality, must master the following tools:

- o Select, sort, filter, and print
 - Select and sort particular values of particular variables
 - Print particular values of particular variables
- One-way descriptive analyses (i.e,. focus on one variable)
 - Descriptive analyses for continuous variables
 - Descriptive analyses for discreet/categorical variables
- Two-way descriptive analyses (relationship between two variables)
 - Categorical by categorical
 - Categorical by continuous
 - Continuous by continuous

Whenever using any of these tools, pay close attention to missing values and how they are coded

I'll focus on the tidyverse approach rather than base R

Tools of EDA: select, sort, filter, and print Let's create a smaller version of the HSLS:09 dataset

```
hsls %>% var_label()
#> $stu id
#> [1] "Student ID"
#>
#> $sch id
#> [1] "School ID"
#>
#> $x2univ1
#> [1] "X2 Sample member status in BY and F1 rounds"
#>
#> $x2univ2a
#> [1] "X2 Base year status and how sample member entered F1 sample"
#>
#> $x2univ2b
#> [1] "X2 Sample member F1 status"
#>
#> $x3univ1
#> [1] "X3 Sample member status in BY, F1, U13, and HS transcript rounds"
#>
#> $x4univ1
#> [1] "X4 Sample member status in BY, F1, U13, HS transcript, and F2 rounds"
#>
#> $s3classes
#> [1] "S3 B01A Taking postsecondary classes as of Nov 1 2013"
#>
```

Tools of EDA: select, sort, filter, and print

We've already know select(), arrange(), filter()

Select, sort, and print specific vars

```
hsls_small %>% arrange(desc(stu_id)) %>%
select(stu_id,x3univ1,x3sqstat,s3classes,s3clglv1)

#print value labels
hsls_small %>% arrange(desc(stu_id)) %>%
select(stu_id,x3univ1,x3sqstat,s3classes,s3clglv1) %>% as_factor()
```

Sometimes helpful to increase the number of observations printed

```
class(hsls_small) #it's a tibble
options(tibble.print_min=50)
# execute this in console
hsls_small %>% arrange(desc(stu_id)) %>% select(stu_id,x3univ1,x3sqstat,s3classe
options(tibble.print_min=10) # set default printing back to 10 lines
```

One-way descriptive stats for continuous vars, Base R approach [SKIP]

```
mean(hsls_small$x2txmtscor)
sd(hsls_small$x2txmtscor)

#Careful: summary stats include value of -8!
min(hsls_small$x2txmtscor)
max(hsls_small$x2txmtscor)
```

Be careful with NA values

```
#Create variable replacing -8 with NA
hsls_small_temp <- hsls_small %>%
  mutate(x2txmtscorv2=ifelse(x2txmtscor==-8,NA,x2txmtscor))
hsls_small_temp %>% filter(is.na(x2txmtscorv2)) %>% count(x2txmtscorv2)
mean(hsls_small_temp$x2txmtscorv2)
mean(hsls_small_temp$x2txmtscorv2, na.rm=TRUE)
rm(hsls_small_temp)
```

One-way descriptive stats for continuous vars, Tidyverse approach

Use ${\tt summarise_at()}$ which is a variation of ${\tt summarise()}$ to calculate descriptive stats

explain .args=list(na.rm=TRUE) on next slide

One-way descriptive stats for continuous vars, Tidyverse approach

"Input vars" in survey data often have negative values for missing/skips
 R includes those negative values when calculating statistics, which you probably don't want

Solution: create version of variable that replaces negative values with NAs

What if you didn't include .args=list(na.rm=TRUE)?

One-way descriptive stats for continuous vars, Tidyverse approach

How to identify these missing/skip values if you don't have a codebook?

 ${\tt count}$ () combined with ${\tt filter}$ () helpful for finding extreme values of continuous vars, which are often associated with missing or skip



One-way descriptive stats for discrete/categorical vars, Tidyverse approach

Use count() to investigate values of discreet or categorical variables

For variables where class==labelled

```
class(hsls_small$s3classes)
#show counts of variable values
hsls_small %>% count(s3classes)
#show counts of value labels
hsls_small %>% count(s3classes) %>% as_factor()
```

o I like count() because the default setting is to show NA values too!

```
hsls_small %>% mutate(s3classes_na=ifelse(s3classes<0,NA,s3classes)) %>% count(s3classes_na)
```

Show both values and value labels on count tables for class==labelled

```
#CRYSTAL'S SOLUTION
x <- hsls_small %>% count(s3classes)
y <- hsls_small %>% count(s3classes) %>% as_factor()
bind_cols(x[,1], y)
```

One-way descriptive stats for discrete/categorical vars, Tidyverse approach

For variables where class==factor [PROBLEM: HOW TO RETAIN FACTOR LEVELS AFTER MUTATE]

```
#use variable from the hsls data frame where vars are factors
class(hsls_f$s3classes)
attributes(hsls_f$s3classes)

#show frequency table
hsls_f %>% count(s3classes)

#frequency table with NAs
    #note: within ifelse() used levels(s3classes)[s3classes]) rather than s3classes
hsls_f %>% mutate(s3classes_f=ifelse(s3classes %in% c("Missing","Unit non-respondence count(s3classes_f))
```

Relationship between variables, categorical by categorical

Two-way frequency table, sometimes called "Cross tabulation", is important for checking data quality - When you create categorical analysis var from single categorical "input" var - Two-way tables show us whether we did this correctly - Two-way tables helpful for understanding skip patterns in surveys

Task: Create a two-way table between s3classes and s3clglvl

```
hsls_small %>% select(s3classes,s3clglvl) %>% var_label()
hsls_small %>% group_by(s3classes) %>% count(s3clglvl)
hsls_small %>% group_by(s3classes) %>% count(s3clglvl) %>% as_factor()
```

What if one of the variables has NAs?

Table created by group_by() and count() shows NAs!

Relationship between variables, categorical by categorical

Tables above are pretty ugly

Use the ${\tt spread}()$ function from tidyr package to create table with one variable as columns and the other variable as rows

The variable you place in spread() will be columns

```
hsls_small %>% group_by(s3classes) %>% count(s3clglvl) %>% spread(s3classes, n)

hsls_small %>% group_by(s3classes) %>% count(s3clglvl) %>% as_factor() %>% spread(s3classes, n)

hsls_small %>% group_by(s3classes) %>% count(s3clglvl) %>% as_factor() %>% spread(s3clglvl, n)
```

Relationship between variables, categorical by continuous

Investigating relationship between multiple variables is a little tougher when at least one of the variables is continuous

One approach is the conditional mean:

- Shows average values of continous variables within groups
- Groups are defined by your categorical variable(s)

Relationship to regression

 Conditional mean is similar to regression with a continuous dependent variable and a categorical X variable

Task:

 Investigate the relationship between math test score, x2txmtscor, and parental education, x2paredu

Relationship between variables, categorical by continuous

Task: Investigate the relationship between math test score, x2txmtscor, and parental education, x2paredu

For checking data quality, helpful to calculate other stats besides mean

Always Investigate presence of missing/skip values

```
hsls_small %>% filter(x2paredu<0) %>% count(x2paredu)
hsls_small %>% filter(x2txmtscor<0) %>% count(x2txmtscor)
```

Replace -8 with NA and re-calculate conditional stats

Relationship between variables, categorical by continuous

[MAKE THIS A STUDENT EXERCISE?]

Can use same approach to calculate conditional mean by multiple $group_by()$ variables

Just add additional variables within group_by()

Task: Calculate mean math test score (x2txmtscor), for each combination of parental education (x2paredu) and sex (x2sex)



Guidelines for "EDA for data quality"

Assme that your goal in "EDA for data quality" is to investigate "input" data sources and create "analysis variables"

 Usually, your analysis dataset will incorporate multiple sources of input data, including data you collect (primary data) and/or data collected by others (secondary data)

While this is not a linear process, these are the broad steps I follow

- 1. Understand how inout data sources were created
 - e.g., when working with survey data, have survey questionnaire and codebooks on hand
- 2. For each input data source, identify the "unit of analysis" and which combination of variables uniquely identify observations
- 3. Investigate patterns in input variables
- 4. Create analysis variable from input variable(s)
- Verify that analysis variable is created correctly through descriptive statistics that compare values of input variable(s) against values of the analysis variable

Always be aware of missing values

"Unit of analysis" and which variables uniquely identify observations

"Unit of analysis" refers to "what does each observation represent" in an input data source

- o If each obs represents a student, you have "student level data"
- If each obs represents a student-course, you have "student-course level data"
- o If each obs represents a school, you have "school-level data"
- o If each obs represents a school-year, you have "school-year level data"

How to identify unit of analysis

- o data documentation
- o investigating the data set

"Unit of analysis" and which variables uniquely identify observations

Identify the variable - or group of vars - that "uniquely identifies" observations

 "uniquely identifies observations": each value of the var has a frequency count of 1

This is important for many data management tasks

o e.g., merging data sources, "tidying" data

How to identify which variable(s) uniquely identify observations

```
data documentationinvestigating the data set
```

```
hsls_small %>% group_by(stu_id) %>% # group_by our candidate

mutate(n_per_id=n()) %>% # calculate number of obs per group

ungroup() %>% # ungroup the data

count(n_per_id==1) # count "true that only one obs per group"

#> # A tibble: 1 x 2

#> `n_per_id == 1` n

#> <lg!> <int>
#> 1 TRUE 23503

#hsls_small %>% count(stu_id) %>% filter(n> 1) ## Patricia's approach

#Karina approaches using asserthat package

#library(assertthat)

#stopifnot(length(unique(uga_pub$ncessch))==nrow(uga_pub))

#assert that(any(duplicated(uga pub, by=c("ncessch", "var2")))==FALSE)
```

Rules for variable creation

Rules I follow for variable creation

- 1. Never modify "input variable"; instead create new variable based on input variable(s)
 - Always keep input variables used to create new variables
- 2. Investigate input variable(s) and relationship between input variables
- 3. Developing a plan for creation of analysis variable
 - ▷ e.g., for each possible value of input variables, what should value of analysis variable be?
- 4. Write code to create analysis variable
- 5. Run descriptive checks to verify new variables are constructed correctly
 - ▶ Can "comment out" these checks, but don't delete them
- 6. Document new variables with notes and labels

Rules for variable creation

Task: Create analysis fir ses qunitile called sesq5 based on x4x2sesq5

```
#investigate input variable
hsls small %>% select(x4x2sesg5) %>% var_label()
hsls small %>% select(x4x2sesq5) %>% val_labels()
hsls_small %>% select(x4x2sesq5) %>% count(x4x2sesq5)
hsls small %% select(x4x2sesq5) %>% count(x4x2sesq5) %>% as_factor()
#create analysis variable
hsls small temp <- hsls small %>%
  mutate(sesq5=ifelse(x4x2sesq5==-8,NA,x4x2sesq5)) # approach 1
hsls small temp <- hsls small %>%
  mutate(sesq5=ifelse(x4x2sesq5<0,NA,x4x2sesq5)) # approach 1</pre>
#verify
hsls_small_temp %>% group_by(x4x2sesq5) %>% count(sesq5)
```

How to be a good researcher, good research assistant (RA) [CUT?]

[keep or cut this section?][hidden curriculum stuff; look at notes from last time you taught Stata data mgt]

- \circ Your advisor/PI getting $\$ to support students depends on previous RAs doing a good job
- Must master certain skills, but beyond that it is mindset and approach that matters



What are skip patterns

Pretty easy to create an analysis variable based on a single input variable

Harder to create analysis variables based on multiple input variables

 When working with survey data, even seemingly simple analysis variables require multiple input variables due to "skip patterns"

What are "skip patterns"? [students answer]

- [?DELETE OR DELAY?] Response on a particular survey item determines whether respondent answers some set of subsequent questions
- What are some examples of this?

Key to working with skip patterns

- Have the survey questionnaire on hand
- Sometimes it appears that analysis variable requires only one input variable, but really depends on several input variables because of skip patterns
 - Don't just blindly turn "missing" and "skips" from survey data to NAs in your analysis variable
 - Rahter, trace why these "missing" and "skips" appear and decide how they should be coded in your analysis variable

Creating analysis variables in the presence of skip patterns

Task: Create a measure of "level" of postsecondary institution attended in 2013 from HSLS:09 survey data

- o "level" is highest award-level of the postsecondary institution
 - ▶ e.g., if highest award is associate's degree (a two-year degree), then 'level==2'
- o The measure, pselev2013, should have following [non-missing] values:
 - 1. Not attending postsecondary education institution
 - 2. Attending a 2-year or less-than-2-year institution
 - 3. Attending 4-year or greater-than-4year institution

Background info:

- o In "2013 Update" of HSLS:09, students asked about college attendance
 - ▶ Variables from student responses to "2013 Update" have prefix s3
- Survey guestionnaire for 2013 update can be found HERE
- o The "online codebook" website HERE has info about specific variables
- o Measure has 3 input variables [usually must figure this out yourself]:
 - 1. x3sqstat: "X3 Student questionnaire status"
 - 2. s3classes: "S3 B01A Taking postsecondary classes as of Nov 1 2013"
 - 3. s3clglvl: "S3 Enrolled college IPEDS level"

hsls_small %>% select(x3sqstat,s3classes,s3clglvl) %>% var_label()

You won't have time to complete this task, but develop a plan for the task and get as far as you can

Creating analysis variables in the presence of skip patterns

Step 1a: Investigate each input variable separately

```
#variable labels
hsls_small %>% select(x3sqstat,s3classes,s3clglvl) %>% var_label()
hsls_small %>% count(x3sqstat)
hsls_small %>% count(x3sqstat) %>% as_factor()
hsls_small %>% count(s3classes)
hsls_small %>% count(s3classes) %>% as_factor()
hsls_small %>% count(s3clglvl)
hsls_small %>% count(s3clglvl) %>% as_factor()
```

Creating analysis variables in the presence of skip patterns

Step 1b: Investigate relationship between input variables

```
#x3sqstate and s3classes
hsls small %>% group_by(x3sqstat) %>% count(s3classes)
hsls_small %>% group_by(x3sqstat) %>% count(s3classes) %>% as_factor()
hsls small %>% filter(x3sqstat==8) %>% count(s3classes)
hsls_small %>% filter(x3sqstat==8) %>% count(s3classes==-8)
hsls small %>% filter(x3sqstat !=8) %>% count(s3classes)
#x3sqstate, s3classes and s3clglvl
hsls small %>% group_by(s3classes) %>% count(s3clglvl)
hsls_small %>% group_by(s3classes) %>% count(s3clglvl) %>% as_factor()
#add filter for whether student did not respond to X3 questionnaire
hsls_small %>% filter(x3sqstat==8) %>% group_by(s3classes) %>% count(s3clglvl)
hsls small %>% filter(x3sqstat !=8) %>% group_by(s3classes) %>% count(s3clglvl)
#add filter for s3classes is "missing" [-9]
hsls small %% filter(x3sqstat !=8,s3classes==-9) %% group_by(s3classes) %>% co
hsls small %>% filter(x3sqstat !=8,s3classes!=-9) %>% group_by(s3classes) %>% co
#add filter for s3classes equal to "no" or "don't know"
hsls small %>% filter(x3sqstat !=8,s3classes!=-9, s3classes %in% c(2,3)) %>% gro
hsls small %% filter(x3sqstat !=8,s3classes!=-9, s3classes %in% c(2,3)) %>% gro
hsls_small %>% filter(x3sqstat !=8,s3classes!=-9, s3classes==1) %>% group_by(s3c
```

Brainstorm next assignment: create GPA from course-level data

Brainstorm in your homework groups for next assignment

Assignment: create GPA variables from student-course level data

Link to assignment HERE

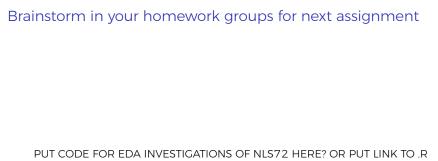
Data source: National Longitudinal Study of 1972 (NLS72)

- o Follows 12th graders from 1972
- Data collection waves
 - ▶ Base year: 1972
 - Follow-up surveys in: 1973, 1974, 1976, 1979, 1986
 - Postsecondary transcripts collected in 1984
- Why use such an old survey for this assignment?
 - NLS72 predates data privacy agreements so postsecondary transcript data are publicly available

#I.OAD DATA

Work on the following in your homework groups:

- Read assignment
- Conduct EDA investigations of input data
 - For homework assignment, you only need to focus on these variables [and can ignore the others]
- Develop a plan for how you will go about creating GPA variables



PUT CODE FOR EDA INVESTIGATIONS OF NLS72 HERE? OR PUT LINK TO .F SCRIPT W/ NLS EDA INVESTIGATIONS HERE?