

# Module 8 Problem Set

INSERT YOUR NAME HERE

## Contents

<b>Instructions</b>	<b>1</b>
Overview . . . . .	1
<b>Load library and data</b>	<b>1</b>
<b>Part I: Conceptual questions</b>	<b>2</b>
<b>Part II: Questions about reshaping long to wide</b>	<b>2</b>
Description of the data . . . . .	2
Overview of the reshaping long to wide tasks . . . . .	2
Load data and create three new data frames . . . . .	2
Questions related to reshaping the dataset <code>agegroup1_obs</code> from long to wide . . . . .	5
Questions related to reshaping the dataset <code>levstudy1_obs</code> from long to wide . . . . .	6
<b>Part III: Questions about reshaping wide to long</b>	<b>6</b>

## Instructions

### Overview

This problem set has three parts.

1. I'll ask you some definition/conceptual questions about the concepts introduced in lecture
2. Tidying untidy data: reshaping from long to wide
  - e.g., dataset has one row for each combination of university ID and enrollment age group, but you want a dataset with one row per university ID and one enrollment variable for each age group
  - for these questions we'll use fall enrollment data from the Integrated Postsecondary Data System (IPEDS), specifically the fall enrollment sub-survey that focuses on enrollment by age group
3. Tidying untidy data: reshaping from wide to long
  - for these questions we'll use data from the NCES digest of education statistics that contains data about the total number of teachers in each state

## Load library and data

```
library(tidyr)
library(tidyverse)
#> -- Attaching core tidyverse packages ----- tidyverse 2.0.0 --
#> v dplyr      1.1.4      v purrr      1.0.2
#> v forcats    1.0.0      v readr      2.1.5
#> v ggplot2    3.5.1      v stringr    1.5.1
#> v lubridate  1.9.4      v tibble     3.2.1
```

```
#> -- Conflicts ----- tidyverse_conflicts() --
#> x dplyr::filter() masks stats::filter()
#> x dplyr::lag() masks stats::lag()
#> i Use the conflicted package (<http://conflicted.r-lib.org/>) to force all conflicts to become errors
library(haven)
library(labelled)
```

## Part I: Conceptual questions

- What is the difference between the terms “unit of analysis” [our term; not necessarily used outside this class] and “observational level” [A Wickham term]?
  - ANSWER:
- What are the three rules of tidy data?
  - ANSWER:

## Part II: Questions about reshaping long to wide

### Description of the data

For these questions, we’ll be using data from the Fall Enrollment survey component of the Integrated Postsecondary Education Data System (IPEDS)

- Specifically, we’ll be using data from the survey sub-component that focuses on enrollment by age-group.
- The dataset we’ll be using data from Fall 2016 (i.e., Fall of the 2016-17 academic year)
- Here is a link to a data dictionary (an excel file) for the enrollment by age dataset: [LINK](#)
- In the dataset you load below:
  - I’ve dropped a few of the variables from the raw enrollment by age data
  - I’ve added a few variables from the “institutional characteristics” survey (e.g., institution name, state, sector) that should be pretty self explanatory if you examine the variable labels and/or value labels
- the variable `unitid` is the ID variable for each college/university
- the dataset has one observation for each combination of the variables `unitid-efbage-lstudy`; in other words the unit of analysis is university per age group per level of study

### Overview of the reshaping long to wide tasks

- We will load the data frame via `read_dta` using the hyperlink and assign it the name `age_f16_allvars_allobs`
- Then, we’ll create two different data frame objects based on the data frame `age_f16_allvars_allobs`
  - A dataframe `agegroup1_obs` that has fewer variables than `age_f16_allvars_allobs` and keeps observations where age-group equals 1 (1. All age categories total)
    - \* this data frame has the simplest structure; we’ll reshape this one first
  - A dataframe `levstudy1_obs` that has fewer variables than `age_f16_allvars_allobs` and keeps observations where “level of study” equals 1 (1. All Students total)
    - \* we’ll reshape this one second
- Questions related to reshaping `agegroup1_obs`
- Questions related to reshaping `levstudy1_obs`

### Load data and create three new data frames

- Load IPEDS data that contains fall enrollment by age

NOTE: IN THIS QUESTION, WE GIVE YOU THE ANSWERS; ALL YOU HAVE TO DO IS RUN THE BELOW CODE CHUNK

```
rm(list = ls()) # remove all objects
#getwd()
#list.files("../../../documents/rclass/data/ipeds/ef/age") # list files in directory w/ NLS data

#Read Stata data into R using read_data() function from haven package
age_f16_allvars_allobs <- read_dta(file="https://github.com/ksalazar3/HED696c_Rclass/raw/master/data/ipeds/age_f16_allvars_allobs.dta")

#rename a couple variables
age_f16_allvars_allobs <- age_f16_allvars_allobs %>% rename(agegroup=efbage, levstudy=lstudy)

#list variables and variable labels
names(age_f16_allvars_allobs)
#> [1] "unitid"      "agegroup"    "levstudy"    "efage01"     "efage02"
#> [6] "efage03"     "efage04"     "efage05"     "efage06"     "efage07"
#> [11] "efage08"     "efage09"     "fullname"    "stabbr"      "sector"
#> [16] "iclevel"     "control"     "hloffer"     "locale"      "merge_age_ic"
age_f16_allvars_allobs %>% var_label()
#> $unitid
#> [1] "Unique identification number of the institution"
#>
#> $agegroup
#> [1] "Age category"
#>
#> $levstudy
#> [1] "Level of student"
#>
#> $efage01
#> [1] "Full time men"
#>
#> $efage02
#> [1] "Full time women"
#>
#> $efage03
#> [1] "Part time men"
#>
#> $efage04
#> [1] "Part time women"
#>
#> $efage05
#> [1] "Full time total"
#>
#> $efage06
#> [1] "Part time total"
#>
#> $efage07
#> [1] "Total men"
#>
#> $efage08
#> [1] "Total women"
#>
#> $efage09
```

```

#> [1] "Grand total"
#>
#> $fullname
#> [1] "Institution (entity) name"
#>
#> $stabbr
#> [1] "State abbreviation"
#>
#> $sector
#> [1] "Sector of institution"
#>
#> $iclevel
#> [1] "Level of institution"
#>
#> $control
#> [1] "Control of institution"
#>
#> $hloffer
#> [1] "Highest level of offering"
#>
#> $locale
#> [1] "Degree of urbanization (Urban-centric locale)"
#>
#> $merge_age_ic
#> NULL

```

- Create two new data frames based on `age_f16_allvars_allobs`

NOTE: IN THIS QUESTION, WE GIVE YOU THE ANSWERS; ALL YOU HAVE TO DO IS RUN THE BELOW CODE CHUNK

```

#Create dataframe that keeps observations where age-group equals `1` (1. All age categories total)
agegroup1_obs <- age_f16_allvars_allobs %>%
  select(fullname,unitid,agegroup,levstudy,efage09,stabbr,locale) %>%
  filter(agegroup==1) %>%
  select(-agegroup)

glimpse(agegroup1_obs)
#> Rows: 7,019
#> Columns: 6
#> $ fullname <chr> "Amridge University", "Amridge University", "Amridge Universi~
#> $ unitid <dbl> 100690, 100690, 100690, 100724, 100724, 100724, 100751, 10075~
#> $ levstudy <dbl+lbl> 1, 2, 5, 1, 2, 5, 1, 2, 5, 1, 2, 1, 2, 5, 1, 2, 5, 1, 2, ~
#> $ efage09 <dbl> 597, 294, 303, 5318, 4727, 591, 37663, 32563, 5100, 1769, 176~
#> $ stabbr <chr> "AL", "AL", "AL", "AL", "AL", "AL", "AL", "AL", "AL", "AL", "AL", "~
#> $ locale <dbl+lbl> 12, 12, 12, 12, 12, 12, 13, 13, 13, 32, 32, 12, 12, 12, 1~

#Create dataframe keeps observations where "level of study" equals `1` (1. All Students total)
levstudy1_obs <- age_f16_allvars_allobs %>%
  select(fullname,unitid,agegroup,levstudy,efage09,stabbr,locale) %>%
  filter(levstudy==1) %>%
  select(-levstudy)

glimpse(levstudy1_obs)

```

```
#> Rows: 36,703
#> Columns: 6
#> $ fullname <chr> "Amridge University", "Amridge University", "Amridge Universi~
#> $ unitid <dbl> 100690, 100690, 100690, 100690, 100690, 100690, 100690, 10069~
#> $ agegroup <dbl+lbl> 1, 2, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 1, 2, ~
#> $ eface09 <dbl> 597, 57, 7, 16, 34, 540, 88, 97, 110, 158, 78, 9, 5318, 4464,~
#> $ stabbr <chr> "AL", "AL", "AL", "AL", "AL", "AL", "AL", "AL", "AL", "AL", "~
#> $ locale <dbl+lbl> 12, 12, 12, 12, 12, 12, 12, 12, 12, 12, 12, 12, 12, 1~
```

## Questions related to reshaping the dataset `agegroup1_obs` from long to wide

- Run whatever investigations seem helpful to you to get to know the data (e.g., list variable names, list variable variable labels, list variable values, tabulations). You may decide to comment out some of these investigations before you knit and submit the problem set so that your pdf doesn't get too long.

Sort and print a few obs

Run some frequencies

- Run the following code, which confirms that there is one row per each combination of `unitid-levstudy`

NOTE: IN THIS QUESTION, WE GIVE YOU THE ANSWERS; BUT TRY TO UNDERSTAND WHAT EACH PART OF THE CODE IS DOING

```
agegroup1_obs %>% group_by(unitid,levstudy) %>% # group by vars
  summarise(n_per_group=n()) %>% # create a measure of number of observations per group
  ungroup %>% # ungroup (otherwise frequency table [next step] created) separately for each group
  count(n_per_group) # frequency of number of observations per group
#> `summarise()` has grouped output by 'unitid'. You can override using the
#> `.groups` argument.
#> # A tibble: 1 x 2
#>   n_per_group      n
#>   <int> <int>
#> 1         1  7019
```

Using code from previous question as a guide, confirm that the object `agegroup1_obs` has more than one observation for each value of `unitid`

- Diagnose whether the data frame `agegroup1_obs` meets each of the three criteria for tidy data
  - YOUR ANSWERS HERE:
    - \* Each variable must have its own column:
    - \* Each observation must have its own row:
    - \* Each value must have its own cell:
  - What changes need to be made to `agegroup1_obs` to make it tidy?
    - YOUR ANSWER HERE:
  - With respect to “reshaping long to wide” to tidy a dataset, define the “names\_from” parameter.
    - YOUR ANSWER HERE:
  - What should the “names\_from” column be in the data frame `agegroup1_obs`?
    - YOUR ANSWER HERE:
  - With respect to “reshaping long to wide” to tidy a dataset, define the “values\_from” parameter.
    - YOUR ANSWER HERE:
  - What should the “values\_from” column be in the data frame `agegroup1_obs`?
    - YOUR ANSWER HERE:

Tidy the data frame `agegroup1_obs` and create a new object `agegroup1_obs_tidy`, then print a few

observations

Confirm that the new object `agegroup1_obs_tidy` contains one observation for each value of `unitid`

Create a new object `agegroup1_obs_tidy_v2` from the object `agegroup1_obs` by performing the following steps in one line of code with multiple pipes:

- Create a variable `level` that is a character version of the variable 'levstudy'
- Drop the original variable `levstudy`
- Tidy the dataset

Print a few observations of `agegroup1_obs_tidy_v2`; Why is this data frame preferable over `agegroup1_obs_tidy`?

- YOUR ANSWER HERE:

### Questions related to reshaping the dataset `levstudy1_obs` from long to wide

- Run whatever investigations seem helpful to you to get to know the data frame `levstudy1_obs` (e.g., list variable names, list variable labels, list variable values, tabulations). You may decide to comment out some of these investigations before you knit and submit the problem set so that your pdf doesn't get too long.

Sort and print a few obs

Run some frequencies

- Confirm that there is one row per each combination of `unitid`-`agegroup`

Using code from previous question as a guide, confirm that the object `levstudy1_obs` has more than one observation for each value of `unitid`

- Why is the data frame `levstudy1_obs` not tidy?
  - YOUR ANSWER HERE:
- What changes need to be made to `levstudy1_obs` to make it tidy?
  - YOUR ANSWER HERE:

Tidy the data frame `levstudy1_obs` and create a new object `levstudy1_obs_tidy` (it is up to you whether you want to create character version of the variable `agegroup` prior to tidying) then print a few observations

Confirm that the new object `levstudy1_obs_tidy` contains one observation for each value of `unitid`

## Part III: Questions about reshaping wide to long

Here, we load a table from NCES digest of education statistics that contains data about the total number of teachers in each state for particular years.

```
load(url("https://github.com/ksalazar3/HED696C_Rclass/raw/master/data/nces_digest/nces_digest_table_208_30"))
```

```
#convert character variables for teacher totals to integers
```

```
table208_30[2:6] <- data.frame(lapply(table208_30[2:6], as.integer))
```

```
table208_30
```

```
#> # A tibble: 51 x 6
```

```
#>   state   tot_fall_2000 tot_fall_2005 tot_fall_2009 tot_fall_2010 tot_fall_2011  
#>   <chr>         <int>         <int>         <int>         <int>         <int>  
#> 1 Alabam~      48194      57757      47492      49363      47722  
#> 2 Alaska~       7880       7912       8083       8170       8087  
#> 3 Arizon~     44438     51376     51947     50030     50800  
#> 4 Arkans~     31947     32997     37240     34272     33982
```

```
#> 5 Califo~      298021      309222      316298      260806      268688
#> 6 Colora~      41983       45841       49060       48542       48077
#> 7 Connec~      41044       39687       43592       42951       43804
#> 8 Delawa~       7469        7998        8639        8933        8587
#> 9 Distri~       4949        5481        5854        5925        6278
#> 10 Florid~     132030     158962     183827     175609     175006
#> # i 41 more rows
```

- Why is the data frame `table208_30` not tidy?  
– YOUR ANSWER HERE:
- What changes need to be made to `table208_30` to make it tidy?  
– YOUR ANSWER HERE:

Tidy the data frame `table208_30` and create a new object `table208_30_tidy`:

- hint: use the `cols = starts_with()` and `names_prefix=()` options for `pivot_longer()`
- after you tidy the data, print a few observations

Once finished, knit to (pdf) and upload both .Rmd and pdf files to class website under the week 6 tab  
*Remeber to use this naming convention "lastname\_firstname\_ps6"*