# Data Transformation: Part III

## Data Structure Conversion: Wide -> Long data

**Cultivate Learning Innovation Lab Workshop**
**August 31, 2020 | Monday | 2:00 - 3:30 p.m.**

**[황보 민] Min Hwangbo**

**PhD Candidate in Learning Sciences & Human Development**
**Graduate Certificate Recipient in Demographic Methods**

# Learning Agenda

- Why do we convert wide data to long data?
- Data transformation: `gather`
- Next steps: Challenge!

# [Intro] Inspiration

- Land
  - *"The University of Washington & Cultivate Learning acknowledges that it sits on Indigenous Land, which touches the shared waters of all tribes and bands within the Duwamish, Suquamish, Tulalip, and Muckleshoot Tribes."*
- People
  - Aimée Dechter | Affiliate Assistant Professor & Former Research Coordinator at Center for Studies in Demography & Ecology (CSDE)
  - [Chuck] Charles C. Lanfear | PhD Candidate & R Guru | 2020 Distinguished Teaching Award Recipient
  - [Jose] Jose Hernandez | Data Science Fellow & Research Staff @ eScience Institute
  - [Liz] Elizabeth Sanders | Associate Professor & Quantitative Researcher @ College of Education
  - [RStudio] Hadley Wickham & Garrett Grolemund | Authors of R for Data Science | RStudio Chief Scientist, Data Scientist & Statistician | Creators of RStudio
  - [이근열] Keun Yeol, Lee. | Professor in Busan National University & Qualitative/Dialect/Linguistics Researcher
  - [本橋智光] Motohashi, Tomomitsu. (2018). Maeshoritaizen data bunseki no tame no SQL/R/Python jissen technique (데이터 전처리 대전. 2019).
  - [Nicolas] Nicolas Pröllochs | Tenure-track Professor of Data Science in University of Giessen & Social Network Analysis / Text Mining Expert

# Learning Agenda

- Why do we convert wide data to long data?
- Data transformation: `gather`
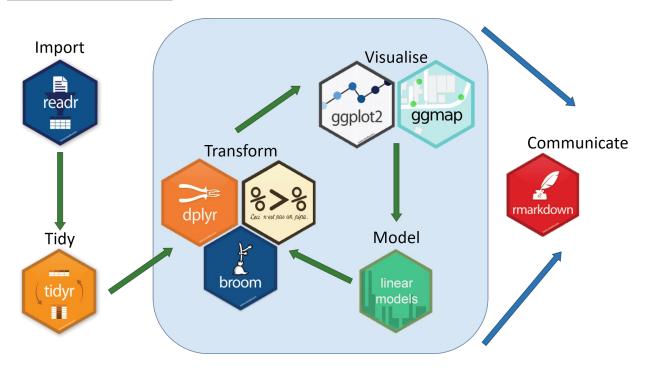- Next steps: Challenge!

# [Intro] Data Transformation?

- *"Process of **converting** data from one format or structure into another format or structure…. **fundamental** aspect of most **data integration** and **data management.**"*[1]
  - Data framing
  - Data manipulation
  - Data wrangling
  - Data management
- Data Transformation Steps
  - Data discovery: Where's my data?
  - Data mapping: Explore your data – What's in it? Using data transformation technique
  - Code generation: Cheatsheet & Stackflow & Pre-Built codes + Ctrl + C / Ctrl + V
  - Code execution: Ctrl + Enter
  - Data review: head(data), tail(data), View(data), missing data, quality of your data, etc…

[1] Source: Wikipedia

# [Intro] Packages: dplyr & tidyr

- **dplyr**: "The grammar of data manipulation."
  - **mutate()**: adds new variables that are functions of existing variables
  - **select()**: picks variables based on their names.
  - **filter()**: picks cases based on their values.
  - **summarise()**: reduces multiple values down to a single summary.
  - **arrange()**: changes the ordering of the rows.
- **tidyr**: "The goal of tidyr is to help you create tidy data."
  - "Every **column** is **variable**."
  - "Every **row** is an **observation**."
  - "Every **cell** is a single **value**."
- Cheatsheet!: Data Transformation Cheatsheet

# [Intro] tidyverse Workflow

# [Why] Wide Data -> Long data?

- *"**Wide data** has a column for each variable. Whereas **long format data** has a column for possible variable types & a column for the values of those variables."* [1]
  - Wide format: Required for Multivariate analysis of variance (MANOVA) or repeated measures in SPSS.
  - Long format:
    - Mixed models (containing fixed & random effects) or most of the survival analysis
    - Also required for most data visualization softwares for survey data analysis (i.e. Qualtrics -> Tableau) - Common @ Cultivate Learning

[1] The analysis factor. (n.d.). *The Wide and Long Data Format for Repeated Measures Data.* https://www.theanalysisfactor.com/wide-and-long-data/

# [Why] Wide Data -> Long data?[1]



[1] Salesforce. (n.d.). *Get your data Tableau-teady.* https://www.tableau.com/learn/get-started/data-structure

# Questions?

# [Prep] Preparation

- Will not cover in this session (*See Session 2: Data Transformation*)
    - Package preparation
        - Installing packages
        - Loading packages
    - Data preparation
        - Loading data set
        - Framing data set
- **Will cover in this session**
    - Checking `Null` values
    - `gather` function to gather multiple columns into one column.
- Will not cover in this session (*See Session 2: Data Transformation*)
    - Saving it as a csv file

# 10 min Break
## Coffee, Tea & Snack Time

# Learning Agenda

- Why do we convert wide data to long data?
- Data transformation: `gather`
- Next steps: Challenge!

# [Prep] Step 1 & 2: Loading Packages & Data Set

```r
18   # Step 1: Loading packages - Data reshape
19   ```{r}
20   # Install packages first using install.packages("pacakgename") on your console!
21   ## Reference: http://www.cookbook-r.com/Manipulating_data/Converting_data_between_wide_and_long_format/
22   ### Inspiration: Dr. Liz Sanders's HLM class
23   #### Inspiration II: Chuck Lanfear Intro to R: https://clanfear.github.io/CSSS508/
24
25   library("dplyr")   # Data reshaping package
26   library("tidyr")   # Data transformation package
27   library("readr")   # CSV loading package
28   ```
```

```r
31   # Step 2: Loading data
32   ```{r, msg = F, warning = F}
33   data.wide <- read_csv("EvalSTARS_CPD_051320.csv") # Loading data from an excel file
34
35   # Quality check descriptive
36   head(data.wide, 5) # First five records
37   tail(data.wide, 5) # Last five records
38   ls(data.wide)       # Vars names
39   ## Data summary: 112 records / 22 vars - check for null values
40   ```
```
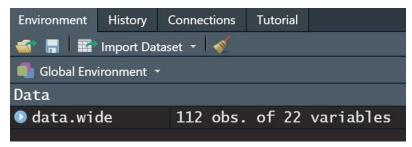
# [Prep] Data exploration: There are NA records!



| | Finished | Recorded Date | Which training did you participate in? | Please mark only one response per line. – Content provided matched the training description. | Please mark only one response per line. – Content provided matched the core competency level indicated in the training description. | Please mark only one response per line. – Examples and illustrations used in the training were relevant to practice. | Please mark only one response per line. – Handouts were useful. | Please mark only one response per line. – Trainer was knowledgeable about the topic. |
|---|---|---|---|---|---|---|---|---|
| 1 | FALSE | 3/24/20 9:49 | NA | NA | NA | NA | NA | NA |
| 2 | FALSE | 3/26/20 9:52 | NA | NA | NA | NA | NA | NA |
| 3 | FALSE | 3/26/20 12:00 | NA | NA | NA | NA | NA | NA |
| 4 | FALSE | 3/26/20 13:18 | NA | NA | NA | NA | NA | NA |
| 5 | FALSE | 3/26/20 14:17 | NA | NA | NA | NA | NA | NA |

# [Prep] Options for NA records

- Option 1: Probably the best option :p
  - Do nothing, who cares...
- Option 2: Delete all NA records in excel... which will take couple hours to figure it out if your data set a ton of records....
- Option 3:
  - Identify number of NA records in R first then
  - `**Filter**` your original data set with no NA records.

# [Prep] Step 3: Checking Null values

sum()                      # Summing number of counts within ()

is.na()                    # Will generate a list of whether (yourdataframe) has a record of NA

                           or null values

- We're going to figure out how many records are available based whether participant(s)'s indication of which training that they participated in?

sum(is.na(data.wide$`Which training did you participate in?`))

# [Prep] Step 3: Checking Null values

sum(is.na(data.wide$`Which training did you participate in?`))

```
# Step 3: Checking `Null` values
```{r}                                                    ⚙ ⏬ ▶
# summary(data.wide)
## Data summary: 112 records / 22 vars - check for null values
sum(is.na(data.wide$`Which training did you participate in?`))
# 26 records = also happens to be something called "False" from the data set.
```
```

```
                                                         ⏏ ⌃ ✕
```

```
 [1] 26
```

# [Prep] Step 4: Filter "Finished" responses

filter (ColumnName == "Value")          # Filter it based on value of the column name

- We're going to only filter responses that are considered "finished" and create the revised data set into data.wide2

data.wide2 <- data.wide %>%
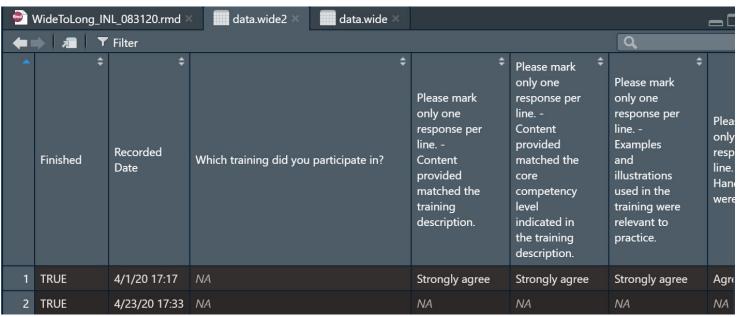
 filter(Finished == "TRUE")

# [Prep] Step 4: Filter "Finished" responses

```r
# Step 4: Data reshape (`Filter` responses that are considered `finished`)
```{r}
data.wide2 <- data.wide %>%
  filter(Finished == "TRUE")
```

| Environment | History | Connections | Tutorial |
|---|---|---|---|

Import Dataset

Global Environment

**Data**

| | |
|---|---|
| data.wide | 112 obs. of 22 variables |
| data.wide2 | 84 obs. of 22 variables |

# [Prep] Step 4: Filter "Finished" responses

- 112 records – 26 records = 84 records?
- Something is going on with 2 responses...

| | Finished | Recorded Date | Which training did you participate in? | Please mark only one response per line. - Content provided matched the training description. | Please mark only one response per line. - Content provided matched the core competency level indicated in the training description. | Please mark only one response per line. - Examples and illustrations used in the training were relevant to practice. | Plea only resp line. Han were |
|---|---|---|---|---|---|---|---|
| 1 | TRUE | 4/1/20 17:17 | NA | Strongly agree | Strongly agree | Strongly agree | Agre |
| 2 | TRUE | 4/23/20 17:33 | NA | NA | NA | NA | NA |

# [Prep] Step 4: Filter "Finished" responses

- 112 records - 26 records = 84 records?
- Something is going on with 2 responses... (don't forget to do your quality check!)

# [Transformation] Step 5: `gather` Wide -> Long

- Using `gather` function to gather questions from multiple columns into one column.

gather(dataset, column1, column2, "start column from data set: end column from data set", factor_key = True)

- Column 1: usually **higher hierarchy** of a column that you'd like to gather it as (i.e. question)
- Column 2: **Value** of the Column 1

# [Transformation] Step 5: `gather` Wide -> Long

data.long <- **gather**(data.wide2, question, response, "Please mark only one response per line. – Content provided matched the training description.":"For future training, what topic(s) are you looking for (Select your top three choices) – Child & Youth Development Competency Areas:", factor_key = T)

```r
# Step 5: Data reshape (Wide -> Long)
* Using `gather` function to consolidate questions into one column.
```{r}
data.long <- gather(data.wide2, question, response, "Please mark only one
response per line. - Content provided matched the training description.":"For
future training, what topic(s) are you looking for (Select your top three
choices) - Child & Youth Development Competency Areas:", factor_key = T)
```

# [Transformation] Step 5: `gather` Wide -> Long

# Summary

- Understanding how data can be reshaped can help you in a long run to feed your data to any software.
- It takes more time to plan these transformation structure than actual execution.
- Try to recycle what you have tried out last time :)

# Next Steps

- Challenge! Institute scenario!

# Institute Scenario

- Cultivate Learning institute manager has collected a data set of "session evaluation" from the last 2020 virtual institute.
- The current data set seems it has each session on each tab, and it's not well organized as we hoped for.
- We communicated to the client that this would take a while to clean up before we attempt to visualize this data on Tableau – the platform where the client wants to visualize the data set.
- **You have 72 hours to think thoroughly and document your plan on how to organize this data set to feed into Tableau.**

# Thank you!

감사합니다!