# Lecture6\_note\_data\_wrangling

### Shukyee Chan

#### 2023-10-12

### Outline

## 2 MEX

## 3 MEX

## # ... with 6,786 more rows

- Reshape (long <-> wide) with pivot\_longer and pivot\_wider
- Stack tables by row or by column with bind\_rows and bind\_cols (or, alternatively, cbind and rbind)
- Merge two tables with inner\_join, full\_join, left\_join, right\_join, semi\_join, and anti\_join
- Save your outputs

```
library(tidyverse)
d <- read_csv("_DataPublic_/vdem/1984_2022/vdem_1984_2022_external.csv")
d |> print(n = 3)
## # A tibble: 6,789 x 211
##
     country_name countr~1 count~2 year historic~3 project histo~4 histn~5 codin~6
##
     <chr>
                  <chr>
                             <dbl> <dbl> <date>
                                                      <dbl>
                                                               <dbl> <chr>
## 1 Mexico
                  MEX
                                 3 1984 1984-12-31
                                                          0
                                                                   1 United~
                                                                                1789
## 2 Mexico
                  MEX
                                    1985 1985-12-31
                                                           0
                                                                   1 United~
                                                                                1789
## 3 Mexico
                  MEX
                                 3 1986 1986-12-31
                                                          0
                                                                   1 United~
                                                                                1789
## # ... with 6,786 more rows, 202 more variables: codingend <dbl>,
       codingstart_contemp <dbl>, codingend_contemp <dbl>, codingstart_hist <dbl>,
## #
       codingend_hist <dbl>, gapstart1 <dbl>, gapstart2 <dbl>, gapstart3 <dbl>,
       gapend1 <dbl>, gapend2 <dbl>, gapend3 <dbl>, gap_index <dbl>,
## #
       COWcode <dbl>, e_v2x_api_3C <dbl>, e_v2x_api_4C <dbl>, e_v2x_api_5C <dbl>,
       e_v2x_civlib_3C <dbl>, e_v2x_civlib_4C <dbl>, e_v2x_civlib_5C <dbl>,
## #
       e_v2x_clphy_3C <dbl>, e_v2x_clphy_4C <dbl>, e_v2x_clphy_5C <dbl>, ...
d_gdp <- d |>
  select(country_text_id, year, e_gdp, e_gdppc) |>
  rename("gdp" = "e_gdp", "gdppc" = "e_gdppc")
d_gdp |> print(n = 3)
## # A tibble: 6,789 x 4
     country_text_id year
                              gdp gdppc
##
     <chr>
                     <dbl> <dbl> <dbl>
## 1 MEX
                      1984 93563. 11.7
```

1985 94259. 11.5

1986 92750. 11.1

# 1. Reshape a Table

d\_gdp\_long <- d\_gdp |>

Wide to Long: pivot\_longer

```
pivot_longer(cols = c("gdp", "gdppc"),
              names_to = "variable", values_to = "value")
d_gdp_long |> print(n = 4)
## # A tibble: 13,578 x 4
##
    country_text_id year variable
                                    value
##
    <chr>
                   <dbl> <chr>
                                  <dbl>
## 1 MEX
                                  93563.
                   1984 gdp
## 2 MEX
                   1984 gdppc
                                   11.7
## 3 MEX
                    1985 gdp
                                  94259.
## 4 MEX
                    1985 gdppc
                                   11.5
## # ... with 13,574 more rows
Long to Wide: pivot_wider
# Reverse the above pivot_long operation
d_gdp_wide_1 <- d_gdp_long |>
 pivot_wider(names_from = "variable", values_from = "value")
d_gdp_wide_1 |> print(n = 4)
## # A tibble: 6,789 x 4
    country_text_id year gdp gdppc
##
    1984 93563. 11.7
## 1 MEX
## 2 MEX
                   1985 94259. 11.5
                    1986 92750. 11.1
## 3 MEX
## 4 MEX
                    1987 93220. 10.9
## # ... with 6,785 more rows
# Make year the column variable
d_gdp_wide_2 <- d_gdp_long |>
 pivot_wider(names_from = "year", values_from = "value")
d_gdp_wide_2 |> print(n = 2)
## # A tibble: 362 x 41
    count~1 varia~2 '1984' '1985' '1986' '1987' '1988' '1989' '1990' '1991' '1992'
                    <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> 
## <chr> <chr>
## 1 MEX
            gdp
                    9.36e4 9.43e4 9.28e4 9.32e4 9.47e4 9.81e4 1.03e5 1.07e5 1.12e5
            gdppc 1.17e1 1.15e1 1.11e1 1.09e1 1.08e1 1.10e1 1.14e1 1.16e1 1.19e1
## # ... with 360 more rows, 30 more variables: '1993' <dbl>, '1994' <dbl>,
## # '1995' <dbl>, '1996' <dbl>, '1997' <dbl>, '1998' <dbl>, '1999' <dbl>,
```

```
'2000' <dbl>, '2001' <dbl>, '2002' <dbl>, '2003' <dbl>, '2004' <dbl>,
       '2005' <dbl>, '2006' <dbl>, '2007' <dbl>, '2008' <dbl>, '2009' <dbl>,
## #
       '2010' <dbl>, '2011' <dbl>, '2012' <dbl>, '2013' <dbl>, '2014' <dbl>,
## #
       '2015' <dbl>, '2016' <dbl>, '2017' <dbl>, '2018' <dbl>, '2019' <dbl>,
## #
       '2020' <dbl>, '2021' <dbl>, '2022' <dbl>, and abbreviated variable ...
## #
# Make country text id the column variable
d_gdp_wide_3 <- d_gdp_long |>
 pivot_wider(names_from = "country_text_id", values_from = "value")
d_gdp_wide_3 |> print(n = 2)
## # A tibble: 78 x 183
                                             CHE
     year variable
                       MEX
                               SUR
                                      SWE
                                                    GHA
                                                           ZAF
                                                                  JPN
                                                                         MMR
                                                                                RUS
##
     <dbl> <chr>
                     <dbl> <
## 1 1984 gdp
                                   2.35e4 2.31e4 3.02e3 3.15e4 2.87e5 4.18e3 3.49e5
                    93563. 286.
                            7.43 2.66e1 3.32e1 2.20e0 9.03e0 2.26e1 1.10e0 1.65e1
## 2 1984 gdppc
                     11.7
## # ... with 76 more rows, and 172 more variables: ALB <dbl>, EGY <dbl>,
      YEM <dbl>, COL <dbl>, POL <dbl>, BRA <dbl>, USA <dbl>, PRT <dbl>,
      SLV <dbl>, YMD <dbl>, BGD <dbl>, BOL <dbl>, HTI <dbl>, HND <dbl>,
      MLI <dbl>, PAK <dbl>, PER <dbl>, SEN <dbl>, SSD <dbl>, SDN <dbl>,
## #
      VNM <dbl>, AFG <dbl>, ARG <dbl>, ETH <dbl>, IND <dbl>, KEN <dbl>,
## #
      PRK <dbl>, KOR <dbl>, XKX <dbl>, LBN <dbl>, NGA <dbl>, PHL <dbl>,
## #
      TZA <dbl>, TWN <dbl>, THA <dbl>, UGA <dbl>, VEN <dbl>, BEN <dbl>, ...
```

## 2. Stack Tables

```
# New data (for stack data vertically)
d_gdp_1945 <-
 read csv(" DataPublic /vdem/1945 1983/vdem 1945 1983 external.csv") |>
  select(country_text_id, year, e_gdp, e_gdppc) |>
 rename("gdp" = "e_gdp", "gdppc" = "e_gdppc")
d_gdp_1906 <-
 read_csv("_DataPublic_/vdem/1906_1944/vdem_1906_1944_external.csv") |>
  select(country_text_id, year, e_gdp, e_gdppc) |>
 rename("gdp" = "e_gdp", "gdppc" = "e_gdppc")
d_gdp_1945 \Rightarrow print(n = 2)
## # A tibble: 6,082 x 4
    country_text_id year gdp gdppc
##
                     <dbl> <dbl> <dbl>
    <chr>
## 1 MEX
                     1945 7827. 3.08
## 2 MEX
                     1946 8331. 3.17
## # ... with 6,080 more rows
# New data (for stack data horizontally)
d edu <- d |>
 select(e_peaveduc, e_peedgini) |>
```

```
rename("edu_15" = "e_peaveduc", "edu_gini" = "e_peedgini")
d_fh <- d |>
  select(starts_with("e_fh")) |>
  rename("fh_CivilLiberty" = "e_fh_cl", "fh_PoliticalRight" = "e_fh_pr",
         "fh_RuleOfLaw" = "e_fh_rol", "fh_Status" = "e_fh_status")
d fh > print(n = 2)
## # A tibble: 6,789 x 4
## fh CivilLiberty fh PoliticalRight fh RuleOfLaw fh Status
##
              <dbl>
                               <dbl> <dbl> <dbl> <dbl>
## 1
                  4
                                   3
                                               NA
## 2
                   4
                                     4
                                                            2
                                                NA
## # ... with 6,787 more rows
bind_rows
d_gdp_1945_2022 <- bind_rows(d_gdp, d_gdp_1945)</pre>
d_gdp_1945_2022 \Rightarrow print(n = 3)
## # A tibble: 12,871 x 4
##
   country_text_id year gdp gdppc
##
     <chr>
                    <dbl> <dbl> <dbl>
                     1984 93563. 11.7
## 1 MEX
## 2 MEX
                    1985 94259. 11.5
                     1986 92750. 11.1
## 3 MEX
## # ... with 12,868 more rows
unique(d_gdp_1945_2022$year) |> sort()
## [1] 1945 1946 1947 1948 1949 1950 1951 1952 1953 1954 1955 1956 1957 1958 1959
## [16] 1960 1961 1962 1963 1964 1965 1966 1967 1968 1969 1970 1971 1972 1973 1974
## [31] 1975 1976 1977 1978 1979 1980 1981 1982 1983 1984 1985 1986 1987 1988 1989
## [46] 1990 1991 1992 1993 1994 1995 1996 1997 1998 1999 2000 2001 2002 2003 2004
## [61] 2005 2006 2007 2008 2009 2010 2011 2012 2013 2014 2015 2016 2017 2018 2019
## [76] 2020 2021 2022
d_gdp_1945_2022_ue_rows <- bind_rows(</pre>
d_gdp |> select(-gdppc),
d_gdp_1945 |> select(-gdp)
d_gdp_1906_2022 <- bind_rows(d_gdp, d_gdp_1945, d_gdp_1906)</pre>
d_gdp_1906_2022 |> print(n = 3)
## # A tibble: 18,559 x 4
## country_text_id year gdp gdppc
##
                    <dbl> <dbl> <dbl>
     <chr>
```

```
## 1 MEX
                      1984 93563. 11.7
## 2 MEX
                      1985 94259. 11.5
## 3 MEX
                      1986 92750.
## # ... with 18,556 more rows
unique(d_gdp_1906_2022$year) |> sort()
     [1] 1906 1907 1908 1909 1910 1911 1912 1913 1914 1915 1916 1917 1918 1919 1920
##
##
    [16] 1921 1922 1923 1924 1925 1926 1927 1928 1929 1930 1931 1932 1933 1934 1935
   [31] 1936 1937 1938 1939 1940 1941 1942 1943 1944 1945 1946 1947 1948 1949 1950
   [46] 1951 1952 1953 1954 1955 1956 1957 1958 1959 1960 1961 1962 1963 1964 1965
   [61] 1966 1967 1968 1969 1970 1971 1972 1973 1974 1975 1976 1977 1978 1979 1980
   [76] 1981 1982 1983 1984 1985 1986 1987 1988 1989 1990 1991 1992 1993 1994 1995
## [91] 1996 1997 1998 1999 2000 2001 2002 2003 2004 2005 2006 2007 2008 2009 2010
## [106] 2011 2012 2013 2014 2015 2016 2017 2018 2019 2020 2021 2022
```

### bind\_cols

```
d_gdp_edu_fh <- bind_cols(d_gdp, d_edu, d_fh)</pre>
d_gdp_edu_fh |> print(n = 3)
## # A tibble: 6,789 x 10
                            gdp gdppc edu_15 edu_g~2 fh_Ci~3 fh_Po~4 fh_Ru~5 fh_St~6
##
     country_te~1 year
                                       <dbl>
                                                <dbl>
                                                        <dbl>
                                                                 <dbl>
##
     <chr>>
                  <dbl>
                         <dbl> <dbl>
                                                                         <dbl>
## 1 MEX
                   1984 93563. 11.7
                                        6.08
                                                 32.7
                                                            4
                                                                     3
                                                                            NA
                                                                                     2
## 2 MEX
                   1985 94259. 11.5
                                                 32.4
                                                                            NA
                                                                                      2
                                        6.22
                                                                                      2
## 3 MEX
                    1986 92750. 11.1
                                        6.36
                                                 31.9
                                                            4
                                                                     4
                                                                            NA
## # ... with 6,786 more rows, and abbreviated variable names 1: country_text_id,
       2: edu_gini, 3: fh_CivilLiberty, 4: fh_PoliticalRight, 5: fh_RuleOfLaw,
## #
       6: fh_Status
names(d_gdp_edu_fh)
                             "year"
    [1] "country_text_id"
                                                  "gdp"
##
    [4] "gdppc"
                             "edu 15"
                                                  "edu gini"
  [7] "fh_CivilLiberty"
                             "fh_PoliticalRight" "fh_RuleOfLaw"
## [10] "fh Status"
```

## These are error-prone operations

- Do bind\_rows and bind\_cols ONLY WHEN you know for sure that there will not be a mismatch!
- If you have any slightest doubt, don't use them.

#### 3. Join Tables

• left\_join: Merge and only keep observations whose identifiers (matching keys) appear in the left-hand-side table.

- right\_join: Merge and only keep observations whose identifiers (matching keys) appear in the right-hand-side table.
- inner\_join: Merge and only keep observations whose identifiers (matching keys) appear in both tables.
- full\_join: Merge and keep observations whose identifiers (matching keys) appear either table.
- anti\_join: Filter out observations whose identifiers (matching keys) appear in the right-hand-side table
- semi\_join: Filter out observations whose identifiers (matching keys) do not appear in the right-hand-side table

#### Task 1: The Case

Join two datasets from the V-Dem data using the above different join\_ functions

• *GDP* data from **2000-2022** 

d\_lj <- d\_gdp\_2000\_2022 |>

d\_lj |> print(n = 2)

• GDP per capita data from 1984 to 2010

```
# Setup
d_gdp_2000_2022 <- d |> filter(year %in% 2000:2022) |>
  select(country_text_id, year, e_gdp) |> rename("gdp" = "e_gdp")
d_gdppc_1984_2010 <- d |> filter(year %in% 1984:2010) |>
  select(country_text_id, year, e_gdppc) |> rename("gdppc" = "e_gdppc")
d_gdp_2000_2022 |> print(n = 2)
## # A tibble: 4,099 x 3
     country_text_id year
##
                               gdp
##
                     <dbl>
                             <dbl>
## 1 MEX
                      2000 145206.
## 2 MEX
                      2001 146993.
## # ... with 4,097 more rows
d_gdppc_1984_2010 |> print(n = 2)
## # A tibble: 4,641 x 3
     country_text_id year gdppc
##
                     <dbl> <dbl>
     <chr>
## 1 MEX
                      1984 11.7
## 2 MEX
                      1985 11.5
## # ... with 4,639 more rows
left_join
```

left\_join(d\_gdppc\_1984\_2010, by = c("country\_text\_id", "year"))

```
## # A tibble: 4,099 x 4
## country_text_id year gdp gdppc
   <chr>
             <dbl> <dbl> <dbl> <
## 1 MEX
                   2000 145206. 13.7
                    2001 146993. 13.6
## 2 MEX
## # ... with 4,097 more rows
unique(d_lj$year) |> sort()
## [1] 2000 2001 2002 2003 2004 2005 2006 2007 2008 2009 2010 2011 2012 2013 2014
## [16] 2015 2016 2017 2018 2019 2020 2021 2022
right_join
d_rj <- d_gdp_2000_2022 |>
 right_join(d_gdppc_1984_2010, by = c("country_text_id", "year"))
d_rj > print(n = 2)
## # A tibble: 4,641 x 4
## country_text_id year gdp gdppc
   <chr> <dbl> <dbl> <dbl>
                    2000 145206. 13.7
## 1 MEX
## 2 MEX
                    2001 146993. 13.6
## # ... with 4,639 more rows
unique(d_rj$year) |> sort()
## [1] 1984 1985 1986 1987 1988 1989 1990 1991 1992 1993 1994 1995 1996 1997 1998
## [16] 1999 2000 2001 2002 2003 2004 2005 2006 2007 2008 2009 2010
inner_join
d_ij <- d_gdp_2000_2022 |>
 inner_join(d_gdppc_1984_2010, by = c("country_text_id", "year"))
d_ij |> print(n = 2)
## # A tibble: 1,951 x 4
   country_text_id year
                           gdp gdppc
    ## 1 MEX
                    2000 145206. 13.7
## 2 MEX
                    2001 146993. 13.6
## # ... with 1,949 more rows
```

```
unique(d_ij$year) |> sort()
## [1] 2000 2001 2002 2003 2004 2005 2006 2007 2008 2009 2010
full join
d_fj <- d_gdp_2000_2022 |>
 full_join(d_gdppc_1984_2010, by = c("country_text_id", "year"))
d_fj > print(n = 2)
## # A tibble: 6,789 x 4
## country_text_id year gdp gdppc
                   <dbl> <dbl> <dbl>
## <chr>
                   2000 145206. 13.7
## 1 MEX
## 2 MEX
                   2001 146993. 13.6
## # ... with 6,787 more rows
unique(d_fj$year) |> sort()
## [1] 1984 1985 1986 1987 1988 1989 1990 1991 1992 1993 1994 1995 1996 1997 1998
## [16] 1999 2000 2001 2002 2003 2004 2005 2006 2007 2008 2009 2010 2011 2012 2013
## [31] 2014 2015 2016 2017 2018 2019 2020 2021 2022
semi_join
d_sj <- d_gdp_2000_2022 |>
  semi_join(d_gdppc_1984_2010, by = c("country_text_id", "year"))
d_{sj} > print(n = 2)
## # A tibble: 1,951 x 3
   country_text_id year
                             gdp
##
   ## 1 MEX
                   2000 145206.
## 2 MEX
                    2001 146993.
## # ... with 1,949 more rows
unique(d_sj$year) |> sort()
## [1] 2000 2001 2002 2003 2004 2005 2006 2007 2008 2009 2010
anti_join
```

## [1] 2011 2012 2013 2014 2015 2016 2017 2018 2019 2020 2021 2022

## 3. Join by Identifiers with Different Variable Names

two options: (1) Rename it beforehand, (2) specify the by = argument differently.

### 4. Many-to-One Join: Repeat!

Calculate each country's average 1984-2010~GDP~per~capita and merge it with our annual GDP data from 2000~to~2022.

```
d_gdppc_1984_2010_avg <- d_gdppc_1984_2010 |> group_by(country_text_id) |>
    summarise(gdppc_1984to2010 = mean(gdppc, na.rm = TRUE))
d_gdppc_1984_2010_avg |> print(n = 2)
## # A tibble: 180 x 2
```

```
d_lj_ManyToOne <- d_gdp_2000_2022 |>
 left_join(d_gdppc_1984_2010_avg, by = "country_text_id")
d_lj_ManyToOne |> print(n = 2)
## # A tibble: 4,099 x 4
    country_text_id year gdp gdppc_1984to2010
##
                   <dbl> <dbl>
##
    <chr>
                                           <dbl>
## 1 MEX
                   2000 145206.
                                            12.8
## 2 MEX
                    2001 146993.
                                           12.8
## # ... with 4,097 more rows
```

# 5. Save Outputs

- .csv "comma-separated values," readable by Excel or a text editor
- .rds "R data serialization," readable by R only

```
# Save to a .csv file
write_csv(d_gdp_1945_2022, "Lecture_06/data/gdp_1945_2002.csv")
# Save to a .rds file
saveRDS(d_gdp_1945_2022, "Lecture_06/data/gdp_1945_2002.rds")
```

# Saving Your Outputs after Data Wrangling

```
# Read a .csv file
d_read_1 <- read_csv("Lecture_06/data/gdp_1945_2002.csv")

# Read a .rds file
d_read_2 <- readRDS("Lecture_06/data/gdp_1945_2002.rds")</pre>
```