# WinteR Statistical Workshop
## Merge

Zack Roman[1,2]

[1]Center for Research Methods and Data Analysis
[2]Department of Psychology

2017

# Outline

# Goals of This Session

**Conceptual:**

- Types of merges
- Merging vocabulary
- When to use merges

**Skill Building:**

- Practicing merging variants
- Different implementations of merging in R
- Dangers associated with improper merging and how to avoid them

# Small Example

```
1   authors
```

```
1       surname nationality deceased
2   1     Tukey          US      yes
3   2  Venables    Australia       no
4   3   Tierney          US       no
5   4    Ripley          UK       no
6   5    McNeil    Australia       no
```

```
1   books
```

```
1          name                  title      other_author
2   1     Tukey Exploratory Data Analysis          <NA>
3   2  Venables Modern Applied Statistics        Ripley
4   3   Tierney                 LISP-STAT          <NA>
5   4    Ripley        Spatial Statistics          <NA>
6   5    Ripley     Stochastic Simulation          <NA>
7   6    McNeil Interactive Data Analysis          <NA>
8   7    R Core    An Introduction to R Venables & Smith
```

## Small Example ...

```
1  merge ( x = authors , y = books , by . x =
       " surname " , by . y = " name " )
```

```
1      surname nationality deceased                          title
2  1    McNeil    Australia       no Interactive Data Analysis
3  2    Ripley           UK       no          Spatial Statistics
4  3    Ripley           UK       no       Stochastic Simulation
5  4   Tierney           US       no                   LISP-STAT
6  5     Tukey           US      yes Exploratory Data Analysis
7  6  Venables    Australia       no Modern Applied Statistics
8      other_author
9  1          <NA>
10 2          <NA>
11 3          <NA>
12 4          <NA>
13 5          <NA>
14 6        Ripley
```

## Merge Arguments

```
1   merge ( x , y , by.x , by.y , by , incomparables ,
        sort , all.x , all.y , all )
```

1. x   Specifies the left data set

2. y   Specifies the right data set

3. by.x, by.y, by   specifies the key as a character string. by   is common to both   x   and   y .

4. incomparables   provides values in the key to not be used for matching, such as NA, blank space, or NaN (not a number).

5. sort   Logical (TRUE or FALSE), sorts the output

6. all.x, all.y, all   Logical, will help us determine the behavior of the merge. We will talk more about this as we go

# Binding is not a merging

- The functions `rbind()` and `cbind()` can be used to "stack" matrices on top of each other (rows bound together), or place them side by side (columns bound together)
- Binding puts data sets together, but if the rows (or columns) are not in exactly the same order, it will corrupt the result. Binding two data sets is not merging
- Merging takes into account a "Key" variable (typically an ID # or Name), so that the correct rows are aligned with each other.

# SQL Terminology

- SQL = "Structured Query Language". Very widely used general purpose data-base framework.
- R merge developed in isolation, used different terminology.
- Next we show that the SQL terms "left join", "inner join" and so forth can be achieved by properly setting the value of the merges `all` parameter ( `all` , `all.x` , and `all.y` )

# Left Join

The "Left Join" is used when the goal data set should **only** have rows
that are present in X. The key variable is used to scan Y for matches,
which are then merged with the X rows.

```
1    dat_legs
```

```
1    animal legs
2  1    dog      4
3  2   cats      4
4  3  human      2
5  4  snake      0
6  5   tree      0
```

```
1    dat_fur
```

```
1    animal    fur
2  1    dog     yes
3  2   cats  Mostly
4  3  human      No
5  4   bird      No
```

# Left Join ...

```
1   merge ( x = dat_legs , y = dat_fur , by =
        " animal " , all.x = TRUE )
```

```
1     animal legs    fur
2   1   cats    4 Mostly
3   2    dog    4    yes
4   3  human    2     No
5   4  snake    0   <NA >
6   5   tree    0   <NA >
```

Setting "all.x" to **TRUE** produces an "Inner Join". The output data will only contain rows that have matching key values on **both** input data sets.

# Left Join

importfigs/left_join_ven.pdf

## Left Join Switched

Let's do a Left Join again, but switch the data sets.

```
1   dat_legs
```

```
1     animal legs
2   1     dog    4
3   2    cats    4
4   3   human    2
5   4   snake    0
6   5    tree    0
```

```
1   dat_fur
```

```
1     animal    fur
2   1     dog    yes
3   2    cats Mostly
4   3   human     No
5   4    bird     No
```

## Left Join Switched ...

```
1   merge ( x = dat_fur , y = dat_legs , by =
        " animal " , all . x = TRUE )
```

```
1     animal     fur legs
2   1    bird      No   NA
3   2    cats  Mostly    4
4   3     dog     yes    4
5   4   human      No     2
```

# Situations calling for Left Join

- You want to investigate the relationship between fur and legs in animals
- You have a data set of the animals you are interested in and their fur status
- You obtain a list of **all** animals legs count
  - Key = Animal Name
  - Output data is the length of the fur data set
- You want to investigate the effect of tuition on retention rate in Florida
- You have Floridian school tuition rates data set
- You obtain a nationwide data set of retention rates
  - Key = School Name
  - Output data is the length of the tuition rates data set

## Inner Join

The "Inner join" is used when the goal data set should only have rows that have keys in both the X and Y data.

```
1  dat_legs
```

```
1     animal legs
2  1     dog    4
3  2    cats    4
4  3   human    2
5  4   snake    0
6  5    tree    0
```
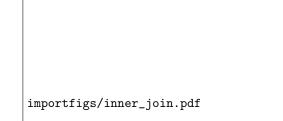
```
1  dat_fur
```

```
1     animal    fur
2  1     dog    yes
3  2    cats Mostly
4  3   human     No
5  4    bird     No
```

## Inner Join ...

```
1    merge ( x = dat_legs , y = dat_fur , by =
         " animal " , all = FALSE )
```

```
1    animal legs    fur
2 1    cats    4 Mostly
3 2     dog    4    yes
4 3   human    2     No
```

Setting "all" to **FALSE** produces an "Inner Join". The output data will
only contain rows that have matching key values on **both** input data sets.

# Inner Join

importfigs/inner_join.pdf

# Qualities of Inner Joins

- Pro, result data set will be more complete than other merges.
- Con, result data set looses more information than other merges.

# Full Join

Full Join keeps all data rows, filling in unmatched rows with missing values.

```
1  dat_legs
```

```
1    animal legs
2  1    dog    4
3  2   cats    4
4  3  human    2
5  4  snake    0
6  5   tree    0
```

```
1  dat_fur
```

```
1    animal    fur
2  1    dog    yes
3  2   cats Mostly
4  3  human     No
5  4   bird     No
```

# Full Join

```
1  merge ( x = dat_legs , y = dat_fur , by =
       " animal " , all = TRUE )
```

```
1     animal legs     fur
2  1    bird   NA      No
3  2    cats    4 Mostly
4  3     dog    4     yes
5  4   human    2      No
6  5   snake    0   < NA >
7  6    tree    0   < NA >
```

# Full Join

importfigs/full_join.pdf

# Properties of Full Joins

- You want an output set with all cases from both data sets
- There will be lots of "missing" values
- You don't loose anything, but working with the data is harder
  - Need to subset before plotting
  - Need to deal with potentially large missing proportion

## Practice

```
1    dat1
```

```
1    Company Earnings
2  |1        A    126345
3  |2        B    492012
4  |3        C    234512
5  |4        D    -28124
6  |5        E    128675
```

```
1    dat2
```

```
1    Company      Region
2  |1        A    Midwest
3  |2        B  Southeast
4  |3        C       West
5  |4        F      North
```

Can you:

- Left Join the data so we have all Earnings in the Output set.
- Left Join the data so we have all Regions in the Output set.
- Inner Join the data so we have no missing data.
- Full Join the data so we have everything in the Output set.

## Practice: Answer 1

```
1   merge ( x = dat1 , y = dat2 , by = " Company " ,
        all . x = TRUE )
```

```
1     Company Earnings     Region
2   1       A    126345     Midwest
3   2       B    492012   Southeast
4   3       C    234512        West
5   4       D    -28124        < NA >
6   5       E    128675        < NA >
```

- Left Join the data so we have all Earnings in the Output set.

## Practice: Answer 2

```
1   merge ( x = dat2 , y = dat1 , by = " Company " ,
         all.x = TRUE )
```

```
1     Company      Region Earnings
2   1        A     Midwest   126345
3   2        B   Southeast   492012
4   3        C        West   234512
5   4        F       North       NA
```

- Left Join the data so we have all Regions in the Output set.

## Practice: Answer 3

```
1   merge ( x = dat1 , y = dat2 , by = " Company " , all
        = FALSE )
```

```
1     Company  Earnings     Region
2   1       A    126345    Midwest
3   2       B    492012  Southeast
4   3       C    234512       West
```

- Inner Join the data so we have no missing data.

## Practice: Answer 4

```
1   merge ( x = dat1 , y = dat2 , by = " Company " , all
        = TRUE )
```

```
1     Company  Earnings     Region
2   1       A    126345    Midwest
3   2       B    492012  Southeast
4   3       C    234512       West
5   4       D    -28124      < NA >
6   5       E    128675      < NA >
7   6       F        NA      North
```

- Full Join the data so we have everything in the Output set.

## Longitudinal Data

- Data comes in 2 typical formats
  1. Wide: Columns that describe units of observation (one row per state, or per school, or per child)

     | state | region |
     |---------|-------|
     | Alabama | south |
     | Alaska | north |

     ⋮

  2. Long: Repeated observations, several times for each unit.

     | year | state | poverty |
     |------|-----------|---------|
     | 2000 | Alabama | 13 |
     | 2001 | Alabama | 12 |

     ⋮

     | 2017 | Wisconsin | 11 |

- We often want to merge the information about the units from the wide format onto the longitudinal data that is in the long format.

## Example: Merging Wide data onto Longitudinal Data

The longitudinal data is about children measured at 3 time points

```
1   dat_long
```

```
1     child_id Time FSIQ
2  1       110    1   98
3  2       110    2  102
4  3       110    3  104
5  4       210    1   89
6  5       210    2   91
7  6       210    3   95
```

Separate data about the education of parents is available for some children

```
1   dat_edu
```

```
1     child_id par_edu
2  1       210      BA
3  2       110      HS
```

# Longitudinal Data: Long

```
1   merge ( x = dat_long , y = dat_edu , by =
        " child_id " , all = TRUE )
```

```
1    child_id Time FSIQ par_edu
2  |1       110   1   98       HS
3  |2       110   2  102       HS
4  |3       110   3  104       HS
5  |4       210   1   89       BA
6  |5       210   2   91       BA
7  |6       210   3   95       BA
```

- This is a full join
- No problems encountered, result *seems* adequate.

## Points of caution in the full join

**1** If information about some families is missing from the wide data,
then missing values will be created in the result
Example:
We change the wide data by removing one child

```
1    child_id par_edu
2  1    210      BA
```

```
1  merge(x = dat_long, y = dat_edu2, by =
       "child_id", all = TRUE)
```

```
1    child_id Time FSIQ par_edu
2  1     110    1   98    <NA>
3  2     110    2  102    <NA>
4  3     110    3  104    <NA>
5  4     210    1   89     BA
6  5     210    2   91     BA
7  6     210    3   95     BA
```

## Points of caution in the full join ...

2 If wide data includes information about children/families that are not tracked in the long data, then the full join will create "extra" all missing lines in the longitudinal part.
Example:
We only change  dat_edu  by inserting additional rows for some children.

```
1     child_id par_edu
2  1      210      BA
3  2      110      HS
4  3      400      ES
5  4      501      HS
```

Why would this happen in real life? Suppose these are child/parent data rows from a different study in which some of the children participated.

## Points of caution in the full join ...

```
1   merge ( x = dat_long , y = dat_edu2 , by =
        " child_id " , all = TRUE )
```

```
1     child_id Time FSIQ par_edu
2   1      110    1   98      HS
3   2      110    2  102      HS
4   3      110    3  104      HS
5   4      210    1   89      BA
6   5      210    2   91      BA
7   6      210    3   95      BA
8   7      400   NA   NA      ES
9   8      501   NA   NA      HS
```

## Points of caution in the full join ...

3. Some users may prefer to think of this as a left join, keeping only rows about children in a study (and omitting rows about families of children who are not in the study)

```
merge ( x = dat_long , y = dat_edu2 , by =
    " child_id " , all.x = TRUE , all.y = FALSE )
```

```
  child_id Time FSIQ par_edu
1    110    1   98     HS
2    110    2  102     HS
3    110    3  104     HS
4    210    1   89     BA
5    210    2   91     BA
6    210    3   95     BA
```

# Longitudinal Data: Long Data by Long Data

```
1    dat_long1
```

```
1     child_id Time FSIQ
2   1       110    1   98
3   2       110    2  102
4   3       110    3  104
5   4       210    1   89
6   5       210    2   91
7   6       210    3   95
```

```
1    dat_long2
```

```
1     child_id Time Reaction
2   1       210    1     0.34
3   2       210    2     0.28
4   3       210    3     0.19
5   4       110    1     0.33
6   5       110    2     0.32
7   6       110    3     0.28
```

Notice here, the dangers are repeating ID's in both data sets.

## Longitudinal Data: Long Data by Long Data

```
1   head(merge(x = dat_long1, y = dat_long2, by =
        "child_id", all.x = TRUE), 12)
```

```
1       child_id Time.x FSIQ Time.y Reaction
2   1        110      1   98      1     0.33
3   2        110      1   98      2     0.32
4   3        110      1   98      3     0.28
5   4        110      2  102      1     0.33
6   5        110      2  102      2     0.32
7   6        110      2  102      3     0.28
8   7        110      3  104      1     0.33
9   8        110      3  104      2     0.32
10  9        110      3  104      3     0.28
11  10       210      1   89      1     0.34
12  11       210      1   89      2     0.28
13  12       210      1   89      3     0.19
```

This is **WRONG!!!** look closely.

## Longitudinal Data: Long Data by Long Data

To solve our problem we provide multiple Keys to the "by" argument:

```
1   merge ( x = dat_long1 , y = dat_long2 , by =
        c ( "child_id" , "Time" ) , all.x = TRUE )
```

```
1     child_id Time FSIQ Reaction
2   1      110    1   98     0.33
3   2      110    2  102     0.32
4   3      110    3  104     0.28
5   4      210    1   89     0.34
6   5      210    2   91     0.28
7   6      210    3   95     0.19
```

That is much better, notice the fix:

```
1   by = c ( "child_id" , "Time" )
```

## Longitudinal Data: Long Data by Long Data

An intuitive way to determine when you need to supply multiple keys to the "by" argument is to ask yourself:

- Can every occurrence of my ID variable be uniquely identified ?
- If not, which other variable is necessary to produce an uniquely identified ID ?

## Longitudinal Data: QUIZ

Which columns together create the proper uniquely identifiable key set?

```
1    dat_nat
```

```
1      ID Year Quarter population illnesses
2   1 USA 1990      Q1  10.585529  97.15840
3   2 USA 1990      Q2  10.709466  90.80678
4   3 USA 1991      Q1   9.890697  98.83752
5   4 USA 1991      Q2   9.546503 118.17312
6   5  UK 1990      Q1  10.605887 103.70628
7   6  UK 1990      Q2   8.182044 105.20216
8   7  UK 1991      Q1  10.630099  92.49468
9   8  UK 1991      Q2   9.723816 108.16900
```

# Longitudinal Data: A Useful way to Identify Keys

```
1   table(dat_nat$ID)
```

```
1   UK  USA
2    4   4
```

Not unique, we need another key

```
1   table(dat_nat$ID, dat_nat$Quarter)
```

```
1         Q1 Q2
2   UK     2  2
3   USA    2  2
```

getting closer

```
1   table(dat_nat$ID, dat_nat$Quarter,
        dat_nat$Year)
```

## Longitudinal Data: A Useful way to Identify Keys ...

```
 1   , ,  = 1990
 2
 3
 4         Q1 Q2
 5    UK    1   1
 6    USA   1   1
 7
 8   , ,  = 1991
 9
10
11         Q1 Q2
12    UK    1   1
13    USA   1   1
```

Winner! Each data point can be uniquely identified as being collected
from a country, during a year, and a quarter.

# Different Key Names

```
1   head(datX)
```

```
1      ID Year Quarter        pop illnesses
2  1 USA 1990      Q1  9.113642  84.02290
3  2 USA 1990      Q2  9.668422 118.05098
4  3 USA 1991      Q1 11.120713  95.18353
5  4 USA 1991      Q2 10.298724 106.20380
6  5  UK 1990      Q1 10.779622 106.12123
7  6  UK 1990      Q2 11.455785  98.37689
```

```
1   head(datY)
```

```
1    Country year Semester percipitation      cars
2  1      USA 1990      Q1    12.049190 111.28511
3  2      USA 1990      Q2    11.632446  76.19642
4  3      USA 1991      Q1    10.254271  89.39734
5  4      USA 1991      Q2    10.491188 109.37141
6  5       UK 1990      Q1     9.675913 108.54452
7  6       UK 1990      Q2     8.337950 114.60729
```

# Different Key Names

```
1   head(datX)
```

```
1      ID Year Quarter       pop illnesses
2   1 USA 1990     Q1 10.583188 106.91171
3   2 USA 1990     Q2  8.693201 108.23795
4   3 USA 1991     Q1  9.459614 121.45065
5   4 USA 1991     Q2 11.947693  76.53056
6   5  UK 1990     Q1 10.053590 101.49592
7   6  UK 1990     Q2 10.351663  86.57469
```

```
1   head(datY)
```

```
1     Country year Semester percipitation      cars
2   1     USA 1990       Q1      9.413120  89.50647
3   2     USA 1990       Q2      8.167623 123.30512
4   3     USA 1991       Q1     10.888139 114.02705
5   4     USA 1991       Q2     11.593488 109.42601
6   5      UK 1990       Q1     10.516855 108.26258
7   6      UK 1990       Q2      8.704328  91.88460
```

## Different Key Names

```
1   merge(x = datX, y = datY, by.x = c("ID",
        "Year", "Quarter"), by.y = c("Country",
        "year", "Semester"),all = TRUE)
```

```
1      ID Year Quarter         pop illnesses percipitation      cars
2  1   UK 1990      Q1 10.053590 101.49592     10.516855 108.26258
3  2   UK 1990      Q2 10.351663  86.57469      8.704328  91.88460
4  3   UK 1991      Q1  9.329023 105.53303     10.054616 104.76248
5  4   UK 1991      Q2 10.277954 115.89963      9.215351 110.21258
6  5  USA 1990      Q1 10.583188 106.91171      9.413120  89.50647
7  6  USA 1990      Q2  8.693201 108.23795      8.167623 123.30512
8  7  USA 1991      Q1  9.459614 121.45065     10.888139 114.02705
9  8  USA 1991      Q2 11.947693  76.53056     11.593488 109.42601
```

# Matching Missing

```
1    datX
```

```
1     ID cars      fear
2  1 111    6  90.61873
3  2 112    5  97.35806
4  3  NA    7  91.15475
5  4 114    6  94.99807
6  5 115    5 106.76902
7  6 116    5 114.09072
8  7  NA    9 109.50524
```

```
1    datY
```

```
1     ID pets
2  1 111    5
3  2  NA    4
4  3 113    4
5  4 114    8
6  5 115    6
7  6  NA    4
8  7 117    7
```

## Matching Missing: The Problem

```
1  merge ( x = datX , y = datY , by = " ID " , all . x =
       TRUE )
```

```
1     ID cars      fear pets
2  1 111    6  90.61873    5
3  2 112    5  97.35806   NA
4  3 114    6  94.99807    8
5  4 115    5 106.76902    6
6  5 116    5 114.09072   NA
7  6  NA    7  91.15475    4
8  7  NA    7  91.15475    4
9  8  NA    9 109.50524    4
10 9  NA    9 109.50524    4
```

**Oops!** That is a dangerous outcome: NA columns were merged together

## Matching Missing: The Remedy

incomparables   to the rescue

```
1   merge ( x = datX , y = datY , by = "ID" , all=
        FALSE , incomparables = "NA" )
```

```
1      ID cars      fear pets
2  |1 111     6  90.61873     5
3  |2 114     6  94.99807     8
4  |3 115     5 106.76902     6
```

That is much better! Always remember to use the   incomparables
argument if you have any missing data on keys.

# Kutils::mergeCheck

```
1   df1
```

```
1    id          x
2   1  1 -0.9806329
3   2  2  0.6873321
4   3  3 -0.5050435
5   4  4  2.1577198
6   5  5 -0.5997976
7   6  6 -0.6945467
8   7  7  0.2239254
```

```
1   df2
```

```
1    id          x
2   1  2 -1.1562233
3   2  3  0.4224185
4   3  4 -1.3247553
5   4  5  0.1410843
6   5  6 -0.5360480
7   6  9 -0.3116061
8   7 10  1.5561096
```

## Kutils::mergeCheck

```
1  library ( kutils )
2  mergeCheck ( df1 , df2 , by = " id " )
```

```
1  Merge difficulties detected
2
3  Unmatched cases from df1 and df2 :
4  df1
5    id        x
6  1  1 -0.9806329
7  7  7  0.2239254
8  df2
9    id        x
10 6  9 -0.3116061
11 7 10  1.5561096
```

- mergeCheck alerts you to potential merging issues
- ID 1 and 7 in the X data frame dont have matching Y IDs
- Further, ID 9 and 10, in the Y data frame dont have matching X IDs

## Kutils::mergeCheck

```
1    df1
```

```
1     idx          x
2   1    1 -0.44803329
3   2    2  0.32112354
4   3    3 -1.23017225
5   4    4 -1.32405869
6   5    5  1.26124227
7   6   NA  1.31923172
8   7  NaN -0.08075376
```

```
1    df2
```

```
1     idy          x
2   1    2 -0.50508981
3   2    3 -0.05215359
4   3    4  0.62886063
5   4    5  2.18000240
6   5    6 -0.06901731
7   6    9  1.54486360
8   7   10  1.32145202
```

## Kutils::mergeCheck

```
1   mergeCheck ( df1 , df2 , by.x = "idx" , by.y =
        "idy" )
```

```
1    Merge difficulties detected
2
3    Unacceptable key values
4    df1
5       idx           x
6    6   NA   1.31923172
7    7  NaN  -0.08075376
8    Unmatched cases from df1 and df2 :
9    df1
10      idx           x
11   1    1  -0.44803329
12   6   NA   1.31923172
13   7  NaN  -0.08075376
14   df2
15      idy           x
16   5    6  -0.06901731
17   6    9   1.54486360
18   7   10   1.32145202
```

# Kutils::mergeCheck ...

- In this situation we are warned of:
    - Unacceptable key values: NA and NaN
    - Again, unmatched IDs: 1,6,7,9,10

# Kutils::mergeCheck

Load `library(kutils)` and run `example(mergeCheck)` to learn more about the function. Our kutils package has much more to offer! check out the kutils help page with `help(package = "kutils")`

# More Information

- The CRMDA has a guide available on merges:
  - https://crmda.ku.edu/guide-41-merge_R_SQL