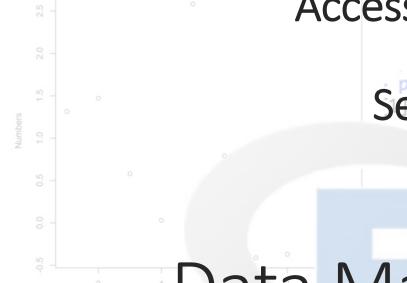
Online Academic Data Analysis Bootcamp Using Open-Access Program R



Session 3

Data Management

> mat<-data.frame(a= |=10, |=10, |=10, |=10, |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10| |=10|

a b c
1 0.68095071 2.62502002 0.2000328
2 -1.41886274 1.77124066 -0.3053134
3 -0.02606923 -0.02631055 0.7104818
4 -0.43837086 1.61360146 0.1967035

Patrick Njage (BSc, MSc, Mstat, PhD) 720 1.32386062 -1.9040575
Technical University of Denmark 61093695 3.40457075 -2.1482692
Academic Data Analysts (WWW.ACADEMICDATAANALYSTS 20RG/)

Institute of Statistics and Data Science (ISADS) (<u>WWW.ISADS-ORG</u>

10 0.78591380 3.31195551 0.3432750

lis a collaborative project with many contributors.

Type 'contributors()' for more information and 'citation()' on how to cite R or R packages in publication

[ype 'demo()' for some demos, 'help()' for on-line help, of
'help.start()' for an HTML browser interface to help.
[ype 'q()' to quit R.

Juring startup - Warning messages:
1: Setting LC_CULATE foiled, using "C"
2: Setting LC_CULATE foiled, using "C"
3: Setting LC_TIME failed, using "C"
4: Setting LC_MUESSAGES failed, using "C"
5: Setting LC_MUETARY failed, using "C"
[R.app GUI 1.70 (7543) x86_64-apple-darwin15.6.0]

WARNING: You're using a non-UTF8 locale, therefore only ASCII characters will work.

Please read R for Mac OS X FAQ (see Help) section 9 and adjust your system preferen
accordingly.

[Workspace restored from /Users/inesfragata/.RData]

Joseph Ndiba (BSc, MSc) Climate Change & Adaptation

norm(n=10, mean=0, sd=1)] 0.5263934 0.4560433 -2.2070991 0.9648372 -1.0221512 0.7316200 -0.1323901 -0.8647938] 0.8693620 -0.2794973

Outline

Part 1: Common Data Management Options in R

Part 2: Transforming Your Data with dplyr

Part 1: Common Data Management Options in R Outline

- Sorting Data
- Merging Data
- Aggregating Data
- Reshaping Data
- Subsetting Data
- Data Type Conversion

Sorting Data

detach(mtcars)

```
order() function used to sort dataframes- default is sorting ASCENDING.
Prepend the sorting variable by a minus sign to indicate DESCENDING order. Here are some examples.
# sorting examples using the mtcars dataset
attach(mtcars)
# sort by mpg
newdata <- mtcars[order(mpg),]</pre>
# sort by mpg and cyl
newdata <- mtcars[order(mpg, cyl),]</pre>
#sort by mpg (ascending) and cyl (descending)
newdata <- mtcars[order(mpg, -cyl),]</pre>
```

Merging Data: Adding Columns

To merge two data frames (datasets) horizontally, use the **merge** function. In most cases, you join two data frames by one or more common key variables (i.e., an inner join).

```
# merge two data frames by ID
total <- merge(data frameA,data frameB,by="ID")</pre>
```

merge two data frames by ID and Country
total <- merge(data frameA,data frameB,by=c("ID","Country"))</pre>

Horizontal joins like this are typically used to add variables to a data frame.

Merging Data: Adding Columns

Examples: Two simple dataframes

```
Dataframe "x" and "y"
#Dataframe "x"
name <- c("John", "Paul", "George", "Ringo", "Stuart", "Pete")
instrument <- c("guitar", "bass", "guitar", "drums", "bass", "drums")
x <- data.frame(name, instrument)
#Dataframe "y"
 name <- c("John", "Paul", "George", "Ringo", "Brian")
band <- c("TRUE", "TRUE", "TRUE", "TRUE", "FALSE")
y <- data.frame(name, band)
Merging Data: Adding Columns
total <- merge(x, y,by="name")
```

```
> X
    name instrument
    John
             guitar
    Pau1
               bass
             guitar
  George
   Ringo
              drums
               bass
5 Stuart
              drums
    Pete
          band
    name
    John
          TRUE
    Paul
          TRUE
3 George
          TRUE
  Ringo
         TRUE
   Brian FALSE
> total
    name instrument band
             guitar TRUE
1 George
             quitar TRUE
    John
    Paul
               bass TRUE
   Ringo
              drums TRUE
```

Merging Data: Adding Rows

To join two data frames (datasets) vertically, use the **rbind** function. The two data frames **must** have the same variables, but they do not have to be in the same order.

```
total <- rbind(data frameA, data frameB)

dataframeA <- mtcars[1:4,]

dataframeB <- mtcars[5:8,]

dataframe_both <- rbind(dataframeA, dataframeB) # combine objects as rows
```

If data frameA has variables that data frameB does not, then either:

- Delete the extra variables in data frameA or
- Create the additional variables in data frameB and set them to NA (missing) before joining them with **rbind()**.

Merging Data: Adding Rows

```
> dataframeA
              mpg cyl disp hp drat wt qsec vs am gear carb
                   6 160 110 3.90 2.620 16.46 0 1
Mazda RX4
             21.0
Mazda RX4 Wag 21.0 6 160 110 3.90 2.875 17.02 0 1
             22.8 4 108 93 3.85 2.320 18.61 1 1 4
Datsun 710
Hornet 4 Drive 21.4 6 258 110 3.08 3.215 19.44 1 0
> dataframeB
                 mpg cyl disp hp drat wt qsec vs am gear carb
Hornet Sportabout 18.7 8 360.0 175 3.15 3.44 17.02 0 0
Valiant
                18.1 6 225.0 105 2.76 3.46 20.22 1 0
             14.3 8 360.0 245 3.21 3.57 15.84 0 0
Duster 360
Merc 240D
                24.4
                      4 146.7 62 3.69 3.19 20.00 1 0
> dataframe_both
                 mpg cyl disp hp drat
                                        wt gsec vs am gear carb
Mazda RX4
                21.0 6 160.0 110 3.90 2.620 16.46
Mazda RX4 Wag
                21.0 6 160.0 110 3.90 2.875 17.02 0 1
                22.8 4 108.0 93 3.85 2.320 18.61 1 1
Datsun 710
Hornet 4 Drive
                21.4
                      6 258.0 110 3.08 3.215 19.44 1 0
Hornet Sportabout 18.7 8 360.0 175 3.15 3.440 17.02 0 0
Valiant
                18.1
                      6 225.0 105 2.76 3.460 20.22 1 0
             14.3 8 360.0 245 3.21 3.570 15.84 0 0
Duster 360
                      4 146.7 62 3.69 3.190 20.00 1 0
Merc 240D
                24.4
```

Aggregating Data

It is relatively easy to collapse data in R using one or more BY variables and a defined function.

Example: aggregate data frame mtcars by cyl and vs, returning means for numeric variables

attach(mtcars) # to make the variables in the dataframe directy available in the workspace aggdata <-aggregate(mtcars, by=list(cyl,vs), FUN=mean, na.rm=TRUE)

print(aggdata)

detach(mtcars)

When using the aggregate() function, the by variables must be in a list (even if there is only one).

Reshaping Data: Transpose

Use the t() function to transpose a matrix or a data frame. In the later case, rownames become variable (column) names.

example using built-in dataset mtcars
t(mtcars)

Hadley Wickham has created a comprehensive package called <u>reshape</u> to massage data. Both an <u>introduction</u> and <u>article</u> are available. There is even a <u>video</u>!

Basically, you "melt" data so that each row is a unique idvariable combination. Then you "cast" the melted data into any shape you would like.

Creating the data frame

```
id <- c(1, 1, 2, 2)
time <- c(1, 2, 1, 2)
measure_1 <- c(5, 3, 6, 2)
measure_2 <- c(6, 5, 1, 4)
mydata <- data.frame(id, time, measure 1, measure 2</pre>
```

```
> mydata
  id time measure_1 measure_2
1 1 1 5 6
2 1 2 3 5
3 2 1 6 1
4 2 2 2 4
>
```

Melt function

library(reshape) mdata <- melt(mydata, id=c("id","time"))

You must specify the variables needed to uniquely identify each measurement (ID and Time)

The variable indicating the measurement variable names (measure_1 or measure_2) is created for you automatically.

Cast function

Now that the data is in a melted form, it can be recast into any shape, using the cast() function.

newdata <- cast(md, formula, FUN)</pre>

Cast the melted data # cast(data, formula, function)

Example 1. Average variable for each subject

subjmeans <- cast(mdata, id~variable, mean)

Example 2. Average variable for each time

timemeans <- cast(mdata, time~variable, mean)

```
> subjmeans
id measure_1 measure_2
1 1 4 5.5
2 2 4 2.5
```

```
> timemeans
  time measure_1 measure_2
1     1     5.5     3.5
2     2     2.5     4.5
>
```

Reshaping a Dataset

With Aggregation Without Aggregation mydata X1 X2 Time cast(md, id+time~variable) Time 3 cast(md, id~variable, mean) 3 2 X1 X2 2.5 (d) (a) cast(md, id+variable~time) md <- melt(mydata, id=c("id", "time")) cast(md, time~variable, mean) Time1 Time 2 Variable Variable X1 X1 5 X2 5 X2 2 X1 3 5.5 X1 2 4.5 X1 6 2.5 X2 X1 2 (b) X2 6 (e) 2 5 X2 cast(md, id~time, mean) X2 cast(md, id~variable+time) X2 Time2 Time1 X1 X1 X2 X2 Time2 Time1 Time2 5.5 3.5 2 (c)

(f)

Subsetting Data

- R has powerful indexing features for accessing object elements.
- •These features can be used to select and exclude variables and observations.
- This section demonstrates:
 - how to keep or delete variables,
 - how to keep or delete observations, and
 - how to take random samples from a dataset.

Subsetting Data: Selecting (Keeping) Variables

```
# select variables v1, v2, v3
myvars <- c("v1", "v2", "v3")
newdata <- mydata[myvars]</pre>
myvars mtcars <- c("mpg", "cyl", "disp")
newdata mtcars <- mtcars[myvars mtcars]</pre>
# select 1st and 5th thru 10th variables
newdata <- mydata[c(1,5:10)]
newdata2_mtcars <- mtcars[c(1,5:10)]</pre>
```

Subsetting Data: Excluding (DROPPING) Variables

```
# exclude variables v1, v2, v3
myvars <- names(mydata) %in% c("v1",
"v2", "v3")
newdata_mtcars <- mydata[!myvars]</pre>
```

exclude 3rd and 5th variable newdata <- mydata[c(-3,-5)]

delete variables v3 and v5 mydata\$v3 <- mydata\$v5 <- NULL

```
Example: exclude variables am, gear, carb

myvars <- names(mtcars) %in% c("am", "gear",
"carb")
newdata_mtcars <- mtcars[!myvars]

# exclude 3rd and 5th variable
newdata_mtcars_2 <- mtcars[c(-3,-5)]

# delete variables am and gear from DataframeA
dataframeA$am <- dataframeA$gear <- NULL
```

Subsetting Data: Selecting Observations

```
# first 5 observations
newdata mtcars 5 <- mtcars[1:5,]
# based on variable values
newdata_mtcars_carb_mpg <- mtcars[ which(mtcars$carb==2 & mtcars$mpg >
20.09), ]
# or
attach(mtcars)
newdata_mtcars_carb_mpg <- mtcars[ which(mtcars$carb==2 & mtcars$mpg >
20.09), ]
detach(mtcars)
```

Subsetting Data: Selection using the Subset Function

The **subset()** function is the easiest way to select variables and observations. In the following example, we select all rows that have a value of weight greater than or equal to 3.610 or wt less then 2.58. We keep the mpg, cyl and disp columns.

```
# using subset function newdata_mtcars_wt <- subset(mtcars, wt >= 3.610 | wt < 2.58, select=c(mpg, cyl, disp))
```

In the next example, we select cars with miles per gallon greater than 20.09 and 2 carburetors but we keep variables horsepower *through* 1/4 mile time (hp, drat, wt and qsec).

```
# using subset function (part 2)
newdata_mtcars_range <- subset(mtcars, mtcars$carb==2 & mtcars$mpg > 20.09,
select=hp:qsec)
```

Subsetting Data: Random Samples

Use the **sample()** function to take a **random sample of size n** from a dataset.

Example: take a random sample of size 15 from the dataset *mtcars* sample without replacement

mysample mtcars <- mtcars[sample(1:nrow(mtcars), 15, replace=FALSE),]

Part 2: Transforming Your Data with dplyr

Outline

Dplyr is aimed at simplifying manipulating, sorting, summarizing, and joining data frames.

Basic dplyr package functions are introduced here including:

select() selects variables

filter() provides basic filtering capabilities

group_by() groups data by categorical levels

summarise() summarizes data by functions of choice

arrange() orders data

join() joins separate dataframes

Dplyr: package and %>% Operator

Package Utilizedinstall.packages("dplyr")library(dplyr)

%>% Operator

- •Although not required, the tidyr and dplyr packages make use of the pipe operator %>%
- •Key advantage: ability to string multiple functions together by incorporating %>%.
- •For instance a function to filter data can be written as:

filter(data, variable == numeric_value)

Same as

data %>% filter(variable == numeric_value)

Dplyr: select() function

Objective: Reduce dataframe size to only desired variables for current task

Description: When working with a sizable dataframe, often we desire to only assess specific variables. The select() function allows you to select and/or rename variables.

Function: select(data, ...)

Same as: data %>% select(...)

Arguments:

data: data frame

...: call variables by name or by function

Dplyr: select() function

Read the provided "expenditure" data into R

Example: our goal is to only assess the 5 most recent years worth of expenditure data. Applying the select() function we can select only the variables of concern.

sub.exp <- expenditures %>% select(Division, State, X2007:X2011)

head(sub.exp) # for brevity only display first 6 rows

```
Division
                 State
                          X2007
                                    X2008
                                             X2009
                                                       X2010
                                                                X2011
               Alabama
                        6245031
                                 6832439
                                           6683843
                                                     6670517
                                                              6592925
                        1634316
2
3
4
5
               Alaska
                                 1918375
                                           2007319
                                                     2084019
                                                              2201270
                        7815720
               Arizona
                                  8403221
                                           8726755
                                                     8482552
                                                              8340211
             Arkansas
                        3997701
                                 4156368
                                                              4578136
           California 57352599 61570555 60080929 58248662
                                                             57526835
             Colorado 6579053
                                 7338766
                                           7187267
                                                     7429302
                                                              7409462
```

Dplyr: select() function

We can also apply some of the special functions within select(). For instance we can select all variables that start with 'X':

head(expenditures %>% select(starts_with("X")))

```
X2000
                             X2001
                                      X2002
                                                X2003
  X1980
           X1990
                                                         X2004
         2275233
                  4176082
                           4354794
                                    4444390
                                              4657643
1146713
                                                       4812479
 377947
          828051
                 1183499
                           1229036
                                    1284854
                                              1326226
                                                       1354846
                 4288739
                                              5892227
        2258660
                           4846105
                                    5395814
                                                       6071785
 949753
         1404545
                  2380331
                           2505179
                                    2822877
                                              2923401
                                                       3109644
9172158 21485782 38129479 42908787 46265544 47983402 49215866
1243049
        2451833
                 4401010
                           4758173
                                    5151003
                                              5551506
                                                       5666191
   X2005
            X2006
                     X2007
                              X2008
                                       X2009
                                                 X2010
                                                          X2011
         5699076
                   6245031 6832439
                                     6683843
                                                        6592925
5164406
                                               6670517
         1529645
                   1634316
                           1918375
                                     2007319
                                                        2201270
 1442269
                                               2084019
                                     8726755
 6579957
          7130341
                   7815720
                            8403221
                                               8482552
                                                        8340211
          3808011
                   3997701
                            4156368
                                     4240839
                                               4459910
                                                        4578136
50918654 53436103 57352599 61570555 60080929 58248662 57526835
 5994440
          6368289
                   6579053 7338766
                                     7187267
                                              7429302
                                                       7409462
```

Dplyr: filter() function

Objective: Reduce rows/observations with matching conditions

Description: Filtering data is a common task to identify/select observations in which a particular variable matches a specific value/condition. The filter() function provides this capability.

Function: filter(data, ...)

Same as: data %>% filter(...)

Arguments:

data: data frame

...: conditions to be met

Dplyr: filter() function

Example: Continuing with the sub.exp dataframe which includes only the recent 5 years worth of expenditures, we can filter by Division:

```
sub.exp %>% filter(Division == 6)
```

```
Division State X2007 X2008 X2009 X2010 X2011 6 Alabama 6245031 6832439 6683843 6670517 6592925
```

We can apply multiple logic rules in the filter() function such as:

- < Less than != Not equal to
- > Greater than %in% Group membership
- == Equal to is.na is NA
- <= Less than or equal to !is.na is not NA</pre>
- >= Greater than or equal to &,|,! Boolean operators

Dplyr: filter() function

Filtering by multiple criteria within a single logical expression (starwars is an R dataset)

```
filter(starwars, hair_color == "none" & eye_color == "black")
```

```
A tibble: 9 x 14
                                                   eye_color birth_year sex
           height mass hair color skin_color
                                                                                gender
  name
            <int> <db1>
                                                                   <db1> <chr>
  <chr>
                         <chr
                                    <chr>>
                                                                                <chr>
                                                   b1ack
                                                                                mascul~
                     68 none
1 Nien Nu~
              160
                                                                         male
                                    grey
                                    white, blue
                                                                                mascul~
              122
                                                   black.
2 Gasgano
                        none
                                                                         male
                     87 none
                                                                                mascul~
                                                                         male
3 Kit Fis∼
              196
                                    green
                                                   black
 Plo Koon
              188
                     80 none
                                                   black.
                                                                      22 male
                                                                                mascul~
                                    orange
                                                   black.
                                                                         male
 Lama Su
              229
                      88\none
                                                                                mascul~
                                    grey
              213
                                                   black
                                                                         female femini~
 Taun We
                        none
                                    grey
                                    red, blue, whA
                                                   black
  Shaak Ti
              178
                     57 none
                                                                         female femini~
8 Tion Me~
                                                                         male
              206
                     80 none
                                    grey
                                                   black
                                                                                mascul~
                                                   black
 BB8
                                                                      NA none
                                                                                mascul~
                         hone
                                    none
  ... with 5 more variables b
                                                    ies∆ chr>, films <list>,
                                ⊲eworld <chr>.
```

Dplyr: group_by() function

Objective: Group data by categorical variables

Description: Often, observations are nested within groups or categories and our goals is to perform statistical analysis both at the observation level and also at the group level.

The group_by() function allows us to create these categorical groupings.

Function: group_by(data, ...)

Same as: data %>% group_by(...)

Dplyr: group_by() function

Example: The group_by() function is a silent function. No observable manipulation of the data is performed after applying the function.

Only change: on top of the actual dataframe, an indicator of what variable the data is grouped by will be provided.

The real importance of the group_by() function comes when we perform summary statistics which we will cover shortly.

group.exp <- sub.exp %>% group_by(Division)

head(group.exp)

```
Division [4]
groups:
 Division State
                         X2007
                                  X2008
                                            X2009
                                                     X2010
                                                              X2011
    <int> <fct>
                         <int>
                                  <int>
                                            <int>
                                                     <int>
                                                               <int>
                       6245031
                                6832439
                                         6683843
        6 Alabama
                                                   6670517
                                                            6592925
        9 Alaska
                                                            2201270
                       1634316
                                1918375
                                         2007319
                                                   2084019
                       7815720
        8 Arizona
                                8403221
                                         8726755
                                                   8482552
                                                            8340211
        7 Arkansas
                       3997701
                                4156368
                                         4240839
                                                            4578136
        9 California 57352599 61570555
                                         60080929 58248662 57526835
        8 Colorado
                       6579053
                                7338766
                                         7187267
                                                            7409462
```

Objective: Perform summary statistics on variables

Description: key goal of data management is to be able to support statistical analysis on the data.

The summarise() function allows us to perform the majority of the initial summary statistics when performing exploratory data analysis.

Function: summarise(data, ...)

Same as: data %>% summarise(...)

Examples

Lets get the mean expenditure value across all states in 2011

sub.exp %>% summarise(Mean_2011 = mean(X2011))

Mean_2011 1 14441473

```
Examples: some more summary stats

sub.exp %>% summarise(Min = min(X2011, na.rm=TRUE),

Median = median(X2011, na.rm=TRUE),

Mean = mean(X2011, na.rm=TRUE),

Var = var(X2011, na.rm=TRUE),

SD = sd(X2011, na.rm=TRUE),

Max = max(X2011, na.rm=TRUE))
```

```
Min Median Mean Var SD Max
1 2201270 7001194 14441473 4.503461e+14 21221360 57526835
```

Previous slide: useful summaries. Comparison of summary statistics at multiple levels reveals important insights.

This is where the group_by() function comes in.

Group by Division and see how the different regions compared in by 2010 and 2011.

Dplyr: arrange() function

Objective: Order variable values

Description: Often, we desire to view observations in rank order for a particular variable(s). The arrange() function allows us to order data by variables in accending or descending order.

Function: arrange(data, ...)

Same as: data %>% arrange(...)

Arguments:

data: data frame

...: Variable(s) to order

^{*}use desc(x) to sort variable in descending order

Dplyr: arrange() function

Examples

Sort mtcars data by cylinder and displacement mtcars[with(mtcars, order(cyl, disp)),]

Same result using arrange: no need to use with(), as the context is implicit
 NOTE: plyr functions do NOT preserve row.names

 arrange(mtcars, cyl, disp)

Let's keep the row.names in this example myCars = cbind(vehicle=row.names(mtcars), mtcars) arrange(myCars, cyl, disp)

Sort with displacement in descending order arrange(myCars, cyl, desc(disp))

Objective: Join two datasets together

Description: Often we have separate dataframes that can have common and differing variables for similar observations and we wish to join these dataframes together.

The multiple xxx_join() functions provide multiple ways to join dataframes.

```
Description: Join two datasets

Function:

inner_join(x, y, by = NULL)

left_join(x, y, by = NULL)

right_join(x, y, by = NULL)

full_join(x, y, by = NULL)

semi_join(x, y, by = NULL)

anti_join(x, y, by = NULL)
```

Arguments:

x,y: data frames to join

by: a character vector of variables to join by. If NULL, the default, join will do a natural join, using all variables with common names across the two tables.

```
Examples: Two simple dataframes
Dataframe "x" and "y"
#Dataframe "x"
name <- c("John", "Paul", "George", "Ringo", "Stuart", "Pete")
instrument <- c("guitar", "bass", "guitar", "drums", "bass", "drums")
x <- data.frame(name, instrument)
#Dataframe "y"
name <- c("John", "Paul", "George", "Ringo", "Brian")
band <- c("TRUE", "TRUE", "TRUE", "TRUE", "FALSE")
y <- data.frame(name, band)
```

```
> y
name band
1 John TRUE
2 Paul TRUE
3 George TRUE
4 Ringo TRUE
5 Brian FALSE
```

```
> X
name instrument
1 John guitar
2 Paul bass
3 George guitar
4 Ringo drums
5 Stuart bass
6 Pete drums
```

inner_join(): Include only rows in both x and y that have a matching value inner_join(x,y)

```
Joining, by = "name"
name instrument band
1 John guitar TRUE
2 Paul bass TRUE
3 George guitar TRUE
4 Ringo drums TRUE
```

```
left_join(): Include all of x, and matching rows of y
left_join(x,y)
```

```
Joining, by = "name"
name instrument band
1 John guitar TRUE
2 Paul bass TRUE
3 George guitar TRUE
4 Ringo drums TRUE
5 Stuart bass <NA>
6 Pete drums <NA>
```

semi_join(): Include rows of x that match y but only keep the columns from x
semi_join(x,y)

```
Joining, by = "name"
name instrument
1 John guitar
2 Paul bass
3 George guitar
4 Ringo drums
```

```
anti_join(): Opposite of semi_join
anti_join(x,y)
```

```
Joining, by = "name"
name instrument
1 Stuart bass
2 Pete drums
```

Thank You