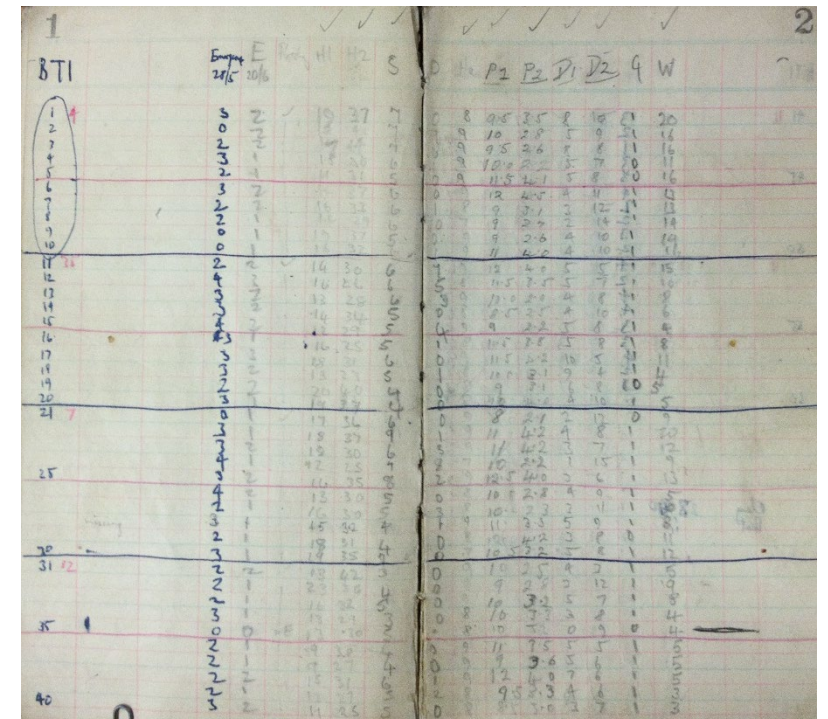# An Introduction to Statistical Computing in R

Heather Kropp

Spring 2019

GEOG 401

# Data & Physical Geography

- Environmental data has rapidly changed with evolving technology
  - Satellites, sensors, sUAS, and citizen science produce more data than ever
  - Many datasets becoming "big data"
- Rapidly improving computational power offers:
  - data management & storage
  - more sophisticated statistics

# Statistical computing

- Interaction between statistics, numerical analysis, and data manipulation

- R: program for statistical computing



The R Project for Statistical Computing

[Home]

**Download**

CRAN

**R Project**

About R
Logo
Contributors
What's New?
Reporting

## Getting Started

R is a free software environment for statistical computing and graphics. It compiles and runs on a wide variety of UNIX platforms, Windows and MacOS. To **download R**, please choose your preferred CRAN mirror.

If you have questions about R like how to download and install the software, or what the license terms are, please read our answers to frequently asked questions before you send an email.

# Why use a programming language?

- Powerful resource for handling large or complex data

- Integrates statistics, data manipulation, and GIS in a single framework

- **Reproducibility**

# Reproducibility crisis?

Home | News & Comment | Research | Careers & Jobs | Current Issue | Archive | Audio & Video | For Authors

Archive > Volume 533 > Issue 7604 > News Feature > Article

NATURE | NEWS FEATURE

## 1,500 scientists lift the lid on reproducibility

Survey sheds light on the 'crisis' rocking research.

Monya Baker

NATURE | COMMENT

## Robust research: Institutions must do their part for reproducibility

C. Glenn Begley, Alastair M. Buchan & Ulrich Dirnagl

01 September 2015

Tie funding to verified good institutional practice, and robust science will shoot up the agenda, say C. Glenn Begle M. Buchan and Ulrich Dirnagl.

Subject terms:   Research management   Institutions   Lab life

Illustration by David Parkins

# Reproducibility issues: two factors

- Statistical uncertainty and limitations
  - Inherent issues with describing the natural world

- Poor data handling and analysis
  - "Fixable"
  - Open data can make analysis transparent
  - Using a programing language allows analyses to be replicated and documented

# Data provenance

- The trajectory of data from its collection to final visualization and analysis

| notebook | raw_data.xlsx | sites_data.xlsx | sites_temp.xlsx | stemp_final.xlsx | st_final_for real.xlsx | |
|---|---|---|---|---|---|---|
| Field notebook | → Digitization into excel | → Average data by site and save file | → Add temperature data to file | → Filter data so that outlier sites are removed and save file | → Sort data for easier plotting | → Plot averaged of data |

- The management of data can readily become unwieldy and easily clouded

# Data provenance

- Using a programming language for data manipulation provides clear documentation and data provenance

notebook       raw_data.csv       Data_script.r

```
[Field notebook] → [Digitization into excel] → [Script that averages, merges temperature data, filters, sorts and plots]
```

- It still requires careful management of data files and scripts

# Reproducibility and R scripts

- A script is a file that stores all of the code for a project or operation

- A good script is written so that it can be run and the results will be the same every time

- It is a good idea to keep data workflows simple. Don't generate a lot of scripts for a project.

- Comments in a script are lines the program ignores. These are used for documentation

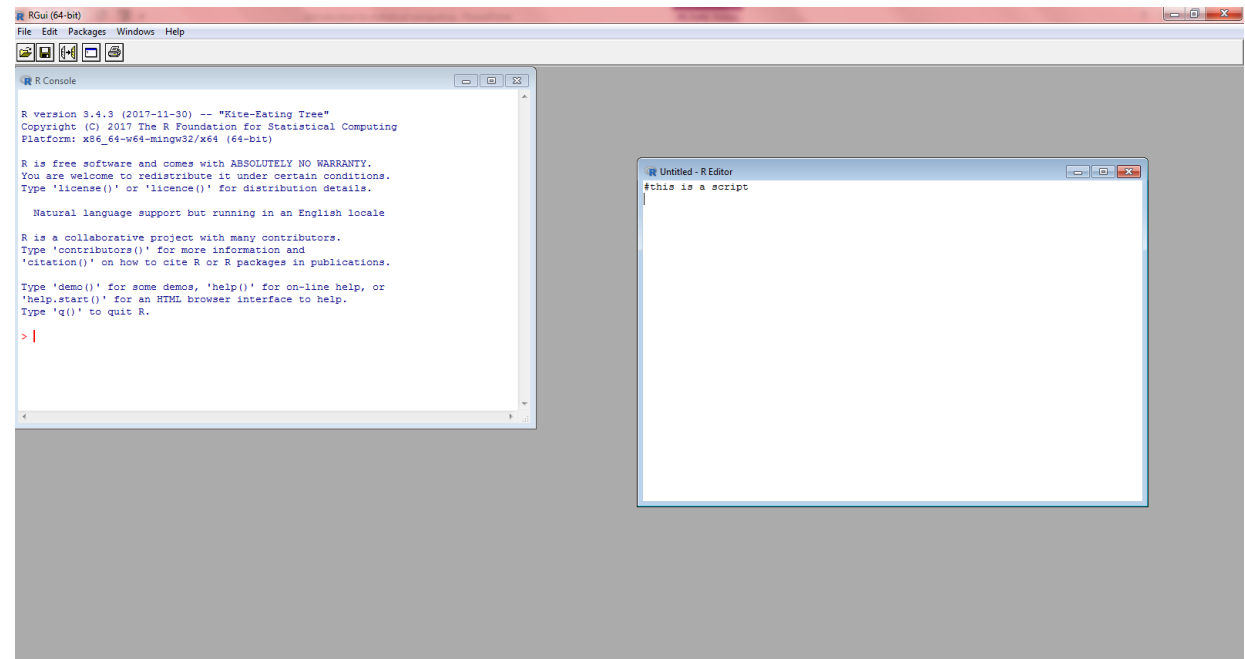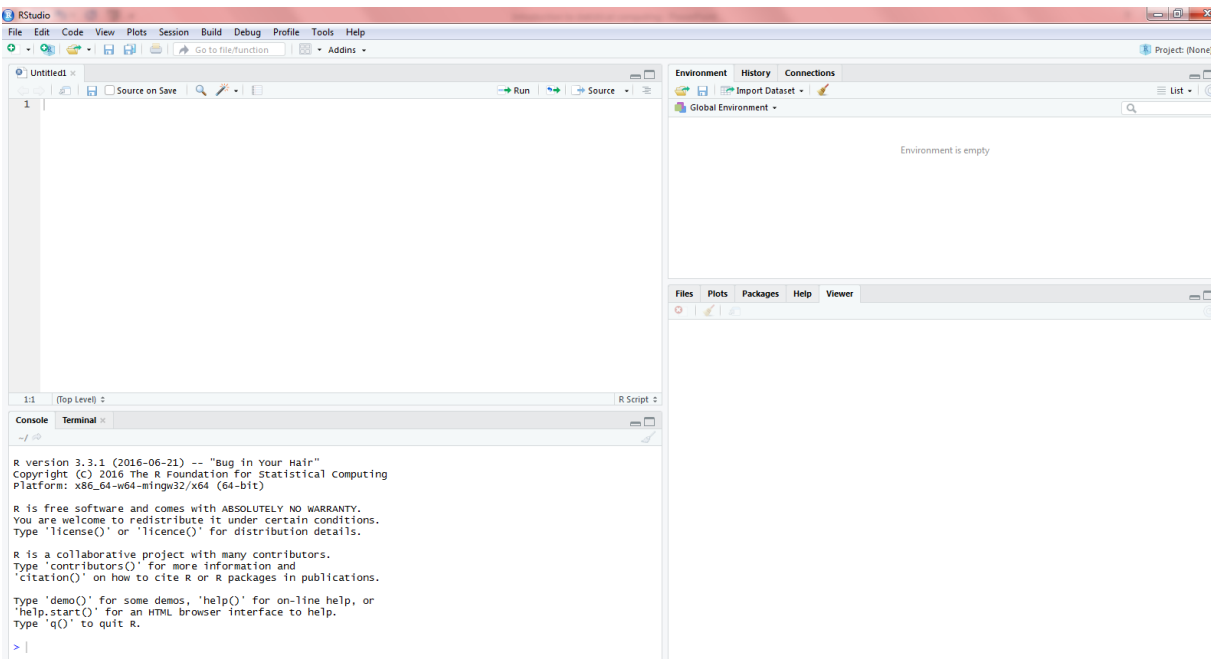- COMMENT, COMMENT COMMENT!!!!! (yes, I am yelling at you in all caps)

```
############################################################
##################This script provides examples for ######
################an introduction to R basic concepts ####
############################################################
#R is just a sophisticated calculator
6^6
5+210
3-10


#Basic data concepts in R
#Assign a name for an object
Ex <- 6
#look at object
Ex
#use the object in a calculation
Ex*5
```

# Many different ways to interface with R
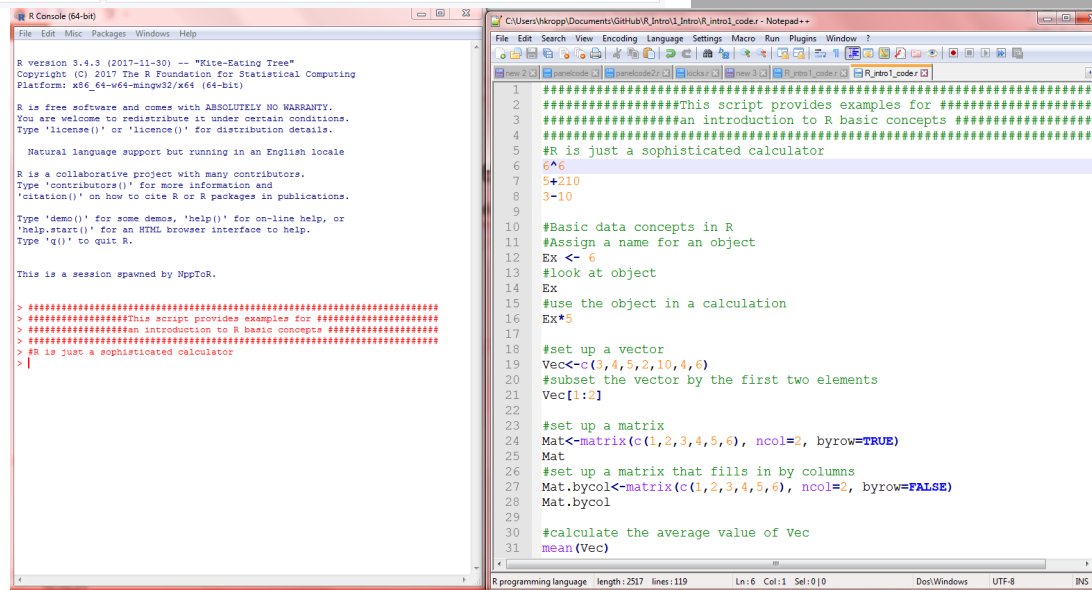
**Rstudio**

**R**

**Text editors that can send code to R**

# Many different ways to interface with R



**Rstudio**

R

**Recommended for beginners who are used to working in excel**

**Text editors that can send code to R**

# The core of R: console and your script

- R console does all of the work. It is basically a giant calculator.

- Your script is where you write code and save it. You need to send it to the console to run it

- The session that you have started in R is called your working environment. It contains everything that you have run in the console.

**script** →

**console** →

# Rstudio basics

- Rstudio also allows you to view the objects you are working with.

- This can provide a more user friendly interface that is more similar to excel

# Working with numbers in R

- R acts like a calculator:

Red=inputs from your script
Blue= output from R

```
R Console

Type 'license()' or 'licence()' for distribution details.

  Natural language support but running in an English locale

R is a collaborative project with many contributors.
Type 'contributors()' for more information and
'citation()' on how to cite R or R packages in publications.

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

> 2*2
[1] 4
> 210/5
[1] 42
> 5^5
[1] 3125
> 3+3
[1] 6
> 1+1
[1] 2
> 4-3
[1] 1
>
```

# Assigning variables in R

- Give an item a name using the <- followed by the item
  - This allows you to refer back to items without having to remember them or write huge amounts of code

- Comment your code using # in front of the line
  - Commenting allows you to keep track of what you are doing and provide reminders for later

# Vectors and matrix in R

- R automatically treats inputs like they are vectors.
- Create a vector using c()

- Set up a ma #set up a vector :rix()
  Vec<-c(3,4,5,2,10,4,6)

```
>
> Vec<-c(3,4,5,2,10,4,6)
> Mat<-matrix(c(1,2,3,4,5,6), ncol=2, byrow=TRUE)
> Mat
     [,1] [,2]
[1,]    1    2
[2,]    3    4
[3,]    5    6
> #set up a matrix that fills in by columns
> Mat.bycol<-matrix(c(1,2,3,4,5,6), ncol=2, byrow=FALSE)
> Mat.bycol
     [,1] [,2]
[1,]    1    4
[2,]    2    5
[3,]    3    6
>
```

```
10   Ex*5
11
12   #set up a vector
13   Vec<-c(3,4,5,2,10,4,6)
14
15   #set up a matrix
16   Mat<-matrix(c(1,2,3,4,5,6), ncol=2, byrow=TRUE)
17   Mat
18   #set up a matrix that fills in by columns
19   Mat.bycol<-matrix(c(1,2,3,4,5,6), ncol=2, byrow=FALSE)
20   Mat.bycol
21
22
23
```

# Functions in R

- Using matrix() is an example of a function in R

- There are a lot of "built in" functions in R that make it easier to work with data or statistics

- For example, calculating an average using mean():

```
> mean(Vec)
[1] 4.857143
>
```

```
21
22    #calculate the average value of Vec
23    mean(Vec)
24
```

# Functions in R

- A **function** typically gives you an **output** based on the inputs you give it

- The inputs needed for a function are often called **arguments**

- R will have a description of the arguments and output for its function
    - google the *function.name*
    - type help(*function.name*) in the console

mean {base} ← *function.name*  R Documentation

Arithmetic Mean

**Description**

Generic function for the (trimmed) arithmetic mean.

**Usage**

mean(x, ...) ←

## Default S3 method:
mean(x, trim = 0, na.rm = FALSE, ...) ←

Way that the function gets used and the
arguments that have some default assumptions

**Arguments**

x
    An R object. Currently there are methods for numeric/logical vectors and date, date-time and time interval objects. Complex vectors are allowed for trim = 0, only. ←

trim
    the fraction (0 to 0.5) of observations to be trimmed from each end of x before the mean is computed. Values of trim outside that range are taken as the nearest endpoint.

na.rm
    a logical value indicating whether NA values should be stripped before the computation proceeds.

...
    further arguments passed to or from other methods.

Detailed
description

- Here *x*, *trim*, and *na.rm* are the names of arguments
- The order here is important because if the names of the arguments aren't used, then they are assumed to be in this default order.
- If the default arguments are sufficient there is no need to include them in the function
- If you don't want to use the default order than you can specify

```
[1] 4.857143
> mean(na.rm=TRUE,x=Vec)
[1] 4.857143
```

```
24    #specify arguments by name and not order
25    mean(na.rm=TRUE,x=Vec)
```

# Working with vectors and matrices:

• Keep in mind that R automatically does vector/matrix math:

```
[3,]    15    18
> #multiply vector by 5
> Vec*5
[1] 15 20 25 10 50 20 30
>
```

```
30
31    #multiply vector by 5
32    Vec*5
```

```
> #set multiply matrix by vector
> Mat.scale<-matrix(c(2,2,2,3,3,3), ncol=2,byrow=TRUE)
> Mat*Mat.scale
      [,1] [,2]
[1,]    2    4
[2,]    6   12
[3,]   15   18
>
```

```
31    # multiply matrix by vector
32    Mat.scale<-matrix(c(2,2,2,3,3,3), ncol=2,byrow=TRUE)
33    Mat*Mat.scale
```

# Packages in R

- People have created thousands of **packages** to add more functions to R

- Packages allow you to download only the ones you want to use (it would take up a lot of space)

- Some functions may have the same name in different packages so be sure to note potential overlap in packages
  - only load the ones you are going to use

# Data Types

- Numeric: number can have any number of decimals

- Character: text

- Factor: text but a short identifying category name

# Reading in Data

- The easiest and most consistent way to read in data in R is through a comma separated text file (.csv)

- You need to tell R where to find the data
  - Set a working directory to always get files from one folder
  - Or specify the file path with the csv name

- File/File paths always need to be in quotes and file paths always have \\ between folders

```
#read in data file
datM<-read.csv("mountain_data.csv")
#check out data
datM
```

- Always think about your names
  - Length
  - Clarity
- Capitalization matters!

```
> datM
   Rank                          Name           Region Elev.m Prom.m Elev.ft Prom.ft
1     1                    Mt Everest      Nepal Tibet   8848   8848   29028   29028
2     2                     Aconcagua        Argentina   6962   6962   22841   22841
3     3             Mt McKinley Denali            US   6194   6138   20320   20138
4     4                   Kilimanjaro         Tanzania   5895   5885   19340   19308
5     5                Cristobal Colon        Colombia   5700   5509   18701   18074
6     6                      Mt Logan          Canada   5959   5250   19550   17224
7     7     Pico de Orizaba Citlaltepetl        Mexico   5636   4922   18491   16148
8     8                 Vinson Massif       Antarctica   4892   4892   16050   16050
9     9                   Puncak Jaya        Indonesia   4884   4884   16023   16023
10   10                   Gora Elbrus          Russia   5642   4741   18510   15554
11   11                    Mont Blanc     France Italy   4809   4696   15777   15406
12   12                      Damavand            Iran   5610   4667   18405   15311
13   13         Klyuchevskaya Volcano          Russia   4750   4649   15584   15252
14   14                  Nanga Parbat         Pakistan   8125   4608   26657   15118
15   15                     Mauna Kea            US   4205   4205   13796   13796
16   16 Jengish Chokusu ex Pik Pobedy Kyrgyzstan China   7439   4148   24406   13609
17   17                    Chimborazo          Ecuador   6267   4122   20561   13523
18   18                    Bogda Shan           China   5445   4122   17864   13523
19   19                  Namcha Barwa           China   7782   4106   25531   13471
20   20                      Kinabalu         Malaysia   4095   4095   13435   13435
21   21                    Mt Rainier            US   4393   4023   14411   13196
22   22                            K2   Pakistan China   8611   4017   28251   13179
23   23                     Ras Dejen         Ethiopia   4533   3980   15092   13090
24   24                Volcan Tajumulco       Guatemala   4220   3980   13845   13058
25   25                   Pico Bolivar        Venezuela   4981   3957   16341   12982
> |
```

# Properties of a data frame

- Typically the columns have names
- All columns are the same length
- There can be different types of data in each column

# Basic Info about a data frame

- Get dimensions

```
> dim(datM)
[1] 25  7
> Mdim<-dim(datM)
>
```

```
55  #get dimensions of the dataset
56  dim(datM)
57  #Note output is a vector of 2 values
58  #we can name this and refer to later
59  Mdim<-dim(datM)
60
61
62
```

- Names of columns

```
> names(datM)
[1] "Rank"    "Name"    "Region"  "Elev.m"  "Prom.m"  "Elev.ft" "Prom.ft"
>
```

```
60
61  #get the column names
62  names(datM)
63
64
65
```

- See what it looks like

```
> head(datM)
  Rank            Name         Region Elev.m Prom.m Elev.ft Prom.ft
1    1      Mt Everest Nepal Tibet    8848   8848   29028   29028
2    2       Aconcagua      Argentina  6962   6962   22841   22841
3    3 Mt McKinley Denali          US  6194   6138   20320   20138
4    4      Kilimanjaro       Tanzania  5895   5885   19340   19308
5    5   Cristobal Colon      Colombia  5700   5509   18701   18074
6    6         Mt Logan        Canada  5959   5250   19550   17224
>
```

```
63
64  #look at the names and first 5 rows
65  head(datM)
66
67
68
69
70
71
72
73
74
75
```

# Referring to data in data frames

- A column can be used by: *data.frame$column*

```
#look at only the name columne
datM$Name
```

- Data frames can also be refered to like matrix where [rows,columns] notation is used

  - Refer to a column without calling its name:

```
#look at name in second column
datM[,2]
```
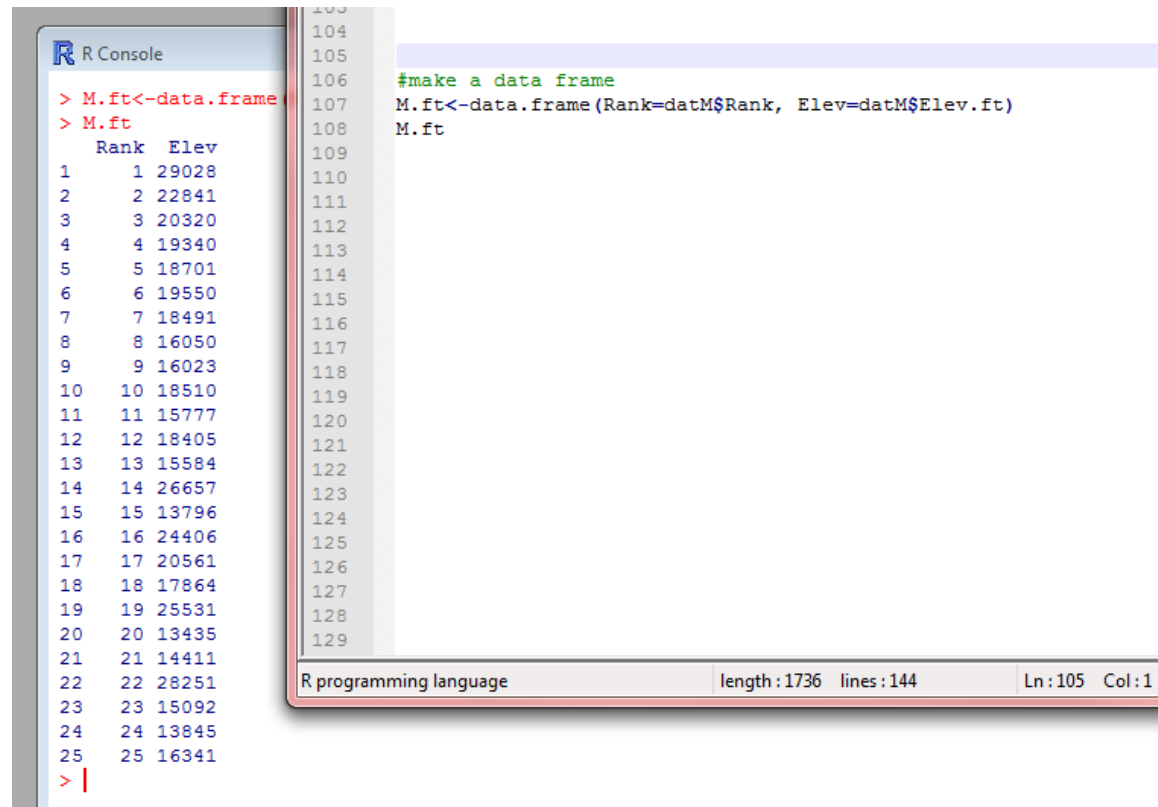
  - Multiple columns

```
#refer to multiple columns
datM[,2:4]
```

  - Rows:

```
#refer to several rows
datM[1:3,]
```

# Creating a data frame

- Use the function data.frame
  - Note: all vectors must be equal lengths

# Subset data

- Data can be subset by a characteristic
- This is done using logical expressions (see R's guide for logical expressions.
  - https://www.r-bloggers.com/logical-operators-in-r/
- Subset with brackets:
  - Mountains in the US

```
> US.M
    Rank              Name Region Elev.m Prom.m Elev.ft Prom.ft
3      3 Mt McKinley Denali     US   6194   6138   20320   20138
15    15         Mauna Kea     US   4205   4205   13796   13796
21    21         Mt Rainier    US   4393   4023   14411   13196
> |
```

```
83   M.ft
84
85   #subset all of the tallest mountains in the US
86   US.M<-datM[datM$Region=="US",]
87
88
89
```

```
> High.M<-datM[datM$Elev.ft>20000,]
> High.M
    Rank                      Name        Region Elev.m Prom.m Elev.ft Prom.ft
1      1               Mt Everest    Nepal Tibet   8848   8848   29028   29028
2      2                Aconcagua      Argentina   6962   6962   22841   22841
3      3        Mt McKinley Denali           US   6194   6138   20320   20138
14    14              Nanga Parbat       Pakistan   8125   4608   26657   15118
16    16 Jengish Chokusu ex Pik Pobedy Kyrgyzstan China   7439   4148   24406   13609
17    17               Chimborazo        Ecuador   6267   4122   20561   13523
19    19             Namcha Barwa          China   7782   4106   25531   13471
22    22                       K2 Pakistan China   8611   4017   28251   13179
> |
```

```
92
93
94
95
96   #subset by mountains above 20,000 ft
97   High.M<-datM[datM$Elev.ft>20000,]
98   High.M
99
100
101
102
103
```

# Missing data

- NA indicates that the data is missing in R

- If there are blank cells in a data file R will automatically fill them in with NA

- You can also designate what marks an NA if it differs in a data file:

# Errors

- Error messages look intimidating at first in R, but they are actually very useful

- Some kinds of examples:
  - Trying to do something where vectors are different lengths

```
> High.M$Elev.ft-US.M$Prom.ft
[1]  8890  9045  7124  6519 10610  7365  5393 14455
Warning message:
In High.M$Elev.ft - US.M$Prom.ft :
  longer object length is not a multiple of shorter object length
> |
```

```
93
94    #look at difference between prominance and elevation
95    High.M$Elev.ft-US.M$Prom.ft
96
97
98
99
```

  - Referring to names incorrectly (capitalization counts!)

```
> mean(High.M$elev.ft)
[1] NA
Warning message:
In mean.default(High.M$elev.ft) :
  argument is not numeric or logical: returning NA
> |
```

- You will get this message when you close out Rstudio.

- The answer is ALWAYS DON'T SAVE!!!!!!!!!!!!!!!!!!!!!!!!!!!
  - I cannot repeat this enough

# In fact:

# Why??

- This saves everything you ran in your previous session and loads it as if you just ran it

- It is easy to loose control of the variables and data in your workspace. You may have things you tested out and didn't work still in the environment.

- It is the easiest way to lead to confusion and problems with reproducing your analyses

- You should be writing nice, clean scripts that can be rerun upon opening. If something works, it's not magic. You wrote your script correctly.