

# Cluster analysis and its application in geochemistry -- A workshop at Goldschmidt 2020

Shuang Zhang (Carnegie Earth and Planets Laboratory)

## Step 1

### 1.1 Download R **first** from <https://cran.r-project.org/>

Note: if you already have R installed on your computer, just ignore this step. However, if you are not able to run the code of this workshop, your R might be very old. Consider updating R. This is a link showing you how to do the update (<https://www.linkedin.com/pulse/3-methods-update-r-rstudio-windows-mac-woratana-ngarmtrakulchol>).

Choose your system, download and install the latest R version (Linux needs to select the right distribution). My R version is 3.6.1. Therefore, if you choose to download and install the latest version, it will be newer than my version but it should work (let us know if it doesn't work). The packages you install afterwards should be compatible with your own R version.

For windows, click base and download the latest version as follows:

**R for Windows**

Subdirectories:

<a href="#">base</a>	Binaries for base distribution. This is what you want to <a href="#">install R for the first time</a> .
<a href="#">contrib</a>	Binaries of contributed CRAN packages (for R $\geq$ 2.13.x; managed by Uwe Ligges). There is also information on <a href="#">third party software</a> available for CRAN Windows services and corresponding environment and make variables.
<a href="#">old contrib</a>	Binaries of contributed CRAN packages for outdated versions of R (for R < 2.13.x; managed by Uwe Ligges).
<a href="#">Rtools</a>	Tools to build R and R packages. This is what you want to build your own packages on Windows, or to build R itself.

Please do not submit binaries to CRAN. Package developers might want to contact Uwe Ligges directly in case of questions / suggestions related to Windows binaries.

You may also want to read the [R FAQ](#) and [R for Windows FAQ](#).

Note: CRAN does some checks on these binaries for viruses, but cannot give guarantees. Use the normal precautions with downloaded executables.

**R-4.0.1 for Windows (32/64 bit)**

[Download R 4.0.1 for Windows](#) (84 megabytes, 32/64 bit)

[Installation and other instructions](#)

[New features in this version](#)

If you want to double-check that the package you have downloaded matches the package distributed by CRAN, you can compare the [md5sum](#) of the .exe to the [fingerprint](#) on the master server. You will need a version of md5sum for windows: both [graphical](#) and [command line versions](#) are available.

**Frequently asked questions**

- [Does R run under my version of Windows?](#)
- [How do I update packages in my previous version of R?](#)
- [Should I run 32-bit or 64-bit R?](#)

## 1.2 Download RStudio from <https://rstudio.com/products/rstudio/download/#download>

Still, choose the right system. The download button will most likely be shown on the website when you click this link, as this website can automatically detect your system and make the recommendation. If not, then just select the right version, download and install.

### RStudio Desktop 1.3.959 - [Release Notes](#)

1. Install R. RStudio requires [R 3.0.1+](#).
2. Download RStudio Desktop. Recommended for your system:



Requires macOS 10.13+ (64-bit)



## All Installers

Linux users may need to [import RStudio's public code-signing key](#) prior to installation, depending on the operating system's security policy.

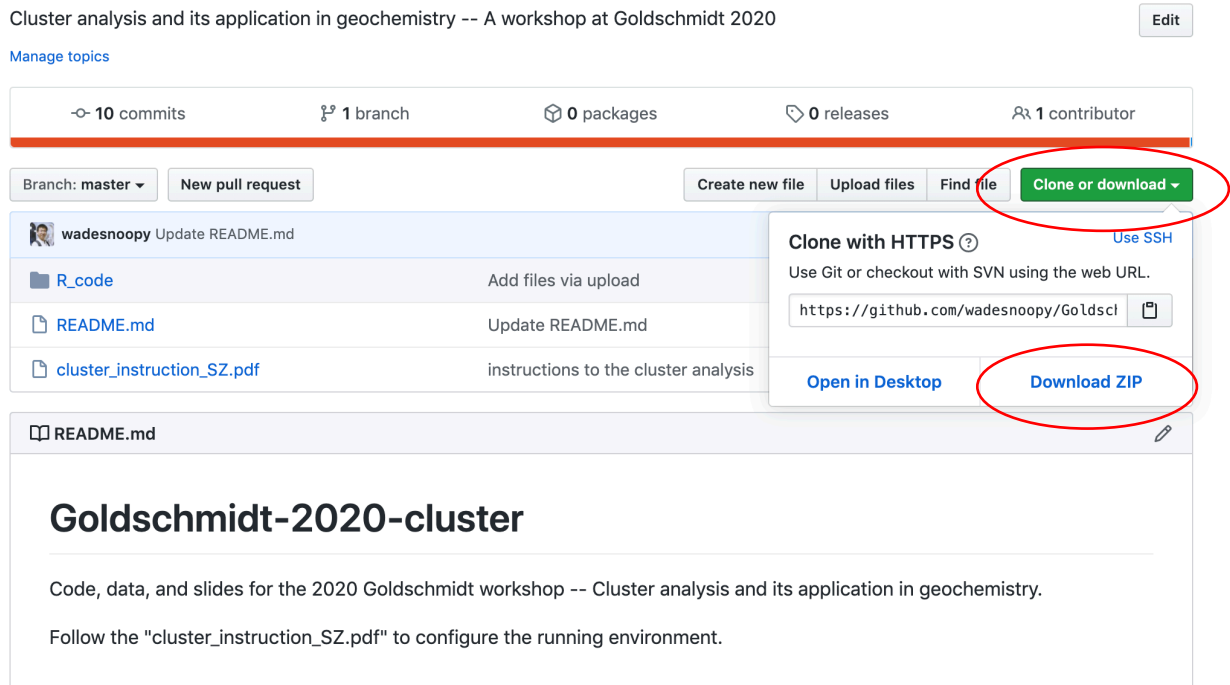
RStudio requires a 64-bit operating system. If you are on a 32 bit system, you can use an [older version of RStudio](#).

OS	Download	Size	SHA-256
Windows 10/8/7	<a href="#">RStudio-1.3.959.exe</a>	171.41 MB	3d493ae5
macOS 10.13+	<a href="#">RStudio-1.3.959.dmg</a>	148.57 MB	7c5b695d
Ubuntu 16	<a href="#">rstudio-1.3.959-amd64.deb</a>	124.57 MB	c2931495
Ubuntu 18/Debian 10	<a href="#">rstudio-1.3.959-amd64.deb</a>	126.11 MB	411ab500
Fedora 19/Red Hat 7	<a href="#">rstudio-1.3.959-x86_64.rpm</a>	146.24 MB	a144e4e6

Open RStudio and make sure everything is working so far. By default, RStudio will automatically detect the R installed in your system and use that R version.

**Step 2:** Download the ZIP file (more straightforward way) that contains the code, data and slides from <https://github.com/wadesnoopy/Goldschmidt-2020-cluster>

Updates (if any) will be uploaded to this repository. To have the latest instruction and material, download a fresh version again on Sunday (June 21, 2020). There might be some tiny changes before June 21.



**Step 3:** Put the ZIP file in your working folder and unzip it. The unzipped folder name will be “Goldschmidt-2020-cluster-master”

**Step 4:** Navigate into the folder of “Goldschmidt-2020-cluster-master”, then enter the folder “R\_code”. In this folder, you will see four files.

- “pyrite\_samples.csv” is the data
- “Goldschmidt-2020-cluster.R” is the code
- “Goldschmidt-2020-cluster.Rmd” is a markdown file that generates the static “Goldschmidt-2020-cluster.html” file.
- “Goldschmidt-2020-cluster.html” is the static HTML file

Open “Goldschmidt-2020-cluster.html” to see the code output.

**Step 5:** Run the code.

“**Goldschmidt-2020-cluster.R**” is the code that we will run. Open it in RStudio. Two ways to run the code.

1) You can run the whole code by first setting the path of working folder to your current “R\_code” folder and then clicking the “Source” button. The packages that are needed in this exercise should be automatically installed if they are not on your computer, based on the code chunk below at the top of the whole code.

```
9 # list the packages that we need
0 packages_need <- c("readr", # read in csv file
1                     "magrittr", # pipe %>%
2                     "dplyr", # manipulate data$
3                     "tidyr", # tidy data
4                     "DataExplorer", # check missing data
5                     "ggplot2", # plotting
6                     "GGally", # pairplot
7                     "corrplot", # correlation plot
8                     "plotly", # interactive 3D plot
9                     "FactoMineR", # do PCA analysis
0                     "factoextra", # facilitate the plotting of PCA outcome and the cluster analysis
1                     "fpc", # clusterboot to test stability
2                     "ClusterR" # for external validation
3 )
4
5
6 # What the code below does is to check if your R has these packages installed (the require() function
, install them (the install.packages() function), and then load the packages using the library() function
7
8 for (package in packages_need){
9   if (!require(package, character.only = TRUE)){
0     install.packages(package)
1   }
2   library(package, character.only = TRUE)
3 }
4
```

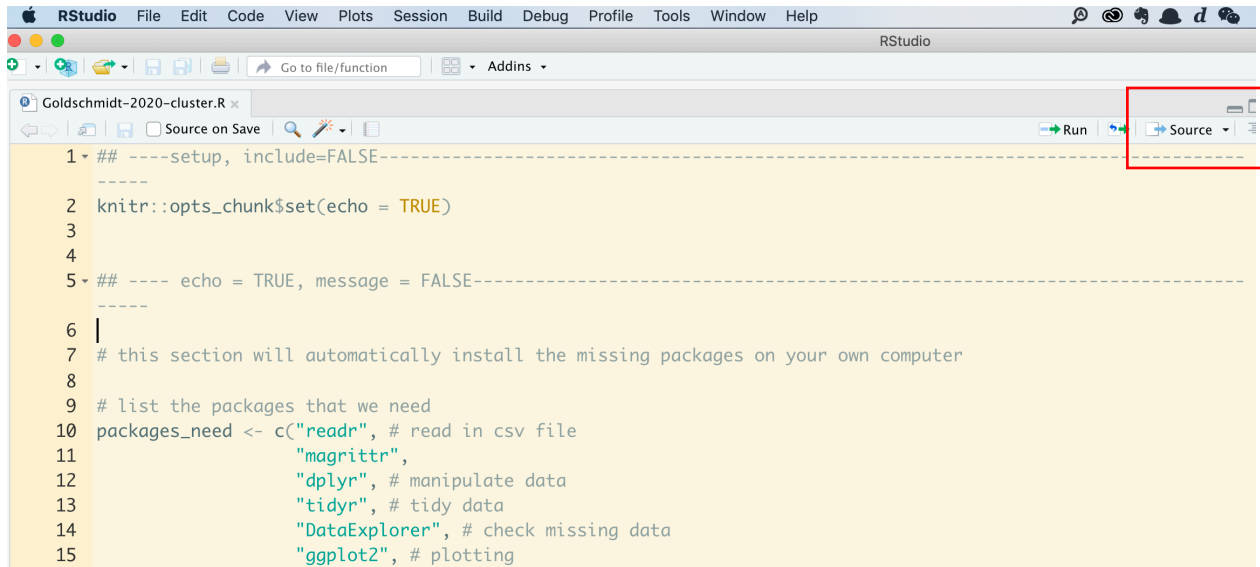
If there are errors occurring during the installation of any of the packages, you need to install the unsuccessful one yourself. Use `install.packages("the package name")` in the RStudio console to install the unsuccessful one, and follow the prompt.

Note that for windows, you also need to use the forward slash “/” instead of the backslash “\” when setting up your own path (shown below)

```
## ---- echo = TRUE, message = FALSE-----
----
# Navigate to your the R_code folder downloaded from github or from this workshop

# change "/Users/shzhang/Documents/Research/Meeting/2020-Goldschmidt/R_code" to your own local path of the
R_code folder downloaded from github or from this workshop
|
setwd("/Users/shzhang/Documents/Research/Meeting/2020-Goldschmidt/R_code")
```

FYI, The whole code can be finished within 2 minutes on my macbook pro (2019) with 2.3 GHz Intel Core i9. I would guess in general, it can be finished within ~8 minutes using a normal laptop. But if you run the code chunk by chunk, time will not be a problem as each chunk runs relatively fast.

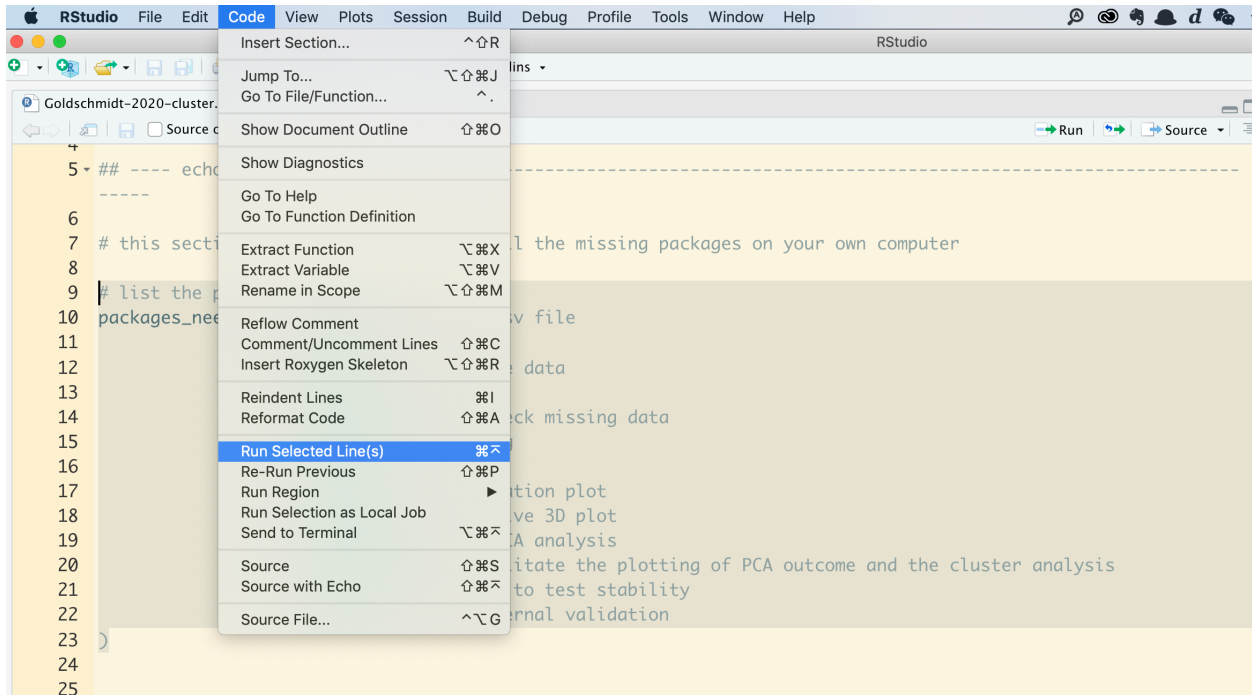


The screenshot shows the RStudio interface. The top menu bar includes File, Edit, Code, View, Plots, Session, Build, Debug, Profile, Tools, Window, and Help. Below the menu bar is a toolbar with icons for file operations and a 'Go to file/function' search bar. The main editor window displays a code chunk titled 'Goldschmidt-2020-cluster.R'. The code within the chunk is as follows:

```
1 ## ----setup, include=FALSE-----
2 knitr::opts_chunk$set(echo = TRUE)
3
4
5 ## ---- echo = TRUE, message = FALSE-----
6 |
7 # this section will automatically install the missing packages on your own computer
8
9 # list the packages that we need
10 packages_need <- c("readr", # read in csv file
11                   "magrittr",
12                   "dplyr", # manipulate data
13                   "tidyr", # tidy data
14                   "DataExplorer", # check missing data
15                   "ggplot2", # plotting
```

On the right side of the code editor, there is a vertical toolbar. The 'Run' button (a green play icon) and the 'Source' button (a blue play icon) are highlighted with a red rectangular box.

2) You can run the code line by line or chunk by chunk from top to bottom. Select the code chunk, then click the “Code” button on the top panel, and click “Run Selected Line(s)”. Tip: to save some time, you can remember the shortcut of “Run Selected Line(s)”. See below.

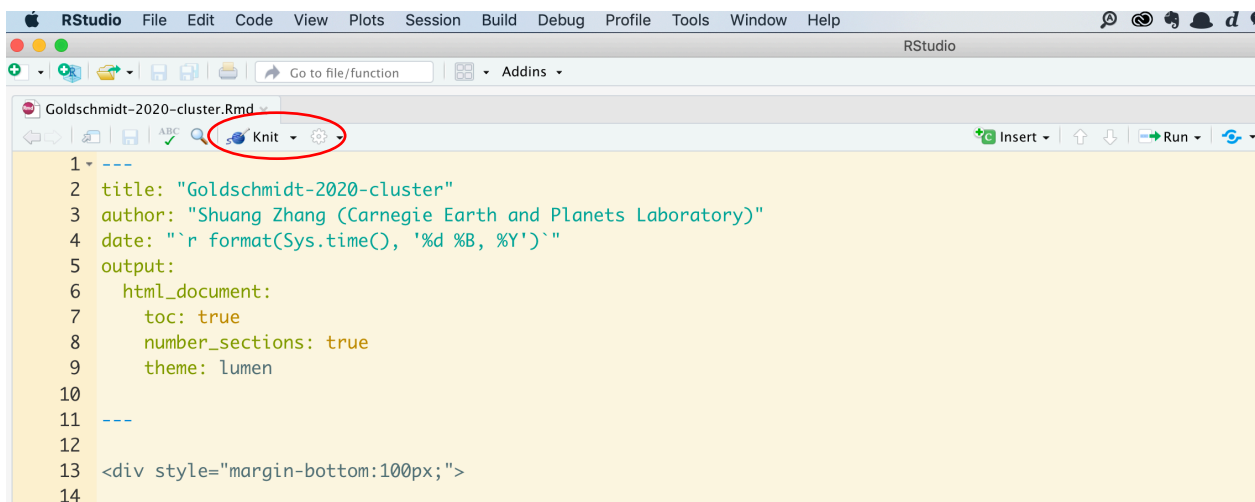


**Step 6:** Bonus step and optional (you can continue if you like)

I “Knit” the code and the code output to the “**Goldschmidt-2020-cluster.html**” file, which is a good way to share the code output among people. This file is generated by the “**Goldschmidt-2020-cluster.Rmd**” file, which is a markdown file and has the same code content with “**Goldschmidt-2020-cluster.R**”.

If you want to generate the “**Goldschmidt-2020-cluster.html**” yourself, you need to open the “**Goldschmidt-2020-cluster.Rmd**” file in RStudio, and then click “Knit”. Note that RStudio will notify you that you need to install “Knit” first if “Knit” is not on your computer. Just install it. And then click “Knit”.

FYI, The whole code can be finished within 2 minutes on my macbook pro (2019) with 2.3 GHz Intel Core i9. I would guess in general, it can be finished within ~8 minutes using a normal laptop.



```
1 ---
2 title: "Goldschmidt-2020-cluster"
3 author: "Shuang Zhang (Carnegie Earth and Planets Laboratory)"
4 date: "`r format(Sys.time(), '%d %B, %Y')`"
5 output:
6   html_document:
7     toc: true
8     number_sections: true
9     theme: lumen
10
11 ---
12
13 <div style="margin-bottom:100px;">
14
```