

# More Docker, Git(hub), and the Unix Shell

## MAS Departmental Disclaimers:

For students trying to take or audit from outside the MAS program.

Taking or auditing 400 courses is simply not permitted because this is a self-supporting program. Sorry, unfortunately, you will NOT be able to take any of the 400 level Stats courses.

There are NO exceptions that can be made by the department. These classes were designed specifically for students who applied directly to the program.

The students of this program are also not allowed to audit or enroll in classes outside of the program as it was created for working professionals.

If you would like to apply for the program, you are welcome to do so: [https://  
master.stat.ucla.edu/admissions/](https://master.stat.ucla.edu/admissions/)

Information is found here: <https://master.stat.ucla.edu/> And here: [https://master.stat.ucla.edu/  
faq/](https://master.stat.ucla.edu/faq/)



A reminder...

class communications

Slack Channel  
uclastat418-class

last time we saw an introduction to docker



today we'll make use of it while also hearing a bit more about what it is.

throughout the rest of course (in most weeks) we will continue to make use of docker environments to incrementally learn about its features

Docker: what is it?

a software to build and run containers with a container being a single instance of an image

an image is a unit of software that allows one to package up code and all its dependencies

## Docker: why is it important?

different software development stacks require different environments;  
especially dependent on task

classical solutions entail virtual environments (anaconda, python  
virtualenv) or virtual machines (vagrant, vmware)

## Docker: so why do data science with it?

it is platform agnostic; as long as a machine has docker we can use the same environment (perfect in a classroom setting where we all have different machines; perfect in a work setting where we might have different machines and certainly different versions)

can share version-controlled environments akin to GitHub for docker images

containers are ephemeral; if you mess up you can start over easily

build on top of other peoples work

ready-made workflow for deployment through Kubernetes or similiar

Docker: running Rstudio

run the following

```
docker run -e PASSWORD='class' --rm -p 8787:8787 rocker/rstudio
```

and then visit in any browser

<http://localhost:8787> or <http://<your-ip-address>:8787>

finally login with the following credentials

Username: rstudio      Password: <what you used above>

Docker: running Rstudio

to exit



We will see another way to exit in a moment

## Docker: running Rstudio with mounting local volume (connect to your computer)

cd to desired directory, then use for Mac

```
docker run -d --rm -e PASSWORD=class -p 8787:8787 -v /Users/nlangholz/Documents/:/home/rstudio/Documents rocker/rstudio
```

chdir to desired directory in Windows command line

```
docker run -d --rm -e PASSWORD=class -p 8787:8787 -v %cd%:/home/rstudio/Documents  
rocker/rstudio
```

cd to desired directory in Windows powershell

```
docker run -d --rm -e PASSWORD=class -p 8787:8787 -v ${PWD}:/home/rstudio/Documents  
rocker/rstudio
```

Again visit in your browser port 8787 either at local host or <your-ip-address> and login with your credentials

## Docker: a few other commands

Now to exit we will need to take a look at what containers are running to find the container id or name

```
docker container ls
```

```
docker kill <container>
```

And finally, to see what docker images we have available

```
docker images
```

## Docker: a few other commands

And finally, to see what docker images exist on our machine

```
docker images
```

If you'd like to remove any of these images

```
docker rmi <image id>
```

# Git

## Git: what is it?

a version control system to track the history of changes as people and teams collaborate on code projects together

Git manages the evolution of a set of files (or in our case a data analysis or model development) in a repository

this provides a transparent history of changes, who made them and how each contribute all while reducing redundancy and unnecessary work

## GIT

< < PREV RANDOM NEXT > >

THIS IS GIT. IT TRACKS COLLABORATIVE WORK  
ON PROJECTS THROUGH A BEAUTIFUL  
DISTRIBUTED GRAPH THEORY TREE MODEL.

COOL. HOW DO WE USE IT?

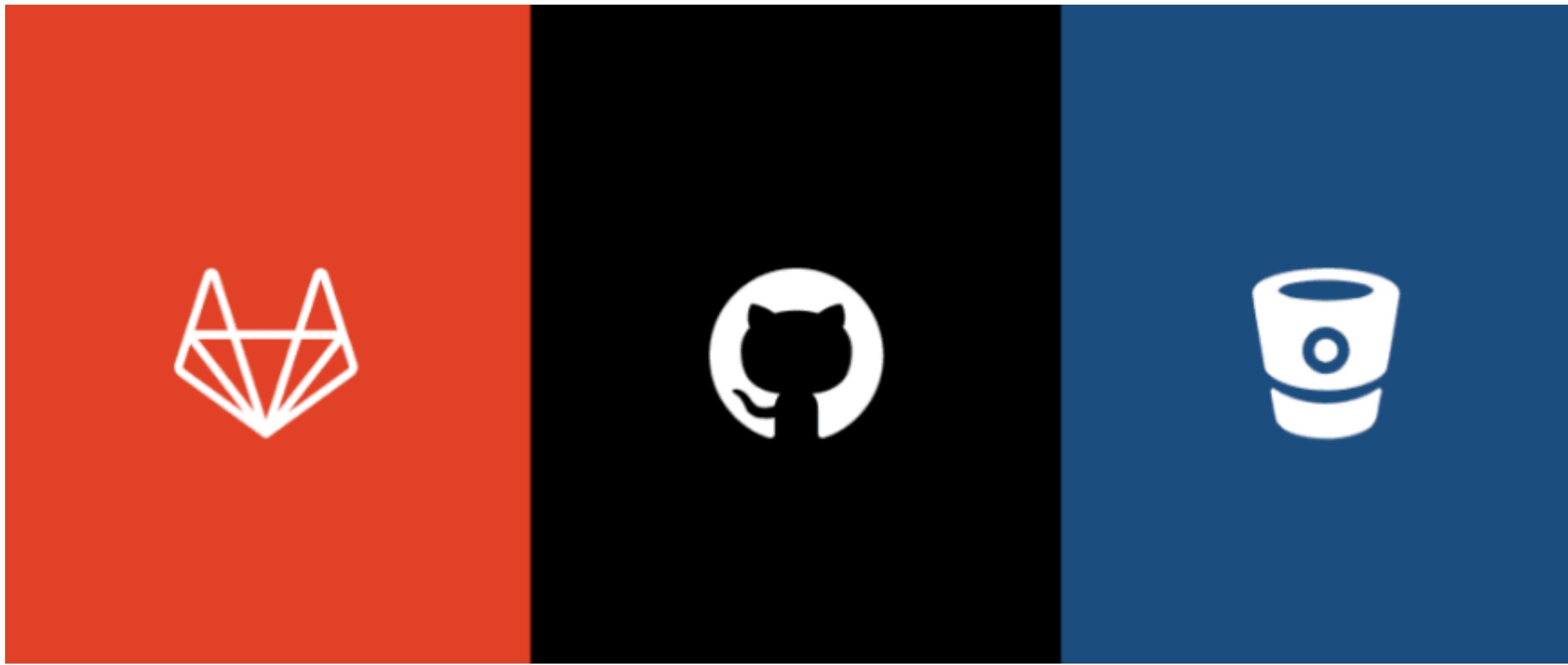
NO IDEA. JUST MEMORIZIZE THESE SHELL  
COMMANDS AND TYPE THEM TO SYNC UP.  
IF YOU GET ERRORS, SAVE YOUR WORK  
ELSEWHERE, DELETE THE PROJECT,  
AND DOWNLOAD A FRESH COPY.



< < PREV RANDOM NEXT > >

PERMANENT LINK TO THIS COMIC: [HTTPS://XKCD.COM/1597/](https://xkcd.com/1597/)  
IMAGE URL (FOR HOTLINKING/EMBEDDING): [HTTPS://IMGS.XKCD.COM/COMICS/GIT.PNG](https://imgs.xkcd.com/comics/git.png)

# Repository management services



# Repository management services



**GitHub**

Code Collaboration & Version Control

Follow

Stacks

**32.5K**

I Use This

Fans    Jobs    Votes  
**28.7K** **3.75K** **10K**



**Bitbucket**

Code Collaboration & Version Control

Follow

Stacks

**9.82K**

I Use This

Fans    Jobs    Votes  
**8.44K** **438** **2.77K**



**GitLab**

Code Collaboration & Version Control

Follow

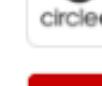
Stacks

**8.35K**

I Use This

Fans    Jobs    Votes  
**7.57K** **357** **2K**

# Repository management services

What companies use GitHub?	What companies use Bitbucket?	What companies use GitLab?
<i>5837 companies on StackShare use GitHub</i>	<i>2270 companies on StackShare use Bitbucket</i>	<i>1584 companies on StackShare use GitLab</i>
 Airbnb	 PayPal	 Alibaba.com
 Netflix	 Salesforce	 Coderus
 Medium	 Zillow	 Webedia
 Instacart	 Starbucks	 Ticketmaster
 reddit	 CircleCI	 WebbyLab
 Lyft	 Tesla Motors	 Edify
 StackShare	 Bitbucket	 Electronic Arts
 Shopify	 Pandora	 Freelancer.com
 9GAG	 Sellsuki	 Citrix
 Asana	 Movielala	 WILD



We're going to use GitHub

<http://github.com>

Make an account if you don't have one

Let it be known...



Microsoft | Stories Our Company ▾ Our Products ▾ Blogs & Communities ▾ Press Tools ▾

## Microsoft to acquire GitHub for \$7.5 billion

June 4, 2018 | Microsoft News Center

[f](#) [t](#) [in](#)

*Acquisition will empower developers, accelerate GitHub's growth and advance Microsoft services with new audiences*

<https://www.theverge.com/2018/6/18/17474284/microsoft-github-acquisition-developer-reaction>



**Mike Coutermash**  
@mscccc

Just googled "undo last commit but keep changes"

And I work for GitHub.

This stuff is hard!

11:10 PM - 18 Jan 2018

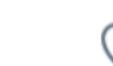
321 Retweets 1,007 Likes



39



321



1.0K



It's not all fun and games. Actually it's quite hard. Initially just working on solo repositories becomes fine but once collaborating you'll start to run into merge conflicts (yikes)

Keep going!!

## GitHub: what tool to use?

Use what makes you feel comfortable and efficient.

“

*“I certainly know Git very well, and honestly think I’m far faster and more efficient in a Git GUI than I could possibly be in the command line – and I’m certainly not slow in the CLI.”*

”

—Dan Clarke, blogger and co-organizer of .Net Oxford group

“

*“I sometimes encounter people who feel it’s “better” to use command line Git, but for very ill-defined reasons. These people may feel like they should work in the shell, even if it leads to Git-avoidance, frequent mistakes, or limiting themselves to a small set of ~3 Git commands. This is counterproductive.”*

”

—Jenny Bryan, Software Engineer at RStudio & Adjunct Professor at the University of British Columbia

Both quotes take from <https://blog.axosoft.com/git-gui-vs-cli/>

## GitHub GUIs



There are a wide variety of choices for GUIs if you are interested.

## GitHub: some (basic) guidelines

Personally, I start all my repositories remotely (on [GitHub.com](https://GitHub.com)). I find it much easier to then pull to my local machine and to then make my first commit

This includes forking others' repos and then cloning to local

Create new branches to make edits (add new features) to the master code base

Commit early and often; typically don't push data (especially when large) and definitely not credentials

(these are all quite simple; you'll surely run into many issues and there are many more resources online that will help you, but don't be afraid to try some new things in this course to learn how it works)

## GitHub: our class repository

Let's create a repo for you to use and push hw to my master repo. cd to the directory where you want to

Fork the class repo on GitHub. Clone your https version of the repo.

Now let's go to docker and do this in a (hopefully) more familiar environment...get an Rstudio container running with a local volume mounted. (if its not still running at the moment)

...will walk through as a class.

A bit more on reproducibility...

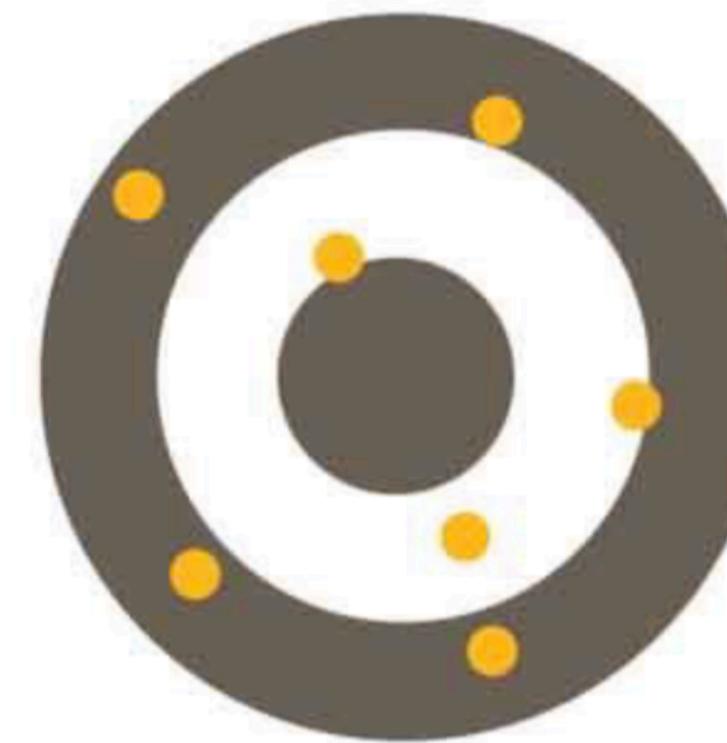
there is a bit of crisis taking place in regards to reproducibility as a whole in science fields. and this true of the data science field as well.

We've seen some of the tools to enable us to do this reproducible work, but what concepts should we be thinking of as we use them?

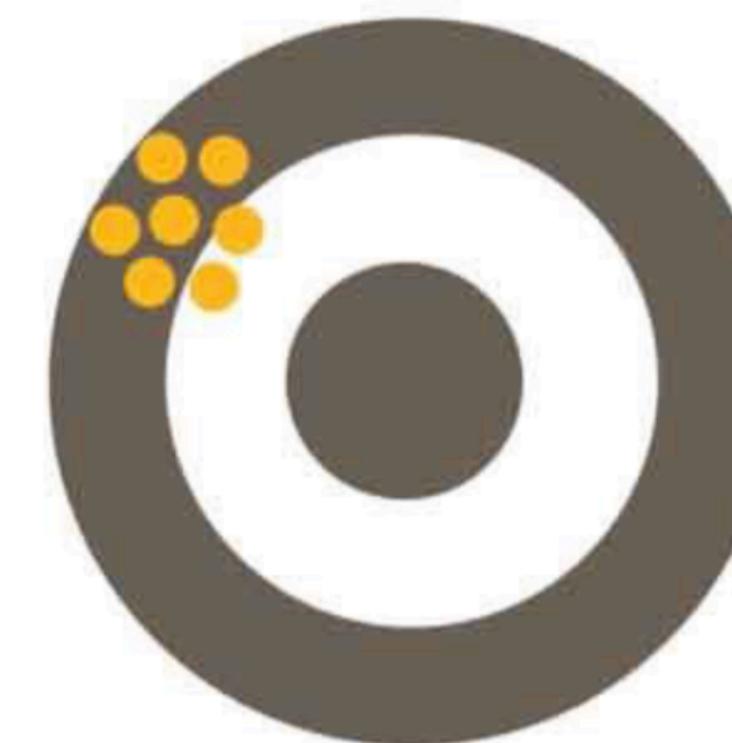
# Repeatability

or test-retest reliability is the variation in measurements taken by a single person or instrument on the same item, under the same conditions, and in a short period of time. A less-than-perfect test-retest reliability causes test-retest variability.

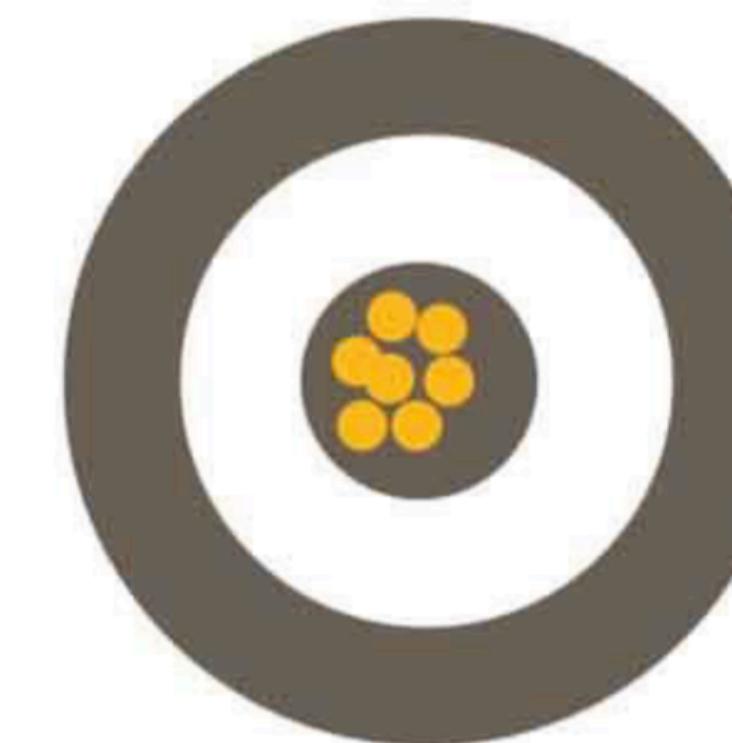
Your systems must be repeatable and reliable, the same query on the same data should return the same results. Otherwise, none of this applies.



Low repeatability, Low accuracy



High repeatability, Low accuracy



High Repeatability, High accuracy

## Reproducibility

an analysis is *reproducible* if there is a specific set of computational functions/analyses (usually specified in terms of code) that exactly reproduce all of the numbers in a published paper from raw data. It is now recognized that a critical component of the scientific process is that data analyses can be reproduced.

## Replicability

but just because a study is reproducible does not mean that it is *replicable*. Replicability is stronger than *reproducibility*. A study is only *replicable* if you perform the exact same experiment (at least) twice, collect data in the same way both times, perform the same data analysis, and arrive at the same conclusions.

# Stodden's Taxonomy of Reproducibility

## Computational

when detailed information is provided about code, software, hardware and implementation details.



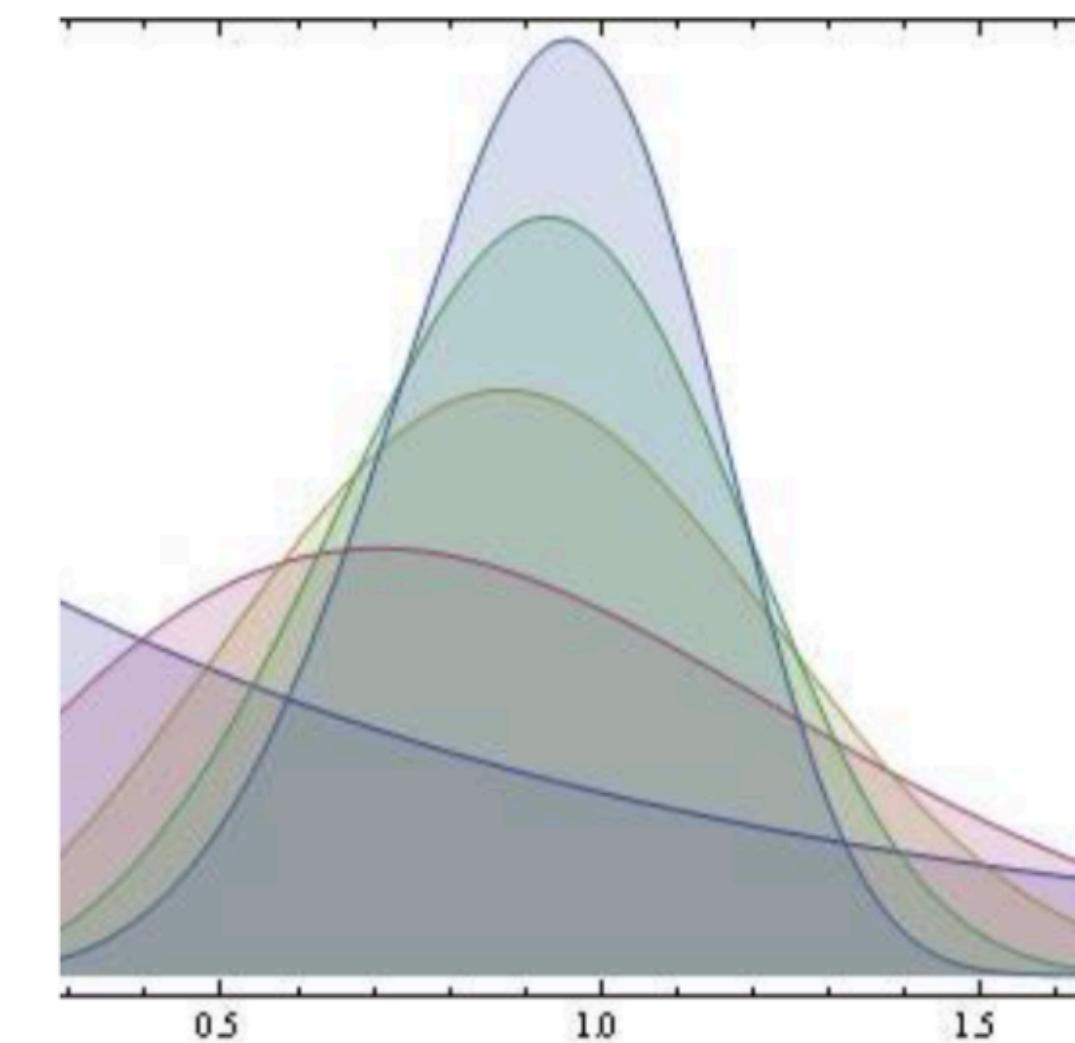
## Empirical

when detailed information is provided about code, software, hardware and implementation details.



## Statistical

when detailed information is provided about code, software, hardware and implementation details.



# Computational reproducibility should be...

## Reviewable

The descriptions of the research methods can be independently assessed and the results judged credible

## Auditable

sufficient records (including data and software) have been archived so that the research can be defended later if necessary or differences between independent confirmations resolved. The archive might be private, as with traditional laboratory notebooks.

## Replicable

tools are made available that would allow one to duplicate the results of the research, for example by running the authors' code to produce the plots shown in the publication.

## Confirmable

the main conclusions of the research, or the model in production can be attained independently without the use of software provided by the author.

## Simple rules for reproducible computational research

track results - whenever a result may be of potential interest, keep track of how it was produced. as a minimum, you should at least record sufficient details on programs, parameters, and manual procedures to allow yourself, in a year or so, to approximately reproduce the results.

script everything - whenever possible, rely on the execution of programs instead of manual procedures to modify data. If manual operations cannot be avoided, you should as a minimum note down which data files were modified or moved, and for what purpose.

store program versions - In order to exactly reproduce a given result, it may be necessary to use programs in the exact versions used originally. As a minimum, you should note the exact names and versions of the main programs you use.

## Simple rules for reproducible computational research

use version control - even the slightest change to a computer program can have large intended or unintended consequences. As a minimum, you should archive copies of your scripts from time to time, so that you keep a rough record of the various states the code has taken during development

store data & intermediate results - in principle, as long as the full process used to produce a given result is tracked, all intermediate data can also be regenerated. In practice, having easily accessible intermediate results may be of great value. As a minimum, archive any intermediate result files that are produced when running an analysis

set a random number seed- many analyses and predictions include some element of randomness, meaning the same program will typically give slightly different results every time it is executed. As a minimum, note which analysis steps involve randomness, so that a level of discrepancy can be anticipated when reproducing the results.

## Simple rules for reproducible computational research

store data viz inputs- from the time a figure is first generated to it being part of an analysis. It is critical to store the data and process that generated it. As a minimum, one should note which data formed the basis of a given plot and how this data could be reconstructed.

allow levels of analysis - in order to validate and fully understand the main result, it is often useful to inspect the detailed values underlying the summaries. Make those fluid and explorable, as a minimum at least once generate, inspect, and validate the detailed values underlying the summaries.

generate text programmatically- there is nothing quite as embarrassing as having a disagreement between your analytical narrative in a writeup and having the tables and figures disagree. Connect them so that there is no chance of disagreement.

.

# Unix Shell

## The shell: why is it important??

to introduce the shell means to have a discussion about the structure of the computer, operating systems, file systems and history

the shell offers programmatic access to a computer's underlying parts, providing the ability to “do” data analysis on directories, on processes, and on their networks

as there are many flavors of the shell, it is the first time we come to a decision about choosing ‘tools’, about evaluating which shell is best for you and the shifting terrain of software development

## The shell: why is it important?

as a practical matter, shell tools are an indispensable part for my own data science practice, data ‘cleaning’, data processing, exploratory analysis, automation; by design they let us deal with data on a scale that can be difficult from with R or Python

it also provides a low level link to other data platforms, data sources, and remote environments

## The shell: how we'll learn

everyone needs access to the shell. everyone with a Mac and Linux already have one. those with Windows machines need something else..

there are instructions on the site for Windows users to get a version of it, but..

as we all have docker installed we can use an image that provides us all with the same environment for learning and exploration. not all functionality works in the image but overall will be useful for us (also we've seem a version in our docker Rstudio environment)

# Operating Systems

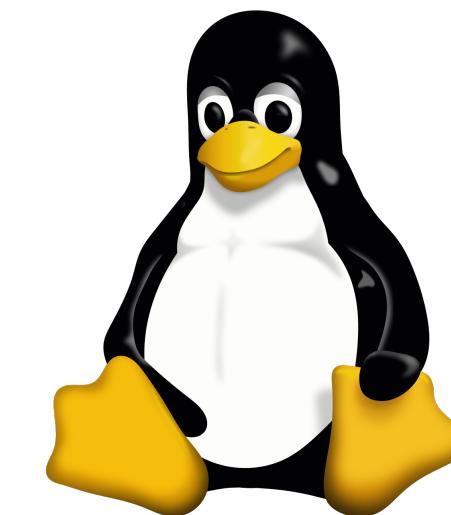
most devices that contain a computer of some kind will have an OS, they tend to emerge when the appliance will have to deal with new applications, complex user-input and possibly changing requirements on its function



# Operating Systems

An operating system is a piece of software that organizes and controls hardware and other software so your computer behaves in a flexible but predictable way

maybe obviously to us all, Window, MacOS and Linux are those most commonly used for home computers. iOS and Android are most common on mobile devices



## A history

In 1964, Bell Labs partnered with MIT and GE to create Multics  
(Multiplexed Information and Computing Service)

“Such systems must run continuously and reliably 7 days a week, 24 hours a day in a way similar to telephone or power systems, and must be capable of meeting wide service demand from multiple man-machine interaction to the sequential processing of absentee-user jobs, from the use of the system with dedicated languages and subsystems to the programming of the stem itself”

## A history

Bell Labs pulled out of the Multics project in 1969, a group of researchers at Bell Labs started work on Unica (uniplexed information and computing system) because initially it could only support one user; as the system matured it was renamed Unix, which isn't an acronym for anything

Ritchie simply says that Unix is a ‘somewhat treacherous pun on Multics’



## The Unix filesystem

Multics brought about the first notion of a hierarchical file system; files were arranged in a tree structure allowing users to have control of them on areas

Unix began (more or less) as a file system and then an interactive shell emerged to let you examine its contents and perform basic operations

## The kernel and the shell

the Unix kernel is the part of the operating system that carries out basic functions like accessing files, handing communication, and others

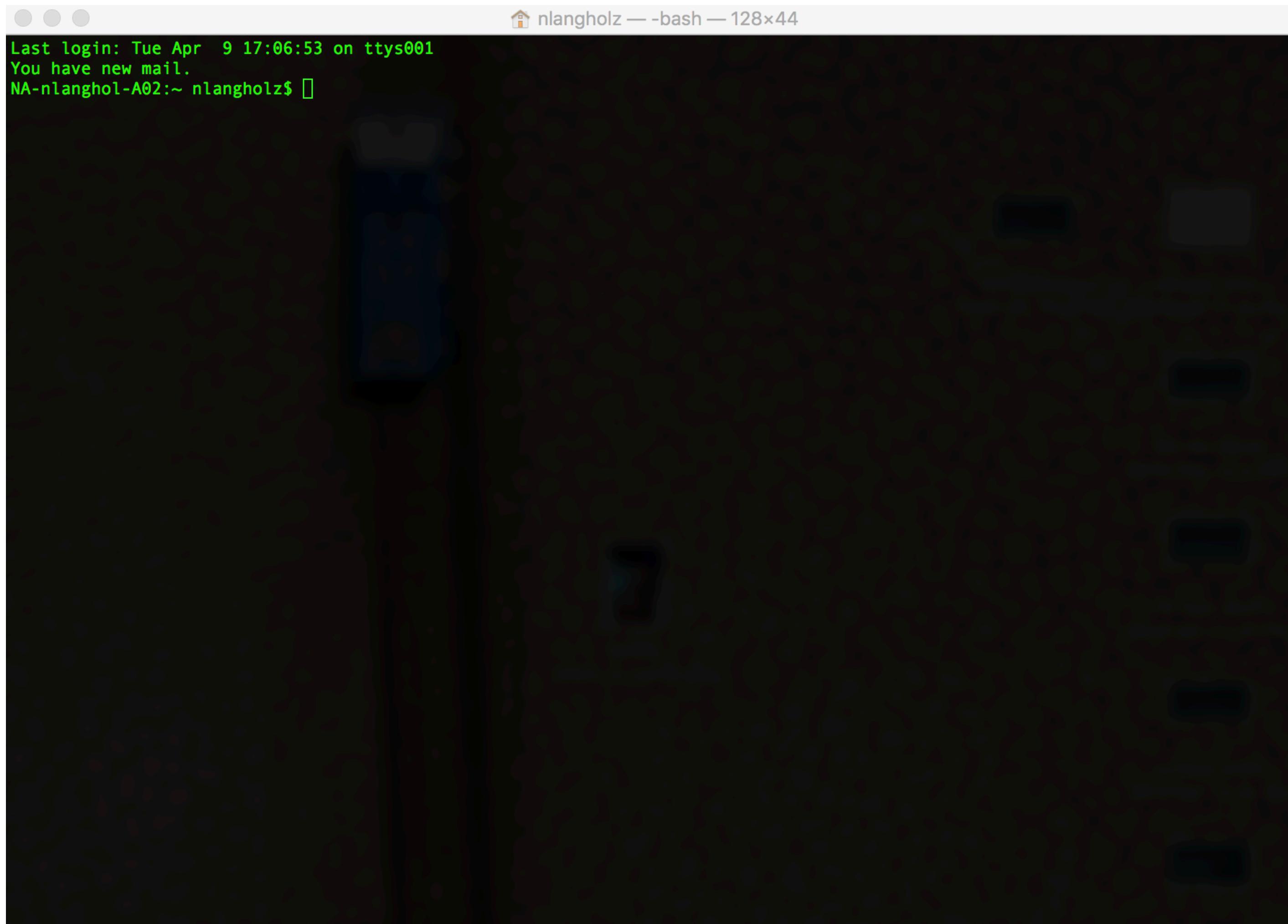
the Unix shell is a user interface to the kernel (keep in mind that Unix was designed for computer scientists and the interface is not optimized for novices)

## Unix shells

A shell is a type of program called an interpreter, think of it as a text-based interface to the kernel

It operates in simple loop: it accepts a command, interprets it, executes the command and waits for another

The shell displays a prompt to tell you that it is ready to accept a command



```
Last login: Tue Apr  9 17:06:53 on ttys001
You have new mail.
NA-nlanghol-A02:~ nlangholz$
```

Not much going on but  
makes us feel like we're in  
the matrix

## Unix shells

The shell is itself a program the the Unix operating system runs for you (a program is referred to as a process when its running)

The kernel manages many processes at once, many of which are the result of user commands (others provide services that keep the computer running)

Some commands are built into the shell, others have been added by users

Either way, the shell waits until the command is executed

```

Processes: 406 total, 2 running, 404 sleeping, 1762 threads
Load Avg: 1.59, 1.85, 1.99 CPU usage: 3.12% user, 5.52% sys, 91.34% idle
SharedLibs: 220M resident, 59M data, 48M linkedit.
MemRegions: 91775 total, 4908M resident, 116M private, 3768M shared. PhysMem: 16G used (2527M wired), 313M unused.
VM: 4303G vsize, 1110M framework vsize, 214622(0) swapins, 285256(0) swapouts.
Networks: packets: 2984993/22G in, 3071920/22G out. Disks: 1046817/18G read, 735692/14G written.

PID  COMMAND      %CPU TIME    #TH  #WQ  #PORT MEM    PURG   CMPRS  PGRP  PPID  STATE   BOOSTS  %CPU_ME
44924 screencaptur 0.7 00:00.19 2   1   49   2628K 36K   0B    604   604   sleeping *0[1]  0.00000
44887 top          3.3 00:01.25 1/1  0   24   5696K 0B    0B    44887  44872  running *0[1]  0.00000
44872 bash         0.0 00:00.01 1   0   19   812K  0B    0B    44872  44871  sleeping *0[1]  0.00000
44871 login        0.0 00:00.02 2   1   29   2292K 0B    0B    44871  44825  sleeping *0[9]   0.00000
44825 Terminal     1.1 00:03.12 8   3   259  65M   396K  0B    44825  1     sleeping *0[20] 0.00000
44602 kcm          0.0 00:00.04 3   3   22   2192K 0B    0B    44602  1     sleeping *0[1]  0.00000
44511 TeamViewer_D 0.1 00:00.57 7   1   92   6052K 0B    0B    44511  1     sleeping *0[1]  0.06399
44507 TeamViewerHo 0.0 00:00.41 8   1   191  16M   8192B 0B    44507  1     sleeping *0[1]  0.00000
44471 mdworker     0.0 00:00.12 4   2   56   5032K 0B    0B    44471  1     sleeping *0[1]  0.00000
44157 com.apple.We 0.0 00:00.22 4   1   142  12M   0B    0B    44157  1     sleeping *0[8]   0.00000
44156 com.apple.We 0.0 00:02.29 5   1   162  129M  0B    0B    44156  1     sleeping *0[140] 0.00000
43783 com.apple.We 0.0 00:20.24 5   1   173  337M  52M   0B    43783  1     sleeping *0[1427] 0.00000
43486 MTLCompilerS 0.0 00:00.04 2   2   22   7500K 0B    0B    43486  1     sleeping 0[3]   0.00000
43277 com.apple.We 0.0 00:03.70 5   1   174  76M   47M   0B    43277  1     sleeping *0[1448] 0.00000
42905 mdworker     0.0 00:00.24 3   1   61   3380K 0B    0B    42905  1     sleeping *0[1]  0.00000
42868 mdworker     0.0 00:00.18 3   1   62   3404K 0B    0B    42868  1     sleeping *0[1]  0.00000
42866 mdworker     0.0 00:00.07 3   1   57   3236K 0B    0B    42866  1     sleeping *0[1]  0.00000
42865 mdworker     0.0 00:00.13 3   1   61   3388K 0B    0B    42865  1     sleeping *0[1]  0.00000
42819 com.apple.We 0.0 00:00.35 5   1   162  22M   0B    0B    42819  1     sleeping *0[204] 0.00000
42810 mdworker     0.0 00:00.09 3   1   57   3260K 0B    0B    42810  1     sleeping *0[1]  0.00000
42809 mdworker     0.0 00:00.09 3   1   61   3236K 0B    0B    42809  1     sleeping *0[1]  0.00000
42566 MTLCompilerS 0.0 00:00.10 2   2   22   10M   0B    0B    42566  1     sleeping 0[11]  0.00000
42563 com.apple.We 0.0 01:29.09 5   1   212  196M  51M   0B    42563  1     sleeping *0[48553] 0.00000
42525 com.apple.We 0.0 00:00.43 5   1   160  19M   0B    0B    42525  1     sleeping *0[240]  0.00000
42241 mdworker     0.0 00:00.20 3   1   61   3448K 0B    0B    42241  1     sleeping *0[1]  0.00000
42084 QuickLookSat 0.0 00:01.07 3  2   47   14M   844K  3344K  42084  1     sleeping 0[0]   0.00000
42082 quicklookd   0.0 00:00.82 4   1   98   2464K 96K   7668K  42082  1     sleeping 0[34]  0.00000
41849 mdworker     0.0 00:00.48 4   1   56   21M   0B    5604K  41849  1     sleeping *0[1]  0.00000
41848 mdworker     0.0 00:00.35 4   1   56   24M   0B    20K   41848  1     sleeping *0[1]  0.00000
41717 mdworker     0.0 00:00.08 3   1   42   3140K 0B    24K   41717  1     sleeping *0[1]  0.00000
41715 mdworker     0.0 00:04.49 3   1   55   7952K 0B    2960K  41715  1     sleeping *0[1]  0.00000
41714 mdworker     0.0 00:00.19 4   1   56   6884K 0B    11M   41714  1     sleeping *0[1]  0.00000
41669 mdworker     0.0 00:00.12 4   1   56   12M   0B    496K   41669  1     sleeping *0[1]  0.00000
41483 CoreServices 0.0 00:00.35 3   1   163  3044K 0B    0B    41483  1     sleeping *0[1]  0.00000
41213 com.apple.We 0.0 00:01.79 5   1   175  57M   0B    4184K  41213  1     sleeping *0[300] 0.00000
40369 com.apple.We 0.0 00:00.37 6   1   150  6776K 0B    7500K  40369  1     sleeping *0[103] 0.00000
38787 netbiosd    0.0 00:00.05 2   2   25   360K  0B    2076K  38787  1     sleeping *0[1]  0.00000
38525 KeychainSync 0.0 00:00.04 2   1   62   52K   0B    1412K  38525  1     sleeping 0[1]  0.00000
34239 VTDecoderXPC 0.0 00:00.57 2   1   46   12M   0B    836K   34239  1     sleeping 0[239] 0.00000
34238 com.apple.au 0.0 00:00.01 2   2   17   8192B 0B    840K   34238  1     sleeping 0[1]  0.00000
34237 MTLCompilerS 0.0 00:00.05 2   2   23   12K   0B    7660K  34237  1     sleeping 0[3]  0.00000
34156 com.apple.We 0.0 02:17.45 9   4   284  367M  32M   47M   34156  1     sleeping *0[7815] 0.00000

```

Result of the command  
top; this is a printout of all  
the processes running on  
your computer

# Operating Systems

## Process Management

Schedules jobs(formally referred to as processes) to be executed by the computer

## Memory and storage management

Allocate space required for each running process in main memory (RAM) or in some other temporary location if space is tight and supervise the storage of data onto disk

# Operating Systems

## Device Management

A program called a driver translates data (files from the filesystem) into signals that

## Application Interface

An API (application programming interface) let programmers use functions of the computer and the operating system without having to know *how* something is done

## User Interface

Finally, the operating system turns and looks at you; the UI is a program that defines how users interact with the computer; some are graphical (Windows is a GUI) and some are text-based (your Unix shell)

# Unix Shell(s)

There are, in fact,  
many different  
kinds of Unix  
shells

The table on the  
right lists a few of  
the most popular

KSH(1)	General Commands Manual	KSH(1)	fish(1)	fish	fish(1)
<b>NAME</b>	<code>ksh, ksh93 - KornShell, a command and programming language</code>		<b>NAME</b>	<code>fish - fish - the friendly interactive shell</code>	
<b>SYNOPSIS</b>	<code>ksh [ +abcefhikmnoprstuvwxyzBCDP ] [ -R file ] [ +o option ] ... [ - ] [ arg ... ] rksh [ +abcefhikmnoprstuvwxyzBCD ] [ -R file ] [ +o option ] ... [ - ] [ arg ... ]</code>		<b>fish - the friendly interactive shell</b>	<code>fish [-h] [-v] [-c command] [FILE [ARGUMENTS...]]</code>	
<b>DESCRIPTION</b>	<p><code>Ksh</code> is a command and programming language that executes commands read from a terminal or a file. <code>Rksh</code> is a restricted version of the command interpreter <code>ksh</code>; it is used to set up login names and execution environments whose capabilities are more controlled than those of the standard shell. <code>Rpksh</code> is a profile shell version of the command interpreter <code>ksh</code>; it is used to execute commands with the attributes specified by the user's profiles (see <code>pexec(1)</code>). See <code>Invocation</code> below for the meaning of arguments to the shell.</p>		<b>Description</b>	<p><code>fish</code> is a command-line shell written mainly with interactive use in mind. The full manual is available in HTML by using the help command from inside <code>fish</code>.</p>	
TCSH(1)	General Commands Manual	TCSH(1)	The following options are available:		
<b>NAME</b>	<code>tcsh - C shell with file name completion and command line editing</code>		<code>-c</code> or <code>--command=COMMANDS</code> evaluate the specified commands instead of reading from the commandline		
<b>SYNOPSIS</b>	<code>tcsh [-bcdefFimnqstvVxX] [-Dname[=value]] [arg ...] tcsh -l</code>		<code>-d</code> or <code>--debug-level=DEBUG_LEVEL</code> specify the verbosity level of <code>fish</code> . A higher number means higher verbosity. The default level is 1.		
<b>DESCRIPTION</b>	<p><code>tcsh</code> is an enhanced but completely compatible version of the Berkeley UNIX C shell, <code>csh(1)</code>. It is a command language interpreter usable both as an interactive login shell and a shell script command processor. It includes a command-line editor (see <code>The command-line editor</code>), programmable word completion (see <code>Completion and listing</code>), spelling correction (see <code>Spelling correction</code>), a history mechanism (see <code>History substitution</code>), job control (see <code>Jobs</code>) and a C-like syntax. The <code>NEW FEATURES</code> section describes major enhancements of <code>tcsh</code> over <code>csh(1)</code>. Throughout this manual, features of <code>tcsh</code> not found in most <code>csh(1)</code> implementations (specifically, the 4.4BSD <code>csh</code>) are labeled with `(+)', and features which are present in <code>csh(1)</code> but not usually documented are labeled with `(u)'.</p>				
BASH(1)	General Commands Manual	BASH(1)			
<b>NAME</b>	<code>bash - GNU Bourne-Again SHell</code>				
<b>SYNOPSIS</b>	<code>bash [options] [command_string   file]</code>				
<b>COPYRIGHT</b>	<p><code>Bash</code> is Copyright (C) 1989-2013 by the Free Software Foundation, Inc.</p>				
<b>DESCRIPTION</b>	<p><code>Bash</code> is an <code>sh</code>-compatible command language interpreter that executes commands read from the standard input or from a file. <code>Bash</code> also incorporates useful features from the <code>Korn</code> and <code>C</code> shells (<code>ksh</code> and <code>csh</code>).</p>				
	<p><code>Bash</code> is intended to be a conformant implementation of the Shell and Utilities portion of the IEEE POSIX specification (IEEE Standard 1003.1). <code>Bash</code> can be configured to be POSIX-conformant by default.</p>				
ZSH(1)	General Commands Manual	ZSH(1)			
<b>NAME</b>	<code>zsh - the Z shell</code>				
<b>OVERVIEW</b>	<p>Because <code>zsh</code> contains many features, the <code>zsh</code> manual has been split into a number of sections:</p>				
	<code>zsh</code> Zsh overview (this section) <code>zshroadmap</code> Informal introduction to the manual <code>zshmisc</code> Anything not fitting into the other sections <code>zshexpn</code> Zsh command and parameter expansion <code>zshparam</code> Zsh parameters <code>zshoptions</code> Zsh options <code>zshbuiltins</code> Zsh built-in functions <code>zshzle</code> Zsh command line editing <code>zshcompswid</code> Zsh completion widgets <code>zshcompsys</code> Zsh completion system <code>zshcomptcl</code> Zsh completion control <code>zshmodules</code> Zsh loadable modules <code>zshcalsys</code> Zsh built-in calendar functions				

## Why the choices?

A shell program was originally meant to take commands, interpret them and then execute some operation

Inevitably, one wants to collect a number of these operations into programs that execute compound tasks at the same time you want to make interaction on the command line as easy as possible (a history mechanism, editing capabilities and so on)

The original Bourne shell is ideal for programming; the C-shell and its variants are good for interactive use; the Korn shell is a combination of both



Steve Bourne, created of sh, 2005 (source wikipedia)

And while we're at it

unix itself comes in different flavors; the 1980s saw an incredible proliferation of Unix versions, somewhere around 100 (System V, AIX, Berkely BSD, SunOS, Linux, ... )

vendors provided (diverging) version of Unix, optimized for their own computer architectures and supporting different features

despite the diversity it was still easier to “port” applications between versions of Unix than it was between different proprietary OS

## A few common commands

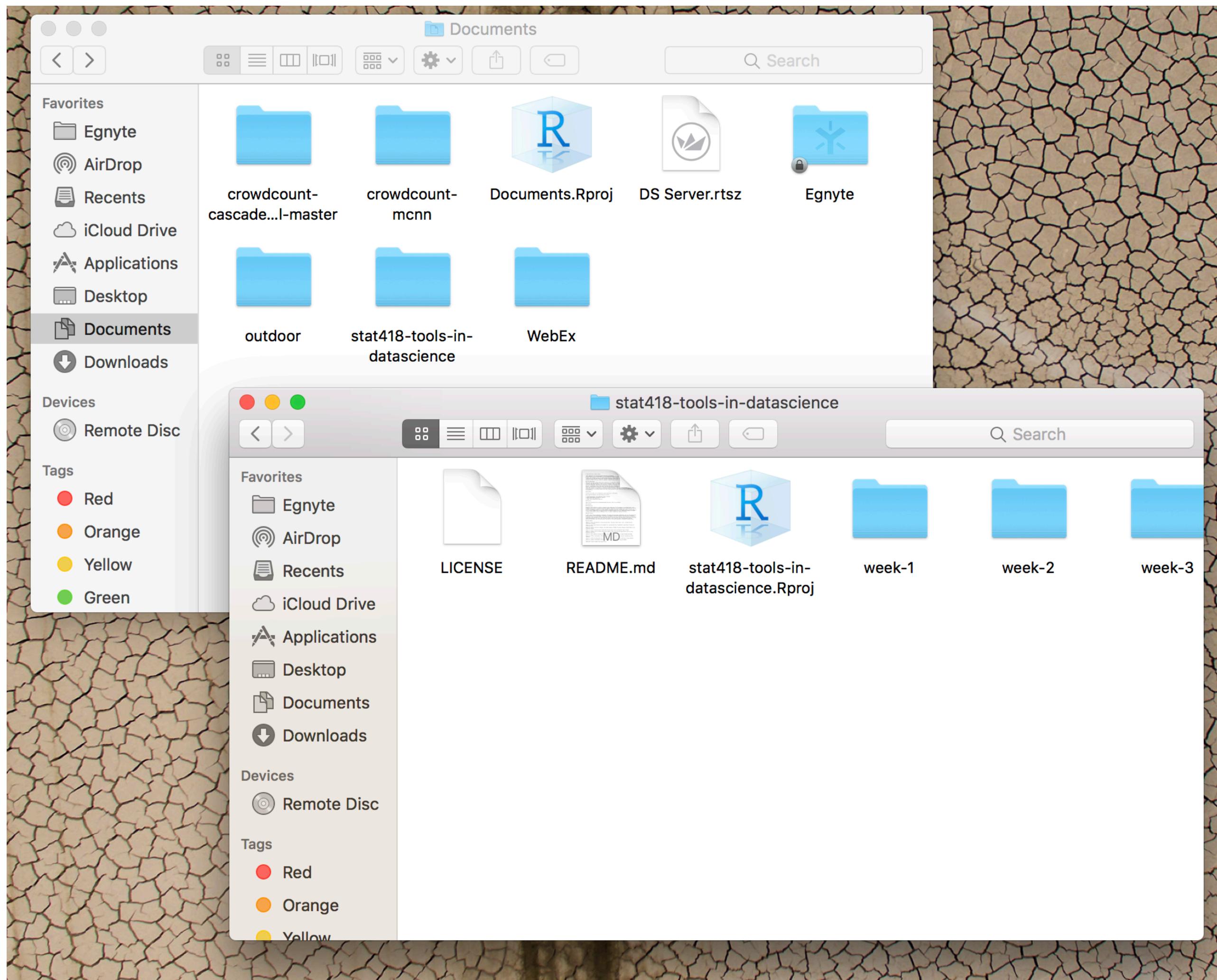
First, commands to explore your file system, walk through directories and list files

`pwd`, `ls`, `cd`

`mkdir`, `rmdir`

`cp`, `mv`, `rm`

```
>Last login: Tue Apr  9 17:29:09 on ttys001
You have new mail.
NA-nlanghol-A02:~ nlangholz$ pwd
/Users/nlangholz
NA-nlanghol-A02:~ nlangholz$ cd Documents/
NA-nlanghol-A02:Documents nlangholz$ pwd
/Users/nlangholz/Documents
NA-nlanghol-A02:Documents nlangholz$ ls
DS Server.rtsz          crowdcount-cascaded-mtl-master
Documents.Rproj          crowdcount-mcnn
Egnyte                   outdoor
WebEx                    stat418-tools-in-datascience
NA-nlanghol-A02:Documents nlangholz$ ls -l
total 32
-rw-r--r--@ 1 nlangholz  staff  9379 Nov  6  2017 DS Server.rtsz
-rw-r--r--  1 nlangholz  staff   205 Jan 31  2018 Documents.Rproj
drwxr-xr-x@ 7 nlangholz  staff   224 Jan 11 11:00 Egnyte
drwxr-xr-x@ 2 nlangholz  staff    64 Jan 18  2018 WebEx
drwxr-xr-x@ 16 nlangholz staff   512 Feb 27  2018 crowdcount-cascaded-mtl-master
drwxr-xr-x  21 nlangholz staff   672 Mar  2  2018 crowdcount-mcnn
drwxr-xr-x  28 nlangholz staff   896 Mar 27 22:18 outdoor
drwxr-xr-x  12 nlangholz staff   384 Apr  9 12:49 stat418-tools-in-datascience
NA-nlanghol-A02:Documents nlangholz$
```



Another view of the filesystem; here your Mac will display directories as folders and of course you navigate by clicking rather than typing commands



stat418-tools-in-datascience — bash — 87x18

```
NA-nlanghol-A02:stat418-tools-in-datascience nlangholz$  
NA-nlanghol-A02:stat418-tools-in-datascience nlangholz$  
NA-nlanghol-A02:stat418-tools-in-datascience nlangholz$  
NA-nlanghol-A02:stat418-tools-in-datascience nlangholz$ ls -al  
total 64  
drwxr-xr-x 12 nlangholz staff 384 Apr  9 12:49 .  
drwx-----+ 15 nlangholz staff 480 Apr  9 17:30 ..  
-rw-r--r--  1 nlangholz staff 8196 Apr  9 11:50 .DS_Store  
drwxr-xr-x  5 nlangholz staff 160 Apr  9 12:49 .Rproj.user  
drwxr-xr-x 16 nlangholz staff 512 Apr  9 12:51 .git  
-rw-r--r--  1 nlangholz staff 653 Apr  9 11:50 .gitignore  
-rw-r--r--  1 nlangholz staff 1057 Apr  9 11:50 LICENSE  
-rw-r--r--  1 nlangholz staff 4803 Apr  9 11:50 README.md  
-rw-r--r--  1 nlangholz staff 205 Apr  9 12:49 stat418-tools-in-datascience.Rproj  
drwxr-xr-x  4 nlangholz staff 128 Apr  9 11:50 week-1  
drwxr-xr-x  4 nlangholz staff 128 Apr  9 12:13 week-2  
drwxr-xr-x  3 nlangholz staff  96 Apr  9 11:50 week-3  
NA-nlanghol-A02:stat418-tools-in-datascience nlangholz$
```

## Kinds of Files

What you'll notice right away is that there are different types of files having different permissions

Unix filesystem conventions places (shared, commonly used) executable files in places like /usr/bin or /usr/local/bin

Different files are opened by different kinds of programs, in OSX, there is a beautiful command called open that decides which program to use

Although we saw a terminal environment in our Rstudio docker container, now let's get another docker image that has some nice data science specific code.

we will use the docker image from the book [Data Science at the Command Line](#) by Jeroen Janssens

<https://www.datascienceatthecli.com/index.html>

First, we need to pull the image using the following

```
docker pull datascienceworkshops/data-science-at-the-command-line
```

# Data Science at the Command Line

To run the image

```
docker run --rm -it datascienceworkshops/data-science-at-the-command-line
```

To leave the environment simply type **exit**

again we didn't mount a local volume so we have no connection to our computer

Now to mount our local directories, cd to desired directory, then use  
on a Mac

```
docker run --rm -it -v`pwd`: /data datascienceworkshops /data-science-at-the-command-line
```

in Windows command line

```
docker run --rm -it -v %cd%: /data datascienceworkshops /data-science-at-the-command-line
```

in Windows powershell

```
docker run --rm -it -v ${PWD}: /data datascienceworkshops /data-science-at-the-command-line
```