

BIOS 6660 course overview

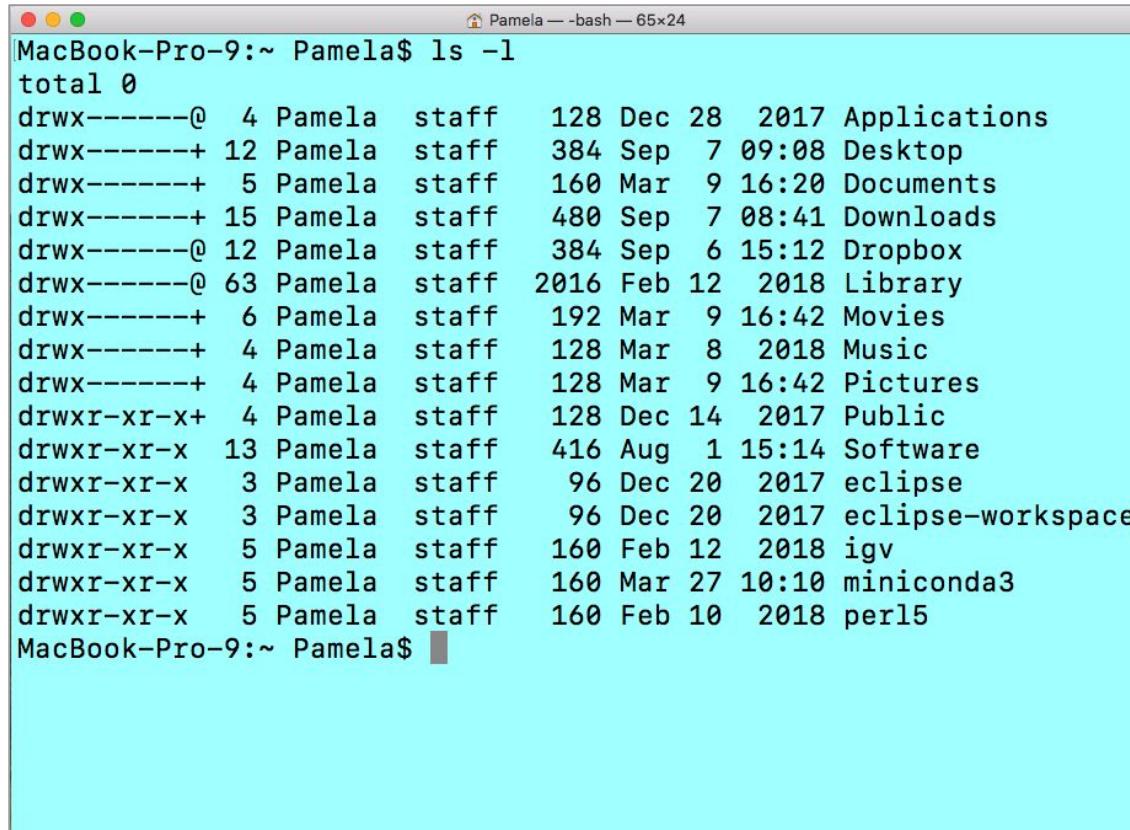
Lecture 1
BIOS 6660, Spring 2019
Instructor: Pam Russell



Basic overview

Weeks	Instructor	Topics
1-5	Pam Russell	Toolkit for reproducible data analysis
6-10	Lauren Vanderlinden	RNA-Seq, ChIP-Seq
11-15	Weiming Zhang	DNA methylation, genetic association

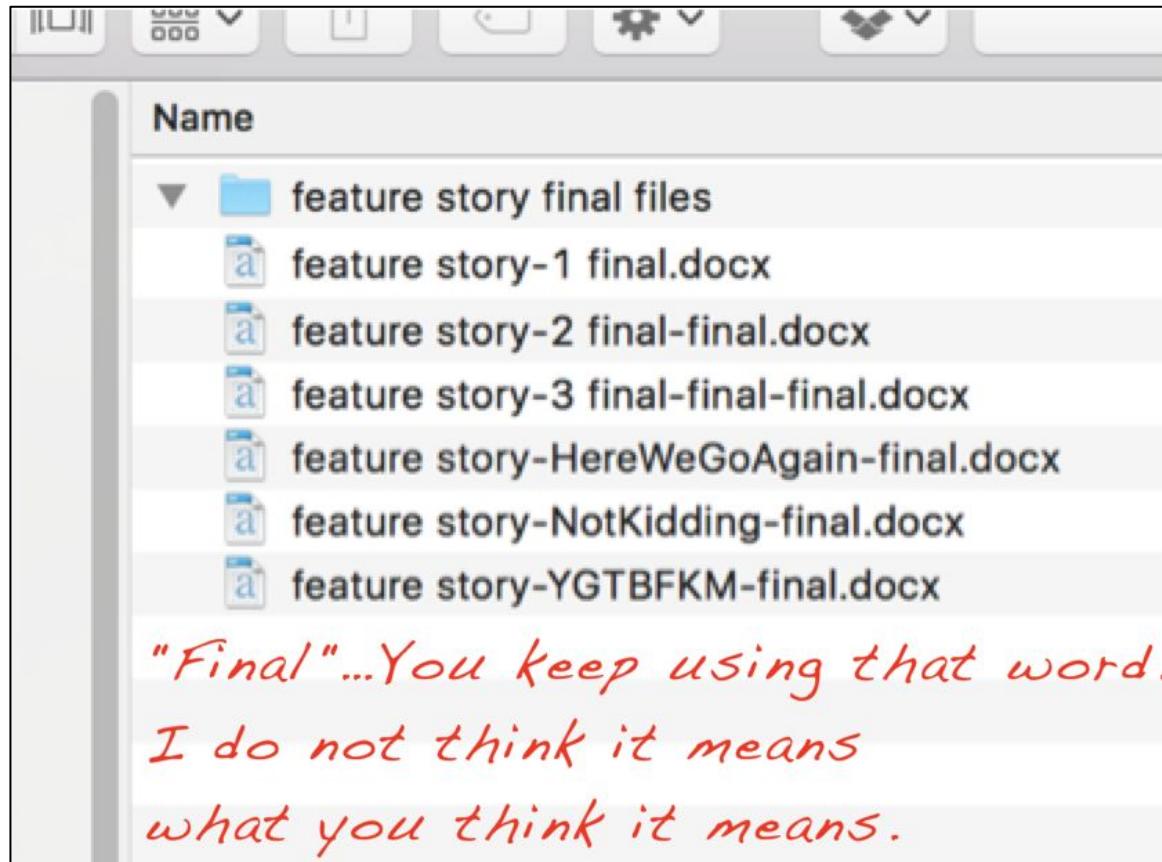
Part 1 topics: Command line



A screenshot of a macOS terminal window titled "Pamela — bash — 65x24". The window shows the output of the command "ls -l". The output lists various directories and files in the current directory, including "Applications", "Desktop", "Documents", "Downloads", "Dropbox", "Library", "Movies", "Music", "Pictures", "Public", "Software", "eclipse", "eclipse-workspace", "igv", "miniconda3", and "perl5". Each entry shows the file mode (e.g., drwxr-xr-x), owner (Pamela), group (staff), size, last modified date, and name.

```
MacBook-Pro-9:~ Pamela$ ls -l
total 0
drwx-----@  4 Pamela  staff   128 Dec 28  2017 Applications
drwx-----+ 12 Pamela  staff   384 Sep  7 09:08 Desktop
drwx-----+  5 Pamela  staff   160 Mar  9 16:20 Documents
drwx-----+ 15 Pamela  staff   480 Sep  7 08:41 Downloads
drwx-----@ 12 Pamela  staff   384 Sep  6 15:12 Dropbox
drwx-----@ 63 Pamela  staff  2016 Feb 12  2018 Library
drwx-----+  6 Pamela  staff   192 Mar  9 16:42 Movies
drwx-----+  4 Pamela  staff   128 Mar  8  2018 Music
drwx-----+  4 Pamela  staff   128 Mar  9 16:42 Pictures
drwxr-xr-x+  4 Pamela  staff   128 Dec 14  2017 Public
drwxr-xr-x+ 13 Pamela  staff   416 Aug  1 15:14 Software
drwxr-xr-x+  3 Pamela  staff    96 Dec 20  2017 eclipse
drwxr-xr-x+  3 Pamela  staff    96 Dec 20  2017 eclipse-workspace
drwxr-xr-x+  5 Pamela  staff   160 Feb 12  2018 igv
drwxr-xr-x+  5 Pamela  staff   160 Mar 27 10:10 miniconda3
drwxr-xr-x+  5 Pamela  staff   160 Feb 10  2018 perl5
MacBook-Pro-9:~ Pamela$
```

Part 1 topics: Version control



Part 1 topics: R programming

The screenshot shows the RStudio interface with the following components:

- Code Editor:** Displays the file `example_r_script.R` containing the following R code:

```
1 x <- 5
2 y <- 7
3 z <- x + y
4
```
- Environment View:** Shows the global environment with variables `x`, `y`, and `z` assigned values 5, 7, and 12 respectively.
- Console View:** Displays the output of running the script, including the R startup message and the assignment of `z`.
- Documentation View:** Shows the R documentation for the `MathFun` package under the `base` namespace, specifically the `Miscellaneous Mathematical Functions` section.

Part 1 topics: R Markdown

RStudio
Project: (None)

First_rmd_doc.Rmd

```
1 ---  
2 title: "First R Markdown Document"  
3 author: "Pam Russell"  
4 date: "10/16/2018"  
5 output: html_document  
6 ---  
7 |  
8 ``{r setup, include=FALSE}  
9 knitr::opts_chunk$set(echo = TRUE)  
10 ...  
11 ## R Markdown  
12  
13 This is an R Markdown document. Markdown is a simple formatting syntax for authoring  
HTML, PDF, and MS Word documents. For more details on using R Markdown see  
http://rmarkdown.rstudio.com.  
14  
15 When you click the Knit button a document will be generated that includes both  
content as well as the output of any embedded R code chunks within the document. You  
can embed an R code chunk like this:  
16  
17 ``{r cars}  
18 summary(cars)  
19 ...  
20  
21 ## Including Plots  
22  
23 You can also embed plots, for example:  
24  
25 ``{r pressure, echo=FALSE}  
26 plot(pressure)  
27 ...  
28 (Top Level) R Markdown
```

Console R Markdown

```
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.
```

Environment History Connections
Files Plots Packages Help Viewer
Publish

First R Markdown Document

Pam Russell
10/16/2018

R Markdown

This is an R Markdown document. Markdown is a simple formatting syntax for authoring HTML, PDF, and MS Word documents. For more details on using R Markdown see <http://rmarkdown.rstudio.com>.

When you click the **Knit** button a document will be generated that includes both content as well as the output of any embedded R code chunks within the document. You can embed an R code chunk like this:

```
summary(cars)
```

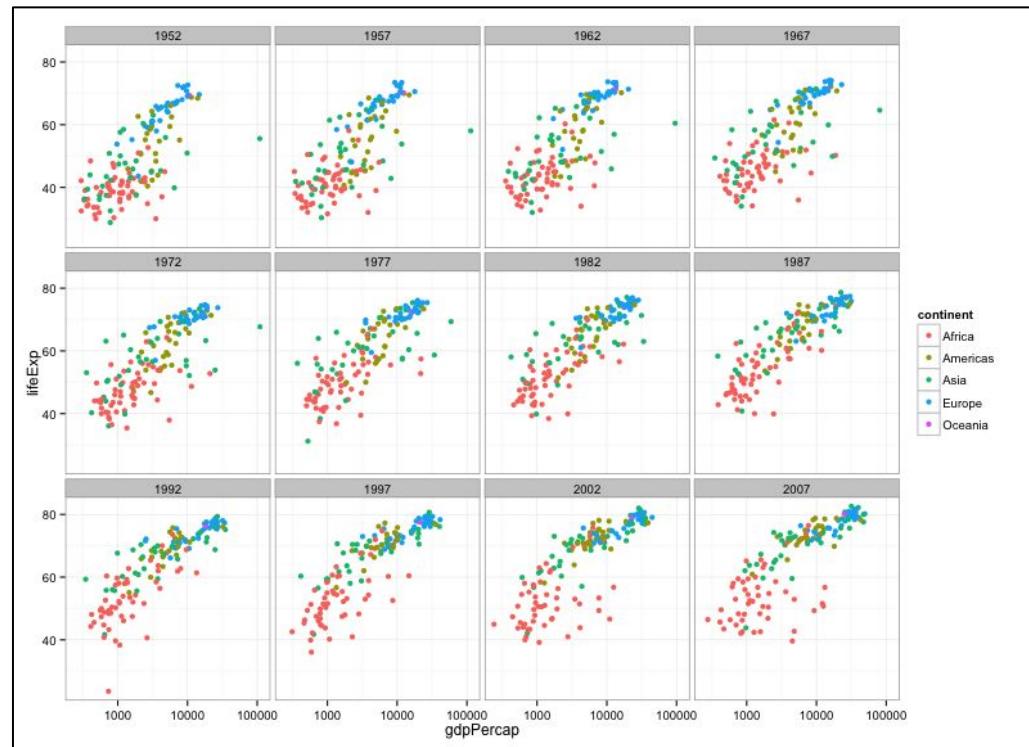
```
##      speed         dist  
##  Min.   : 4.0   Min.   :  2.00  
##  1st Qu.:12.0   1st Qu.: 26.00  
##  Median :15.0   Median : 36.00  
##  Mean   :15.4   Mean   : 42.98  
##  3rd Qu.:19.0   3rd Qu.: 56.00  
##  Max.   :25.0   Max.   :120.00
```

Including Plots

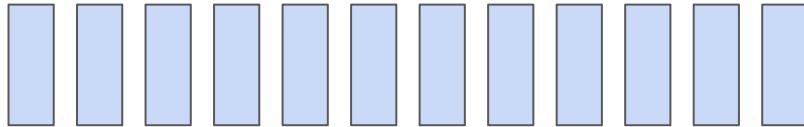
You can also embed plots, for example:



Part 1 topics: R Tidyverse



Part 1 topics: Code organization



Individual statement



Block of code



Whole file

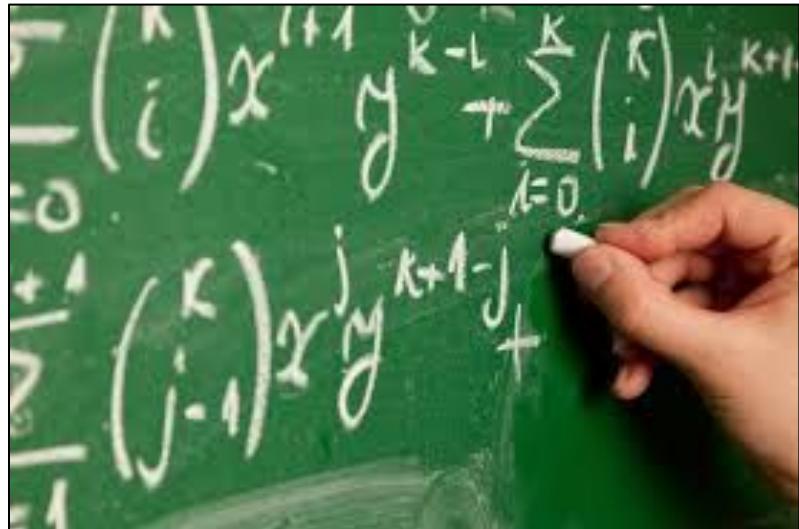


Multiple files



Package

Part 1 topics: Code quality

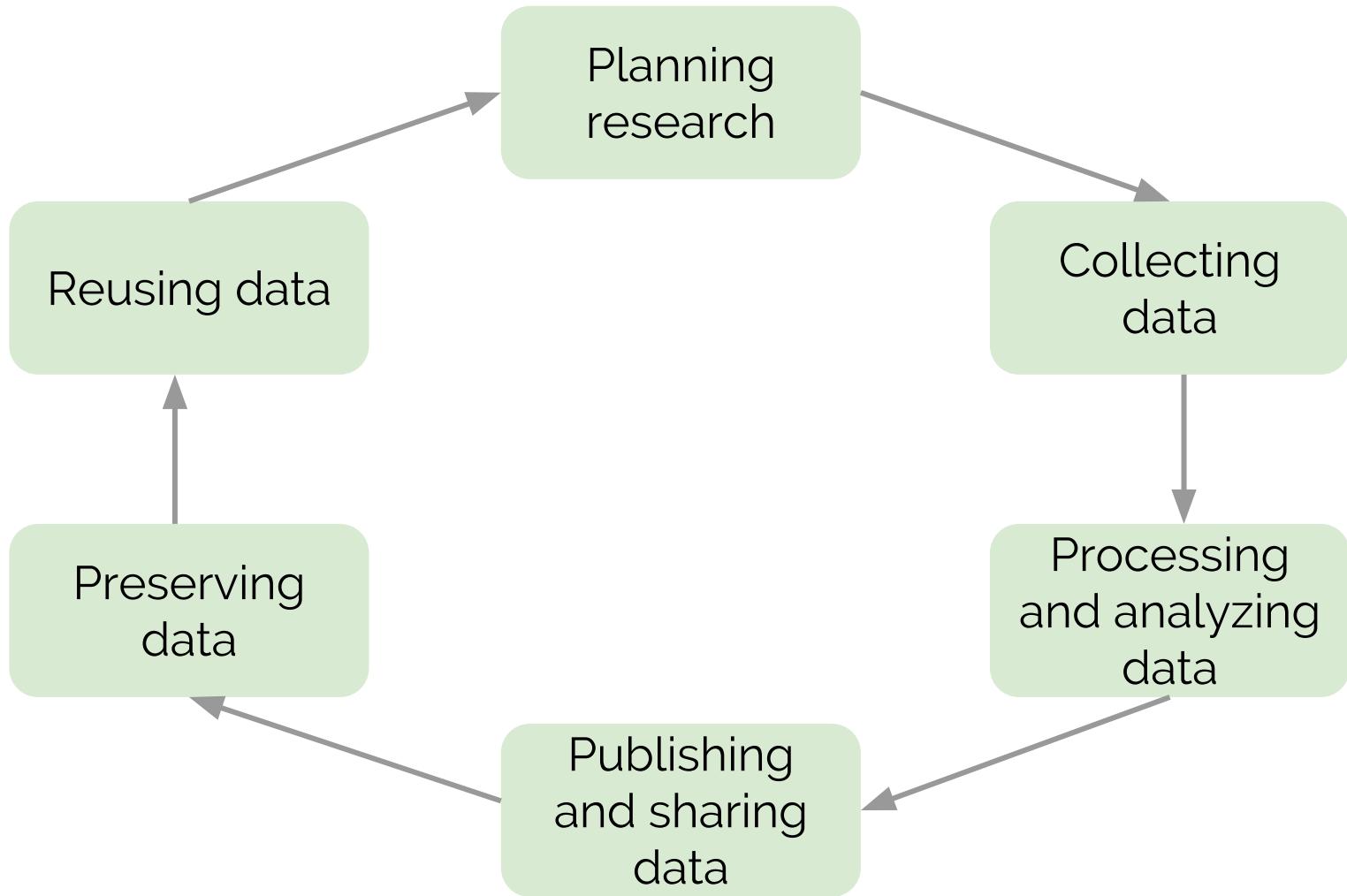


Correctness

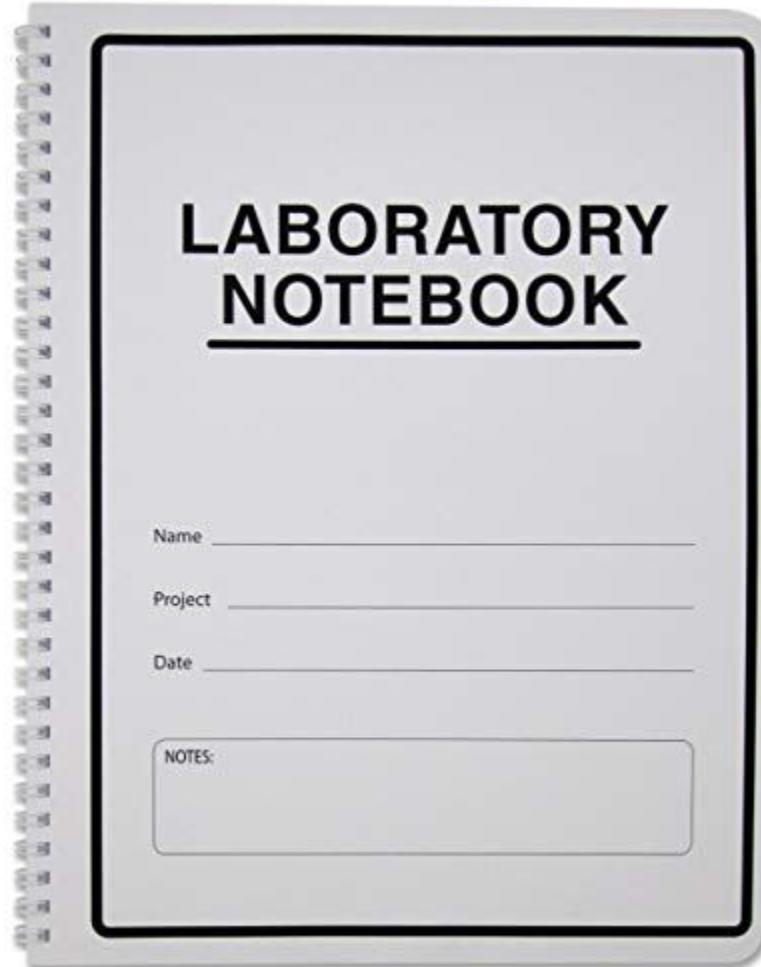


Effective communication

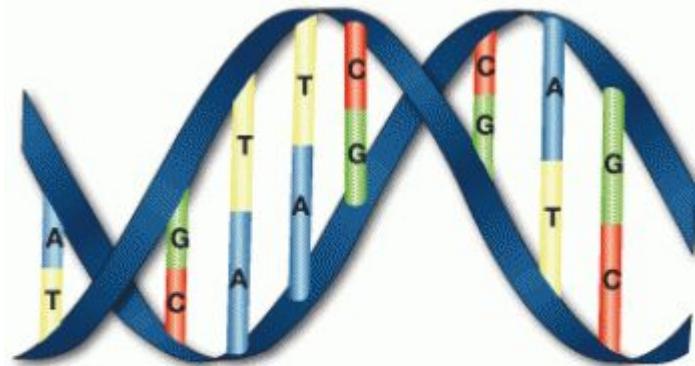
Part 1 topics: Data management



Part 1 topics: Reproducibility



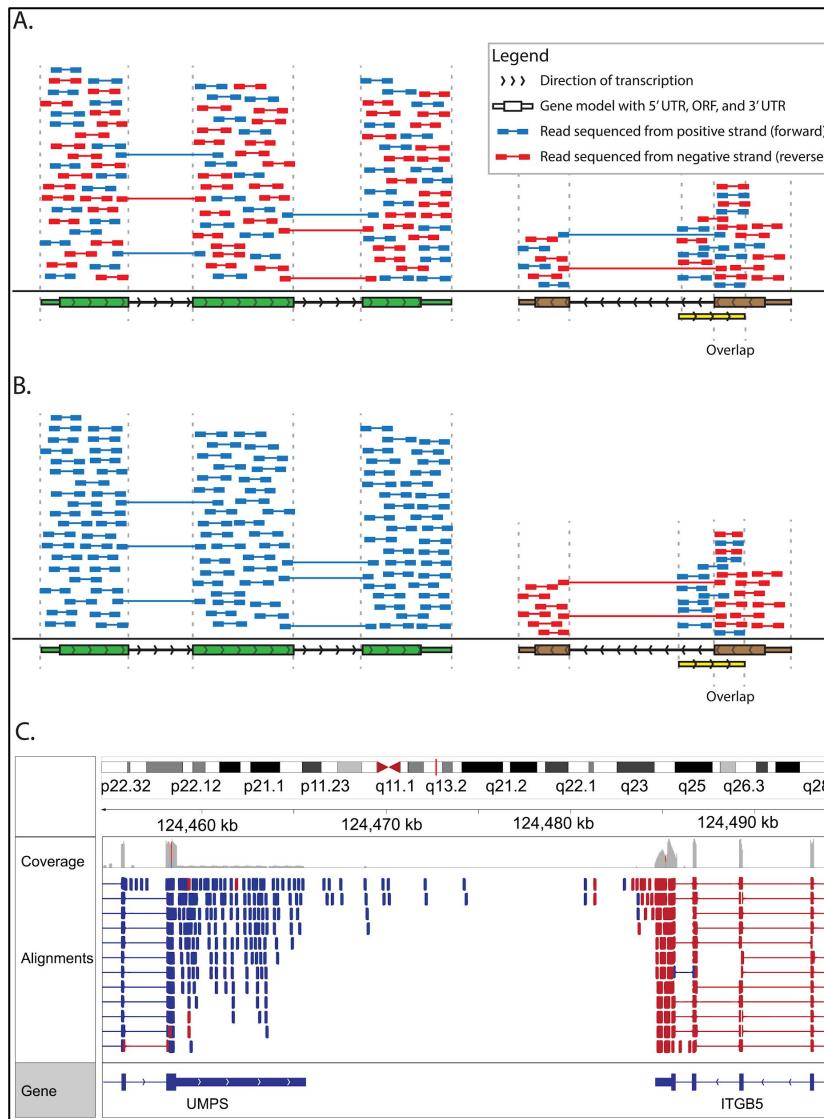
Parts 2 and 3: Genomic data analysis



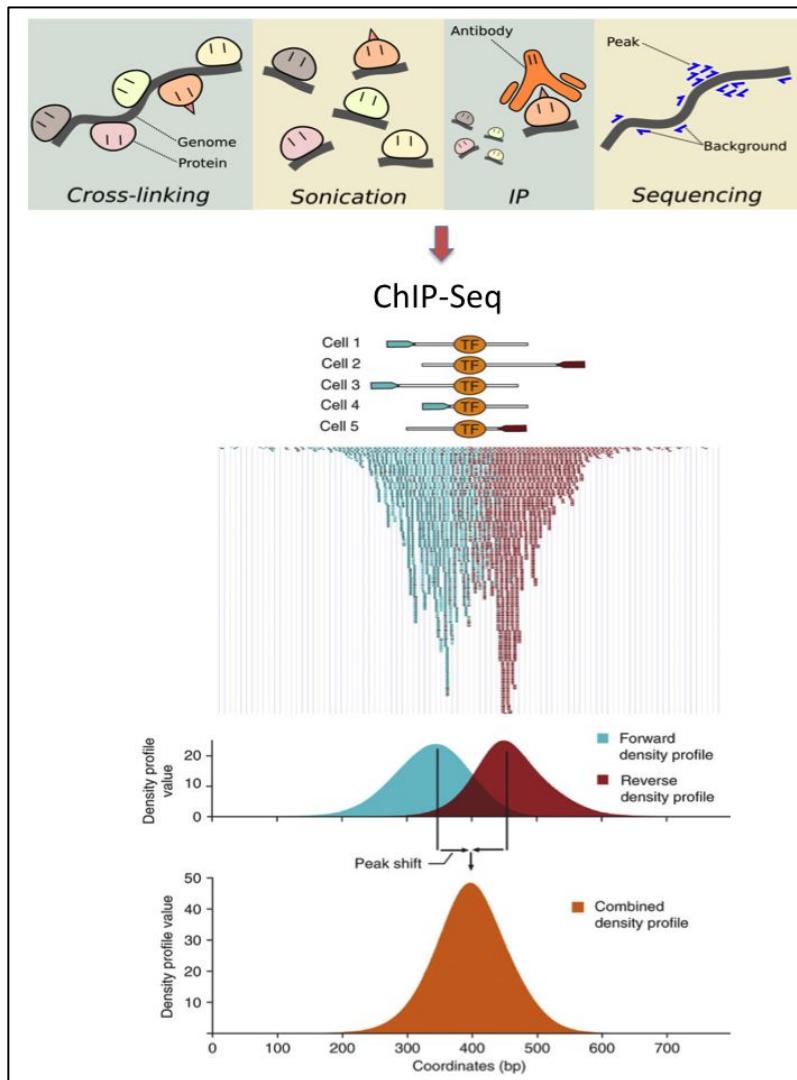
Thymine (Yellow) = T Guanine (Green) = G
Adenine (Blue) = A Cytosine (Red) = C



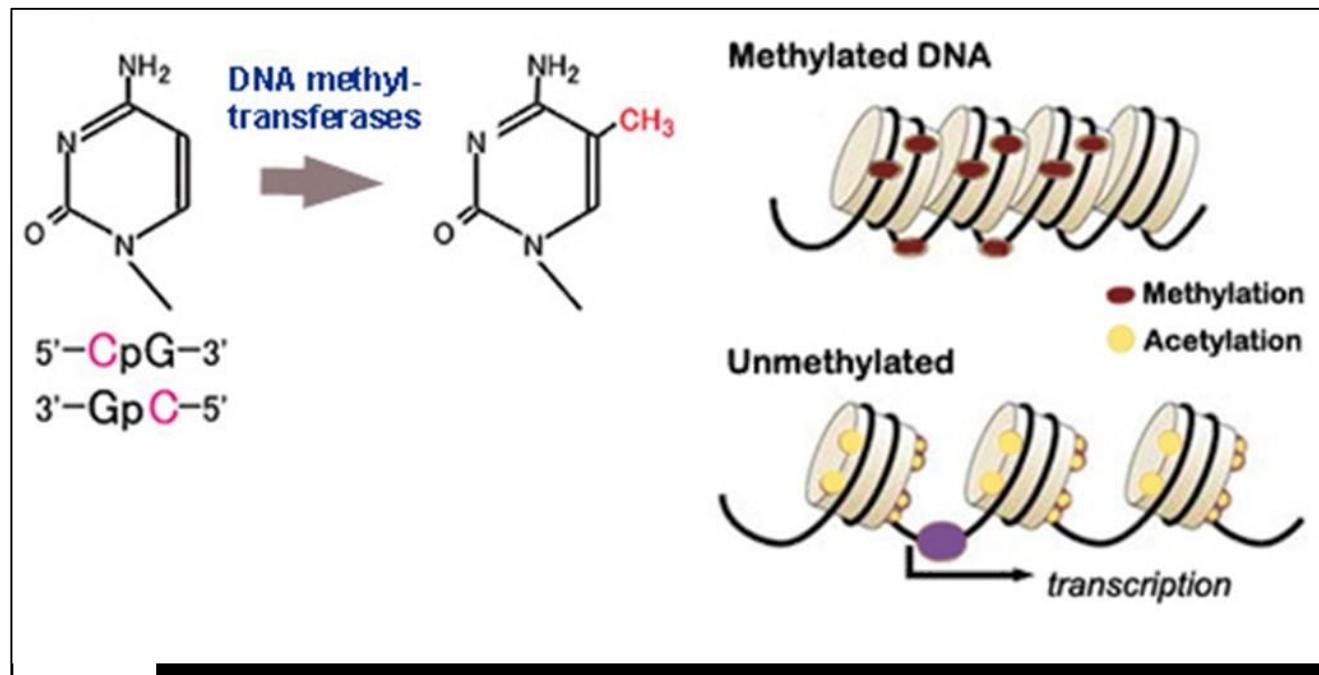
Part 2 topics: RNA-Seq analysis



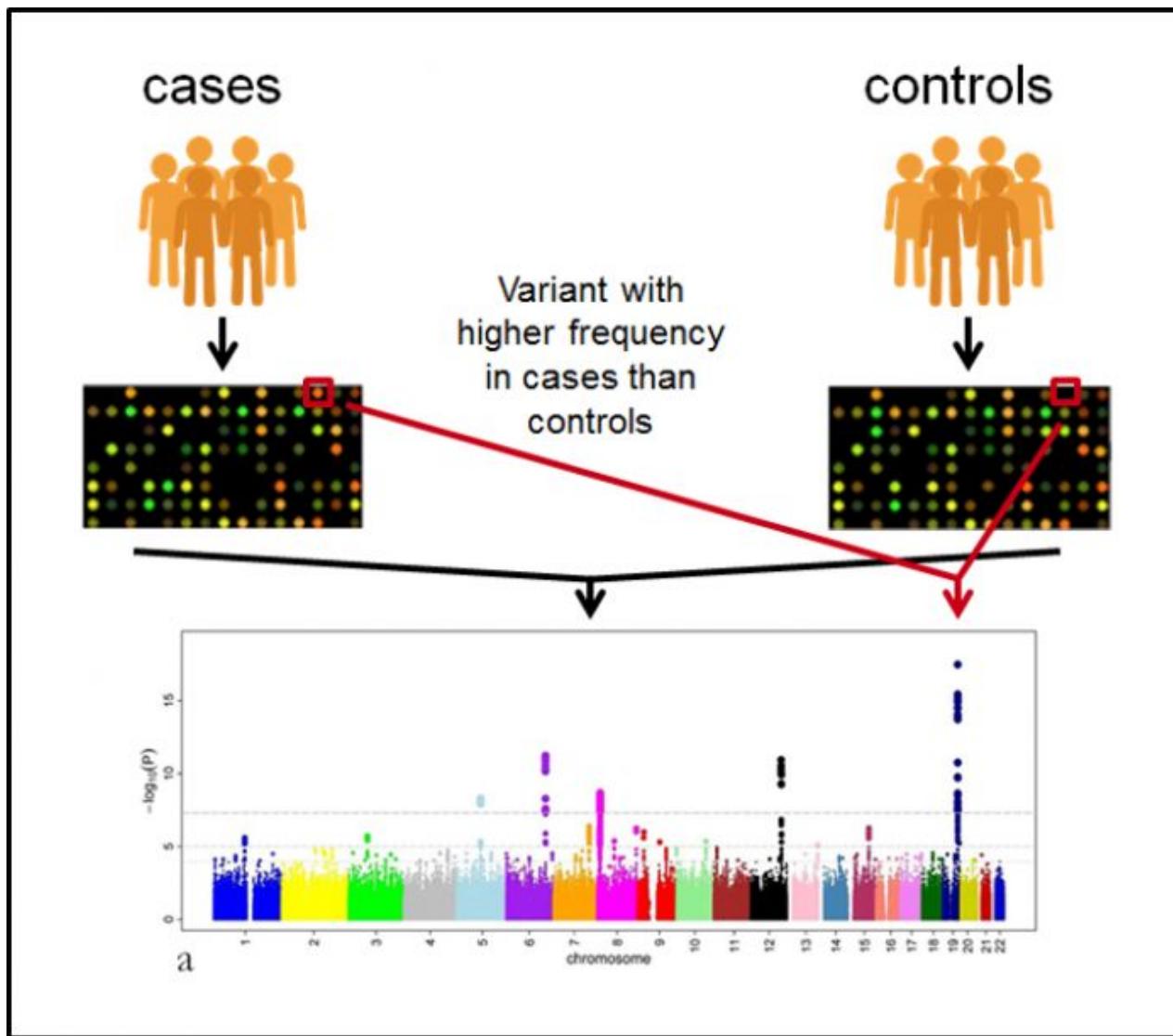
Part 2 topics: ChIP-Seq analysis



Part 3 topics: DNA methylation analysis



Part 3 topics: Genetic association



Everything is on Canvas

- Syllabus
- Schedule
- Lecture slides
- Homework assignments
- Turn in homework

Evaluation

12 homework assignments

No exams

Tools

- ~~Book~~

- R
- RStudio
- R packages
- Command line
- Version control
- Yampa server
- Publically available data

Yampa server details

Some of our work will be done on a remote server called Yampa.

You should have received an email from Pam with login credentials for Yampa. If not, let Pam know ASAP.

Homework 1 first requires you to set up a Linux shell on your computer. From the shell, you can log into Yampa with the command `ssh <username>@yampa.ucdenver.pvt`. You must be on the university network or [VPN](#) to access Yampa.

The first time you log into Yampa, immediately type `passwd+[enter]` to change your password.

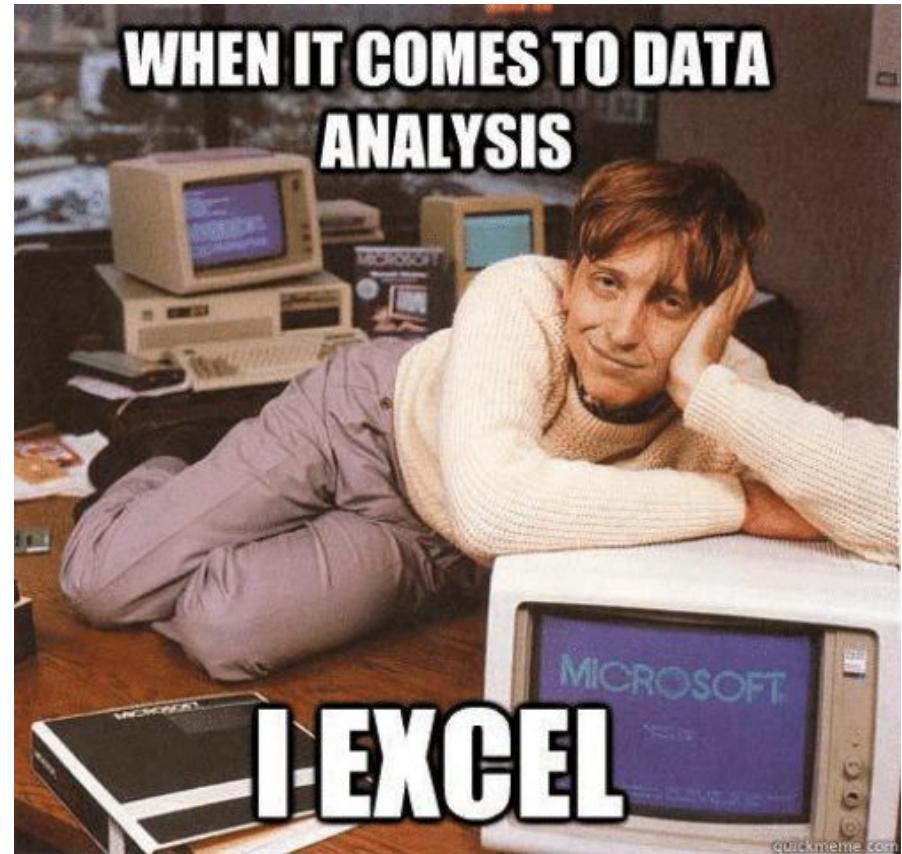
Pam's contact info

Office hours Fridays 10:00-12:00
Or by appointment
Building 406 room 107

pamela.russell@ucdenver.edu

Analysis horror stories

Lecture 1
BIOS 6660, Spring 2019
Instructor: Pam Russell





**KEEP
CALM
AND
GO
PRACTICE**

Data management

Data management

Acquiring, organizing, maintaining,
and protecting data

"Then we saw sequences start to vanish
and we were like, 'Oh my god'"



rm -r -f *

Paper retracted when data can't be found

Retraction: CPAP for the Metabolic Syndrome in Patients with Obstructive Sleep Apnea. N Engl J Med 2011;365:2277-86

TO THE EDITOR:

As the authors of the article entitled “CPAP for the Metabolic Syndrome in Patients with Obstructive Sleep Apnea”¹ published in the Journal on December 15, 2011, we regret to report that transcription errors occurred in the Supplementary Appendix, available with the full text of the article at NEJM.org. There were multiple errors in the table on pages 18 and 19 of the Supplementary Appendix concerning data on the accumulation of abdominal fat as assessed with the use of computed tomography and on carotid intima–media thickness as assessed with the use of ultrasonography. These errors, in turn, changed some values in Table 4 of the article. Although these changes do not alter the conclusions of the article, the primary data could not be located to verify corrections made from secondary tables. Accordingly, we have no way of confirming the correct data and, with regret, wish to retract the article.

R

R

Transparent, automated, reproducible
analysis

Excel error leads to purchase of toxic assets



ACME Sales Numbers.xls

	A	B	C	D	E	F	G	H	I	J
10						\$5,135	\$3,555	\$1,501	\$316	\$7
11	S					\$2,507	\$6,104	\$8,502	\$981	\$87
12						\$10,328	\$14,794	\$16,007	\$2,403	\$1,89
13	Roller SKI					\$170,506	\$183,001	\$355,877	\$187,093	\$276,95
20										
21										
22	Total		Jan	Feb	Mar	Apr	May	Jun	Jul	Aug
23			\$220,889	\$1,120,814	\$339,532	\$413,238	\$333,936	\$471,817	\$278,935	\$387,97
24										
25										

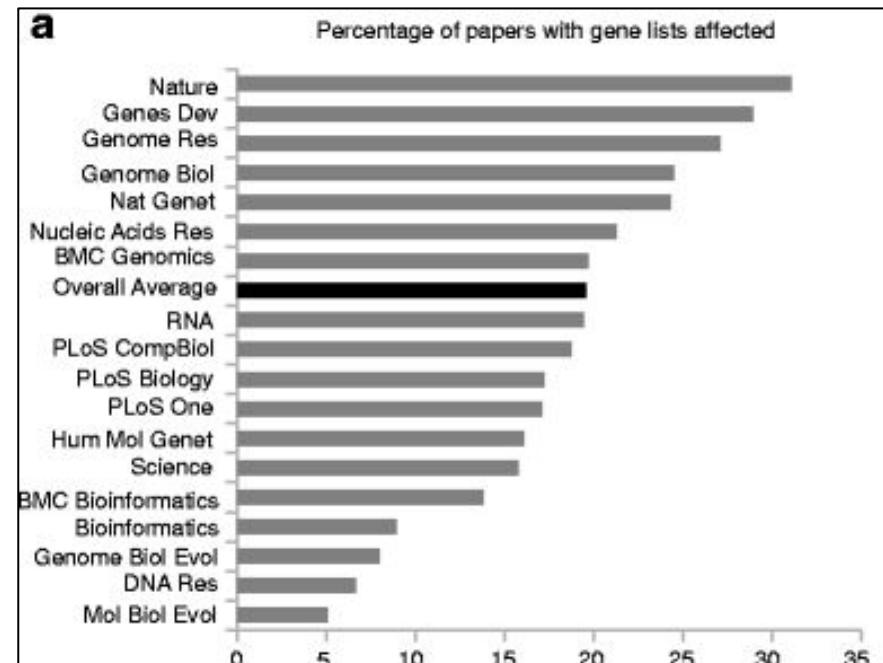
A double line displays in place of the hidden rows and the row numbers are also hidden.

Excel is just trying to help

“*MARCH1*” → “1-Mar”

“*SEPT2*” → “2006/09/02”

“2310009E13” → “2.31E+13”



R for transparent, end-to-end workflows

RStudio Source Editor

```
longrunning.Rmd x
3 ---
4
5 ````{r}
6 # load libraries
7 library(dygraphs)
8 library(leaflet)
9
10 # import cars data
11 source("import.R")
12 cars <- import_data("cars.csv")
13 cars <- cars[order(cars$mpg),]
14 cars <- head(cars, n = 15)
15
16 # import cities data
17 cities <- readr::read_csv("cities.csv")
18 ```
19
20 ## dygraphs
21
22 Dygraphs provides rich facilities for charting time-series data in R and includes support for many interactive features including series/point highlighting, zooming, and panning.
23
24 ````{r}
25 library(dygraphs)
26 dygraph(nhtemp, main = "New Haven Temperatures") %>%
27   dyRangeSelector(dateWindow = c("1920-01-01", "1960-01-01"))
28 ````
```

New Haven Temperatures

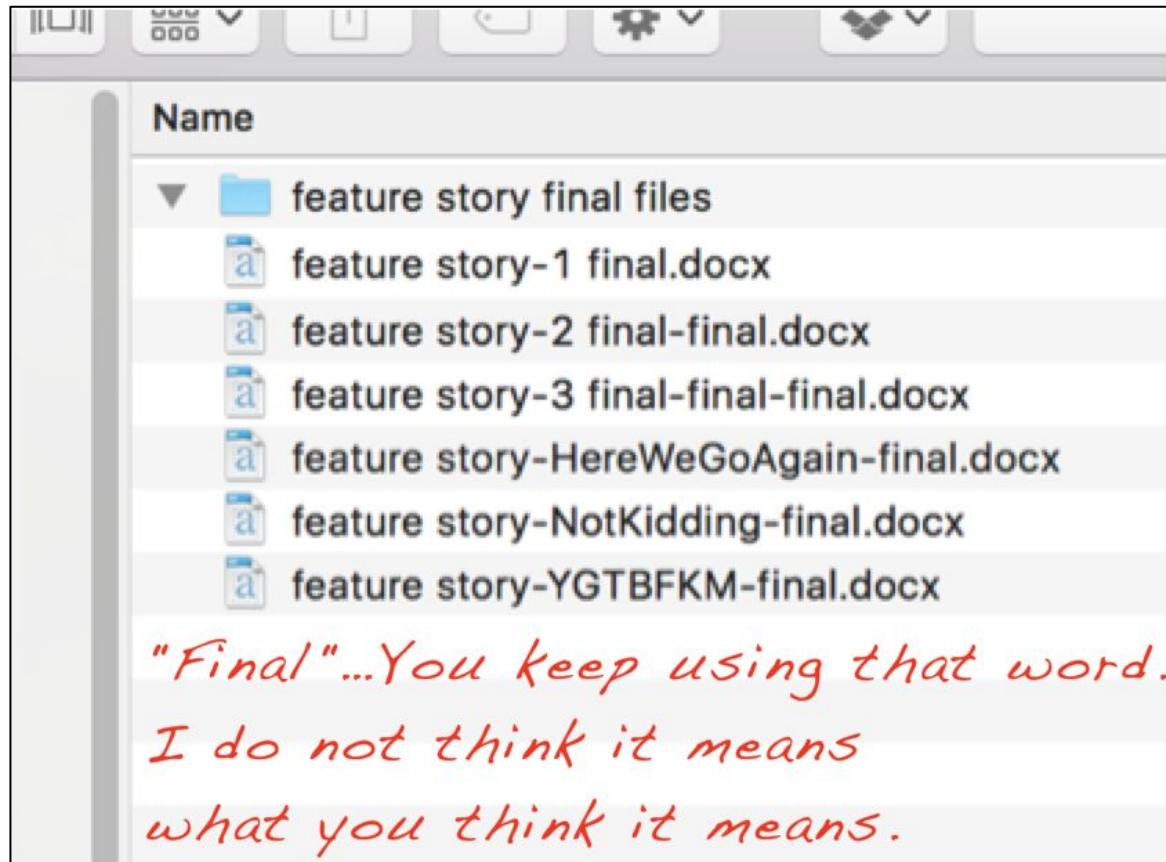
22:1 dygraphs Run All:

Version control

Version control

Maintaining an archive of snapshots throughout project life cycle

DIY version control



Google Docs

several issues: the small size of miRNA-molecules, the small number of unique miRNA-molecules compared to large RNAs, the fact that highly-similar miRNAs are often can be transcribed from multiple genomic loci, and the presence of the presence of miRNA isoforms known as isomiRs. Quantification methods that do not adequately failing to address these issues can return misleading information. We propose a novel quantification method, miR-MaGIC, designed to address addressing these concerns. miR-MaGIC performs highly-stringent mapping to a core region of each miRNA — using individualized miRNA sequences if available — and defines a meaningful set of miRNA sequences by collapsing the miRNA space to the level of “functional groups” of miRNAs. We hypothesize that these two features, mapping stringency and collapsing, provide more optimal quantification to a more meaningful unit (i.e., miRNA family). To test this hypothesis, we evaluate miR-MaGIC and a range of several published methods on 212 mouse whole-brain small RNA-seq libraries, evaluating each method’s ability to accurately reflect global miRNA expression profiles. We reason that methods introduce less bias in subsequent analyses if they report define accuracy as total counts close to the total number of input reads originating from miRNAs. We find that miR-MaGIC, and our modifications to other methods to incorporate both stringency and collapsing, provide the most accurate counts according to this measure. miR-MaGIC software is freely available at <https://github.com/KechrisLab/miR-MaGIC>.

Version history

Only show named versions

Total: 139 edits

December 2017

- December 27, 4:33 PM Current version Pamela Russell
- December 21, 4:03 PM All anonymous users
- December 21, 2:10 PM All anonymous users Pamela Russell

November 2017

- November 6, 4:32 AM Pamela Russell
- November 6, 3:55 AM Pamela Russell
- November 5, 9:25 PM All anonymous users
- November 3, 12:01 PM Pamela Russell
- November 3, 11:09 AM All anonymous users
- November 3, 10:21 AM All anonymous users

November 1, 2:49 PM Show changes

Find and replace

10 → ten

A clearer head would have foreseen the carnage. 110 became **1ten**, 5.10 became **5.ten** and 101010 became, well, you get the idea. One of the things I learnt that day is that banking websites have a lot of numbers!

ten → 10

This was clearly pure genius as it immediately fixed everything I'd just broken. It also brought another problem to my **at10tion**.

Code quality

Code quality

Effective code design => fewer errors

NASA loses Mars Climate Orbiter

Confusion leads to Mars failure

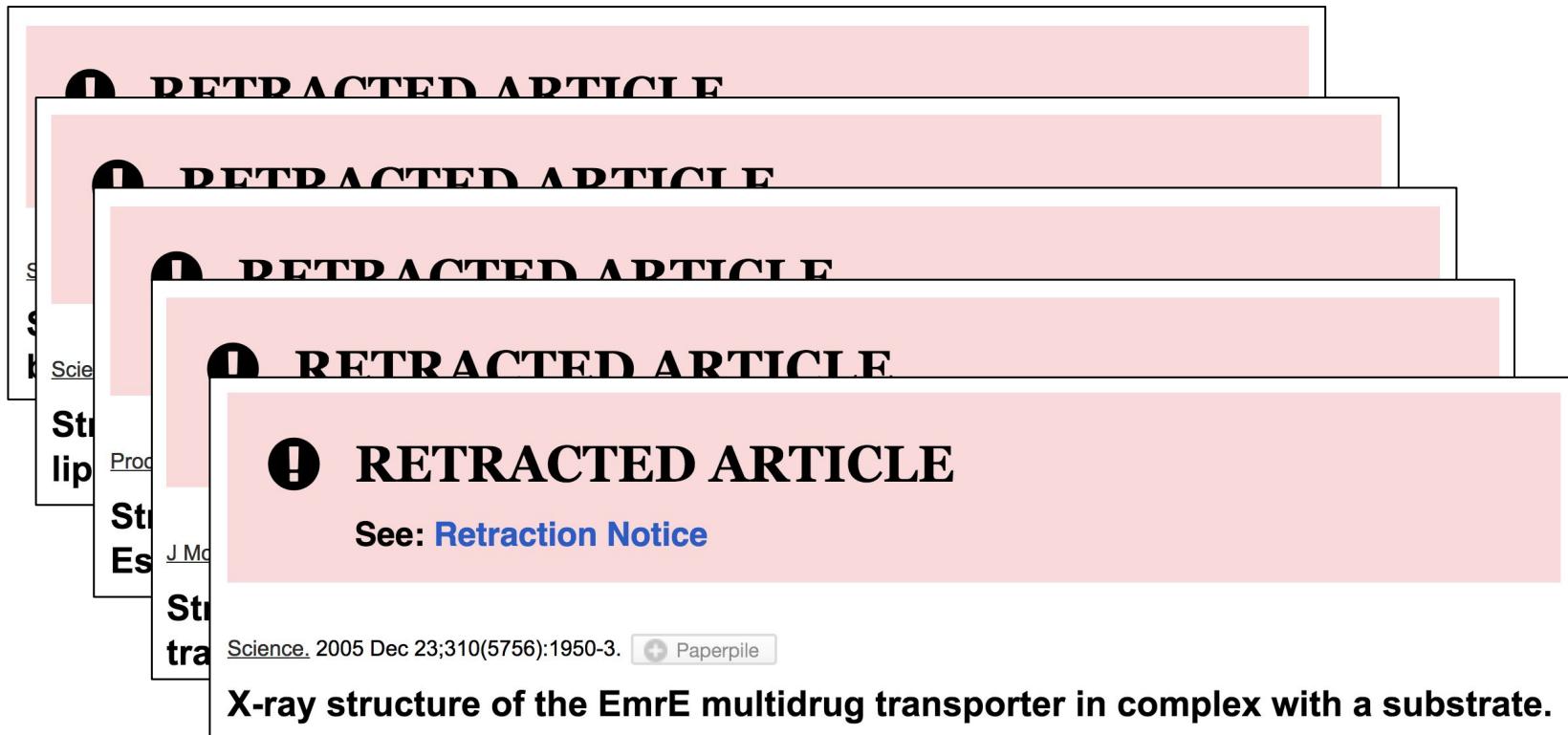


The Mars Climate Orbiter: Now in pieces on the planet's surface

The Mars Climate Orbiter Spacecraft was lost because one Nasa team used imperial units while another used metric units for a key spacecraft operation.

"People sometimes make errors," said Dr. Edward Weiler, NASA's Associate Administrator for Space Science. **"The problem here was not the error,** it was the failure of NASA's systems engineering, and the checks and balances in our **processes to detect the error.** That's why we lost the spacecraft."

“I hope people will understand that it was a mistake, and I’m very sorry for it.”

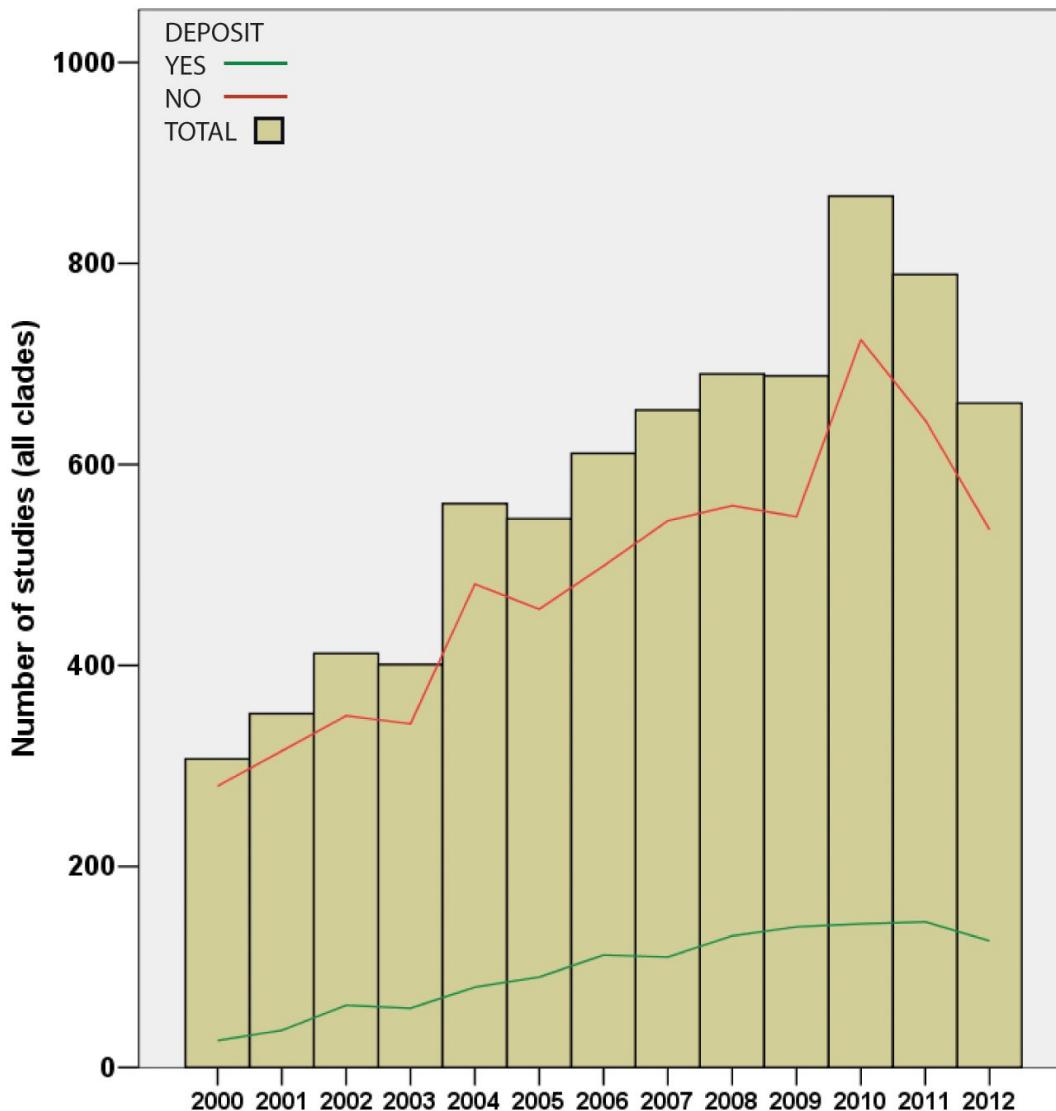


Reproducibility

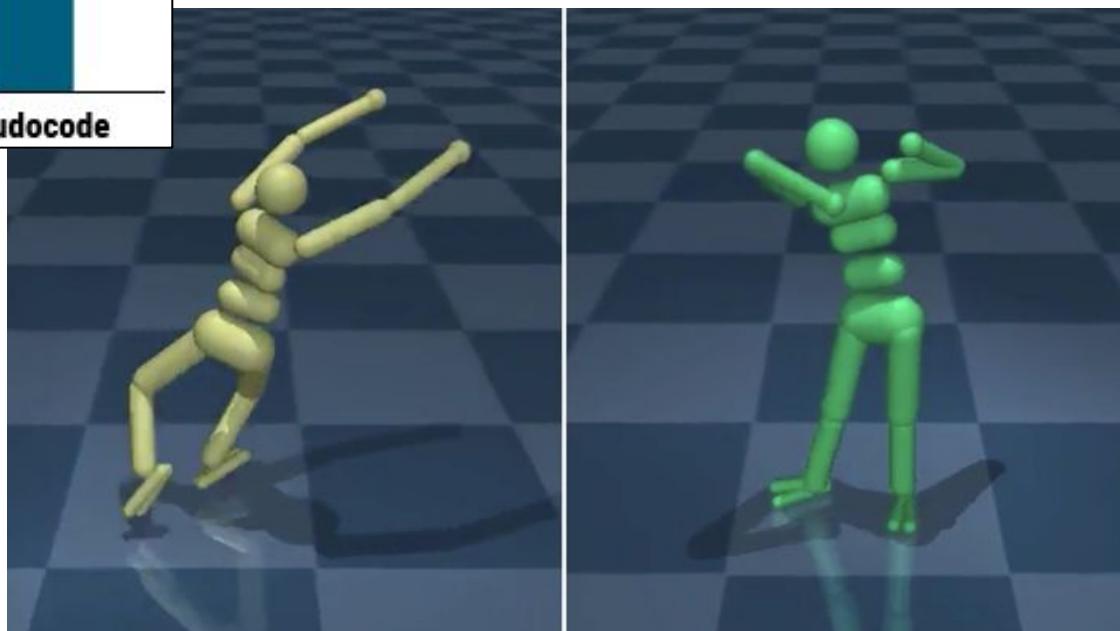
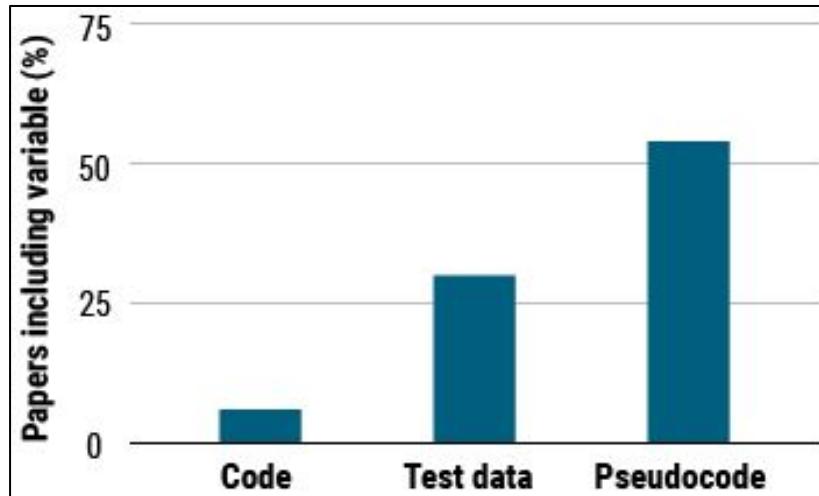
Reproducibility

The ability to recapitulate the results in a paper from publicly available data, code, and workflows

Lost branches on the tree of life



**“We tried for two months and we
couldn’t get anywhere close.”**

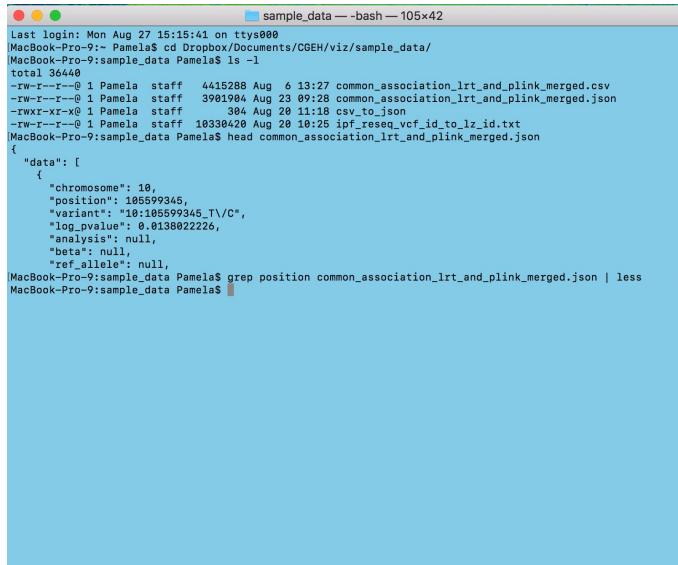


Command line

Command line

Interact with computer through
commands

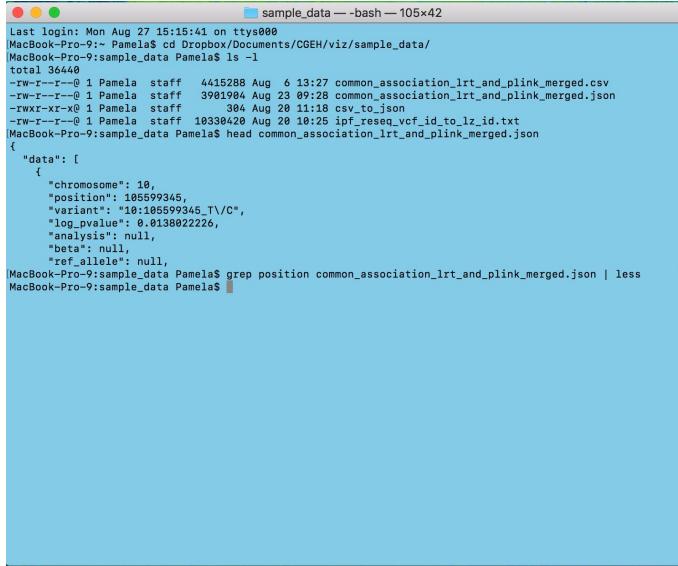
Text processing one liners



```
sample_data -- bash -- 105x42
Last login: Mon Aug 27 15:15:41 on ttys000
MacBook-Pro-9:~ Pamela$ cd Dropbox/Documents/CGEH/viz/sample_data/
MacBook-Pro-9:sample_data Pamela$ ls -l
total 36448
-rw-r--r--@ 1 Pamela  staff  4415288 Aug  6 13:27 common_association_lrt_and_plink_merged.csv
-rw-r--r--@ 1 Pamela  staff  3981904 Aug 23 09:28 common_association_lrt_and_plink_merged.json
-rwxr-xr-x@ 1 Pamela  staff    304 Aug 28 11:18 csv_to_json
-rw-r--r--@ 1 Pamela  staff  10330420 Aug 28 10:25 ipf_reseq_vcf_id_to_lz_id.txt
MacBook-Pro-9:sample_data Pamela$ head common_association_lrt_and_plink_merged.json
{
  "data": [
    {
      "chromosome": 18,
      "position": 186599345,
      "variant": "10:186599345_T\\C",
      "log_pvalue": 0.0138022226,
      "analysis": null,
      "beta": null,
      "ref_allele": null,
      "alt_allele": null
    }
  ]
}
MacBook-Pro-9:sample_data Pamela$ grep position common_association_lrt_and_plink_merged.json | less
MacBook-Pro-9:sample_data Pamela$
```

- Print first 10 lines of file
- Find and replace string inside file
- Count lines or characters in file
- Sort file in decreasing order by column 5
- Make a new file with columns 7, 3, 6 of original file
- Find differences between two files
- Delete all blank lines from file
- Remove duplicate lines from file
- View file with long lines chopped

Running analyses



```
Last login: Mon Aug 27 15:18:41 on ttys008
MacBook-Pro-9:~ Pamela$ cd Dropbox/Documents/CGEH/viz/sample_data/
MacBook-Pro-9:sample_data Pamela$ ls -l
total 36440
-rw-r--r-- 1 Pamela  staff  4415288 Aug  6 13:27 common_association_lrt_and_plink_merged.csv
-rw-r--r-- 1 Pamela  staff  3981984 Aug 23 09:28 common_association_lrt_and_plink_merged.json
-rwxr-xr-x  1 Pamela  staff    384 Aug 20 11:18 csv_to_json
-rw-r--r-- 1 Pamela  staff 10330428 Aug 20 10:25 ipf_reseq_vcf_id_to_lz_id.txt
MacBook-Pro-9:sample_data Pamela$ head common_association_lrt_and_plink_merged.json
{
  "data": [
    {
      "chromosome": 18,
      "position": 105599345,
      "variant": "10:105599345_T\\C",
      "log_pvalue": 0.0138022226,
      "analysis": null,
      "beta": null,
      "ref_allele": null,
      "alt_allele": null
    }
  ]
}
MacBook-Pro-9:sample_data Pamela$ grep position common_association_lrt_and_plink_merged.json | less
MacBook-Pro-9:sample_data Pamela$
```

- Most genomics tools work from command line
- Keep sets of commands in scripts for later reference and reuse
- Keep scripts under version control
- Manage version control from command line

Preview of what's next

Lecture 2

Command line and version control

Homework 1

- Currently posted; due next Tuesday
- Some problems are just tool setup that you can begin working on now
- Some problems use material from lecture 2

Getting help with homework

Office hours on Fridays are timed with assignments due Tuesdays

I won't answer last-minute emails on Monday night! :)