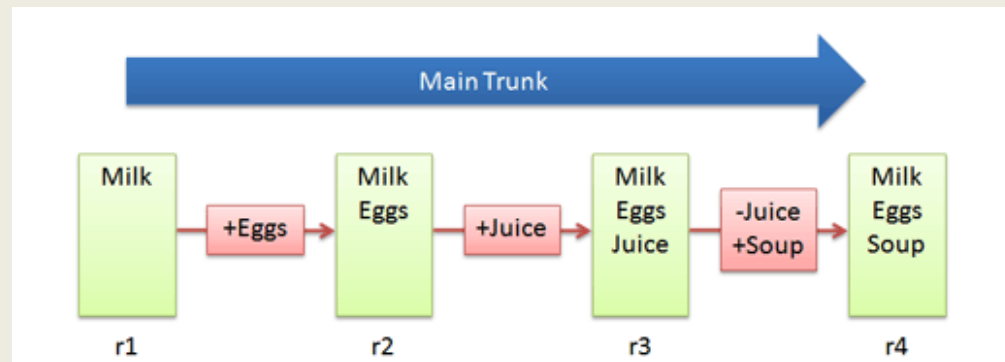


Version control and documentation

Lecture 8, CPSC 499 Fall 2018

What is version control?

- A computer program that tracks changes you have made to a set of files
- Lets you easily revert to old versions
- Lets you document why you made the changes
- Allows multiple users working on the same project





Git

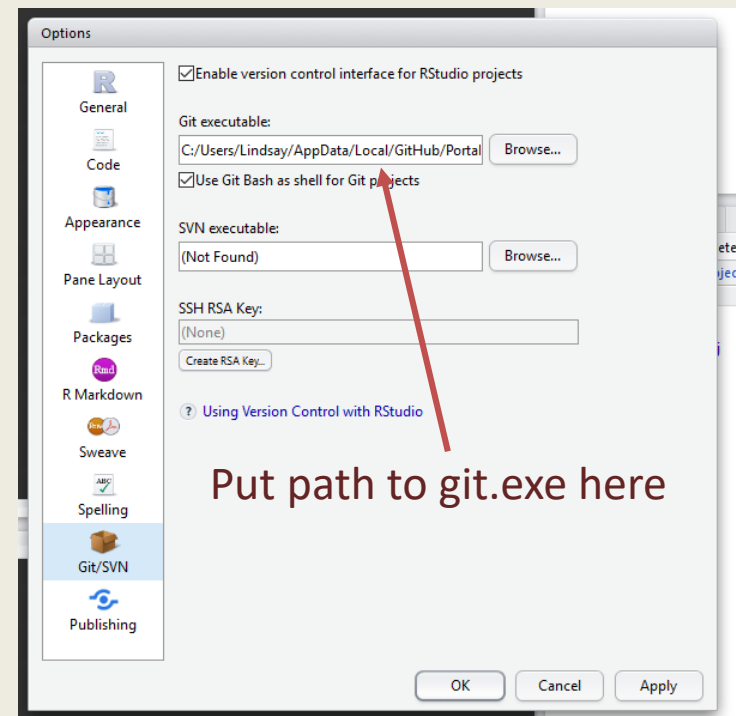
- Popular software for version control
- Integrates with RStudio to track changes to R scripts and source code on your computer
- Allows branching so you can try new things out before permanently incorporating changes into software

Useful info for integrating Git with RStudio

In RStudio Tools -> Global Options

At git-scm.com you can get a "portable" version of git.exe that doesn't require admin access. Use this to put Git on the lab computers.

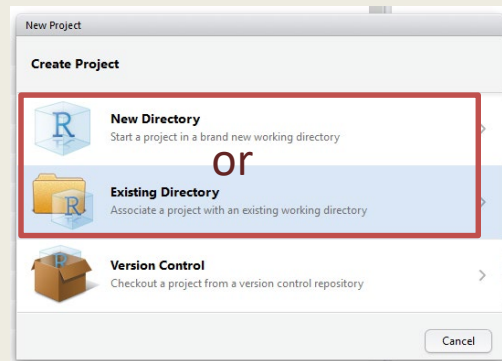
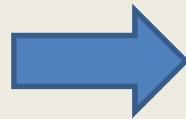
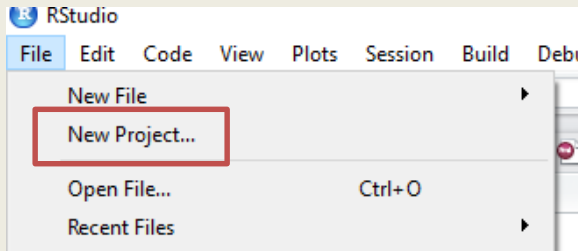
For your own computer, I recommend installing the regular version of Git.



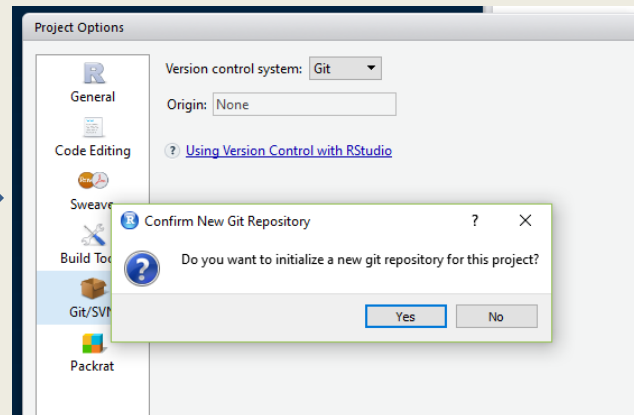
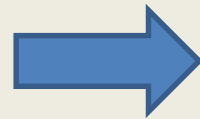
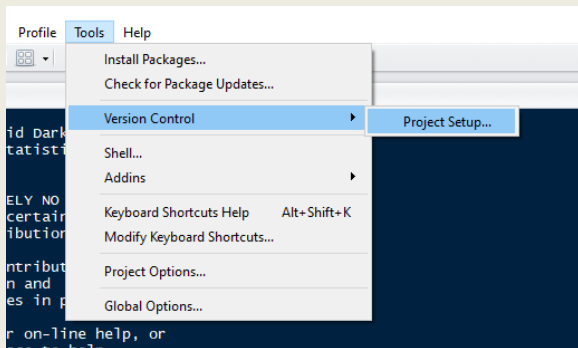
Restart RStudio after adding git.exe path.
Edit .gitconfig in your user directory with your name and email.

Setting up version control for an R project stored locally on your computer

- For analysis etc. that you don't want to share with the world (yet)



Remaining
dialogs to set up
the project

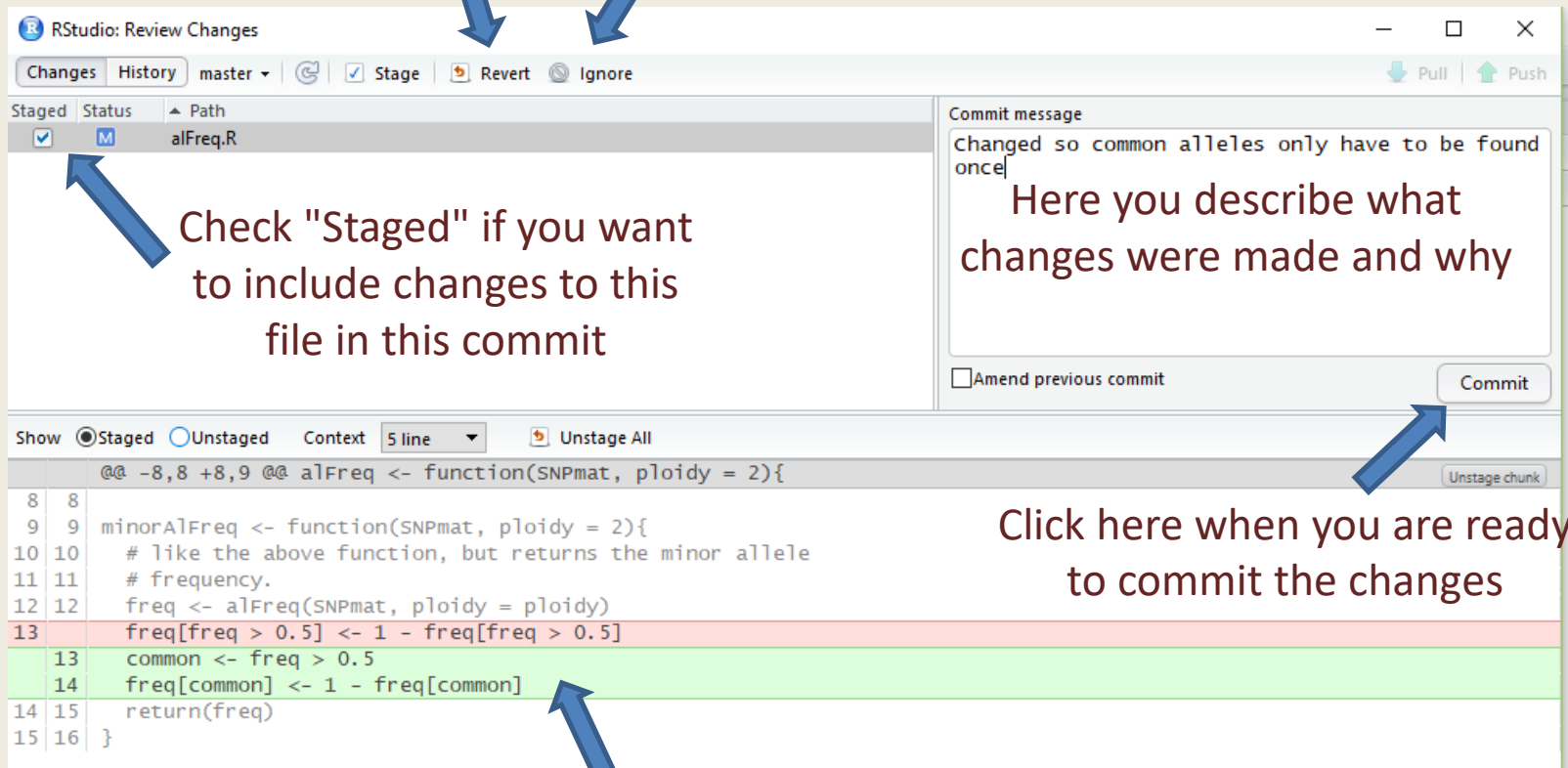


Reviewing changes and making commits

 ← This button

To "undo" changes

To stop tracking changes to a file



The screenshot shows the RStudio 'Review Changes' interface. At the top, there are tabs for 'Changes' and 'History', and a dropdown for the current branch, 'master'. Below these are buttons for 'Stage', 'Revert', and 'Ignore'. On the right side, there are 'Pull' and 'Push' buttons. The main area is divided into two panes. The left pane, titled 'Staged', 'Status', and 'Path', shows a list of files. The file 'alFreq.R' is listed with a checked box in the 'Staged' column and a blue 'M' icon in the 'Status' column. A blue arrow points from the text 'Check "Staged" if you want to include changes to this file in this commit' to this checked box. The right pane, titled 'Commit message', contains a text area with the message 'Changed so common alleles only have to be found once' and a larger text area below it with the text 'Here you describe what changes were made and why'. Below the text areas are checkboxes for 'Amend previous commit' and a 'Commit' button. A blue arrow points from the text 'Click here when you are ready to commit the changes' to the 'Commit' button. At the bottom, there is a 'Show' section with radio buttons for 'Staged' (selected) and 'Unstaged', a 'Context' dropdown set to '5 line', and an 'Unstage All' button. Below this is a code editor showing the contents of 'alFreq.R'. The code is as follows:

```
@@ -8,8 +8,9 @@ alFreq <- function(SNPmat, ploidy = 2){
8 8
9 9 minorAlFreq <- function(SNPmat, ploidy = 2){
10 10 # like the above function, but returns the minor allele
11 11 # frequency.
12 12 freq <- alFreq(SNPmat, ploidy = ploidy)
13 13 freq[freq > 0.5] <- 1 - freq[freq > 0.5]
14 14 common <- freq > 0.5
15 15 freq[common] <- 1 - freq[common]
16 16 return(freq)
17 }
```

Line 13 is highlighted in red, and line 14 is highlighted in green. A blue arrow points from the text 'Shows what changes have been made to each file' to the green highlight on line 14.

Shows what changes have been made to each file

"History" to browse through all commits

The screenshot shows a Git web interface. At the top, there are tabs for "Changes" and "History", with "History" selected. Below the tabs, there's a dropdown menu showing "master" and "(all commits)". To the right, there's a search bar and a "Pull" button. The main content area displays a list of commits. The first commit is highlighted in blue. Below the list, there's a section for the selected commit, showing its SHA, Author, Date, Subject, and Parent. The commit message is "Add minorAlFreq function". Below this, there's a link to the file "alFreq.R". The file content is displayed in a code editor, showing R code for the "alFreq" and "minorAlFreq" functions. The code is highlighted with a green background.

Subject	Author	Date	SHA
HEAD -> refs/heads/master Changed so common alleles only have to be found once	Lindsay Clark <dragonzuela04@gmail.com>	2017-11-11	fc1d6bf0
Add minorAlFreq function	Lindsay Clark <dragonzuela04@gmail.com>	2017-11-11	2b1d3cf3
First commit	Lindsay Clark <dragonzuela04@gmail.com>	2017-11-11	9fb4c238

SHA 2b1d3cf3
Author Lindsay Clark <dragonzuela04@gmail.com>
Date 2017-11-11 18:30
Subject Add minorAlFreq function
Parent 9fb4c238

[alFreq.R](#)

[alFreq.R](#) [View file @ 2b1d3cf3](#)

```
@@ -4,4 +4,12 @@ alFreq <- function(SNPmat, ploidy = 2){
4 4 # vector of allele frequencies in the population.
5 5 freq <- colMeans(SNPmat, na.rm = TRUE)/ploidy
6 6 return(freq)
7 }
No newline at end of file
7 }
8
9 minorAlFreq <- function(SNPmat, ploidy = 2){
10 # like the above function, but returns the minor allele
11 # frequency.
12 freq <- alFreq(SNPmat, ploidy = ploidy)
13 freq[freq > 0.5] <- 1 - freq[freq > 0.5]
14 return(freq)
15 }
```

Mini exercise

- Open a project from a previous lab or lecture
- Add version control using Git
- Make a couple commits and then view them in history



GitHub

- Free website where you can host your Git repositories
- Good for sharing code, collaborating
- Social networking features, wikis, web hosting
- Can give other people feedback on their software
- Integrates with RStudio

Setting up a GitHub repository

If you know you will want to host a repository on GitHub, it is easiest to set it up on GitHub then import it to RStudio, rather than the other way around. (see happygitwithr.com)

On github.com:

GitHub navigation bar: Pull requests, Issues, Marketplace, Explore, + button, user profile.

Create a new repository

A repository contains all the files for your project, including the revision history.

Owner: lvclark / Repository name: testproj ✓

Description (optional): A test project for class

Public (selected): Anyone can see this repository. You choose who can commit.

Private: You choose who can see and commit to this repository.

☒ Initialize this repository with a README

This will let you immediately clone the repository to your computer. Skip this step if you're importing an existing repository.

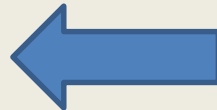
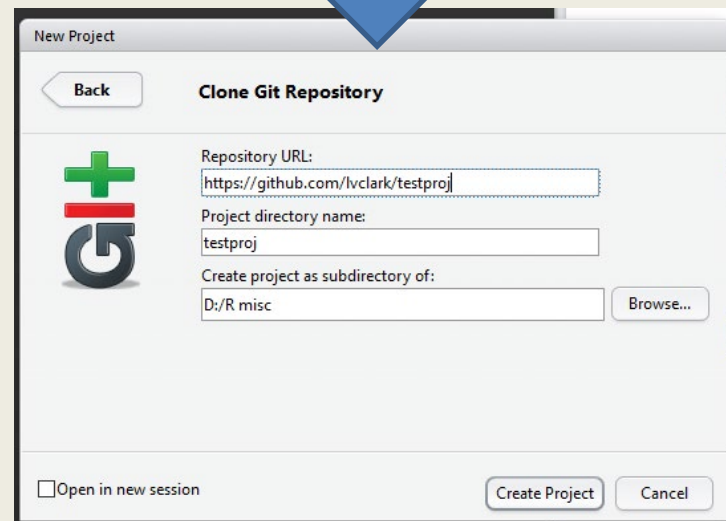
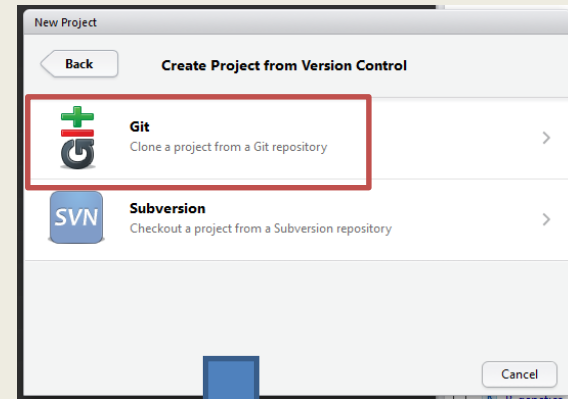
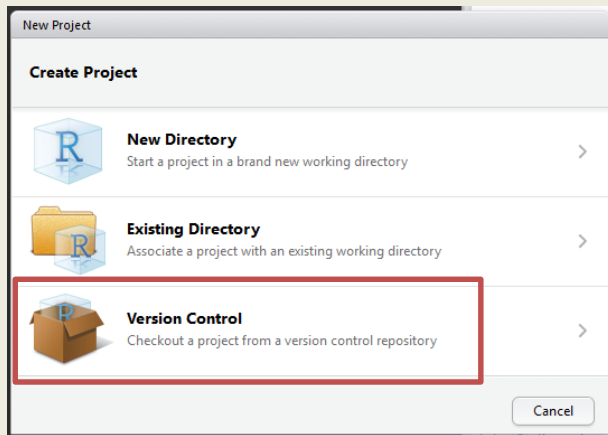
Add .gitignore: None | Add a license: GNU General Public License v3.0 ⓘ

Create repository

+ button to bring up this dialog

GPL 3 is good for scientific software.
You are also not liable if your software causes any harm.

Bringing your GitHub repository into RStudio

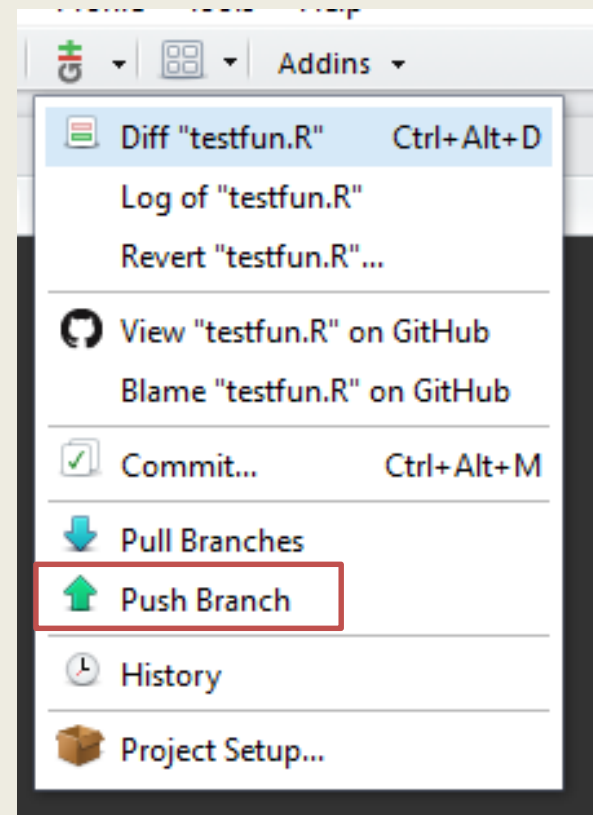
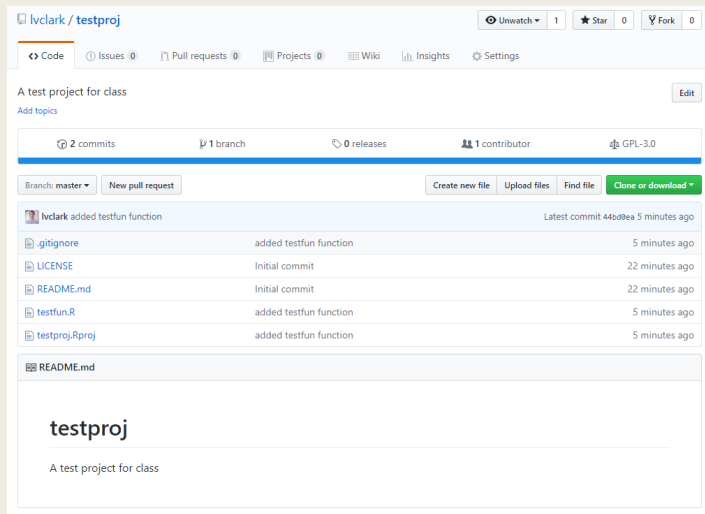


Now you can add and edit files and directories and make commits as demonstrated earlier

To sync your changes back to GitHub

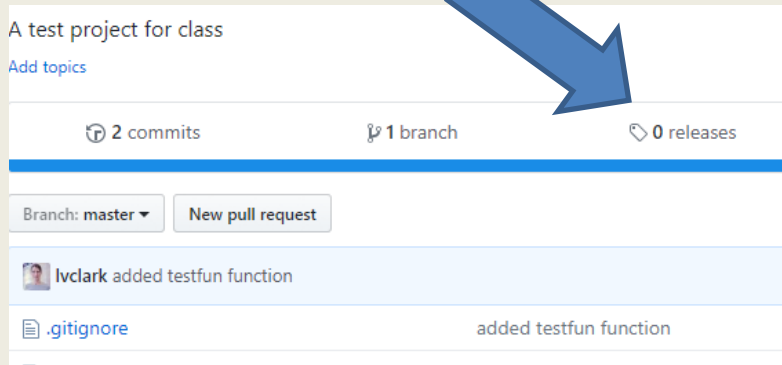
- Make a commit in RStudio as demonstrated earlier
- "Push" the changes

Now you should see your changes on the website



Making a release

- A snapshot of the repository at a given point in time
- Generally a fully-working version of the software
- If the current version isn't working, users can always download last release instead
- To turn in final project, make a release and send me the link



A screenshot of the GitHub "Create new release" form. The form includes a version number "v0.0", a target branch "master", a release title "early development release of testproj", a description "This release is just to test out making a release", and buttons for "Publish release" and "Save draft".

Semantic versioning

- E.g. 2.14.1, 1.0, 0.2-1
- First number is MAJOR version
 - 0 when software still in development
 - 1 for first stable release to public
 - Increase further when there are backwards-incompatible changes
- Second number is MINOR version
 - Starts at 0, increase when new functionality added
- Third number is PATCH version
 - Starts at 0, increase for bug fixes or code optimization

Getting a DOI for a release

- See <https://guides.github.com/activities/citable-code/>
- Link to Zenodo account
- Then every time you make a release of your repository in GitHub, it gets permanently archived on Zenodo (hosted at CERN)

The Zenodo logo, consisting of the word "zenodo" in white lowercase letters on a blue rectangular background.


Using GitHub for collaboration

Reporting bugs and problems

BenLangmead / bowtie2

Watch 30 Star 196 Fork 69







Code Issues 35 Pull requests 7 Projects 0 Wiki Insights

 **Want to submit an issue to BenLangmead/bowtie2?** [Dismiss](#)

If you have a bug or an idea, browse the open issues before opening a new one. You can also take a look at the [Open Source Guide](#).

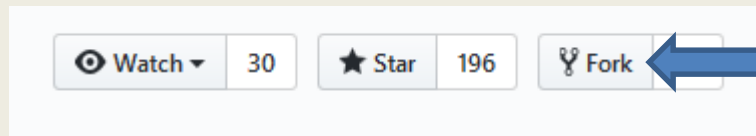
Filters Labels Milestones [New issue](#)

35 Open ✓ 146 Closed Author Labels Projects Milestones Assignee Sort

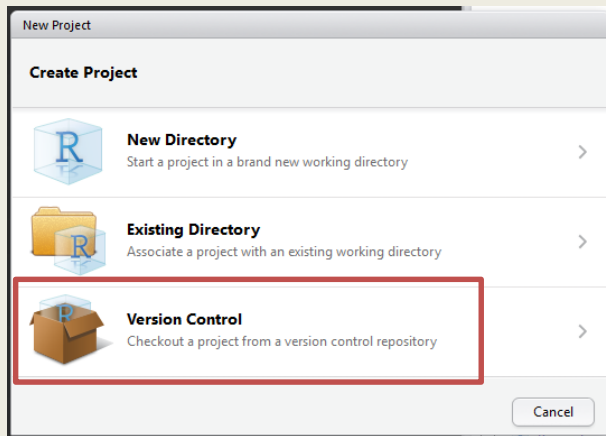
-  **bowtie2 runs in one core despite the --threads and despite the TBB libraries properly installed** 9
#209 opened 12 days ago by anagtz
-  **TLEN is negative for both reads in a properly aligned read pair**
#208 opened 14 days ago by brianyu2010
-  **Adding gcc-8 case on Travis CI.**
#207 opened 27 days ago by junaruga
-  **(ERR): bowtie2-align died with signal 11 (SEGV) (core dumped)** 1
#206 opened on Sep 7 by a-kroh
-  **Fully masked chromosome with Ns gets merged to the previous chromosome** 1
#205 opened on Sep 6 by ccwang002
-  **TLEN = 0 in some cases of soft clipped reads** 3
#203 opened on Aug 21 by mcjmigdal

Better option than emailing software creator. Can see if your issue was already reported.

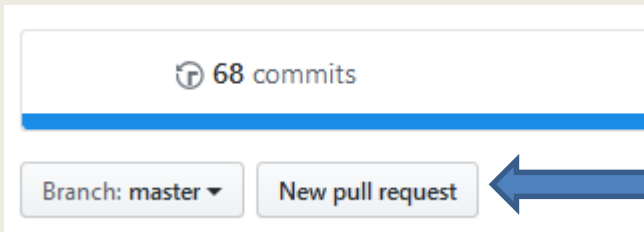
Suggesting changes to someone else's software



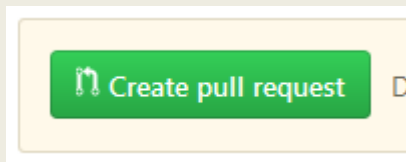
First click **"fork"** to make your own copy



Then **clone** your fork to your computer using RStudio (or Git clone). Make any changes you want, make **commits**, and **push** back to GitHub.

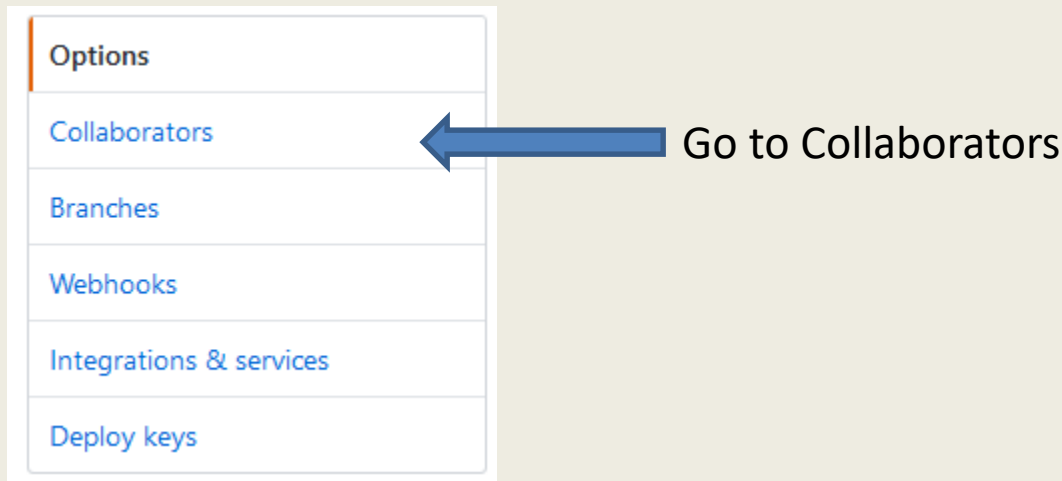
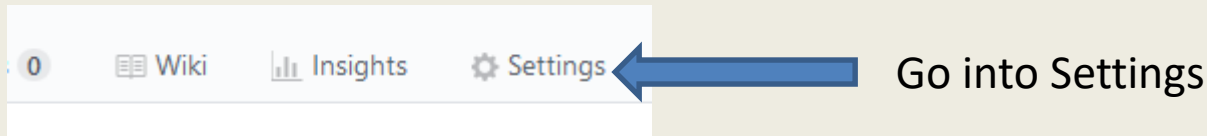


Go back to your fork and start a **pull request**

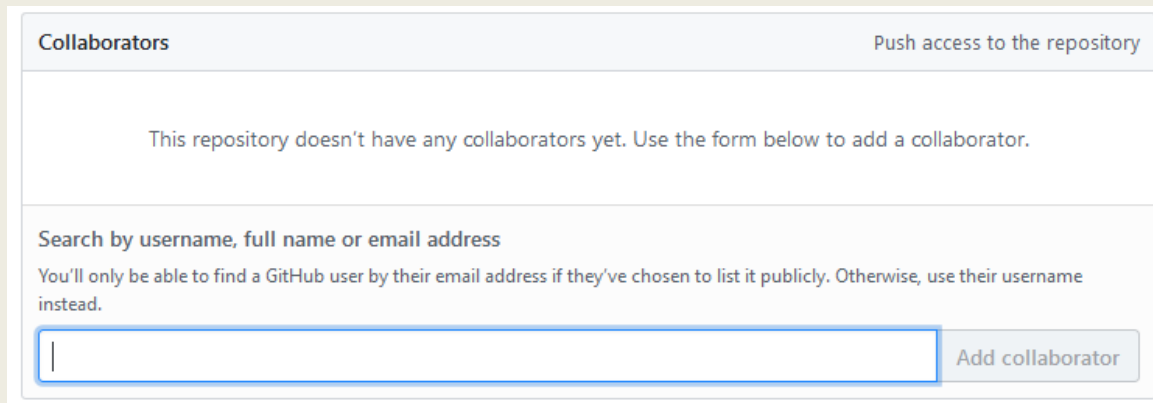


Check that you are using the right fork, and create the pull request. Include a description.

Collaborating on a repository



I recommend doing this for group work on the final projects.



A screenshot of the GitHub 'Collaborators' page. The page title is 'Collaborators' with a link to 'Push access to the repository'. The main content area states: 'This repository doesn't have any collaborators yet. Use the form below to add a collaborator.' Below this is a search instruction: 'Search by username, full name or email address' and a note: 'You'll only be able to find a GitHub user by their email address if they've chosen to list it publicly. Otherwise, use their username instead.' There is a text input field and an 'Add collaborator' button.

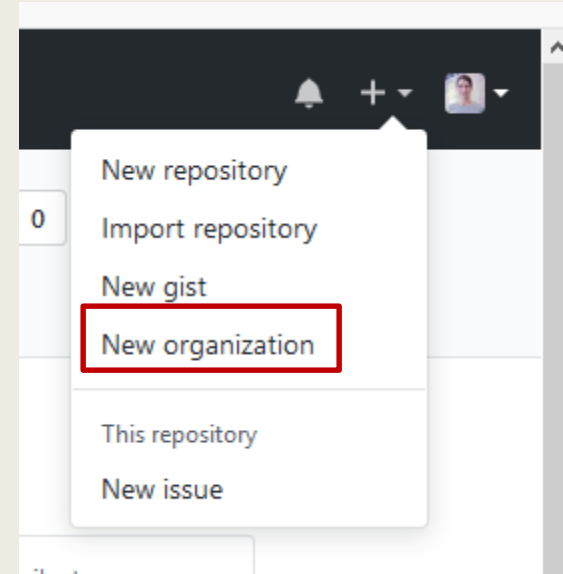
Add usernames.

These people now can **push** to the repository the same as you can.

Remember to **pull** before you edit to avoid conflicts!

GitHub Organizations

- For if you have multiple repositories that you want to work on with the same group of people
- Repositories will belong to the organization rather than any one individual
- Can organize people into teams with different privileges



Might make an organization for your lab, etc.

Tools for incorporating R code
and output into documents

R Markdown format

- Available with "knitr" package
- Based on Markdown
- Quick and easy to write
- Renders nicely to HTML
- .Rmd extension

Graphics in R

One of the major draws of the R programming language is its graphical capabilities. As you prepare manuscripts for publication, you may often find that R is the best choice of (free) software for producing publication-quality figures, even if you did the analysis using some other software.

In this lab we'll use graphical functions available with the standard installation of R, which will give you some opportunities to practice indexing and other skills that we learned earlier in class. There is also a popular package called ggplot2 for making plots in R, which I recommend looking into on your own time.

A basic scatter plot

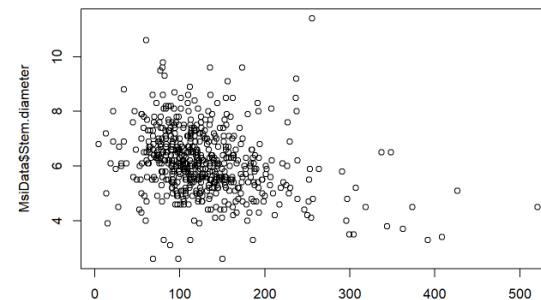
We'll read in some data to practice on. Each row is a genotype (AKA an individual or clone) of the grass *Miscanthus sinensis*, and the columns have some data about those genotypes, including genetic group, latitude and longitude of collection for wild plants, and the least-squared means of some phenotypes measured in a field trial.

```
MsiData <- read.csv("Miscanthus_sinensis_yield.csv", header = TRUE,  
  stringsAsFactors = FALSE)  
str(MsiData)
```

```
## 'data.frame': 580 obs. of 7 variables:  
## $ Genetic.group : chr "S Japan" "SE China/tropical" "US naturalized" "S Japan" ...  
## $ Latitude : num NA 23.7 38.8 NA 44.2 ...  
## $ Longitude : num NA 121 -77 NA 142 ...  
## $ Biomass.yield : num 2436.3 81.9 613.5 2388.9 530.3 ...  
## $ Plant.height : num 292 165 163 272 184 ...  
## $ Number.of.stems : num 144.4 24.7 97.2 143.4 94.2 ...  
## $ Stem.diameter : num 8.3 5.9 4.6 7 5.3 6.1 4.8 5.6 6.8 6 ...
```

Let's explore some of this data with a scatter plot. If we call the `plot` function with integer or numeric vectors for both the first and second arguments, we get a scatter plot, with the first argument on the x-axis.

```
plot(MsiData$Number.of.stems, MsiData$Stem.diameter)
```



Sweave format

- Available with base R installation
- Based on LaTeX
- Makes very nicely typeset PDFs
- Renders mathematical equations well
- Lots of LaTeX utilities for citations, tables, figures, links within document, etc.
- .Rnw extension

4 Getting Started: A Tutorial

4.1 Creating a dataset

As with any genetic software, the first thing you want to do is import your data. For this tutorial, go into the “extdata” directory of the polysat package installation, and find a file called “GeneMapperExample.txt”. Open this file in a text editor and inspect its contents. This file contains simulated genotypes of 300 diploid and tetraploid individuals at three loci. Move this text file into the R working directory. The working directory can be changed with the `setwd` function, or identified with the `getwd` function:

```
> getwd()
```

```
[1] "D:/GitHub/polysat/vignettes"
```

Then read the file using the `read.GeneMapper` function, and assign the dataset a name of your choice (`simgen` in this example) by typing:

```
> simgen <- read.GeneMapper("GeneMapperExample.txt")
```

The dataset now exists as an object in R. The following commands display, respectively, some basic information about the dataset, the sample and locus names, a subset of the genotypes, and a list of which genotypes are missing.

```
> summary(simgen)
```

```
Dataset with allele copy number ambiguity.  
Insert dataset description here.  
Number of missing genotypes: 5  
300 samples, 3 loci.  
1 populations.  
Ploidies: NA  
Length(s) of microsatellite repeats: NA
```

```
> Samples(simgen)
```

Uses for R markdown and Sweave

- Software tutorials
- Reproducible analysis; can write entire paragraphs about what you are doing
- Dynamically-generated reports; say if you want to make some figures and an explanation of what they mean, and also store the code with it

Documentation used in R packages

All user-level functions in an R package must have an Rd file

- These are used to generate the help pages
- Use the `prompt` function to create a skeleton Rd file for a function in your global environment
- Open the file in RStudio and fill in the sections
- Hit "Preview" to see what it will look like in RStudio after the package is installed

Anatomy of an Rd file

```
1 \name{testfun}
2 \alias{testfun}
3 \title{
4 Add One to a Number
5 }
6 \description{
7 This function takes a number and adds one.
8 }
9 \usage{
10 testfun(x)
11 }
12 % this is a comment
13 \arguments{
14   \item{x}{
15 A number.
16 }
17 }
18 \details{
19 This function is pretty self-explanatory and doesn't need
20 more detail, but if it did this is where it would go.
21 }
22 \value{
23 The value of \code{x}, plus one.
24 }
25 \references{
26 Podunk, B. B. and Podunk, B. J. (2017) Adding one to stuff.
27 \emph{Nature} \bld{551}, 13507--13508.
28 }
29 \author{
30 Lindsay Clark
31 }
32
33 \seealso{
34 \code{\link{Arithmetic}}
35 }
36 \examples{
37 testfun(5)
38 }
39 \keyword{arith }% use one of RShowDoc("KEYWORDS")
40
41
```

- `\xxxx{yyyy}` or `\xxxx{yyyy}{zzzz}` to delineate sections and format text
- `%` to indicate comments
- Examples should be executable without needing anything else in global environment.
- Can put multiple functions into one Rd file: use multiple `\alias` sections
- See "Writing R documentation files" in the "Writing R Extensions" manual

Package vignettes

- For when the package needs a tutorial above and beyond the info in the Rd files
- A package can have any number of vignettes (including zero)
- All code in the vignette must be executable (no errors) for the package to pass its check
- R can extract all the code from a vignette into a script
- Can use R markdown or Sweave

Midterm extra credit

- Contribute to study guide:
<https://github.com/lvclark/CPSC499-Fall2018-studyguide>
- Give me your GitHub username on Thursday (or sooner) and will add you
- One point per commit, up to five points
- (Edit the study guide before the exam, but not during!!)