

Git

What is version control?

A version control system is a tool that **manages changes** made to the files and directories in a project. Many version control systems exist; this lesson focuses on one called Git, which is used by many of the data science tools covered in our other lessons. Its strengths are:

- **Keep track of changes to files:** Nothing that is saved to Git is ever lost, so you can always go back to see which results were generated by which versions of your programs.
- **Notice conflicts between changes made by different people:** Git automatically notifies you when your work conflicts with someone else's, so it's harder (but not impossible) to accidentally overwrite work.
- **Synchronize files between different computers:** Git can synchronize work done by different people on different machines, so it scales as your team does.

Version control isn't just for software: books, papers, parameter sets, and anything that changes over time or needs to be shared can and should be stored and shared using something like Git.

Where does Git store information?

Each of your Git projects has two parts:

1. the files and directories that you create and edit directly, and
 2. the extra information that Git records about the project's history.
- The combination of these two things is called a **repository**.
 - Git stores all of its extra information in a directory called `.git` located in the **root directory of the repository**. Git expects this information to be laid out in a very precise way, so you should never edit or delete anything in `.git`.

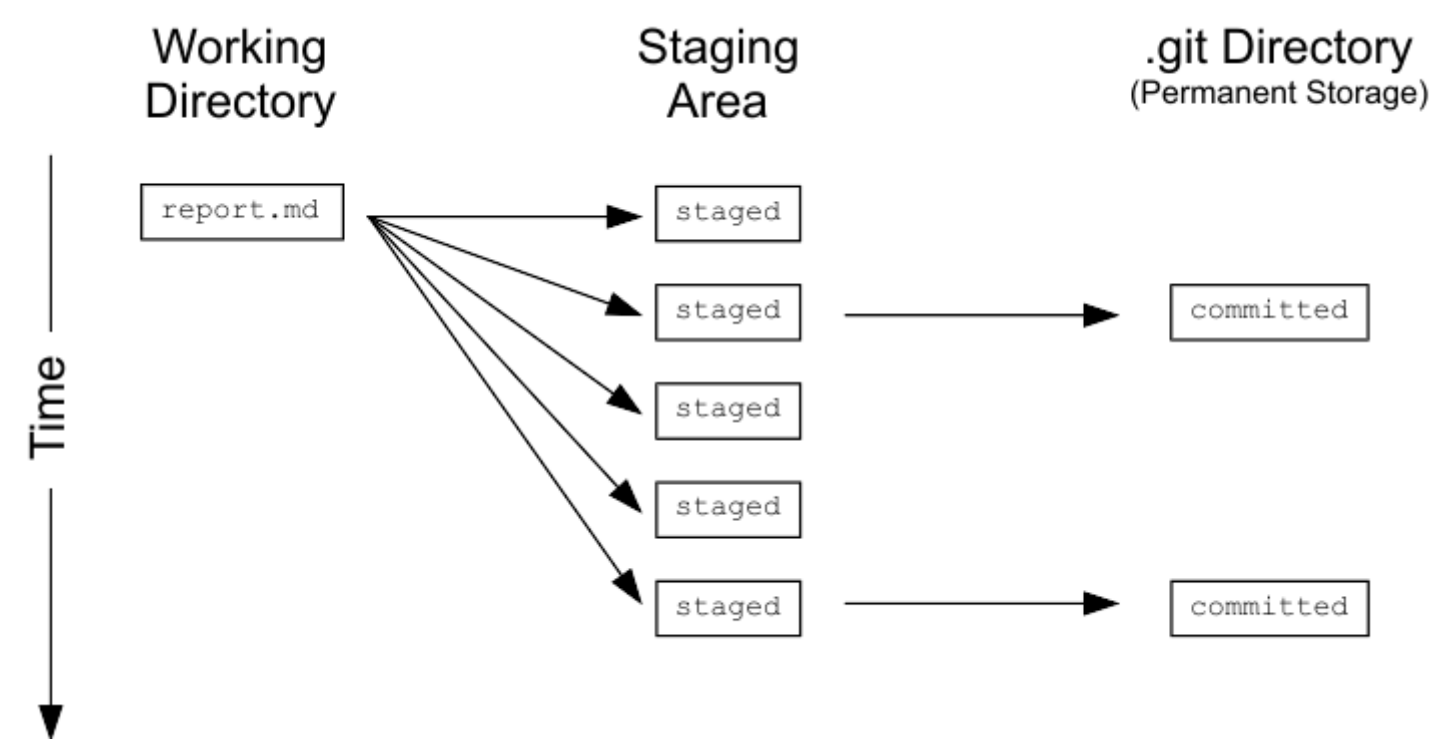
Suppose your home directory `/home/repl` contains a repository called `dental`, which has a sub-directory called `data`. Where is information about the history of the files in `/home/repl/dental/data` stored?

Answer:

`/home/repl/dental/.git`

How can I tell what I have changed?

Git has a **staging area** in which it **stores files with changes you want to save that haven't been saved yet**. Putting files in the staging area is like putting things in a box, while **committing** those changes is like putting that box in the mail: you can add more things to the box or take things out as often as you want, but once you put it in the mail, you can't make further changes.



Desktop programming tools like RStudio can turn diffs like this into a more readable side-by-side display of changes; you can also use standalone tools like DiffMerge or WinMerge.

How can I edit a file?

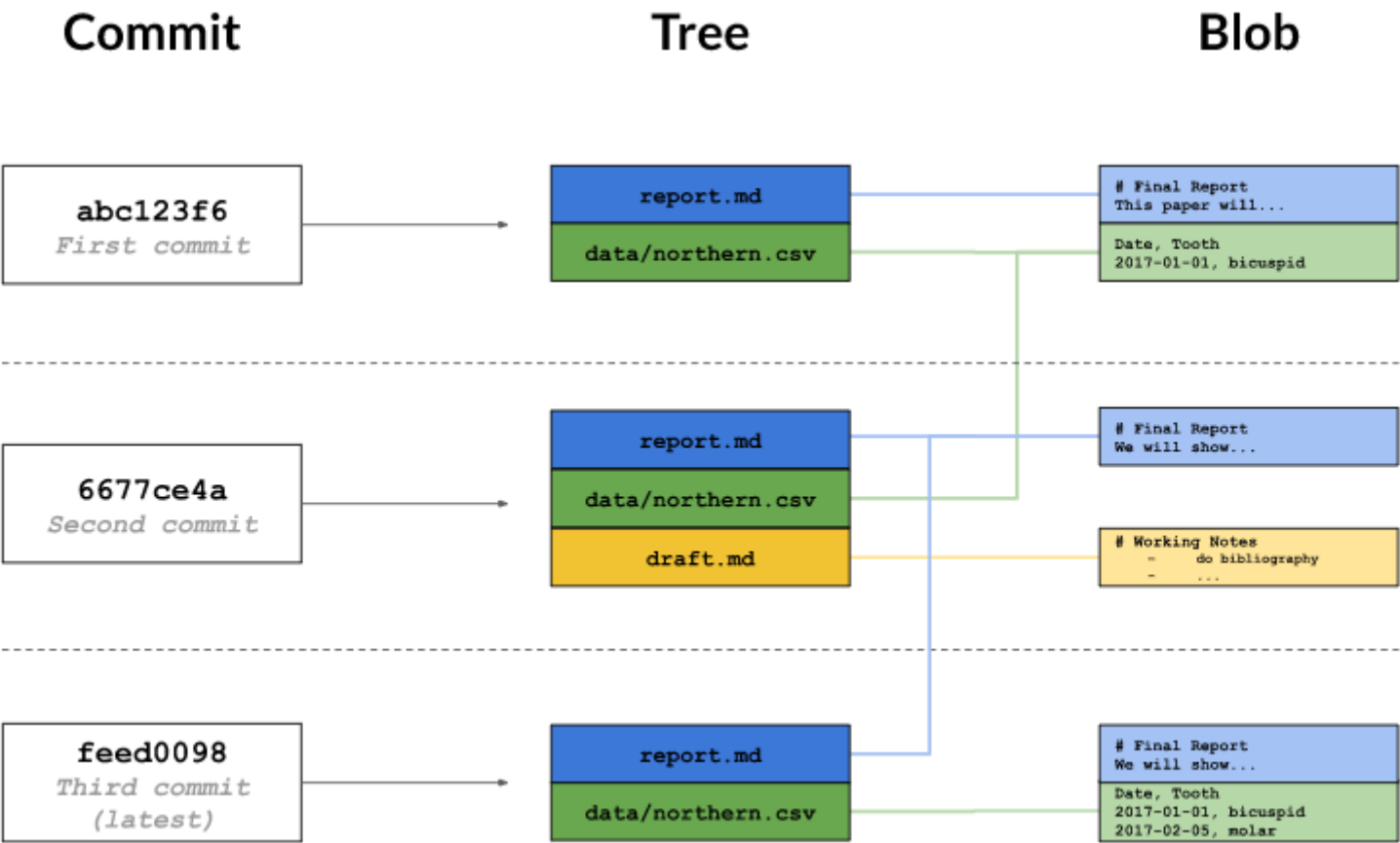
- one Unix text editor s called Nano
- if you type `nano filename`, it will open `filename` for editing (or create it if it doesn't already exist). You can then move around with the arrow keys, delete characters with the backspace key, and so on. You can also do a few other operations with control-key combinations:

- Ctrl-K: delete a line.
- Ctrl-U: un-delete a line.
- Ctrl-O: save the file ('O' stands for 'output').
- Ctrl-X: exit the editor.

How does Git store information?

Git uses a three-level structure for this.

1. A **commit** contains **metadata** such as the **author**, the **commit message**, and the **time** the commit happened. In the diagram below, the most recent commit is at the bottom (`feed0098`), underneath its parent commits.
2. Each commit also has a **tree**, which tracks the **names and locations** in the repository when that commit happened. In the oldest (top) commit, there were two files tracked by the repository.
3. For each of the files listed in the tree, there is a **blob**. This contains a compressed **snapshot of the contents of the file when the commit happened** (blob is short for *binary large object*, which is a SQL database term for "may contain data of any kind"). In the middle commit, `report.md` and `draft.md` were changed, so the blobs are shown next to that commit. `data/northern.csv` didn't change in that commit, so the tree links to the blob from the previous commit. Reusing blobs between commits help make common operations fast and minimizes storage space.



<ul style="list-style-type: none">You can restrict the results to a single file or directory	
<div><div><div>Saving changes</div><div><div><div>1. add file to staging area</div><div>2. commit everything in staging area</div><div><div><div></div><div></div></div><div>log message required to commit changes</div><div><div></div><div></div></div><div>can ammend</div><div><div><div></div><div></div></div><div>Git does not track files by default, it waits until you have used <code>git add</code> at least once</div><div><div></div><div></div></div><div>the untracked files won't have a blob, and won't be listed in a tree, i.e. won't benefit from version control</div><div><div></div><div></div></div><div>use git status to check files that are in your repository but aren't (yet) being tracked.</div><div><div></div><div></div></div><div>good practice to save periodically</div></div></div></div></div></div></div>	<div><div>cd dental</div><div>(git status)</div><div>git add report.txt</div><div>git commit -m "Adding a reference"</div><div></div><div>git commit --amend - m "new message"</div></div>
<div><div><div>Edit a file</div><div><div><div></div><div></div></div><div>when run git commit without message, git launches a text editor</div></div></div></div>	<div><div>nano names.txt</div><div>(type info)</div><div>(Ctrl-O to write the file out and save what is written)</div><div>(Ctrl-X and Enter to exit the editor)</div></div>
<div>Repositories</div>	
<div><div><div>View a repository's history</div><div><div><div></div><div></div></div><div><div>Log entries are shown most recent first</div><div>The commit line displays a unique ID for the commit called a hash</div></div></div></div></div>	<div><div>git log (path)</div><div>(Press the space bar to go down a page or the 'q' key to quit.)</div></div>
<div><div><div>View a specific commit</div><div><div><div></div><div></div></div><div>git log + git diff</div></div></div></div>	<div><div>git show 0da2f7</div><div>(q to quit)</div><div></div><div>for relative path - commit made just before the most recent one:</div><div>git show HEAD~1</div></div>
<div><div><div>See who changed what in a file</div><div><div><div><div></div><div></div></div><div>Each line contains five elements,</div><div><div><div>1. The first eight digits of the hash, <code>04307054</code> .</div><div>2. The author, <code>Rep Loop</code> .</div><div>3. The time of the commit, <code>2017-09-20 13:42:26 +0000</code> .</div><div>4. The line number, <code>1</code> .</div><div>5. The contents of the line, <code># Seasonal Dental Surgeries (2017) 2017-18</code> .</div></div></div></div></div></div></div>	<div><div>git annotate report.txt</div></div>
<div><div><div>See what changed between two commits</div></div></div>	<div><div>git diff ID1..ID2</div><div>git diff HEAD~1..HEAD~3</div></div>
<div><div><div>Tell Git to ignore certain files</div><div><div><div></div><div></div></div><div>for temporary or intermediate files that you don't want to save</div></div></div></div>	<div><div><div><div><div></div><div></div></div><div>by creating a file in the root directory of your repository called <code>.gitignore</code> and storing a list of wildcard patterns that specify the files you don't want Git to pay attention to</div><div><div></div><div></div></div><div>to ignore any file or directory called <code>build</code> (and, if it's a directory, anything in it), as well as any file whose name ends in <code>.mpl</code> . THEN <code>.gitignore</code> should contain:</div></div></div><div><div>build</div><div>*.mpl</div><div>pdf (should be *.pdf)</div><div>*.pyc -> ignore bin/analyze.pyc</div><div>backup -> ignore backup/northern.csv</div></div></div>
<div><div><div>Remove unwanted (untracked) files</div><div><div><div></div><div></div></div><div>git clean -n will show you a list of files that are in the repository, but whose history Git is not</div></div></div></div>	<div><div>cd dental</div><div>git status</div><div>git clean -f</div></div>

<div>currently tracking</div> <ul style="list-style-type: none">• git clean -f will then delete those files.• <code>git clean</code> only works on untracked files• ls to list the files in your current working directory. backup.log should no longer be there!	<div>>removing backup.log</div> <div>ls</div> <div>>bin data report.txt results</div>
<div>See how Git is configured</div> <div><code>git config --list</code> with one of three additional options:</div> <ul style="list-style-type: none">• <code>--system</code> : settings for every user on this computer.• <code>--global</code> : settings for every one of your projects.• <code>--local</code> : settings for one specific project. <div>Each level overrides the one above it, so local settings (per-project) take precedence over global settings (per-user), which in turn take precedence over system settings (for all users on the computer).</div>	<div>git config --list --local</div>
<div>Change my Git configuration</div> <ul style="list-style-type: none">• there are two you should set on every computer you use: your name and your email address• sed to identify the authors of a project's content in order to give credit (or assign blame, depending on the circumstances).	<div>git config --global user.email rep.loop@datacamp.com</div> <div>To change the email address (user.email) configured for the current user for all projects to rep.loop@datacamp.com.</div>
<div>Undo</div>	
<div>Unstage a file you shouldn't have</div>	<div>git reset HEAD</div>
<div>Undo changes to unstaged files</div> <ul style="list-style-type: none">• <code>--</code> must be there to separate• once you discard changes in this way, they are gone forever.	<div>git checkout -- data/northern.csv</div>
<div>Undo changes to staged files</div> <ul style="list-style-type: none">• step 1: unstage the file• step 2: undo changes since the last commit	<div>git reset HEAD path/to/file</div> <div>git checkout -- path/to/file</div>
<div>Restore an old version of a file</div> <ul style="list-style-type: none">• restore versions of that file from a commit• committing as saving your work, checking out as loading that saved version• takes two arguments: the hash that identifies the version you want to restore, and the name of the file.• Restoring a file doesn't erase any of the repository's history. Instead, the act of restoring the file is saved as another commit, because you might later want to undo your undoing.• Passing <code>-</code> followed by a number restricts the output to that many commits.• cat print the content of a file onto the standard output stream	<div>cat data/western.csv</div> <div>git log -2 data/western.csv</div> <div>git checkout 2242bd data/western.csv</div> <div>cat data/western.csv</div> <div>(to display the updated contents)</div> <div>git commit -m "restored"</div>
<div>Undo all of the changes I have made</div> <ul style="list-style-type: none">• give <code>git reset</code> a directory. For example, <code>git reset HEAD data</code> will unstage any files from the <code>data</code> directory• if you don't provide any files or directories, it will unstage everything• <code>HEAD</code> is the default commit to unstage, so you can simply write <code>git reset</code> to unstage everything	<div>git reset</div> <div>To remove all files from the staging area</div> <div>git checkout -- .</div> <div>To put those files back in their previous state</div>

- `git checkout -- data` will then restore the files in the `data` directory to their previous state.
- So `git checkout -- .` will revert all files in the current directory.

What is a branch?

- allows you to have multiple versions of your work, and lets you track each version systematically.
- instead of creating different subdirectories to hold different versions of your project in different states which could cause confusion on which is the right version of each file in the right subdirectory, and you risk losing work.
- one branch do not affect other branches (until you **merge** them back together).
- Branches are the reason Git needs both trees and commits: a **commit will have two parents** when **branches are being merged**.
 - *blobs* for files, *trees* for the **saved states of the repositories**, and *commits* to record the changes
- By default, every Git repository has a branch called `master`

What are conflicts?

- for example, bug fixes might touch the same lines of code, or analyses in two different branches may both append new (and different) records to a summary data file.
- In this case, Git relies on you to reconcile the conflicting changes.
- Example
 - file `todo.txt`

A) Write report.
B) Submit report.

- create a branch called `update` and modify the file to be:

A) Write report.
B) Submit final version.
C) Submit expenses.

- switch back to the `master` branch and delete the first line, so that the file contains:

B) Submit report.

When you try to merge `update` and `master` , what conflicts does Git report?

- B) Submit report.
- +A) Write report.
- +B) Submit final version.
- +C) Submit expenses.
 - Git can merge the deletion of line A and the addition of line C automatically.

How can I find out where a cloned repository originated?

- Git remembers where the original repository was. It does this by storing a **remote** in the new repository's configuration. A remote is like a browser bookmark with a name and a URL.
- If you use an online git repository hosting service like GitHub or Bitbucket, a common task would be that you clone a repository from that site to work locally on your computer. **Then the copy on the website is the remote.**

Pull and Push changes from a remote repository?

- remote repository is often a repository in an online hosting service like GitHub.
- A typical workflow is that you pull in your collaborators' work from the remote repository so you have the latest version of everything, do some work yourself, then push your work back to the remote so that your collaborators have access to it.

Branch	
List all of the branches in a repository <ul style="list-style-type: none">• The branch you are currently in will be shown with a <code>*</code> beside its name.	git branch
View the differences between branches <ul style="list-style-type: none">• same as revisions• to check what conflict Git report before merging	git diff branch-1..branch-2

<div>Switch from one branch to another + Remove file<ul style="list-style-type: none">• same as restoring version using hash• Git will only let you do this if all of your changes have been committed - can get around this• <code>git rm</code> removes the file then stages the removal of that file with <code>git add</code> , all in one step.</div>	<div><code>cd dental</code> <code>git branch</code> <code>#switch to summary-statistics branch</code> <code>git checkout summary-statistics</code> <code>#delete report.txt</code> <code>git rm report.txt</code> <code>git commit -m "Removing report"</code> <code>#check that it's gone</code> <code>ls</code> <code>#switch back to master branch</code> <code>git checkout master</code> <code>ls</code> <code>-report.txt in master branch should not be affected</code></div>
<div>Create a branch<ul style="list-style-type: none">• using <code>-b</code> flag allows you to create a branch then switch to it in one step• The contents of the new branch are initially identical to the contents of the original. Once you start making changes, they only affect the new branch.</div>	<div><code>cd dental</code> <code>git branch</code> <code>git checkout -b deleting-report</code> <code>git rm report.txt</code> <code>git commit -m "report deleted"</code> <code>#to compare the master branch with the new state of the deleting-report branch</code> <code>git diff master..deleting-report</code></div>
<div>Merge two branches<ul style="list-style-type: none">• When you merge source branch into destination source, Git incorporates the changes made to the source branch into the destination branch• If those changes don't overlap, the result is a new commit in the destination branch that includes everything from the source branch• Git automatically opens an editor so that you can write a log message for the merge</div>	<div><code>cd dental</code> <code>git merge source destination</code> <code>Ctrl+O</code> and then <code>Ctrl+X</code> to exit this. If you want to avoid this, use the <code>--no-edit</code> flag in the <code>git merge</code> command.</div>
<div>Merge two branches with conflicts<ul style="list-style-type: none">• running <code>git status</code> after the merge reminds you which files have conflicts that you need to resolve by printing <code>both modified:</code> beside the files' names.• In many cases, the destination branch name will be <code>HEAD</code> because you will be merging into the current branch.• To resolve the conflict, edit the file to remove the markers and make whatever other changes are needed to reconcile the changes, then commit those changes.</div>	<div><code>cd dental</code> <code>git merge alter-report-title master</code> <code>#to check which file has conflict</code> <code>git status</code> <code>#edit text file to reconcile changes</code> <code>nano report.txt</code> <code>Ctrl + K</code> to remove entire lines <code>Ctrl + O</code> to Save then enter <code>Ctrl + X</code> to exit <code>git add report.txt</code> <code>git commit -m "remove lines"</code></div>
<div>Repository</div>	
<div>Create a brand new repository<ul style="list-style-type: none">• One thing you should <i>not</i> do is create one Git repository inside another.• Nested repositories may be used in very large projects but no ideal</div>	<div><code>git init project-name</code></div>
<div>Turn an existing project into a Git repository<ul style="list-style-type: none">• after initializing the folder into a repo, Git immediately notices that there are a bunch of changes that can be staged (and committed afterwards)</div>	<div><code>git init</code> (in the project's root directory) or <code>git init /path/to/project</code></div>
<div>Example - in directory dental with is not yet a Git repository</div>	<div><code>pwd</code> #'present working directory' <code>git init</code> <code>git status</code></div>

	git status
<div>Create a copy of an existing repository</div> <div><ul style="list-style-type: none">When? Join a project that is already running, inherit a project from someone else, or continue working on one of your own projects on a new machineWhen you clone a repository, Git uses the name of the existing repository as the name of the clone's root directory, for example:</div> <div>git clone /existing/project</div> <div>will create a new directory called <code>project</code> inside your home directory.</div> <div>git clone /existing/project newprojectname</div> <div>to add new name</div>	<div>pwd</div> <div>ls</div> <div>git clone /home/thunk/repo dental</div> <div>git clone</div> <div>https://github.com/datacamp/project.git</div>
<div>List names of remotes (cloned repos)</div> <div><ul style="list-style-type: none">-v flag shows the remote's URLs.</div>	git remote -v (for "verbose")
<div>Define remotes</div> <div><ul style="list-style-type: none">When you clone a repository, Git automatically creates a remote called <code>origin</code> that points to the original repository.You can connect any two Git repositories this way, but in practice, you will almost always connect repositories that share some common ancestry.</div>	
<div>Add more remotes</div>	<div>git remote add remote-name URL</div> <div>git remote add thunk /home/thunk/repo</div>
<div>Remove existing remotes</div>	<div>git remote rm remote-name</div>
<div>Pull in changes from a remote repository</div> <div><ul style="list-style-type: none">pull changes from those repositories and push changes to themgets everything in <code>branch</code> in the remote repository identified by <code>remote</code> and merges it into the current branch of your local repositoryget changes from <code>latest-analysis</code> branch in the repository associated with the remote called <code>thunk</code> and merge them into your <code>quarterly-report</code> branch.</div>	<div>git pull remote branch</div> <div>cd dental</div> <div>git pull thunk latest-analysis</div>
<div>Pull with unsaved changes</div> <div><ul style="list-style-type: none">Just as Git stops you from switching branches when you have unsaved workpulling in changes from a remote repository when doing so might overwrite things you have done locallyFix: either commit your local changes or revert them, and then try to pull again.commit all your local changes if you want your git pull to run smoothly</div>	<div>cd dental</div> <div>git pull origin master</div> <div>> overwrite warning</div> <div>#to discard changes in the repo</div> <div>git checkout -- report.txt</div> <div>git pull origin master</div>
<div>Push changes to a remote repository</div> <div><ul style="list-style-type: none">push local changes - contents of your branch <code>branch-name</code> into a branch with the same name in the remote repository associated with <code>remote-name</code>not good practice to use different branch names</div>	<div>cd dental</div> <div>git add data/northern.csv</div> <div>git commit -m "Added more northern data"</div> <div>git push origin master</div>
<div>Push with unsaved changes</div>	<div>git push origin master</div>

Push with unsaved changes

- Git does not allow you to push changes to a remote repository unless you have merged the contents of the remote repository into your own work.
- Fix: bring your repository up to date with origin with git pull

git push origin master
>error: failed to push some refsto '/home/thunk/repo'
hint: Updates were rejected because the tip of your
current branch is behind
hint: its remote counterpart. Integrate the remote
changes (e.g.hint: 'git pull ...') before pushing again.
hint: See the 'Note about fast-forwards' in 'git push --
help' for details.
Ctrl + X
git push origin master