

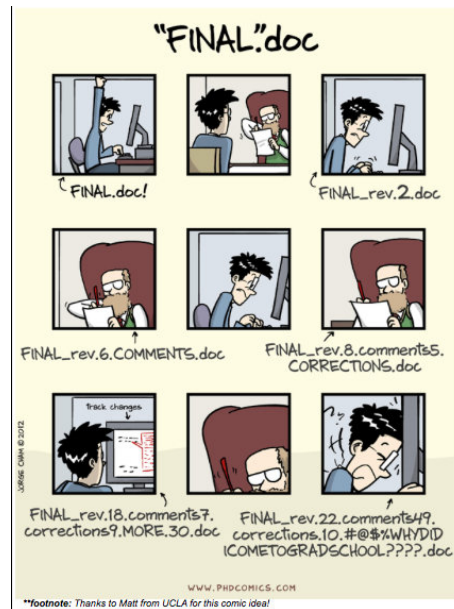
W4995 Applied Machine Learning

Git, Github and Testing

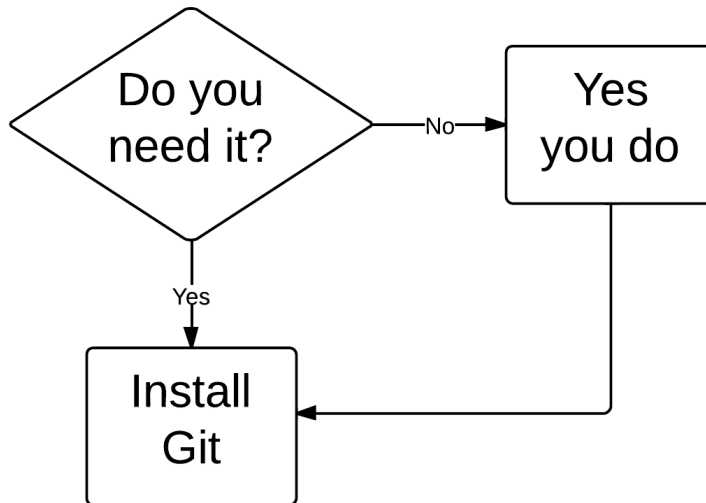
01/28/19

Andreas C. Müller

Why do I care?



Version Control Flowchart



Git and Github

<https://guides.github.com/>

<http://rogerdudler.github.io/git-guide/>





Configuration

```
git config --global user.name "Andreas Mueller"  
git config --global user.email "acm2248@columbia.edu"  
git config --global color.ui "auto"  
git config --global core.editor "vim"
```

Show your configuration:

```
git config --list
```

Creating a repository

```
$ mkdir homework1  
$ cd homework1
```

```
$ git init
```

```
> Initialized empty Git repository in /tmp/homework1/.git/
```

```
$ ls -a
```

```
.  ..  .git
```

```
$ git status
```

```
On branch master
```

```
No commits yet
```

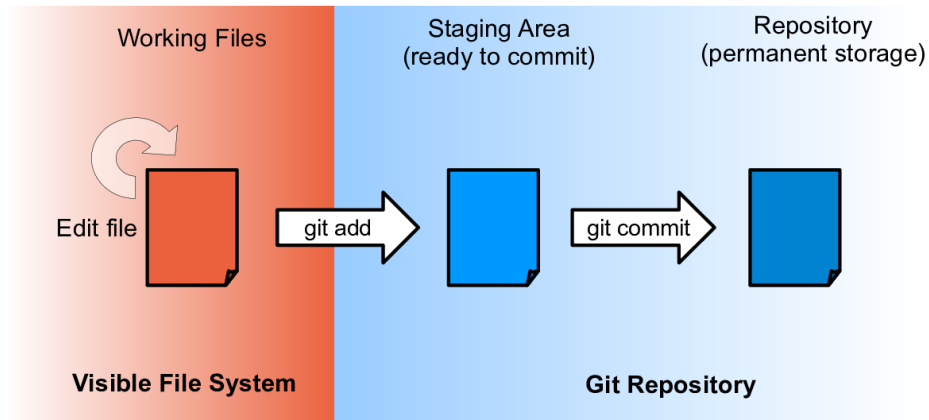
```
nothing to commit (create/copy files and use "git add" to track)
```


Cloning a repository

The screenshot shows the GitHub interface for the repository 'aml-spring-19 / homework-1'. At the top, there are buttons for 'Unwatch', 'Star' (1), and 'Fork' (0). Below this is a navigation bar with 'Code', 'Issues' (0), 'Pull requests' (0), 'Projects' (0), 'Wiki', 'More', and 'Settings'. The main content area shows 'Base Repository for Homework 1' with an 'Edit' button. Below this, it says 'Manage topics'. A summary bar indicates '1 commit', '1 branch', '0 releases', and '1 contributor'. A 'Branch: master' dropdown is visible. On the left, a list of files shows '.gitignore' and 'README.md', both marked as 'Initial commit'. On the right, a 'Clone or download' button is expanded, showing options to 'Clone with SSH' (with a link to 'Use HTTPS') and a text input field containing 'git@github.com:aml-spring-19/homework-1'. Below the input field is a 'Download ZIP' button. The repository name 'homework-1' is displayed prominently, followed by the description 'Base Repository for Homework 1'.

```
$ git clone git@github.com:aml-sprint-19/homework-1
```

Workflow?



Changing a File

```
$ echo "print('Hello world!')" >> task1.py
```

```
$ git status
```

On branch master

No commits yet

Untracked files:

(use "git add <file>..." to include in what will be committed)

task1.py

nothing added to commit but untracked files present (use "git add" to track)

```
$ git add task1.py
```

Viewing your history

```
$ git commit -m "say hello"
```

```
[master (root-commit) 7d139fb] say hello
1 file changed, 1 insertion(+)
create mode 100644 task1.py
```

```
$ git log
```

```
commit 7d139fb317ecfa7d629654b709747e91ecec444 (HEAD -> master)
Author: Andreas Mueller <andreas.mueller@columbia.edu>
Date:   Mon Jan 28 12:37:29 2019 -0500
```

say hello

A note on viewing changes

Changes between working directory and what was last staged

```
git diff
```

Changes between staging area and last commit

```
git diff --staged
```

Referencing different versions

- Shorthand for different versions of a repository (refers to commits)
 - Current Version (most recent commit): HEAD
 - Version before current: HEAD~1
 - Version before that: HEAD~2
- Each of these also has a commit hash
 - use `git log` to get appropriate hash

Exploring History

Changes made in the last commit

```
git diff HEAD~1
```

Changes made in the last 2 commits

```
git diff HEAD~2
```

Changes made since commit hash...

```
git diff 0b0d55e
```

Recovering Older Versions of Files

- Overwrite task1.py:

```
$ echo "print('goodbye, cruel world!')" > task1.py  
$ cat task1.py
```

```
print('goodbye, cruel world!')
```

Recover last recorded version:

```
$ git checkout HEAD task1.py  
$ cat task1.py
```

```
print('Hello world!')
```


Recovering Older Versions of whole repo

```
$ echo "print('goodbye, cruel world!')" > task1.py
$ git add task1.py
$ git commit -m "saying goodbye"
```

```
[master 95f3b39] saying goodbye
1 file changed, 1 insertion(+), 1 deletion(-)
```

```
$ git log
```

```
commit 95f3b39402451c474e3887be94d2cc37a11be511 (HEAD -> master)
Author: Andreas Mueller <andreas.mueller@columbia.edu>
Date: Mon Jan 28 12:44:52 2019 -0500
```

```
    saying goodbye
```

```
commit 7d120fb31733ef7d620654b700747e01cccc444
```

Resetting:

```
$ git reset --hard HEAD~1
```

HEAD is now at 7d139fb say hello

Branches

```
$ git checkout -b "new_feature"
```

Switched to a new branch 'new_feature'

make some changes, add, commit...

Moving between branches:

```
$ git checkout master
```

Switched to branch 'master'

changes are not present in master..

Merge

- Fast-forward merge:

```
* 6cec4ed (HEAD -> another_one) E
* 1e96a3b D
* 513ced1 (master) C
* b5ba00a B
* db928a0 A

/tmp/git_graphs [git::master] [andy@dsi-amueller] [15:38]
> git merge another_one
Updating 513ced1..6cec4ed
Fast-forward
 D | 0
 E | 0
2 files changed, 0 insertions(+), 0 deletions(-)
 create mode 100644 D
 create mode 100644 E

* 6cec4ed (HEAD -> master, another_one) E
* 1e96a3b D
* 513ced1 C
* b5ba00a B
* db928a0 A
```

- Merge-commits:

```
* 43563a5 (HEAD -> master) G
* c3ea8c8 F
* 6cec4ed (another_one) E
* 1e96a3b D
*/
* 513ced1 C
* b5ba00a B
* db928a0 A

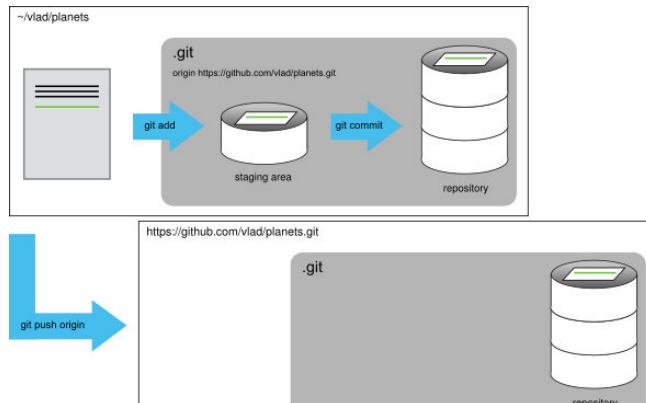
/tmp/git_graphs [git::master] [andy@dsi-amueller] [15:43]
> git merge another_one
Merge made by the 'recursive' strategy.
 D | 0
 E | 0
2 files changed, 0 insertions(+), 0 deletions(-)
 create mode 100644 D
 create mode 100644 E

* d6fedb0 (HEAD -> master) Merge branch 'another_one'
* 6cec4ed (another_one) E
* 1e96a3b D
* 43563a5 G
* c3ea8c8 F
*/
* 513ced1 C
* b5ba00a B
* db928a0 A
```

Collaborating using git & github

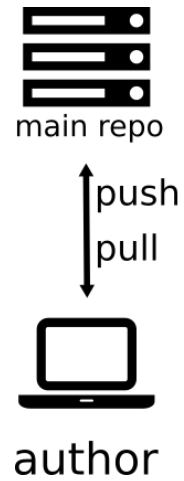
Remote repository

- central location everyone can see
- requires network ?
- github, bitbucket
- public vs private

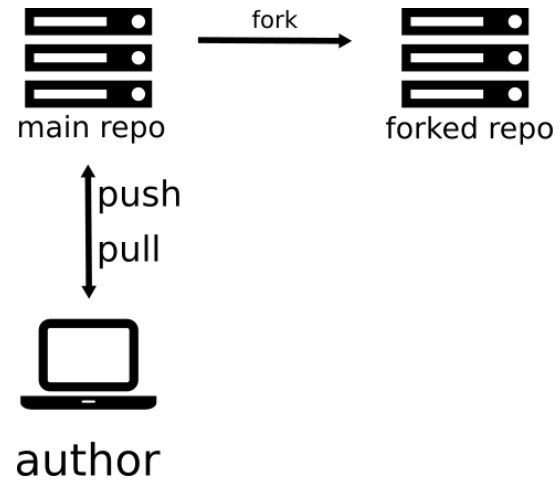


Github

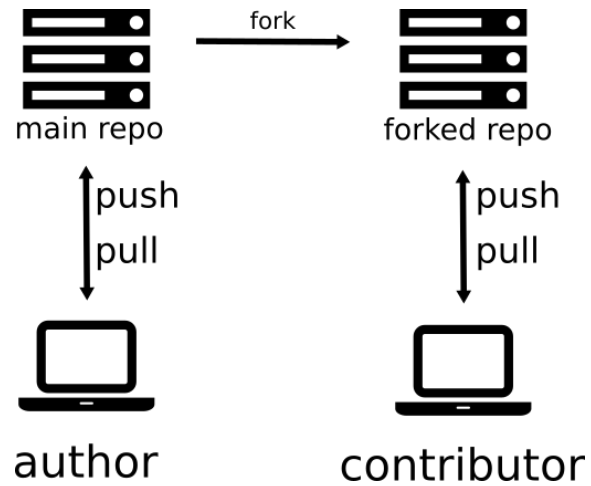
Github pull request workflow



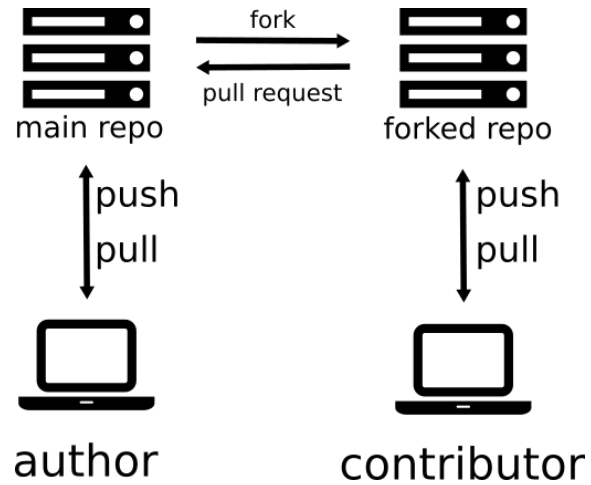
Github pull request workflow



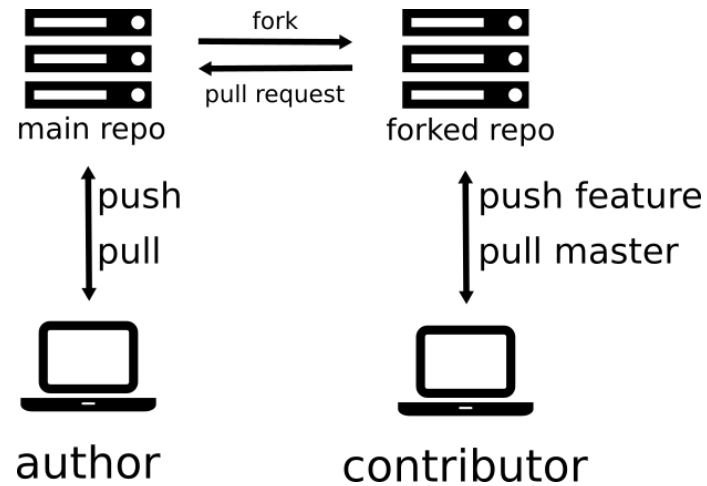
Github pull request workflow



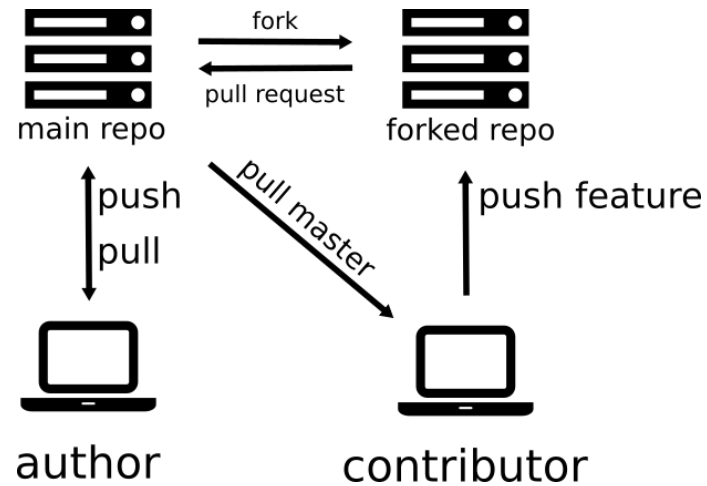
Github pull request workflow



Github pull request workflow



Github pull request workflow



Typical Workflow

- Clone
- Branch
- Add, commit, add, commit, add, commit, ...
- Merge / rebase
- Push

Some tips

- `$ git status`
- Install shell plugins for status and branch (oh-my-zsh)

```
/home/andy/checkout/scikit-learn [git::master *] [andy@dsi-amueller] [16:35]  
> █
```

- Set editor, pager and diff-tool (check out meld!)
- Use `.gitignore`

Git log

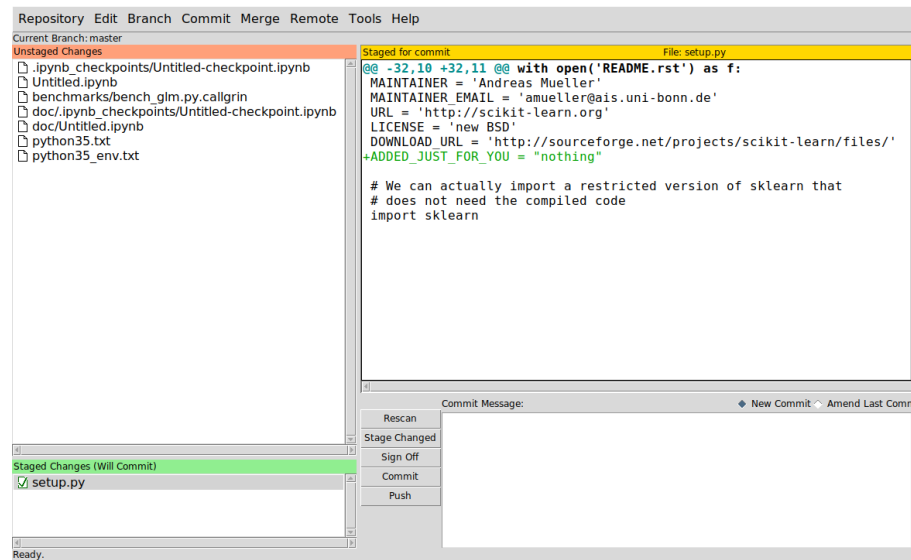
```

/home/andy/checkout/scikit-learn [git:master *] [andy@dsi-umel1eri [43]
$ git log --oneline --decorate --all --graph -n 50
c9868f1 HEAD upstream/master: CI remove obsolete comment
7978119 [MRG-] #2818: in FAQ, link deep learning question to GPU question (#2280)
3b503ce TST/FIX Add check for estimator: parameters not modified by 'fit' (#7746)
906868c TST/FIX Issue #5173 - pass n_neighbors in ML computation (#8181)
a113e11 all sorted on lfu folder path contents (#7548)
a940709 [MRG-3] FIX Memory leak in MAE; Use safe realloc; Acquire GIL only when raising; Propagate all errors to python interpreter level
08772cd [FIX] Ensure coref is an ndarray when fitting Lassosars (#8616)
556469d MNT/BLD Use Github's merge refs to test PRs on CircleCI (#2101)
66443aa [MRG-1] Add prominent mention of Laplacian Eigenmaps (#8155)
b0dda7b [MRG-+ 2] [MAINT] Update to Sphinx-Gallery 0.1.7 (#7986)
a0e4462 [MRG-3] Fixes #8198 - error in datasets, make moons (#8199)
90333ad TST/FIX fix flake8 diff, sh check files (#8208)
e6abef3 DOC add missing parentheses in TfidfTransformer docstring
3c8d82f DOC additional fixes to 20 newsgroups to prevent TypeError (#8204)
c0b059d removed stray space in __main__ (#8203)
c6d70a Upgrade html documentation to JQuery v3.1.1 (#8145)
c46c7b Clarify error message for min_samples_split. (#8167)
2a1408a fixing typo in cs_mse_path deprecation (#8176)
168455e [MRG-+ 1] Fix the cross_val_predict function for methods='predict_proba' (#7889)
491b11c DOC Fix Link (#8172)
898985e Fix Ridge floating point instability (#8154)
2775f91 [MRG-+ 1] Add fowlkes-mallows and other supervised cluster metrics to SCORERS dict so it can be used in hyper-param search (#8171)
8062125 [MRG-] #1 add fowlkes fit to multisetoutput module (#8054)
08e2125 [MRG] FIX Avoid default mutable argument in constructor of AgglomerativeClustering (#8153)
d8ce49e [MRG-2] Avoid failure in first iteration of RANSAC regression (#7914)
0874398 [MRG-1] DOC: complete list of online learners (#8152)
085691e FIX sphinx gallery rendering of plot digits pipe example
8349497 [MRG-1] Deprecate ridge alpha param on SparsePCA.transform() (#8137)
c9c9499 DOC: updating GridSearchCV's n_jobs parameter (#8106)
2c0f6d7 [MRG-1] fowlkes_mallows score: more unit tests (Fixes #8101) (#8140)
543095e [MRG-1] Add DSCM support for additional metric params (#8139)
07777ce [MRG-1] Fix "cite us" link in sidebar (#8142)
288827b [MRG-1] update copyright years for 2017 (#8138)
e2ab0db DOC Fix typo in FAQ (#8132)
380e4b0 DOC Split data using safe_split in permutation test score (#5697)
all4c44 [MRG-1] Catch cases for different class size in MLPClassifier with warm start (#7976) (#8035)
dc72405 (origin/repr_give_up, repr_give_up) i have no idea what I'm doing
c25f4de pep8
8b30325 playing around, then giving up
* 0d1398a (upstream/ignore_lambda_to_diff_errors) MNT Ignore E731: Use a def instead of lambda
//
d97d13e DOC add sklearn-crfutils to related projects (#7878)
fc0786a [MRG-3] Fused types for MultiTaskElasticNet (#8061)
695649c [MRG-1] MAINT Python 3.6 fixes (#8123)
e0808d0 DOC Fix indentation errors and username links (#8121)
473c60e [MRG-2] FIX IsolationForest(max_features=0.8).predict(X) fails input validation (#7575)

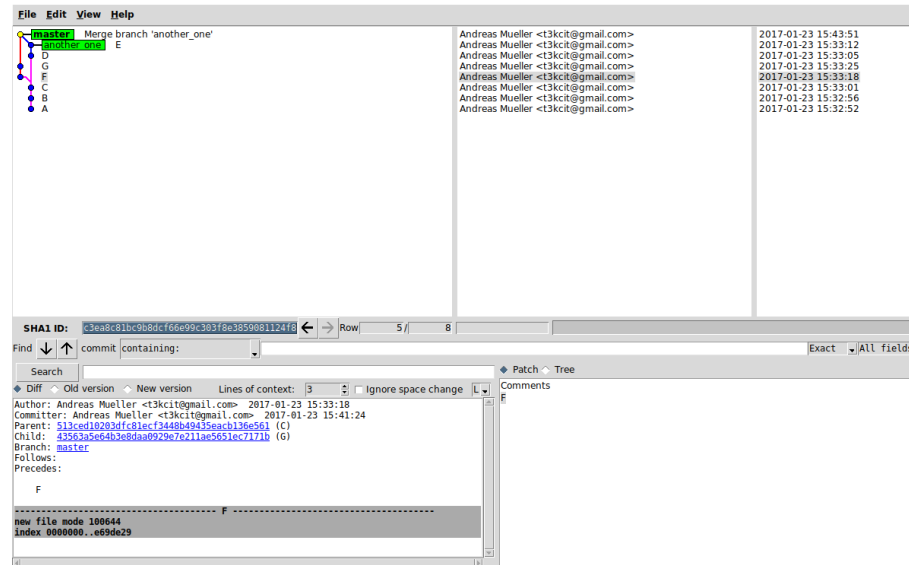
```




git gui



gitk



Understanding Git

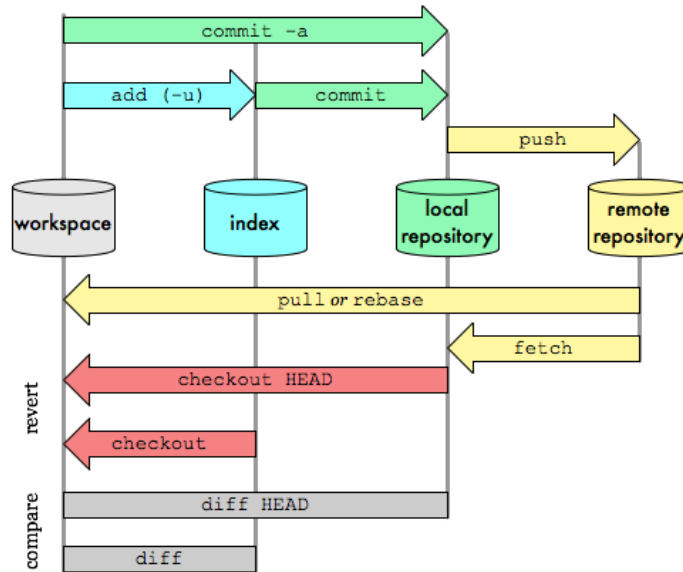
- Working directory
- Repository (Commit graph, history)
- Index (Staging Area)
- Branches
- Head

git Commands

- **git add**
puts files from working director into staging area (index) If not tracked so far, adds them to tracked files.
- **git commit**
commits files from staging area (index) to repository, moves current branch with HEAD
- **git checkout [<commit>] [<file>]**
Set <file> in working directory to state at <commit> and stages it.
- **git checkout [-b] <branch>**
moves HEAD to <branch> (-b creates it), changes content of working dir
- **git reset --soft <commit>**
moves HEAD to <commit> (takes the current branch with it)
- **git reset --mixed <commit>**
moves HEAD to <commit>, changes index to be at <commit> (but not working directory)
- **git reset --hard <commit>**
moved HEAD to <commit>, changes index and working tree to <commit>.

Git Data Transport Commands

<http://osteele.com>



reflog

\$ git reflog

```
d6fedb0 HEAD@{0}: merge another_one: Merge made by the 'recursive' strategy.
43563a5 HEAD@{1}: rebase finished: returning to refs/heads/master
43563a5 HEAD@{2}: rebase: G
c3ea8c8 HEAD@{3}: rebase: F
513ced1 HEAD@{4}: rebase: checkout 513ced1
4db426c HEAD@{5}: merge feature: Fast-forward
b5ba00a HEAD@{6}: checkout: moving from master to master
b5ba00a HEAD@{7}: reset: moving to HEAD~3
6cec4ed HEAD@{8}: merge another_one: Fast-forward
513ced1 HEAD@{9}: checkout: moving from another_one to master
6cec4ed HEAD@{10}: reset: moving to 6cec4ed
4db426c HEAD@{11}: checkout: moving from feature to another_one
4db426c HEAD@{12}: checkout: moving from master to feature
513ced1 HEAD@{13}: reset: moving to HEAD~4
4db426c HEAD@{14}: checkout: moving from feature to master
4db426c HEAD@{15}: checkout: moving from master to feature
4db426c HEAD@{16}: commit: G
125e957 HEAD@{17}: commit: F
6cec4ed HEAD@{18}: commit: E
1e96a3b HEAD@{19}: commit: D
513ced1 HEAD@{20}: commit: C
b5ba00a HEAD@{21}: commit: B
db928a0 HEAD@{22}: commit (initial): A
```

```
> git log --oneline --graph --decorate --all
* d6fedb0 (HEAD -> master) Merge branch 'another_one'
| \
| * 6cec4ed (another_one) E
| * 1e96a3b D
| * | 43563a5 G
| * | c3ea8c8 F
|/
* 513ced1 C
* b5ba00a B
* db928a0 A
```

Git for ages 4 and up:

<https://www.youtube.com/watch?v=1ffBJ4sVUb4>

(with play-doh!)

Unit Tests and integration tests

Why test?

- Ensure that code works correctly.
- Ensure that changes don't break anything.
- Ensure that bugs are not reintroduced.
- Ensure robustness to user errors.
- Ensure code is reachable.

Test-driven development?

Types of tests

- Unit tests – function does the right thing.
- Integration tests – system / process does the right thing.
- Non-regression tests – bug got removed (and will not be reintroduced).

How to test?

- pytest – <http://doc.pytest.org>
- Searches for all test * methods.
- Reports nice errors!
- Dig deeper: <http://pybites.blogspot.com/2011/07/behind-scenes-of-pytests-new-assertion.html>

Example

```
(x):
```

```
    x + 2
```

```
inc
```

```
inc
```

```
():
```

```
    inc(3) == 4
```

```
> py.test test_sample.py
===== test session starts =====
platform linux -- Python 3.5.2, pytest-2.9.2, py-1.4.31, pluggy-0.3.1
rootdir: /tmp, inifile:
plugins: cov-2.4.0, timeout-1.2.0
collected 1 items

test_sample.py F

===== FAILURES =====
_____ test_answer _____
```

Example

```
(x):  
    x + 1
```

```
inc      inc
```

```
():  
    inc(3) == 4
```

```
> py.test test_sample.py  
===== test session starts =====  
platform linux -- Python 3.5.2, pytest-2.9.2, py-1.4.31, pluggy-0.3.1  
rootdir: /tmp, inifile:  
plugins: cov-2.4.0, timeout-1.2.0
```


Test coverage

```
(x):  
x < 0:  
    0  
    x + 1  
  
(x):  
    x - 1
```

```
inc      inc  
  
():  
    inc(3) == 4
```

```
> py.test  
===== test session starts =====  
platform linux -- Python 3.5.2, pytest-2.9.2, py-1.4.31, pluggy-0.3.1  
rootdir: /tmp/myproj, inifile:  
plugins: cov-2.4.0, timeout-1.2.0  
collected 1 items  
  
test_inc.py .  
  
===== 1 passed in 0.01 seconds =====
```

Test coverage

```
(x):  
x < 0:  
    0  
    x + 1  
  
(x):  
    x - 1
```

```
inc      inc  
  
():  
inc(3) == 4
```

```
/tmp/myproj [andy@dsi-amueller] [10:51]  
> py.test --cov inc  
===== test session starts =====  
platform linux -- Python 3.5.2, pytest-2.9.2, py-1.4.31, pluggy-0.3.1  
rootdir: /tmp/myproj, inifile:  
plugins: cov-2.4.0, timeout-1.2.0  
collected 1 items  
  
test_inc.py .  
  
----- coverage: platform linux, python 3.5.2-final-0 -----  
Name      Stmts  Miss  Cover  
-----  
inc.py      6      2   67%
```

HTML report

```
$ pytest --cov inc --cov-report=html
```

Coverage report: 67%

Module ↓	statements	missing	excluded	coverage
inc.py	6	2	0	67%
Total	6	2	0	67%

coverage.py v4.2, created at 2017-01-23 10:53

Coverage for **inc.py** : 67%

6 statements 4 run 2 missing 0 excluded

```
1 # inc.py
2 def inc(x):
3     if x < 0:
4         return 0
5     return x + 1
6
7
8 def dec(x):
9     return x - 1
```

Continuous integration (with GitHub)

What is Continuous integration?

- Run command on each commit (or each PR).
- Unit testing and integration testing.
- Can act as a build-farm (for binaries or documentation).
- requires clear declaration of dependencies.
- Build matrix: Can run on many environments.
- Standard serviced: TravisCI, Appveyor, Azure Pipelines and CircleCI

Benefits of CI

- Can run on many systems
- Can't forget to run it
- Contributor doesn't need to know details
- Can enforce style
- Can provide immediate feedback
- Protects the master branch (if run on PR)

What does it do?

- Triggered at each commit / push
- Sets up a virtual machine with your configuration.
- Pulls the current branch.
- Runs command. Usually: install, then test.
- Reports success / Failure to github.

Setting up TravisCI

- Create account linked to your github account: <http://travis-ci.org>



AppliedMachineLearning/Homework-I-starter

- Check out docs at <https://docs.travis-ci.com>

```
language: python
python:
  - "2.7"
  - "3.4"
  - "3.5"
  - "3.6"
  - "3.7"
  - "nightly"

install:
  - pip install -r requirements.txt

script:
  - pytest
```


Using Travis

- Triggered any time you push a change
- Integrated with Pull requests
- Try a pull request on your own repository!

Questions ?