

BETTER SCIENCE CODE

Eric Denovellis

Presentation: **<https://edeno.github.io/Better-Science-Code>**

Repository: **<https://github.com/edeno/Better-Science-Code>**

Why should you care about producing good code

REASON 1. Doing good science!

Why should you care about producing good code

Following good coding practices reduces the chance of making mistakes.

Why should you care about producing good code

REASON 2. Want to remember what the code does months later

Why should you care about producing good code

REASON 3. Want to be able to share it with other people

Why should you care about producing good code

REASON 4. Avoid introducing new errors

Why should you care about producing good code

REASON 5. Can serve as a resume for future employers

Exercise in managing complexity:

- break problems down into smaller components
- eliminate unnecessary dependencies
- keep track of what you did (be organized)

How to write good code???

STEP 1. Decompose programs into small, well-defined functions

```
import numpy as np

def bad_function():
    X = np.load('/tmp/123.npy', mmap_mode='r')
    y, x1, x2 = X[:, 0], X[:, 1], X[:, 2]
    z1 = (x1 - x1.mean()) / x1.std()
    Q1, R1 = np.linalg.qr(z1, mode='reduced')
    b1 = np.linalg.solve(R1, np.dot(Q1.T, y1))
    z2 = (x2 - x2.mean()) / x2.std()
    Q2, R2 = np.linalg.qr(z1, mode='reduced')
    b2 = np.linalg.solve(R2, np.dot(Q2.T, y2))
    b = b1 - b2
    np.save('ans.npy', b)
```

```
import numpy as np

def better_function():
    y, x1, x2 = load_data('/tmp/123.npy')
    b1 = linear_regression(zscore(x1), y)
    b2 = linear_regression(zscore(x2), y)
    b = b1 - b2
    np.save('ans.npy', b)

def load_data(data_name):
    X = np.load(data_name, mmap_mode='r')
    return X[:, 0], X[:, 1], X[:, 2]

def zscore(x):
    return (x - x.mean()) / x.std()

def linear_regression(design_matrix, response):
    Q, R = np.linalg.qr(design_matrix, mode='reduced')
    return np.linalg.solve(R, np.dot(Q.T, response))
```

How to write good code???

Try to keep functions to less than 60 lines (small)

How to write good code???

Try to keep what the function does as simple as possible (well-defined)

How to write good code???

Be ruthless about eliminating duplication of code.

Small, well-defined, without duplicates

```
import numpy as np

def bad_function():
    X = np.load('/tmp/123.npy', mmap_mode='r')
    y, x1, x2 = X[:, 0], X[:, 1], X[:, 2]
    z1 = (x1 - x1.mean()) / x1.std()
    Q1, R1 = np.linalg.qr(z1, mode='reduced')
    b1 = np.linalg.solve(R1, np.dot(Q1.T, y1))
    z2 = (x2 - x2.mean()) / x2.std()
    Q2, R2 = np.linalg.qr(z1, mode='reduced')
    b2 = np.linalg.solve(R2, np.dot(Q2.T, y2))
    b = b1 - b2
    np.save('ans.npy', b)
```

Small, well-defined, without duplicates

```
import numpy as np

def better_function():
    y, x1, x2 = load_data('/tmp/123.npy')
    b1 = linear_regression(zscore(x1), y)
    b2 = linear_regression(zscore(x2), y)
    b = b1 - b2
    np.save('ans.npy', b)

def load_data(data_name):
    X = np.load(data_name, mmap_mode='r')
    return X[:, 0], X[:, 1], X[:, 2]

def zscore(x):
    return (x - x.mean()) / x.std()

def linear_regression(design_matrix, response):
    Q, R = np.linalg.qr(design_matrix, mode='reduced')
    return np.linalg.solve(R, np.dot(Q.T, response))
```


Use good variable/function names

```
import numpy as np

def bad_function():
    X = np.load('/tmp/123.npy', mmap_mode='r')
    y, x1, x2 = X[:, 0], X[:, 1], X[:, 2]
    z1 = (x1 - x1.mean()) / x1.std()
    Q1, R1 = np.linalg.qr(z1, mode='reduced')
    b1 = np.linalg.solve(R1, np.dot(Q1.T, y1))
    z2 = (x2 - x2.mean()) / x2.std()
    Q2, R2 = np.linalg.qr(z1, mode='reduced')
    b2 = np.linalg.solve(R2, np.dot(Q2.T, y2))
    b = b1 - b2
    np.save('ans.npy', b)
```

Use good variable/function names

```
import numpy as np

def better_function():
    y, x1, x2 = load_data('/tmp/123.npy')
    b1 = linear_regression(zscore(x1), y)
    b2 = linear_regression(zscore(x2), y)
    b = b1 - b2
    np.save('ans.npy', b)

def load_data(data_name):
    X = np.load(data_name, mmap_mode='r')
    return X[:, 0], X[:, 1], X[:, 2]

def zscore(x):
    return (x - x.mean()) / x.std()

def linear_regression(design_matrix, response):
    Q, R = np.linalg.qr(design_matrix, mode='reduced')
    return np.linalg.solve(R, np.dot(Q.T, response))
```

Use good variable/function names

```
import numpy as np

def better_function():
    response, design_matrix1, design_matrix2 = load_data(
        '/tmp/123.npy')
    coefficient1 = linear_regression(
        zscore(design_matrix1), response)
    coefficient2 = linear_regression(
        zscore(design_matrix2), response)
    coefficient_difference = coefficient1 - coefficient2
    np.save('ans.npy', coefficient_difference)

def load_data(data_name):
    X = np.load(data_name, mmap_mode='r')
    return X[:, 0], X[:, 1], X[:, 2]

def zscore(x):
    return (x - x.mean()) / x.std()

def linear_regression(design_matrix, response):
    Q, R = np.linalg.qr(design_matrix, mode='reduced')
    return np.linalg.solve(R, np.dot(Q.T, response))
```


Document your functions

Easy thing: brief sentence describing the function without using the name of the function*

**this is the most important*

Document your functions

```
def zscore(x):  
    return (x - x.mean()) / x.std()  
  
def linear_regression(design_matrix, response):  
    Q, R = np.linalg.qr(design_matrix, mode='reduced')  
    return np.linalg.solve(R, np.dot(Q.T, response))
```

Document your functions

```
def zscore(x):  
    '''Number of standard deviations from the mean'''  
    return (x - x.mean()) / x.std()  
  
def linear_regression(design_matrix, response):  
    Q, R = np.linalg.qr(design_matrix, mode='reduced')  
    return np.linalg.solve(R, np.dot(Q.T, response))
```

Document your functions

```
def zscore(x):  
    '''Number of standard deviations from the mean'''  
    return (x - x.mean()) / x.std()  
  
def linear_regression(design_matrix, response):  
    '''Calculate a linear least-squares regression for  
    two sets of measurements'''  
    Q, R = np.linalg.qr(design_matrix, mode='reduced')  
    return np.linalg.solve(R, np.dot(Q.T, response))
```

Document your functions

- additional detail about what the function does or method it implements
- description of the parameters
- description of the outputs
- examples if you can

Document your functions

```
def linear_regression(design_matrix, response):  
    '''Calculate a linear least-squares regression for  
    two sets of measurements  
  
    Uses the QR decomposition to avoid numerical instability  
    in taking the inverse.  
  
    Parameters  
    -----  
    design_matrix, response : array_like  
        Two sets of measurements. Both arrays should have  
        the same length.  
  
    Returns  
    -----  
    coefficients : array_like  
        Parameters estimated from the model.  
  
    Examples  
    -----  
>>> design_matrix = np.random.random(10)  
>>> response = np.random.random(10)  
>>> coefficients = linear_regression(design_matrix, response)  
  
    '''
```


Test your code

```
import numpy as np

def zscore(x):
    '''Number of standard deviations from the mean'''
    return (x - x.mean()) / x.std()

def test_zscore():
    pass
```

Test your code

```
import numpy as np

def zscore(x):
    '''Number of standard deviations from the mean'''
    return (x - x.mean()) / x.std()

def test_zscore():
    test_values = np.asarray([1, 3])
    expected_values = np.asarray([-1, 1])

    assert np.allclose(zscore(test_values), expected_values)
```

Test your code

Unit tests test a small component of your code (usually a small function) and makes sure it works like you think it works

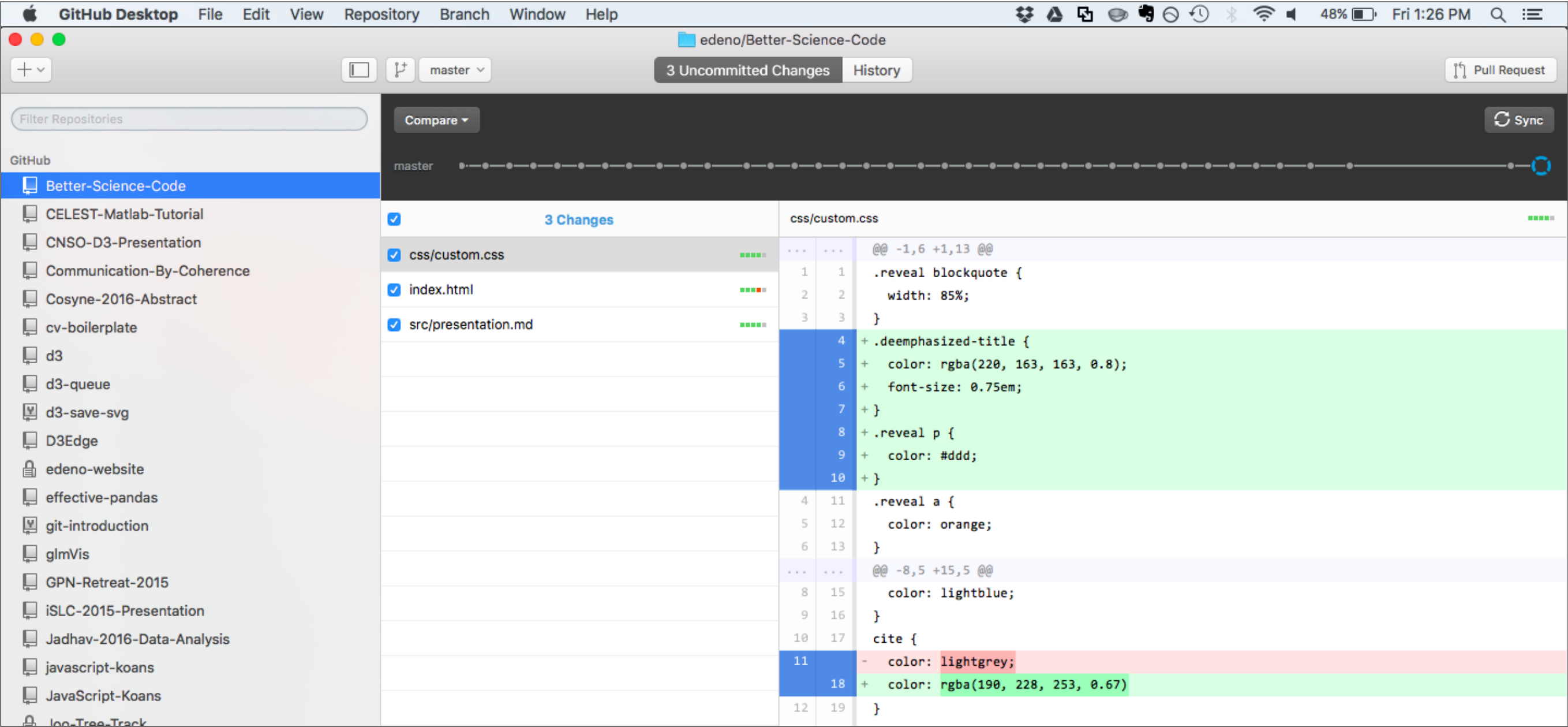
Matlab, Python, R have unit test packages

- **Matlab unit test framework**
- **Python unit test**
- **Pytest**
- **R: testthat**

Use version control

Sophisticated way to track change in your code over time

Use version control



Github Desktop

GitHub Desktop

FileEditViewRepositoryBranchWindowHelp

edeno/Better-Science-Code

4 Uncommitted ChangesHistory

Pull Request

Filter Repositories

GitHub

Better-Science-Code

CELEST-Matlab-Tutorial

CNSO-D3-Presentation

Communication-By-Coherence

Cosyne-2016-Abstract

cv-boilerplate

d3

d3-queue

d3-save-svg

D3Edge

edeno-website

effective-pandas

git-introduction

glmVis

Compare

Commit History

Sync

master

Add resource for file naming

23 hours ago by edeno

Fix example code

1 month ago by edeno

3

Add another reference

1 month ago by edeno

2

Add annotation

1 month ago by edeno

Emphasize that the workflows are from various scie...

1 month ago by edeno

Update citation

1 month ago by edeno

Add resource for file naming

edeno 4702910 23 hours ago

README.md

... @@ -35,5 +35,7 @@ A presentation on best coding/data management practices based on my own experien

35 35 [Millman, K.J., and Pérez, F. (2014). Developing open-source scientific practice. Implementing Reproducible Research 149.

36 36](https://osf.io/h9gsd/)

37 37

38 + [Best practices for file naming](http://library.stanford.edu/research/data-management-services/data-best-practices/best-practices-file-naming)

39 +

38 40 # Requirements

39 41 * [pandoc](http://pandoc.org/)

Use version control

Helps you back up your work

Use version control

Go back to previous versions of your code

GitHub Desktop

FileEditViewRepositoryBranchWindowHelp

edeno/Better-Science-Code

4 Uncommitted ChangesHistory

Pull Request

Filter Repositories

GitHub

Better-Science-Code

CELEST-Matlab-Tutorial

CNSO-D3-Presentation

Communication-By-Coherence

Cosyne-2016-Abstract

cv-boilerplate

d3

d3-queue

d3-save-svg

D3Edge

edeno-website

effective-pandas

git-introduction

glmVis

Compare

Commit History

Sync

master

Add resource for file naming

23 hours ago by edeno

Fix example code

1 month ago by edeno

3

Add another reference

1 month ago by edeno

2

Add annotation

1 month ago by edeno

Emphasize that the workflows are from various scie...

1 month ago by edeno

Update citation

1 month ago by edeno

Add resource for file naming

edeno 4702910 23 hours ago

README.md

... @@ -35,5 +35,7 @@ A presentation on best coding/data management practices based on my own experien

35 35 [Millman, K.J., and Pérez, F. (2014). Developing open-source scientific practice. Implementing Reproducible Research 149.

36 36](https://osf.io/h9gsd/)

37 37

38 + [Best practices for file naming](http://library.stanford.edu/research/data-management-services/data-best-practices/best-practices-file-naming)

39 +

38 40 # Requirements

39 41 * [pandoc](http://pandoc.org/)

Use version control

Reduce code clutter and confusion

Use version control

Experiment with different versions of code (branches)

Use version control

Makes it easier to work with others

Use version control

Commit early and often

Refactor your code

Always leave the code in a better state than when you first found it.

How to write good code???

Exercise in managing complexity:

- break problems down into smaller components
- eliminate unnecessary dependencies
- keep track of what you did (be organized)

Summary:

1. Write small well-defined, well-named functions
2. Use good function and variable names
3. Document your functions
4. Test your code
5. Refactor your code
6. Use version control
7. Always search for well-maintained software libraries that do what you need.

break problems down into smaller components

1. Write small well-defined, well-named functions
2. Use good function and variable names
3. Document your functions
4. Test your code
5. Refactor your code
6. Use version control
7. Always search for well-maintained software libraries that do what you need.

keep track of what you did (be organized)

1. Write small well-defined, well-named functions
2. Use good function and variable names
3. Document your functions
4. Test your code
5. Refactor your code
6. Use version control
7. Always search for well-maintained software libraries that do what you need.

Tidy Data:

- Each variable forms a column.
- Each observation forms a row.
- Each type of observational unit forms a table
- flat is better than nested

File naming:

- Don't use spaces in file names
- Use leading zeros (001 vs. 1)