

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!

We want code that **works** (it does what you say it does) and is **reproducible** (you can get to the same result every time using the same data and code):

Don't want to have to retract papers because the code had bugs

Following good coding practices reduces the chance of making mistakes

IT'S TOO EASY TO MAKE MISTAKES

“As the complexity of a software program increases, the likelihood of undiscovered bugs quickly reaches certainty” – Poldrack et al. 2017

We are writing *complex code*

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

“The single biggest reason you should write nice code is so that your future self can understand it.” – Greg Wilson

“All code has at least one collaborator and that is future you.” – Hadley Wickham

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

Particularly important with statistical methods development

REASON 4. Avoid introducing new errors

REASON 5. Can serve as a resume for future employers

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)

Goal: Want to form good habits

Don't be overwhelmed *and not do any of these things*

Don't beat yourself up *if you don't do all these things all the time*

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))
```

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

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

Be ruthless about eliminating duplication of code

```
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))
```

Small, well-defined functions are more *maintainable*

Small, well-defined functions are more *composable*

Small, well-defined functions are more *readable*

* if you give them good names

STEP 2. Use good variable/function names to clarify what things do

```
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))
```

```
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))
```

You don't need comments if the variable or function already tells you what it does (self-documenting)

Use the naming conventions of your language of choice (`snake_case` or `camelCase`) and be consistent

Avoid using abbreviations that are not commonly used

(jmi vs. joint_mark_intensity)

Prefer whole words

(elec_poten vs. electric_potential)

STEP 3. Document your functions

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

**this is the most important*

```
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))
```

More complicated thing:

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

```
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)

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

STEP 4. Test your code

Make sure your code works like you think it does

Think about how your code can fail

Small, well-defined, well-named functions are easy to test!

```
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)
```

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

Unit tests prevent regression of your code

If you change your code, you want to know what still works and what has broken

Functions should be simple to test

If you find a bug, write a test.

Use unit tests to define the requirements of your code

You can use programs called **test runners** to run a group of unit tests automatically.

Matlab, Python, R have unit test packages

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

There are also libraries available that will work with your version control system to run these tests every time you commit a new piece of code ([continuous integration](#))

STEP 5. Use version control

Sophisticated way to track change in your code over time

GitHub Desktop

FileEditViewRepositoryBranchWindowHelp

eden0/Better-Science-Code

3 Uncommitted ChangesHistory

Pull Request

Filter Repositories

Sync

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

eden0-website

git-introduction

glmVis

GPN-Retreat-2015

iSLC-2015-Presentation

Jadhav-2016-Data-Analysis

javascript-koans

JavaScript-Koans

Joo-Tree-Track

Kalman-and-Bayesian-Filters-in-Python

lorenlab-matlab

matorque

modern-pandas

MusicMappr

playground

Prospectus

Prospectus-Presentation

pydata-book

python-koans

Compare

master

2

Use markdown syntax highlighting

1 hour ago by eden0

2

Revert "Use the reveal.js highlighting"

1 hour ago by eden0

Fix mistake in directory

1 hour ago by eden0

Add in highlight.js manually

1 hour ago by eden0

Update

1 hour ago by eden0

Use the reveal.js highlighting

1 hour ago by eden0

Separate Reasons from Steps

1 hour ago by eden0

Remove type and shape from description

14 hours ago by eden0

2

Add more about commenting code

14 hours ago by eden0

2

Add to docstring example

14 hours ago by eden0

2

Eliminate knuth quote

15 hours ago by eden0

2

Turn into lists

15 hours ago by eden0

2

Formatting

15 hours ago by eden0

2

Turn into links

15 hours ago by eden0

Use markdown syntax highlighting

eden0 8289959 1 hour ago

index.html

... 10 11 12

@@ -10,6 +10,44 @@

<meta name="viewport" content="width=device-width, initial-scale=1.0, maximum-scale=1.0, user-scalable=no, minimal-ui">

<link rel="stylesheet" href="revealjs/css/reveal.css">

<style type="text/css">code{white-space: pre;}</style>

+ <style type="text/css">

+ div.sourceCode { overflow-x: auto; }

+ table.sourceCode, tr.sourceCode, td.lineNumbers, td.sourceCode {

+ margin: 0; padding: 0; vertical-align: baseline; border: none; }

+ table.sourceCode { width: 100%; line-height: 100%; background-color: #303030; color: #cccccc; }

+ td.lineNumbers { text-align: right; padding-right: 4px; padding-left: 4px; }

+ td.sourceCode { padding-left: 5px; }

+ pre, code { color: #cccccc; background-color: #303030; }

+ code > span.kw { color: #f0dfaf; } /* Keyword */

+ code > span.dt { color: #dfd7bf; } /* DataType */

+ code > span.dv { color: #dcdccc; } /* DecVal */

+ code > span.bn { color: #dca3a3; } /* BaseN */

+ code > span.fl { color: #c0bed1; } /* Float */

+ code > span.ch { color: #dca3a3; } /* Char */

+ code > span.st { color: #cc9393; } /* String */

+ code > span.co { color: #7f9f7f; } /* Comment */

+ code > span.ot { color: #efef8f; } /* Other */

+ code > span.al { color: #ffcfaf; } /* Alert */

+ code > span.fu { color: #efef8f; } /* Function */

+ code > span.er { color: #c3bf9f; } /* Error */

+ code > span.wa { color: #7f9f7f; font-weight: bold; } /* Warning */

+ code > span.cn { color: #dca3a3; font-weight: bold; } /* Constant */

+ code > span.sc { color: #dca3a3; } /* SpecialChar */

Commit History

Version control stores the whole history of your project

GitHub Desktop

FileEditViewRepositoryBranchWindowHelp

edeno/Better-Science-Code

3 Uncommitted ChangesHistory

Pull Request

Filter Repositories

CompareSync

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

git-introduction

glmVis

GPN-Retreat-2015

iSLC-2015-Presentation

Jadhav-2016-Data-Analysis

javascript-koans

JavaScript-Koans

Joo-Tree-Track

Kalman-and-Bayesian-Filters-in-Python

lorenlab-matlab

matorque

modern-pandas

MusicMappr

playground

Prospectus

Prospectus-Presentation

pydata-book

python-koans

master

2

Use markdown syntax highlighting

1 hour ago by edeno

Revert "Use the reveal.js highlighting"

1 hour ago by edeno

Fix mistake in directory

1 hour ago by edeno

Add in highlight.js manually

1 hour ago by edeno

Update

1 hour ago by edeno

Use the reveal.js highlighting

1 hour ago by edeno

Separate Reasons from Steps

1 hour ago by edeno

Remove type and shape from description

14 hours ago by edeno

Add more about commenting code

14 hours ago by edeno

Add to docstring example

14 hours ago by edeno

Eliminate knuth quote

15 hours ago by edeno

Turn into lists

15 hours ago by edeno

Formatting

15 hours ago by edeno

Turn into links

15 hours ago by edeno

Use markdown syntax highlighting

edeno 8289959 1 hour ago

index.html

@@ -10,6 +10,44 @@

10 10 <meta name="viewport" content="width=device-width, initial-scale=1.0, maximum-scale=1.0, user-scalable=no, minimal-ui">

11 11 <link rel="stylesheet" href="revealjs/css/reveal.css">

12 12 <style type="text/css">code{white-space: pre;}</style>

13 + <style type="text/css">

14 + div.sourceCode { overflow-x: auto; }

15 + table.sourceCode, tr.sourceCode, td.lineNumbers, td.sourceCode {

16 + margin: 0; padding: 0; vertical-align: baseline; border: none; }

17 + table.sourceCode { width: 100%; line-height: 100%; background-color: #303030; color: #cccccc; }

18 + td.lineNumbers { text-align: right; padding-right: 4px; padding-left: 4px; }

19 + td.sourceCode { padding-left: 5px; }

20 + pre, code { color: #cccccc; background-color: #303030; }

21 + code > span.kw { color: #f0dfaf; } /* Keyword */

22 + code > span.dt { color: #dfd7bf; } /* DataType */

23 + code > span.dv { color: #dcdccc; } /* DecVal */

24 + code > span.bn { color: #dca3a3; } /* BaseN */

25 + code > span.fl { color: #c0bed1; } /* Float */

26 + code > span.ch { color: #dca3a3; } /* Char */

27 + code > span.st { color: #cc9393; } /* String */

28 + code > span.co { color: #7f9f7f; } /* Comment */

29 + code > span.ot { color: #efef8f; } /* Other */

30 + code > span.al { color: #ffcfaf; } /* Alert */

31 + code > span.fu { color: #efef8f; } /* Function */

32 + code > span.er { color: #c3bf9f; } /* Error */

33 + code > span.wa { color: #7f9f7f; font-weight: bold; } /* Warning */

34 + code > span.cn { color: #dca3a3; font-weight: bold; } /* Constant */

35 + code > span.sc { color: #dca3a3; } /* SpecialChar */

Commit History

Helps you back up your work

Go back to previous versions of your code

Reduce code clutter and confusion

Experiment with different versions of code (branches)

Makes it easier to work with others

Commit early and often

STEP 6. Refactor your code

“Whenever I have to think to understand what the code is doing, I ask myself if I can refactor the code to make that understanding more immediately apparent.” – Martin Fowler, Refactoring: Improving the Design of Existing Code

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

STEP 7. Always search for well-maintained software libraries that do what you need.

Don't rewrite functions that are already implemented as part of the core language.

Use other software libraries if they are well-maintained

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.

Conclusion: Writing good code takes work

We have a scientific obligation to ensure the correctness of our programs.

Bonus: Data Management

Put different projects in different folders/repositories

Use relative paths

Separate the data from the code

Processed Data should be separated from Raw Data to avoid accidentally changing the data

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

If original data is not in a good form, convert it to a good form (but don't overwrite the original data)

Don't hand-edit data files.

All aspects of data cleaning should be in scripts