

Data Science is Software: Developer #lifehacks for the Jupyter Data Scientist

PETER BULL
[@drivendataorg](https://drivendata.org) / [@pjbull](https://twitter.com/pjbull)

bit.ly/jupyter-lifehacks

Screenshot of a GitHub repository page for `pjbull / data-science-is-software`.

The repository has 5 commits, 1 branch, 0 releases, and 1 contributor.

Key UI elements highlighted with red boxes:

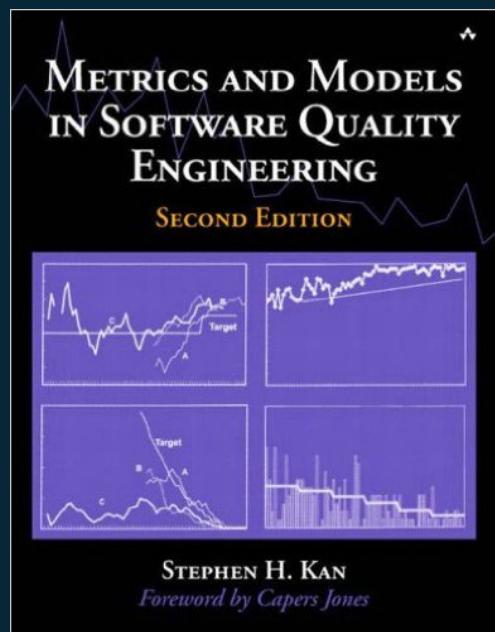
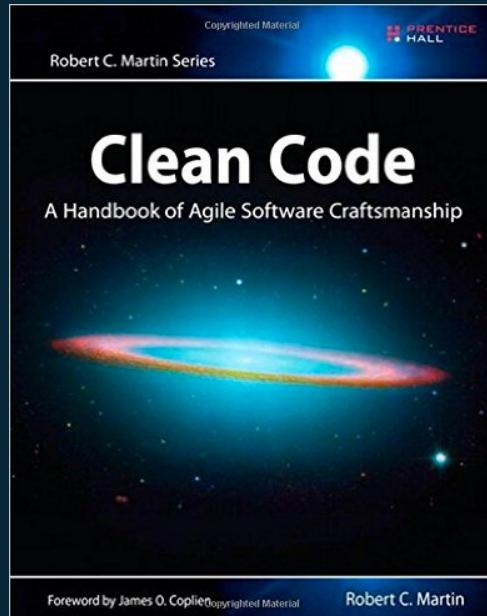
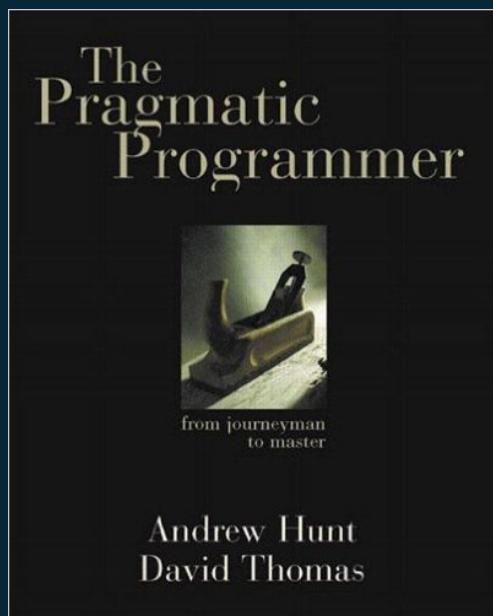
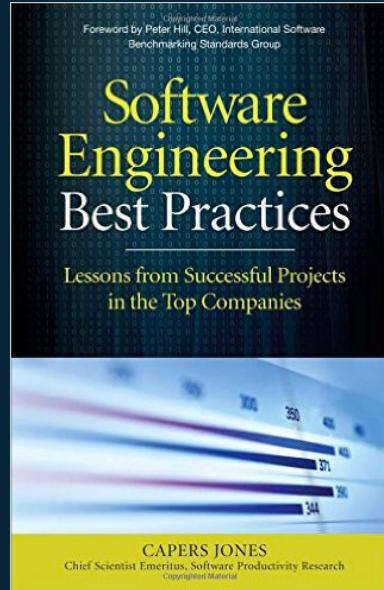
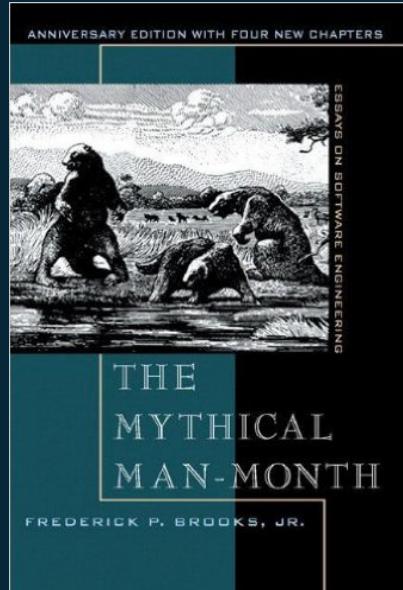
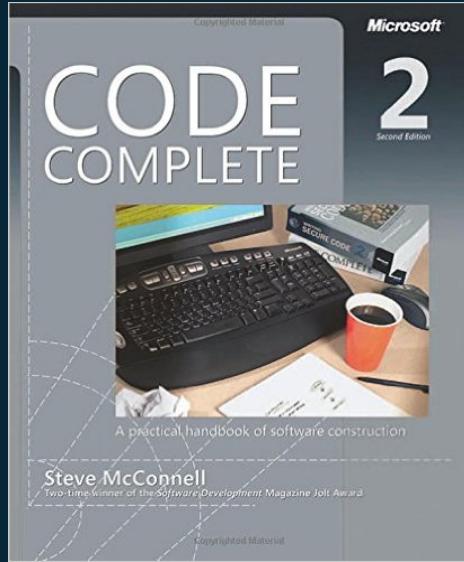
- HTTPS** dropdown and URL (`https://github.com/pjbull`)
- OR** link
- Download ZIP** button

git clone is displayed prominently above the commit list.

File	Message	Time
<code>data</code>	Add BDM Notebook	2 days ago
<code>notebooks</code>	Materials update	an hour ago
<code>slides</code>	Materials update	an hour ago
<code>src</code>	Materials update	an hour ago
<code>.gitignore</code>	Talk updates!	12 hours ago
<code>LICENSE</code>	Initial commit	2 days ago
<code>README.md</code>	Talk updates!	12 hours ago
<code>requirements.txt</code>	Materials update	an hour ago

“It's easy to enhance a
FORTRAN compiler to
compile COBOL as well; it's
just a SMOP.”

– SMOP entry in The Jargon File
(a comprehensive
compendium of
hacker slang)



1

This is my house

Spot 7 differences between
these images with piglets.

<http://www.everydayok.com>







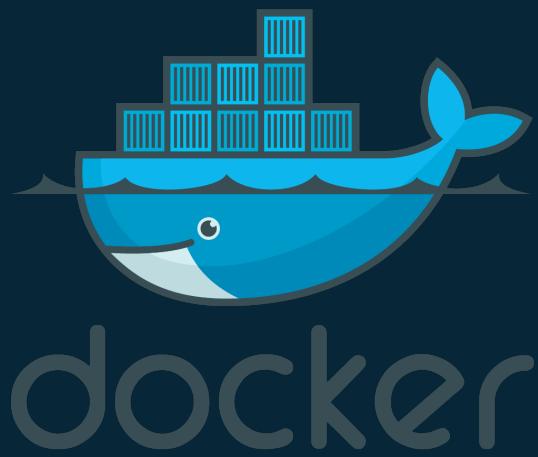
WATERMARK

VIRTUALENV

VIRTUALENVWRAPPER

PIP REQUIREMENTS.TXT

Other options (for more complex environments)



2

The Life-Changing Magic of Tidying Up



#1
NEW YORK TIMES
BEST SELLER
—
3 MILLION
COPIES SOLD

the life-changing magic of tidying up

the Japanese art of decluttering
and organizing

marie kondo

```
.  
├── Inspection_count_min.jpeg  
├── README.Rmd  
├── README.html  
├── dd_dictionary.csv  
├── mallet.rar  
├── scripts\ and\ data  
│   ├── AllViolations.csv  
│   ├── PhaseIISubmissionFormat.csv  
│   ├── build_rev_tm.R  
│   ├── docsAsTopicsProbs_noStopwords.txt  
│   ├── feature_eng.R  
│   ├── features_test_phase2.csv  
│   ├── features_train_phase2.csv  
│   ├── learning_final.R  
│   ├── negative-words.txt  
│   ├── positive-words.txt  
│   ├── rand_neg.txt  
│   ├── restaurant_ids_to_yelp_ids.csv  
│   ├── rev_tm.txt  
│   ├── review_sentiscored.csv  
│   ├── run.R  
│   ├── sentiment_script.R  
│   ├── sub_2_PhaseII_h20.csv  
│   ├── yelp.stops  
│   └── yelp_academic_dataset_business.json  
└── varimp_gbm1.jpeg  
└── varimp_gbm2.jpeg  
└── varimp_sev.jpeg
```

```
.  
├── AllViolations.csv  
├── BusinessClass.py  
├── GenLearningData.py  
├── GenTestingData.py  
├── InspectionClass.py  
├── LearnTest.py  
├── PhaseIISubmissionFormat.csv  
├── PhaseIISubmissionFormat_final.csv  
├── PhaseIISubmissionFormat_test.csv  
├── README.txt  
├── ReviewClass.py  
├── restaurant_ids_to_yelp_ids.csv  
├── yelp_boston_academic_dataset  
└── yelp_duplicate_ids.csv
```

```
├── Step\ 1\ -\ install\ necessary\ software\ and\ packages.txt  
├── Step\ 2\ -\ one-off\ step\ to\ create\ postgresql\ server\ instance\ and\ a\ database.txt  
├── Step\ 3\ -\ one-off\ step\ to\ create\ tables\ and\ views\ in\ postgresql.py  
└── Step\ 4\ -\ The\ only\ file\ to\ run\ when\ you\ want\ to\ run\ models\ and\ generate\ new\ scores.py
```

Ruby on Rails Application

```
.  
|   ├── Gemfile  
|   ├── Gemfile.lock  
|   ├── Guardfile  
|   ├── LICENSE  
|   ├── README.md  
|   ├── README.nitrous.md  
|   └── Rakefile  
|       ├── app  
|       |   ├── assets  
|       |   ├── controllers  
|       |   ├── helpers  
|       |   ├── mailers  
|       |   ├── models  
|       |   └── views  
|       ├── bin  
|       |   ├── bundle  
|       |   ├── rails  
|       |   └── rake  
|       ├── config  
|       |   ├── application.rb  
|       |   ├── boot.rb  
|       |   ├── cucumber.yml  
|       |   ├── database.yml.example  
|       |   ├── environment.rb  
|       |   ├── environments  
|       |   ├── initializers  
|       |   ├── locales  
|       |   └── routes.rb  
|       ├── config.ru  
|       └── db  
|           ├── migrate  
|           ├── schema.rb  
|           └── seeds.rb  
|       ├── features  
|       |   ├── signing_in.feature  
|       |   ├── step_definitions  
|       |   └── support  
|       ├── lib  
|       |   ├── assets  
|       |   └── tasks  
|       ├── log  
|       ├── public  
|       |   ├── 404.html  
|       |   ├── 422.html  
|       |   ├── 500.html  
|       |   ├── assets  
|       |   ├── favicon.ico  
|       |   └── robots.txt  
|       ├── script  
|       |   └── cucumber  
|       ├── spec  
|       |   ├── controllers  
|       |   ├── factories.rb  
|       |   ├── helpers  
|       |   ├── models  
|       |   ├── requests  
|       |   ├── spec_helper.rb  
|       |   └── support  
|       └── vendor  
|           └── assets
```

Django Application

```
.  
|   ├── README.md  
|   └── media  
|       └── init.txt  
|   └── projectname  
|       ├── __init__.py  
|       └── home  
|           ├── __init__.py  
|           ├── models.py  
|           ├── tests.py  
|           └── views.py  
|       ├── manage.py  
|       ├── settings  
|       |   ├── __init__.py  
|       |   ├── default.py  
|       |   └── local.template.py  
|       ├── urls.py  
|       └── wsgi.py  
|   └── requirements.txt  
└── static-assets  
    ├── apple-touch-icon.png  
    ├── css  
    |   └── main.css  
    ├── favicon.ico  
    ├── humans.txt  
    ├── images  
    |   └── init.txt  
    ├── js  
    |   ├── main.coffee  
    |   ├── main.js  
    |   └── main.map  
    └── libs  
        ├── bootstrap-3.3.5  
        ├── font-awesome-4.3.0  
        ├── html5shiv.js  
        ├── jquery  
        └── modernizr  
    └── media -> ../media/  
    └── robots.txt  
└── templates  
    ├── 404.html  
    ├── 500.html  
    ├── base.html  
    └── home.html
```

Cookiecutter Data Science

[Why use this project structure?](#)

[Other people will thank you](#)

[You will thank you](#)

[Nothing here is binding](#)

[Getting started](#)

[Requirements](#)

[Starting a new project](#)

[Example](#)

[Directory structure](#)

[Opinions](#)

[Data is immutable](#)

[Notebooks are for exploration and communication](#)

[Analysis is a DAG](#)

[Build from the environment up](#)

[Keep secrets and configuration out of version control](#)

[Be conservative in changing the default folder structure](#)

[Contributing](#)

[Links to related projects and references](#)

Cookiecutter Data Science

A logical, reasonably standardized, but flexible project structure for doing and sharing data science work.

Why use this project structure?

We're not talking about bikeshedding the indentation aesthetics or pedantic formatting standards — ultimately, data science code quality is about correctness and reproducibility.

When we think about data analysis, we often think just about the resulting reports, insights, or visualizations. While these end products are generally the main event, it's easy to focus on making the products *look nice* and ignore the *quality of the code that generates them*. Because these end products are created programmatically, **code quality is still important!** And we're not talking about bikeshedding the indentation aesthetics or pedantic formatting standards — ultimately, data science code quality is about correctness and reproducibility.

It's no secret that good analyses are often the result of very scattershot and serendipitous explorations. Tentative experiments and rapidly testing approaches that might not work out are all part of the process for getting to the good stuff, and there is no magic bullet to turn data exploration into a simple, linear progression.

That being said, once started it is not a process that lends itself to thinking carefully about the structure of your code or project layout, so it's best to start with a clean, logical structure and stick to it throughout. We think it's a pretty big win all around to use a fairly standardized setup like this one. Here's why:

drivendata.github.io/cookiecutter-data-science

```
LICENSE  
Makefile  
README.md      <- Makefile with commands like `make data` or `make train`  
                  <- The top-level README for developers using this project.  
data  
  external      <- Data from third party sources.  
  interim       <- Intermediate data that has been transformed.  
  processed     <- The final, canonical data sets for modeling.  
  raw           <- The original, immutable data dump.  
  
docs          <- A default Sphinx project; see sphinx-doc.org for details  
  
models        <- Trained and serialized models, model predictions, or model summaries  
  
notebooks     <- Jupyter notebooks. Naming convention is a number (for ordering),  
                  the creator's initials, and a short '-' delimited description, e.g.  
                  `1.0-jqp-initial-data-exploration`.  
  
references    <- Data dictionaries, manuals, and all other explanatory materials.  
  
reports  
  figures      <- Generated analysis as HTML, PDF, LaTeX, etc.  
                  <- Generated graphics and figures to be used in reporting  
  
requirements.txt <- The requirements file for reproducing the analysis environment, e.g.  
                  generated with `pip freeze > requirements.txt`  
  
src  
  __init__.py   <- Source code for use in this project.  
                  <- Makes src a Python module  
  
  data          <- Scripts to download or generate data  
    make_dataset.py  
  
  features      <- Scripts to turn raw data into features for modeling  
    build_features.py  
  
  models        <- Scripts to train models and then use trained models to make  
                  predictions  
    predict_model.py  
    train_model.py  
  
  visualization <- Scripts to create exploratory and results oriented visualizations  
    visualize.py  
  
tox.ini       <- tox file with settings for running tox; see tox.testrun.org
```



3

Edit-run-repeat:
Stopping the cycle of pain





AUTORELOAD, PDB, DEBUG

Q, MODULES

ASSERT, UNITTEST

ENGARDE, NUMPY.TESTING



4

Next-level code
inspection

8200

19

Town

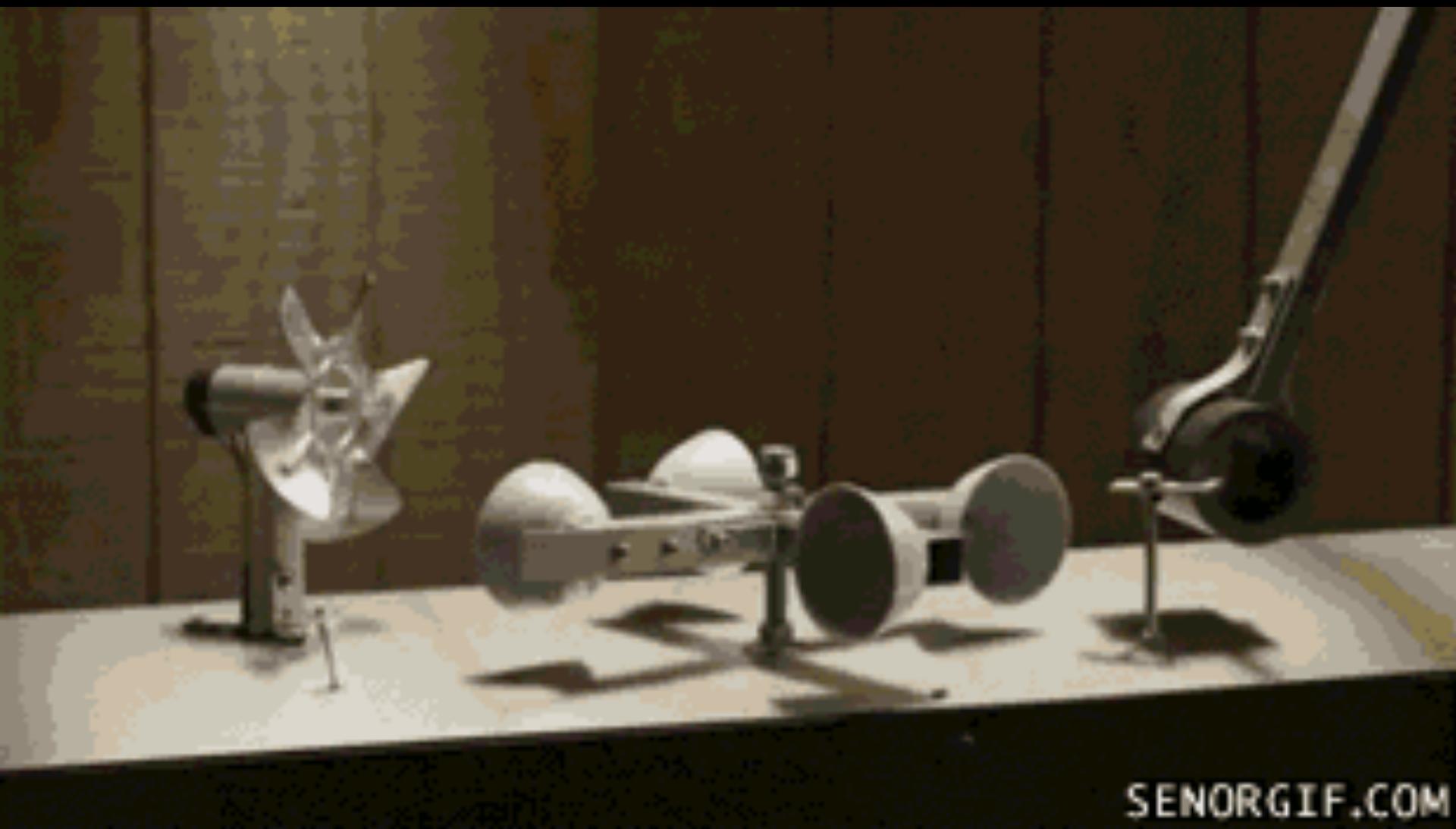


Save

Load







SENORGIF.COM







COVERAGE.PY

%PRUN

PYCHARM

FLAKE8 LINTING FOR SUBLIME

Not Covered

Version Control `git`

Code review `git branch + pull requests`

Branching strategy [GitHub Flow](#)

Issue tracking `GitHub Issues + waffle.io`

Automating workflows [Make](#)

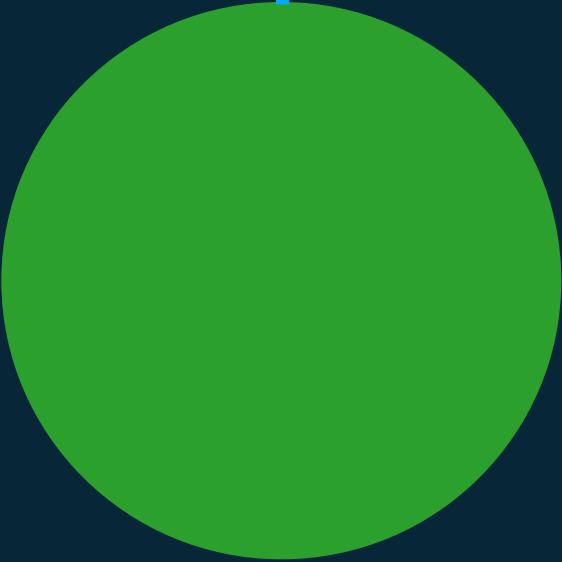
Documentation Generator [Sphinx](#)

Docstrings [Sphinx ReST](#)

Style guide [flake8](#)

Notebook differencing `nbconvert to .py on save`

Configuration isolation `config.py`



Questions?

DRIVENDATA
[@drivendataorg](https://twitter.com/drivendataorg) / [@pjbull](https://twitter.com/pjbull)
peter@drivendata.org