BETTER CODE FOR DATA SCIENCE

ALEXANDER CS HENDORF @ PYDATA GLOBAL 2020

ALEXANDER C. S. HENDORF

MANAGING PARTNER & PRINCIPAL CONSULTANT DATA SCIENCE & AI AT KÖNIGSWEG.

PYTHON SOFTWARE FOUNDATION FELLOW, PYTHON SOFTWAREVERBAND CHAIR, PYCONDE & PYDATA BERLIN CHAIR, LOCAL PYDATA COMMUNITY ORGANIZER



KŌNIGSWEG

We do digital excellence.

STRATEGY & INNOVATION DATA & ARTIFICIAL INTELLIGENCE BUSINESS TRANSFORMATION & OPERATIONS

Get in touch with our specialists.

PYDATA FRANKFURT



WEDNESDAY MAY 22 18:00 TECHQUARTIER FRANKFURT

- SciKit-Learn Keynote

Olivier Griesel INRIA





- Jupyter Keynote

Sylvain Corlay Quantstack



KÖNIGSWEG

PYDATA SÜDWEST

THU, NOV 21 18:00 - 21:30 X-HOUSE, HEIDELBERG (@MAIN STATION)

- A.I. Caramba: How to Maintain Sanity

Vincent D. Warmerdam GoDataDriven



- 2119: A Data Science Escape Room

Konrad Heimpel

Getsafe



KÖNIGSWEG







PyConDE & PyData Berlin 2021

Oct 13-15 2021 bcc Berlin Congress Center









Why do Programming Skills Matter?

- Quality Assurance
- Focus on value, not bugs
- productivity my vary by 10x*
- Maintainability
- Reusability

What about Personalities?



Some Personas

Stormy Simone	Coding to get things done studying, constantly overexcited about the latest paper
Abstract Alec	Just left academia, high level of problem abstraction, lacks production experience
Steady Saša	Skilled in another programming language, looking forward to learning new skills
Sceptical Stéphane	Skilled in another programming language, sceptical to new things, constantly comparing
Determined Dylan	Self-taught programmer with strong domain expertise, feels constantly not good enough as programmer
Leading Lian	Superior who used to program back in the days tends to throw too many good ideas in the room
Master Michi	Experienced coder and architect, knows personal limits und unknowns



What Workflows Do We Need?

- Clear
- Simple
- Minimal
- Stable



How Should We Set-Up Our Code?

- Have Principles!
- Coding
- Architecture
- Stack



Principles I PEP20 The Zen of Python, by Tim Peters

>>> import this

Beautiful is better than ugly. Explicit is better than implicit.

Simple is better than complex. Complex is better than complicated.

Flat is better than nested. Sparse is better than dense.

Readability counts.

Special cases aren't special enough to break the rules. Although practicality beats purity.

Errors should never pass silently. Unless explicitly silenced.

In the face of ambiguity, refuse the temptation to guess.

There should be one - and preferably only one - obvious way to do it.

Although that way may not be obvious at first unless you're Dutch.

Now is better than never.

Although never is often better than *right* now.

If the implementation is hard to explain, it's a bad idea.

If the implementation is easy to explain, it may be a good idea.

Namespaces are one honking great idea -- let's do more of those!



Principles II Architecture

- Clear
- Simple
- Minimal
- Stable
- Scalable
- Independent

"Saw a post on LinkedIn about that framework..."



"What's wrong with about our currently used framework?..."

"New!!! Quack!! New better!!! Quack!!! Quaaaack!!!...



Principles III Software Stack

- Defined
- Minimal
- Stable
- Set up for experimentation
- Extendable



Stack: Python Standard Library

- Comes with Python
- Many useful libraries
- Look at standard lib first!





~ % cat .bash_profile # for brew search/ export HOMEBREM_GITHUB_API_TOKEN="<

Anaconda export PATH="\$PATH:/Users/hendorf/anaconda3/bin"

Jupyter
export PATH="\$PATH:/Users/hendorf/.local/bin"

The next line updates PATH for the Google Cloud SDK.
if [-f '/Users/hendorf/Downloads/googlecloud-sdk/path.bash.inc'; fi
cloud-sdk/path.bash.inc'; fi

>>> conda initialize >>> # !! Contents within this block are managed by 'conda init' !! _conda_setup="\$('/Users/hendorf/anaconda3/bin/conda' 'shell.bash' 'hook' 2> /dev/null)" if [\$? -eq 0]; then eval "\$ conda setup" else -f "/Users/hendorf/anaconda3/etc/profile.d/conda.sh"]; then if "/Users/hendorf/anaconda3/etc/profile.d/conda.sh" else PATH="/Users/hendorf/anaconda3/bin:\$PATH" expo fi fi unset __conda_setup # <<< conda initialize <<<</pre>

What About Environments?

- Spaces to program within
- Environment variables
- Config files in many places
- Bash, zsh vars setups on load
- Management with conda, pip,...
- Reproducibility
- IDE
- Environment clones to try new stuff

Menu of the Day: Spaghetti Code





What Makes Code Good or Bad?

- Readable
- Maintainable
- Shareable
- Reliable
- Documented



What's Code Readability?

- Concise, descriptive names
- Explicit code
- Split long lines, 2 pages on one screen
- Break complex or long into smaller functions
- No fancy stuff because you can
- Add line comments to provide context
- Linted with autopep8/black et al.
- Define your own style guide for stability



Tip: Dataclasses

- Be explicit if more than one return value
- Dataclasses provide useful interfaces with benefits





What Makes Code Maintainable?

- Consistent style and approaches
- Self-explanatory
- Single points of failure / DRY Principle
- Well origanised
- Tests



How Can I Share My Code?

- Use git
- Commit often
- Use a repository (service)
- .gitignore secrets
- Guidelines how to make pull requests
- Team reviews are good
- CI/CD with auto tests and publishing



What Makes Code Reliable?

- Testing
- Systems are not reliable use explicit Exceptions to handle
- Type hints
- Simplicity



What Does Documented Mean?

- READMEs
- Docstrings
- Line Comments
- Working Example Code
- Do not state the obvious
- Concise, not chatty
- Constantly update and refactor



Breaking the Rules!

- Know when to break the rules
- Have a good reason
- Be explicit about it





Jupyter Notebooks to Production



What About Jupyterlab?

- Exploratory Data Analysis
- Exploratory Code Execution
- Open Source
- Visualisation
- Formats / MD-Texts
- Ecosystem

Director of Artificial Intelligence and Autopilot Vision at Tesla



Andrej Karpathy 🤣 @karpathy · 28 Min.

so I accidentally held down something and deleted all cells in this jupyter notebook I've been building for ~2 months, and the "undo delete cell" isn't bringing them back. Lol.

(7) 319

226



1] 23



1

, ^,

I'm still shook. Some jupyter hotkey, somehow held down with my left palm, just iteratively deletes everything and undo doesn't bring them back (it creates an empty cell only). Hug your favorite notebooks and keep them safe 🤎

Q 24



Andrej Karpathy 🤣 @karpathy · 5 Min.

1] 6

thanks everyone, I was luckily able to find a snapshot in the .ipynb_checkpoints/ folder. You know that annoying thing you always add to the top of your .gitignore? turns out it can actually be useful :)

 \bigcirc 3 ① 1 0 75

Aug 30 2020, https://twitter.com/karpathy/status/1299972064426651650



No Comment.

KÖNIGSWEG



Any Downsides Using Notebooks?

- Code Completion
- Pointing Code Smells
- Code Versioning
- Notebooks are scripts
- Execution order
- Introspection
- Debugging

On my day job I move my hand up and down

Ah, easy! You're a Data Scientist scrolling notebooks!

Guessing the Occupation



What Common Antipatterns in Notebooks?

- Wall of code
- Non-local variables
- Redundancies
- Lack of versioning
- Missing documentation



What Can I Do About Wall of Code?

- Use multiple cells
- Break down into function
- Break down into multiple functions
- Classes
- Use modules



What About Non-local Variables?

- Hard to refactor
- Hard to move code with non-local variables
- Use function parameters
- Use classes and attributes





What Can I Do About Redundancies?

- Avoid them ;)
- Use functions
- DRY-principle
- Constant Refactoring



How Can I Debug?

- Use Introspection
- IDE provide good debuggers
- Print statements are ok
- Logging is better



	D 🔴	Debug	
Deb	oug: 🔄 🤤 NeuralStyleTrans	fer ×	* -
đ	Frames Variables Console		=
	Frames Va	iables 🛛 💫 Console	
	Frames Val	iables Console iables Console iables torch_content_image = {NoneType} None iables torch_style_image = {NoneType} None iables use_cuda = {bool} True iables vgg = {VGG} VGG(\n (conv1_1): Conv2d(3, 64, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1)) iables iable iables vgg = {VGG} VGG(\n (conv1_1): Conv2d(3, 64, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1)) iables iables iables conv1_1 = {Conv2d} Conv2d(64, 64, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1)) iables conv2_1 = {Conv2d} Conv2d(64, 64, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1)) iiiiiiiiiiiiiiiiiiiiiiiiiiiiiiiiiiii	
		<pre> [i] is_cdd = (bool) False [i] is_clad = (bool) False [i] is_clad = (bool) False [i] is_cquantized = (tople: 2) (3, 3) [i] is_cquantized = (tople: 2) (3, 3) [i] outcut_padding = (tople: 2) (0, 0) [i] is_adding = (tople: 2) (0, 0)</pre>	

Frames Variables Console = 2 ± ± ± ± * 1 = Frames Variables Console = 2 ± ± ± 1 * 1 = Frames Variables Console = 2 ± ± ± 1 * 1 = Frames Variables Console = 2 Console M < + > = targets = {list: 0} [] M < + > = targets = {list: 0} [] I corch_content_image = {NoneType} None I corch_style_image = {NoneType} None I corch_style_image = {NoneType} None I use_cuda = {bool} True * = vgg = {VGG} VGG(\n (conv1_1): Conv2d(3, 64, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1))\n (conv1_2): Conv2d(64, 64, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1)) * = conv1_1 = {Conv2d} Conv2d(64, 64, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1)) * = conv2_1 = {Conv2d} Conv2d(64, 128, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1)) * = conv2_1 = {Conv2d} Conv2d(64, 128, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1)) * = conv2_1 = {Conv2d} Conv2d(64, 128, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1)) * = conv2_1 = {Conv2d} Conv2d(64, 128, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1)) * = conv2_1 = {Conv2d} Conv2d(64, 128, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1))	=
<pre>Frames Variables Q Console Frames Variables Q Console Variables Q Variables Q Variables Vari</pre>	
 M < + > = targets = {list: 0} [] torch_content_image = {NoneType} None torch_style_image = {NoneType} None torch_style_image = {NoneType} None use_cuda = {bool} True v = vgg = {VGG} VGG (\n (conv1_1): Conv2d(3, 64, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1))\n (conv1_2): Conv2d(64, 64, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1)) conv1_1 = {Conv2d} Conv2d(3, 64, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1)) conv1_2 = {Conv2d} Conv2d(64, 64, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1)) conv2_1 = {Conv2d} Conv2d(64, 128, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1)) 	
 use_cuda = {bool} True use_cuda = {bool}	
Conv1_2 = {Conv2d} Conv2d(64, 64, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1)) Conv2_1 = {Conv2d} Conv2d(64, 128, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1)) Conv2_1 = {Conv2d} Conv2d(64, 128, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1)) Conv2_1 = {Conv2d} Conv2d(64, 128, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1)) Conv2_1 = {Conv2d} Conv2d(64, 128, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1)) Conv2_1 = {Conv2d} Conv2d(64, 128, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1)) Conv2_1 = {Conv2d} Conv2d(64, 128, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1)) Conv2_1 = {Conv2d} Conv2d(64, 128, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1))	(3, 3), stride=(1, 1), padding= View
> = conv2_1 = {conv2a} Conv2a(64, 128, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1)	
Conv2d {Conv2d {Conv{Conv{Conv2d {Conv2d {Conv2d {Conv{Conv2d {Conv{Conv{Conv{Conv{Con	0.0064, -0.1856, 0.1297, View
Code fragment: Code fr	297, -0.0177, -0.2785, -0.019 View
Result: 	
<pre>> = dtype = {dtype} forch.rioat32 0 grad = {NoneType} None 0 grad_fn = {NoneType} None > = 0 is_cuda = {bool} False > 0 is leaf = {bool} True</pre>	0.0477 0.2795 0.0404
<pre>> * > * Close Evaluate Outp_patcnes = {DU0} raise</pre>	, -0.0177, -0.2700, -0.0104, VIEW
groups = {int} 1 [in_channels = {int} 128	
<pre> 21 out_channels = {int} 128 22 output_padding = {tuple: 2} (0, 0) 23 padding = {tuple: 2} (1, 1) 24 25 25 25 25 25 25 25 25 25 25 25 25 25</pre>	





Jupyter and an IDE!

- Best of both
- Very productive
- Focus on the strong suits



How Can I Have Shorter, Concise Notebooks?

- Functions
- Classes
- Modules
- Magics!
- Relevant code only notebooks





Anything Else?

- Type Hints
- Decomposition
- Parallelizing execution
- Parameters as config
- Sub-classing
- Multi modules
- Descriptors
- Decorators



What Are the Main Take-Aways?

- Good coding is as good engineering an essential success factor
- Python offers several ways to achieve the required quality level for a successful implementation
- Beyond writing quality code, there are other aspects to consider such as personas, standardization and architecture
- Mentors will speed up the process a lot

ALEXANDER C. S. HENDORF

MANAGING PARTNER & PRINCIPAL CONSULTANT DATA SCIENCE & AI AT KÖNIGSWEG.

Thank you!

