

A Data Scientist's Guide to Code Reviews

PyCon/PyData 2022



Doing code reviews

Pros

- They'll improve the code clarity
- They might uncover errors
- I'll probably learn something while doing it
- Literally everyone says that I should

Cons

- I don't want to

⇒ I will ****not**** do code reviews

Doing code reviews on data science work

Pros

- They'll improve the code clarity
- They might uncover errors
- I'll probably learn something while doing it

Cons

- Hardly anyone says that I should
- The code will likely not end up in a live system as-is
- The work is a one-off thing

⇒ I will ****not**** do code reviews

However, having someone else review your work is as important in data science as in software engineering.

What are code reviews for?

- Verifying that the specified goal is achieved
- Uncovering errors and misunderstandings
- Knowledge transfer
- Feedback for architectural or design decisions
- Improving your code & coding practice

The traditional code review practice is not applicable to “typical” data science work.

Different focus

Software Engineering

- Is the artifact functional?
- Are there bugs?
- Are coding guidelines & quality standards met?
- Can someone else than the author work on the artifact?

⇒ Code Review

Data Science

- Is the chosen approach comprehensible & clear?
- Have data peculiarities been taken into account?
- Are the results plausible?
- Can someone else than the author explain the concept?

⇒ Peer Review

What are code reviews for in data science?

- Verifying that the specified goal is achieved ✓
- Uncovering errors and misunderstandings
- Knowledge transfer
- Feedback for architectural or design decisions
- Improving your code & coding practice

What are code reviews for in data science?

- Verifying that the specified goal is achieved ✓
- Uncovering **logical** errors and misunderstandings ✓
- Knowledge transfer ✓
- Feedback for architectural or design decisions
- Improving your code & coding practice

What are code reviews for in data science?

- Verifying that the specified goal is achieved ✓
- Uncovering **logical** errors and misunderstandings ✓
- Knowledge transfer ✓
- Feedback for ~~architectural or design decisions~~ **approach** ✓
- Improving your code & coding practice

What are code reviews for in data science?

- Verifying that the specified goal is achieved ✓
- Uncovering **logical** errors and misunderstandings ✓
- Knowledge transfer ✓
- Feedback for ~~architectural or design decisions~~ **approach** ✓
- ~~Improving your code & coding practice~~ **Reproducibility** ↺

Code review checklist

- ☐ Overview over present files and the task
 - changelist
 - MR's description
 - accompanying ticket (when working with a ticket system, e.g. JIRA)
- ☐ Run the code and reproduce the results
 - ☐ [optional] if GitLab CI is used it might be worth checking the pipeline
 - ! fixing the pipeline is the author's responsibility
- ☐ Ensure comprehension: ask, ask, ask

Why has the author decided to do XY, chosen package A instead of B, selected model 42 as baseline,...?

github.com/awoerner92/talks/pydata2022/checklist.md

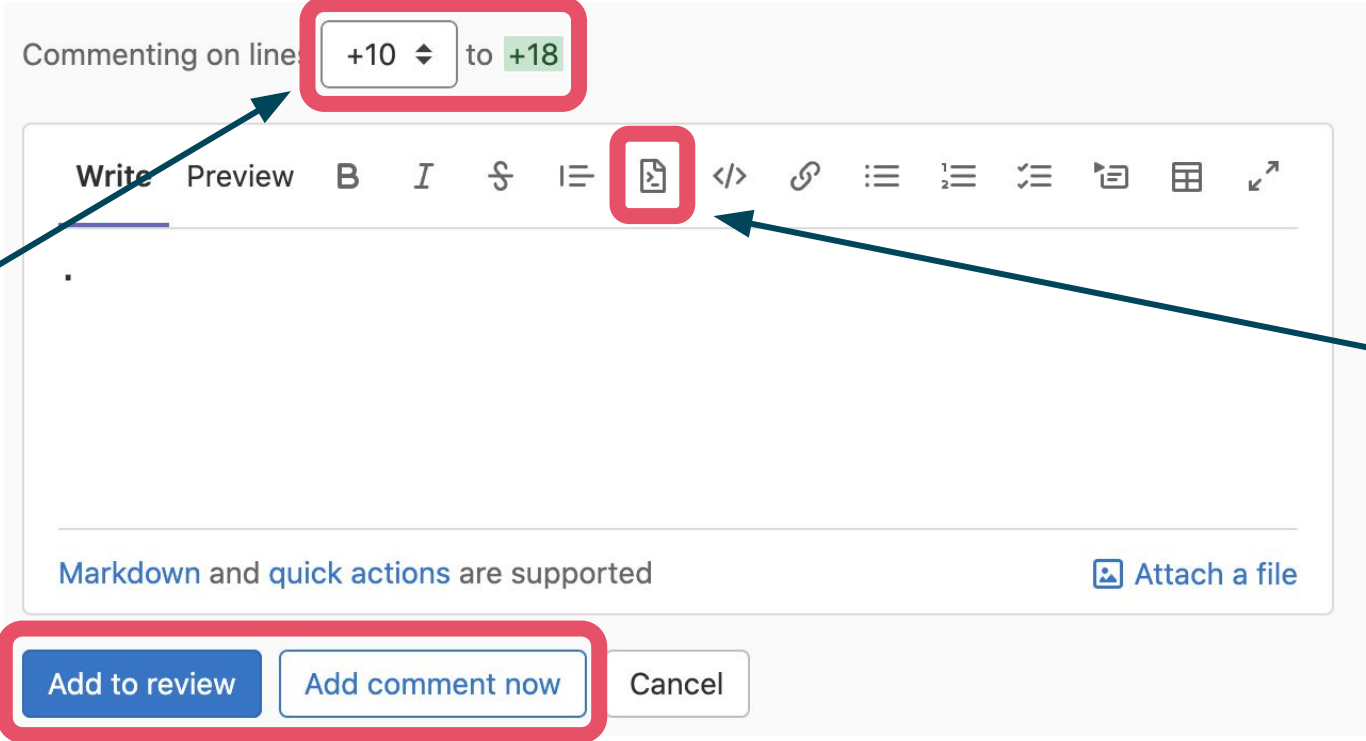
How to do code reviews?

- How to Do Code Reviews Like a Human
 - <https://mtlynch.io/human-code-reviews-1/>
 - <https://mtlynch.io/human-code-reviews-2/>
- How to Make Your Code Reviewer Fall in Love With You
 - <https://mtlynch.io/code-review-love/>

Useful Git functionality: pre-commit hooks

- Runs pre-defined set of tools with every commit
- Tools:
 - Jupyter notebook conversion: nbconvert
 - Code formatter: black, isort
 - Linter: flake8

Useful GitLab functionalities: Comment field



The screenshot shows the GitLab comment field interface. The 'Commenting on line' dropdown is set to '+10 to +18'. The 'Write' tab is active, showing a rich text editor with a toolbar containing icons for bold, italic, strikethrough, link, code suggestion, and other formatting options. The 'code suggestion' icon (a document with a code symbol) is highlighted with a red box and an arrow pointing to the text '(GitHub: ```suggestion```)' on the right. The 'Add to review' and 'Add comment now' buttons are highlighted with a red box and an arrow pointing to the text 'add comment to batch or comment immediately' at the bottom. An arrow points from the text 'comment & mark multiple lines' to the line range dropdown.

comment & mark multiple lines

code suggestion

(GitHub: ````suggestion````)

add comment to batch or comment immediately

Useful GitLab functionalities: Mark viewed

Overview 55 Commits 19 Pipelines 14 **Changes 43**

✓ All threads resolved

Compare version 1 and latest version Show latest version 39 files +939 -909 Expand all files

> tests/test_model.py 0 → 100644 +252 -0 **Viewed**

- Collapses the file
- Helps to keep an overview

Thank you!



Alexandra Wörner

Data Scientist

alexandra.woerner@scieneers.de



@alex_woerner

Questions?