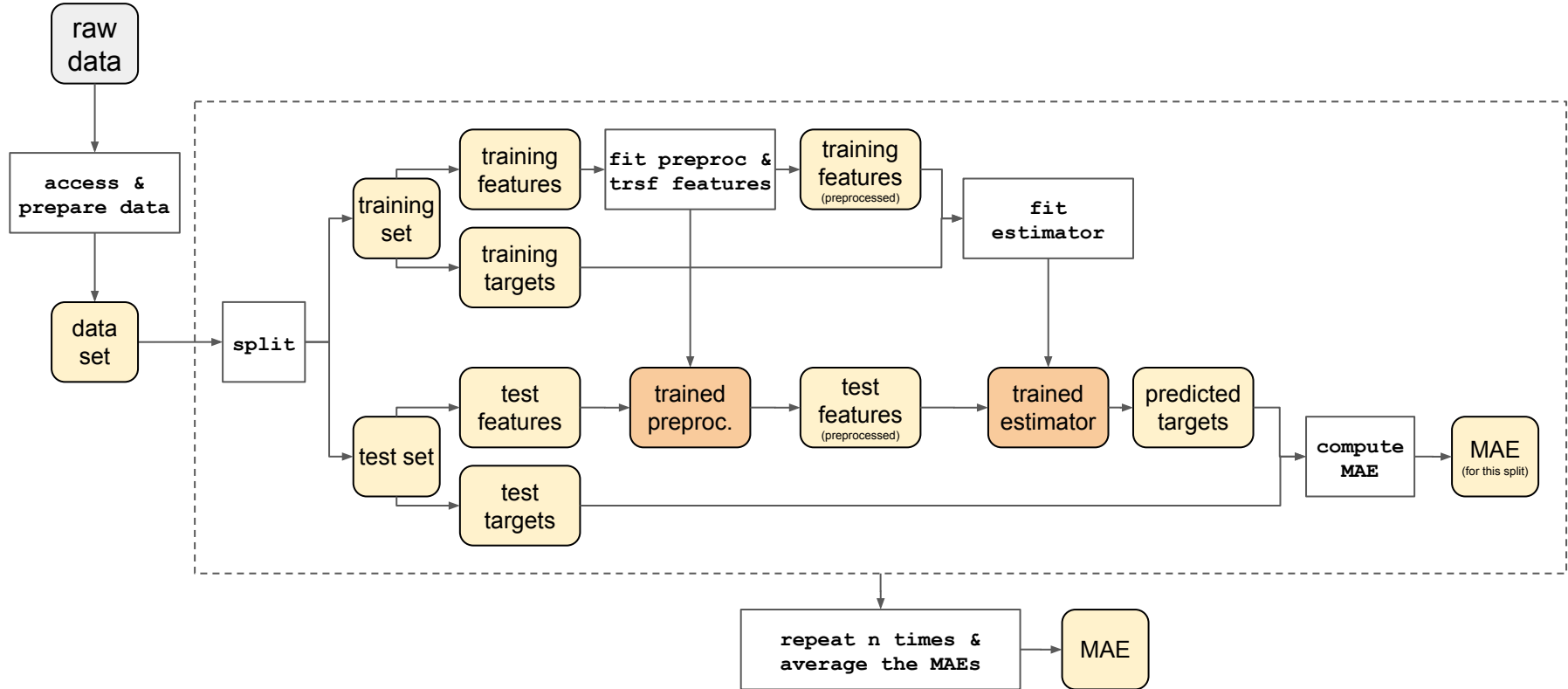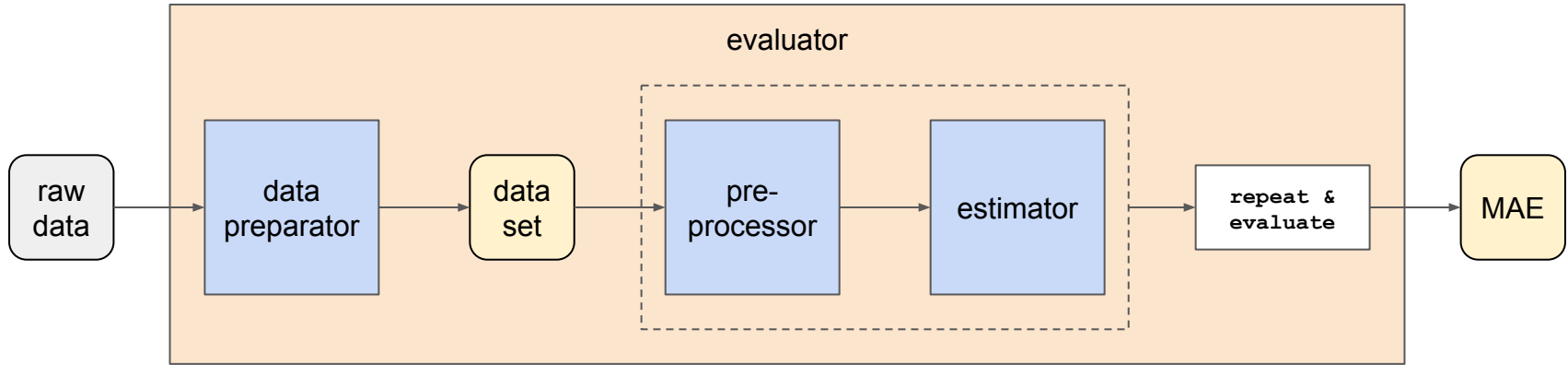# M05 mini-project

UniDistance - Master AI - SS2023

Benjamin Décaillet
Valentin Décaillet
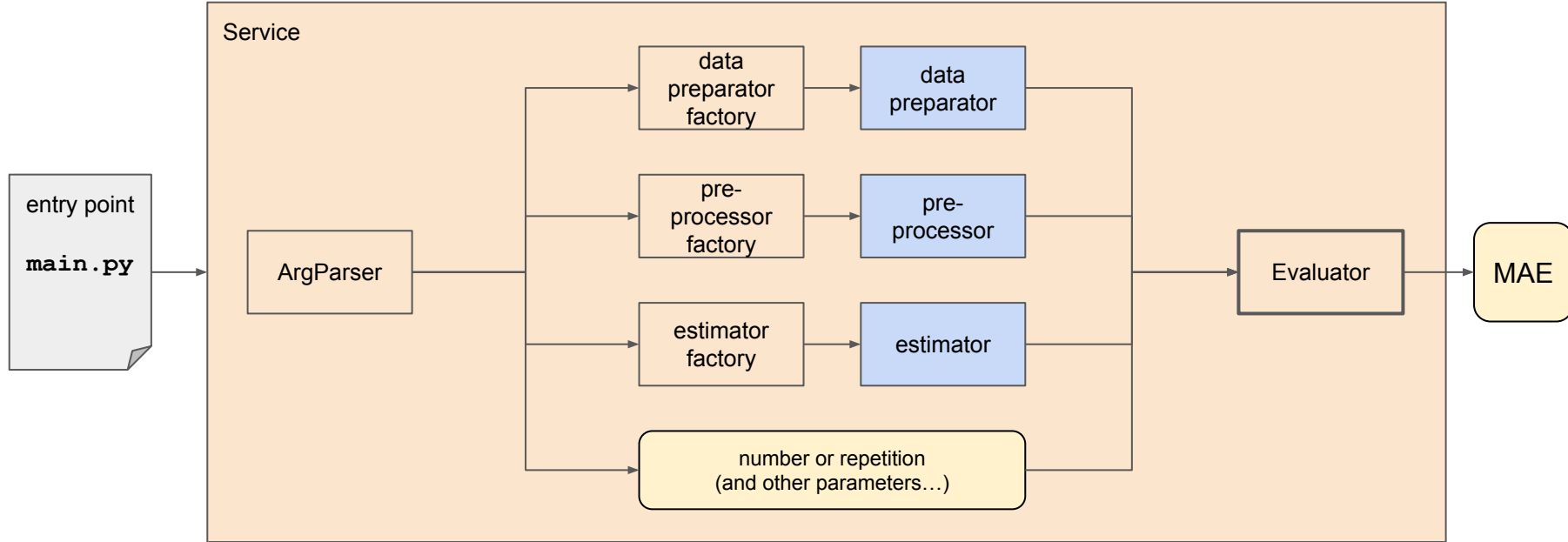
# Workflow - functional overview

# Workflow - class diagram (high level)



- advantage: the evaluator provides flexibility (swap/customize/fine-tune the "blue blocs")
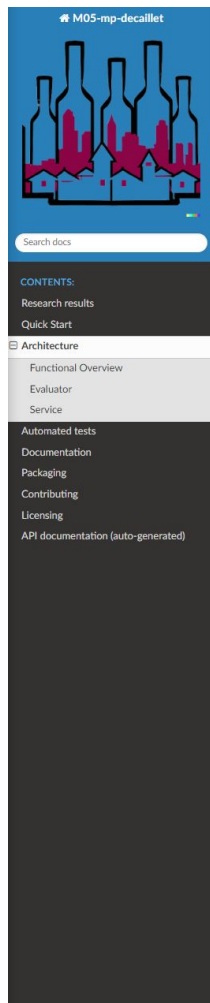- disadvantage: constructing an evaluator is technically complex

# Workflow - construction of the evaluator



- Service is easy to call from any entry point (typically a one-liner: `Service().run()`)
- The entire point is to instantiate an Evaluator (which we can then run…)

# Documentation

- ## Use sphinx
  - ### Accept .rst and .md files
  - ### Theme "read the docs"
  - ### Auto-generated from code and docstrings
- ## Links from readme.md
  - ### Badges to sphinx doc
- ## Sphinx deployed to github pages
  - ### CI-deploy, branch main only
  - ### GitHub shenanigans
    - #### not a jekyll site
    - #### permissions for github bot

# Documentation

- Use sphinx
  - Accept .rst and .md files
  - Theme "read the docs"
  - Auto-generated from code and docstrings
- Links from readme.md
  - Badges to sphinx doc
- Sphinx deployed to github pages
  - CI-deploy, branch main only
  - GitHub shenanigans
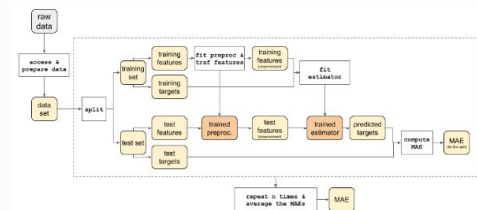    - not a jekyll site
    - permissions for github bot

# Documentation

- Use sphinx
  - Accept .rst and .md files
  - Theme "read the docs"
  - Auto-generated from code and docstrings
- Links from readme.md
  - Badges to sphinx doc
- Sphinx deployed to github pages
  - CI-deploy, branch main only
  - GitHub shenanigans
    - not a jekyll site
    - permissions for github bot

```yaml
12  jobs:
13    main-ci-action:
14      runs-on: ubuntu-latest
15
16      permissions:
17        # required to be able to push sphinx doc
18        # to the branch 'gh-pages'
19        contents: write
20
21      steps:
22        - uses: actions/checkout@v3

      [...]

49        - name: Build documentation
50          run: |
51            rm -rf ./doc/apidoc
52            sphinx-apidoc src/ -o ./doc/apidoc --no-toc --separate --module-first
53            sphinx-build doc sphinx
54            touch sphinx/.nojekyll  # Dear Github, this is not a jekyll site. Chill.
55
56        - name: Deploy documentation to GitHub Pages
57          if: success()  &&  github.ref == 'refs/heads/main'
58          uses: crazy-max/ghaction-github-pages@v3
59          with:
60            target_branch: gh-pages
61            build_dir: sphinx
62          env:
63            GITHUB_TOKEN: ${{ secrets.GITHUB_TOKEN }}
```

# Version control

- GIT
  - feature branches
  - rebase before pull requests
    inspired from Gitflow
  - commit message convention
    inspired from Angular
- *It's all explained in the doc*

# Version control



**M05-mp-decaillet**

Search docs

**CONTENTS:**

🏠 / Contributing

View page source

## Contributing

### Commit messages convention

As a general rule, follow this syntax for commit messages. (Emphasis on the header-line, not so much on body and footer)

### Branching strategy

As a general rule:

- avoid pushing directly to branch `dev`
- push your changes to a feature branch (named `feature/name-in-kebab-case`) and create a pull request when you're done.

Branch `main` **is for releases only**; never push directly to main, always merge from `dev`.
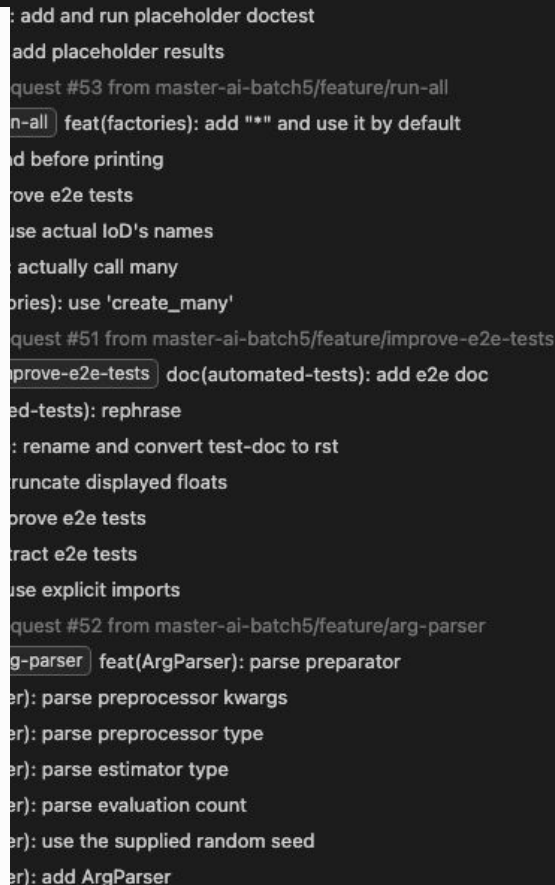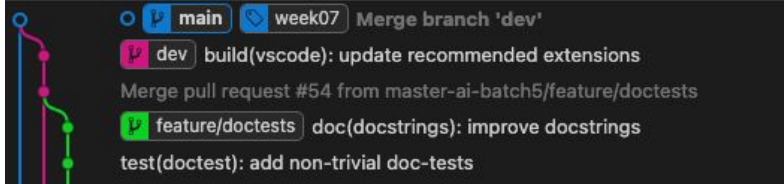In principle a (merge) commit to `main` *is* a release.

### Linter

The python code in this project must match autopep8 and isort linting/formatting rules.

GitHub actions will enforce these rules.

### Reformat from command line

- activate your virtualenv: `conda activate m05-mp-decaillet`
- apply autopep8 to all local files: `autopep8 --max-line-length=120 --recursive . -aaa --in-place`
- apply isort to all local files: `isort .`



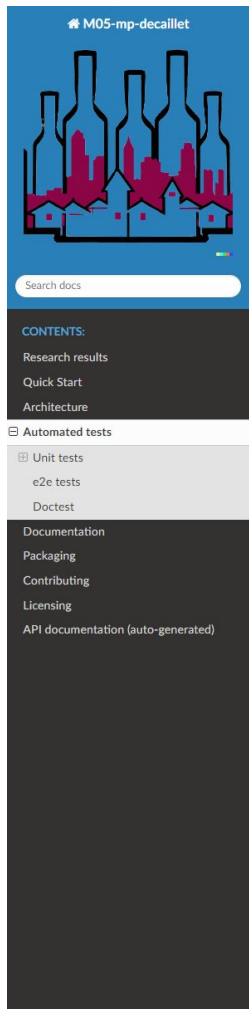| | | | |
|---|---|---|---|
| ○ 🔵 main 🏷 week07 | Merge branch 'dev' | | 19 Mar 2023 14:26 |
| 🔵 dev | build(vscode): update recommended extensions | | 19 Mar 2023 14:22 |
| | Merge pull request #54 from master-ai-batch5/feature/doctests | | 19 Mar 2023 12:48 |
| 🟢 feature/doctests | doc(docstrings): improve docstrings | | 16 Mar 2023 01:02 |
| | test(doctest): add non-trivial doc-tests | | 16 Mar 2023 00:59 |
| | ...: add and run placeholder doctest | | 16 Mar 2023 00:56 |
| | ...add placeholder results | | 15 Mar 2023 15:37 |
| | ...quest #53 from master-ai-batch5/feature/run-all | | 19 Mar 2023 12:44 |
| ...n-all | feat(factories): add "*" and use it by default | | 16 Mar 2023 00:36 |
| | ...d before printing | | 16 Mar 2023 00:36 |
| | ...rove e2e tests | | 16 Mar 2023 00:36 |
| | ...use actual IoD's names | | 16 Mar 2023 00:36 |
| | ... actually call many | | 16 Mar 2023 00:36 |
| | ...ories): use 'create_many' | | 16 Mar 2023 00:36 |
| | ...quest #51 from master-ai-batch5/feature/improve-e2e-tests | | 19 Mar 2023 12:33 |
| ...prove-e2e-tests | doc(automated-tests): add e2e doc | | 15 Mar 2023 22:25 |
| | ...ed-tests): rephrase | | 15 Mar 2023 22:18 |
| | ...: rename and convert test-doc to rst | | 15 Mar 2023 22:04 |
| | ...runcate displayed floats | | 15 Mar 2023 16:25 |
| | ...rove e2e tests | | 15 Mar 2023 16:25 |
| | ...tract e2e tests | | 15 Mar 2023 16:25 |
| | ...use explicit imports | | 15 Mar 2023 16:25 |
| | ...quest #52 from master-ai-batch5/feature/arg-parser | | 19 Mar 2023 12:28 |
| ...g-parser | feat(ArgParser): parse preparator | | 13 Mar 2023 15:37 |
| | ...er): parse preprocessor kwargs | | 13 Mar 2023 15:23 |
| | ...er): parse preprocessor type | | 5 Mar 2023 23:50 |
| | ...er): parse estimator type | | 13 Mar 2023 15:11 |
| | ...er): parse evaluation count | | 5 Mar 2023 23:04 |
| | ...er): use the supplied random seed | | 13 Mar 2023 11:00 |
| | ...er): add ArgParser | | 13 Mar 2023 10:38 |

# Code Sharing (GitHub)

- Organization created: "master-ai-batch5"
- Repository under the organization → Role management for Repository
- Protection on branches:
  - main:
    - Require linear history
    - Can only be updated from dev
  - dev: Set up as default
    - Require a pull request before merging
    - Require approvals
    - Dismiss stale pull request approvals when new commits are pushed
    - Require conversation resolution before merging
    - Allow force pushes from certain member only
  - feature branches: no security
- Issues used to track tasks and assignments

# Automated testing

- unit-tests
- e2e tests
- doctests

# Automated testing - unit-tests

- native 'unittest' framework
- coverage 100%

| Name | Stmts | Miss | Cover | Missing |
|---|---|---|---|---|
| decm05\__init__.py | 3 | 0 | 100% | |
| decm05\arg_parser.py | 61 | 0 | 100% | |
| decm05\estimating\__init__.py | 5 | 0 | 100% | |
| decm05\estimating\contract.py | 10 | 0 | 100% | |
| decm05\estimating\decision_tree_estimator.py | 9 | 0 | 100% | |
| decm05\estimating\factory.py | 21 | 0 | 100% | |
| decm05\estimating\linear_regression_estimator.py | 8 | 0 | 100% | |
| decm05\estimating\sklearn_estimator_base.py | 23 | 0 | 100% | |
| decm05\evaluator.py | 37 | 0 | 100% | |
| decm05\preparator\__init__.py | 5 | 0 | 100% | |
| decm05\preparator\base_preparator.py | 17 | 0 | 100% | |
| decm05\preparator\boston_preparator.py | 14 | 0 | 100% | |
| decm05\preparator\contract.py | 8 | 0 | 100% | |
| decm05\preparator\factory.py | 34 | 0 | 100% | |
| decm05\preparator\wine_preparator.py | 36 | 0 | 100% | |
| decm05\preprocessing\__init__.py | 6 | 0 | 100% | |
| decm05\preprocessing\contract.py | 10 | 0 | 100% | |
| decm05\preprocessing\factory.py | 31 | 0 | 100% | |
| decm05\preprocessing\min_max_preprocessor.py | 8 | 0 | 100% | |
| decm05\preprocessing\polynomial_preprocessor.py | 13 | 0 | 100% | |
| decm05\preprocessing\sklearn_preprocessor_base.py | 24 | 0 | 100% | |
| decm05\preprocessing\standard_preprocessor.py | 8 | 0 | 100% | |
| decm05\service.py | 24 | 0 | 100% | |
| TOTAL | 415 | 0 | 100% | |

91/91 tests passed (100%)

- M05-mp-decaillet
  - unit_tests
    - estimating
      - test_decision_tree_estimator.py
        - TestDecisionTreeEstimator
          - test__name
          - test__happy_path
          - test__fails_it_not_fit
          - test__fails_if_columns_dont_match
      - test_factory.py
      - test_linear_regression_estimator.py
    - preparator
    - preprocessing
    - test_arg_parser.py
    - test_evaluator.py
    - test_service.py

# Automated testing - e2e tests

- a bash script

```bash
#!/usr/bin/env bash
set -euo pipefail
IFS=$'\n\t'
# why do I bother with the above ?
# http://redsymbol.net/articles/unofficial-bash-strict-mode/

#===============================================================
#title          : e2e.sh
#description     : run e2e tests
#date           : 2022/09/15
#usage          : e2e.sh
#notes          :
#===============================================================

echo "Running e2e tests..."


python main.py --seed=42 \
               --dataset=boston \
               > output.log 2>> error.log
diff -q <(cat <<EOF
dataset preprocessor       estimator  evaluation count   MEAN ABSOLUTE ERROR
 boston      min-max linear-regression             3             3.4924
 boston      min-max     decision-tree             3             2.7763
 boston     standard linear-regression             3             3.4944
 boston     standard     decision-tree             3             3.1092
 boston   polynomial linear-regression             3             5.0997
 boston   polynomial     decision-tree             3             2.9869
EOF
) output.log || (echo "Output 1 does not match expected output" \
                        && cat output.log \
                        && exit 1)

python main.py --seed=42 \
               --dataset=wines \
               --preprocessor-type=min-max \
               --estimator-type=decision-tree \
               > output.log 2>> error.log
diff -q <(cat <<EOF
dataset preprocessor       estimator  evaluation count  MEAN ABSOLUTE ERROR
  wines      min-max decision-tree                 3             0.5765
EOF
) output.log || (echo "Output 2 does not match expected output" \
                        && cat output.log \
                        && exit 2)
```

# Automated testing - doctests

- sphinx.ext.doctest

```
Wine Quality
------------

.. testcode::

    from decm05 import Service
    Service(["--dataset=wines", "--seed=42"]).run()

.. testoutput::

    dataset preprocessor        estimator  evaluation count  MEAN ABSOLUTE ERROR
      wines       min-max linear-regression               3               0.5695
      wines       min-max    decision-tree               3               0.5800
      wines      standard linear-regression               3               0.5729
      wines      standard    decision-tree               3               0.5850
      wines    polynomial linear-regression               3               0.5582
      wines    polynomial    decision-tree               3               0.5790


Boston House Prices
-------------------

.. testcode::

    from decm05 import Service
    Service(["--dataset=boston", "--seed=42"]).run()

.. testoutput::

    dataset preprocessor        estimator  evaluation count  MEAN ABSOLUTE ERROR
     boston       min-max linear-regression               3               3.4924
     boston       min-max    decision-tree               3               2.7763
     boston      standard linear-regression               3               3.4944
     boston      standard    decision-tree               3               3.1092
     boston    polynomial linear-regression               3               5.0997
     boston    polynomial    decision-tree               3               2.9869
```
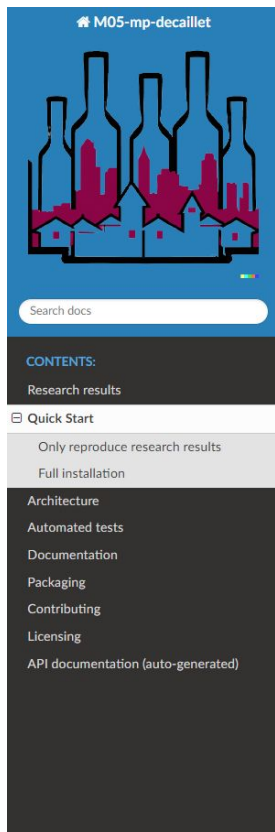
# Environment management

- Python 3.11
  - using conda virtual environments
  - dependencies pinned in `build-requirements.txt` and `dev-requirements.txt`
- "**quick start**" section in documentation:
  - project-specific command lines
  - links to 3rd parties

# Continuous Integration (CI)

- run for each
  - push (branch / version tags)
  - pull request
- one job, several steps
  - cf image
  - abort on failure
- conditional deploy steps
  - deploy doc only for branch `main`
  - deploy to PyPI only for version tags

# Packaging

- Standard python package repository PyPi
  - only deployed on test https://test.pypi.org/
  - package name decm05
- setuptools
- dependencies pinned with two requirements files
  - "build-requirements.txt" contains runtime dependencies
  - "dev-requirements.txt" contains all other dependencies
- Semantic version numbering i.e: "1.0.0"
  - tag must match [0-99].+[0-99].+[0-99]
  - without "alpha,beta, etc" version candidates
- CI uploading new version directly on test PyPi

# Licensing

- MIT licensing
- Licenses of dependencies
  - Auto-generated with [pip-licenses](pip-licenses)
  - See our documentation for generation
- Licenses of data
  - *.names file containing information about licensing and content included in package

# Sources

M05: Open Science and Ethics (Dr. André Anjos, Flavio Tarsetti, Joël Dumoulin)

See [our documentation](#) for more information

Any questions?

1. Dr. Anjos, A., Tarsetti, F. & Dumoulin, J. (2023). M05: Open Science and Ethics [Teaching Course]. IDIAP Research Institute. Martigny, Switzerland.
2. "Question Mark Question Response" by OpenClipart-Vectors from Pixabay. Available at: https://pixabay.com/illustrations/question-mark-question-response-1020165/
3. Github pages by Décaillet V. & Décaillet B. at "https://master-ai-batch5.github.io/M05-mp-decaillet/index.html"