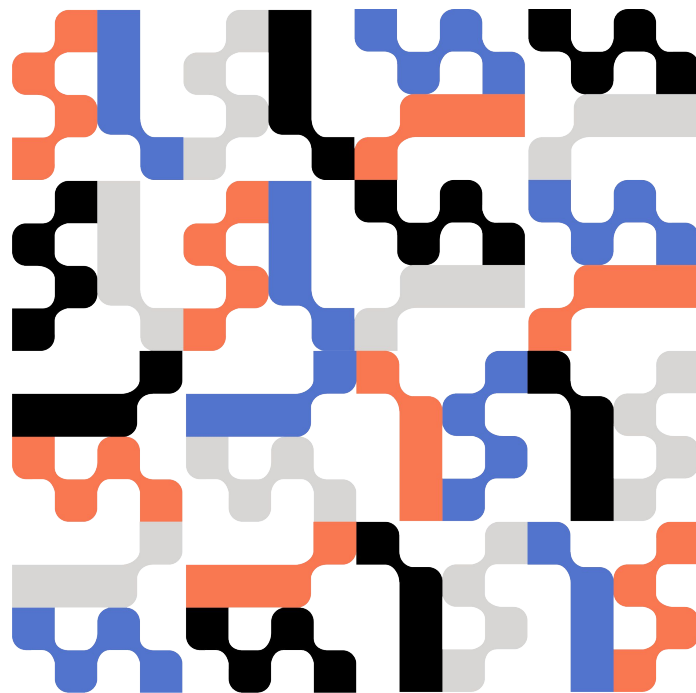


# FawltyDeps

Finding **undeclared** and **unused** dependencies  
in your notebooks and projects

Johan Herland – [@jherland](#)

FOSDEM 2025



# Replication crisis

18 languages

- Contents [hide]
- (Top)
- > Background
- > Prevalence
- > Causes
- > Consequences
- > Remedies
- See also
- Notes
- References
- Further reading

## Data science

46 languages

- Contents [hide]
- (Top)
- > Foundations
  - Relationship to statistics
- > Etymology
  - Early usage
  - Modern usage
- Data Science and Data Analysis
- History
- See also
- References

Article Talk Read Edit View history Tools

From Wikipedia, the free encyclopedia

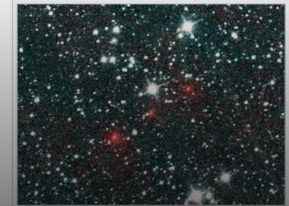
*Not to be confused with information science.*

**Data science** is an interdisciplinary academic field<sup>[1]</sup> that uses statistics, scientific computing, scientific methods, processes, algorithms and systems to extract or extrapolate knowledge and insights from noisy, structured, and unstructured data.<sup>[2]</sup>

Data science also integrates domain knowledge from the underlying application domain (e.g., natural sciences, information technology, and medicine).<sup>[3]</sup> Data science is multifaceted and can be described as a science, a research paradigm, a research method, a discipline, a workflow, and a profession.<sup>[4]</sup>

Data science is a "concept to unify statistics, data analysis, informatics, and their related methods" to "understand and analyze actual phenomena" with data.<sup>[5]</sup> It uses techniques and theories drawn from many fields within the context of mathematics, statistics, computer science, information science, and domain knowledge.<sup>[6]</sup> However, data science is different from computer science and information science. Turing Award winner Jim Gray imagined data science as a "fourth paradigm" of science (empirical, theoretical, computational, and now data-driven) and asserted that "everything about science is changing because of the impact of information technology" and the data deluge.<sup>[7][8]</sup>

A **data scientist** is a professional who creates programming code and combines it with statistical knowledge to create insights from data.<sup>[9]</sup>



The existence of Comet NEOWISE (here depicted as a series of red dots) was discovered by analyzing astronomical survey data acquired by a space telescope, the Wide-field Infrared Survey Explorer.

### Foundations

# Computational reproducibility of Jupyter notebooks from biomedical publications

Sheeba Samuel<sup>1,2\*†</sup> and Daniel Mietchen<sup>3,4,5†‡</sup>

<sup>1</sup>Heinz-Nixdorf Chair for Distributed Information Systems, Friedrich Schiller University Jena, Germany and

<sup>2</sup>Michael Stifel Center Jena, Germany and <sup>3</sup>Ronin Institute, Montclair, New Jersey, United States and <sup>4</sup>Institute for Globally Distributed Open Research and Education (IGDORE) and <sup>5</sup>FIZ Karlsruhe — Leibniz Institute for Information Infrastructure, Berlin, Germany

\*sheeba.samuel@uni-jena.de

†These authors contributed equally to this work

‡daniel.mietchen@ronininstitute.org

## Abstract

**Background** Jupyter notebooks facilitate the bundling of executable code with its documentation and output in one interactive environment, and they represent a popular mechanism to document and share computational workflows, including for research publications. The reproducibility of computational aspects of research is a key component of scientific reproducibility but has not yet been assessed at scale for Jupyter notebooks associated with biomedical publications.

**Approach** We address computational reproducibility at two levels: (1) Using fully automated workflows, we analyzed the computational reproducibility of Jupyter notebooks associated with publications indexed in the biomedical literature repository PubMed Central. We identified such notebooks by mining the article's full text, trying to locate them on GitHub and attempting to

333v1 [cs.DL] 11 Aug 2023

22,578 Python notebooks

70% had declared dependencies

46% could install dependencies

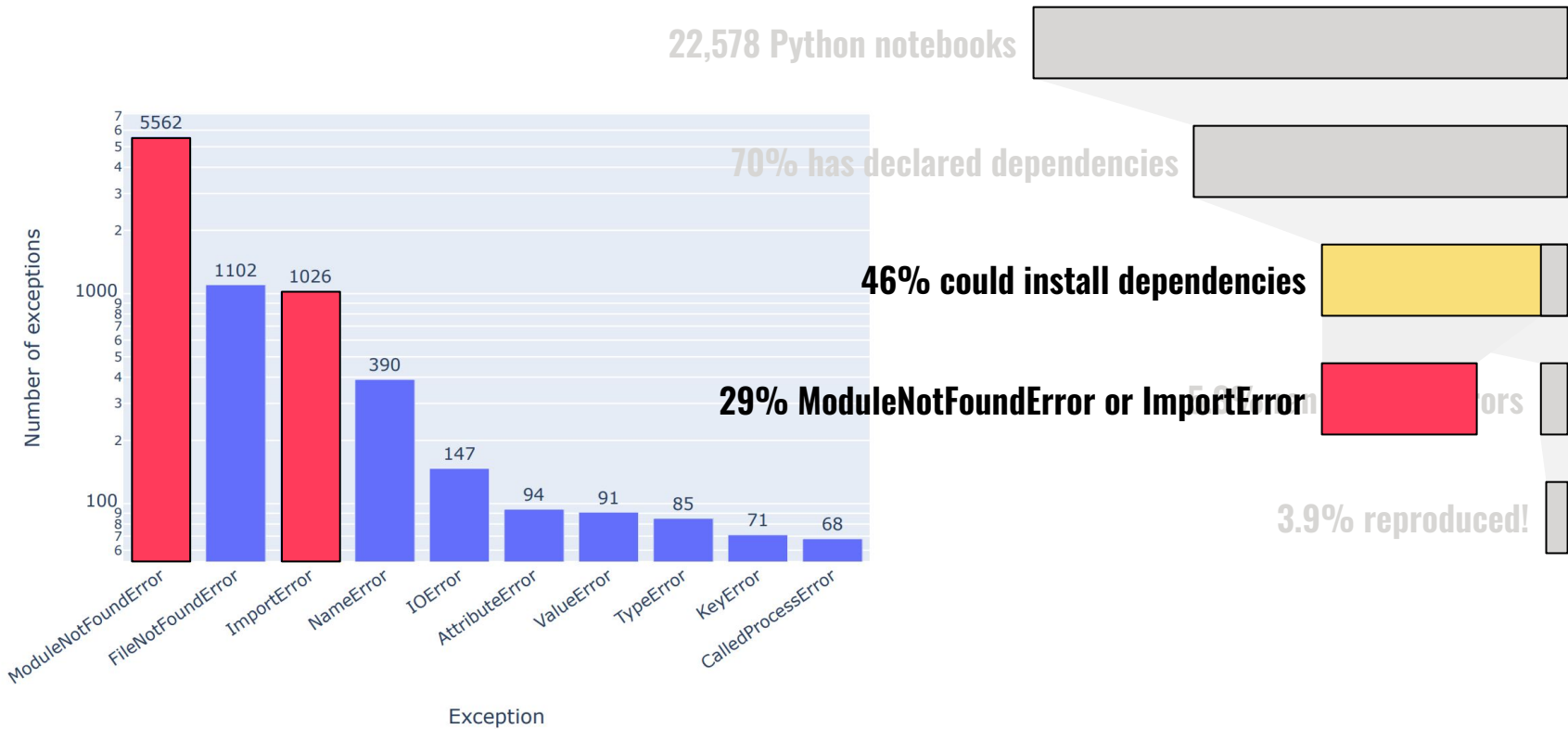
5.3% ran without errors

3.9% reproduced!

**Results** Out of 27,271 Jupyter notebooks from 2,660 GitHub repositories associated with 3,467 publications, 22,578 notebooks were written in Python, including 15,817 that had their dependencies declared in standard requirement files and that we attempted to re-run automatically. For 10,388 of these, all declared dependencies could be installed successfully, and we re-ran them to assess reproducibility. Of these, 1,203 notebooks ran through without any errors, including 879 that produced results identical to those reported in the original notebook, and 324 for which our results differed from the originally reported ones. Running the other notebooks resulted in exceptions.

notebooks resulted in exceptions

reproduced! 879 that produced results identical to those reported in the original notebook, and 324 for which our results differed from the originally reported ones.



**Figure 19.** Exceptions occurring in Jupyter notebooks in our corpus. See Table 5 for information about the nature of these errors and potential fixes.

**30% failed to declare dependencies**

**24% failed to install dependencies**

**29% failed to import dependencies**

**>80% fail because of missing or incorrect dependency declarations.**

22,578 Python notebooks

70% has declared dependencies

46% could install dependencies

29% ModuleNotFoundError or ImportError

3.9% reproduced!

PAPER

## Computational reproducibility of Jupyter notebooks from biomedical publications

Sheeba Samuel<sup>1,2\*†</sup> and Daniel Mietchen<sup>3,4,5†‡</sup>



# Less is More? An Empirical Study on Configuration Issues in Python PyPI Ecosystem

Yun Peng

The Chinese University of Hong Kong  
Hong Kong, China  
ypeng@cse.cuhk.edu.hk

Ruida Hu

Harbin Institute of Technology  
Shenzhen, China  
200111107@stu.hit.edu.cn

Ruohe Wang

Harbin Institute of Technology  
Shenzhen, China  
200110930@stu.hit.edu.cn

Cuiyun Gao\*

Harbin Institute of Technology  
Shenzhen, China  
gaocuiyun@hit.edu.cn

Shuqing Li

The Chinese University of Hong Kong  
Hong Kong, China  
sqli21@cse.cuhk.edu.hk

Michael R. Lyu

The Chinese University of Hong Kong  
Hong Kong, China  
lyu@cse.cuhk.edu.hk

[cs.SE] 19 Oct 2023

### ABSTRACT

Python is the source code for diverse third-party libraries. However, the utility of these libraries is often hindered by configuration conflicts in their dependencies. We study the dependency conflicts to automatically infer version-level checks and inference, based on the assumption that configurations of libraries in the PyPI ecosystem are correct. However, our study reveals that this assumption is not universally valid, and relying solely on version-level checks proves inadequate in ensuring compatible run-time environments.

Only 5,371 (65%) libraries, comprising 131,720 (39%) releases, successfully pass all three checks of PYCON. This indicates that approximately 30% of libraries and 51% of releases on the PyPI platform can be installed but may encounter source-level compatibility problems.

### PFRL.

in fostering its prosperity. The accessibility and utility of Python are further amplified by the public libraries available on the Python Package Index (PyPI) platform. With over 470 thousand Python

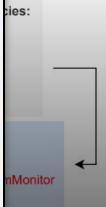


Table 2: Configuration issues detected by PyCON. There may be multiple issues occurring in one release.

Category	Issue	Check	#Releases	Fatal?	Possible Reasons
Incomplete Configuration	<u>Missing configuration files</u>	Installation Check	251	✓	Missing required information
	<u>Missing required libraries for setup</u>	Installation Check	3,318	✓	
	Missing Python versions	Dependency Check	55,138	✗	
	<u>Missing required libraries for direct imports</u>	Import Validation	142,521	✗	
Incorrect Configuration	Dependency conflicts in setup	Installation Check	6,318	✓	Unsolvable constraints
	Incorrect Python versions	Installation Check	4,155	✓	Incorrect dependencies
	Other run-time Errors in setup	Installation Check	3,464	✓	Missing files
	Inconsistent configurations with metadata	Dependency Check	592	✗	Naming error
	Inconsistent version numbers with release dates	Dependency Check	12,018	✗	Confusing version orders
	<u>Missing required modules for indirect imports</u>	Import Validation	11,023	✗	Incorrect dependencies
	Inconsistent modules in direct imports with installed dependencies	Import Validation	6,678	✗	
Other run-time Errors in imports	Import Validation	8,178	✗		
Incorrect Code	Missing source code	Dependency Check	2,588	✓	Creating placeholders
	Parsing error	Dependency Check	431	✓	Invalid syntax/encoding
	Multiple version control failure	Import Validation	15,507	✗	Incorrect dependencies

sues in

ng  
Technologyhina  
it.edu.cnLyu  
of Hong Kong  
China  
edu.hk

alled Dependencies:

```

==1.13.1
==0.26.2
py==1.21.6
ck==3.12.2
w==9.5.0

```

Monitor as \_GymMonitor

ird-party library

d utility of Python  
lable on the Python  
thousand Python



# Dependencies and dependency declarations

code.py:

```
import numpy as np  
...
```

pyproject.toml

```
[project]  
name = "my_project"  
dependencies = ["numpy"]
```

setup.cfg:

```
[options]  
install_requires =  
    numpy==1.26.2
```

setup.py:

```
from setuptools import setup  
  
setup(  
    name="my_project",  
    install_requires=["numpy"]  
)
```

```
$ python code.py  
Traceback (most recent call last):  
  File ".../code.py", line 1, in <module>  
    import numpy as np  
ModuleNotFoundError: No module named 'numpy'
```

# Dependencies and dependency declarations

code.py:

```
import numpy as np  
...
```

requirements.txt:

```
numpy
```

# Undeclared and unused dependencies

code.py:

```
import numpy as np
import pandas
...
```

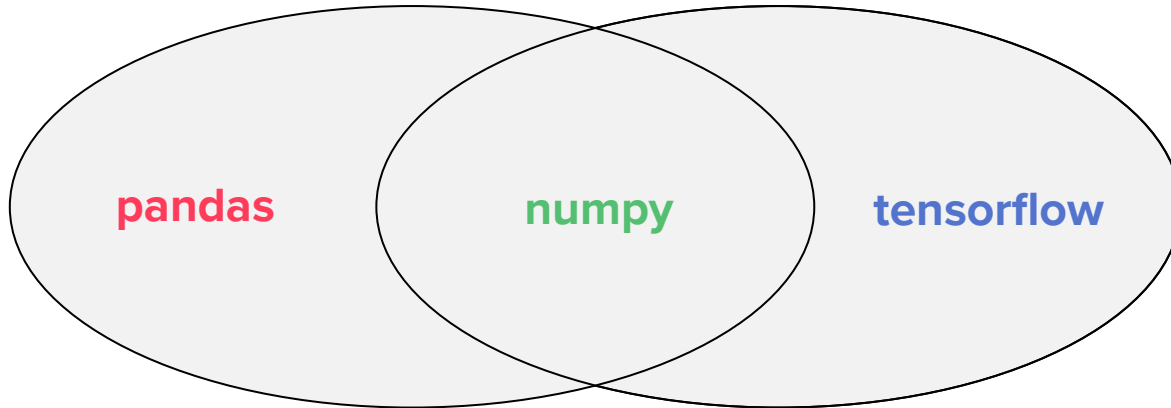
requirements.txt:

```
numpy
tensorflow
```

# Undeclared and unused dependencies

Dependencies  
in your code:

Declared  
dependencies:



# Can we automate this analysis?

# TWEAG

by Modus Create



Nour El Mawass



Maria Knorps



Johan Herland



Zhihan Zhang



Richard  
Bullington-McGuire



# FawltyDeps

code.py:

```
import numpy as np
import pandas
...
```

requirements.txt:

```
numpy
tensorflow
```

```
$ fawltydeps --detailed
$ fawltydeps --detailed
T These imports appear to be undeclared dependencies:
T - 'pandas' imported at:
T   code.py:2
T
T These dependencies appear to be unused (i.e. not imported):
F - 'tensorflow' declared in:
  requirements.txt
```



# FawltyDeps

## Python code and compatibility:

- + Python scripts (.py)
- + Jupyter notebooks (.ipynb)
- + Python v3.8 – v3.13
- + Linux, MacOS, and Windows

## Availability:

- + via PyPI:
  - + `pip install fawltydeps`
  - + As a dev dependency: `uv add fawltydeps`
  - + As a stand-alone tool: `uvx fawltydeps`
- + as a [pre-commit hook](#) or [Github Action](#)

## Dependency declarations:

- + `requirements.txt`
- + `pyproject.toml` (PEP 621, Poetry, Pixi)
- + `setup.py` (simple)
- + `setup.cfg`
- + `environment.yml` (Conda)
- + `pixi.toml` (Pixi)

## Configurable:

- + `fawltydeps --help`
- + Command-line options override configuration
- + `$fawltydeps_<directive>` in environment
- + `[tool.fawltydeps]` in `pyproject.toml`
- + `fawltydeps --generate-toml-config`

# FawltyDeps output

code.py:

```
import numpy as np
import pandas
...
```

requirements.txt:

```
numpy
tensorflow
```

```
$ fawltydeps --list-imports --detailed
code.py:1: numpy
code.py:2: pandas
```

```
$ fawltydeps --list-deps --detailed
requirements.txt: numpy
requirements.txt: tensorflow
```

```
$ fawltydeps --list-sources --detailed
```

Sources of Python code:

code.py (using ./ as base for 1st-party imports)

Sources of declared dependencies:

requirements.txt (parsed as a requirements.txt file)

# FawltyDeps

code.py:

```
import numpy as n
import pandas
...
```

```
$ fawltydeps --li
code.py:1: numpy
code.py:2: pandas
```

```
$ fawltydeps --json
{
  "settings": { ... },
  "sources": [
    {
      "source_type": "CodeSource",
      "path": "code.py",
      "base_dir": "."
    },
    {
      "source_type": "DepsSource",
      "path": "requirements.txt",
      "parser_choice": "requirements.txt"
    }
  ],
  "imports": [
    {
      "name": "numpy",
      "source": {"path": "code.py", "lineno": 1}
    },
    {
      "name": "pandas",
      "source": {"path": "code.py", "lineno": 2}
    }
  ],
  "declared_deps": [
    {
      "name": "numpy",
      "source": {"path": "requirements.txt", "lineno": 1}
    },
    {
      "name": "tensorflow",
      "source": {"path": "requirements.txt", "lineno": 2}
    }
  ]
}
```

```
],
"resolved_deps": {
  "tensorflow": {
    "package_name": "tensorflow",
    "import_names": ["tensorflow"],
    "resolved_with": "IdentityMapping",
    "debug_info": null
  },
  "numpy": {
    "package_name": "numpy",
    "import_names": ["numpy"],
    "resolved_with": "IdentityMapping",
    "debug_info": null
  }
},
"undeclared_deps": [
  {
    "name": "pandas",
    "references": [
      {"path": "code.py", "lineno": 2}
    ]
  }
],
"unused_deps": [
  {
    "name": "tensorflow",
    "references": [
      {"path": "requirements.txt"}
    ]
  }
],
"version": "0.18.0"
}
```

# FawltyDeps as pre-commit hook

<https://pre-commit.com>



```
$ pip install pre-commit
[...]
$ vim .pre-commit-config.yaml
$ pre-commit install
pre-commit installed at .git/hooks/pre-commit
```

.pre-commit-config.yaml:

```
repos:
- repo: https://github.com/tweag/FawltyDeps
  rev: v0.18.0
  hooks:
  - id: check-undeclared
  - id: check-unused
```

# FawltyDeps as pre-commit hook

<https://pre-commit.com>



```
$ vim code.py
$ git commit -am "Add pandas"
[...]
FawltyDeps-undeclared.....Failed
- hook id: check-undeclared
- exit code: 3
```

These imports appear to be undeclared dependencies:

```
- 'pandas' imported at:
  code.py:2
```

```
FawltyDeps-unused..(no files to check)Skipped
```

code.py:

```
import numpy as np
+import pandas
...

```

# FawltyDeps as pre-commit hook

<https://pre-commit.com>



```
$ vim code.py requirements.txt
$ git commit -am "Add pandas"
FawltyDeps-undeclared.....Passed
FawltyDeps-unused.....Passed
[main 6ece52e] Add pandas
 2 files changed, 2 insertions(+)
$ vim requirements.txt
$ git commit -am "Declare tensorflow dependency"
FawltyDeps-undeclared..(no files to check)Skipped
FawltyDeps-unused.....Failed
- hook id: check-unused
- exit code: 4
```

These dependencies appear to be unused (...):

- '**tensorflow**' declared in:  
requirements.txt

code.py:

```
import numpy as np
import pandas

...
```

requirements.txt:

```
numpy
pandas
+tensorflow
```



# FawltyDeps in CI (GitHub Actions)

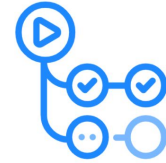
[github.com/tweag/FawltyDeps-action](https://github.com/tweag/FawltyDeps-action)

.github/workflows/fawltydeps.yaml:

```
name: FawltyDeps

on: [push, pull_request]

jobs:
  fawltydeps:
    runs-on: ubuntu-latest
    steps:
      - name: checkout
        uses: actions/checkout@v4
      - name: lint - FawltyDeps
        uses: tweag/FawltyDeps-action@v0.2.0
        with:
          options: --detailed
```



GitHub Actions

# FawltyDeps in CI (GitHub Action)

The screenshot shows a GitHub Actions workflow run titled "Add undeclared use of pandas #1" that has failed. The failure message is "All checks have failed" with 2 failing checks. One of the failing checks is "FawltyDeps". A dropdown menu is open for the "FawltyDeps" check, showing options: Summary (highlighted with a red arrow), Jobs, fawltydeps (with a red 'x'), Run details, Usage, and Workflow file. The "Summary" view is expanded, showing the error message: "Error: 'fawltydeps --detailed' found issues: These imports appear to be undeclared dependencies: - 'pandas' imported at: code.py:2". The code diff shows the addition of "import pandas" on line 2.

Add undeclared use of pandas #1

Open

All checks have failed  
2 failing checks

FawltyDeps

FawltyDeps

This branch  
Merging can be

Merge pull request

Summary

Jobs

fawltydeps

Run details

Usage

Workflow file

Annotations

Details

fawltydeps summary

Error: 'fawltydeps --detailed' found issues:

These imports appear to be undeclared dependencies:  
- 'pandas' imported at:  
code.py:2

```
1 1 import numpy as np
2 2 + import pandas
2 3
3 4 ...
```

---

---

# How to draw an Owl.

*“A fun and creative guide for beginners”*

---

---

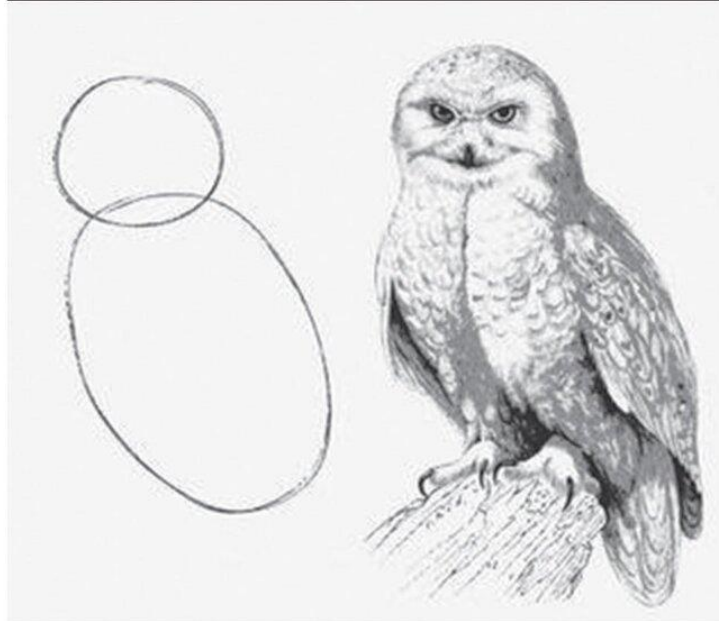


Fig 1. Draw two circles

Fig 2. Draw rest of the damn Owl

# Matching imports names to package names

code.py:

```
import numpy as np
import sklearn
import setuptools
import pkg_resources
...
```

requirements.txt:

```
numpy
scikit-learn
setuptools
```

Imports are found

Dependencies are installed

Your project's  
Python environment

# Resolving dependencies into import names

- + Python environments inside the project
  - + Search the entire project?
  - + Search within the paths given to FawltyDeps
  - + Python environments given with `--pyenv`
- + Python environment where FawltyDeps is installed
- + Custom mapping: `[tool.fawltydeps.custom_mapping]`
  - + Takes precedence, if given
- + Auto-install from PyPI? `--install-deps`
  - + Expensive, disabled by default
- + Final resort: Identity mapping

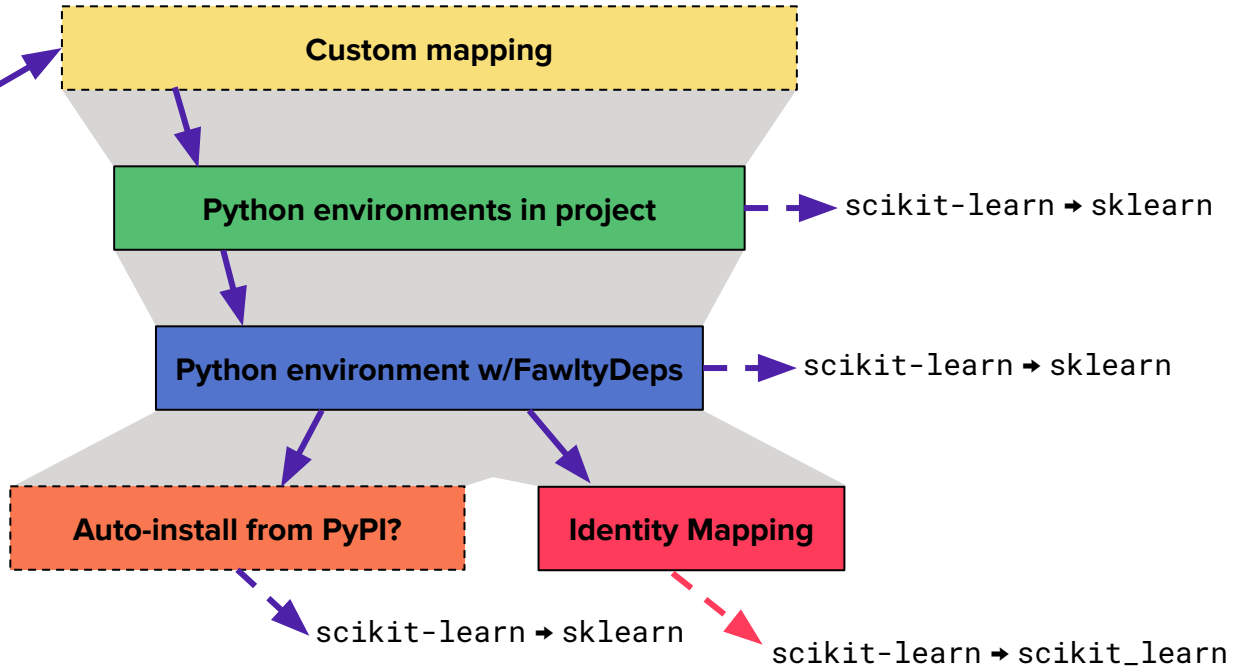
# Resolving dependencies into import names

code.py:

```
import sklearn
```

requirements.txt:

```
scikit-learn
```





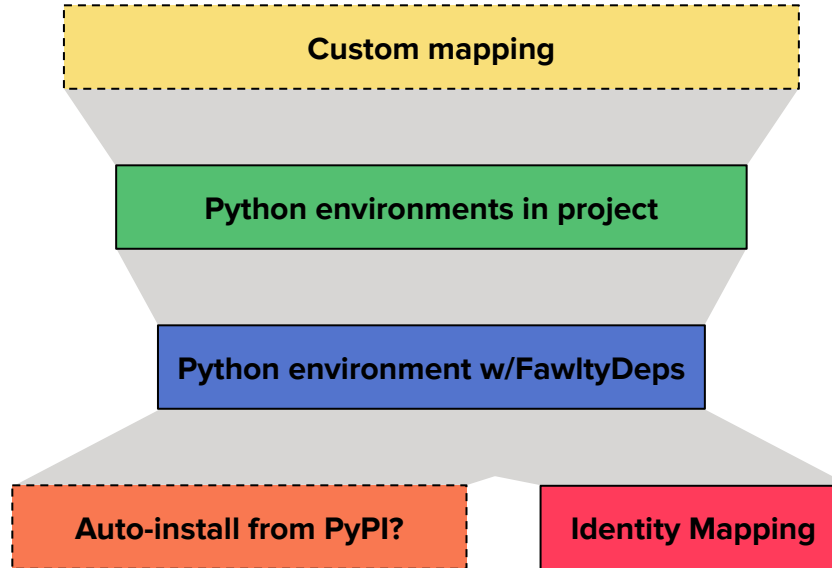
# Resolving dependencies into import names

code.py:

```
import sklearn
```

requirements.txt:

```
scikit-learn
```



# More complexities and advanced use cases

- + Exclude parts of your project from analysis:
  - + `--exclude tests/`
  - + `--exclude-from .gitignore`
- + Some dependencies are never meant to be imported:
  - + Tools: tox, black, flake8, ruff, pylint, mypy, etc.
  - + Type stubs
  - + Indirect imports or plugin-style package collections
  - + `--ignore-unused`
- + Conditional imports, alternative imports, e.g.: `try: import foo; except: import bar`
  - + `--ignore-undeclared`



# Work in progress

- + Dynamic imports, e.g.: `importlib.import_module(<some dynamic expression>)`
- + Required dependencies vs. various degrees of optional dependencies
  - + `install-requires` vs. `extras-require` in `setup.py`
  - + `[project.dependencies]` vs `[project.optional-dependencies]` in `pyproject.toml`
  - + PEP 735: `[dependency-groups]` in `pyproject.toml`
  - + Poetry: `[tool.poetry.group.$NAME.dependencies]` in `pyproject.toml`
  - + Development dependencies: `dev-requirements.txt` vs `requirements.txt`



# Summary

- + Reproducibility matters
- + Missing or incorrect dependency declarations is the greatest culprit
- + **FawlyDeps** finds **undeclared** and **unused** dependencies in your project
  - + Works with Python scripts and Jupyter notebooks, across Python version and platforms
  - + Understands the many ways of declaring dependencies
  - + Works out of the box on most projects, and can be configured to suit *all* projects
  - + Can be an automated part of your workflow, either pre-commit or in your CI

# QUESTIONS?



Slides

- + [FawltyDeps on GitHub](#)
- + [FawltyDeps on PyPI](#)
- + [FawltyDeps on Discord](#)
  
- + Tweag blog: [tweag.io/blog](https://tweag.io/blog)
- + Tweag: [tweag.io](https://tweag.io)
- + Modus Create: [moduscreate.com](https://moduscreate.com)



@jherland

# THANK YOU!



MODUS CREATE



by Modus Create

