

The Tools of Data-Driven Engineering

Handout Edition
Niels Skovgaard Jensen

This is the handout version of the lecture slides.

- Most references to tools have links to their respective websites.
- There are links to all slides in the Contents slide
- I reference a lot of different tools so you can investigate them yourself. My recommended tools are usually marked in **bold**.

My hope for this lecture:

- Show you where to look for good SciML development practices
- Get you thinking about what fits your use-cases
- Give you a heads-up of what is expected in industry

What you **should** do after this lecture:

- Investigate tools and reconsider your workflows
- Understand that a little work up front can make a lot of things easier

What you **should not** do:

- Panic that you don't use *all* these tools
- Spend all your time trying to learn *all* these tools

Tools are supposed to make your life easier

(But it will take a little work and interest)

Contents

The Tools of Data-Driven

Engineering	1
Handout README	2
My hope for this lecture:	3
.....	4
Analysis as a DAG	5
Each Step Should Only Depend on the Previous	6
How Notebooks Often Look	7
How to ensure Acyclicity	8
Environment Reproducibility	9
Reproducibility in CompSci	10
Python Environments	11
Basic uv commands	12
Reuse code when possible	13
Share code in libraries, use the code in scripts	14
Coding Practices	15
Linting and Formatting	16
Formatting Example	17

Linting Example	18
Python type-hinting	19
Static Type Checkers	20
Static Type Checker Example ..	21
How it can look in your IDE ...	22
Dataclasses	23
Organising Configuration in Dataclasses	24
Using classmethods as a flexible __init__	25
Dataclasses in Python – Structs in Julia	26
Making Your Scripts Reconfigurable	27
Non-Reconfigurable	28
Reconfigurable	29
Version Control	30
Git – Really quick	31
Git is your friend	32
GUIs for GIT	33

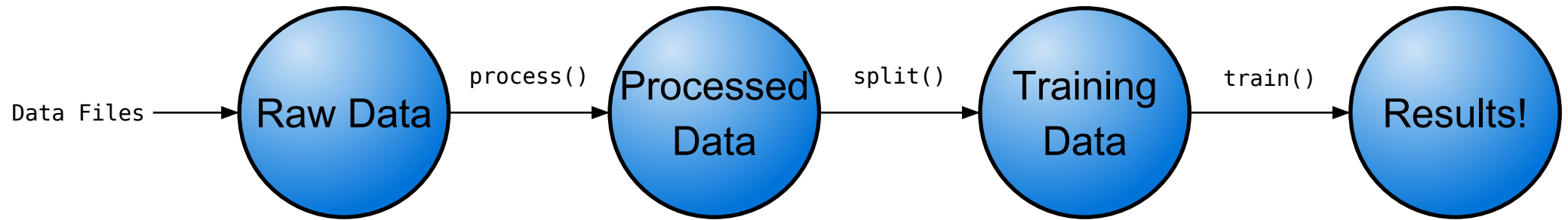
Git for Data-Driven

Engineering	34
.gitignore – What not to Commit	35
Pre-Commit	36
Exercises	37
Exercises	38
2.1 – Environment	39
2.2 – Formatting and Linting ...	40
2.3 – Typing	41
2.4 – Git	42
2.5 – Reproducibility	44
2.6 – Investigate Your IDE	45
Advanced Tools	46

Analysis as a DAG

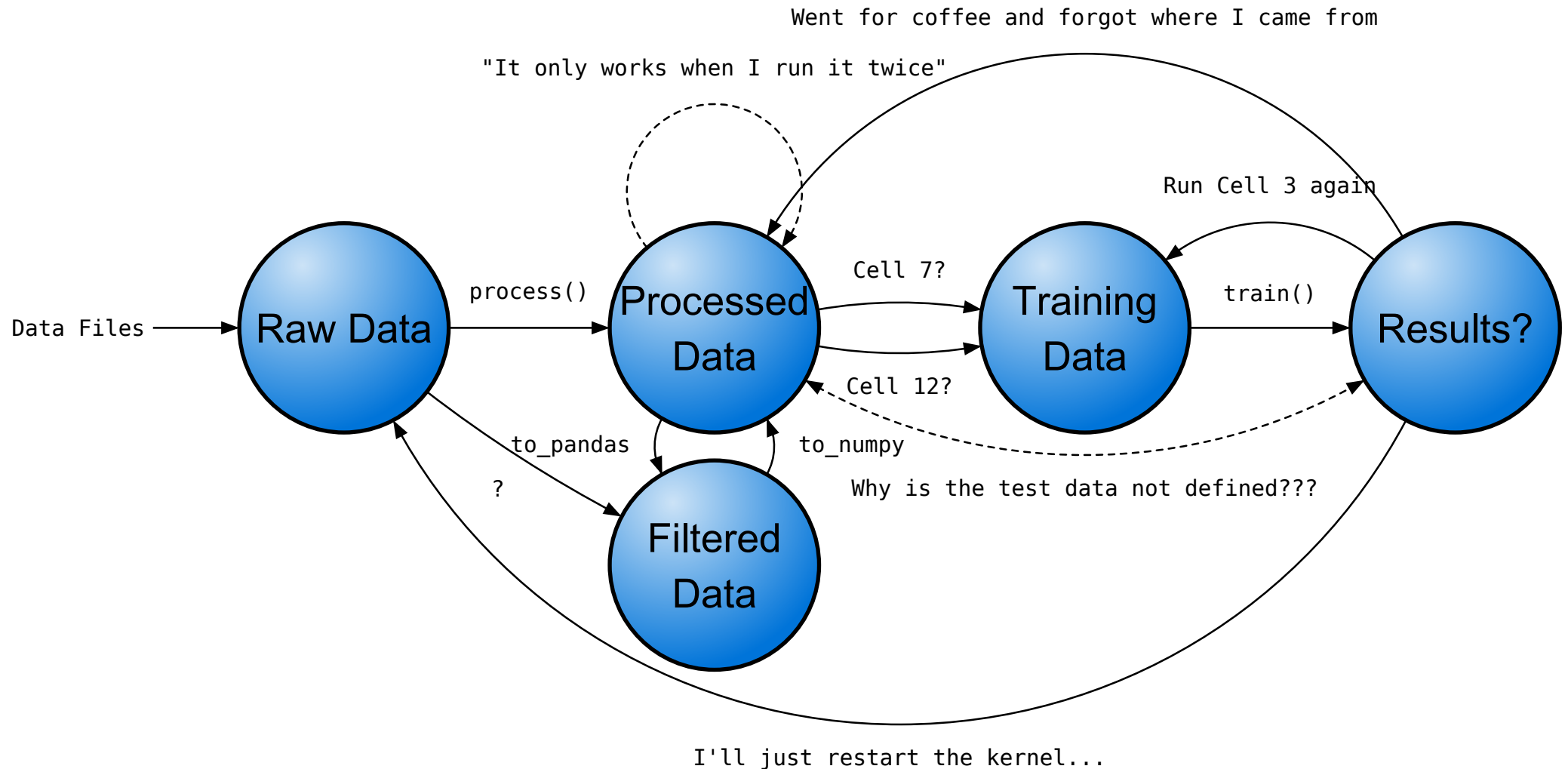
Each Step Should Only Depend on the Previous

Analysis as a DAG



How Notebooks Often Look

Analysis as a DAG



How to ensure Acyclicity

- Prefer scripts or library code for things you want reproduced
 - Or use reproducible notebooks such as [Marimo](#)
- Consider raw data *immutable*
- Design a data pipeline into your project

```
1 |— data
2 |   |— interim  <- Intermediate data that has been transformed.
3 |   |— processed <- The final, canonical data sets for modeling.
4 |   └— raw      <- The original, immutable data dump
```

- TIP: You can use command-run files, such as [make](#) or [just](#), to string together processing and experimentation scripts.

Environment Reproducibility

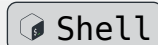
- We are in one of the few sciences with almost entirely controlled experimentation
 - Reproducibility is (relatively) easy
- SciML is a very experiment-driven field
 - Reproducibility is harder
 - But much more important

There are many ways to use python and install packages:

- pip, venv, [conda](#), [uv](#)
 - uv is quickly becoming the modern industry standard
- Reproducibility when using uv
 - pyproject.toml – Defines requirements
 - uv.lock – locks requirements
 - uv creates isolated environments per project
- uv makes it easy to:
 - Share and reproduce environments across computers
 - Ensure you are using the right environment when running code
- It is also *really* fast

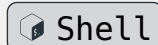
Create a new project

```
1 uv init --lib # Create a library
2 uv init --script # Standalone script
```



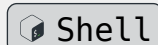
Add a dependency

```
1 uv add numpy #numpy, torch, jax, etc...
```



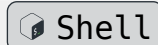
Run a file

```
1 uv run myscript.py # Uses the venv uv has created
```



Reproduce the project


```
1 uv sync # Reproduce the environment on a new computer
```



Reuse code when possible

- uv makes it easy to create a shared library of code

```
1 uv init --lib my_library
```


 Shell

- This creates a folder structure like:

```
1 my_library/  
2 |— pyproject.toml      # Project configuration file  
3 |— src  
4   |— my_library        # Library source code  
5   |— __init__.py  
6   |— hello.py          # Example module
```

Which can then be imported into a script

```
my_script.py
```

 Python

```
1 from my_library.hello import hello
```

Share code in libraries, use the code in scripts

```
1  awesome_sciml_project/
2  └─ pyproject.toml      # Project configuration file
3  └─ src
4      └─ my_library      # Library source code
5          └─ __init__.py
6              └─ data.py  # Example module
7                  └─ train.py
8                      └─ models.py
9  └─ data/
10      └─ processed
11      └─ interim
12      └─ raw
13  └─ scripts/ # Create seperate scripts with some agreed-on naming scheme
14      └─ 1.0-nsj-first_visualisation.py
15      └─ 2.0-apek-fit_data.py
```

For more inspiration: [Cookiecutter Datascience Template](#)

Coding Practices

Linting and Formatting

linting – Automated coding style guides

formatting – Automated code aesthetics

- Modern Python has converged on a style
- It is much easier to read consistently styled-code
- Formatting tools make it easy to follow the style!

The industry standard formatter/linter for python is [Ruff](#)

Formatting Example

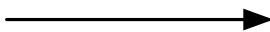
ugly_code.py

Python

```
1 import numpy as np
2 import torch
3 def foo(
4     x,y ,
5     z, u, items=[]
6 ):
7     q = 2; t = np.array([1, 2, 3]); z
8     = 4
9     assert (x > 0, "must be
10    positive")
11    if x == None or y == True:
12        return x - y
13    if u > 0:
14        return y+z
15    else:
16        return x+ y
17 if __name__ == "__main__":
18     result = foo(1, 2 ,3,
19     4)
20     print("result:", result )
```

```
1 uvx ruff format
   ugly_code.py
```

Shell



formatted.py

Python

```
1 import numpy as np
2
3
4 def foo(x, y, z, u, items=[]):
5     q = 2
6     t = np.array([1, 2, 3])
7     z = 4
8     assert (x > 0, "must be
9     positive")
10    if x == None or y == True:
11        return x - y
12    if u > 0:
13        return y + z
14    else:
15        return x + y
16
17 if __name__ == "__main__":
18     result = foo(1, 2, 3, 4)
19     print("result:", result)
```

Linting Example

```

formatted.py Python
1  import numpy as np
2
3
4  def foo(x, y, z, u, items=[]):
5      q = 2
6      t = np.array([1, 2, 3])
7      z = 4
8      assert (x > 0, "must be
9          positive")
10     if x == None or y == True:
11         return x - y
12     if u > 0:
13         return y + z
14     else:
15         return x + y
16
17 if __name__ == "__main__":
18     result = foo(1, 2, 3, 4)
19     print("result:", result)

```

```

1
2 uvx ruff check formatted.py

```



```

1  LOTS OTHER ERRORS
2  ...
3
4  E712 Avoid equality comparisons to
   `True`; use `y:` for truth checks
5  --> unformatted.py:8:21
6  |
7  6 |         q = 2; t = np.array([1, 2, 3]);
   z = 4
8  7 |         assert (x > 0, "must be
   positive")
9  8 |         if x == None or y == True:
10 |                                     ~~~~~~
11  9 |             return x - y
12 10 |         if u > 0:
13 |
14 help: Replace with `y`
15
16 Found 7 errors.
   No fixes available (4 hidden fixes can be
17 enabled with the `--unsafe-fixes`
   option).
18

```

```
1 x: float = 1.0
2 y: np.ndarray = np.array([2.0])
3 t: Tensor = torch.tensor([2.0])
4
5 def foo(x: np.ndarray) -> float:
6     return float(x.item())
7 my_func: Callable = foo
```

Python is a *dynamically* typed language.

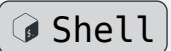
In Python, type hints are structured comments (they do not affect runtime). But they can be analyzed by external software

This allows for *static type checking* – A powerful tool to catch bugs before you run your code

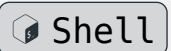
Static Type Checkers

- Examples: [mypy](#), [pyright](#), [basedpyright](#), ([ty](#), [pyrefly](#) – In Beta)
- Can easily be run in your IDE
- Use uv tool (uvx) to experiment with them:

```
1 uvx mypy check <my_code>.py
```



```
1 uvx basedpyright check <my_code>.py
```



```
1 uvx ty check <my_code>.py
```



```
static_broken.py Python
1 def add(a: int, b: int) -> int:
2     return a + b
3
4 def greet(name: str) -> str:
5     return "Hello, " + name
6
7 result = add("hello", 5)
8 message = greet(42)
```

```
1
2 uvx mypy static_broken.py Shell
```




```
typing_example/
static_broken.py:7: error:
1 Argument 1 to "add" has incompatible type "str";
  expected "int" [arg-type]
  Text
typing_example/static_broken.py:8:
2 error: Argument 1 to "greet" has incompatible type "int"; expected
  "str" [arg-type]
3 Found 2 errors in 1 file (checked 1
  source file)
```

Caught **before** running the code!


How it can look in your IDE

static_broken.py


 Python

```
1 def add(a: int, b: int) -> int:
2     return a + b
3
4 def greet(name: str) -> str:
5     return "Hello, " + name
6
7 result = add("hello" expected int, 5)
8 message = greet(42 expected str)
```

```
1  from dataclasses import dataclass
2
3  # Init example from python documentation
4  class Complex:
5      def __init__(self, realpart, imagpart):
6          self.r = realpart
7          self.i = imagpart
8  z = Complex(1.0,1.0)
9
10 # Dataclass equivalent
11 @dataclass
12 class Complex:
13     r: float
14     i: float
15 z = Complex(1.0,1.0) # Create in the same way
16
17 @dataclass
18 class Dataset: # They are excellent for packing data
19     mat: jnp.array
20     coordinate: list[float]
21     desc: str
22     ... # Whatever your heart desires
```


 Python

train.py

 Python


```
1 @dataclass
2 class TrainingConfig:
3     model: torch.nn.Module
4     train_data: Data
5     test_data: Data
```

train.py

 Python

```
1 @dataclass
2 class TrainingResult:
3     model_final: torch.nn.Module
4     train_loss: list[float]
5     test_loss: list[float]
```

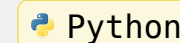
train.py

 Python

```
1 config = TrainingConfig(model = ..., ... )
2
3 def train(config: TrainingConfig) -> TrainingResult:
4     ...
5     return TrainingResult(
6         model_final=model,
7         train_loss=train_loss,
8         ...
9     )
10 result = train(config)
11 result.model_final # Access trained model or other results easily
12
13 result2 = train(config2) # Easy to keep track of multiple runs!
```


Using classmethods as a flexible `__init__`

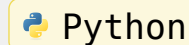
```
1  @dataclass
2  class Data:
3      mat: torch.Tensor
4      R: torch.Tensor
5      # name, version, metadata... Whatever
6
7      @classmethod
8      def from_file(cls, path: Path) -> 'Data':
9          # load data from file
10         mat, R = load_file(path)
11         return cls(mat = mat, R = R)
12
13     @classmethod
14     def from_numpy(cls, array: np.ndarray) -> 'Data':
15         # convert numpy array to Data instance
16         mat = torch.from_numpy(array)
17         R = torch.eye(mat.size(1)) # Example: identity matrix as R
18         return cls(mat = mat, R = R)
```



Python


- Julia has built-in type system
- Types are enforced at runtime
- Static typing is not as mature as Python
- Dataclasses are (roughly) equivalent to Julia structs
- Multiple-dispatch with structs is one of Julia's superpowers

- Bad Practice: Changing code to alter variables
 - Each experiment changes code state
- Good Practice: Making your scripts reconfigurable
 - Save the command, and reproduce your result in the future
- Tip: LLM Coding tools are excellent at making your scripts reconfigurable.




```
1
2 @dataclass
3 class TrainConfig:
4     lr: float = 1e-3
5     epochs: int = 1000
6     batch_size: int = 32
7
8 def train(config: TrainConfig):
9     ...
10
11 if __name__ == "__main__":
12     config = TrainConfig()
13     train(config)
```

```
1 import argparse
2 parser = argparse.ArgumentParser()
3 parser.add_argument("--lr", type=float, default=0.001)
4 parser.add_argument("--epochs", type=int, default=10)
5 parser.add_argument("--batch-size", type=int, default=32)
6
7 ... # Same as previous example
8 if __name__ == "__main__":
9     args = parser.parse_args()
10    config = TrainConfig(lr=args.lr, epochs=args.epochs, batch_size=args.batch_size)
11    train(config)
```

 Python

Now you can change the script when running it

```
1 uv run <myscript.py> --lr 0.002 --epochs 100 --batch-size 64
```

 Shell

Advanced Tool: [Hydra](#)

Version Control

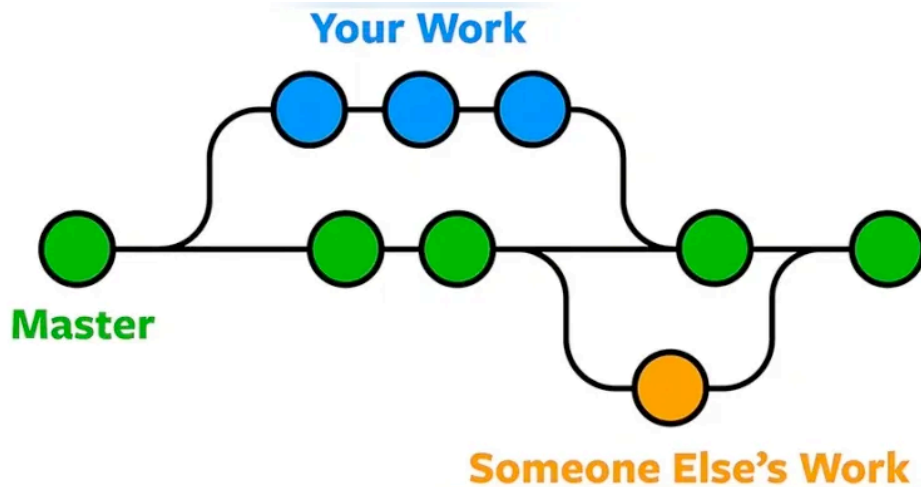
Git – Really quick

Git is a distributed version-control system.

- Enables collaboration
- And *restoration* of code state

Each change is made in a *commit*

Multiple chained *commits* form *branches*



Git is your friend

External git services:

- [GitHub](#)
- [GitLab](#)
- *Makes collaboration easy*

Remember: Git \neq github

Local git repos are also very useful

Advanced Tool: [Jujutsu](#)

GUIs for GIT

- [GitGraph](#) in VSCode – Free GUI for using git
- [GitHub Desktop](#) – Free GUI for using GitHub
- [Smartgit](#)

But most power comes from the command line

Some pragmatic guidelines:

- Use commits to restore code state
- When something works – commit
- When you have a result – commit
- Before you experiment – commit
- Forgot to commit before experimenting? Use `git stash`
 - Saves your uncommitted changes and restores a clean state
 - `git stash pop` to bring them back

.gitignore – What not to Commit

- Git tracks *everything* by default – you must tell it what to ignore
- Create a .gitignore file in your project root:

```
1 __pycache__/      # Python cache files
2 .venv/            # Virtual environment
3 data/             # Large datasets
4 *.pt              # PyTorch model checkpoints
5 *.ckpt            # Other checkpoints
6 .env              # API keys and secrets
7 wandb/            # Experiment logs
```

[Pre-commit](#) is a tool to make sure everyone follows the same linter/formatter guidelines.

Example: Before each commit run:

- Linting
- Formatting
- uv sync (To make uv.lock)

Pre-commit ensures a commit cannot be made if the tools fail

It is especially useful for collaboration

Advanced: CI/CD pipelines.z

Exercises

Exercises

Like the lecture, these exercises are made to give an introduction to the tools that have been discussed. Thus you are free to choose and experiment for yourself. The goal is to create a forum for discussing how we work, use tools and do the best possible SciML.

Comprehension questions: Discuss these questions in groups

- What is meant by treating your analysis as a DAG? Why is this desirable?
- Why should raw data be treated as immutable?
- What is the difference between a lockfile and a requirements.txt?
- What is the difference between *linting* and *formatting*?
- Which style of typing does Python have, and what is the advantage, what is the disadvantage?
- Which types of errors can static type checking fix?
- What is a dataclass, and why is it useful for organizing experiment configurations and data?
- Why is it better to make scripts reconfigurable (e.g. with argparse) rather than editing variables in the code?
- What is git, and how does it aid in reproducibility?
- Why is it useful to commit before experimenting with your code?

The following exercises all assume that uv is installed on your computer. All other used tools are installed as uv tools through the uvx command.

If you use another programming language, spend your effort on understanding how you can do these things in your language to replicate the effects and coding standards.


If you are a tool wizard and all this seems trivial, there is a list of more advanced tools at the bottom you can investigate.

Installing uv Guide:

<https://docs.astral.sh/uv/getting-started/installation/>

Note that you can use conda to install uv if you do not want to use the standalone installer.

```
1 conda install conda-forge::uv
```


 Shell

2.1 – Environment

uv documentation: <https://docs.astral.sh/uv/>


1. Create a new uv project with a src library

```
1 uv init --lib myproject
```

 Shell

2. Use `cd myproject` to go into the folder and read the `pyproject.toml`
3. Add a dependency of your choice (e.g. `numpy`)

```
1 uv add numpy
```

 Shell


4. Read the `pyproject.toml` again and notice the changes. Also investigate the `uv.lock` file.
5. Write a small script that imports your added dependency and run it with `uv run`

2.2 – Formatting and Linting

Ruff documentation: <https://docs.astral.sh/ruff/>


1. Find a short piece of code you have written previously
2. Run the ruff formatter on it and observe the changes

```
1 uvx ruff format <my_code>.py
```

 Shell

3. Run the ruff linter and try to understand the warnings it gives (if any)

```
1 uvx ruff check <my_code>.py
```

 Shell


4. Try fixing the linting errors, either manually or with the `--fix` flag

2.3 – Typing

Python typing documentation: <https://docs.python.org/3/library/typing.html>


1. Find a piece of code you have previously written that does not have type hints.
2. Type-hint the code as best you can
3. Run mypy on the type-hinted code to see what errors come up:

```
1 uvx mypy <your_code>.py
```

 Shell

4. Try to understand the errors (if any) and think about what types of bugs they might introduce, then fix all the errors.
5. (Extra) Now run the stricter basedpyright type checker and investigate and fix the errors it points to.


```
1 uvx basedpyright <your_code>.py
```

 Shell

2.4 – Git


1. Initialize a git repository in your project folder

```
1 git init
```

 Shell


2. Add your files and make an initial commit

```
1 git add .  
2 git commit -m "Initial commit"
```

 Shell


3. Make a small change to your code (e.g. add a print statement)
4. Use `git diff` to see what changed, then stage and commit the change

```
1 git add <my_code>.py  
2 git commit -m "Describe the change"
```

 Shell

5. Use `git log` to view your commit history
6. Create a repository on [GitHub](https://github.com) and push your project to it

```
1 # Create a repo on github.com, then:  
2 git remote add origin https://github.com/<your-username>/<your-repo>.git  
3 git push -u origin main
```

 Shell

7. (Extra) Create a new branch, make a change on it, and merge it back into main

```
1 git checkout -b experiment  
2 # make changes and commit
```

 Shell

```
3 git checkout main
```

```
4 git merge experiment
```

2.5 – Reproducibility

1. Download the code that was made during the lecture demonstration on GitHub using the GitHub CLI, GitHub Desktop or git directly.
2. Run `uv sync` to reproduce my environment
3. Run my example script with `uv run`

This is how easy reproducibility should be!

2.6 – Investigate Your IDE

Most of the tools you have just used have plugins for VS Code and many other IDEs, investigate these plugins and how to use them.

Advanced Tools

If all of these exercises seem trivial and you are already a tool wizard. Then here are some advanced tools you can look into that are great for more complicated projects and users.

Version Control: Jujutsu A git-compatible version control system with superpowers. (for an amazing Jujutsu TUI: jjui)

Argparsing: Hydra Create configurable and reproducible experiments directly from your python dataclasses.

einops – Self documenting tensor operations in Jax, Torch and Numpy For a deep-dive on einstein-like notation: <https://grrddm.notion.site/Einsums-in-the-wild-bd773f01ba4c463ca9e4c1b5a6d90f5f?v=4b86eb72820943c083c84c2674f67ac4>

Jaxtyping: Runtime type checking for tensor shapes in Jax, Torch, Numpy and more....

Typst: An amazing LaTeX and Overleaf alternative.

Oxidise Your Terminal: A list of really cool Rust-based terminal tools <https://www.youtube.com/watch?v=rWMQ-g2QDsI&themeRefresh=1> Honorable mentions: fish, zoxide, fzf, eza,

Writing Python like it is Rust: <https://kobzol.github.io/rust/python/2023/05/20/writing-python-like-its-rust.html> A cool perspective on how to write maintainable and well-typed Python code.