

# STAT243 Lecture 4 Good Practices

## 1 Good Coding Practices

### 1.1 Editors

- 选择合适的编辑器/IDE

使用支持目标语言的现代编辑器 (VS Code、Emacs、Sublime、vim、RStudio 等)。

优点包括：

- 语法高亮与代码补全；
- 自动缩进与括号匹配；
- 行号显示 (便于定位错误)；
- 集成运行/调试功能；
- 集成 AI 辅助编程 (如 GitHub Copilot、Gemini Code Assist 等)。

#### Logic

在 VS Code 或 RStudio 中使用 Quarto，可以直接运行、调试并格式化代码，适合课程作业与科研分析。

### 1.2 Code Syntax and Style

#### 1.2.1 General Rules

- 参考标准**：遵循 Python 的官方风格指南 [PEP 8](#)。
- 文件头部信息**：用注释说明作者、日期、用途、Python 版本、主要依赖库。
- Docstrings**：为公共模块、类、函数编写文档字符串；非公共项使用行内注释。
- 缩进**：每层 4 个空格 (避免使用制表符)。
- 空白与可读性**：
  - 操作符两侧加空格；
  - 函数参数与元素之间加空格；
  - 使用空行分隔逻辑块；
  - 使用括号分组长表达式以便换行与注释。

#### Remark

避免超过 79 字符的代码行 和 72 字符的注释行。

可用 `ruff format` 或 `black` 自动格式化。

- 命名约定 (PEP8)**：
  - 类名 → `UpperCamelCase`
  - 函数、变量 → `snake_case`
  - 非公共方法 → `_leading_underscore`
  - 函数名使用动词：如 `calc_mean()`、`compute_error()`。
- 注释原则**：
  - 说明复杂逻辑或魔法常数；
  - 简洁明确，完整句式；
  - 例如：



Python

```
1 newdf = (
2     pd.read_csv('file.csv')
3     .rename(columns={'STATE': 'us_state'}) ## adjust column names
4     .dropna() ## remove NA rows
5 )
```

### Logic ↴

编写可读代码的核心：**清晰的结构 + 自文档化命名 + 简明注释。**

## | 1.2.2 Coding Style Suggestions

- **模块化设计：**
  - 拆解为小函数，每个函数执行单一任务。
  - 避免重复代码，核心功能应封装复用。
  - 每个函数需：
    - 明确职责；
    - 有 docstring；
    - 具备单元测试。
- **可重用与透明：**
  - 函数应通过输入参数获取数据，而非操作全局变量。
  - 函数输入输出明确定义，方便复现分析。
- **可维护性与防御性编程：**
  - 检查输入有效性，提供默认参数；
  - 对潜在错误发出警告或抛出异常；
  - 使用断言 (assert) 检测假设。
- **其他建议：**
  - 不硬编码常数（建议使用如 speed\_of\_light = 3e8）。
  - 用 list/tuple 打包相关数据；
  - 避免平台相关代码；
  - 让他人审阅 (code review)。

### Logic ↴

小函数、显式参数、单一职责 → 易调试、易测试、易复用。

## | 1.2.3 Linting

- **Linting 定义：**自动检测代码风格与潜在错误的工具。
- **推荐工具：** ruff (轻量快速) 或 black。



Shell

```
1 pip install ruff
2 ruff check test.py ## 检查语法
3 ruff format test.py ## 自动格式化
```

ruff 自动修复空格、缩进与注释风格。

### ⚠ Remark ▾

Linting 不仅改善可读性，也能防止潜在逻辑错误。

## 1.3 Assertions, Exceptions, and Testing

### 1.3.1 Exceptions

- 异常 (Exception): 运行时错误的处理机制。
- 可通过 `raise` 主动抛出异常：



Python

```
1 def myfun(val):
2     if val <= 0:
3         raise ValueError(`val` should be positive)
```

- **try-except 结构**: 捕获并处理异常。



Python

```
1 try:
2     with open("file.txt", "r") as f:
3         text = f.read()
4 except Exception as err:
5     print(f"{err}\nCheck path: {os.getcwd()}")
6     return None
```

- **re-raise 异常**:



Python

```
1 try:
2     requests.get(url)
3 except Exception as err:
4     print("Problem accessing URL.")
5     raise
```

### ⌚ Logic ▾

异常处理逻辑应：

- ① 捕获可预见错误；
- ② 提供明确提示；
- ③ 在必要时重新抛出异常。

### 1.3.2 Assertions

- **断言 (assert)** 用于验证程序状态是否符合预期：



Python

```
1 number = -42
2 assert number > 0, f"Expected positive, got {number}"
```

若条件不满足，抛出 `AssertionError`。

- **常用断言形式**：



Python

```
1 assert x in y
2 assert isinstance(x, float)
3 assert all(arr)
```

### ⌚ Logic ▾

断言主要用于开发阶段调试和“sanity checks”；  
在生产环境可禁用以提高性能。

## 1.3.3 Testing

- **单元测试 (Unit Tests)**: 验证单个函数的行为。
- 推荐使用 `pytest` :



Python

```
1 import pytest
2 import dummy
3
4 def test_numeric():
5     assert dummy.add_one(3) == 4
6
7 def test_bad_input():
8     with pytest.raises(TypeError):
9         dummy.add_one('hello')
```

运行测试：



Shell

```
1 pytest test_dummy.py
```

### ⌚ Logic ▾

测试不仅验证正确输出，也应验证错误处理行为。

## 1.3.4 Automated Testing (CI)

- **持续集成 (Continuous Integration, CI)**：  
自动在代码变更时运行测试。
- **GitHub Actions 示例：**

`.github/workflows/test.yml`



YAML

```
1 on:
2   push:
3     branches: [main]
4
5 jobs:
6   CI:
7     runs-on: ubuntu-latest
8     steps:
9       - uses: actions/checkout@v4
```

```
10      - uses: actions/setup-python@v5
11      with:
12          python-version: 3.12
13      - run: |
14          pip install pytest
15          pip install --user .
16          pytest
```

#### ⚠ Remark ▾

CI 确保测试在远程环境中自动运行，  
防止“本地能跑、服务器报错”的问题。

## 1.4 Version Control

- 使用版本控制（Git/GitHub）记录修改历史。
- 建议：
  - 使用 GitHub Issues 管理任务；
  - 维护清晰的 commit 信息；
  - 保留运行日志或开发笔记。

## 1.5 AI-Assisted Coding Tools

- **集成式 AI 编程助手**（如 Copilot、Gemini、Cursor）：
  - 代码自动补全与建议；
  - Chat/Agent 模式可自动修改代码；
  - 可基于上下文（文件/目录）生成代码。

#### ⚠ Remark ▾

- 对 AI 生成的代码保持警惕：需理解其逻辑。
- 可用于简单语法或可验证任务（如绘图格式化）。
- 对复杂算法仍需手动验证。

#### ⌚ Logic ▾

适度使用 AI 辅助工具以加速开发，但**理解优先于生成**。

## 2 Debugging and Recommendations for Avoiding Bugs

#### ⌚ Logic: Overview ▾

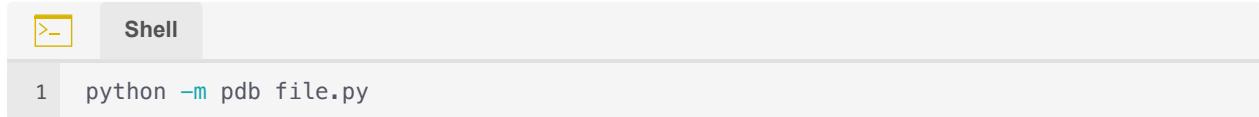
- **Common Debuggers in Python**
  - `pdb`：Python 标准调试器（内置模块）。
  - `ipdb`：IPython 封装的增强版本，支持自动补全与高亮。
  - **JupyterLab**：自带交互式调试器。
  - **VS Code**：最常用的图形化调试器，集成度高，推荐使用。

## 2.1 基本调试策略

- 阅读 traceback
  - 从错误信息底部开始阅读；底部显示触发错误的实际行。
  - 用双引号包裹错误信息进行网页搜索（尤其在 Stack Overflow 上）常能快速找到解决方案。
- 理解调用栈（call stack）
  - 错误通常发生在多层函数调用中。
  - 调试重点是：出错时被执行的函数及其参数。
  - 若错误出现在外部库函数中，应检查你自己代码传入的参数。
- 修错顺序
  - 当有多个错误时，从最早出现的错误开始修复，后续错误往往由它引起。
- 检查可复现性
  - 重新启动 Python 环境，看错误是否仍存在，以排查全局变量或作用域污染问题。
- 隔离错误源
  - 逐步删除或添加代码（类似“二分查找”策略）以定位问题所在。
  - 模块化设计（函数化）可单独测试各部分逻辑。
- 传统方法
  - 插入 `print()` 输出变量值以追踪执行流。
  - 尽管有更好工具，这一策略仍在简单情况下有效。
- 逐行执行
  - Python 可逐行执行代码，适合快速排查简单问题。
  - 若错误涉及复杂嵌套函数或变量作用域，建议使用调试器。

### 2.1.1 Using pdb

- 激活调试器的方法
  1. 插入 `breakpoint()` 或 `import pdb; pdb.set_trace()`。
  2. 运行后出错时使用 `pdb.pm()` 进入调试模式。
  3. 在 Jupyter 中使用 `%debug` 魔法命令。
  4. 用 `pdb.run("function_call")` 控制函数执行。
  5. 启动时直接进入调试模式：



```
1 python -m pdb file.py
```

- 示例：使用 `breakpoint()`



```
1 import run_with_break as run
2 run.fit(run.data, run.n_cats)
```

运行后进入交互模式（Pdb），可使用：

- `n`：执行当前行并到下一行
- `c`：继续到下一个断点
- `p var`：打印变量

### 2.1.2 Post-mortem Debugging（事后调试）

- 方法

Python

```
1 import pdb
2 import run_no_break as run
3 run.fit(run.data, run.n_cats)
4 pdb.pm()
```

在错误发生点自动进入调试器。

- **导航堆栈**

- `u`: 上移一层到用户代码。
- `l`: 查看周围代码。
- `p var`: 打印变量值。

### 2.1.3 常用 `pdb` 命令速查表

命令	含义
<code>h / help</code>	显示所有命令
<code>l / list</code>	查看当前行附近代码
<code>p / print</code>	打印变量
<code>n / next</code>	执行当前行
<code>s / step</code>	进入函数内部
<code>r / return</code>	跳出当前函数
<code>c / continue</code>	继续执行至下个断点
<code>b / break</code>	设置断点
<code>tbreak</code>	临时断点
<code>u / d</code>	在调用栈中上/下移动
<code>where</code>	查看调用栈
<code>q</code>	退出调试
<code>&lt;Enter&gt;</code>	重复上一命令

## 2.2 常见错误来源

- 括号不匹配 (parenthesis mismatch)
- `==` 与 `=` 混淆
- 浮点比较 (`==` 不可靠)
- 返回值类型或形状与预期不符
- 错用函数或变量名
- 无名参数顺序错误
- 变量未定义, Python 从外层作用域取到错误值 (典型作用域错误)
- Python 自动降维导致维度不一致

## 2.3 防止与捕获错误的技巧

### 2.3.1 Defensive Programming (防御式编程)

- 检查函数输入的有效性。

- 提供合理默认值。
- 对输入/输出进行范围验证。
- 用 `assert`、`try` 或 `raise` 抛出带信息的错误。
- 示例：

```
Python
1 import warnings
2
3 def mysqrt(x):
4     if isinstance(x, str):
5         raise TypeError(f"What is the square root of '{x}'?")
6     if isinstance(x, (float, int)):
7         if x < 0:
8             warnings.warn("Input value is negative.", UserWarning)
9             return float('nan')
10        return x**0.5
11    else:
12        raise ValueError(f"Cannot take the square root of {x}")
```

## 2.3.2 try/except 捕获运行时错误

```
Python
1 try:
2     model = statsmodels.api.OLS(sub['y'], statsmodels.api.add_constant(sub['x']))
3     fit = model.fit()
4     params[i, :] = fit.params.values
5 except Exception as error:
6     print(f"Regression cannot be fit for stratum {i}.")
```

- 若某分层无数据导致回归失败，程序仍能继续运行。

## 2.3.3 维度保持 (Dimensionality Handling)

- Numpy 常会自动压缩多余维度，导致错误：

```
Python
1 mat = np.array([[1, 2], [3, 4]])
2 mat2 = mat[1, :]      ## mat2 维度变为 (2,)
```

- 修复：

```
Python
1 if len(mat2.shape) != 2:
2     mat2 = mat2.reshape(1, -1)
```

## 2.3.4 避免使用全局变量

- 全局变量使代码不可预测，易受污染。
- 清理环境变量或在新 session 中运行函数以检测依赖。

```
1 del x
2
3 def f(z):
4     y = 3
5     print(x + y + z)
6 ## NameError: x is not defined
```

## 2.3.5 其他调试与编码建议

- 优先使用标准库或成熟算法包。
- 模块化设计：小函数易测试与复用。
- 先保证正确性与可读性，再考虑优化。
- 逐步构建代码并测试中间结果。
- 为 `if`、`for` 等逻辑写健壮条件判断。
- 避免硬编码数值（如 `3e8` 应定义为 `speed_of_light`）。
- 提早编写测试（unit tests）。