

Lesson 46 Summary: Project Work 2

Data Science with Python – Key Concepts

Data Science Program

ML Project Best Practices

Reproducibility

- Random seeds
- Version control
- Requirements.txt

Validation

- Train/test split
- Cross-validation
- Holdout test set

Features

- Feature selection
- Importance ranking
- Domain knowledge

Model Choice

- Compare multiple
- Ensemble methods
- Hyperparameter tuning

Code Quality

- Modular functions
- Clear comments
- Unit tests

Results

- Clear metrics
- Visualizations
- Business impact

Essential for trust:

- **Seeds:** Set random seeds everywhere
- **Version:** Git for code, DVC for data
- **Environment:** requirements.txt/conda.yml

Anyone should get the same results

Professional standards:

- **Modular:** Functions, not scripts
- **Documented:** Docstrings, comments
- **Tested:** Unit tests for key functions

Good code = maintainable code

Systematic approach:

- **Compare:** Try multiple algorithms
- **Validate:** Use cross-validation
- **Tune:** Grid/random search

The best model depends on the problem

Complete the story:

- **README:** Project overview, setup
- **Notebooks:** Exploratory analysis
- **Report:** Methodology, results

Document as you go, not at the end

Project Essentials:

Component	Purpose
requirements.txt	Reproducible environment
README.md	Project documentation
src/	Production code
notebooks/	Exploration/analysis
tests/	Unit tests
models/	Saved model artifacts

Standard project structure enables collaboration