

# Outils

Environnement, installation, IDE, versionning, qualité, DQ, exécution

## Programme Data Engineer

Public intermédiaire — Travaux pratiques

23 octobre 2025

# Pile logicielle du TP

- **Langage** : Python 3.10+ (numpy, pandas, matplotlib).
- **IDE/Notebooks** : VS Code (extensions Python, Jupyter) ou JupyterLab.
- **Stockage local** : CSV/Parquet ; option **DuckDB** embarqué pour SQL.
- **Versionning** : Git (Github/GitLab), branches feature.
- **Qualité** : black, isort, ruff, pre-commit ; tests pytest.
- **DQ (option)** : Great Expectations ou Soda.

# Installation rapide (venv/pip) et alternative Conda

## Option A — venv + pip

```
1 python3 -m venv .venv
2 source .venv/bin/activate    # Windows: .venv\\Scripts\\activate
3 python -m pip install --upgrade pip
4 pip install numpy pandas matplotlib duckdb pytest black isort ruff
   pre-commit
```

## Option B — Conda/Miniconda

```
1 conda create -n dw_tp python=3.11 -y
2 conda activate dw_tp
3 pip install numpy pandas matplotlib duckdb pytest black isort ruff
   pre-commit
```

# Structure de projet recommandée

```
1 mini_dw_tp/  
2   data/  
3     raw/           # fichiers sources (CSV)  
4     staging/       # nettoyages intermédiaires  
5     dim/           # dimensions persistes  
6     fact/          # faits persistés  
7 notebooks/  
8 src/  
9   etl/  
10    __init__.py  
11    extract.py      # lecture CSV/Parquet  
12    transform.py    # nettoyage, assemblage  
13    load.py         # persist CSV/Parquet/DuckDB  
14    utils/  
15      io.py, dq.py, logging_conf.py  
16 tests/  
17 Makefile  
18 requirements.txt (ou pyproject.toml)
```

# Makefile : commandes pratiques

```
1 .PHONY: setup format lint test run
2 setup:
3     python -m pip install -r requirements.txt
4     pre-commit install
5
6 format:
7     black src notebooks
8     isort src notebooks
9
10 lint:
11     ruff check src notebooks
12     ruff format --check
13     ruff fix --select I --unsafe-fixes
14     ruff check --select F401 # imports inutiles
15     ruff check --select E --select W
16     ruff check --select I
17     ruff check --statistics
18     ruff check --show-source
```

# Git : bonnes pratiques & .gitignore

**Flux** : main (protégée) → *branchesfeature/..* → *PR&review*.

## DuckDB (option) : moteur SQL embarqué

```
1 import duckdb
2 con = duckdb.connect("mini_dw.duckdb")
3 con.execute("CREATE SCHEMA IF NOT EXISTS dw")
4 con.execute("CREATE TABLE IF NOT EXISTS dw.dim_product AS SELECT *
5             FROM read_csv_auto('data/dim/dim_product.csv')")
6 # Requeter le CA par mois
7 con.sql(
8     """
9     SELECT strftime(order_date, '%Y-%m') AS ym, SUM(amount) AS
10        revenue
11    FROM read_csv_auto('data/fact/fact_sales.csv')
12    GROUP BY ym ORDER BY ym
13    """
14 ).df()
```

# Tests (pytest) & Data Quality

## pytest — exemple minimal

```
1 # tests/test_amount.py
2 import pandas as pd
3
4 def test_amount_non_negative():
5     df = pd.read_csv('data/fact/fact_sales.csv')
6     assert (df['amount'] >= 0).all()
```

## Great Expectations — CLI (option)

```
1 # initialisation dans le repo
2 great_expectations init
3 # cr er une suite et valider un fichier
4 great_expectations suite new
5 great_expectations checkpoint run my_checkpoint
```



# Formatage, lint & hooks pre-commit

```
1 # .pre-commit-config.yaml (extrait)
2 repos:
3   - repo: https://github.com/psf/black
4     rev: 24.8.0
5     hooks: [{id: black}]
6   - repo: https://github.com/charliermarsh/ruff-pre-commit
7     rev: v0.5.5
8     hooks: [{id: ruff}]
9   - repo: https://github.com/PyCQA/isort
10     rev: 5.13.2
11     hooks: [{id: isort}]
```

# Logs & configuration (.env)

```
1 # src/utis/logging_conf.py
2 import logging, os
3 LOG_LEVEL = os.getenv("LOG_LEVEL", "INFO")
4 logging.basicConfig(
5     level=LOG_LEVEL,
6     format="%(asctime)s %(levelname)s [%(name)s] %(message)s",
7 )
8 logger = logging.getLogger("dw_tp")
```

```
1 # .env (exemple)
2 LOG_LEVEL=INFO
3 DATA_DIR=./data
```

# Exécution : scripts, notebooks, VS Code

- **Scripts** : `python -m src.etl.load` (réutilisable, testable).
- **Notebooks** : exploration/validation (sauver en `.py` pour versionning).
- **VS Code** : tasks, debug config, intégration Git, extensions Python/Jupyter.
- **Make** : `make setup` | `make run` | `make test` | `make lint`.

# Schéma d'outillage du TP

