

Pourquoi utiliser Python pour des projets data

Pourquoi choisir Python pour les projets data ?

Python est devenu le langage de référence dans le domaine de la data, que ce soit pour la data analysis, le machine learning, le traitement de données massives ou encore le déploiement d'APIs de données. Voici pourquoi.

Simplicité et lisibilité

- Python a une syntaxe claire, proche du langage naturel.
- Il est facile à lire, à apprendre et à maintenir, ce qui est crucial quand plusieurs personnes (scientifiques, analystes, ingénieurs) collaborent sur les mêmes scripts.

Écosystème riche et mature

- Python dispose d'une immense bibliothèque d'outils pour la data :
 - **NumPy** et **Pandas** : traitement et manipulation de données.
 - **Matplotlib**, **Seaborn**, **Plotly** : visualisation.
 - **Scikit-learn**, **XGBoost** : machine learning.
 - **TensorFlow**, **PyTorch** : deep learning.
 - **SQLAlchemy**, **Psycopg2**, **Pymongo** : connexion aux bases de données.
 - **Airflow**, **Luigi** : orchestration de pipelines.
 - **Dask**, **PySpark** : traitement distribué.

Interfaçage facile avec d'autres outils

- Python s'intègre facilement avec des bases de données, des APIs REST, des outils cloud, des scripts shell ou du code C/C++.
- Il est compatible avec tous les grands services cloud (AWS, Azure, GCP).

Un langage adapté à toutes les étapes de la data

- Du **nettoyage des données** à la **modélisation**, en passant par la **visualisation** et le **déploiement**, Python couvre l'ensemble de la chaîne de valeur data.
- Il est utilisé aussi bien par les **data analysts**, **data engineers** que par les **data scientists**.

Une grande communauté

- Des milliers de ressources, tutoriels, notebooks et projets open-source sont disponibles.
- Des réponses quasi instantanées sur StackOverflow, GitHub ou Medium.

Exemple : un petit projet Python de data exploration

Imaginons qu'on veuille charger des données CSV, les analyser, les nettoyer, puis les enregistrer dans une base PostgreSQL. Voici un exemple simple avec `pandas` et `sqlalchemy`.

Structure du projet

```
mon-projet-data/  
├── app/  
│   ├── main.py  
│   └── requirements.txt  
├── data/  
│   └── ventes.csv  
├── Dockerfile  
└── docker-compose.yml
```

Contenu du `ventes.csv`

```
1  date,produit,quantite,prix_unitaire  
2  2024-01-01,Stylo,10,1.5  
3  2024-01-01,Cahier,5,2.0  
4  ...
```

Contenu du `main.py`

```
1  import pandas as pd  
2  from sqlalchemy import create_engine  
3  
4  # Chargement du CSV  
5  df = pd.read_csv("data/ventes.csv")  
6  
7  # Nettoyage rapide  
8  df.dropna(inplace=True)  
9  df['total'] = df['quantite'] * df['prix_unitaire']  
10  
11 # Connexion à PostgreSQL via SQLAlchemy  
12 engine = create_engine("postgresql://myuser:mypassword@db:5432/mydb")  
13 df.to_sql("ventes", engine, if_exists="replace", index=False)  
14  
15 print("Données enregistrées avec succès.")  
16
```

Contenu du requirements.txt

```
1 pandas
2 sqlalchemy
3 psycopg2-binary
4
```

Contenu du Dockerfile

```
1 FROM python:3.11-slim
2
3 WORKDIR /app
4
5 COPY app/requirements.txt .
6 RUN pip install --no-cache-dir -r requirements.txt
7
8 COPY app/ .
9 COPY data/ ../data/
10
11 CMD ["python", "main.py"]
12
```

Contenu du docker-compose.yml

```
1  version: '3.9'
2
3  services:
4    app:
5      build: .
6      depends_on:
7        - db
8      volumes:
9        - ./data:/data
10     networks:
11       - backend
12
13     db:
14       image: postgres:15
15       environment:
16         POSTGRES_DB: mydb
17         POSTGRES_USER: myuser
18         POSTGRES_PASSWORD: mypassword
19       volumes:
20         - pgdata:/var/lib/postgresql/data
21       networks:
22         - backend
23
24     volumes:
25       pgdata:
26
27     networks:
28       backend:
29
```

Lancement du projet

```
git:(master) docker-compose up --build
```

Conclusion

Python s'impose dans les projets data pour sa **simplicité**, sa **puissance**, et son **écosystème extrêmement riche**. Il est à la fois un outil d'exploration, d'analyse, de traitement, de modélisation et de déploiement. Il permet de transformer rapidement une idée en solution opérationnelle, même à grande échelle.