

# **Documentation Technique**

Pipeline d'Ingestion Crypto vers Azure Data Lake

**Projet Data Engineering - Niveau 1**

**Yoann LEHONG CHEFFSON**

19 décembre 2025

## Table des matières

<b>1</b>	<b>Introduction</b>	<b>3</b>
1.1	Objectifs du projet . . . . .	3
<b>2</b>	<b>Architecture Technique</b>	<b>3</b>
2.1	Flux de Données (Data Flow) . . . . .	3
2.2	Stack Technologique . . . . .	3
<b>3</b>	<b>Installation et Configuration</b>	<b>3</b>
3.1	Prérequis . . . . .	3
3.2	Mise en place de l'environnement . . . . .	4
3.3	Configuration de la Sécurité (.env) . . . . .	4
<b>4</b>	<b>Structure du Code (Extract &amp; Load)</b>	<b>4</b>
4.1	Chargement Sécurisé de la Configuration . . . . .	4
4.2	Simulation des Données (Mocking) . . . . .	4
4.3	Ingestion vers Azure (Upload) . . . . .	5
<b>5</b>	<b>Exécution et Déploiement</b>	<b>5</b>
5.1	Lancement du Pipeline . . . . .	5
5.2	Gestion de version (Git) . . . . .	5
<b>6</b>	<b>Conclusion</b>	<b>5</b>

## 1 Introduction

Ce document technique décrit l'architecture et le fonctionnement du pipeline d'ingestion de données (ETL) conçu pour le suivi des marchés de cryptomonnaies.

Dans un contexte où la fiabilité des données est critique, ce projet simule une récupération de données de marché en temps réel et assure leur stockage sécurisé et pérenne dans un environnement Cloud Enterprise (\*\*Azure Data Lake Gen2\*\*).

### 1.1 Objectifs du projet

- **Extraction** : Simuler des données de marché réalistes (Bitcoin, Ethereum, Solana, BNB) pour pallier les limitations des API publiques.
- **Transformation** : Structurer et nettoyer les données via Pandas.
- **Chargement (Load)** : Automatiser l'envoi sécurisé vers le Cloud Microsoft Azure.
- **Sécurité** : Implémenter les bonnes pratiques DevSecOps (gestion des secrets via .env).

## 2 Architecture Technique

### 2.1 Flux de Données (Data Flow)

Le pipeline suit une architecture linéaire de type ETL :

1. **Source** : Générateur de données Python (Simulation API).
2. **Processing Local** : Script Python + Pandas (Création du CSV).
3. **Transport** : SDK Azure Storage via protocole HTTPS sécurisé.
4. **Destination** : Azure Storage Account (Conteneur `raw-data`).

### 2.2 Stack Technologique

**Langage** Python 3.x

**Cloud Provider** Microsoft Azure

**Librairies Clés**

- `pandas` : Manipulation et structuration des données.
- `azure-storage-blob` : Interface client pour Azure Data Lake.
- `python-dotenv` : Gestion de la sécurité et des variables d'environnement.

**Infrastructure** Azure Data Lake Storage Gen2 (Hierarchical Namespace activé).

## 3 Installation et Configuration

### 3.1 Prérequis

- Python 3.9 ou supérieur installé.
- Un compte Microsoft Azure avec une souscription active.
- Un compte de stockage Azure créé avec un conteneur nommé `raw-data`.

### 3.2 Mise en place de l'environnement

Il est recommandé d'utiliser un environnement virtuel pour isoler les dépendances du projet.

```

1 # 1. Cr ation de l'environnement virtuel
2 python -m venv venv
3
4 # 2. Activation (Windows Powershell)
5 .\venv\Scripts\Activate
6
7 # 3. Installation des librairies requises
8 pip install pandas azure-storage-blob python-dotenv

```

Listing 1 – Installation des dépendances

### 3.3 Configuration de la Sécurité (.env)

**Avertissement de sécurité :** Les clés d'accès Azure (Connection Strings) sont des informations sensibles. Elles ne doivent jamais être écrites en dur dans le code source ni committées sur GitHub.

Nous utilisons un fichier `.env` (exclu par `.gitignore`) pour stocker ces secrets localement.

```

1 AZURE_CONNECTION_STRING="DefaultEndpointsProtocol=https;AccountName=votrecompte;
  AccountKey=votre_cle_secrete_tres_longue==;EndpointSuffix=core.windows.net"

```

Listing 2 – Format du fichier `.env`

## 4 Structure du Code (Extract & Load)

Le script `main.py` orchestre l'ensemble du pipeline. Voici les composants critiques.

### 4.1 Chargement Sécurisé de la Configuration

Le code initialise l'environnement en chargeant les variables cachées avant toute opération.

```

1 from dotenv import load_dotenv
2 import os
3
4 # Chargement des variables d'environnement depuis .env
5 load_dotenv()
6
7 # R cup ration s curis e de la cha ne de connexion
8 AZURE_CONNECTION_STRING = os.getenv("AZURE_CONNECTION_STRING")
9 CONTAINER_NAME = "raw-data"

```

### 4.2 Simulation des Données (Mocking)

Cette fonction remplace l'appel API pour garantir la stabilité du pipeline lors du développement, indépendamment des quotas d'API ou des pannes réseau.

```

1 def get_mock_data():
2     """G n ration de donn es al atoires r alistes"""
3     fake_data = []
4     base_prices = {'bitcoin': 95000, 'ethereum': 3500}
5
6     # Logique de variation al atoire (+/- 5%)
7     # ...
8     return fake_data

```

### 4.3 Ingestion vers Azure (Upload)

La fonction d'envoi utilise le SDK Azure Blob Storage pour transférer le fichier CSV généré.

```

1 def upload_to_azure(local_file_name):
2     try:
3         # Cr ation du client de service
4         blob_service_client = BlobServiceClient.from_connection_string(
5             AZURE_CONNECTION_STRING)
6
7         # Ciblage du conteneur et du fichier
8         blob_client = blob_service_client.get_blob_client(
9             container=CONTAINER_NAME,
10            blob=local_file_name
11        )
12
13        # Envoi des donn es binaires
14        with open(local_file_name, "rb") as data:
15            blob_client.upload_blob(data, overwrite=True)
16
17        print("Succ s ! Fichier transf r sur Azure.")
18
19    except Exception as e:
20        print(f"Erreur Critique Azure : {e}")

```

## 5 Exécution et Déploiement

### 5.1 Lancement du Pipeline

Pour exécuter le script manuellement depuis le terminal :

```
1 python main.py
```

### 5.2 Gestion de version (Git)

Le projet respecte les standards DevOps. Le fichier `.gitignore` est configuré pour exclure les fichiers sensibles et temporaires :

```

1 venv/          # Environnement virtuel
2 .env           # Secrets (CRUCIAL)
3 *.csv          # Donn es locales
4 --pycache__/  # Cache Python
5 *.db           # Bases de donn es locales

```

## 6 Conclusion

Ce projet valide les compétences fondamentales d'un Data Engineer :

1. Capacité à scripter un pipeline ETL en Python.
2. Maîtrise de l'interaction avec un Cloud Provider (Azure).
3. Application stricte des normes de sécurité des données.

Le pipeline est désormais prêt pour être étendu avec des outils de visualisation (Power BI) ou de transformation avancée (Databricks).