

Mastere 2 Data Engineer

Expert informatique et système d'information - Titre certifié niveau 7

ENERGIEWATCH



ADIB Melissa
DIALLO Alimou M.
GBE Keagnon Grâce H.

SOMMAIRE

1. Une plongée dans Energiewatch
2. Energiewatch, quelle réponse ?
3. Productions et objectifs
4. Notre Business model
5. Gestion de projet
6. Architecture
7. Flux de Travail : De l'ETL au déploiement cloud
8. Modèle Machine Learning
9. Empreinte carbone
10. Intégration CI
11. Interface utilisateur
12. Difficultés rencontrés et évolutions
13. Maintenance de l'application
14. Démo

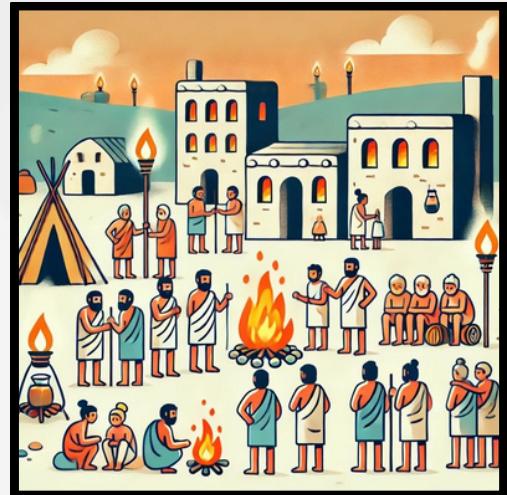


1.

UNE PLONGÉE DANS **ENERGIEWATCH**

UN PEU D'HISTOIRE...

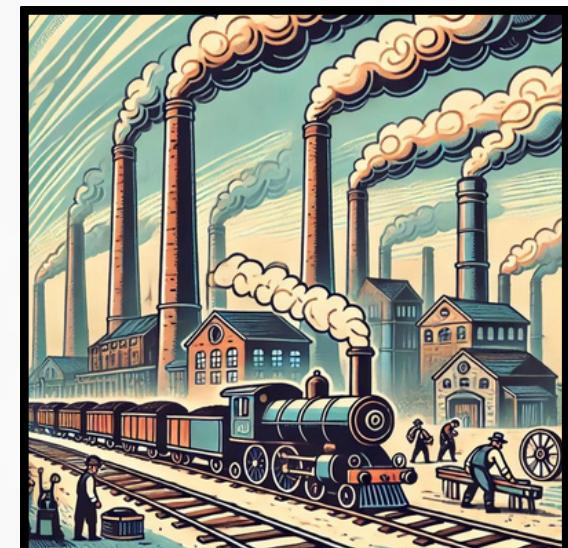
**3000 av. J.-C. à
500 apr. J.-C.**



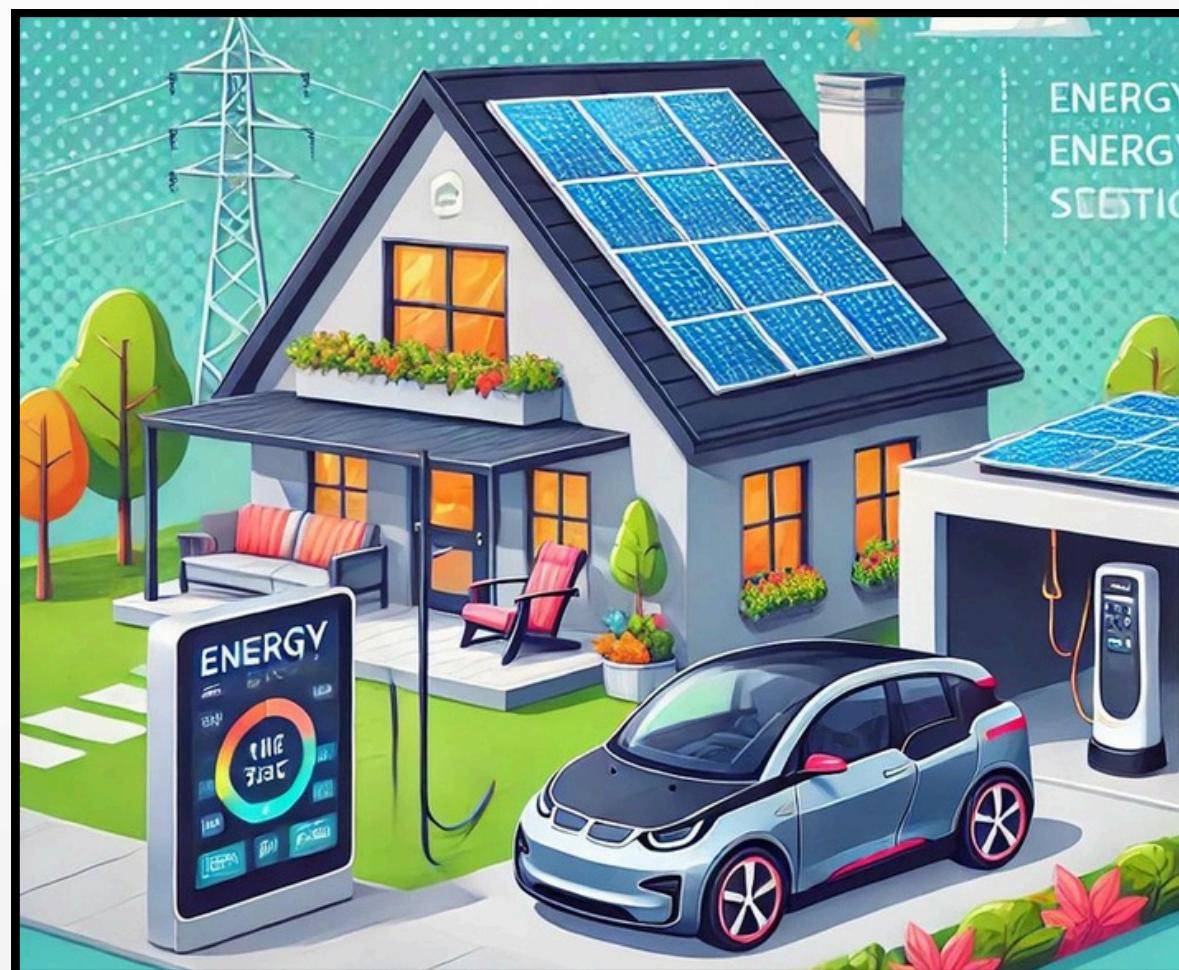
500 à 1500 apr. J.-C.



1750 à 1850



2000 à nos jours



CONTEXTE

 Hausse des Prix de l'Énergie

 Complexité Croissante des Réseaux

 L'Ère des Données

 Urgence Environnementale
Transition Énergétique



QUELS ENJEUX ?

Détection des Anomalies

Optimisation de la Consommation

Prévision de la Demande

Sécurité des Données



2.

ENERGIEWATCH, QUELLE RÉPONSE ?

QUELLE RÉPONSE ?

Construire un avenir énergétique plus durable



COMMENT ?

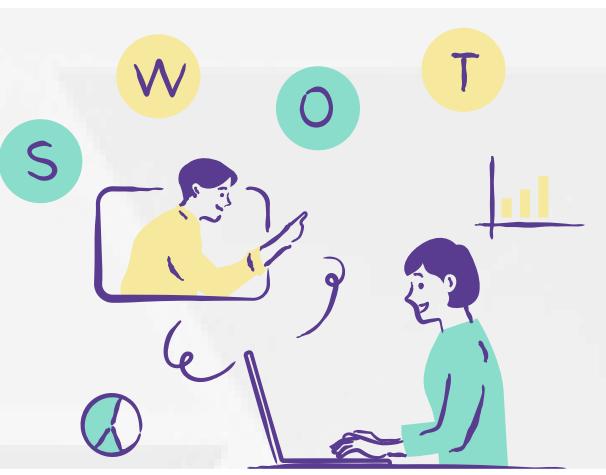
Optimisation gestion énergétique

Surveillance continue

Détection d'anomalies en temps réel

Anticipation des besoins énergétiques

ANALYSE



S

- Innovation et Flexibilité
- Facilité d'Utilisation
- Adaptabilité
- Orientation Durable

W

- Manque d'Isolation des Données
- Complexité de Gestion des Mises à Jour

O

- Marché en Croissance
- Partenariats Stratégiques
- Éducation du Marché
- Évolution Technologique

T

- Dépendance aux Données
- Vulnérabilité aux Cyberattaques
- Absence de Système de Prévention
- Scalabilité Limitée

LA CONCURRENCE

Quelques solutions existantes dans le domaine
de la détection d'anomalies



Anodot



Amazon Lookout for Metrics

3.

OBJECTIFS ET LIVRABLES

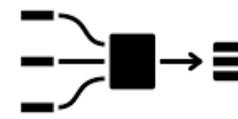
OBJECTIFS

Mise en place d'un outil qui permet :

-  La prédition de la consommation
-  La détection des Anomalies
-  Le calcul et le suivi de l'empreinte carbone

LIVRABLES

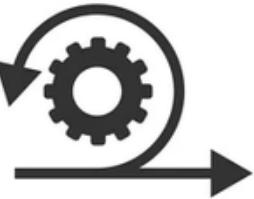
Pour nos équipes



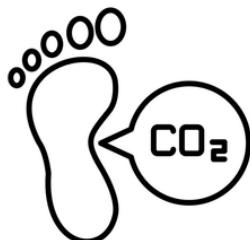
Pipeline automatisé



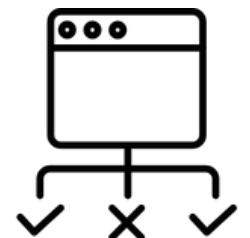
Modèles ML



Scripts pour le ré-entraînement



Logs empreinte carbone



Tests unitaires



Documentation Technique

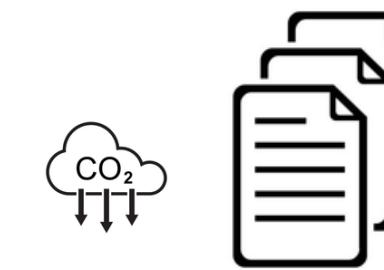
Pour les clients



Tableau de bord interactif



Rapports de prédition



Rapport d'empreinte carbone



Modèles prédictifs

CONSTRAINTES

Plusieurs contraintes peuvent être rencontrées à différents niveaux :



% Accuracy



Disponibilité



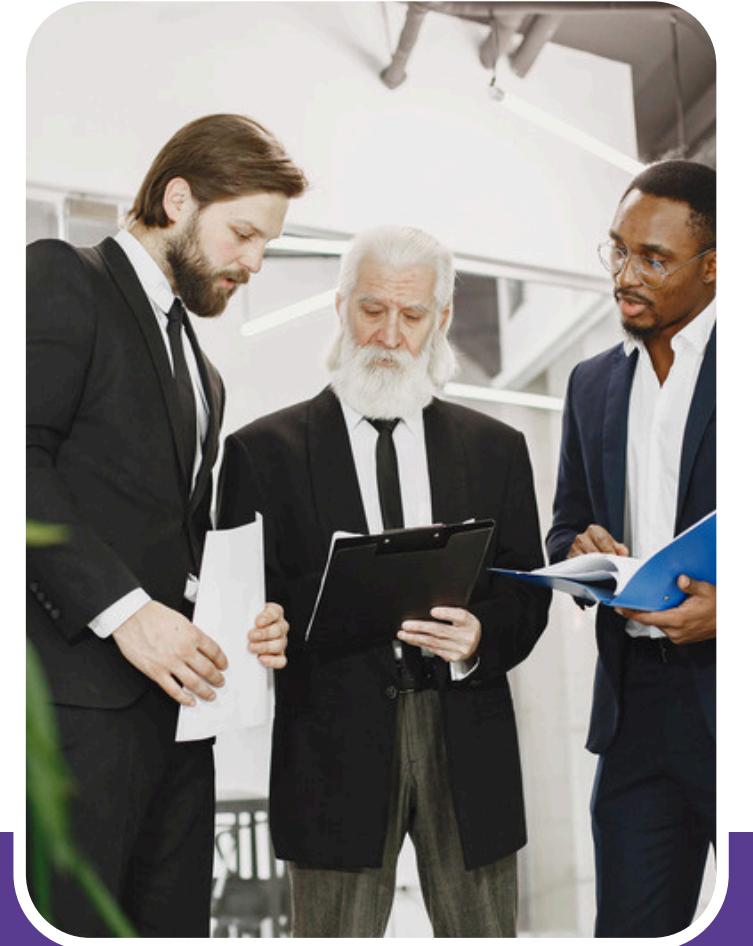
Environnementales



Stockage des logs



Anomalie critique



4.

NOTRE BUSINESS MODEL



MODÈLE ÉCONOMIQUE

OBJECTIF DE RENTABILITÉ



- **Abonnements Mensuels/Annuels :**
 - **Basique** : 700 EUR/mois
 - **Avancé** : 1500 EUR/mois
 - **Entreprise** : 3000 EUR/mois
- **Services Additionnels :**
 - Consulting
 - Formations
- **Premium**



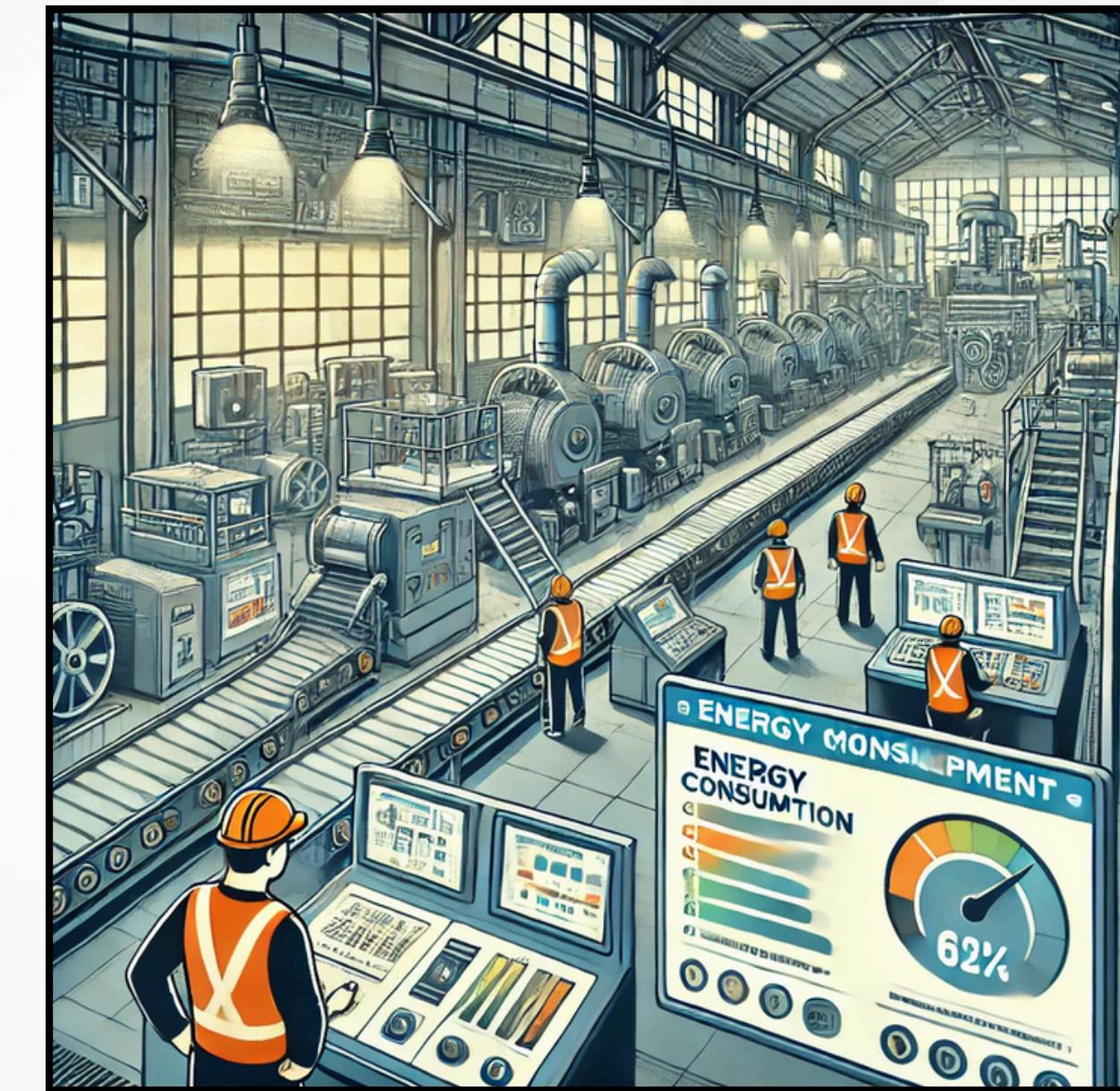
- **Objectif de Rentabilité :**
 - **Abonnement basique** : 8 400 EUR/an
 - **Abonnement Niveau Avancé** : 18 000 EUR/an
 - **Abonnement Niveau Entreprise** : 36 000 EUR/an

CIBLE MARKET

Entreprises d'Électricité



Industries à Forte Consommation





HELENA 40 ans

Profession : Directrice des Opérations d'une Entreprise Énergétique

Contexte : Gère la production et la distribution d'énergie dans une grande entreprise d'électricité, avec la nécessité de s'adapter à des fluctuations imprévisibles de la demande causées par des événements sociaux et des changements météorologiques.

Objectifs :

- Optimiser l'ajustement de la production énergétique en fonction de la demande tout en réduisant les coûts et en améliorant l'efficacité.
- Répondre rapidement aux perturbations causées par les conditions météorologiques extrêmes ou les mouvements sociaux.
- Suivre et réduire l'empreinte carbone associée aux opérations énergétiques.

Besoins :

- Une solution qui intègre des variables sociales et météorologiques pour mieux prédire la demande énergétique.
- Un outil qui permet de détecter rapidement des anomalies dans les données de consommation et de production.
- Un suivi précis de l'empreinte carbone des traitements et des infrastructures IT pour mieux aligner les opérations avec les objectifs environnementaux.

Ce que EnergieWatch peut apporter :

- Prédiction fiables basées sur des facteurs sociaux et climatiques pour mieux planifier la production.
- Détection des anomalies dans les données en temps réel pour prévenir les problèmes avant qu'ils n'impactent les opérations.
- Calcul et suivi de l'empreinte carbone des traitements énergétiques, aidant l'entreprise à réduire son impact environnemental.

5.

GESTION DE PROJET

APPROCHE

Identification du besoin

Etablir les exigences

Planification et suivi



- KANBAN
- Planification des tâches



- Canal de discussion
- Planification des sprints
- Partage de documents



Apache
Airflow



 MongoDB
Atlas



 kedro



 + 
elasticsearch kibana

 Streamlit

 mlflow™



 S3
amazon
web services

RESSOURCES LOGICIELLES

BUDGETS ALLOUÉ AUX RESSOURCES



- **Coût mensuel :** 500 --- 1 000 EUR/mois
- **Coût annuel :** 6 000 et 12 000 EUR.



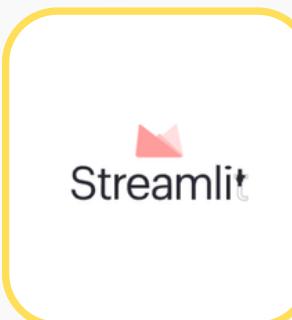
- **Coût mensuel :** 21 EUR/utilisateur.
- **Coût annuel (pour 5 utilisateurs) :** 1 260 EUR.



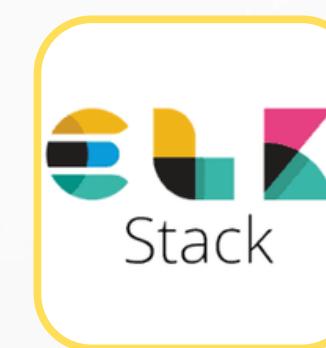
- **Coût mensuel :** 200 --- 500 EUR/mois
- **Coût annuel :** 2 400 et 6 000 EUR.



- **Coût mensuel :** Entre 200 et 1 000 EUR/mois en fonction du niveau de ressources alloué (espace, nombre de clusters).
- **Coût annuel :** Entre 2 400 et 12 000 EUR



- **Coût mensuel :** Entre 200 et 1 000 EUR/mois en fonction du niveau de ressources alloué (espace, nombre de clusters).
- **Coût annuel :** Entre 2 400 et 12 000 EUR



- **Coût mensuel :** Environ 100 à 500 EUR/mois pour une utilisation standard.
- **Coût annuel :** Entre 1 200 et 6 000 EUR.



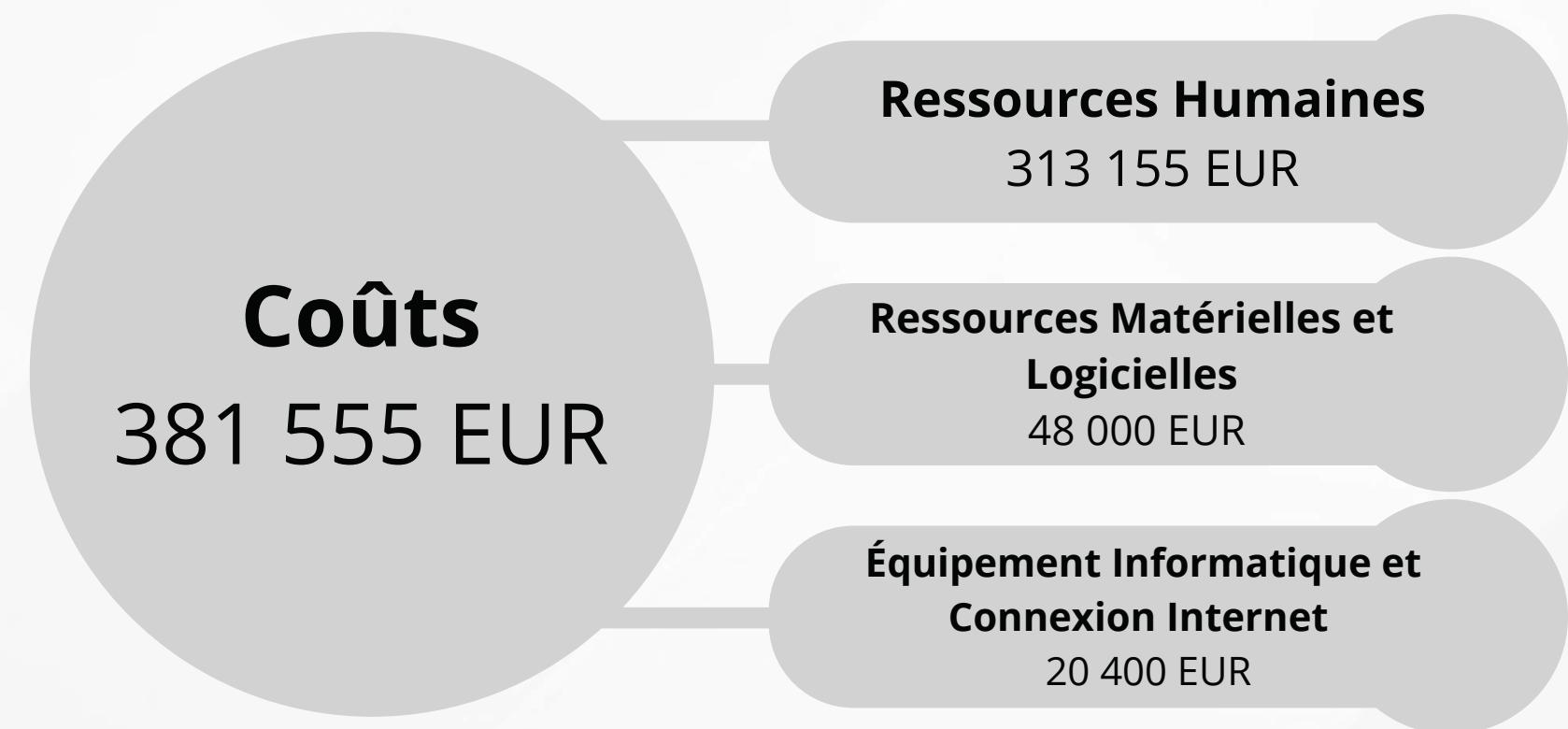
- **Coût mensuel :** Varie grandement selon les ressources utilisées. On peut estimer environ 500 à 1 500 EUR/mois.
- **Coût annuel :** Entre 6 000 et 18 000 EUR.

Coût minimum :
21 660 EUR
Coût maximum :
67 260 EUR

RESSOURCES HUMAINES

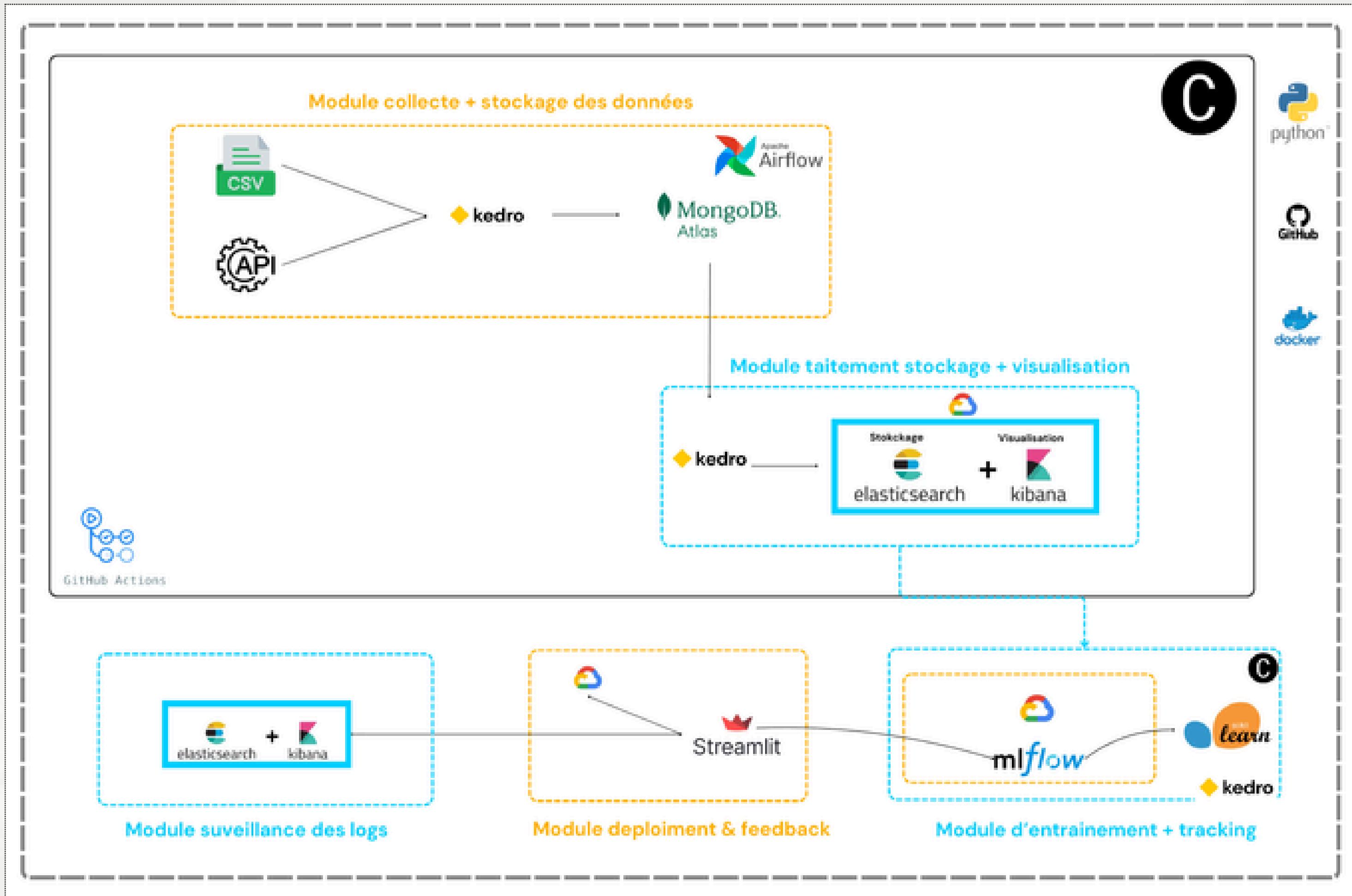
RÔLE	RESPONSABILITÉS	COMPÉTENCES	COÛT ANNUEL (EUR)
Data Analyst	<ul style="list-style-type: none"> Collecter et analyser données Préparer des rapports et des visualisations 	<ul style="list-style-type: none"> Python (pandas, NumPy) Tableau, Power BI ou Kibana 	60,269
Data Engineer	<ul style="list-style-type: none"> Concevoir, développer et maintenir les pipelines de données Assurer la qualité et l'intégrité des données 	<ul style="list-style-type: none"> MongoDB, SQL Kedro, Airflow, oozie GCP, AWS Python, Java, Docker 	60,269
Ingénieur DevOps	<ul style="list-style-type: none"> Automatisation des Processus de Déploiement Gestion des Infrastructures 	<ul style="list-style-type: none"> Jenkins, GitHub action AWS, GCP, Prometheus Docker, ELK, Grafana 	67,519
Data scientist	<ul style="list-style-type: none"> Développer, entraîner, et surveiller les modèles ML Gérer les expérimentations Mettre en production les modèles 	<ul style="list-style-type: none"> scikit-learn, TensorFlow Docker, Mlflow, Airflow Deep learning 	67,519
Coordinateur / PO	<ul style="list-style-type: none"> Gestion de projet, coordination des équipes Interaction avec les parties prenantes 	<ul style="list-style-type: none"> Méthode Agiles 	60,269

STRUCTURE DES COÛTS



6.

ARCHITECTURE



ROADMAP

	PHASE 1			PHASE 2			PHASE 3		
	Janvier	Fevrier	Mars	Avril	Mai	Juin	Juillet	Aout	Septembre
Découpage des tâches et deadline									
Récupération et Ingestion des Donnée									
Création du Pipeline ETL avec Kedro									
Test unitaire et intégration CI									
Création du pipeline fusion des données avec Kedro									
Développement des modèles de ML									
Mise en place d'un serveur Mlflow, Airflow et optimisation des modèles									
Mise en place de tests CI avancés									
Création des dashboards avec Streamlit/Kibana									
Création du Makefile et finalisation des CI									
Mise en place d'un système de monitoring (logs, métriques, ELK)									
Calcul de l'empreinte carbone									

7.

FLUX DE TRAVAIL : DE L'ETL AU DÉPLOIMENT CLOUD

Do not use SQLite as metadata DB in production – it should only be used for dev/testing. We recommend using Postgres or MySQL. [Click here](#) for more information.

Do not use the SequentialExecutor in production. [Click here](#) for more information.

DAGs

All (2)	Active (1)	Paused (1)	Running (1)	Failed (0)	Filter DAGs by tag	Search DAGs	Auto-refresh	Links
DAG	Owner	Runs	Schedule	Last Run	Next Run	Recent Tasks	Actions	Links
Prix_des_carburants_en_France	airflow	1 (1)	*/10 * * * *	2024-09-11, 21:26:32	2024-09-11, 22:00:00	1 (1)		
consommation_gaz_industriel	airflow	1 (1)	*/10 * * * *	2024-09-11, 22:15:51	2024-09-10, 08:00:00	1 (1)		

Showing 1-2 of 2 DAGs

AIRFLOW

Airflow DAGs Cluster Activity Datasets Security Browse Admin Docs

22:09 CEST (+02:00) AD

DAG: Prix_des_carburants_en_France Un DAG pour charger les données de l'API dans MongoDB

Schedule: "00 * * * *"
Next Run ID: 2024-09-11, 22:09:00 CEST

09/11/2024 10:05:03 PM All Run Types All Run States Clear Filters

Duration: 00:00:15 API to MongoDB Task

Run History: 25 runs

Run States: Failed (1), Pending (1), Running (1), Scheduled (1), Skipped (1), Success (20), Up for reschedule (1), Up for retry (1), Upstream failed (1), No plan (1)

DAG: Prix_des_carburants_en_France

Details, Flow Graph, Gantt, Code, Event Log, Run Duration, Task Duration, Calendar

DAG Runs Summary

Total Runs Displayed: 25

Total success: 20

First Run Start: 2024-09-11, 21:39:05 CEST

Last Run Start: 2024-09-11, 22:04:34 CEST

Max Run Duration: 00:00:15

Mean Run Duration: 00:00:06

Min Run Duration: 00:00:05

DAG Summary

Total Tasks: 1

PythonOperator: 1

DAG Details

GOOGLE COLAB

Colab Enterprise

NOTEBOOKS ENVIRONNEMENTS D'EXÉCUTION MODÈLES D'EXÉCUTION EXÉCUTIONS PROGRAMMATIONS

Fichiers + ⚡ C 🔍 gbegrace_14) X +

MES NOTEBOOKS PARTAGÉS AVEC MO

Région us-central1 (Iowa)

Filtre Filtrer par notebook

gbegrace2018 (23 août 2024, 14:40:34) gbegrace2018 (23 août 2024, 16:20:14) gbegrace2018 (23 août 2024, 21:47:14) gbegrace2018 (24 août 2024, 09:31:52) gbegrace2018 (24 août 2024, 09:32:21)

```
+ Code + Texte + Commandes
```

```
import pandas as pd
df_carbu = pd.read_csv("gs://engiedatastorage/prix-carburants-quotidien@opendatamef.csv",sep=";")
print(df_carbu.head())
df_carbu.columns
```

	id	Code postal	pop	adresse	ville	horaires	rupture	fermeture	geom	Mise à jour des prix	... Nom Officiel EPCI
0	17250063	17250	R	Rue du gros chêne	Saint-Porchaire	{"@automate-24-24": "1", "jour": [{"@id": "1", ...}}	("gid": "4", "@nom": "GPLc", "@debut": "2020-0...	NaN	45.814, -0.774	2024-08-06T16:44:05+02:00	... CC Cœur de Saintonge
1	17290091	17290	R	37 BIS PLACE DE LA REPUBLIQUE	AIGREFEUILLE-D'AUNIS	{"@automate-24-24": "", "jour": [{"@id": "1", ...}}	("gid": "4", "@nom": "GPLc", "@debut": "2018-0...	NaN	46.117, -0.934	2024-07-23T11:41:52+02:00	... CC Aunis Sud
2	17290091	17290	R	37 BIS PLACE DE LA REPUBLIQUE	AIGREFEUILLE-D'AUNIS	{"@automate-24-24": "", "jour": [{"@id": "1", ...}}	("gid": "6", "@nom": "SP98", "@debut": "2018-0...	NaN	46.117, -0.934	2024-07-23T11:41:52+02:00	... CC Aunis Sud
3	17300003	17300	R	Impasse du 11 Novembre 1918	ROCHEFORT	{"@automate-24-24": "1", "jour": [{"@id": "1", ...}}	("gid": "2", "@nom": "SP95", "@debut": "2017-0...	NaN	45.929, -0.958	2024-08-05T02:01:00+02:00	... CA Rochefort Océan
4	17310004	17310	R	2 Route des Mirouelles	SAINT-PIERRE-D'OLÉRON	{"@automate-24-24": "1", "jour": [{"@id": "1", ...}}	("gid": "4", "@nom": "GPLc", "@debut": "2017-0...	NaN	45.943, -1.321	2024-08-07T10:59:10+02:00	... CC de l'Île d'Oléron

	Numéro	Département	Département	Code Officiel	Région
0		17	Charente-Maritime	75.0	
1		17	Charente-Maritime	75.0	
2		17	Charente-Maritime	75.0	
3		17	Charente-Maritime	75.0	
4		17	Charente-Maritime	75.0	

PIPELINES

```
data-collection-kedro > src > data_collection_kedro > pipelines > etl_pipeline > pipeline.py > create_pipeline
1 """
2     pipeline 'etl_pipeline'
3 """
4 from kedro.pipeline import Pipeline, node
5
6 from .nodes import etl_api_data, etl_csv_data
7
8 def create_pipeline(**kwargs) -> Pipeline:
9 """
10     Create a Kedro pipeline for the ETL process with multiple API sources.
11     Returns:
12         Pipeline: The constructed Kedro pipeline.
13 """
14
15     return Pipeline(
16         [
17             node(
18                 func=etl_api_data,
19                 inputs=["params:api_urls", "params:db_name", "params:collection_names"],
20                 outputs=None,
21                 name="process_api_data_node",
22             ),
23             node(
24                 func=etl_csv_data,
25                 inputs=[
26                     "params:csv_file_paths",
27                     "params:db_name",
28                     "params:collection_names",
29                 ],
30                 outputs=None,
31                 name="process_csv_data_node",
32             ),
33         ]
34     )
```

```
data-collection-kedro > src > data_collection_kedro > pipelines > data_fusion_pipeline > pipeline.py > ...
1 """Data fusion pipeline module."""
2 from kedro.pipeline import Pipeline, node
3
4 from .nodes import (
5     display_dataframes,
6     load_collections,
7     merge_data_store_in_elastic,
8     normalize_columns,
9     select_columns,
10 )
11
12 def create_pipeline(**kwargs) -> Pipeline:
13     return Pipeline(
14         [
15             node(
16                 func=load_collections,
17                 inputs=dict(
18                     collection_names="params:collection_names",
19                     db_name="params:db_name1",
20                 ),
21                 outputs="loaded_dataframes",
22                 name="load_collections_node",
23             ),
24             node(
25                 func=select_columns,
26                 inputs=dict(
27                     dataframes="loaded_dataframes",
28                     columns_to_select="params:columns_to_select",
29                 ),
30                 outputs="selected_dataframes",
31                 name="select_columns_node",
32             ),
33             node(
34                 func=display_dataframes,
35                 inputs="selected_dataframes",
36                 outputs="displayed_selected_data",
37                 name="display_selected_data_node",
38             ),
39             node(
40                 func=normalize_columns,
41                 inputs="selected_dataframes",
42                 outputs="normalized_dataframes",
43                 name="normalize_columns_node",
44             ),
45             node(
46                 func=merge_data_store_in_elastic,
47                 inputs=["normalized_dataframes"],
48                 outputs=["merged_meteo_courbe", "merged_courbe_mouvement"],
49                 name="merge_data_node",
50             ),
51         ]
52     )
```

ETL PIPELINE

```

PROBLEMS OUTPUT DEBUG CONSOLE TERMINAL PORTS COMMENTS
side. If you want to opt out, set the 'KEDRO_DISABLE_TELEMETRY' or 'DO_NOT_TRACK' environment variables, or create a '.travis.yml' file in the current working directory with the contents 'consent: false'. Read more at https://docs.kedro.org/en/stable/configuration/telemetry.html
(base) grace@grace-HP-ENVY-x360-Convertible:~/Projects/training_CDI/detectionAnomalie/data-collection-kedros$ kedro run --pipeline=etl_pipeline
[09/23/24 21:53:19] INFO Kedro project data-collection-kedros
[09/23/24 21:53:21] INFO Kedro is sending anonymous usage data with the sole purpose of improving the product. No personal data or IP addresses are sent. If you want to opt out, set the 'KEDRO_DISABLE_TELEMETRY' or 'DO NOT TRACK' environment variables, or create a '.travis.yml' file in the current working directory with the contents 'consent: false'. Read more at https://docs.kedro.org/en/stable/configuration/telemetry.html
[09/23/24 21:53:22] INFO Using synchronous mode for loading and saving data. Use the --async flag for potential performance gains.
https://docs.kedro.org/en/stable/nodes_and_pipelines/run_a_pipeline.html#load-and-save-asynchronously
    INFO Loading data from params:api.urls (MemoryDataset)...
    INFO Loading data from params:db_name (MemoryDataset)...
    INFO Loading data from params:collection_names (MemoryDataset)...
    INFO Running node: process_api_data node: process_api_data([params:api.urls;params:db_name;params:collection_names]) -> None
        date operateur secteur_d_activite region 00_00_00 ... 22_00_00 23_00 conso_journaliere_mwh_pcs_0degc
0 2022-05-29 GRTgaz GRS/ELD Grand Est 915.000000 ... 1064.328188 915.448799 27238.120793 def
1 2022-05-29 GRTgaz GRS/ELD Bretagne 387.198921 ... 424.963041 400.326384 13436.980963 def
2 2022-05-29 GRTgaz GRS/ELD Bourgogne-Franche-Comté 483.575449 ... 547.148815 509.835779 15277.687807 def
3 2022-05-29 GRTgaz GRS/ELD Pays de la Loire 686.634681 ... 784.503555 731.180790 23482.329316 def
4 2022-05-31 GRTgaz GRS/ELD Grand Est 1153.074509 ... 1155.682929 1262.429416 35586.643889 def
[5 rows x 31 columns]
Index(['date', 'operateur', 'secteur_d_activite', 'region', '00_00_00',
       '01_00', '02_00', '03_00', '04_00', '05_00_00', '06_00_00', '07_00',
       '08_00', '09_00', '10_00_00', '11_00_00', '12_00_00', '13_00', '14_00',
       '15_00', '16_00', '17_00', '18_00', '19_00_00', '20_00', '21_00',
       '22_00', '23_00', 'conso_journaliere_mwh_pcs_0degc', 'statut',
       'code_region'],
      dtype='object')
    periode source elec_tt gaz_tt propane fed_ct1 sp95 ... baseload_elec_2026 baseload_elec_2027 prix_tonne_co2 gazole_100km sp98_100km
0 2021-07-28 EDF/Powernet None None NaN NaN None ... None None None None None
1 2021-08-04 EDF/Powernet None None NaN NaN None ... None None None None None
2 2021-08-18 EDF/Powernet None None NaN NaN None ... None None None None None
3 2021-09-08 EDF/Powernet None None NaN NaN None ... None None None None None
4 1983-05-01 Edis None None 4.2323 3.9841 None ... None None None 2.632 None
[5 rows x 29 columns]
Index(['periode', 'source', 'elec_tt', 'gaz_tt', 'propane', 'fed_ct1', 'sp95', 'baseload_elec_2026', 'baseload_elec_2027', 'prix_tonne_co2', 'gazole_100km', 'sp98_100km'],
      dtype='object')
INFO Completed 1 out of 2 tasks
INFO Loading data from params:csv_file_paths (MemoryDataset)...
INFO Loading data from params:db_name (MemoryDataset)...
INFO Loading data from params:collection_names (MemoryDataset)...
INFO Running node: process_csv_data node: process_csv_data([params:csv_file_paths;params:db_name;params:collection_names]) -> None
Processing file data/prix_carburant_data_2018_2024.csv for collection csv.prix_carburant_data_2018_2024
    id code_postal_pcp ... adresse ... carburant_en_rupture ... début_rupture fin_rupture automate 24-24_(oui/non)
[09/12/24 22:19:24] INFO Completed 2 out of 2 tasks
    INFO Pipeline execution completed successfully.
(kedro-dataeng-env) (base) grace@grace-HP-ENVY-x360-Convertible:~/Projects/training_CDI/Projet_Master2/DetectionAnomalie/data-collection-kedros$ 

```

```

PROBLEMS OUTPUT DEBUG CONSOLE TERMINAL PORTS COMMENTS
Batch 42 with 500 records stored successfully in test_csv_courbe_de_charge_data_2018_2024.
Batch 43 with 500 records stored successfully in test_csv_courbe_de_charge_data_2018_2024.
Batch 44 with 500 records stored successfully in test_csv_courbe_de_charge_data_2018_2024.
Batch 45 with 500 records stored successfully in test_csv_courbe_de_charge_data_2018_2024.
[codecarbon INFO @ 22:18:39] Energy consumed for RAM : 0.000325 kWh. RAM Power : 4.333893299102783 W
[codecarbon INFO @ 22:18:39] Energy consumed for all CPUs : 0.000563 kWh. Total CPU Power : 7.5 W
[codecarbon INFO @ 22:18:39] 0.000888 kWh of electricity used since the beginning.
Batch 46 with 500 records stored successfully in test_csv_courbe_de_charge_data_2018_2024.
Batch 47 with 500 records stored successfully in test_csv_courbe_de_charge_data_2018_2024.
Batch 48 with 500 records stored successfully in test_csv_courbe_de_charge_data_2018_2024.
Batch 49 with 500 records stored successfully in test_csv_courbe_de_charge_data_2018_2024.
Batch 50 with 500 records stored successfully in test_csv_courbe_de_charge_data_2018_2024.
Batch 51 with 500 records stored successfully in test_csv_courbe_de_charge_data_2018_2024.
Batch 52 with 500 records stored successfully in test_csv_courbe_de_charge_data_2018_2024.
Batch 53 with 500 records stored successfully in test_csv_courbe_de_charge_data_2018_2024.
[codecarbon INFO @ 22:18:54] Energy consumed for RAM : 0.000343 kWh. RAM Power : 4.333893299102783 W
[codecarbon INFO @ 22:18:54] Energy consumed for all CPUs : 0.000594 kWh. Total CPU Power : 7.5 W
[codecarbon INFO @ 22:18:54] 0.000937 kWh of electricity used since the beginning.
Batch 54 with 500 records stored successfully in test_csv_courbe_de_charge_data_2018_2024.
Batch 55 with 500 records stored successfully in test_csv_courbe_de_charge_data_2018_2024.
Batch 56 with 500 records stored successfully in test_csv_courbe_de_charge_data_2018_2024.
Batch 57 with 500 records stored successfully in test_csv_courbe_de_charge_data_2018_2024.
Batch 58 with 500 records stored successfully in test_csv_courbe_de_charge_data_2018_2024.
Batch 59 with 500 records stored successfully in test_csv_courbe_de_charge_data_2018_2024.
Batch 60 with 500 records stored successfully in test_csv_courbe_de_charge_data_2018_2024.
[codecarbon INFO @ 22:19:09] Energy consumed for RAM : 0.000361 kWh. RAM Power : 4.333893299102783 W
[codecarbon INFO @ 22:19:09] Energy consumed for all CPUs : 0.000625 kWh. Total CPU Power : 7.5 W
[codecarbon INFO @ 22:19:09] 0.000986 kWh of electricity used since the beginning.
Batch 61 with 500 records stored successfully in test_csv_courbe_de_charge_data_2018_2024.
Batch 62 with 500 records stored successfully in test_csv_courbe_de_charge_data_2018_2024.
Batch 63 with 500 records stored successfully in test_csv_courbe_de_charge_data_2018_2024.
Batch 64 with 500 records stored successfully in test_csv_courbe_de_charge_data_2018_2024.
Batch 65 with 500 records stored successfully in test_csv_courbe_de_charge_data_2018_2024.
Batch 66 with 500 records stored successfully in test_csv_courbe_de_charge_data_2018_2024.
Batch 67 with 500 records stored successfully in test_csv_courbe_de_charge_data_2018_2024.
Batch 68 with 220 records stored successfully in test_csv_courbe_de_charge_data_2018_2024.
Processing file data/movement_data_2018_2024.csv for collection test_csv_mouvements_sociaux_data_2018_2024
    date_debut date_de_fin motif organisations_syndicales metiers_cibles population nombre_grevistes taux_grevistes date_debut
0 2018-04-03 NaN Reforme de la SNCF CGT,UNSA,SUD Rail,CFDT,FO NaN NaN NaN 34.0 2018-04-03
1 2018-04-13 NaN Reforme de la SNCF CGT,UNSA,SUD Rail,CFDT,FO NaN NaN NaN 22.5 2018-04-13
2 2018-04-19 NaN Reforme de la SNCF CGT,UNSA,SUD Rail,CFDT,FO NaN NaN NaN 22.7 2018-04-19
3 2018-06-07 NaN Reforme de la SNCF CGT,UNSA,SUD Rail,CFDT,FO NaN NaN NaN 14.3 2018-06-07
4 2018-06-23 NaN Reforme de la SNCF CGT,UNSA,SUD Rail,CFDT,FO NaN NaN NaN 20.7 2018-06-23
Index(['date_debut', 'date_de_fin', 'motif', 'organisations_syndicales',
       'metiers_cibles', 'population', 'nombre_grevistes', 'taux_grevistes',
       'date_debut'],
      dtype='object')
Batch 1 with 158 records stored successfully in test_csv_mouvements_sociaux_data_2018_2024.
[codecarbon INFO @ 22:19:24] Energy consumed for RAM : 0.000379 kWh. RAM Power : 4.333893299102783 W
[codecarbon INFO @ 22:19:24] Energy consumed for all CPUs : 0.000656 kWh. Total CPU Power : 7.5 W
[codecarbon INFO @ 22:19:24] 0.001035 kWh of electricity used since the beginning.
[09/12/24 22:19:24] INFO Completed 2 out of 2 tasks
    INFO Pipeline execution completed successfully.
(kedro-dataeng-env) (base) grace@grace-HP-ENVY-x360-Convertible:~/Projects/training_CDI/Projet_Master2/DetectionAnomalie/data-collection-kedros$ 

```

MONGODB

Energy

Overview Real Time Metrics **Collections** Atlas Search Performance Advisor Online Archive

DATABASES: 1 COLLECTIONS: 16

[+ Create Database](#)[Search Namespaces](#)

DBEnergy

[Consommation quotid...](#)[Consommation-quotidi...](#)[Données-de-consomm...](#)[Niveaux-de-prix-par-c...](#)[consommation_horaire...](#)[csv_combined_meteo_...](#)[csv_courbe_de_charg...](#)[csv_meteo_data_2018...](#)[csv_mouvements_soci...](#)[csv_prix_du_carburant...](#)[donnée-local-energie](#)[ecomix](#)[energiedata](#)[prix-des-carburants](#)[prod-region-annuelle...](#)

DBEnergy.Consommation quotidienne brute régionale

STORAGE SIZE: 224KB LOGICAL DATA SIZE: 1.37MB TOTAL DOCUMENTS: 3720 INDEXES TOTAL SIZE: 116KB

[Find](#)[Indexes](#)[Schema Anti-Patterns](#)[Aggregation](#)[Search Indexes](#)[Generate queries from natural language in Compass](#)[Filter](#)

Type a query: { field: 'value' }

QUERY RESULTS: 1-20 OF MANY

```
_id: ObjectId('667fc8de65b11852abad0128')
date_heure : "2021-06-28T11:30:00+00:00"
date : "2021-06-28"
heure : "13:30"
code_insee_region : "52"
region : "Pays de la Loire"
consommation_brute_gaz_grtgaz : null
statut_grtgaz : null
consommation_brute_gaz_terega : null
statut_terega : null
consommation_brute_gaz_totale : null
consommation_brute_electricite_rte : 2958
statut_rte : "Consolidé"
consommation_brute_totale : null
```

```
_id: ObjectId('667fc8de65b11852abad0129')
date_heure : "2021-06-28T11:30:00+00:00"
date : "2021-06-28"
heure : "13:30"
```

DATA FUSION PIPELINE

```
PROBLEMS OUTPUT TERMINAL PORTS COMMENTS
[09/16/24 12:53:06] INFO Saving data to loaded_dataframes (MemoryDataset)...
0.091263 INFO Completed 1 out of 3 tasks
[09/16/24 12:53:09] INFO Loading data from loaded_dataframes (MemoryDataset)...
0.091269 INFO Loading data from params:columns_to_select (MemoryDataset)...
0.091339 INFO Running node: select_columns_node: select_columns((loaded_dataframes,params:columns_to_select)) -> (selected_dataframes)
0.091364
Colonnes disponibles dans le DataFrame: Index(['_id', 'Unnamed: 0', 'TempMax_Deg', 'TempMin_Deg', 'Wind_kmh',
   'Wet_percent', 'Visibility_km', 'CloudCoverage_percent',
   'Dayduration_hour', 'region', 'day'],
  dtype='object')
Colonnes disponibles dans le DataFrame: Index(['_id', 'date', 'opérateur', 'secteur_d'activité', 'région', '00:00',
   '01:00', '02:00', '03:00', '04:00', '05:00', '06:00', '07:00', '08:00',
   '09:00', '10:00', '11:00', '12:00', '13:00', '14:00', '15:00', '16:00',
   '17:00', '18:00', '19:00', '20:00', '21:00', '22:00', '23:00',
   'consommation_journalière_(mwh_pc_0°C)', 'statut', 'code_région'],
  dtype='object')
Colonnes disponibles dans le DataFrame: Index(['_id', 'date_de_debut', 'date_de_fin', 'motif',
   'organisations_syndicales', 'metiers_cibles', 'population',
   'nombre_grevistes', 'taux_grevistes', 'date_debut'],
  dtype='object')
[0.091363] INFO Saving data to selected_dataframes (MemoryDataset)...
0.091363 INFO Completed 2 out of 3 tasks
[09/16/24 12:53:23] INFO Loading data from selected_dataframes (MemoryDataset)...
0.091369 INFO Running node: display_selected_data_node: display_dataframe(selected_dataframes) -> (displayed_selected_data)
0.091364
DataFrame 1:
  TempMax_Deg TempMin_Deg Wind_kmh Wet_percent Visibility_km CloudCoverage_percent Dayduration_hour region day
0      9.0       5.0      25      74     0.250          20.0          8:23:0  alsace 2018/01/01
1     11.0       7.0      47      78     0.750          42.0          8:12:0  bretagne 2018/01/01
2      8.0       6.0      38      83     0.500          91.0          8:18:0  lorraine 2018/01/01
3     14.0      11.0      23      95     7.825          96.0          8:53:0  aquitaine 2018/01/02
4     12.0       5.0      39      93     0.000          78.0          8:29:0    centre 2018/01/02

DataFrame 2:
  date opérateur secteur_d'activité régions 00:00 ... 21:00 22:00 23:00 consommation_journalière_(mwh_pc_0°C) statut
0 2018-11-15 6Ktgaz GDF/EDF Hauts-de-France 6880.895246 ... 6378.954872 6741.885649 6939.526476 148113.988853 04finale
```

Total carbon emissions for this run: 1.343783220306978e-05 kgCO₂eq

```
...ions for this run. 1.543/032202009/00-00 Hydrology
INFO Saving data to merged_meteo_courbe (MemoryDataset)...
INFO Saving data to merged_courbe_mouvement (MemoryDataset)
INFO Completed 5 out of 5 tasks
INFO Pipeline execution completed successfully.
INFO Loading data from merged_meteo_courbe (MemoryDataset)
INFO Loading data from merged_courbe_mouvement (MemoryDataset)
INFO Loading data from dissolved_selected_data (MemoryDataset)
```

DATA FUSION PIPELINE

ELASTIC SEARCH

elastic

Find apps, content, and more.

Setup guides

Stack Management Index Management Indices

Management

Ingest (3)

- Ingest Pipelines
- Logstash Pipelines

Data (10)

- Index Management**
- Index Lifecycle Policies
- Data Set Quality
- Snapshot and Restore
- Rollup Jobs
- Transforms
- Cross-Cluster Replication
- Remote Clusters

Alerts and Insights (5)

- Alerts
- Rules
- Cases
- Connectors
- Reporting

Machine Learning

Watcher

Maintenance Windows

Security (3)

- Users
- Roles
- API keys

Console

Notebooks

Index Management docs

Index Management

Indices Data Streams Index Templates Component Templates Enrich Policies

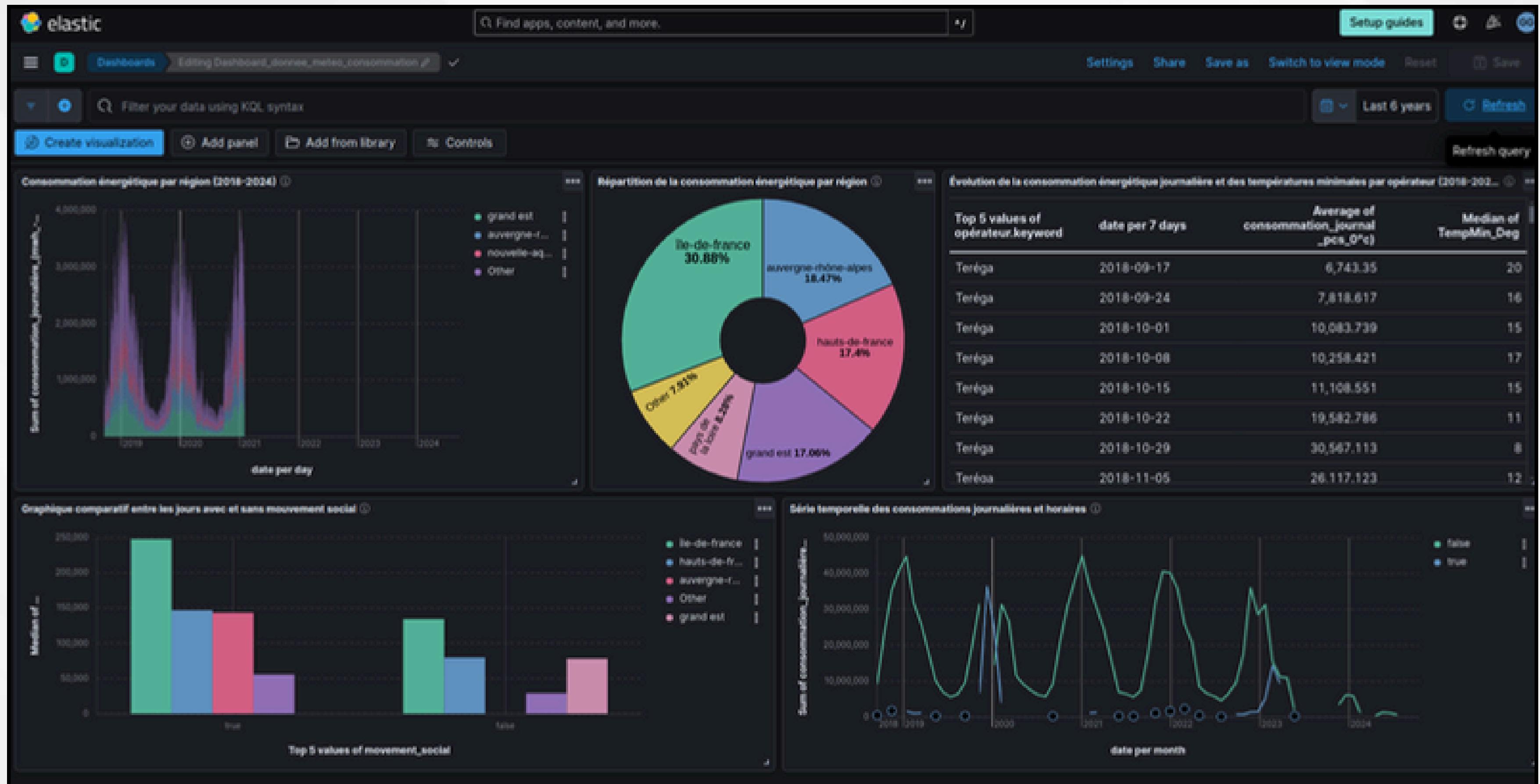
Update your Elasticsearch indices individually or in bulk. [Learn more.](#) ⓘ

Include hidden indices Include rollup indices

Name	Health	Status	Primaries	Replicas	Docs count	Storage size	Data stream
course_mouvement_index	green	open	1	1	22 318	25.05mb	
logs_englewatch	green	open	1	1	191	229.87kb	
meteo_course_index	green	open	1	1	29 986	33.17mb	
metrics-endpoint.metadata.current_default	green	open	1	1		4996	

Rows per page: 10 ⏪ ⏴ ⏵ ⏹

KIBANA



RESULTATS TEST UNITAIRES

```
===== test session starts =====
platform linux -- Python 3.11.5, pytest-8.3.2, pluggy-1.5.0
rootdir: /home/grace/Projects_training_CDI/Projet_Master2/DetectionAnomalie/data-collection-kedro
configfile: pyproject.toml
plugins: cov-5.0.0
collected 14 items

tests/pipelines/data_fusion_pipeline/test_pipeline.py ...
tests/pipelines/etl_pipeline/test_pipeline.py ..... [ 21%]
tests/pipelines/etl_pipeline/test_transform.py ... [ 78%]
[100%]

===== warnings summary =====
tests/pipelines/data_fusion_pipeline/test_pipeline.py::test_store_in_elasticsearch
  /home/grace/Projects training_CDI/DetectionAnomalie/data-collection-kedro-dataeng-env/lib/python3.11/site-packages/numpy/_core/fromnumeric.py:57: FutureWarning:
  'DataFrame.swapaxes' is deprecated and will be removed in a future version. Please use 'DataFrame.transpose' instead.
    return bound(*args, **kwds)

-- Docs: https://docs.pytest.org/en/stable/how-to/capture-warnings.html
===== 14 passed, 1 warning in 9.40s =====
```

8.

MODÈLE DE MACHINE LEARNING

NOTRE SELECTION

Modèle	Algorithm	Objectif
Détection des Anomalies	Isolation Forest	Identifier les anomalies dans les habitudes de consommation
Clustering par Région	DBSCAN & KMeans	Regrouper les régions en fonction de profils de consommation
Prédiction de la Consommation avec Données de Mouvements	Random Forest	Prédire la consommation en fonction des mouvements humains
Prédiction de la Consommation avec Données Météorologiques	CatBoost	Prédire la consommation en fonction de la météo

DÉTECTION D'ANOMALIES

💬 Détection d'anomalies dans les données

Choisissez un fichier CSV

Drag and drop file here
Limit 200MB per file - CSV

[Browse files](#)

 df_with_anomalies.csv 4.3MB

X

Nombre d'anomalies : 320

Nombre de lignes sans anomalies : 10920

	date	région	consommation_moyenne_journalière	statut	anomaly
67	2018-12-26 00:00:00	bourgogne-franche-comté	3739.191501	Définitif	-1
68	2018-12-29 00:00:00	bourgogne-franche-comté	3488.272348	Définitif	-1
69	2018-12-30 00:00:00	Île-de-france	11126.150450	Définitif	-1
70	2019-01-02 00:00:00	Île-de-france	30958.687564	Définitif	-1
71	2019-01-03 00:00:00	nouvelle-aquitaine	2829.525728	Définitif	-1
72	2019-01-05 00:00:00	occitanie	1214.573892	Définitif	-1
73	2019-01-06 00:00:00	provence-alpes-côte d'azur	3505.725034	Définitif	-1
74	2019-01-07 00:00:00	Île-de-france	32547.795176	Définitif	-1
75	2019-01-07 00:00:00	hauts-de-france	6963.340698	Définitif	-1
76	2019-01-07 00:00:00	provence-alpes-côte d'azur	3739.080995	Définitif	-1
77	DATE AT naanumna	new caledonia, new caledonian	7165.766968	Potentiel	-1

Tableau des anomalies détectées

CYCLE DE VIE DES MODÈLES

Non sécurisé 14.70.175.234:5000/#/experiments/6?searchFilter=&orderByKey=attributes.start_time&orderByAsc=false&startTIme=A&utilLifecycleFilter=Active&modelVersionFilter=All+Runs&datasetFilter=W10%10

mlflow 2.18.0 Experiments Models GitHub Docs

Experiments

Energy Consumption Prediction - METEO

Provide Feedback Add Description Share

Run Evaluation Experimental Traces Experimental

metrics.rmse < 1 and params.model = "tree"

Time created State: Active Datasets Sort: Created Columns + New run

Group by

Run Name	Created	Dataset	Duration	Source	Models
CatBoost Regressor (Tu...)	4 days ago	dataset (eadbc3c9) Train	45.1min	catboost...	-
CatBoost Regressor (Tu...)	5 days ago	dataset (eadbc3c9) Train	37.7min	catboost...	-
Random Forest Regressor	5 days ago	dataset (eadbc3c9) Train	24.8s	rm.py	-
Random Forest Regressor	5 days ago	dataset (eadbc3c9) Train	21.8s	rm.py	-
CatBoost Regressor (Tu...)	5 days ago	dataset (eadbc3c9) Train	37.6min	catboost...	-
CatBoost Regressor (Tu...)	5 days ago	dataset (eadbc3c9) Train	2.3min	catboost...	-
CatBoost Regressor	5 days ago	dataset (eadbc3c9) Train	1.1min	nodes.py	-
LightGBM Regressor	5 days ago	dataset (eadbc3c9) Train	1.8min	nodes.py	-
XGBoost Regressor	5 days ago	dataset (eadbc3c9) Train	34.5s	nodes.py	-
GridSearch Gradient Bo... Gradient Boosting Regr...	5 days ago	dataset (eadbc3c9) Train	34.9min	nodes.py	-
Gradient Boosting Regr...	17 days ago	dataset (eadbc3c9) Train	17.8s	nodes.py	-
GridSearch Gradient Bo... Gradient Boosting Regr...	17 days ago	dataset (eadbc3c9) Train	28.1min	nodes.py	-
Gradient Boosting Regr...	17 days ago	dataset (eadbc3c9) Train	15.5s	nodes.py	-
Linear Regression	17 days ago	dataset (eadbc3c9) Train	13.7s	nodes.py	-
Random Forest Regressor	17 days ago	dataset (eadbc3c9) Train	23.4s	nodes.py	-
Logistic Regression	17 days ago	dataset (5f64ee57) Train	0.7s	nodes.py	-
Random Forest Regressor	17 days ago	dataset (5f64ee57) Train	23.7s	nodes.py	-
rimble-parch-M43	17 days ago	-	1.4s	nodes.py	-
Linear Regression	17 days ago	dataset (5f64ee57) Train	13.2s	nodes.py	-

29 matching runs Show more columns (85 total)

STOCKAGE ARTEFACTS MLFLOW

Google Cloud EngieWatchProject Tapez / pour rechercher des ressources, des documents, des produits, etc Recherche

Cloud Storage Cloud Storage Informations sur le bucket

Zone : eu (plusieurs régions dans l'Union européenne) Classe de stockage : Standard Accès public : Non public Protection : Supprimer de façon réversible

ACCÉDER AU CHEMIN ACTUALISER DÉMARRER

Buckets

OBJET CONFIGURATION AUTORISATIONS PROTECTION CYCLE DE VIE OBSERVABILITÉ RAPPORTS SUR L'INVENTAIRE OPÉRATIONS

Navigateur de dossiers Buckets > engiedatastorage

Créer un dossier Importer Transférer les données Autres services

Filtrer par préfixe de nom uniquement Filtrer les objets et dossiers Afficher Objets actifs uniquement

Objet	Taille	Type	Création	Classe de stockage	Dernière modification	Accès public	Historique des
couche-de-chage-mlflow-regional...	11,1 Mo	text/csv	23 août 2024, 13:01:31	Standard	23 août 2024, 13:01:31	Non public	
data/	—	Dossier	—	—	—	—	
mlflow_experiment/	—	Dossier	—	—	—	—	
models/	—	Dossier	—	—	—	—	
mouvements-vecteur-depuis-200...	105 Ko	text/csv	23 août 2024, 12:59:48	Standard	23 août 2024, 12:59:48	Non public	
prix-carburants quotidien@open...	74,9 Mo	text/csv	23 août 2024, 13:02:37	Standard	23 août 2024, 13:02:37	Non public	

Marktplace Notes de version

9.

EMPREINTE CARBONE

CALCUL EMPREINTE CARBONE

- **Globalement faible** : Empreinte carbone totale de 0.000593 kgCO2eq
- **Modèle le plus énergivore** : Modèle CatBoost.

```
scripts/empreinte_carbone_projet/emissions/emissions_catboosting.csv: Dernier run: 0.0004381058554062 kgCO2eq
scripts/empreinte_carbone_projet/emissions/emissions_data_fusion_pipeline.csv: Dernier run: 1.3437832202069784e-05 kgCO2eq
scripts/empreinte_carbone_projet/emissions/emissions_dbscan.csv: Dernier run: 4.407894442148034e-06 kgCO2eq
scripts/empreinte_carbone_projet/emissions/emissions_etl_pipeline.csv: Dernier run: 5.8006733084674365e-05 kgCO2eq
scripts/empreinte_carbone_projet/emissions/emissions_isolationforest.csv: Dernier run: 6.892161395051047e-06 kgCO2eq
scripts/empreinte_carbone_projet/emissions/emissions_random_forest.csv: Dernier run: 1.603075068442143e-05 kgCO2eq
scripts/empreinte_carbone_projet/emissions/emissions_ridge_model.csv: Dernier run: 2.3823826691643307e-06 kgCO2eq

Emprinte carbone totale du projet: 0.00059263609883729 kgCO2eq
```

10.

INTEGRATION CI

MAKEFILE

```
M Makefile M X
M Makefile
1 MYPY_PATHS = data-collection-kedro/src/data_collection_kedro/pipelines/etl_pipeline \
2 |   |   |   |   data-collection-kedro/src/data_collection_kedro/pipelines/data_fusion_pipeline \
3 |
4 PYLINT_PATHS = data-collection-kedro/src/data_collection_kedro/pipelines/etl_pipeline \
5 |   |   |   |   data-collection-kedro/src/data_collection_kedro/pipelines/data_fusion_pipeline \
6 |   |   |   |   dashboard_ui/
7
8 ISORT_PATHS = data-collection-kedro/src/data_collection_kedro/pipelines/etl_pipeline \
9 |   |   |   |   data-collection-kedro/src/data_collection_kedro/pipelines/data_fusion_pipeline \
10 |   |   |   |   dashboard_ui/
11
12 TEST_PATH = tests/
13
14 mypy:
15     mypy $(MYPY_PATHS)
16 pylint:
17     pylint $(PYLINT_PATHS)
18 isort:
19     isort $(ISORT_PATHS)
20 black:
21     black $(PYLINT_PATHS)
22 test:
23     pytest $(TEST_PATH) --disable-warnings
24
25 # Commande pour exécuter tous les outils (mypy, pylint, isort, black, pytest)
26 check: mypy pylint isort black test
27
28 # Nettoyage des fichiers temporaires
29 clean:
30     rm -rf __pycache__
31     rm -rf .mypy_cache
32     rm -rf .pytest_cache
33
```

GITHUB ACTION

```
page_anomalie_detection.py          ci.yml  M X  
.github > workflows > ci.yml  
1  name: Python Lint, Type Check, and Tests  
2  on: [push, pull_request]  
3  jobs:  
4    build:  
5      runs-on: ubuntu-latest  
6      steps:  
7        - uses: actions/checkout@v2  
8        - name: Set up Python  
9          uses: actions/setup-python@v2  
10         with:  
11           python-version: '3.11.5'  
12        - name: Install dependencies  
13          run:  
14            python -m pip install --upgrade pip  
15            pip install -r requirements_test.txt  
16            pip install -r data-collection-kedro/requirements.txt  
17            pip install boto3 botocore  
18            python3 -m pip install types-requests  
19            pip install pytest pytest-cov # Ajout de pytest et pytest-cov  
20            pip install pytest-mock  
21        - name: Check import formatting with isort  
22          run: isort .  
23        - name: Check code formatting with black  
24          run: black .  
25        - name: Lint with pylint  
26          run: pylint --rcfile=.pylintrc .  
27        - name: Run tests with pytest  
28          run:  
29            pytest --disable-warnings --maxfail=5  
30        - name: Generate full coverage report  
31          run:  
32            pytest --cov=. data-collection-kedro/tests/ --cov-report=xml  
33            pytest --cov=. data-collection-kedro/tests/ --cov-report=html  
34        - name: Upload coverage report artifact  
35          uses: actions/upload-artifact@v3  
36          with:  
37            name: coverage-report  
38            path: htmlcov/  
39
```

11.

INTERFACE UTILISATEUR



Prédiction d'énergie consommée

Dans cette section, nous faisons la prédiction en prenant en compte des conditions météorologiques.

- Prédiction Météo
- Anomalie Détection
- Prédiction Consommation
- Clustering
- Feedback
- Tracking

Température Max (°C)	25,00	- +	Humidité (%)	50	- +
Température Min (°C)	15,00	- +	Visibilité (km)	10	- +
Vitesse du vent (km/h)	10	- +	Couverture Nuageuse (%)	50	0 100

Prédire la Consommation Énergétique

12.

DIFFICULTES RENCONTRÉES ET EVOLUTIONS

DIFFICULTES RENCONTREES

- Qualité du dataset
- Technologies diverses et variées à prendre en main
- Gestion des Données Massives
- Contraintes Budgétaires
- Données sur l'Empreinte Carbone



EVOLUTIONS



**Prédiction
des risques
d'incidents**



**Anticipation des
Impacts**

**Optimisation
de la
Planification**

13.

MAITENANCE DE L'APPLICATION

MONITORING

Screenshot of the Elastic Stack Monitoring interface showing various dashboards and metrics for different machine learning models.

Top Left Panel: Statistics des Modèles par Application

Top 5 values of application_name.keyword	Top 5 values of model_name.keyword	Median of response_time	Median of cpu_usage	Median of memory_usage	Unique count of model_version.keyword	Count of status.keyword
PredictionConsoleAppWithMetrics	CatBoost	-	13	62.5	1	6
PredictionConsoleApp	RandomForest	27.257	2	60	1	14
ClusteringApp	DBSCAN	1	14.5	60.5	1	6
AnomalyDetectionApp	IsolationForest	0.223	11	58	1	30

Top Right Panel: Répartition des Exécutions par Modèle

Modèle	Pourcentage
IsolationForest	43.97%
CatBoost	34.84%
DBSCAN	9.61%
RandomForest	12.38%

Middle Left Panel: Nombre d'Exécutions Réussies et Echouées par Modèle

Modèle	Nombre d'exécutions réussies	Nombre d'exécutions échouées
IsolationForest	25	0
CatBoost	18	0
RandomForest	12	0
DBSCAN	5	0

Middle Right Panel: Historique des Erreurs

Top 3 values of status.keyword	Top 3 values of error_message.keyword	Top 3 values of error_type.keyword	Top 3 values of model_name.keyword	Minimum of timestamp
failed	name 'psutil' is not defined	NameError	IsolationForest	Sep 17, 2024 @ 12:57:56.46
failed	name 'psutil' is not defined	NameError	CatBoost	Sep 17, 2024 @ 12:55:58.44

Bottom Panel: Top 4 values of model_name.keyword

Modèle	Nombre d'exécutions
IsolationForest	25
CatBoost	18
RandomForest	12
DBSCAN	5

14.

DEMO



MERCI
