

Projet d'évaluation de Data Science

Auteurs : Marc-Antoine Bernard et Abdoul Aziz Toure

Cursus : ESTIA, Master BIHAR 2025-26

Date : 31 décembre 2025

Table des matières

Projet d'évaluation de Data Science.....	1
1 INFORMATIONS GÉNÉRALES.....	3
2 PROBLÉMATIQUE CLIENT ET CAHIER DES CHARGES.....	4
2.1 Informations clients (échange avec Yvan Le Bihan du 9/12/25) :.....	4
2.2 Objectifs opérationnels.....	4
2.3 Données à disposition.....	4
3 CADRAGE DU PROJET.....	5
4 DESCRIPTION DÉTAILLÉE DES TRAVAUX RÉALISÉS (Marc-Antoine Bernard).....	7
4.1 Préparation des données.....	7
4.2 Étude des saisonnalités.....	9
4.3 Évaluation des modèles de prédiction.....	10
4.3.a Caractéristiques des données.....	10
4.3.b Analyse des modèles disponibles.....	10
4.3.c Déploiement des modèles.....	11
4.4 Baseline.....	12
4.5 Baseline ajustée de la tendance.....	13
4.6 Conclusion.....	13
4.7 Bibliographie.....	14
5 DESCRIPTION DÉTAILLÉE DES TRAVAUX RÉALISÉS (Abdoul Aziz Toure).....	15
5.1 Préparation des données.....	15
5.1.a Visualisation globale de la consommation.....	15
5.1.b Suppression des outliers extrêmes (99 ^e percentile).....	15
5.2 Mise en évidence de la saisonnalité annuelle.....	16
5.2.a Saisonnalité forte et structurée.....	16
5.3 Identification d'une période de référence « saine ».....	17
5.3.a Justification du découpage temporel.....	17
5.4 Cahier des charges fonctionnel et technique.....	18
5.4.a Contraintes fonctionnelles.....	18
5.4.b Contraintes mathématiques et data science.....	18
5.5 Positionnement méthodologique retenu.....	18
5.5.a Conclusion Partielle.....	19
5.6 Utilisation du modèle et détection des anomalies.....	19
5.6.a Objectif de la phase de modélisation.....	19
5.6.b Choix méthodologique : Holt-Winters global normalisé.....	20
5.7 Prétraitement des données.....	20
5.7.a Tri temporel et transformation logarithmique.....	20
5.7.b Calcul de la moyenne par compteur (train uniquement).....	20
5.7.c Jointure et centrage par ID.....	21

5.8 Apprentissage du modèle Holt-Winters.....	21
5.8.a Construction de la série temporelle d'apprentissage.....	21
5.8.b Entraînement du modèle Holt-Winters.....	21
5.9 Prédiction sur la période test.....	22
5.9.a Reprojection dans l'espace réel.....	22
5.9.b Interprétation.....	22
5.10 Détection robuste des anomalies.....	22
5.10.a Calcul des résidus.....	22
5.10.b Seuil robuste basé sur la MAD.....	22
5.10.c Détection des anomalies.....	23
5.11 Visualisation des résultats globaux.....	23
5.12 Analyse par compteur.....	24
5.12.a Extraction de caractéristiques par ID.....	24
5.12.b Visualisation multi-compteurs.....	24
5.13 Évaluation quantitative du modèle.....	25
5.13.a Résultats obtenus.....	25
5.13.b Interprétation.....	25
5.14 Conclusion partielle.....	26

1 INFORMATIONS GÉNÉRALES

Pour un opérateur public de distribution d'eau, la détection rapide des fuites sur le réseau est un enjeu stratégique majeur.

Les équipes terrain interviennent quotidiennement sur des milliers de kilomètres de canalisation, où les fuites peuvent passer inaperçues et entraîner des pertes importantes d'eau, des coûts élevés et des perturbations du service.

HUPI travaille depuis plusieurs mois afin de développer un Assistant Virtuel de détection de fuites, basé sur des modèles de Machine Learning.

Cet assistant prend en compte :

- les caractéristiques du réseau,
- les données historiques de consommation et de pression,
- ainsi que le profil et les habitudes de fonctionnement du réseau, afin d'alerter automatiquement et judicieusement lorsqu'un risque de fuite est détecté.

L'objectif est de concevoir des modèles auto-apprenants, capables de s'adapter aux spécificités de chaque zone géographique et de chaque type d'infrastructure, afin de générer des alertes personnalisées et adaptées à chaque contexte de réseau.

Nous disposons d'un ensemble de données contenant les mesures quotidiennes de consommation d'eau pour 502 compteurs différents. Chaque ligne correspond à une observation d'un compteur à une date donnée.

Les variables sont les suivantes :

- valeur_active : consommation mesurée (en m³)
- valeur_date : date de la mesure (quotidienne)
- libelle : identifiant du compteur (502 valeurs différentes)

L'objectif du projet est d'analyser et de caractériser la consommation d'eau des différents compteurs.

Pour répondre à cette problématique, trois tâches principales peuvent être menées :

- Prédire la consommation future des compteurs à partir de leurs historiques.
- Identifier la tendance de la consommation (hausse, stabilité, baisse).
- Classer les niveaux de consommation (faible, moyen, fort) pour caractériser les comportements des compteurs.

Après une analyse de l'attendu et une exploration des données en binôme, chaque membre s'est concentré sur un modèle en particulier.

2 PROBLÉMATIQUE CLIENT ET CAHIER DES CHARGES

2.1 Informations clients (échange avec Yvan Le Bihan du 9/12/25) :

Le client est capable d'identifier les fuites importantes (gros débit) car des résurgences apparaissent sur le terrain. A l'inverse, les fuites plus petites (petit débit) sont rarement visibles rapidement d'où la demande du client de pouvoir les identifier sur les relevés de compteur.

Pour localiser une fuite sur le terrain, le client utilise des équipements de détection (Par exemple à base d'ultra son) qui nécessitent de parcourir en surface le long de la canalisation. Ce temps de recherche est relativement long et incompressible, d'où l'attendu du client d'identifier rapidement une fuite sur les courbes (ordre de grandeur : quelques semaines).

Par ailleurs, si l'utilisateur est généralement en mesure d'identifier une fuite à l'échelle d'un compteur individuel, il ne dispose ni de la capacité ni des ressources nécessaires pour analyser exhaustivement l'ensemble des compteurs. De plus, les moyens d'intervention sur le terrain étant limités, l'outil de détection doit prioritairement fournir un indicateur de fuite dont la **précision augmente à mesure que le seuil de décision se resserre**, au détriment éventuel du rappel. Cette stratégie vise à concentrer l'analyse sur des situations présentant une probabilité élevée de fuite avérée, afin de renforcer la confiance de l'utilisateur dans l'outil et de permettre une mobilisation ciblée et efficace des ressources opérationnelles.

2.2 Objectifs opérationnels

L'outil doit donc :

- se focaliser sur les compteurs à débit relativement faible
- proposer un metric de probabilité de fuite
- l'identifier en quelques semaines.

2.3 Données à disposition

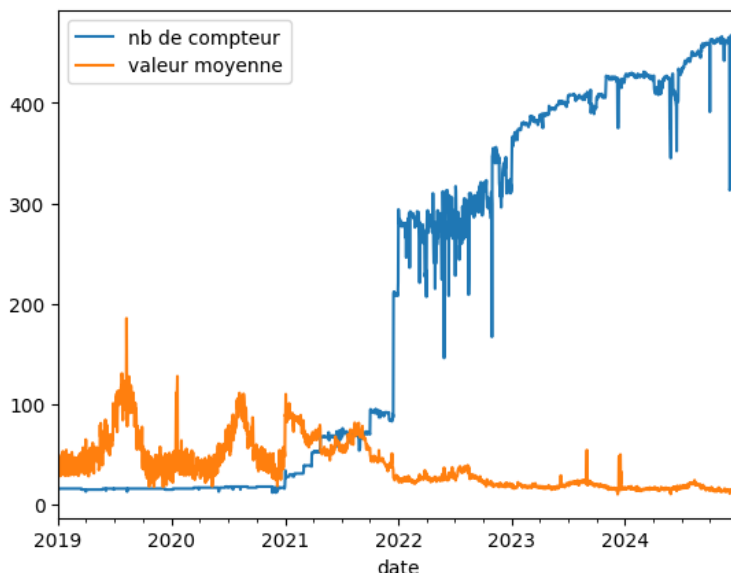
Les données fournies sont des télé relevés quotidiens de 502 compteurs d'eau sur une plage s'étalant de 2019 à 2024

D'après le client :

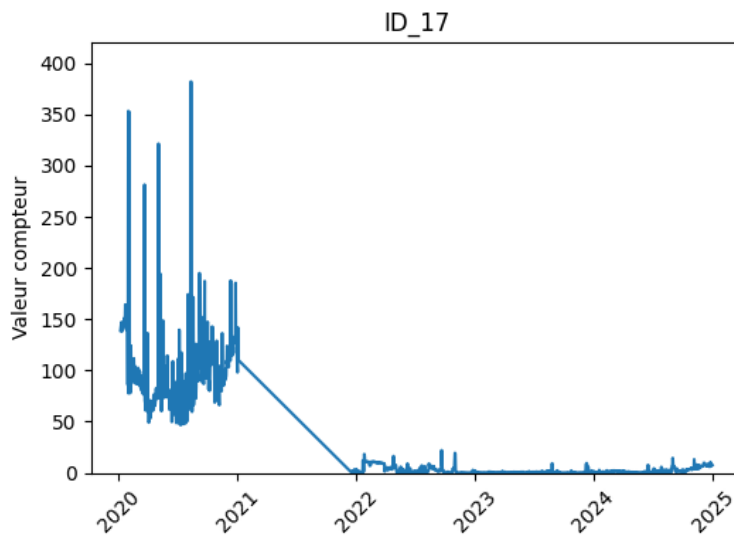
- les données de la plage du COVID (2020-2021) ne sont pas toujours représentatives.
- Les données sont brutes, sans pré-traitement. Des artefacts de mesures sont présents. Les données ne sont pas labellisées (présence de fuite non renseignée).
- la consommation d'eau dans le Sud-Ouest de la France est très liée aux saisons climatiques. En effet, la population augmente fortement l'été du fait de l'activité touristique.

3 CADRAGE DU PROJET

A l'origine du projet en 2021, Hupi avait reçu les données de 20 compteurs initiaux (ID_1 à ID_20) sur la plage 2019-2021. La plage de temps et le nombre de compteurs a été étendue dans un deuxième temps à partir de 2022, ce qui est illustré ci dessous.



Cela a conduit à parfois créer des plages non continues dans les données de certains compteurs. Ces plages sont aussi parfois non cohérentes, tel qu'illustré ci dessous pour le compteur ID_17



Afin de disposer d'un maximum de compteurs sur une période représentative, il a été convenu de focaliser notre étude sur la plage 2022-2024.

Les données sont des séries temporelles indépendantes, c'est à dire que les données sont à considérer individuellement par compteur et dans un ordre chronologique.

L'objectif est de déterminer un modèle capable d'identifier en quelques semaines ce qui se rapproche le plus d'une signature typique de fuite et de le soumettre au client pour vérification sur le terrain.

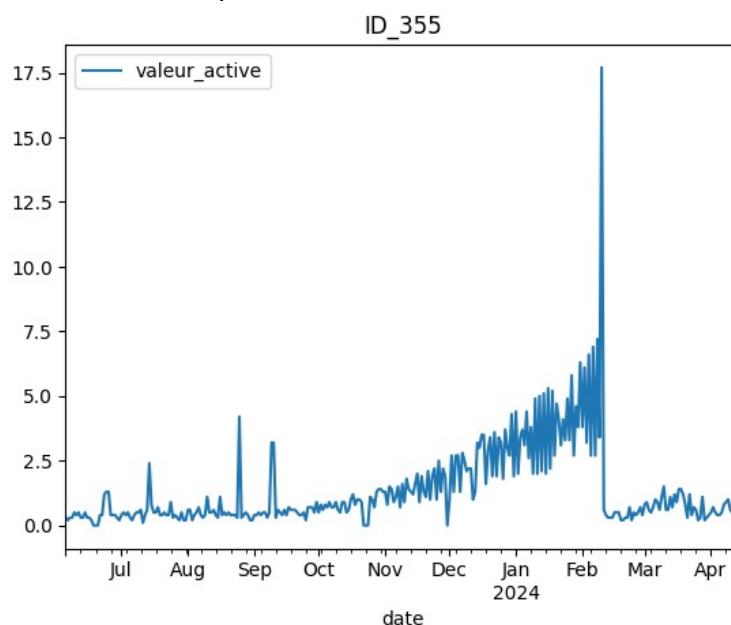
Par définition, les fuites recherchées n'étant pas connues, nous pouvons exclure ce qui relève d'évènements identifiés : conduite arrachée sur des travaux de VRD, dégâts des eaux chez le consommateur. Ces évènements sont généralement :

- soudains = montée rapide de la consommation
- violents = montée élevée de la consommation
- courts = retour à la normale rapidement après l'évènement.

Les fuites recherchées sont des tendances de consommation qui se caractérisent par :

- non saisonnier car inhabituel
- une consommation en augmentation progressive et éventuellement en accélération car la fuite ne peut qu'aller en s'aggravant sur un réseau d'adduction d'eau.
- sur un intervalle de temps de plusieurs mois car non détectée.

Un exemple typique est le cas du compteur ID_355.



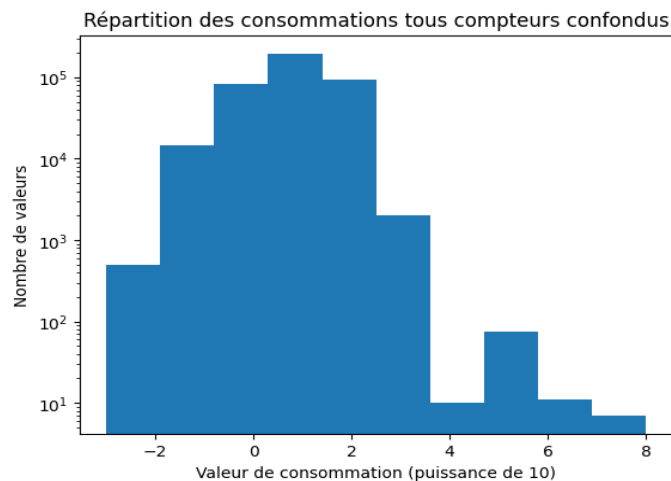
On peut imaginer que la fuite a commencé en octobre 2023 et a été réparée en février 2024.

La démarche consiste à prédire sur quelques semaines la consommation future des compteurs à partir de leurs historiques en tenant compte des saisonnalités et des tendances long terme, et de comparer cette prédiction avec la réalité pour qualifier et quantifier un résiduel.

4 DESCRIPTION DÉTAILLÉE DES TRAVAUX RÉALISÉS (Marc-Antoine Bernard)

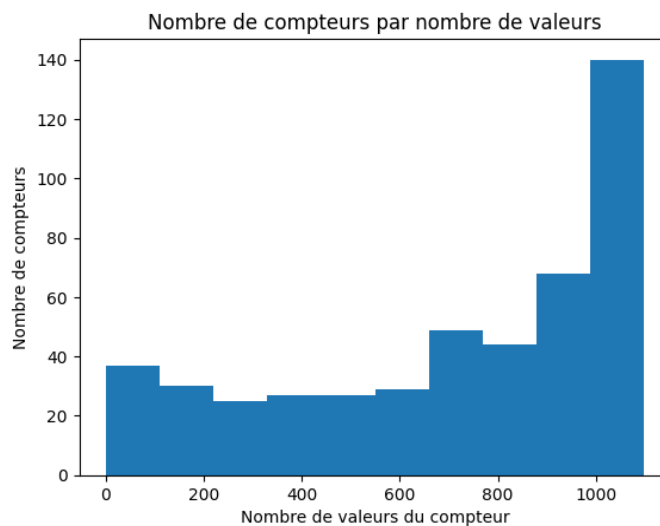
4.1 Préparation des données

Un nettoyage préalable a été menée suite à l'exploration des données :

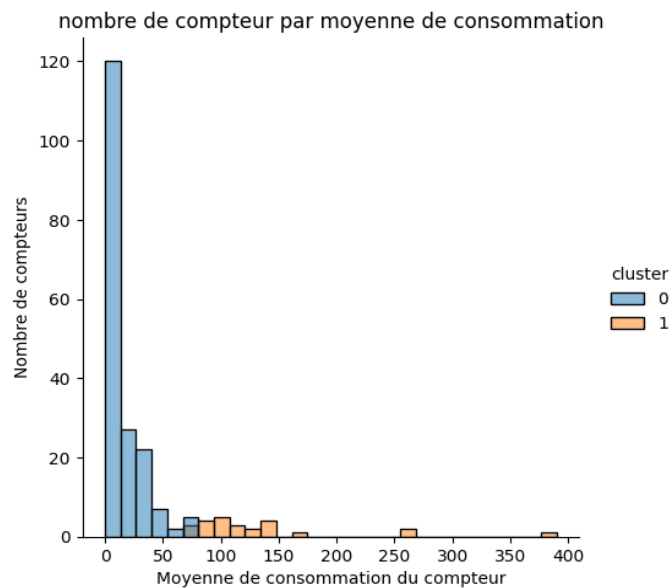


- Suppression des valeurs négatives ou globalement aberrantes ($>1e3$).
- Suppression par compteur des valeurs manquantes, nulles ou quasi-nulles (percentile 1%), et des pics (percentile 99%).

A la suite de ce nettoyage, seuls les compteurs avec au moins 80 % de données exploitables ont été conservés, soit 208 compteurs.



Une clusterisation sur la moyenne de consommation des compteurs a ensuite été menée afin de retirer les compteurs à gros débit et conserver un cluster homogène de 183 compteurs.



Enfin, un resampling a été réalisé afin d'obtenir des séries temporelles continues à pas quotidien aptes à être traitées par les algorithmes de prédictions :

- regroupement par moyenne des données d'un même jour
- ajout de donnée par continuité sur les jours manquants

4.2 Étude des saisonnalités

Le client a indiqué que les données présentaient une saisonnalité annuelle. Cette caractéristique peut impacter fortement le choix du modèle de prédiction et doit donc être quantifiée.

Pour cela, nous utilisons le Seasonal Variance Ratio (SVR) issu du SNR en traitement de signal.

Le Signal-to-Noise Ratio (SNR) analyse le bruit résiduel dans un signal :

$$SNR = \frac{Puissance_{signalUtile}}{Puissance_{noise}}$$

La puissance est la moyenne du carré d'un signal centré sur zéro, soit sa variance.

Dans un signal avec des composantes décorrélées, les variances s'additionnent :

$$Var_{signal} = Var_{season} + Var_{trend} + Var_{noise}$$

Le SNR est défini comme la part de puissance totale issue de la saisonnalité, soit :

$$SVR = \frac{Var_{saisonnalité}}{Var_{totale}}$$

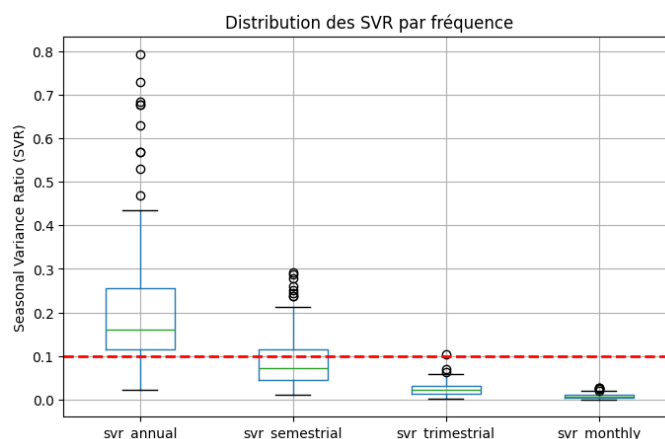
Le seuil significatif statistique est de 10%.

Pour une saisonnalité de période donnée, le SVR est estimé en projetant l'ensemble des observations sur leur phase intra-période, en agrégeant les valeurs par bins de phase sur tous les cycles disponibles, et en calculant la variance des bins rapportée à la variance du signal.

Le nombre de bins contrôle la résolution intra-période. c'est un compromis biais/variance :

- peu de bins = profil très lissé
- beaucoup de bins = profil détaillé mais bruité

Nous avons retenu un bins de 12 qui correspond à un partitionnement mensuel dans une saisonnalité annuelle.



La saisonnalité annuelle est effectivement bien marquée sur une majorité de compteurs. On constate aussi l'existence d'une saisonnalité semestrielle sur quelques compteurs.

4.3 Évaluation des modèles de prédiction

4.3.a Caractéristiques des données

Les séries temporelles analysées présentent les propriétés suivantes :

- Les observations sont strictement positives, la consommation d'eau ne pouvant, par définition, prendre de valeurs négatives.
- Les séries exhibent une saisonnalité annuelle marquée.
- Les séries sont de longueur réduite (moins de 3 ans)
- Elles présentent des évolutions interannuelles caractérisées par des tendances parfois non linéaires, incluant des variations abruptes. Ces évolutions se traduisent par des changements proportionnels du niveau de consommation, impliquant que les phases de baisse s'accompagnent d'un écrasement progressif de la série vers zéro.

Dans ce contexte, la consommation peut être modélisée de manière pertinente comme une décomposition multiplicative de composantes positives :

$$\text{Consommation} = \text{Tendance} * \text{Saisonnalité} * \text{Bruit}$$

4.3.b Analyse des modèles disponibles

Plusieurs familles de modèles de séries temporelles ont été testées et comparées :

Baseline - Comparaison directe $x(t) - x(t - 365)$: Cette approche fournit une référence simple et aisément implémentable pour la détection de variations annuelles. Toutefois, elle se limite à une comparaison ponctuelle entre deux instants et ne permet pas de prendre en compte la tendance.

Modèles fondés sur une fenêtre glissante (par exemple ARIMA, lissage LOESS) : ces approches sont adaptées à des séries globalement stationnaires présentant une saisonnalité dominante et stable. Elles requièrent toutefois un nombre suffisant de cycles saisonniers pour estimer de manière fiable les paramètres du modèle, typiquement de l'ordre de 5 à 7 périodes complètes. Cette condition n'est pas satisfaite dans notre cas, les données disponibles ne couvrant que 3 ans.

Méthodes de machine learning : le jeu de données ne comporte qu'une variable explicative explicite (la date), ce qui restreint l'espace des modèles exploitables à des approches de régression simples. Ces modèles ne sont pas en mesure de capturer de manière adéquate les dynamiques saisonnières et structurelles du phénomène étudié, mais ont été testé pour déterminer les tendances, en particulier avec des fonctions logistiques qui sont bien adaptées à ce profil de consommation.

Approches de deep learning : les séries temporelles étant indépendantes les unes des autres, l'apprentissage nécessiterait l'entraînement d'un modèle distinct par série. Compte tenu de la longueur limitée des séries, cette stratégie conduit à un risque élevé de under-fitting et n'est pas envisageable dans ce contexte.

Décomposition fréquentielle par transformée de Fourier : La transformée de Fourier permet d'isoler efficacement des composantes saisonnières, y compris non multiples, à partir d'un nombre réduit d'observations. Toutefois, dans le cas présent, où la saisonnalité annuelle est dominante et déjà capturée par la baseline, cette approche n'apporte pas de gain significatif supplémentaire. Par ailleurs, la transformée de Fourier ne permet pas de modéliser explicitement les tendances, celles-ci étant assimilées à des cycles de très basse fréquence. En conséquence, cette méthode n'est pas retenue comme solution principale dans notre cadre d'étude.

4.3.c Déploiement des modèles

Les caractéristiques des séries n'ont pas permis de déterminer simplement un modèle. Par ailleurs, la baseline semblait déjà une bonne approche. Pour autant, quel que soit le modèle utilisé, son déploiement devra répondre aux caractéristiques suivantes :

Pour une date de prédiction donnée, le modèle n'exploite exclusivement que les informations disponibles antérieurement à cette date, afin de reproduire un cadre d'utilisation réaliste dans lequel l'utilisateur interroge l'outil à l'instant présent pour la détection d'un événement de fuite en cours.

L'historique de consommation associé à un compteur est alors partitionné en deux segments temporels distincts :

- Période d'évaluation : fenêtre temporelle la plus récente de l'historique, de durée paramétrable, s'achevant à la date de prédiction et sur laquelle une prédiction va être produite par le modèle. Cette période constitue le support principal du calcul des indicateurs utilisés pour la détection.

- Période de référence : ensemble des données historiques antérieures à la période d'évaluation, servant de référence au modèle pour la caractérisation du comportement nominal du compteur.

Sur la période d'évaluation, un résiduel est calculé par la différence entre la consommation réelle et la consommation prédite par le modèle. Ce résiduel peut être positif ou négatif car il est une différence de consommation et non une consommation.

Ce résiduel va être Qualifié (hausse, stabilité, baisse) et Quantifié (faible, moyen, fort) à l'aide de deux metrics :

Trend score : balance d'énergie résiduelle positive vs négative

$$\frac{Var_{Positive} - Var_{Negative}}{Var_{Positive} + Var_{Negative}}$$

borné entre -1 pour un résiduel entièrement baissé et +1 entièrement haussé.

Level score : niveau moyen (absolu) du résiduel vs niveau moyen de consommation

$$\frac{|Moy_{Residual}|}{Moy_{Référence}}$$

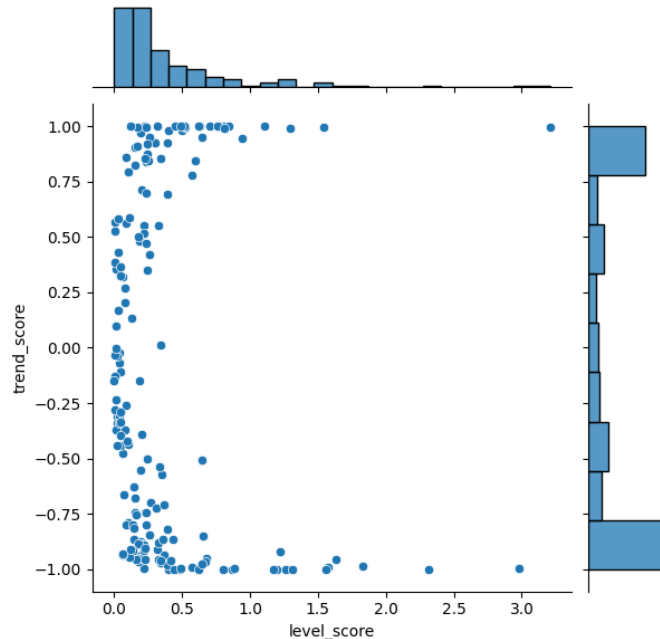
positif, 0 pour un résiduel nul.

La probabilité de fuite va augmenter avec ces deux metrics.

4.4 Baseline

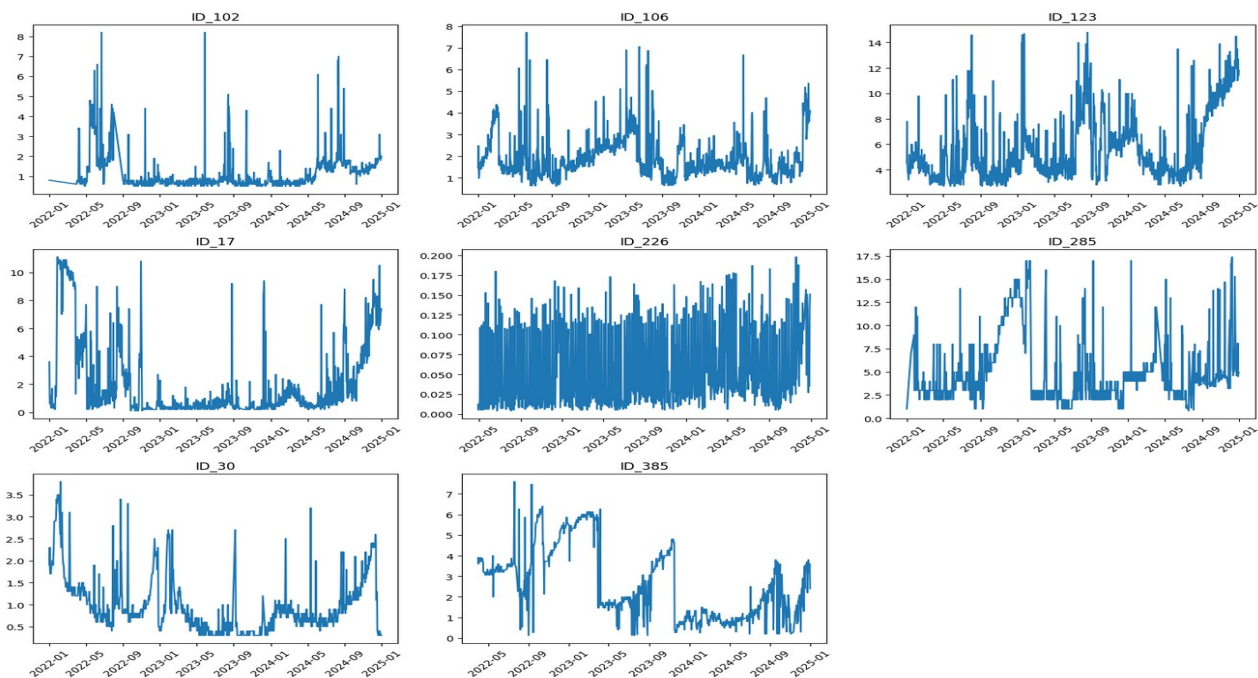
Avec le modèle baseline, la prédiction est la valeur de la fenêtre d'évaluation de l'année précédente.

Résultats obtenus sur les metrics au dernier jour de chaque série (généralement le 31/12/24) et une période d'évaluation de 30 jours (soit 4 semaines afin de répondre au cadrage de projet) :



Observation des 8 compteurs ayant un `trend_score` > 0.9 et un `level_score` > 0.8 :

- Les compteurs sans tendances ou avec des tendances baissières présentent effectivement des profils de fuite : ID_106, ID_123, ID_285, ID_385
- Les autres compteurs ayant des tendances haussières, la détection confond une fuite avec la tendance de consommation.



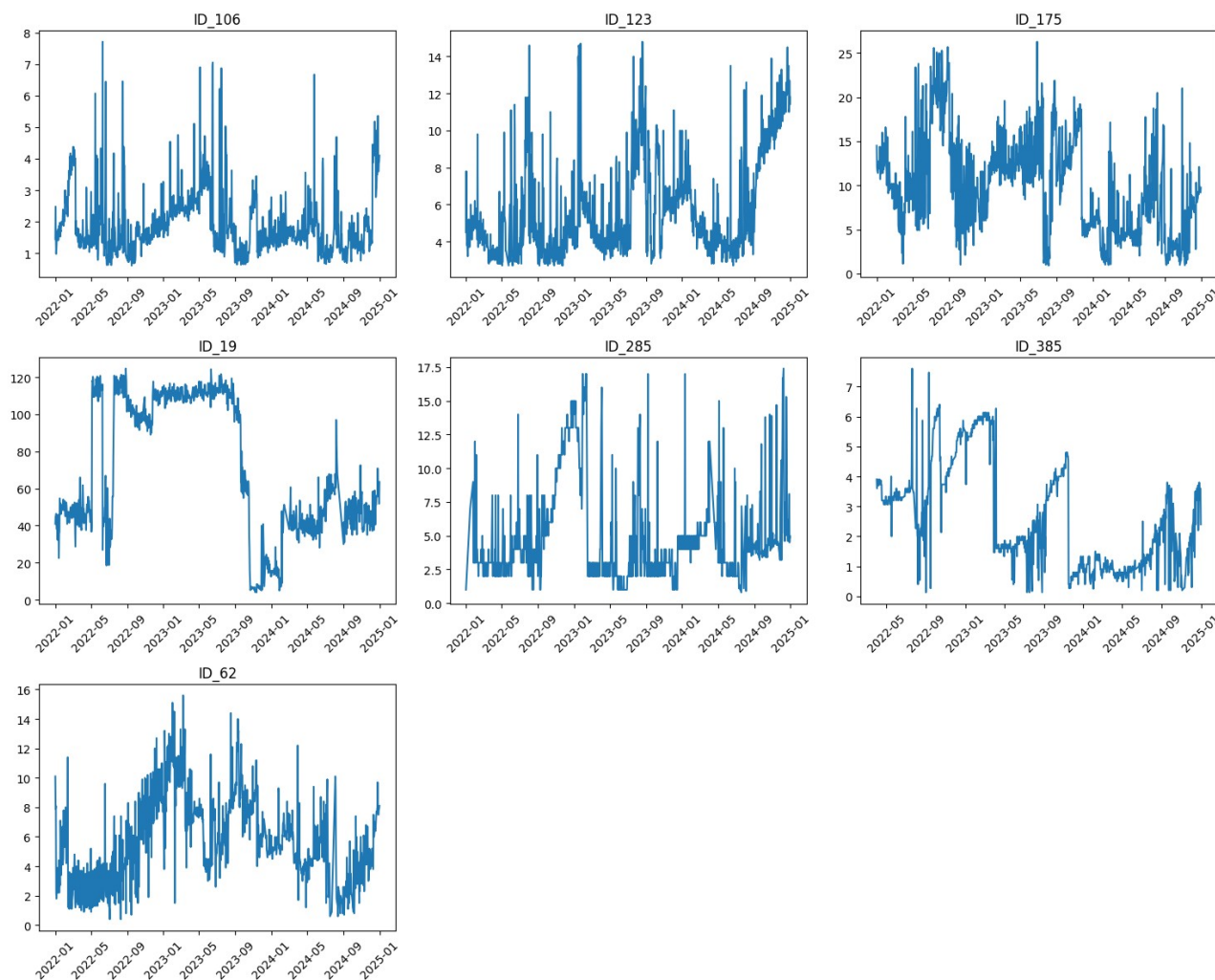
La précision du modèle baseline peut donc être évaluée à 50 %.

4.5 Baseline ajustée de la tendance

En corrigeant la prédiction de la baseline d'une tendance calculée en comparant les 180 jours précédents la période d'évaluation aux 180 jours équivalents de l'année précédente, le modèle isole 7 compteurs ayant un `trend_score` > 0.9 et un `level_score` > 0.7.

On retrouve les 4 compteurs correctement identifiés par la baseline : ID_106, ID_123, ID_285, ID_385.

Les 3 autres compteurs ont des tendances de consommation fortement baissières présentant un mois de décembre 2024 en hausse importante alors que traditionnellement, c'est un mois de consommation basse.



La précision est relativement proche de 100 % qui était l'objectif recherché.

4.6 Conclusion

Un modèle relativement simple inspiré de la baseline a permis d'obtenir un premier résultat proche de l'attendu avec des séries de taille réduite (575 jours nécessaires). Des améliorations sont possibles, sur le calcul de la tendance par exemple qui se limite pour le moment à une constante. Un feedback du client est nécessaire avant d'aller plus loin.

4.7 Bibliographie

Moyenne mobile :

- <https://blog.statoscop.fr/timeseries-4.html>

Pandas échantillonnage et interpolation :

- <https://stackoverflow.com/questions/30530001/python-pandas-time-series-interpolation-and-regularization>
- <https://pandas.pydata.org/docs/reference/api/pandas.Series.interpolate.html>

Détection des ruptures :

- <https://centre-borelli.github.io/ruptures-docs/code-reference/detection/kernelcpd-reference/>

Transformée rapide de Fourier :

- https://fr.wikipedia.org/wiki/Transformation_de_Fourier_rapide
- <https://docs.scipy.org/doc/scipy/tutorial/fft.html>
- <https://numpy.org/devdocs/reference/generated/numpy.fft.fftfreq.html>

auto corrélation, décomposition saisonnière :

- <https://www.statsmodels.org/stable/index.html>
- https://en.wikipedia.org/wiki/Decomposition_of_time_series

traitement des séries temporelles avec un RNN :

- https://www.tensorflow.org/tutorials/structured_data/time_series?hl=fr

Time serie Classifier avec un RNN:

- https://keras.io/examples/timeseries/timeseries_classification_from_scratch/

Traitement de signal

- https://fr.wikipedia.org/wiki/%C3%89nergie_d%27un_signal
- https://fr.wikipedia.org/wiki/Rapport_signal_sur_bruit

Saisonnalité

- <https://en.wikipedia.org/wiki/Seasonality>

Fonctions logistiques

- [https://fr.wikipedia.org/wiki/Fonction_logistique_\(Verhulst\)](https://fr.wikipedia.org/wiki/Fonction_logistique_(Verhulst))

5 DESCRIPTION DÉTAILLÉE DES TRAVAUX RÉALISÉS (Abdoul Aziz Toure)

5.1 Préparation des données

5.1.a Visualisation globale de la consommation

La première visualisation (Figure 1) présente l'évolution de la consommation d'eau entre **2019 et 2025**, après suppression des valeurs extrêmes manifestement aberrantes.

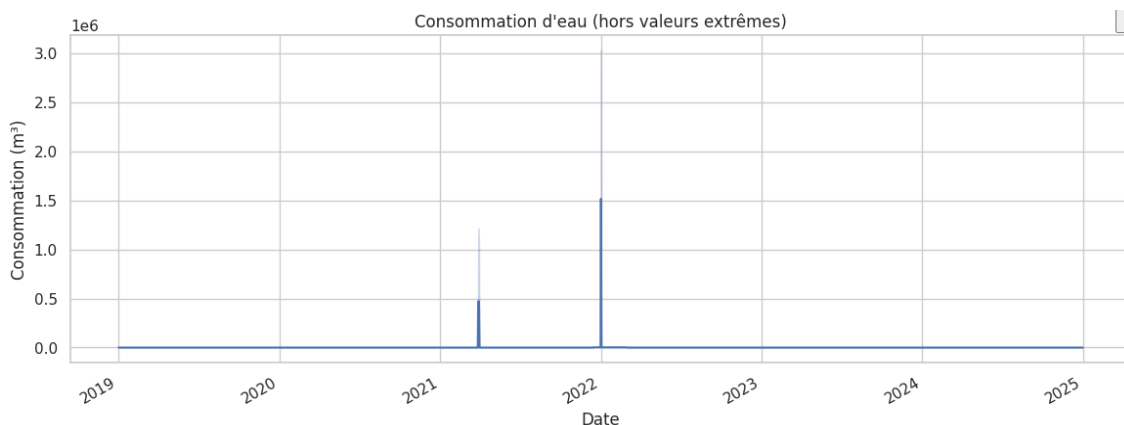


Figure 1 : Consommation d'eau avec valeur extremes

Constats majeurs :

- La série est dominée par quelques pics très élevés (jusqu'à plusieurs centaines de milliers de m³)
- Ces valeurs écrasent totalement l'échelle et masquent le comportement réel du réseau
- Une analyse directe sur les données brutes est donc non exploitable

Conclusion : un prétraitement robuste est indispensable avant toute modélisation.

5.1.b Suppression des outliers extrêmes (99^e percentile)

Afin de révéler la structure réelle de la consommation, un filtrage basé sur le 99^e percentile a été appliqué.

Effets observés (Figure 2) :

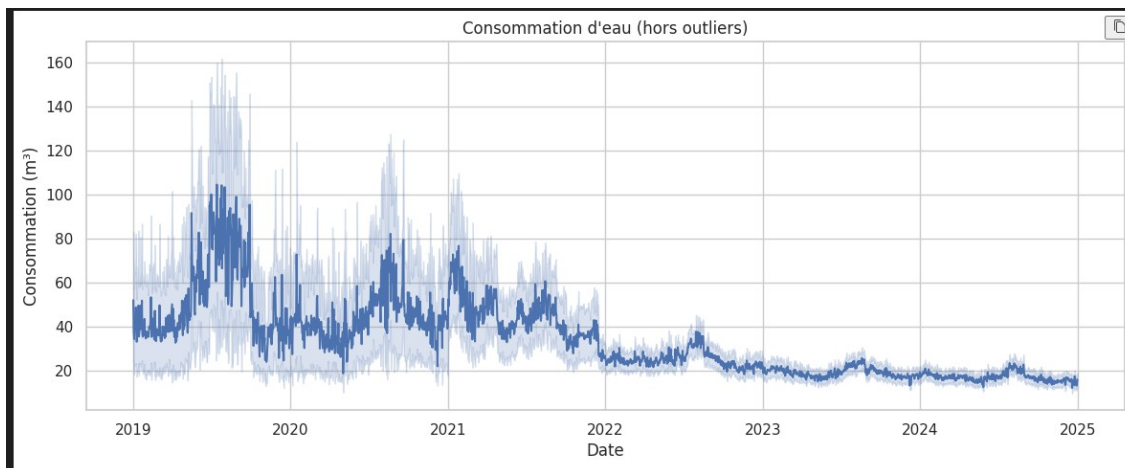


Figure 2 :

Consommation d'eau sans les outliers

- Apparition claire d'un niveau moyen de consommation
- Mise en évidence de variations régulières au cours de l'année
- Disparition des artefacts liés à des erreurs de mesure ou événements exceptionnels

Cette étape permet d'obtenir une vision réaliste et exploitable du fonctionnement hydraulique.

5.2 Mise en évidence de la saisonnalité annuelle

5.2.a Saisonnalité forte et structurée

Les figures obtenues après filtrage montrent clairement :

- Une saisonnalité annuelle marquée
- Des pics récurrents correspondant aux périodes de forte consommation (été)
- Des creux en période hivernale

Point clé métier :

Une augmentation de consommation en été n'est pas une anomalie, mais un comportement normal du réseau (arrosage, usage domestique accru, chaleur).

Cette observation impose une contrainte forte au cahier des charges :

Le modèle ne doit jamais confondre saisonnalité et fuite.

5.3 Identification d'une période de référence « saine »

5.3.a Justification du découpage temporel

Un découpage temporel a été effectué autour du 1er janvier 2022, conduisant à :

- Période 1 (avant 2022) : historique long, hétérogène, potentiellement contaminé par des fuites

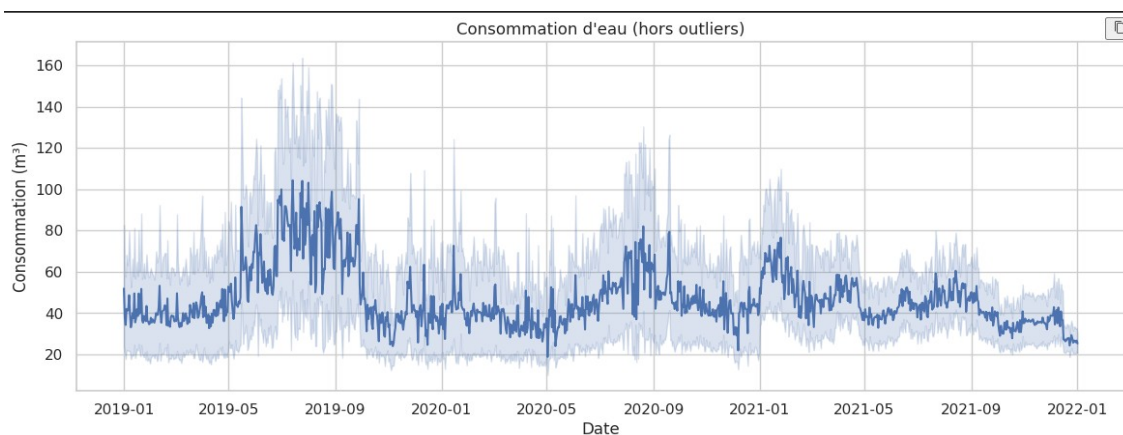


Figure 3 :

Consommation d'eau de la période de 2019-2022

- Période 2 (2022–2024) : période stable, visuellement plus régulière

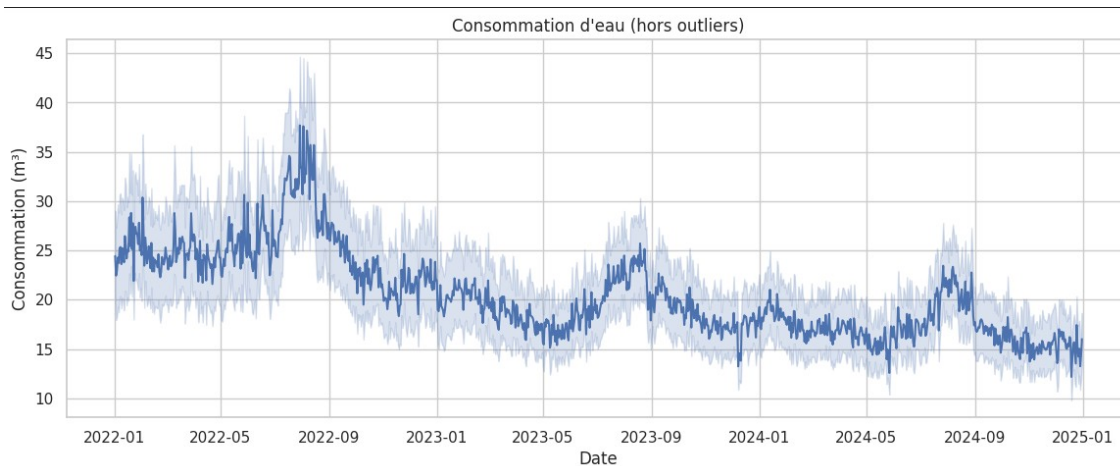


Figure 4 :

Consommation d'eau de la période de 2022-2025

Les visualisations montrent que la période 2022–2024 présente :

- Une consommation plus lisse
- Une variabilité plus faible
- Une saisonnalité cohérente et répétable

Cette période est donc considérée comme une période propre (baseline).

Décision experte :

La période 2022–2024 est utilisée comme référence de comportement normal du réseau.

5.4 Cahier des charges fonctionnel et technique

À partir des constats précédents, le cahier des charges du système de détection est défini comme suit.

5.4.a Contraintes fonctionnelles

Le système doit :

1. Apprendre le comportement normal à partir d'une période saine
2. Intégrer explicitement :
 - o La saisonnalité annuelle
 - o Les niveaux moyens propres à chaque compteur
3. Détecter :
 - o Les fuites brutales (pics anormaux)
 - o Les fuites lentes (dérive progressive)
4. Fournir une visualisation claire et interprétable pour les équipes métier

5.4.b Contraintes mathématiques et data science

Le modèle doit :

- Être robuste aux outliers
- Ne pas nécessiter de données labellisées (approche non supervisée)
- Produire une séparation explicite entre :
 - o Comportement normal
 - o Comportement anormal
- Être scalable à plusieurs compteurs (IDs)

5.5 Positionnement méthodologique retenu

Compte tenu :

- de la saisonnalité forte,
- de l'absence de labels fiables,
- du besoin d'interprétabilité,

l'approche retenue repose sur :

- Transformation logarithmique (stabilisation de variance)
- Centrage par compteur
- Apprentissage du comportement normal sur période saine
- Détection d'écarts persistants dans le temps

Ce cadre justifie pleinement l'utilisation ultérieure de :

- Holt-Winters (approche temporelle)

5.5.a Conclusion Partielle

Cette analyse exploratoire démontre que :

- La problématique de détection de fuites ne peut être abordée sans une compréhension fine de la saisonnalité
- La période 2022–2024 constitue une référence fiable du comportement normal
- Toute solution crédible doit intégrer temps, saison et structure globale des données

Cette section pose ainsi les fondations scientifiques et métier nécessaires à la conception d'un système de détection de fuites robuste et innovant.

5.6 Utilisation du modèle et détection des anomalies

5.6.a Objectif de la phase de modélisation

L'objectif de cette phase est de modéliser le comportement normal de la consommation d'eau afin d'identifier automatiquement les écarts anormaux, interprétés comme des fuites potentielles.

Contrairement à une approche supervisée, aucune étiquette "fuite / non-fuite" n'est disponible. Le problème est donc formulé comme une détection d'anomalies temporelles non supervisée, en tenant compte :

- de la saisonnalité annuelle,
- des différences de niveau entre compteurs,
- de la robustesse face aux valeurs extrêmes.

5.6.b Choix méthodologique : Holt-Winters global normalisé

Le modèle retenu repose sur :

- une transformation logarithmique (stabilisation de variance),
- un centrage par compteur (normalisation inter-ID),
- un modèle Holt-Winters additif (tendance + saisonnalité annuelle),
- une détection robuste des anomalies via MAD.

Ce choix permet :

- d'expliquer la saisonnalité (été \neq fuite),
- de séparer clairement comportement normal **et** résidus anormaux,
- d'obtenir un modèle interprétable et industriellement exploitable.

5.7 Prétraitement des données

5.7.a Tri temporel et transformation logarithmique

```
# Tri chronologique
data_plot_filtered = data_plot_filtered.sort_values("valeur_date")

# Transformation logarithmique (stabilisation de la variance)
data_second_half["log_valeur"] = np.log(data_second_half["valeur_active"])
data_first_half["log_valeur"] = np.log(data_first_half["valeur_active"])
```

Justification

La transformation logarithmique permet :

- de réduire l'influence des pics,
- de rendre les écarts relatifs comparables,
- d'améliorer la stabilité du modèle Holt-Winters.

5.7.b Calcul de la moyenne par compteur (train uniquement)

```
# Moyenne log par ID calculée uniquement sur la période saine (train)
mu_id = (
    data_second_half
    .groupby("libelle")["log_valeur"]
    .mean()
    .rename("mu_id")
)
```

Principe clé :

La moyenne par compteur est calculée uniquement sur la période de référence saine (**2022–2024**) afin d'éviter toute fuite d'information depuis la période test.

5.7.c Jointure et centrage par ID

```
# Suppression d'éventuelles colonnes existantes
for df in [data_second_half, data_first_half]:
    if "mu_id" in df.columns:
        df.drop(columns="mu_id", inplace=True)

# Ajout de mu_id
data_second_half = data_second_half.join(mu_id, on="libelle")
data_first_half = data_first_half.join(mu_id, on="libelle")

# Centrage
data_second_half["log_centered"] = (
    data_second_half["log_valeur"] - data_second_half["mu_id"]
)
data_first_half["log_centered"] = (
    data_first_half["log_valeur"] - data_first_half["mu_id"]
)
```

Pourquoi le centrage est essentiel ?

- Chaque compteur a un niveau de consommation propre.
- Le centrage permet au modèle d'apprendre une dynamique commune, indépendante du volume absolu.
- Sans centrage, les compteurs à forte consommation domineraient le modèle.

5.8 Apprentissage du modèle Holt-Winters

5.8.a Construction de la série temporelle d'apprentissage

```
train_series = (
    data_second_half
    .sort_values("valeur_date")
    .set_index("valeur_date")["log_centered"]
)
```

La série d'apprentissage représente le comportement normal agrégé du réseau.

5.8.b Entraînement du modèle Holt-Winters

```
hw = ExponentialSmoothing(
    train_series,
    trend="add",
    seasonal="add",
    seasonal_periods=365
).fit()
```

Interprétation mathématique

Le modèle apprend :

- une tendance additive (évolution lente du réseau),
- une saisonnalité annuelle (usage estival / hivernal),
- un niveau moyen nul (grâce au centrage).

5.9 Prédiction sur la période test

```
forecast_log_centered = hw.forecast(len(data_first_half))
```

```
forecast_log_centered.index = (
    data_first_half
    .sort_values("valeur_date")
    .index
)
```

5.9.a Reprojection dans l'espace réel

```
data_first_half["prediction"] = np.exp(
    forecast_log_centered.values + data_first_half["mu_id"].values
)
```

5.9.b Interprétation

On reconstruit la consommation attendue :

1. moyenne du compteur
2. exponentiation inverse du log

5.10 Détection robuste des anomalies

5.10.a Calcul des résidus

```
data_first_half["residual"] = (
    data_first_half["valeur_active"] -
    data_first_half["prediction"]
)
```

5.10.b Seuil robuste basé sur la MAD

```
from statsmodels.robust.scale import mad
```

```
# Écart-type robuste
sigma = mad(data_first_half["residual"])
```

```
# Seuil robuste (3-sigma)
threshold = 3 * sigma
```

Pourquoi MAD plutôt que l'écart-type ?

- Insensible aux valeurs extrêmes

- Adapté aux distributions non gaussiennes
- Standard en détection d'anomalies industrielles

5.10.c Détection des anomalies

```
data_first_half["anomaly"] = (
    data_first_half["residual"].abs() > threshold
).astype(int)
```

5.11 Visualisation des résultats globaux

```
plt.figure(figsize=(14,5))

plt.plot(
    data_first_half["valeur_date"],
    data_first_half["valeur_active"],
    label="Consommation réelle"
)

plt.plot(
    data_first_half["valeur_date"],
    data_first_half["prediction"],
    "--",
    label="Comportement normal"
)

plt.scatter(
    data_first_half[data_first_half["anomaly"]==1]["valeur_date"],
    data_first_half[data_first_half["anomaly"]==1]["valeur_active"],
    color="red",
    label="Anomalie"
)

plt.title("Détection d'anomalies – Holt-Winters global (log + normalisation ID)")
plt.legend()
plt.show()
```

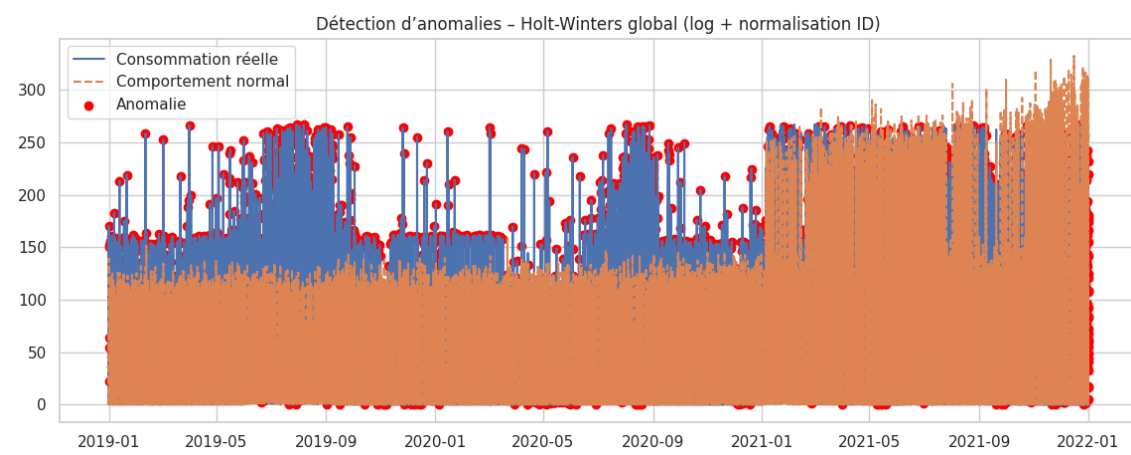


Figure 5 : Détection d'anomalies globale

Lecture métier

- Courbe bleue : consommation observée
- Courbe pointillée : consommation attendue
- Points rouges : anomalies potentielles (fuites)

5.12 Analyse par compteur

5.12.a *Extraction de caractéristiques par ID*

```
features_by_id = (
    data_first_half
    .groupby("libelle")
    .agg(
        mean_residual=("residual", "mean"),
        mad_residual=("residual", mad),
        max_residual=("residual", lambda x: np.max(np.abs(x))),
        anomaly_rate=("anomaly", "mean")
    )
)
```

Ces indicateurs permettent :

- le **classement des compteurs à risque**,
- la priorisation des interventions terrain.

5.12.b *Visualisation multi-compteurs*

```
top_ids = (
    data_first_half["libelle"]
    .drop_duplicates()
    .sample(n=5)
    .values
)

fig, axes = plt.subplots(5, 1, figsize=(14, 20), sharex=True)

for ax, libelle in zip(axes, top_ids):
    data_id = data_first_half[data_first_half["libelle"] == libelle]

    ax.plot(data_id["valeur_date"], data_id["valeur_active"], label="Réal")
    ax.plot(data_id["valeur_date"], data_id["prediction"], "--", label="Normal")
    ax.scatter(
        data_id[data_id["anomaly"]==1]["valeur_date"],
        data_id[data_id["anomaly"]==1]["valeur_active"],
        color="red"
    )
    ax.set_title(f"Compteur {libelle}")
    ax.legend()

plt.tight_layout()
plt.show()
```

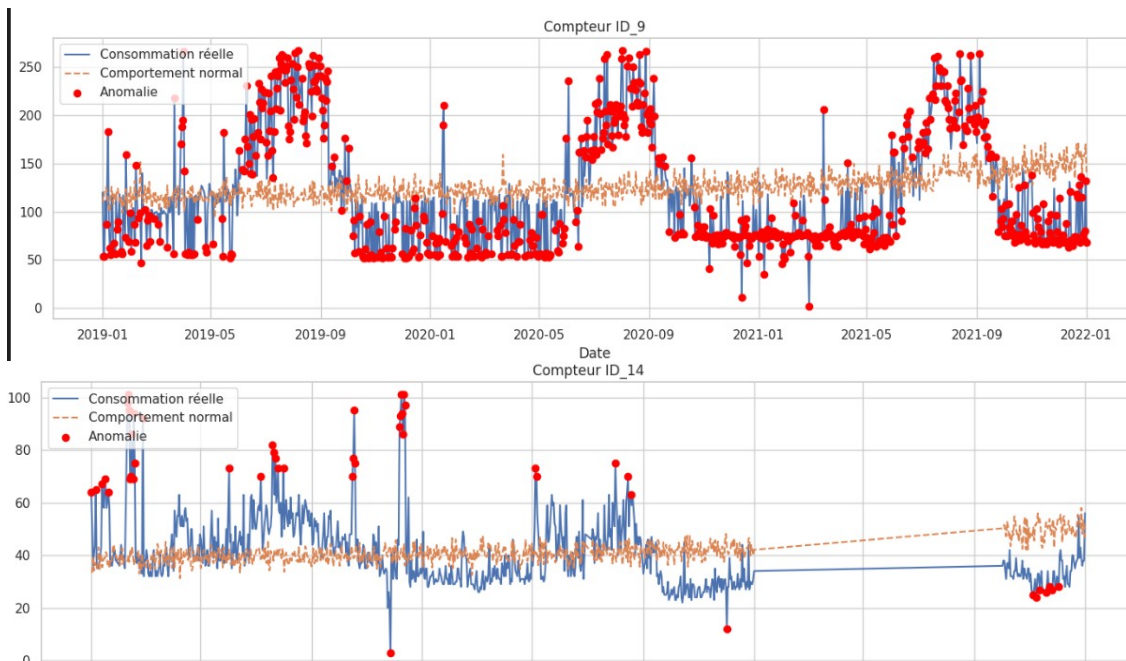



Figure 6 : Détection d'anomalie par compteur

5.13 Évaluation quantitative du modèle

L'évaluation est réalisée **uniquement sur les points considérés comme normaux**, afin de mesurer la qualité de modélisation du comportement sain.

```
from sklearn.metrics import mean_absolute_error, mean_squared_error, r2_score
```

```
mask_normal = data_first_half["anomaly"] == 0
```

```
y_true = data_first_half.loc[mask_normal, "valeur_active"]
```

```
y_pred = data_first_half.loc[mask_normal, "prediction"]
```

```
mae = mean_absolute_error(y_true, y_pred)
```

```
rmse = np.sqrt(mean_squared_error(y_true, y_pred))
```

```
r2 = r2_score(y_true, y_pred)
```

```
mape = np.mean(np.abs((y_true - y_pred) / y_true)) * 100
```

5.13.a Résultats obtenus

- **MAE** : 5.67
- **RMSE** : 8.50
- **MAPE** : 69.56 %
- **R²** : 0.936

5.13.b Interprétation

- Le R² élevé confirme une excellente modélisation du comportement normal
- Le MAPE élevé est attendu sur de faibles consommations (effet relatif)
- Le modèle est suffisamment précis pour la détection d'écarts anormaux

5.14 Conclusion partielle

Cette approche Holt-Winters normalisée permet :

- une détection fiable des anomalies,
- une prise en compte explicite de la saisonnalité,
- une interprétabilité métier immédiate,
- une base solide pour une industrialisation à grande échelle.