

Exploration des données

CM3 : Séries temporelles

Mickaël Martin Nevot

V1.0.0



Cette œuvre de Mickaël Martin Nevot est mise à disposition sous licence Creative Commons Attribution - Utilisation non commerciale - Partage dans les mêmes conditions.

Exploration des données

- I. Présentation
- II. Introduction
- III. Détection d'anomalies
- IV. Extraction de caractéristiques
- V. Séries temporelles

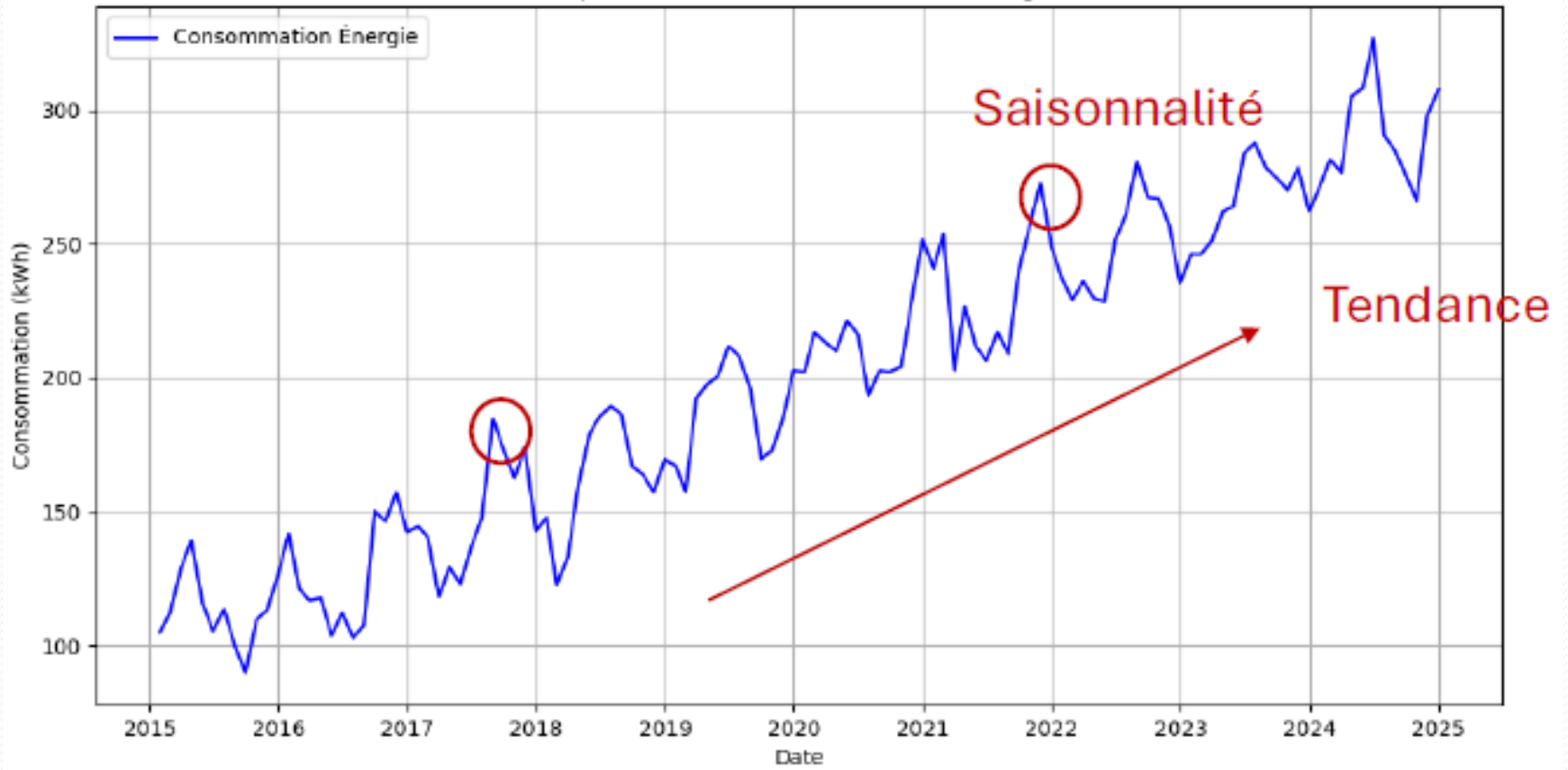
Séries temporelles

- Suite d'observations répétées d'un même phénomène à des dates différentes
- Exemple :
 - Économie : chômage, croissance, actions en bourse, etc.
 - Environnement : pression, température, niveau d'eau, etc.
 - Santé : électrocardiogramme, etc.
- Utilité :
 - Compréhension du phénomène représenté par la série
 - Prédictions

Dépendance dans le temps et comportements saisonniers

Exemple

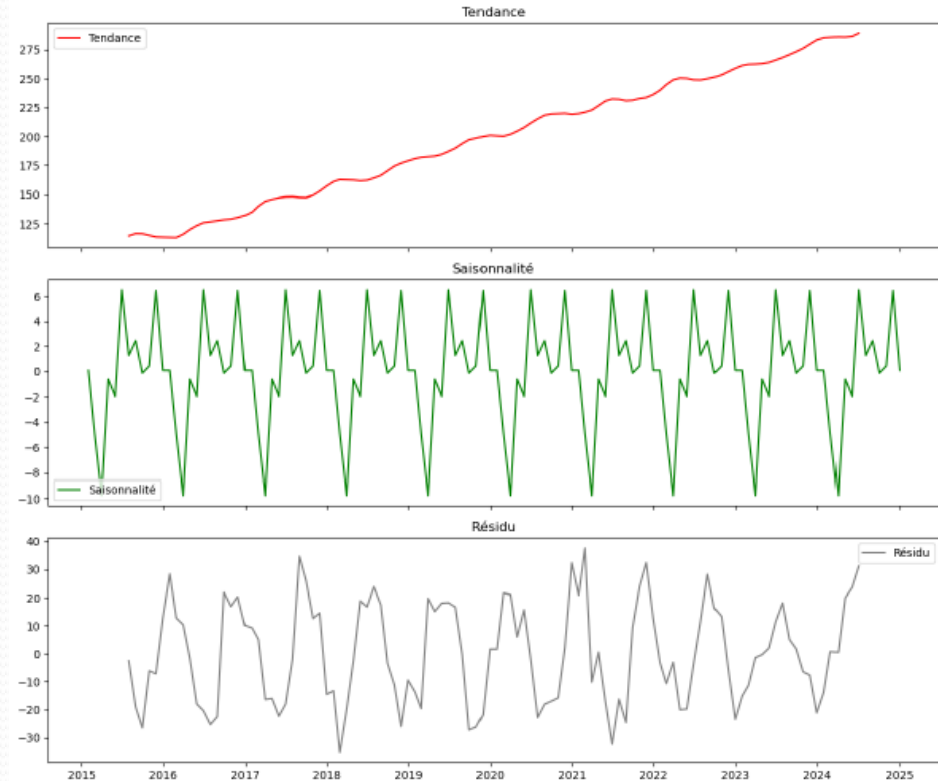
Série Temporelle de la Consommation d'Énergie (kWh)



Décomposition

En ses composantes

- Composants :
 - **Tendance (T)** : stabilité / augmentation / diminution dans le temps
 - **Saisonnalité (S)** : variations périodiques (facteurs saisonniers ?)
 - **Résidu (R)** : variations aléatoires inexplicables par T / S (bruit) / erreur entre valeur observée et valeur prédite

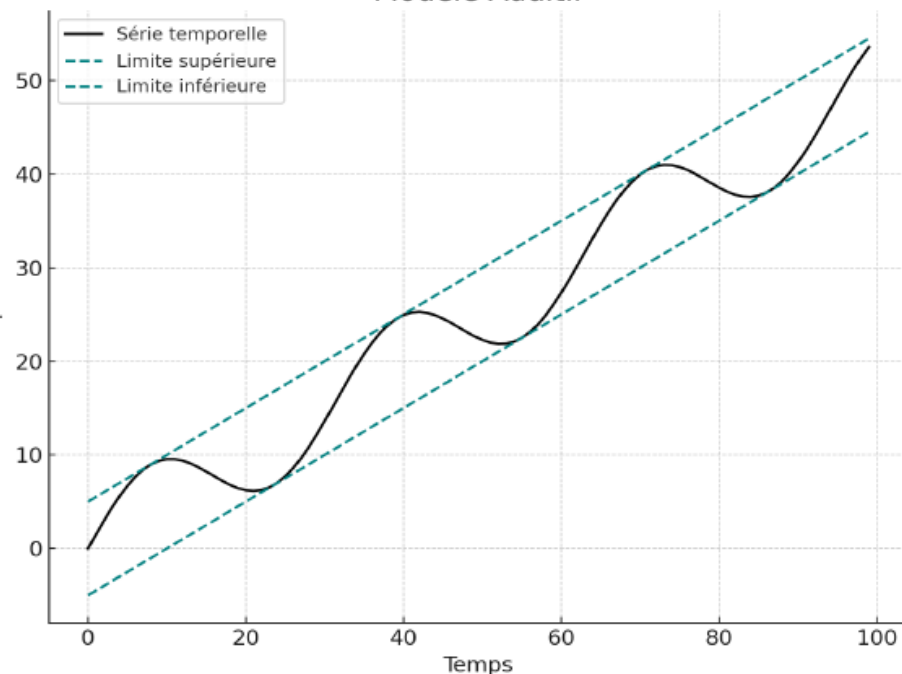


Modèles de décomposition

Modèle additif

- Amplitude : $Y_t = T_t + S_t + R_t$
- Ex. : temp. journalière en hiver

Modèle Additif

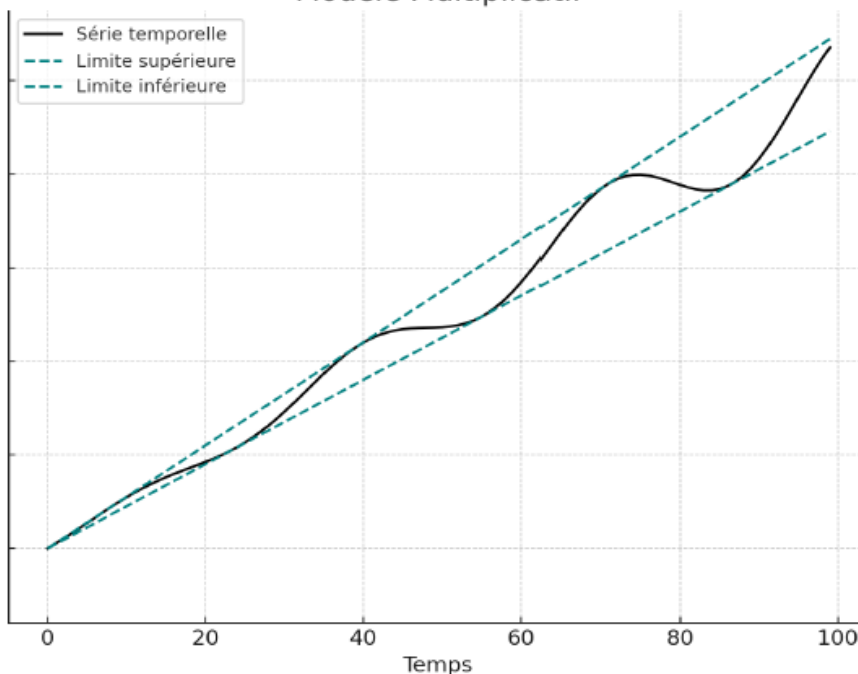


La saisonnalité a une amplitude constante, indépendamment de la tendance

Modèle multiplicatif

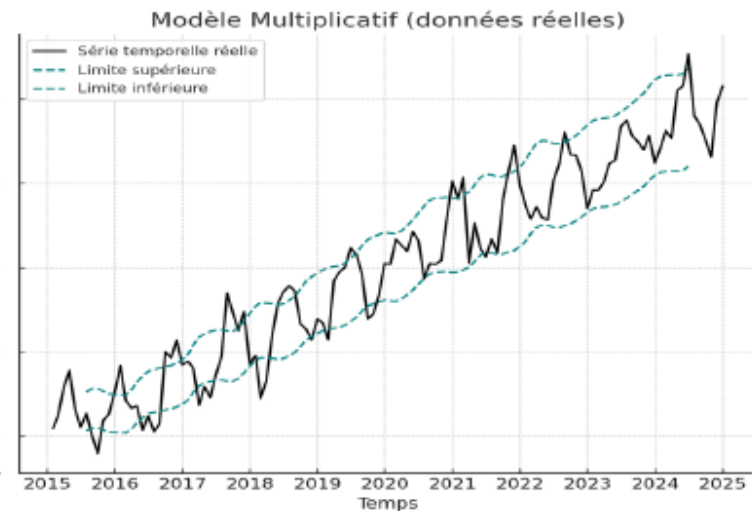
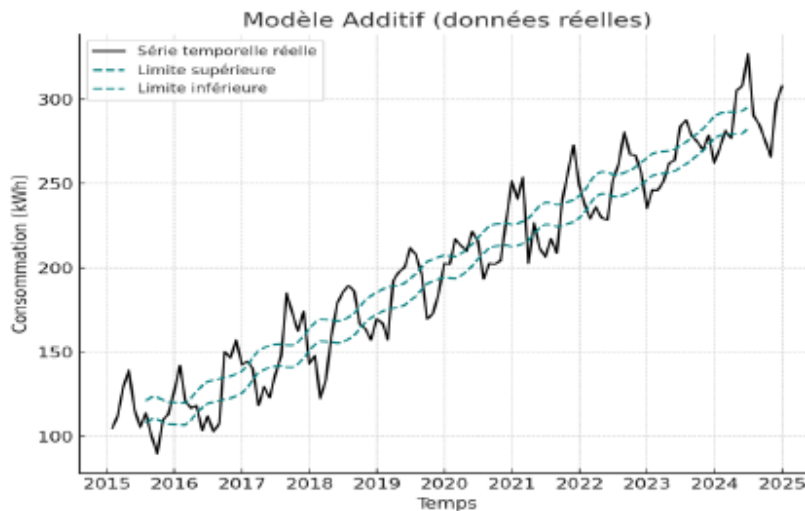
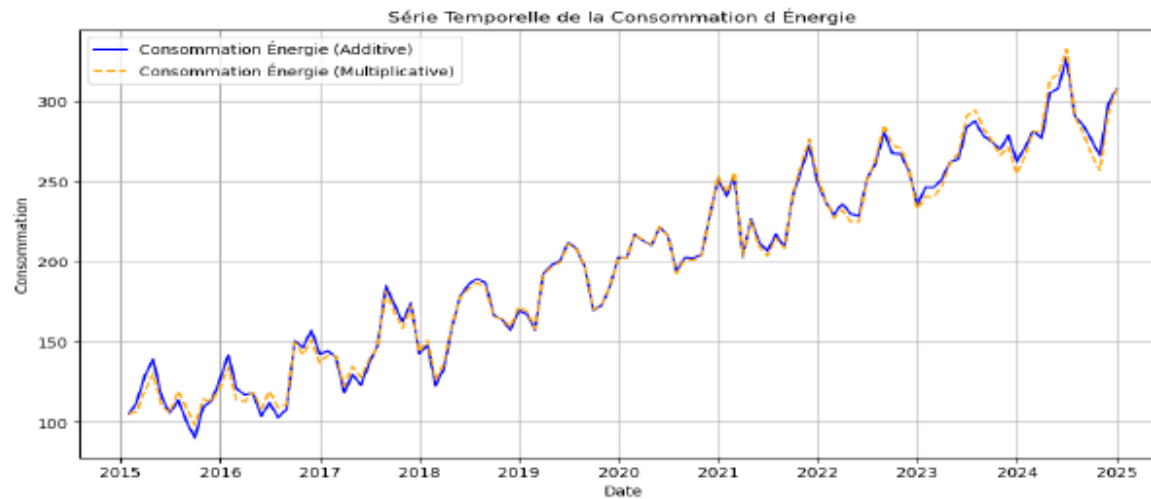
- Amplitude : $Y_t = T_t \times S_t \times R_t$
- Ex. : taux de pluie par mois

Modèle Multiplicatif



L'amplitude des variations augmente avec la tendance

Modèles de décomposition

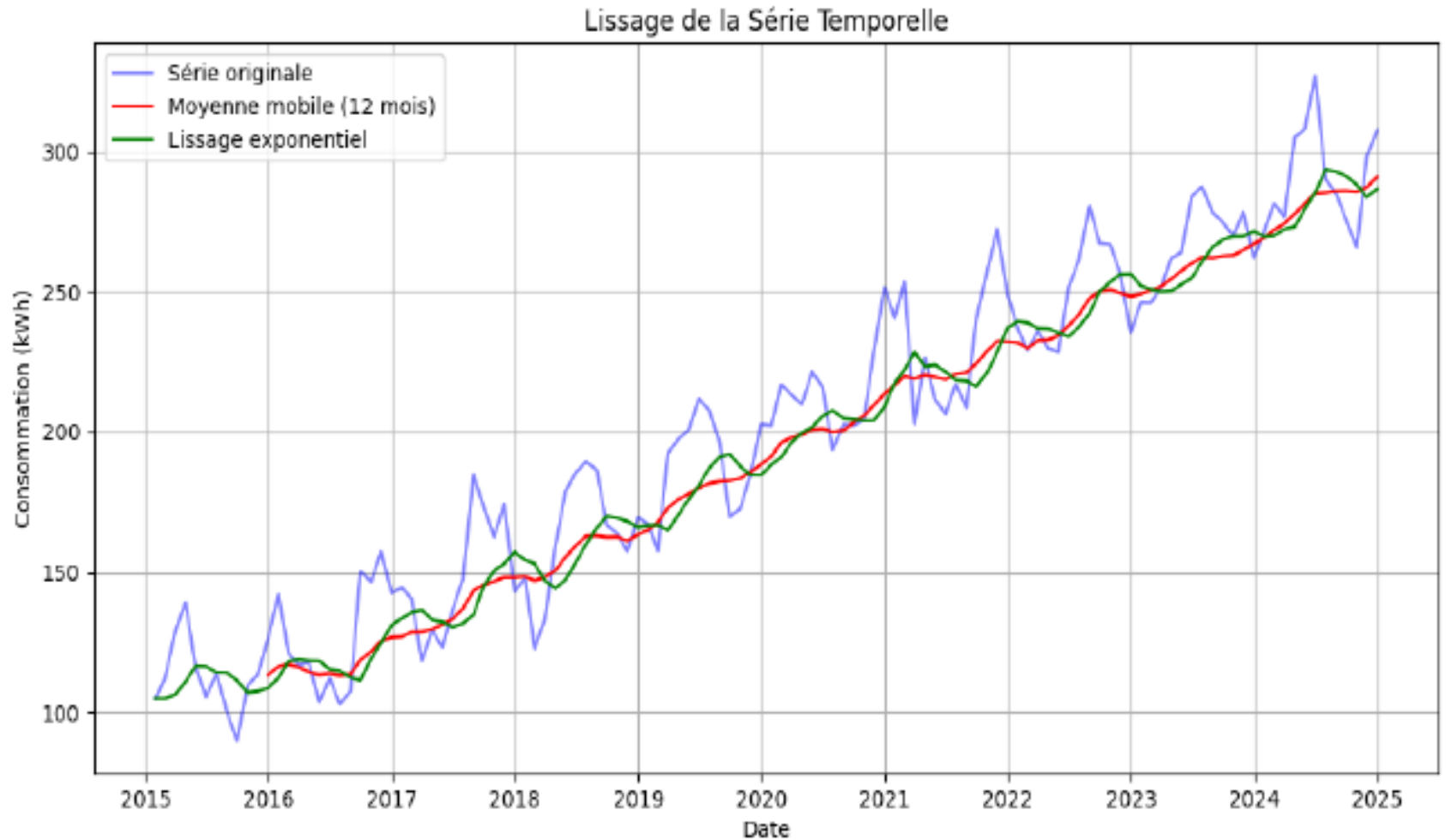


Lissage / filtrage

- Application d'une transformation pour :
 - Réduire le bruit
 - Déterminer la structure (T et S)
 - Révèle les composantes non-stationnaires :
 - Tendance : non-stationnarité en moyenne
 - Variation d'amplitude : non-stationnarité en variance
- Méthodes de lissage :
 - **Moyenne mobile** : remplace chaque valeur par la moyenne des k valeurs précédentes
 - **Lissage exponentiel** : observations récentes ont plus de poids que les plus anciennes (sensible aux fluctuations)

Méthode de Holt-Winters : lissage exponentiel avec prise en compte de T et S

Lissage / filtrage



Stationnarité

- Problèmes de la non-stationnarité :
 - Résultats biaisés
 - Mauvaise interprétation
 - Incohérences dans les résultats des tests statistiques
 - **Imprédictibilité**
- Tester la stationnarité :
 - Méthodes visuelles : graphique, **ACF**, **PACF**
 - Méthodes statistiques : **ADF**, **KPSS**, **PP**

Les propriétés statistiques (moyenne, variance, autocovariance, autocorrélation) d'une série temporelle stationnaire ne varient pas dans le temps

Une série temporelle non-stationnaire devrait être transformée

Bruit blanc : série temporelle purement aléatoire

Tester la stationnarité

1. Visualisation

- Graphique, ACF

2. Vérifier la stationnarité en moyenne (test 1)

- Test ADF, KPSS

3. Vérifier la saisonnalité

- ACF/PACF, décomposition

4. Détection de la non stationnarité en variance

- Test de ARCH

5. Correction de la non stationnarité

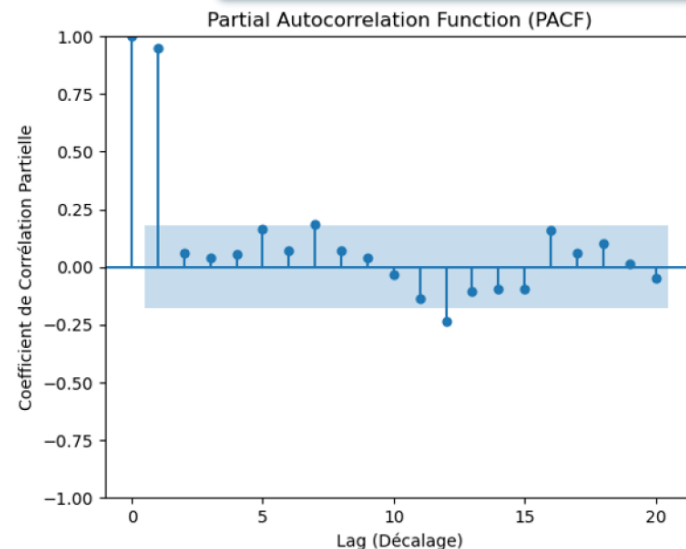
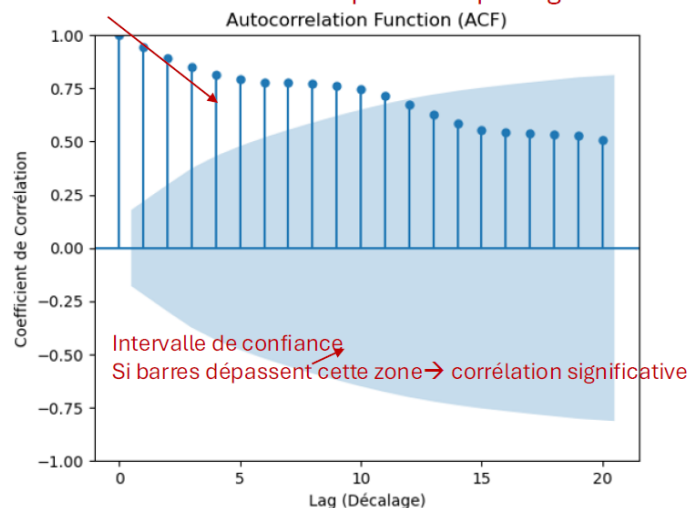
- Différenciation, log, ..

Stationnarité : visuellement

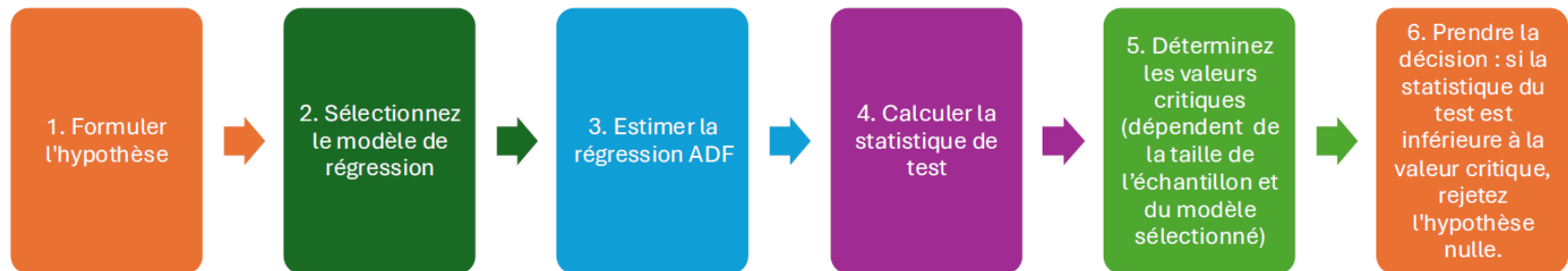
Type non-stationnarité	ACF	PACF
Tendance (moyenne)	Diminue lentement, ne coupe pas à zéro	Forte valeur au lag n°1 Décroit lentement
Saisonnière (périodicité)	Pics périodiques	Pics périodiques aux mêmes lags
Variance (hétéroscédasticité)	Instable	Variable selon les périodes

Lag : nb de périodes de retard

Indicateurs de corrélation pour chaque lag



ADF

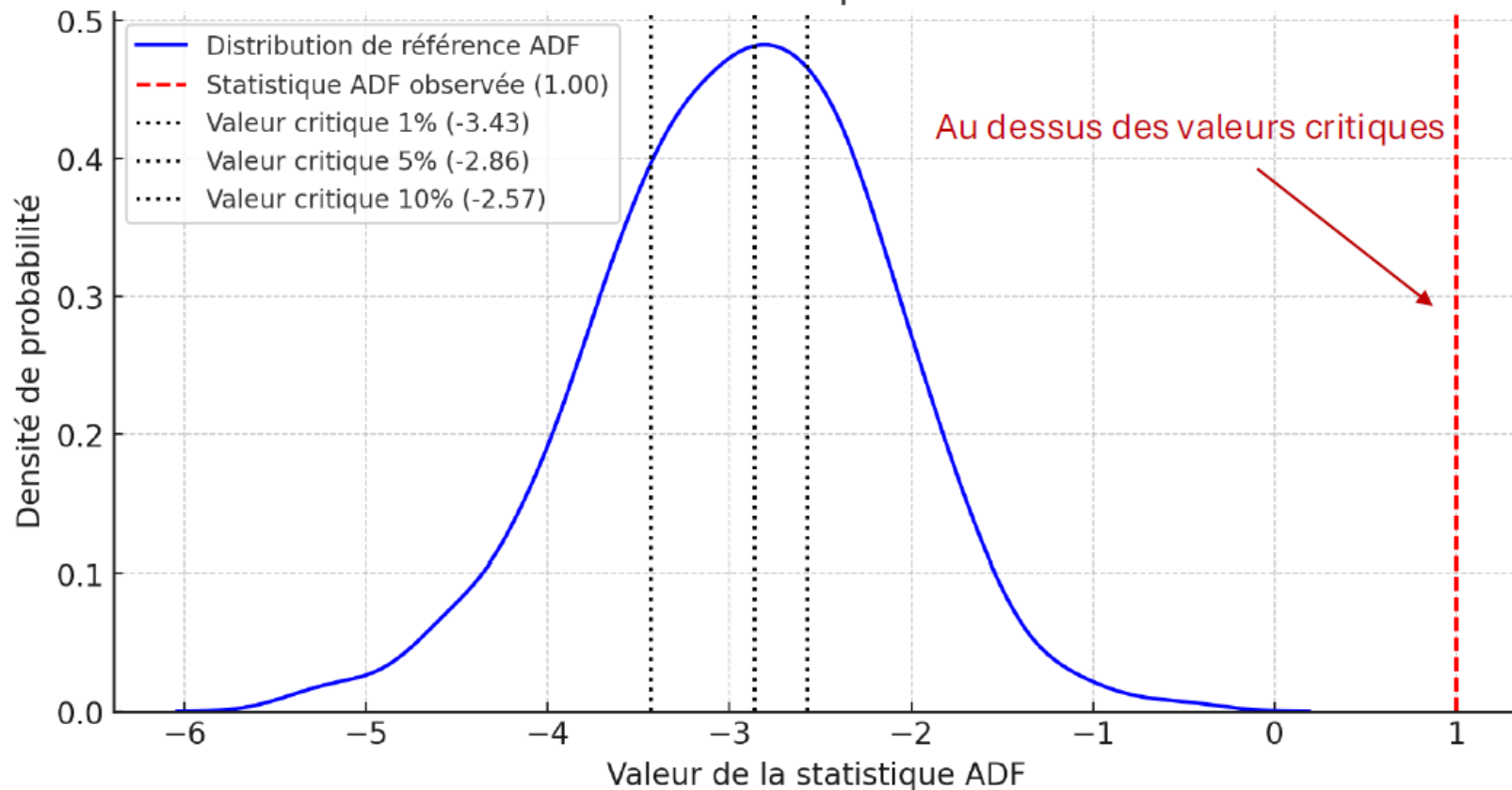


- **Hypothèse nulle** : série **non-stationnaire**
- Résultats : Générée par un processus présentant une racine unitaire
 - Statistique ADF : si $<$ valeurs critiques, série stationnaire
 - *p-value* (probabilité d'erreur) : si ≤ 0.05 , série stationnaire
 - Valeurs critiques : seuils de référence pour décision

ADF sert aussi à tester la stationnarité des résidus

ADF

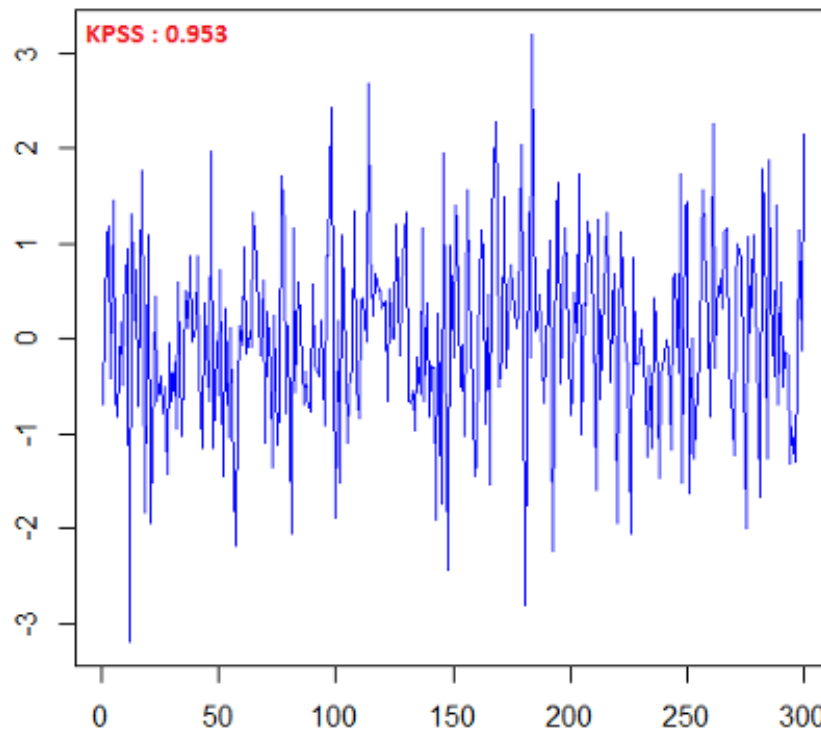
Illustration du calcul de la p-value dans le test ADF



Stat. ADF	<i>p-value</i>	Valeur critique 1%	Valeur critique 5%	Valeur critique 10%
1.0018	0.994	-3.4918	-2.888	-2.5811

KPSS

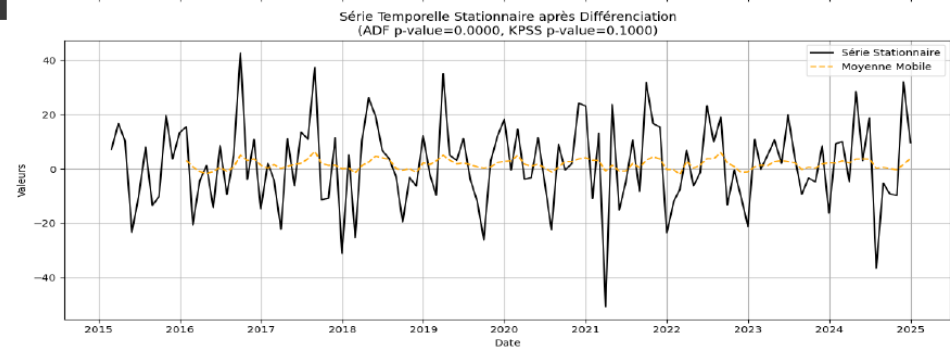
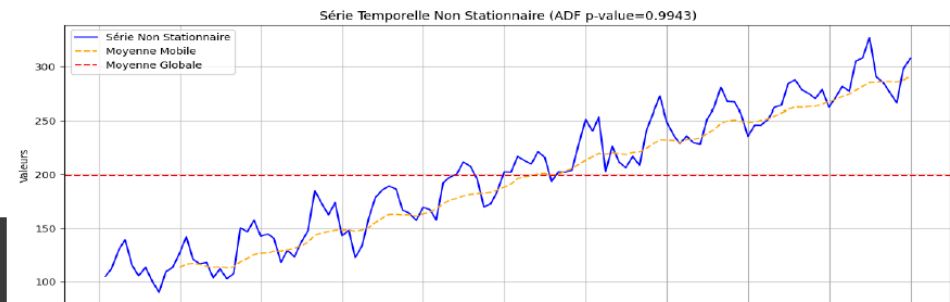
- Similaire à ADF avec **hypothèse nulle : série stationnaire**
- Résultats :
 - *p-value* : si > 0.05 : série stationnaire



ADF + KPSS

ADF	KPSS	Interprétation
Non-stationnaire	Stationnaire	Stationnaire (autour d'une tendance)
Non-stationnaire	Non-stationnaire	Série à transformer
Stationnaire	Stationnaire	Stationnaire
Stationnaire	Non-stationnaire	Incohérent, à compléter avec ACF/PACF

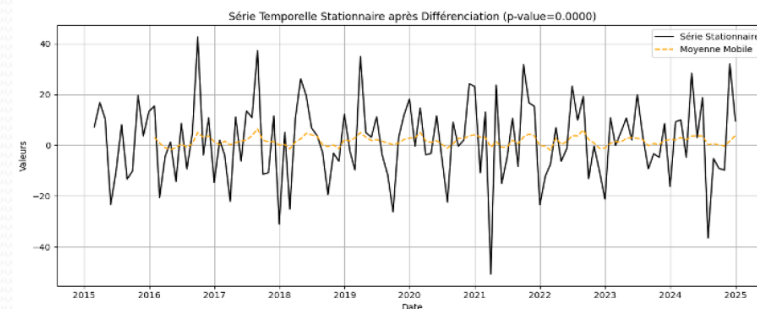
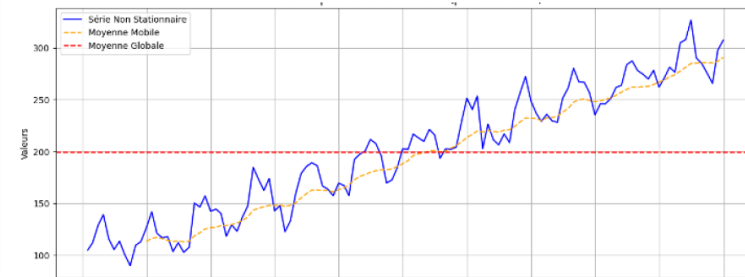
Test	Statistic (Original)	p-value (Original)
ADF	1.001846	0.994285
KPSS	1.790778	0.010000



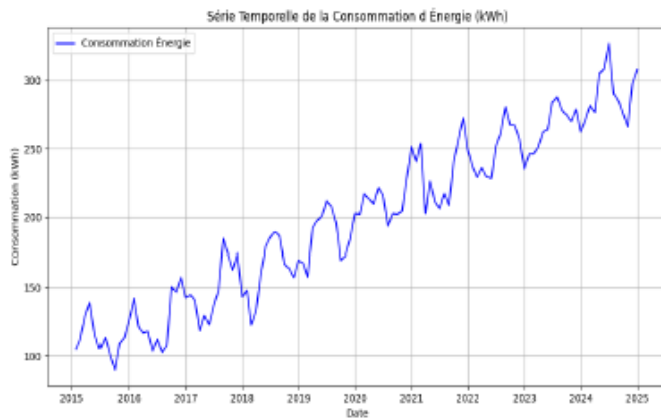
Transformations pour stationnarité

- Transformation pour non-stationnarité :
 - En moyenne (tendance) : différenciation simple
 - $Y_t - Y_{t-1}$ ← Supprime les tendances linéaires
 - De saisonnalité : différenciation saisonnière
 - $Y_t - T_{t-s}$ ← Supprime les cycles saisonniers
 - En variance : trans. logarithmique
 - $Y'_t = \log Y_t$

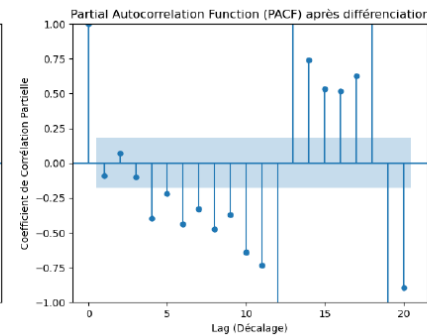
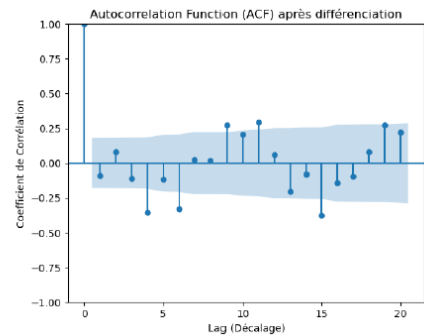
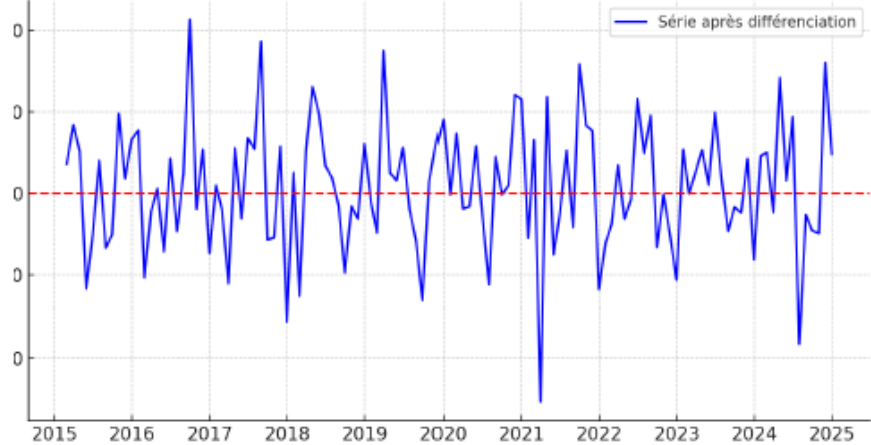
Plusieurs transformations à la suite peuvent être nécessaires



Transformations pour stationnarité



Série après différenciation simple (Suppression de la tendance linéaire)



Méthode	Statistique	p -value	Stats. (diff.)	p -value (diff.)
ADF	1.0018	0.9942	-9.6598	0
KPSS	1.79077	0.01	0.0123	0.1

Prédiction de séries temporelles

- Série stationnaire : **AR, MA, ARMA** Univariées
- Série non-stationnaire (avant) : **ARIMA, SARIMA**
- Critères d'information : ← Mesure qualité d'un modèle
 - AIC (*Akaike info. criterion*) : $AIC = -2 \log L + 2k$
 - BIC (*Bayesian info. criterion*) : $BIC = -\log L + k \log N$

Prévision court terme quand infos corrélées dans le temps

Un bon modèle doit générer des résidus stationnaires : sans tendance ni saisonnalité restante

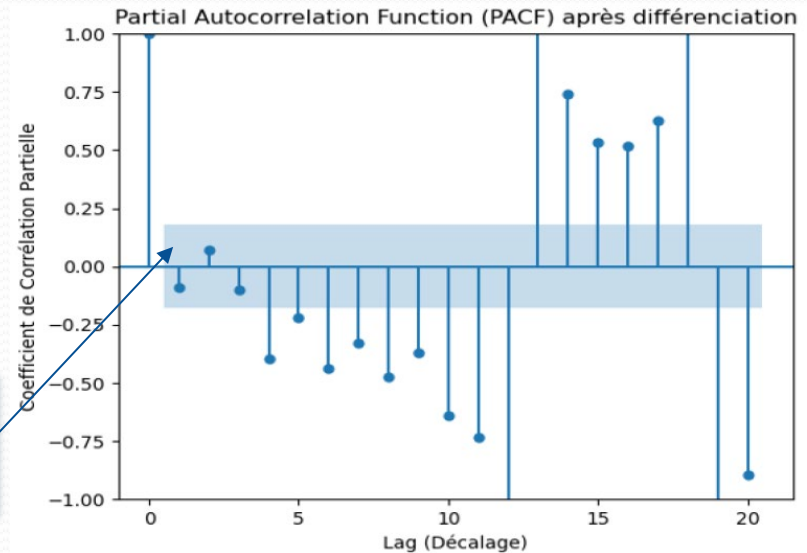
Les résidus doivent ressembler à un bruit blanc : aléatoire et structure détectable

Modèle AR

- **Autorégressif :**

- Prédire futures valeurs avec valeurs passées
- $AR(p)$: valeur actuelle dépend de p valeurs (T) passées
- $Y_t = c + \phi_1 Y_{t-1} + \dots + \phi_p Y_{t-p} + \epsilon_t$
- Choix de p : PACF
- Estimation des ϕ_p :
 - OLS, Lasso, Ridge

Avec $AR(p)$ la PACF coupe nettement après le *lag* p (ici : $p = 1$)

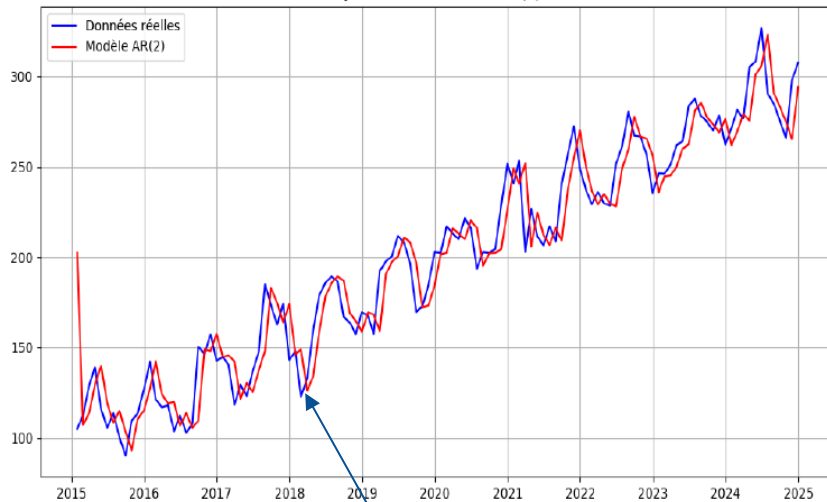


AR ne prend pas en compte les effets saisonniers, sensible aux bruits et données manquantes

Modèle AR

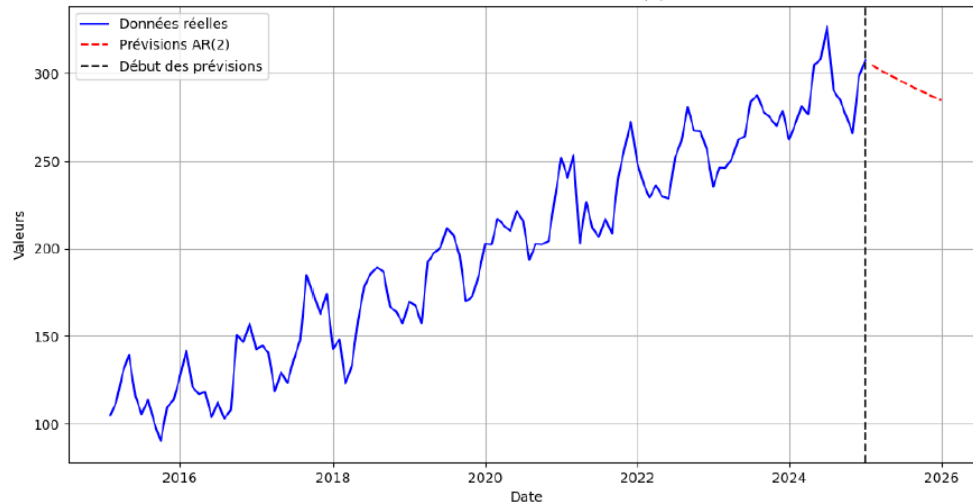
Série temporelle transformée avec différenciation

Ajustement du modèle AR(2)



AR(2) est un bon modèle car
courbe rouge suit la bleue

Prévisions sur 12 mois avec AR(2)



Évaluation du modèle AR(2) :

AIC : 1011.81

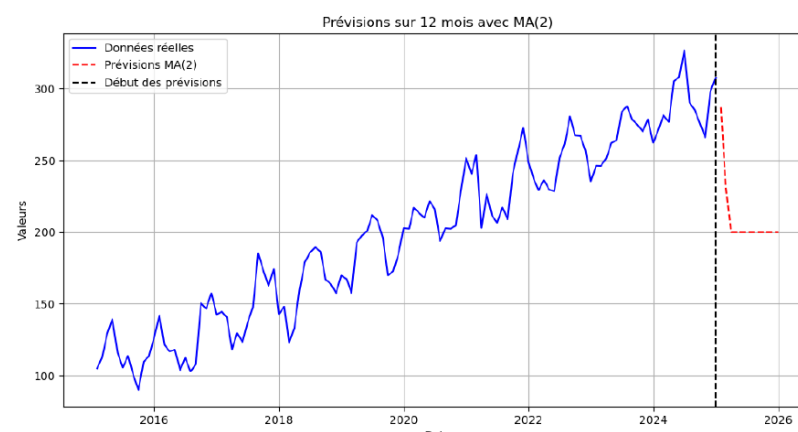
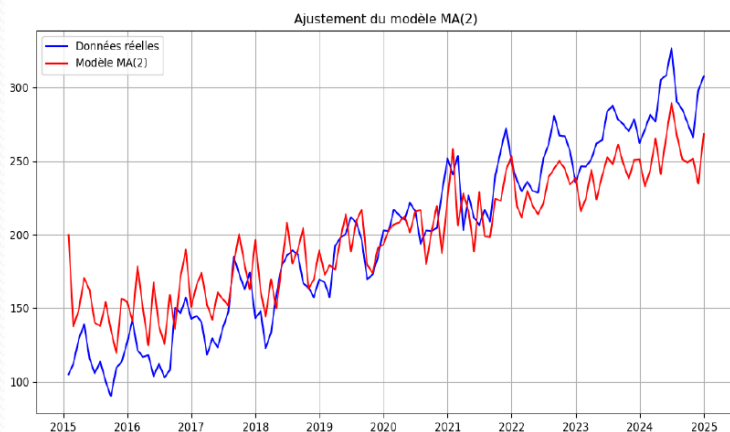
BIC : 1022.96

MAE (Erreur absolue moyenne) : 13.07

RMSE (Erreur quadratique moyenne) : 17.90

Modèle MA

- **Moyenne mobile** ← **≠ outil de lissage**
 - Prédire futures valeurs avec erreurs (R) passés
 - $MA(q) : Y_t = c + \epsilon_t + \theta_1\epsilon_{t-1} + \dots + \theta_q\epsilon_{t-q}$
 - Choix de q : ACF



Métrique	Valeur
AIC	1144.25706
BIC	1155.407027
MAE	23.116971
RMSE	28.275765

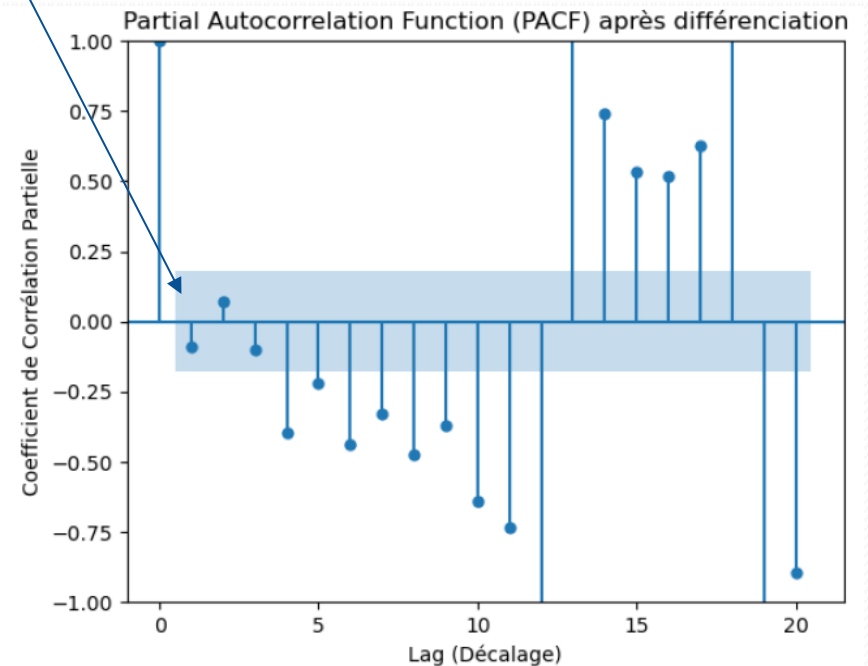
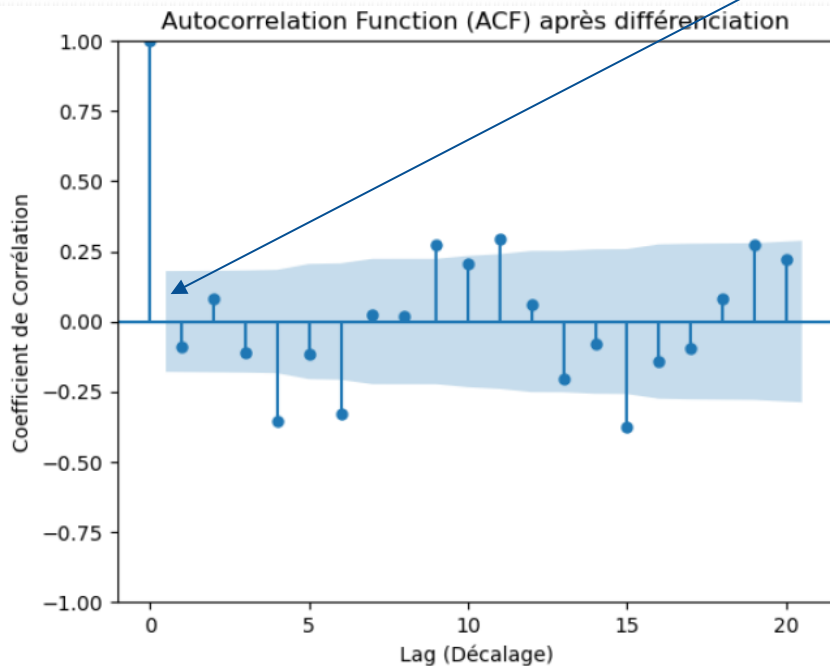
Stationnarité des résidus Non

Prédit des valeurs qui convergent rapidement vers une moyenne constante

Modèle ARMA

- $AR(p) + MA(q)$

$q=1, p=1$



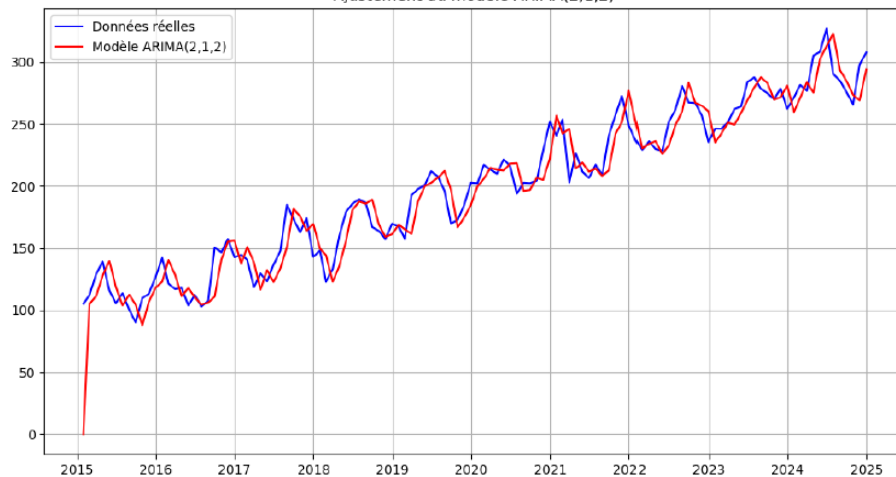
Modèle ARIMA

- *Autoregressive integrated moving average*
- $AR(p) + I(d) + MA(q)$
- **Intégration** : Saisonnalité
 - $I(d)$: d = nb de différenciations pour rendre stationnaire
 - Choix de d : ADF, KPSS

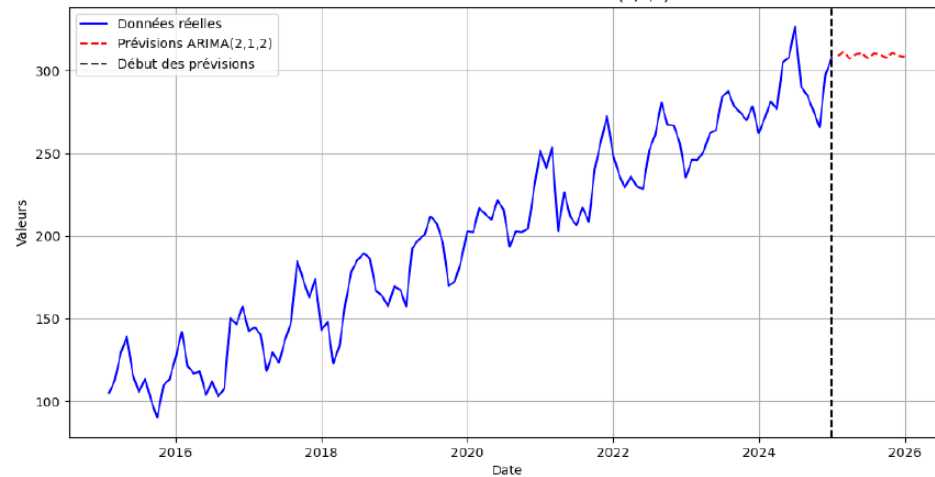


Modèle ARIMA

Ajustement du modèle ARIMA(2,1,2)



Prévisions sur 12 mois avec ARIMA(2,1,2)



Métrique	Valeur
AIC	996.18597
BIC	1010.081588
MAE	12.931678
RMSE	17.95582
Stationnarité des résidus	oui

Comparaison des modèles

Modèle	AIC	BIC	MAE	RMSE	Stationnarité
<i>AR</i> (2)	1011.81	1022.96	13.07	17.9	Non-stationnaire
<i>MA</i> (2)	1144.26	1115.41	23.12	28.28	Non-stationnaire
<i>ARIMA</i> (2, 1, 2)	996.19	1010.08	12.93	17.96	Stationnaire

ARIMA(2, 1, 2) est le modèle le plus performant

Choix d'un modèle de prédiction

- Minimiser AIC, BIC
- Méthode Box-Jenkins :
 - Protocole statistique de modélisation, d'analyse et d'estimation de modèle

Phase I : Identification

- Stabiliser la variance
- Stabiliser la moyenne

- Examiner les données, AVF, PACF
- Identifier les modèles potentiels

Phase II : Estimation et test

- Estimer les paramètres des modèles potentiels
- Choisir le meilleur modèle selon un critère (AIC ou BIC)

- Examiner les données, ACF, PACF et identifier les modèles potentiels
- Tester la stationnarité des erreurs résiduels

Phase III : Application

- Faire le forecasting avec le modèle retenu

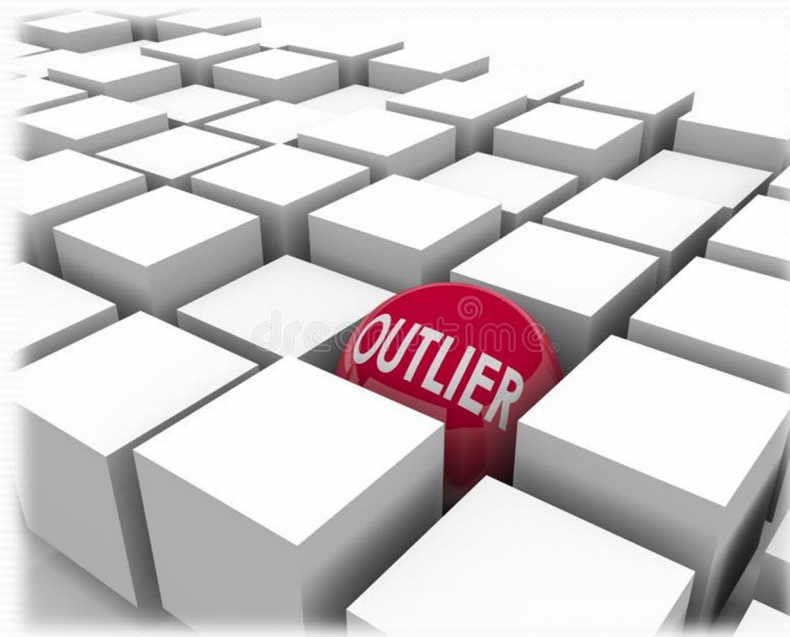
Projet de contributions

- *Outliers et skylines* :

Skylines sensibles aux outliers

- *Outliers* \Rightarrow *skyline** \Rightarrow classement avec méthode de SkyRank
- *Outliers* + *skyline* + SkyRank \leftarrow Gain d'efficacité
- *Outlier cube* : toutes les combinaisons de dimensions

Connaître dimension(s)
à l'origine d'un *outlier*



Aller plus loin

- LCOF (*local correlation outlier factor*)
- Sélection des variables
- Isolation Forest
- SVM
- HQC (*Hannan-Quinn information criterion*)
- SARIMA
- XGBoost
- Random Forest
- LSTM
- Séries temporelles multivariées

Liens

- Document classique :
 - Sana Sellami. *Exploration de données*.
 - Ricco Rakotomalala. *Détection des anomalies*.
 - Breunig , Kriegel, Ng and Sander. *LOF: identifying density based local outliers*.

Crédits

Auteur

Mickaël Martin Nevot

mmartin.nevot@gmail.com



Carte de visite électronique

Relecteurs

Cours en ligne sur : www.mickael-martin-nevot.com

