

## NUAGE DE POINTS :

- bien qu'imparfait, on constate une légère relation linéaire croissante entre les deux variables
- un modèle linéaire estimé par les MCO pourrait donc être approprié
- on peut s'attendre à un  $\beta$  positif mais inférieur à 1 et à un  $\alpha$  négatif

## CORRELATION :

- nous avons un coefficient de corrélation positif et relativement élevé
  - la relation linéaire entre les deux variables semble se confirmer
- ### MCO :
- si la valeur de t-stat est inférieure à la valeur tabulée (table de student) donc n'est pas significative
  - Nous obtenons un  $\beta$  de 1.5977, c'est à dire positif (ce que nous attendions). Il est néanmoins supérieur à 1.
  - Ainsi, une variation de 1% de la variable ersand (variable explicative) induira une variation de 1.5977% de erfod. (variable expliquée)

- La constante n'est pas négative, ce à quoi nous nous attendions à partir du nuage de points. Cela signifie que s'il n'y a aucune variation du marché, le rendement des actions de Ford sera quand même légèrement positif.

- La p-valeur associée au  $\beta$  est égale à 0.000 : on peut dire que la variable ersand est statistiquement significative au seuil de 5%
- l'hypothèse de non significativité des paramètres est :

- non-rejetée pour  $\alpha$  : -rejetée pour  $\beta$
- Significativité globale du modèle :
- le  $R^2$  ajusté vaut 0.287 : donc environ 30% des variations de erfod sont expliquées par le modèle (ce qui est tout à fait raisonnable).

## MCO « HAC » :

- on remarque que l'option "HAC" permet de mener une correction de la matrice de variance-covariance de l'erreur à la Newey-West -cette correction semble avoir modifié les écarts-types des estimateurs (à la hausse ce qui est normal puisque l'hétéroscédasticité et l'autocorrélation ont tendance à mener à sous-estimer la variance des estimateurs).
- on pourra rappeler que ces deux problèmes sont un soucis dans la mesure où ils entraînent un estimateur erroné de la matrice de variance-covariance des erreurs, ce qui fausse l'inférence statistique (c'est à dire les tests de significativité).

## RESIDUS :

- Les résidus semblent être de moyenne constante et égale à 00
  - En revanche leur variance n'est pas constante
  - Difficile de dire s'il existe une structure d'autocorrélation.
- ### Breusch-Pagan / white (nuage de 4 éléments) :
- 1er élément : Statistique du test LM (se référer à la table)
  - 2ème : p-valeur associée à la statistique de test LM (inférieure à 0.1 (white : 0.01), on rejette l'hypothèse nulle d'homoscédasticité à 10%)
  - 3ème : Statistique du test de Fischer (significativité globale du modèle de régression des résidus sur X)
  - 4ème : p-valeur associée à ce test (même interprétation puisque si la p-valeur < 0.01, le modèle est globalement significatif, et donc la relation entre les résidus et les variables explicatives est significative)

Ici nous rejetons l'hypothèse nulle d'homoscédasticité à 10%

## 2-VIF :

- 1 explicatives = [ersand, roif, d\_crisis, d\_covid]
- for n, var in enumerated(1 explicatives):  
print(le VIF de la variable', var, ' est :',  
'\nVIF[d', explicatives].dropna().values,n))  
le VIF de la variable erfod est : 1.0591829579271606  
le VIF de la variable roif est : 1.0473746769468446  
le VIF de la variable d\_crisis est : 1.029664691134221  
le VIF de la variable d\_covid est : 1.003964217666122

## Commentaires :

- Nous avons calculé les Variance Inflation Factor de toutes les variables explicatives
  - Sans avoir de règle de décision précise, nous pouvons dire qu'avec des VIF inférieur à 2 pour toutes les variables, la présence de multicollinéarité est faible.
  - L'hypothèse MCO selon laquelle la matrice des variables explicatives  $X \cdot X'$  est de plein rang, semble respecter (hypothèse qui assure de pouvoir estimer les paramètres).
  - import numpy as np  
from numpy.linalg import det  
D = det(dif1 explicatives.com(0))  
# calcul de la T : T=244 et k=4  
FG = (-244 / (1/6) \* (2\*(4+1)) + 5)\*np.log(D)  
print(FG)  
**Commentaires :**
    - Nous menons ici un test de multicollinéarité de Farrar-Glober
    - Avec une valeur tabulée pour le  $\chi^2(1/2(k+1))$  de 18.307 à 5%, l'hypothèse nulle d'absence de multicollinéarité est rejetée.
    - Ce test confirme donc nos commentaires précédents.
  - échantillons liés pour voir l'absence de séries et éventuellement en déduire de l'information sur la stationnarité des séries**
    - autocorélogramme** : permet de voir la structure autocorrégressive (de long terme) des séries
    - autocorrélations partielles** : permet de voir l'autocorrélation entre deux périodes après avoir purgé l'effet des autres périodes
  - ADE** : si la statistique est supérieure à la valeur tabulée nous ne sommes pas H0 la série contient une racine unitaire. La série est non stationnaire de type stochastique ; il faut la différencier (on prend le taux de croissance pour la stationnariser. Nous devons, bien entendu vérifier que la variable ainsi transformée est elle-même stationnaire. rford est donc stationnaire et la série rford est dite intégrée d'ordre 1 puisqu'il a fallu la différencier une fois pour la stationnariser. D'où on travaillera sur la série de rendement Si p-val est inférieure à 1% le test rejette H0 de présence d'une racine unitaire. La série
- ### Exercice 1 :
- Question 1 : A l'aide de quel(s) statistique(s) peut-on comparer la significativité globale de 2 modèles estimés par la même méthode ?
- Réponse : Pour comparer la significativité globale de deux modèles estimés par la même méthode, plusieurs statistiques peuvent être utilisées en fonction du type de modèle utilisé : test  $\chi^2$ , test de Fisher, le critère d'information d'Akaike (AIC) ou le critère d'information bayésien (BIC).
- Question 2 : L'estimateur MCO est-il « BLUE » ? Que signifie cet acronyme et à quel renvoie-t-il ?
- Réponse : BLUE : best linear unbiased estimator. Propriété1 « Linéarité de l'estimateur ». Propriété2 « Estimateur dans biais ». Propriété3 «  $\beta$  est de variance minimale ». Propriété totale : «  $\beta$  est un estimateur BLUE de  $\beta$  ». On dit de l'estimateur MCO «  $\beta$  est BLUE (Best Linear Unbiased Estimator).
- Question 3 : Quelle propriété des variables est-il impératif d'étudier lorsque l'on analyse des séries temporelles ? Pourquoi ?
- Réponse : Lors de l'analyse de séries temporelles, il est crucial d'étudier la stationnarité des variables. La stationnarité, qui mesure que les propriétés statistiques de la série restent constantes au fil du temps. Cette propriété est essentielle pour garantir la validité des résultats d'analyse, assurer la stabilité des modèles, faciliter l'interprétation des résultats et améliorer la fiabilité des prévisions à long terme. Des méthodes graphiques et des tests statistiques, tel que le test Augmented Dickey-Fuller (ADF) est utilisé pour évaluer la stationnarité. En cas de non-stationnarité, des transformations ou ajustements peuvent être nécessaires avant d'appliquer des modèles d'analyse.
- **Modèle de régression simple** :  $Y = \alpha + \beta X + \epsilon$
- Y est la variable expliquée (observée)
  - X est la variable explicative (observée)

- $\alpha$  (la constante) et  $\beta$  (le coefficient de pente) sont les paramètres à estimer (inconnus)
- c'est le terme d'erreur (inconnu)
- Dans un premier temps, on considère que X est observé sans erreur, autrement dit, X est une variable certaine. De ce fait, elle est indépendante du terme d'erreur  $\epsilon$ . La variable Y est elle en revanche, une variable aléatoire (puisque'elle dépend précisément d' $\epsilon$ ).
- Si les variables X et Y sont des variables qui dépendent de T observations, on écrit le modèle de régression simple comme suit :  $Y_t = \alpha + \beta X_t + \epsilon_t$
- $\forall t \in [1; T]$  où T peut désigner deux dimensions différentes :
  - le temps (séries temporelles)
  - des individus (coupes instantanées)
- **propriété du terme d'erreur** : On ne peut pas mesurer ou prévoir le terme d'erreur  $\epsilon$  pour chaque observation. On émet un certain nombre d'hypothèses sur ce terme...
  - 1- **Nullité de l'erreur moyenne** :  
 $E(\epsilon) = 0$  ou  $E(\epsilon) = 0$ , mais sans plus de chance d'avoir l'un ou l'autre  
→  $E(\epsilon) = 0$   
-Revient à dire que, en moyenne, le modèle est bien spécifié
  - 2- **L'absence d'autocorrélation des erreurs** :  
-La valeur de l'erreur en t (période ou individu) ne dépend pas de celle en t'  
-Il n'y a pas de corrélation entre deux erreurs à deux dates différentes ou entre deux individus  
→  $E(\epsilon_t \epsilon_{t'}) = 0, \forall t \neq t'$
  - 3- **Homoscédasticité des erreurs**  
-La variance de l'erreur est constante (au cours du temps ou entre les individus)  
→  $E(\epsilon_t^2) = \sigma^2, \forall t$   
→  $\sigma^2$  représente la variance de l'erreur et ne dépend pas de t
  - 4- **Normalité des erreurs**  
-On suppose (théorème central limite) que le terme d'erreur suit une loi normale centrée et de variance constante  
-Cette hypothèse se vérifie d'autant plus que le nombre d'observations est grand  
→  $\epsilon \sim N(0, \sigma^2 \epsilon)$
- **Principe des MCO** :  $Y = T \cdot \alpha + \beta X_t + \epsilon_t$   
L'objectif des MCO est d'estimer les paramètres  $\alpha$  et  $\beta$ . Une telle estimation nous donne l'équation d'une droite, appelée **droite de régression**, donnée par :  
 $Y = T \cdot \alpha + \beta X_t$   
où :
  - $\alpha$  et  $\beta$  sont les estimateurs des paramètres  $\alpha$  et  $\beta$
  - $Y$  est la valeur estimée de Yt et donné par le modèle (une valeur de Yt\* correspond à un point sur l'axe des ordonnées (Yt) pour un Xt et obtenu via la droite de régression MCO).
- **Propriétés des estimateurs** :
  - Propriété 1** : la droite de régression passe par le point moyen  $(\bar{X}, \bar{Y})$
  - Propriété 2** : la variable observée Yt et estimée Yt\* sont de même moyenne :  $Y = Y^*$
  - Propriété 3** : en moyenne, les résidus sont nuls (modèle en moyenne correctement spécifié)  $e = 0$
  - Propriété 4** : la covariance entre la variable explicative X et les résidus, ainsi qu'entre la variable expliquée estimée Yt\* et les résidus, est nulle (pas de corrélation avec le résidu :  $Cov(X, \epsilon) = Cov(Y, \epsilon) = 0$ )
  - Propriété 5** : les changements d'origine et d'échelle ne changent pas  $\beta$
  - Propriété 6** : les estimateurs sont linéaires (on peut les exprimer comme des fonctions linéaires de la variable expliquée)
- Très important - Propriété 7** : ce sont des estimateurs sans biais (en moyenne, ils sont égaux à leur vraies valeurs)  $E(\alpha) = \alpha$  et  $E(\beta) = \beta$
- Propriété 8** : les estimateurs MCO sont convergents et de variance minimale. Autrement dit, lorsque le nombre d'observations T tend vers l'infini, la variance converge vers 0.

- les estimateurs MCO sont BLUE : Best Linear Unbiased Estimator
- **Estimateur MCO : Hypo du modèle de régression multiple** :  
 $Y_t = \alpha + \beta_1 X_{1t} + \beta_2 X_{2t} + \dots + \beta_k X_{kt} + \epsilon_t$
- H1 : la matrice X est non aléatoire** : les valeurs de X sont observées sans erreur (X et le terme d'erreur sont indépendants.  $E(\epsilon|X) = 0$ )
- H2 : la matrice X est de plein rang** : Rang(X) = k+1. Cette hypothèse revient à dire que les colonnes sont linéairement indépendantes, ou que l'on ne peut pas écrire une colonne de X comme une combinaison linéaire des autres. Nous faisons cette hypothèse pour assurer l'inversibilité de la matrice  $X'X$  (dont nous avons besoin pour le calcul de  $\beta$ ). Une hypothèse sous-jacente est que le nombre d'observations ne peut pas être inférieur au nombre de variables : T > k+1. Sans cela, X ne peut pas être de plein rang.
- H3 : nullité de l'espérance de l'erreur** : Comme pour le modèle de régression simple, nous supposons que le terme d'erreur est nul en moyenne :  $E(\epsilon) = 0$ . On en déduit que  $E(Y) = X\beta$  c'est à dire que le modèle est en moyenne bien spécifié.
- H4 : homoscedasticité et absence d'autocorrélation des erreurs** : Soit  $E(\epsilon \epsilon')$  la matrice de variance covariance du terme d'erreur, avec  $\epsilon'$  la transposée du vecteur des termes d'erreur. L'hypothèse H4 s'écrit comme suit :  $E(\epsilon \epsilon') = \sigma^2 I$  où I désigne la matrice identité et  $\sigma^2$ , la variance du terme d'erreur.
- H5 : normalité des erreurs** : nous supposons (en vertu du théorème central limite) que les erreurs suivent une loi normale d'espérance nulle et de variance constante :  $\epsilon \sim N(0, \sigma^2 \epsilon)$  A noter que cette hypothèse ne sert pas au calcul des paramètres mais à la validation des résultats statistiques et à la construction des tests statistiques.
- **Estimateur de la variance des erreurs**  
On ne connaît pas les erreurs dans un modèle. Pour estimer la variance des erreurs  $\sigma^2$ , il convient d'utiliser les résidus  $e$  :  
 $\sigma^2 = e'e / (T - k - 1)$   $\sigma^2$  est un estimateur sans biais de la variance des erreurs  $\sigma^2 \epsilon$
- **Test sur un coefficient de régression : test de student** : A partir de la statistique t définie à la section précédente, nous pouvons déduire un intervalle de confiance à 100(1 - p)% pour  $\beta$  :  
 $\beta \pm t_{p/2, n-k-1} \sqrt{\sigma^2 \epsilon}$

- On peut donc tester si  $\beta$  s'écarte plus ou moins une certaine valeur  $\beta_0$  : H0, l'hypothèse nulle :  $\beta = \beta_0$  - H1, l'hypothèse alternative :  $\beta \neq \beta_0$
- en pratique, on utilise ce test comme un test de significativité, c'est à dire qu'on prend  $\beta_0 = 0$  : H0, l'hypothèse nulle :  $\beta = 0$  - H1, l'hypothèse alternative :  $\beta \neq 0$
  - **Test sur plusieurs coefficients : test de Fisher** : le test de Student nous permet de tester la significativité statistique d'une variable seule. Nous procédons maintenant à la description du test de Fisher, qui permet de tester la significativité de plusieurs coefficients simultanément. Idée générale du test : poser une contrainte sur plusieurs coefficients en même temps,  $R\beta = 0$ . Où R est une matrice de plein rang et de taille (k,1) avec k le nbr de contraintes et l le vecteur des contraintes de taille (q,1)
  - **Equation de l'analyse de la variance** :  $Y = X\beta + \epsilon$  nous avons obtenu, après estimation, une expression pour le résidu :  $\epsilon = Y - X\hat{\beta}$ , où l'on note X et Y en minuscule pour signifier qu'elles sont centrées (on leur retire leur moyenne).
  - **coefficient de détermination (R²)** : mesure la part de la variance totale expliquée par les variables explicatives et permet de juger de la qualité de l'ajustement du modèle.  $R^2 = SCE/SCT = 1 - SCR/SCT$  où  $0 \leq R^2 \leq 1$ . → **coefficient de détermination corrigé (ajusté R²)** : le coefficient de détermination classique augmente avec le nombre de variables explicatives intégrées dans le modèle (puisque l'introduction de nouvelles variables ne peut pas faire diminuer la SCR) on peut donc le corriger de ce biais :  $R^2 = 1 - ((T-1)/(T-k-1)) * (1-R^2)$

→ **Critère d'information** : critère d'information d'Akaike (AIC) : plus ce critère est faible plus le modèle a un pouvoir explicatif important. Critère d'information de Schwarz (SIC) : ce critère est aussi un critère qu'il faut minimiser.

→ **Prévision** : estimation d'un modèles : - permet d'expliquer ce qui est observé : pouvoir explicatif (en échantillon) - permet d'émettre des prévisions : pouvoir prédictif du modèle (hors échantillon) : on estime le modèle sur une partie de l'échantillon, puis on regarde l'erreur de prévision sur la partie restante de l'échantillon).

→ **Propriété des estimateurs en présence d'autocorrélation et/ou hétéroscédasticité** : modèle :  $Y = X\beta + \epsilon$  où X est non aléatoire de plein rang. L'estimateur des MCO donne  $\hat{\beta} = (X'X)^{-1}X'Y$  : en présence d'hétéroscédasticité ou d'autocorrélation, la variance de l'estimateur MCO n'est plus minimale. Les MCO ne sont alors plus "Best".

→ **Hétéroscédasticité des erreurs** : On rappelle que l'hétéroscédasticité renvoie au fait que les termes situés sur la diagonale de la matrice de variance-covariance des erreurs sont non-égaux entre eux. En série temporelle : la variance des erreurs est non constante dans le temps.

→ **D'où peut venir l'hétéroscédasticité ?**

- d'hétérogénéité** de l'échantillon (pays différents ou époques différentes)
- oubli d'une variable explicative dans le modèle
- asymétrie dans la distribution des certaines variables
- une mauvaise transformation des variables (par exemple : prendre le modèle linéaire alors que le bon modèle aurait été le modèle log-linéaire).

→ **Tests d'hétéroscédasticité : 1-test graphique : pour détecter**

- il existe divers tests permettant de la détecter
- première intuition graphique est toujours utile :
  - estimation du modèle comme s'il n'y avait pas d'hétéroscédasticité
  - graphique (nuage de points) de la variable expliquée estimée (Y) en fonction des résidus au carré
  - s'il existe une relation apparente entre les deux séries, nous pouvons avoir présomption de présence d'hétéroscédasticité

-il est aussi possible de faire le nuage de points entre les résidus et une variable explicative : de même si une relation apparaît, nous avons présomption d'hétéroscédasticité.

-une autre analyse graphique consiste à faire le graphique de l'évolution des résidus selon les observations : la variance doit apparaître constante.

**2-quelques tests d'hétéroscédasticité** : Le test de Goldfeld-Quandt, développé en 1965, sert à détecter l'hétéroscédasticité due à l'une des variables explicatives dans une régression linéaire. Le processus implique le classement des valeurs de la variable suspecte par ordre croissant, la suppression des valeurs centrales, l'estimation par les MCO sur les deux sous-échantillons résultants, le calcul des sommes des carrés des résidus (SCR), et enfin la comparaison des SCR. Une différence significative suggère la présence d'hétéroscédasticité dans le modèle. Test de Glejser (1969) : a pour but de détecter la présence ET la forme de l'hétéroscédasticité.

**3-test de Breusch-Pagan (1979)** : -test général couvrant un grand nombre de cas d'hétéroscédasticité

- valable que pour des échantillons assez grands (beaucoup d'observation) : test asymptotique
- n'implique pas de savoir quelle variable explicative est à l'origine de l'hétéroscédasticité
- test simple : ne se base que sur les résidus issus de la régression MCO du modèle

On considère le modèle de régression multiple suivant :  
 $Y_t = \alpha + \beta_1 X_{1t} + \beta_2 X_{2t} + \dots + \beta_k X_{kt} + \epsilon_t$

- $BP \sim \chi^2_p$ , l'hypothèse nulle d'homoscédasticité n'est pas rejetée

- $BP \sim \chi^2_p$  l'hypothèse nulle d'homoscédasticité est rejetée

**4-test WHITE** : ne repose pas sur l'hypothèse de normalité des erreurs

-ne s'intéresse pas à la forme de l'hétéroscédasticité.

- si  $TR^2 \sim \chi^2_k$ , on ne rejette pas l'hypothèse nulle d'homoscédasticité
- si  $TR^2 \sim \chi^2_k$ , on rejette l'hypothèse nulle d'homoscédasticité

→ **Procédure d'estimation en présence d'hétéro** : Problème : cette variance est utilisée dans le calcul des test de significativité. Solution : -estimation du modèle par les MCG, il faut connaître la matrice de variance covariance du terme d'erreur.

-Correction des MCO par l'utilisation d'un estimateur de la matrice de var-cov.

**1-estimateur de la matrice de var-cov de white :**

$\Omega = e'(e'(T-k-1)) * (X'X)^{-1} * (\sum e_t e_t X_t X_t' / (T-k-1))$

**2-estimateur de la matrice de var-cov de Newey-West** : -La proposition de correction de White repose sur l'hypothèse de non-autocorrélation des résidus.

-La proposition faite par Newey et West (1987) non

→ **Autocorrélation des erreurs** : nous travaillons sur l'autocorrélation en série temporelle, le terme d'erreur à un instant, dépend du même terme d'erreur à un autre instant/ en coupe instantanée : autocorrélation spatiale.

→ **Tests d'autocorrélation : 1-détection graphique de la présence d'autocorrélation** : de la même manière que pour l'hétéroscédasticité, il est possible de repérer la présence d'autocorrélation dans les erreurs en nous intéressant à la structure des résidus, de manière graphique

**2-test de Durbin-Watson** : Permet de tester l'autocorrélation à l'ordre 1 (en série temporelle, cela l'autocorrélation à l'ordre 1 renvoie à une autocorrélation du terme d'erreur entre 2 périodes). Test basé sur une hypothèse de distribution des résidus

**3-test de Breusch-Godfrey** : permet de détecter d'autocorrélation à des ordres supérieurs à 1. Le test se base sur le processus autorégressif d'ordre p (AR(p)).  $BG \sim \chi^2_p$ , l'hypothèse nulle d'absence d'autocorrélation n'est pas rejetée.  $BG \sim \chi^2_p$ , l'hypothèse nulle d'absence d'autocorrélation est rejetée

→ **Procédure d'estimation en présence d'autocorrélation** :

- Même remarques que pour l'hétéroscédasticité :
  - autocorrélation  $\Rightarrow$  variance de l'estimateur non minimale
  - variance non minimale  $\Rightarrow$  tests de significativité biaisés
  - modèle non interprétable.

-Il est donc important de corriger de l'autocorrélation

-Pour effectuer cette correction : même problème que pour l'hétéroscédasticité concernant les MCG.

-On préférera donc procéder à une correction du modèle en utilisant un estimateur de la matrice de var-cov des erreurs.

-Nous avons vu que la correction de Newey-West permettrait est aussi adaptée à la correction de modèles en présence d'autocorrélation : c'est donc la correction que nous utiliserons.

→ **Problème lié aux variables explicatives : variables aléatoires** : l'hypothèse,  $X \cdot X'$  est non aléatoire :  $E(\epsilon|X) = 0$

renvoie à l'hypothèse d'exogénéité des variables explicatives : l'espérance conditionnelle de l'erreur étant nulle, les variables explicatives ne sont pas déterminées par le modèle et lui sont donc exogènes.

-Lorsque cette hypothèse est non vérifiée : même si la taille de l'échantillon croît indéfiniment, les estimateurs MCO ne se rapprochent plus de la vraie valeur.

-Problème : trouver un estimateur convergent

-l'idée :

- variables explicatives ne sont pas réellement "exogènes"
- solutions : les remplacer par des variables qui le sont
- enjeux : trouver des variables exogènes, représentative de ce que l'on souhaite mesurer

-Pour détecter le problème : étude de la corrélation entre les variables explicatives et le résidu de la régression

-origine du problème : mauvaise mesure des variables ou manque d'une variable explicative importante dans le modèle ou la variable expliquée est aussi variable explicative du modèle.

→ **Si la corrélation est forte** : -Passer le modèle en dynamique (avec des retards). Vérifier si une variable explicative importante ne manque pas. Etudier les sens de causalité de la relation.

→ **Problème lié aux variables explicatives : multicollinéarité** : Nous avons étudié l'hypothèse selon laquelle X est une matrice non aléatoire. Nous nous intéressons maintenant, plus en détails à l'hypothèse qui établit que X est de plein rang.

→ **Détection multicollinéarité-VIF (variance inflation factor)** :

-Cette statistique est appelée facteur d'influence de la variance (Variance Inflation Factor).

-Des valeurs élevées pour cette statistique pour une variable explicative donnée est le signe de la présence éventuelle d'une situation de multicollinéarité presque parfaite l'impliquant.

-Il n'existe cependant pas une règle établissant un seuil pour cette statistique au-delà duquel on peut prendre une décision. La statistique est néanmoins utile car donnant des indications sur les paramètres  $\beta$  pouvant être estimés avec peu de précision.

-Parfois considéré que :  $VIF < 2$  est parfait, et qu'à partir de 10, la multicollinéarité est trop forte.

→ **Test de FG : corrélation entre variables explicatives** :

-Test basé sur la matrice des coefficients de corrélation entre les variables explicatives.

-Si la matrice est linéairement corrélée, le déterminant de cette variable est nul

-La statistique de test est donnée par :

$FG = (T-1) - 1/6((k+1)*5) * \log(D)$

où D est le déterminant de la matrice de corrélation des variables explicatives.

-Sous l'hypothèse nulle d'orthogonalité (non-multicollinéarité) des variables explicatives :

$FG \sim \chi^2_{1/2(k+1)}$

-La règle de décision est donnée par :

- Si  $FG \sim \chi^2_{1/2(k+1)}$ , l'hypothèse nulle n'est pas rejetée (pas de multicollinéarité)
- Si  $FG \sim \chi^2_{1/2(k+1)}$ , l'hypothèse nulle est rejetée (présomption de multicollinéarité)

→ **Solutions à la multicollinéarité** : - Augmenter la taille de l'échantillon : augmenter le nombre d'observations pour diminuer la part des observations corrélées entre elles. Plus T est supérieur à k, moins la multicollinéarité a de chance d'être forte.

-La régression ridge : l'une des colonnes de X est égale à une combinaison linéaire des autres. En ajoutant une constante sur la diagonale de X, cette multicollinéarité disparaît.

→ **Test de stabilité des coefficients** : -En robustesse de l'estimation menée : tester la stabilité des coefficients (en série temporelle principalement).

-S'assurer que les coefficients ne varient pas selon le contexte

-Étudier l'impact d'un changement structurel.

-Parmis d'autres approches, nous étudions ici le test de Chow.

→ **3 types de structure de données** : -Données en coupe instantanée (plusieurs individus à un instant donné)

-Données en série temporelle (plusieurs périodes pour un individu donné)

-Données de panel (plusieurs individus sur plusieurs périodes)

→ **Pourquoi introduire les retards** : une variable que l'on cherche à expliquer a de grandes chances de dépendre de valeurs passées des variables explicatives : modèles à retards échelonnés

-la variable expliquée peut dépendre de ses propres retards : modèles autorégressifs

→ **Dans quel cadre utiliser les retards ?** -analyse financière, prévision autorégressive, prévision échelonnée, persistance ou effet structurel d'un choc.

→ **Nombre de retard et estimation** : -Regarder la significativité des coefficients -Estimer de plusieurs modèles avec plusieurs retards : on retient le modèle qui a les meilleurs critères

d'information (minimiser AIC, BIC, SIC) et le plus petit  $R^2$  ajusté.

→ **Modèles à retards échelonnés : PGL** :

$Y_t = \alpha + \beta_0 X_t + \beta_1 X_{t-1} + \beta_2 X_{t-2} + \dots + \beta_k X_{t-k} + \epsilon_t$

→ **Modèles à retards échelonnés autorégressifs** : Dans de nombreux cas (séries financières par exemple) : la variable endogène dépend de ses propres retards. On parle de processus autorégressifs :  $Y_t = \mu + \phi_1 Y_{t-1} + \phi_2 Y_{t-2} + \dots + \phi_p Y_{t-p} + \epsilon_t$

Le cas le plus connu étant celui de la marche aléatoire : la valeur d'aujourd'hui est la même que celle d'hier à un bruit prêt (finance comprise).

→ **Caractéristiques des séries temporelles :**

**Saisonnalité** : caractéristique des séries temporelles où les données décrivent des variations régulières et prévisibles.

On parle aussi de périodicité : tout schéma se répétant avec un intervalle de temps régulier. (Climat, ventes, finance).

→ **Stationnarité au second ordre** : Modélisation en série temporelle : étude obligatoire de la stationnarité !!

-Nous nous intéressons qu'à la stationnarité du second ordre / stationnarité faible : le plus courant

-Définitions de la stationnarité au second ordre d'un processus Yt :

- Le moment d'ordre 2 est fini et constant au cours du temps (homoscédasticité) :  $E(Y_t^2) < \infty, \forall t$
- La moyenne est constante :  $E(Y_t) = m, \forall t$
- La covariance entre t et t+h ne dépend que de h :  $Cov(Y_t, Y_{t+h}) = \gamma_h, \forall t$

→ **Fonction d'autocovariance, d'autocorrélation** : La fonction d'autocovariance est définie, pour un processus Yt de variance finie, par :

$\gamma_h = Cov(Y_t, Y_{t+h}) = E(Y_t - E(Y_t))(Y_{t+h} - E(Y_{t+h}))$

c'est la covariance d'une même variable entre deux dates. De là, nous pouvons définir la fonction d'autocorrélation :

$\rho_h = \gamma_h / \sigma^2$

elle mesure les liaisons temporelles d'une série. On se sert de la

