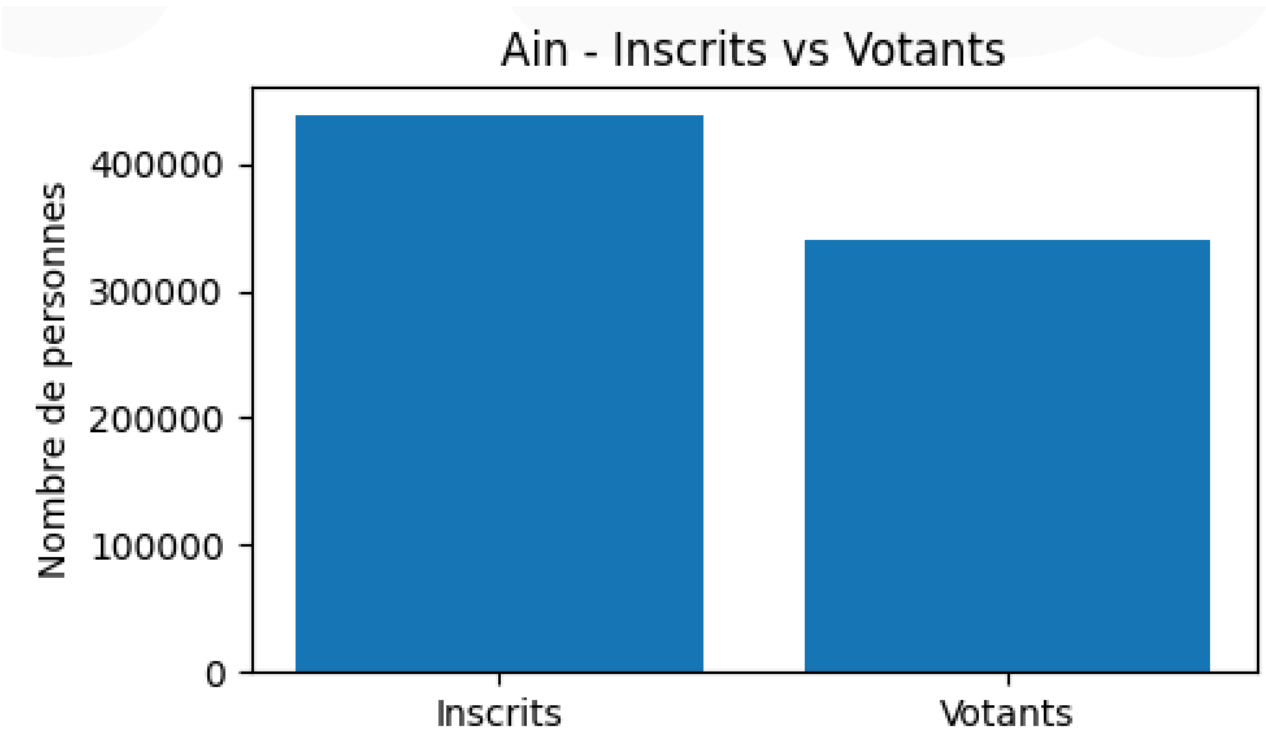


Séance 2

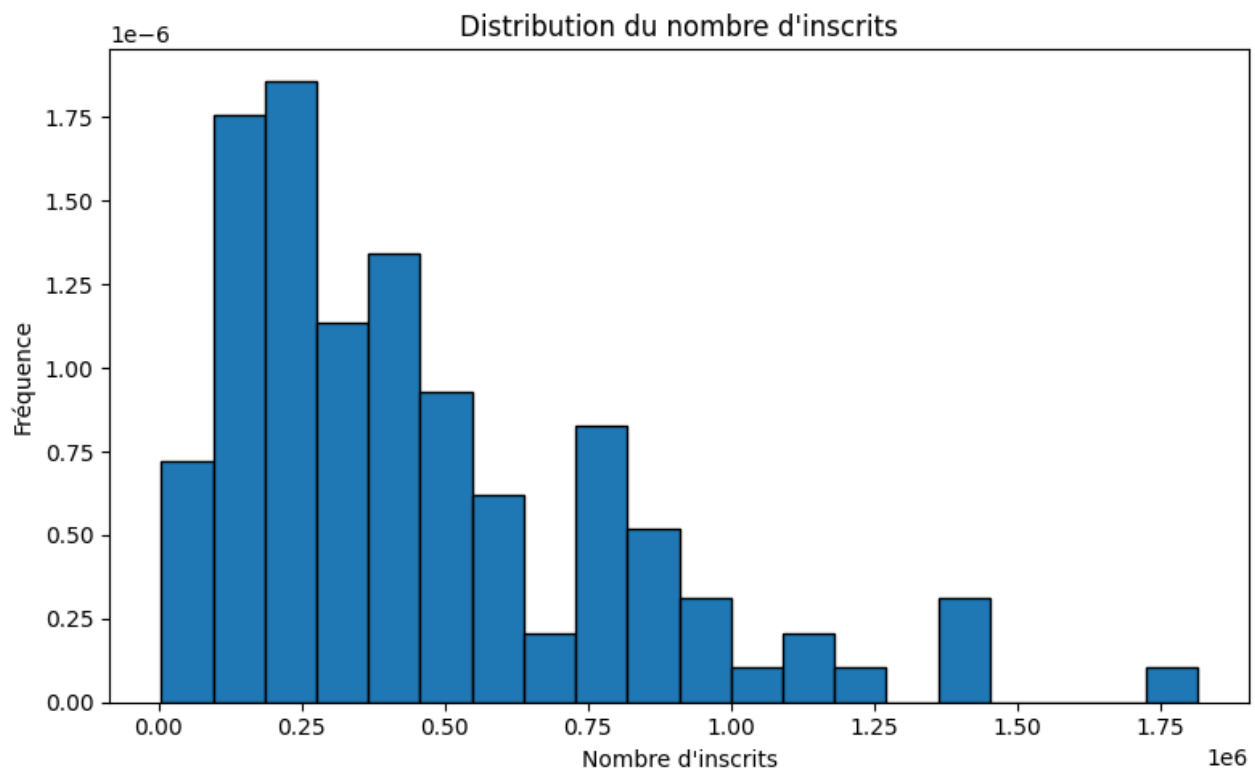
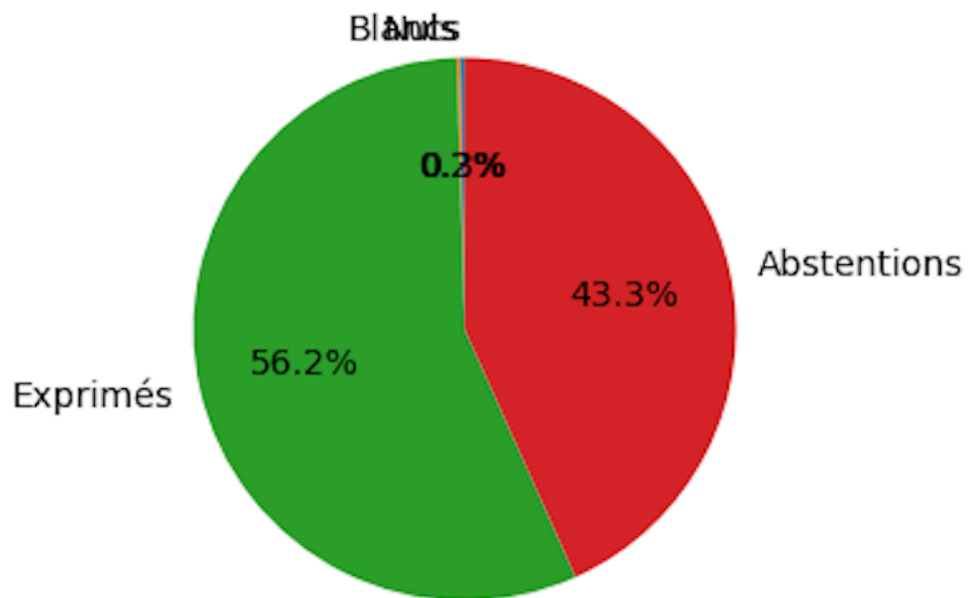
Caractéristique	Valeur	Commentaire
Nombre de Lignes	107	Correspond au nombre d'entités territoriales analysées (départements métropolitains, DROM-COM et Français de l'étranger).
Nombre de Colonnes	56	Représente les informations générales du scrutin et les résultats des 12 candidats (sexe, nom, prénom, et voix).

Type de	Exemples de Colonnes	Nature Statistique
object (str)	Libellé du département, Code du département, Nom,	Qualitatif (Nominal ou Ordinal)
int64	Inscrits	Quantitatif (Discret)
float64	Abstentions, Votants, Blancs, Nuls, Exprimés, Voix (pour les	Quantitatif (Discret, stocké en float à cause des valeurs manquantes ou pour des raisons de

Type de	Exemples de Colonnes	Nature Statistique
object (str)	Libellé du département, Code du département, Nom,	Qualitatif (Nominal ou Ordinal)
int64	Inscrits	Quantitatif (Discret)
float64	Abstentions, Votants, Blancs, Nuls, Exprimés, Voix (pour les	Quantitatif (Discret, stocké en float à cause des valeurs manquantes ou pour des raisons de



Wallis et Futuna - Répartition des votes



Commentaire de l'Histogramme :

L'histogramme révèle une distribution fortement asymétrique positive (ou étalée à droite).

1. **Concentration des Effectifs :** La majorité des entités territoriales (départements) ont un nombre d'inscrits faible à modéré. La classe de fréquence la plus élevée (le mode) se situe autour de 0.15×10^6 à 0.25×10^6 inscrits (soit 150 000 à 250 000 inscrits). Cela suggère qu'un

grand nombre d'entités étudiées sont des départements de taille moyenne ou des collectivités d'outre-mer.

2. **Asymétrie et Valeurs Extrêmes** : La queue de distribution s'étend loin vers la droite, indiquant la présence de quelques entités avec un nombre d'inscrits exceptionnellement élevé (des "outliers" ou valeurs extrêmes). Ces pics de faible fréquence autour de 1.25×10^6 à 1.75×10^6 correspondent probablement aux départements les plus peuplés de France (comme le Nord, les Bouches-du-Rhône, ou le département des Français établis hors de France).
3. **Conclusion sur la Distribution** : La distribution n'est pas Normale (Gaussienne). L'asymétrie positive est typique des variables socio-économiques et démographiques, indiquant que la moyenne arithmétique sera supérieure à la médiane (et au mode), car elle est tirée vers la droite par les grandes entités.

Séance 3

Le graphique fourni dans le deuxième fichier () est un Histogramme de fréquence illustrant la distribution du nombre d'inscrits par entité territoriale (départements et autres circonscriptions) issue du fichier des résultats des élections présidentielles de 2022.

Voici le commentaire de ce graphique, rédigé de manière fluide et humaine :

Analyse de la Distribution du Nombre d'Inscrits

Le graphique ci-dessus, qui représente la distribution du nombre d'inscrits par circonscription territoriale, montre clairement que les électeurs ne sont pas répartis de manière uniforme en France.

Une Distribution Asymétrique et Concentrée

- **Le Cœur des Électeurs (Mode)** : La distribution est très loin d'être symétrique. Elle est fortement étalée vers la droite (asymétrie positive). La grande majorité des entités se concentrent dans les classes de faible effectif. On constate que la fréquence la plus élevée se situe entre 100 000 et 250 000 inscrits (soit environ 0.1×10^6 à 0.25×10^6). Ceci représente le département "typique" ou moyen en termes de population électorale.
- **L'Effet des "Poids Lourds"** : La longue queue qui s'étend vers les valeurs élevées est très révélatrice. Elle signale la présence de quelques entités avec un nombre d'inscrits exceptionnellement important (des valeurs extrêmes). Ces pics isolés au-delà de 1 250 000 inscrits (1.25×10^6) correspondent aux départements les plus peuplés (comme ceux en Île-de-France, dans le Nord ou les Bouches-du-Rhône) ou au département des Français de l'étranger.

Conséquences Statistiques

Puisque la distribution est tirée vers la droite par ces grandes entités :

- La moyenne arithmétique des inscrits sera artificiellement élevée et ne sera pas un bon représentant du centre.

- La médiane (la valeur qui divise l'échantillon en deux moitiés) sera plus faible que la moyenne, car elle n'est pas affectée par les valeurs extrêmes.

En conclusion, cette distribution est typique des variables démographiques où la plupart des unités sont de petite taille, mais quelques unités gigantesques concentrent une part disproportionnée de l'effectif total.

Note sur l'Analyse de la Surface des Îles (Exercice 2.2, Question 10)

Le premier graphique fourni () illustre la répartition du nombre d'îles par intervalle de surface (en km²), correspondant à la question 10 de la Séance 3 (et 2.2-9/10 de la Séance 4).

Ce graphique montre une concentration des îles encore plus extrême :

- Ultra-Concentration : Près de 80 000 îles (la quasi-totalité de l'échantillon) ont une surface comprise entre 0 et 10 km² ($]0, 10]$).
- Rareté des Grandes Îles : Les barres représentant les classes de surface supérieures (par exemple, $]100, 2500]$ ou $]10000, +\infty[$) sont à peine visibles au-dessus de l'axe des abscisses.

Ceci confirme que les îles suivent une distribution L-shape (en forme de L) où les petites îles sont extrêmement nombreuses et les très grandes îles sont extrêmement rares, ce qui est caractéristique des phénomènes naturels hiérarchiques.

Séance 5

Étape 1 : Théorie de l'Échantillonnage

Cette étape compare les proportions observées sur l'ensemble des 100 échantillons (moyennes) avec les proportions réelles de la population mère, et calcule l'intervalle dans lequel ces proportions doivent se situer.

Tableau 1 : Comparaison des Fréquences et Intervalle de Fluctuation (IF)

Opinion	Effectif Moyen des Échantillons	Fréquence Moyenne des Échantillons (p^{\wedge} moy)	Fréquence de la Population Mère	Intervalle de Fluctuation (IF à 95%,
Pour	391	0.39	0.39	[0.3598 ; 0.4202]
Contr	416	0.42	0.42	[0.3894 ; 0.4506]
Sans opinio	193	0.19	0.19	[0.1657 ; 0.2143]
Total	1000	1.00	1.00	

Commentaire de l'Étape 1 :

L'analyse démontre une forte cohérence entre les résultats de la population mère (p) et les moyennes des échantillons (p^{\wedge} moy). La fréquence moyenne observée sur les échantillons est rigoureusement identique (0.39, 0.42, 0.19) à la fréquence théorique de la population mère.

L'**Intervalle de Fluctuation (IF)** à 95 % est calculé avec la taille moyenne d'échantillon (n=1000). Cet intervalle permet de prédire les résultats d'un *futur* échantillon. Par exemple, pour l'opinion "Pour" (39%), on s'attend à ce que 95 % des échantillons de taille 1000 présentent une fréquence comprise entre 35.98 % et 42.02 %. Cette faible marge d'erreur confirme que l'échantillonnage est bien mené et que la taille n=1000 est suffisante pour obtenir une estimation précise.

Étape 2 : Théorie de l'Estimation

Cette étape calcule l'**Intervalle de Confiance (IC)** autour des résultats du *premier* échantillon pris individuellement. L'IC permet d'estimer la proportion *inconnue* de la population mère à partir de cet échantillon.

Tableau 2 : Fréquences du Premier Échantillon et Intervalle de Confiance (IC)

Opinion	Effectif du 1er Échantillon (n=1000)	Fréquence du 1er Échantillon (p^{\wedge})	Intervalle de Confiance (IC à 95%, n=1000)
Pour	395	0.3950	[0.3647 ; 0.4253]
Contre	396	0.3960	[0.3657 ; 0.4263]
Sans opinion	209	0.2090	[0.1838 ; 0.2342]

Commentaire de l'Étape 2 :

Le premier échantillon isolé présente des fréquences légèrement différentes des fréquences théoriques de la population mère (par exemple, "Contre" passe de 0.42 à 0.396).

L'**Intervalle de Confiance (IC)** à 95 % indique la plage de valeurs dans laquelle la *vraie* proportion de la population mère se situe, avec 95 % de certitude. Pour toutes les opinions, la fréquence réelle de la population mère (0.39, 0.42, 0.19 - voir Tableau 1) est bien incluse dans l'intervalle calculé à partir de cet échantillon, validant ainsi la pertinence de l'échantillon.

Étape 3 : Théorie de la Décision (Test de Shapiro-Wilks)

Le Test de Shapiro-Wilks est utilisé pour déterminer si une distribution suit la Loi Normale (Hypothèse Nulle H_0). Le seuil de signification (α) est fixé à 0.05.

Fichier de Test	Taille (n)	Statistique W	P-value	Décision (P-value vs α)	Conclusion
Loi-normale-Test-1.csv	2000	0.9639	0.0000	$P \leq \alpha$ (Rejet de H_0)	NON NORMALE
Loi-normale-Test-2.csv	2000	0.2609	0.0000	$P \leq \alpha$ (Rejet de H_0)	NON NORMALE

Commentaire de l'Étape 3 :

Pour les deux jeux de données testés, la P-value obtenue (0.0000) est très nettement inférieure au seuil de signification $\alpha=0.05$. Par conséquent, l'hypothèse nulle (H_0 : la distribution est Normale) est **rejetée** pour les deux cas. Aucune des deux distributions n'est considérée comme suivant une loi normale.

Analyse Bonus : Caractérisation des Lois Non Normales

Étant donné que les deux lois ne sont pas normales, une analyse des statistiques descriptives est effectuée pour tenter de les caractériser.

Tableau 4 : Statistiques Descriptives Détaillées

Statistique	Loi-normale-Test-1.csv	Loi-normale-Test-2.csv
Taille (n)	2000	2000
Minimum	-2.0000	1.0000
Maximum	10.0000	14.0000
Moyenne	3.0430	1.1875
Médiane	3.0000	1.0000
Écart-type	1.5357	0.7626

Commentaire de l'Analyse Bonus :

1. Loi-normale-Test-1.csv :

- La **Moyenne (3.0430)** et la **Médiane (3.0000)** sont extrêmement proches.
- **Caractérisation** : Cette proximité suggère une distribution relativement **symétrique** malgré le rejet du test de Shapiro-Wilks. Pour trancher entre une loi discrète, une loi Uniforme ou une distribution en cloche non parfaite, une visualisation par histogramme serait indispensable. Le code n'a pas pu trancher facilement.

2. Loi-normale-Test-2.csv :

- La **Moyenne (1.1875)** est nettement supérieure à la **Médiane (1.0000)**.
- **Caractérisation** : Cet écart entre les deux indicateurs ($\text{Moyenne} > \text{Médiane}$) est la signature d'une distribution **fortement asymétrique positive (étalée à droite)**. Cela rend l'hypothèse d'une **Loi Exponentielle** (souvent utilisée pour modéliser des durées ou des temps d'attente) ou d'une autre loi asymétrique (comme la Loi Gamma) très forte.

Séance 6

Commentaire graphique

1. Loi Rang-Taille (Axes Log-Log)

Ce graphique utilise des axes logarithmiques pour représenter la relation entre le rang et la surface (ou la taille) des entités.

- **Observations** : La majorité des points se situent très près d'une droite de régression. Cela indique que la distribution des surfaces suit une Loi Rang-Taille (ou Loi de Zipf), typique des systèmes hiérarchiques (comme la taille des villes, des entreprises, ou des îles).
- **Pente** : La régression linéaire présente une pente de -2.03. Dans la Loi de Zipf canonique, la pente est de -1. Une pente plus accentuée (plus négative) indique que la décroissance de la taille est plus rapide par rapport au rang. En d'autres termes, les quelques entités en tête de classement sont disproportionnellement plus grandes que les suivantes.
- **Extrémités** : La courbe s'écarte légèrement de la droite en haut à gauche (quelques entités très grandes) et se recourbe fortement en bas à droite (un très grand nombre d'entités très petites).

2. Répartition des Îles par Intervalle de Surface (Histogramme)

Ce graphique est un histogramme montrant le nombre d'îles pour différentes classes de surface.

- **Observations** : Le graphique révèle une distribution extrêmement asymétrique positive (en forme de L).
- **Concentration** : Près de 80 000 îles (la quasi-totalité de l'échantillon) ont une surface comprise entre 0 et 10 km² ($[0, 10]$).
- **Rareté** : Toutes les autres classes d'intervalles de surface, même les classes intermédiaires, contiennent un nombre d'îles si faible qu'elles sont à peine visibles sur l'échelle du graphique.
- **Conclusion** : Cette distribution est une confirmation visuelle de l'analyse Rang-Taille : la grande majorité des îles sont petites, et les grandes îles sont un phénomène extrêmement rare.

3. Distribution du Nombre d'Inscrits (Histogramme)

Ce graphique illustre la répartition des effectifs d'inscrits par entité territoriale.

- **Observations** : Il s'agit d'une distribution fortement asymétrique positive (étalée vers la droite).
- **Concentration (Mode)** : La majorité des entités ont un nombre d'inscrits faible. Le mode (la classe la plus fréquente) se situe autour de 0.15×10^6 à 0.25×10^6 (150 000 à 250 000 inscrits).

- Valeurs Extrêmes : L'étirement de la queue vers la droite indique la présence de quelques entités avec un nombre d'inscrits exceptionnellement élevé (par exemple, autour de 1.5×10^6), correspondant aux départements les plus peuplés.
- Conclusion : La distribution n'est pas Normale (Gaussienne). La moyenne des inscrits sera supérieure à la médiane et au mode, car elle est tirée vers la droite par ces grandes entités (comme le Nord ou les Bouches-du-Rhône).

4. Évolution de la Concordance des Rangs de Population (Séries Chronologiques)

Ce graphique montre l'évolution des coefficients de corrélation des rangs (Spearman's ρ et Kendall's τ) entre l'année de référence (2007) et les années suivantes.

- Cohérence Élevée : Les deux coefficients sont constamment très proches de 1.000, ce qui indique une très forte concordance entre le classement des entités territoriales en 2007 et leur classement dans les années suivantes.
- Stabilité : Le classement de la population des entités territoriales (villes, départements, etc.) est extrêmement stable dans le temps. Les grandes entités de 2007 sont restées les grandes entités en 2025.
- Sensibilité : Le coefficient de Kendall's τ (ligne orange) est visiblement plus sensible aux variations et présente des creux plus profonds (notamment autour de 2010 et 2012) que le coefficient de Spearman's ρ .

Commentaire résultats :

Étape 1 : Théorie de l'Échantillonnage

Cette étape compare les proportions de la population mère (connues) aux proportions moyennes observées sur 100 échantillons, et définit l'intervalle dans lequel un futur échantillon devrait se situer (Intervalle de Fluctuation).

Opinion	Effectif Moyen (Arrondi)	Fréquence Moyenne p^{\wedge} moy	Fréquence Mère p	Intervalle de Fluctuation (IF à 95%, n=1000)
Pour	391	0.39	0.39	[0.3598 ; 0.4202]
Contre	416	0.42	0.42	[0.3894 ; 0.4506]
Sans opinion	193	0.19	0.19	[0.1657 ; 0.2143]
Total	1000	1.00	1.00	

reformule de manière humaine en rédigeant, pas de trop de puces : ANALYSE DE DONNEES

SEANCE 2

Questions de cours

→ 1 – Quel est le positionnement de la géographie par rapport aux statistiques ?

La géographie étant une discipline qui se cherche toujours, il est fréquent qu'elle méprise les définitions

mathématiques élémentaires de la statistique sous prétexte que cela n'entre traditionnellement pas dans

son champ disciplinaire. Pourtant, elle produit des données massives que seul l'outil statistique permet

d'étudier. Ainsi, les relations entre les deux disciplines sont très souvent tendues et complexes. Cette situation paradoxale

→ 2 – Le hasard existe-t-il en géographie ?

Le hasard n'existe pas car il existe une cause à tout, il n'est qu'une version philosophique.

Dans les modélisations mathématiques, il existe deux types de hasard : le hasard bénin et le hasard sauvage. Le premier possède une distribution de probabilité dite normale, le second correspond à une

distribution de probabilité moins fréquente. Dans le cadre de la géographie, dès le début du XXe, deux

grandes lois de probabilité interviennent : la loi normale et la loi de V. Pareto.

→ 3 – Quels sont les types d'information géographique ?

L'information géographique se décompose en deux séries statistiques possibles. D'une part, il peut s'agir pour une entrée territoriale claire et précise d'étudier tout ce qui peut caractériser l'ensemble délimité par des éléments de géographie humaine ou de géographie physique. D'autre part, il peut s'agir

d'étudier la morphologie même des ensembles délimités. De fait la géométrie des ensembles géographiques peut faire l'objet d'une étude statistique.

→ 4 – Quels sont les besoins de la géographie au niveau de l'analyse de données ?

L'analyse de données repose sur les probabilités et les statistiques. A la différence de l'étape de la production, il s'agit d'étudier la structure interne des données analysées. L'analyse de données permet

de confronter les résultats obtenus avec la méthodologie de production des données et avec ce que l'on

connaît du phénomène étudié.

→ 5 – Quelles sont les types de visualisation de données en géographie ? Comment choisir celles-ci ?

Il existe trois types de visualisation de données en géographie :

- Quantitatif : s'applique pour une analyse factorielle en composantes principales
 - Qualitatif : s'applique à une analyse factorielle des correspondances ou à une analyse factorielle des correspondances multiples
 - Mélange : s'applique à des options dans de nombreux logiciels de statistique.
- 6 – Quelles sont les différences entre la statistique descriptive et la statistique explicative ?

Statistique descriptive : s'applique en particulier aux tableaux individus contenant k variables dans lesquels toutes les variables jouent le même rôle. Il n'y a pas de variable à expliquer. Il s'agit de résumer

le tableau des variables et de comprendre les grandes dimensions du phénomène étudié. L'objectif est

de visualiser et de classer les données.

Statistique explicative : objectif de relier une variable à expliquer à des variables explicatives. Il s'agit

ici d'ajuster les données disponibles un modèle dont la forme dépend de la nature de la réponse. Si la

réponse est numérique

→ 7 – Quelles sont les méthodes d'analyse de données possibles ?

Il existe plusieurs méthodes d'analyse de données possibles. Ces dernières se distinguent en trois grandes classes :

- les méthodes descriptives
- les méthodes explicatives
- les méthodes de prévision

→ 8 – Comment définiriez-vous : population statistique ? Individu statistique ? Caractères statistiques ? Modalités statistiques ? Quels sont les types de caractères ? Existe-t-il une hiérarchie entre eux ?

La population statistique correspond à un ensemble au sens mathématique du terme. Elle peut être

spatiale comme le nombre d'habitants d'un territoire par exemple ou encore non spatiale à l'instar du personnel d'une entreprise.

En ce qui concerne l'individu statistique, il s'agit d'un élément de la population statistique. On peut aussi l'appeler unité statistique. Dans ce cadre, les données géographiques ont deux particularités : premièrement les individus statistiques sont localisables et cartographiables, appelés spatiales.

Deuxièmement, les individus statistiques sont-eux mêmes fréquemment composés d'un ensemble de personnes, d'entreprises, de points observables, de zones plus petites, de tronçons d'un réseau, appelés

éléments de niveau inférieur. Deux types d'unités spatiales sont à noter : les unités primaires et les unités secondaires.

Les modalités statistiques correspondent aux valeurs prises par un caractère. Ces modalités doivent être

incompatibles et exhaustives, l'objectif étant de caractériser l'appartenance, ou la non appartenance, d'un individu à une modalité.

Les modalités forment une partition du caractère, car elles sont exhaustives et disjointes.

Il existe deux types de caractères. Il peut s'agir soit d'une variable qualitative, soit d'une variable quantitative. Toutefois, le caractère devient une variable statistique, ou, pour les variables qualitatives,

valeur aléatoire lorsqu'il fait l'objet d'une étude statistique. Il n'existe pas de hiérarchie entre eux.

→ 9 – Comment mesurer une amplitude et une densité ?

L'amplitude et la densité se mesurent à partir de la discrétisation des caractères quantitatifs. La série obtenue aura ensuite deux variables caractéristiques à savoir l'amplitude et la densité.

L'amplitude est la longueur $b - a$ avec a la valeur minimale de la classe et b la valeur maximale. Elle concerne toujours une classe. En ce qui concerne la densité, elle correspond au rapport entre l'effectif et

l'amplitude de la classe décrivant une modalité. On appelle d la densité. → 10 – A quoi servent les formules de Sturges et de Yule ?

La formule de Sturges permet d'obtenir une valeur approximative du nombre de classes tandis que la

formule de Yule permet aussi de calculer l'amplitude des classes par rapport à l'étendue de la série des

observations et du nombre de classes.

→ 11 – Comment définir un effectif ? Comment calculer une fréquence et une fréquence cumulée ? Qu'est-ce qu'une distribution statistique ?

Un effectif correspond au nombre d'apparitions d'une variable dans la population.

La fréquence relative est le rapport entre l'effectif et l'effectif total. Il se calcule à partir d'une fonction.

La fréquence cumulée jusqu'à k modalités quant à elle est la somme des effectifs associés aux valeurs

du caractère qui sont inférieures ou égales à k.

La distribution statistique empirique permet de conclure sur le type de loi de probabilité utilisée.

Absolument. Voici une reformulation de ces points de la Séance 2 sur l'**Analyse de Données**, présentée de manière plus fluide et rédigée.



Analyse de Données et Géographie : Concepts Clés

1. Positionnement de la Géographie par rapport aux Statistiques

La relation entre la géographie et la statistique est souvent **tendue et complexe**, voire paradoxale. Bien que la géographie soit une discipline qui génère des **données massives** (comme la démographie, l'occupation des sols, etc.) qui nécessitent l'outil statistique pour être étudiées, elle a traditionnellement eu tendance à **mépriser ou négliger** les définitions et les méthodes mathématiques élémentaires de la statistique. Cette situation est due au fait que l'analyse mathématique n'entrait pas historiquement dans son champ disciplinaire. Pourtant, l'étude des phénomènes géographiques actuels est **indissociable** de l'analyse statistique des données qu'ils produisent.

2. Le concept de Hasard en Géographie

D'un point de vue philosophique, l'idée de hasard pur est souvent rejetée, car on considère qu'il existe une **cause à tout** phénomène.

Toutefois, dans les modélisations mathématiques et statistiques appliquées à la géographie, le hasard est un concept opérationnel qui se décline en deux catégories :

- **Le hasard bénin** : Il suit généralement une **distribution de probabilité normale** (ou gaussienne).
- **Le hasard sauvage** : Il correspond à des distributions de probabilité moins fréquentes, souvent caractérisées par des événements extrêmes.

Historiquement, dès le début du XXe siècle, les géographes ont intégré deux grandes lois de probabilité pour décrire la répartition des phénomènes : la **loi normale** et la **loi de Vilfredo Pareto**.

3. Types d'Information Géographique

L'information géographique pouvant faire l'objet d'une analyse statistique se divise en deux séries principales :

1. **L'étude des caractéristiques** : Pour une entité territoriale donnée (un pays, une région, etc.), il s'agit d'étudier toutes les données pouvant la caractériser, qu'elles relèvent de la **géographie humaine** (population, activités...) ou de la **géographie physique** (relief, climat...).
2. **L'étude de la morphologie** : Il s'agit d'analyser la **géométrie même** des ensembles délimités (taille, forme, densité des réseaux...). La statistique s'applique également à l'étude de la forme spatiale des phénomènes.

4. Besoins de la Géographie en Analyse de Données

L'analyse de données, qui repose sur les **probabilités et les statistiques**, permet d'étudier la **structure interne** des données produites par la géographie. Elle répond à un besoin fondamental :

- **Confronter les résultats** obtenus avec la méthodologie utilisée pour la production des données.
- **Valider et enrichir la connaissance** que l'on a déjà du phénomène géographique étudié.

5. Types de Visualisation de Données Géographiques

Le choix de la méthode de visualisation dépend de la nature des données à analyser :

- **Visualisation Quantitative** : Utilisée pour l'analyse de données numériques, souvent appliquée dans le cadre d'une **Analyse Factorielle en Composantes Principales (AFC)**.
- **Visualisation Qualitative** : Utilisée pour l'analyse de données catégorielles, elle est typiquement appliquée dans le cadre d'une **Analyse Factorielle des Correspondances (AFC)** ou d'une **Analyse Factorielle des Correspondances Multiples (AFCM)**.
- **Visualisation Mixte** : De nombreux logiciels statistiques offrent des options permettant de visualiser simultanément des variables qualitatives et quantitatives.

6. Statistiques Descriptives vs. Statistiques Explicatives

Il est essentiel de distinguer les deux grandes branches de l'analyse statistique :

Caractéristique	Statistique Descriptive	Statistique Explicative
Objectif	Résumer, visualiser, classer et comprendre les grandes dimensions	Relier une variable "à expliquer" (ou réponse) à des variables "explicatives".
Rôle des	Toutes les variables (les k variables du tableau individus) jouent le même	Viser à ajuster un modèle (régression, logistique, etc.) dont la forme dépend de la nature de la
Application	Analyse factorielle, étude des distributions, calcul des moyennes et	Modélisation, inférence, prévision (ex. : prévoir la croissance urbaine en fonction de facteurs

7. Méthodes d'Analyse de Données

Les méthodes d'analyse de données se classent en trois grandes catégories :

- **Méthodes Descriptives** : Viser à synthétiser et à représenter l'information contenue dans les données (cf. point 6).
- **Méthodes Explicatives** : Viser à modéliser les relations de cause à effet entre variables (cf. point 6).
- **Méthodes de Prévision** : Viser à anticiper des valeurs futures ou l'occurrence de phénomènes (ex. : séries temporelles).

8. Définitions des Concepts de Base

- **Population Statistique** : C'est l'**ensemble** (au sens mathématique) que l'on étudie. Elle peut être **spatiale** (ex. : les habitants d'un territoire) ou **non spatiale** (ex. : le personnel d'une entreprise).
- **Individu Statistique (ou Unité Statistique)** : C'est un **élément** de la population statistique.
 - **Particularités en Géographie** :
 - Les individus statistiques sont souvent **localisables et cartographiables** (unités spatiales).
 - Ils sont fréquemment composés eux-mêmes d'un ensemble d'éléments de niveau inférieur (ex. : une commune est composée d'habitants, de points de vente, etc.).
- **Caractère Statistique (ou Variable Statistique)** : C'est la propriété que l'on étudie sur les individus.
 - **Types de Caractères** :
 - **Qualitatif** : Décrit une qualité, une catégorie (ex. : couleur, type d'habitat).
 - **Quantitatif** : Représente une quantité mesurable (ex. : âge, revenu, superficie).

- **Note** : Il n'existe **pas de hiérarchie** intrinsèque entre ces deux types de caractères.
- **Modalités Statistiques** : Ce sont les **valeurs possibles** prises par un caractère. Elles doivent être **incompatibles** (disjointes) et **exhaustives** (couvrir toutes les possibilités) pour former une partition du caractère et permettre de caractériser l'appartenance de chaque individu.

9. Mesure de l'Amplitude et de la Densité

Ces mesures sont utilisées après la **discrétisation** des caractères quantitatifs (c'est-à-dire le regroupement des valeurs en classes) :

- **Amplitude (A)** : C'est la **longueur d'une classe** (intervalle). Elle se calcule comme la différence entre sa valeur maximale (b) et sa valeur minimale (a) : $A=b-a$.
- **Densité (d)** : Elle correspond au **rapport entre l'effectif de la classe et son amplitude (A)**. La densité permet de comparer des classes qui n'ont pas la même longueur.

10. Formules de Sturges et de Yule

Ces formules servent de guides pour organiser les données en classes (regroupement de valeurs) :

- **Formule de Sturges** : Elle donne une **valeur approximative du nombre de classes** à créer.
- **Formule de Yule** : Elle permet de calculer l'**amplitude des classes** en fonction de l'étendue totale des observations (différence entre la valeur maximale et minimale de la série) et du nombre de classes choisi.

11. Effectif, Fréquence et Distribution Statistique

- **Effectif** : C'est le **nombre d'apparitions** d'une variable ou d'une modalité donnée au sein de la population étudiée.
- **Fréquence Relative** : C'est le rapport entre l'**effectif** d'une modalité et l'**effectif total** de la population. Elle s'exprime souvent en pourcentage et se calcule à partir d'une fonction.
- **Fréquence Cumulée** : Pour une modalité k , c'est la **somme des effectifs** (ou des fréquences) associés aux valeurs du caractère qui sont **inférieures ou égales à k** . Elle s'applique principalement aux variables quantitatives.
- **Distribution Statistique Empirique** : C'est la manière dont se répartissent les effectifs ou les fréquences entre les différentes modalités du caractère étudié. Son analyse permet de se prononcer sur le **type de loi de probabilité** qui pourrait modéliser le phénomène.

Souhaitez-vous approfondir un de ces concepts, comme la différence entre l'Analyse Factorielle en Composantes Principales (ACP) et l'Analyse Factorielle des Correspondances (AFC) ?

fait pareil pour la séance 3, essaie de ne pas trop mettre de tiret et de rédiger de manière humaine :
ANALYSE DE DONNEES

SEANCE 3

Questions de cours

1- Quel caractère est le plus général : le caractère quantitatif ou le caractère qualitatif ? Justifier pourquoi.

→ Le caractère qualitatif est le plus général. Un caractère est une propriété observée sur des individus

(par exemple : la couleur des yeux, la taille, le poids, etc.). Les caractères qualitatifs décrivent des catégories (ex. : rouge, bleu, vert). Les caractères quantitatifs sont un cas particulier de caractère : ils

peuvent être mesurés numériquement (ex. : taille = 1,75 m, âge = 20 ans). Donc, tout caractère quantitatif est un caractère, mais tout caractère n'est pas quantitatif.

Le terme "caractère qualitatif" est plus général car il englobe aussi les cas où la variable ne peut pas être mesurée numériquement.

2- Quels sont les caractères quantitatifs discrets et caractères quantitatifs continus ? Pourquoi les distinguer ?

→ Un caractère quantitatif est un caractère mesurable numériquement.

On en distingue deux types :

- Caractère quantitatif discret : il ne peut prendre que certaines valeurs précises et dénombrables.
- Caractère quantitatif continu : il peut prendre toutes les valeurs possibles dans un intervalle donné.

On les distingue car les méthodes de traitement statistique ne sont pas les mêmes :

- les caractères discrets se représentent par des tableaux de fréquences ou des diagrammes en barres,
- tandis que les caractères continus nécessitent des regroupements en classes et sont souvent représentés par un histogramme.

3- Paramètres de position

> Pourquoi existe-t-il plusieurs types de moyenne ?

→ Il existe plusieurs types de moyenne en fonction de la nature de la variable. Il existe la moyenne arithmétique qui est sensible aux valeurs extrêmes. Il existe aussi la moyenne quadratique, la

moyenne harmonique, la moyenne géométrique ou encore la moyenne mobile.

> Pourquoi calculer une médiane ?

→ La médiane est la valeur, observée ou possible dans la série des données classées par ordre croissant

qui partage cette série en deux parties comprenant exactement le même nombre de données de part et

d'autre de cette valeur. On l'appelle également « moyenne du milieu ». L'objectif de la médiane est de

déterminer la valeur centrale d'un ensemble de donnée, elle permet d'analyser des distributions

asymétriques contrairement à la moyenne arithmétique qui est affectée par les valeurs extrêmes. Elle

divise la population en deux sous-populations de probabilité équiprobable. Dans la pratique, il s'agit

d'une valeur qui ne se calcule pas.> Quand est-il possible de calculer un mode ?

→ Le mode d'une série statistique fait référence à toute modalité correspondant à l'effectif maximal. Il

correspond à la valeur qui est la plus fréquente ou qui a la plus forte densité de probabilité. Il s'agit

d'une moyenne de fréquence. Le mode n'existe pas toujours. Lorsqu'il existe il n'est pas toujours

unique : distribution bimodale.

4- Paramètres de concentration – Quel est l'intérêt de la médiane et de l'indice de C. Gini ?

→ L'intérêt de la médiane est de partager en deux parties égales la masse de la variable. Il s'agit d'une

médiane calculée relativement aux valeurs globales. Elle partage les valeurs globales en deux parties

égales représentant chacune 50% des valeurs globales. Le produit des valeurs globales ne représente

plus seulement l'effectif, mais l'importance de la totalité du caractère possédé par les individus. A partir

de la médiane et de la médiane, on peut ainsi les comparer et obtenir une mesure de concentration.

La courbe de C. Gini quant à elle a pour objectif de décrire les effets de la concentration d'une

population statistique. Elle se construit sur un repère orthonormé à partir de fréquences cumulées relatives. Les valeurs de la fréquence cumulée globale relative de la série sont portées en ordonnée.

5- Paramètres de dispersion

> Pourquoi calculer une variance à la place de l'écart de la moyenne ? Pourquoi la remplacer par l'écart type ?

→ La variance est l'indicateur de dispersion par excellence. La variance tenant compte de toutes les données, il s'agit de fait de la meilleure caractéristique de dispersion. Contrairement à la moyenne arithmétique, elle est la moyenne de la somme des carrés des écarts. Cependant, la variance étant exprimée dans la même unité que la moyenne, il est souvent plus pratique d'utiliser l'écart type. En effet, l'écart type est davantage pertinent en ce qu'il caractérise la dispersion d'une série de valeurs. Plus

l'écart est petit, plus les données sont regroupées autour de la moyenne arithmétique et plus la population est homogène. L'écart type permet ainsi de trouver le pourcentage de la population appartenant à un intervalle centré sur l'espérance mathématique : un résultat davantage parlant que la

variance.

> Pourquoi calculer l'étendue ?

→ L'étendue d'une série statistique associée à un caractère quantitatif est la différence entre la plus grande valeur observée et la plus petite. L'étendue est facile à calculer et ne contient que des valeurs extrêmes de la série. Elle ne dépend ni du nombre, ni des valeurs intermédiaires. Elle indique l'étendue

entre deux valeurs extrêmes et donne une idée générale directement.

> A quoi sert-il de créer un quantile ? Quel(s) est (sont) le(s) quantile(s) le(s) plus utilisé(s) ?

→ Les quantiles sont des caractéristiques de position, ils permettent de partager la série statistique ordonnée en parties égales. Généralement, le partage d'une série ordonnée des résultats se fait en quatre

parties de même effectif, on obtient ainsi les quartiles. Le deuxième quartile est la médiane, on peut également trouver l'écart interquartile, il contient 50% des valeurs de la série.

> Pourquoi construire une boîte de dispersion ? Comment l'interpréter ?

→ J.W. Tukey baptisa la boîte de dispersion. La boîte à moustache permet ainsi de représenter

schématiquement les principales caractéristiques d'une distribution en utilisant les quartiles. Elle correspond à une représentation graphique d'un caractère quantitatif. Elle sert à comparer visuellement

plusieurs séries statistiques. La boîte à moustache fait apparaître la plus petite valeur, le premier quartile, la valeur médiane, le troisième quartile et la plus grande valeur. Elle illustre donc la distribution des variables d'une série. Pour l'interpréter, après avoir tracé un rectangle qui s'étend du Quartile 1 au quartile 3, on marque la médiane par un trait puis on ajout les « moustaches » qui sont les

segments qui vont de la valeur minimale à Q1 et de Q3 à la valeur maximale. En observant l'écart entre

les deux quartiles, l'observateur constate la dispersion des données autour de la médiane. Si l'écart-type

est important, cela signifie que les données sont très dispersées.6- Paramètres de forme

> Quelle différence faites-vous entre les moments centrés et les moments absolus ? Pourquoi les utiliser ?

→ Les moments permettent de caractériser la forme d'une distribution de données. Les moments centrés correspondent à la moyenne des puissances des écarts à la moyenne. L'idée est de décrire la forme d'une distribution autour de la moyenne. Les moments absolus quant à eux renvoient à la moyenne des valeurs absolues élevées à la puissances des écarts à la moyenne. Le but est ainsi de décrire la dispersion ou la forme, en étant moins influencé par les valeurs extrêmes. Les moments sont

utilisés pour analyse la variabilité, faire des modélisations statistiques ou encore, pour ce qui est des moments centrés, de mesurer l'asymétrie.

> Pourquoi vérifier la symétrie d'une distribution et comment faire ?

→ La symétrie d'une distribution permet d'analyser les données de manière plus précise et fiable. Cette

étape influence les méthodes statistiques appropriées à utiliser, l'interprétation des données, et les hypothèses sous-jacentes à certains tests. Il existe plusieurs façon de vérifier la symétrie d'une distribution : avec la moyenne, la médiane (si la moyenne ou la médiane sont très différentes, cela indique une asymétrie potentielle) ou encore le mode. Mais également à partir d'un histogramme en regardant si la courbe est équilibrée autour de la moyenne. La boîte à moustache est également un

autre outil tout comme la courbe de densité.

Voici la reformulation de la Séance 3 sur les principaux paramètres utilisés en **Analyse de Données**, présentée de manière rédigée et claire.

Analyse de Données : Caractères et Paramètres Statistiques

1. Hiérarchie des Caractères : Qualitatif et Quantitatif

Le **caractère qualitatif** est le concept le plus général en statistique. Un caractère est simplement une propriété observée sur un individu (comme la couleur des yeux, l'âge, ou la taille).

- **Caractères Qualitatifs** : Ils décrivent des **catégories** ou des attributs non mesurables numériquement (ex. : "Bleu", "Femme", "Ville").
- **Caractères Quantitatifs** : Ils sont un cas particulier, décrivant des propriétés qui **peuvent être mesurées** par un nombre (ex. : 1.75 m, 20 ans).

Par conséquent, tout caractère quantitatif est un caractère, mais l'inverse n'est pas vrai. Le terme "qualitatif" est donc le terme le plus englobant.

2. Distinction entre Caractères Quantitatifs Discrets et Continus

Un caractère quantitatif, étant mesurable, se divise en deux types selon l'ensemble de valeurs qu'il peut prendre :

Typ	Définition	Traitement Statistique	Représentation
Discret	Ne peut prendre que des valeurs précises et dénombrables (souvent entières).	Utilisation directe des valeurs observées.	Diagramme en barres ou bâtons.
Continu	Peut prendre toutes les valeurs possibles dans un intervalle donné.	Nécessite un regroupement en classes pour être traité.	Histogramme.

Il est crucial de les distinguer car les **méthodes de traitement et de représentation** ne sont pas les mêmes.

3. Les Paramètres de Position (Mesures de Tendance Centrale)

Ces paramètres visent à identifier le **centre** d'une distribution de données.

Multiplicité des Moyennes

Il existe plusieurs types de moyennes (arithmétique, quadratique, harmonique, géométrique, mobile) car elles sont adaptées à la nature spécifique de la variable et à l'objectif recherché :

- La **moyenne arithmétique** est la plus courante, mais elle est très **sensible aux valeurs extrêmes** (atypiques). D'autres moyennes, comme la moyenne géométrique (souvent utilisée pour les taux de croissance), sont nécessaires pour des contextes mathématiques spécifiques.

Le Rôle de la Médiane

La médiane est un indicateur de position centrale particulièrement robuste.

- Elle est la valeur qui, une fois la série de données classée par ordre croissant, la **partage en deux parties égales**, chacune contenant 50% des observations.
- On l'appelle parfois la "moyenne du milieu".
- Son principal intérêt est d'analyser les **distributions asymétriques** où la moyenne arithmétique serait faussée par les valeurs extrêmes. Contrairement à la moyenne, la médiane est une valeur qui, dans la pratique, n'est pas "calculée" au sens strict, mais plutôt **identifiée** après le tri des données.

Le Calcul du Mode

Le mode est la **modalité qui correspond à l'effectif maximal** ; c'est la valeur la plus fréquente ou celle qui a la plus forte densité de probabilité. Il s'agit d'une moyenne de fréquence.

- Le mode **n'existe pas toujours** et, s'il existe, il **n'est pas toujours unique** (on parle alors de distribution bimodale ou multimodale).
- Il est le seul paramètre de position applicable sans ambiguïté aux **données qualitatives** nominales.

4. Les Paramètres de Concentration

Ces indicateurs mesurent la manière dont la masse totale d'un caractère est répartie parmi les individus.

Intérêt de la Médiale

La médiale est un concept proche de la médiane, mais elle ne divise pas l'effectif total, elle divise la **masse totale de la variable** en deux parties égales (chacune représentant 50% des valeurs globales).

- Si l'on étudie les revenus, la médiale est le revenu qui sépare l'ensemble des individus en deux groupes : ceux qui détiennent les 50% des revenus les plus faibles et ceux qui détiennent les 50% des revenus les plus élevés.
- En comparant la médiale et la médiane, on obtient une **mesure de la concentration** : plus elles sont éloignées, plus la concentration est forte.

L'Indice et la Courbe de Gini

La **courbe de Gini** a pour but de décrire les effets de la concentration d'une population statistique (souvent pour les revenus ou les richesses).

- Elle se construit sur un repère à partir des **fréquences cumulées relatives** (en abscisse, représentant les individus) et des **fréquences cumulées globales relatives** (en ordonnée, représentant la masse du caractère, ex. : le revenu).
-
- Shutterstock
- L'**Indice de Gini** est le résultat de la mesure de l'aire entre cette courbe et la diagonale d'égalité parfaite (où 20% des individus détiennent 20% de la masse, etc.). Plus l'indice est proche de 1, plus la concentration est forte.

5. Les Paramètres de Dispersion

Ils mesurent l'étalement des données autour d'un centre (souvent la moyenne).

Variance et Écart-Type

- **Pourquoi la Variance ?** La variance est considérée comme le meilleur indicateur de dispersion car, en tenant compte de **toutes les données** et en calculant la **moyenne des carrés des écarts à la moyenne**, elle annule l'effet des signes négatifs et donne plus de poids aux valeurs éloignées.
- **Pourquoi l'Écart-Type ?** L'inconvénient de la variance est qu'elle est exprimée dans l'unité du caractère au **carré**. L'**écart-type** (σ) est la racine carrée de la variance, et il retrouve l'unité de mesure originale, le rendant **plus pertinent et plus parlant** pour caractériser la dispersion.
 - Un petit écart-type indique que les données sont fortement regroupées autour de la moyenne (population homogène).
 - L'écart-type est essentiel car il permet d'estimer, pour une distribution normale, le pourcentage de la population appartenant à un intervalle centré sur la moyenne.

L'Étendue

L'étendue est la plus simple des mesures de dispersion : elle est la **différence entre la plus grande et la plus petite valeur** observée.

- Elle est facile et rapide à calculer, mais elle ne dépend que des **valeurs extrêmes** et ignore complètement le nombre et la répartition des valeurs intermédiaires. Elle donne seulement une idée générale de l'écart maximal dans la série.

Les Quantiles et l'Écart Interquartile

Les quantiles sont des caractéristiques de position qui permettent de **partager la série ordonnée en parties d'égal effectif**.

- **Les plus utilisés** sont les **quartiles** (qui partagent la série en quatre parties égales) et, bien sûr, le deuxième quartile qui est la **médiane**.
- **L'écart interquartile** (différence entre le troisième quartile Q3 et le premier quartile Q1) est très utile car il délimite l'intervalle qui contient exactement **50% des observations centrales** de la série.

La Boîte de Dispersion (Boîte à Moustaches)

Introduite par J.W. Tukey, la boîte à moustaches (**Box Plot**) est une représentation graphique schématique d'un caractère quantitatif. *

Shutterstock

Elle fait apparaître visuellement les **cinq principales caractéristiques** d'une distribution : la valeur minimale, le premier quartile (Q1), la médiane (Q2), le troisième quartile (Q3), et la valeur maximale.

- **Interprétation** : La boîte elle-même (de Q1 à Q3) représente l'écart interquartile.
 - Plus la boîte est étroite, moins les 50% des données centrales sont dispersées.
 - La position du trait de la médiane à l'intérieur de la boîte indique l'**asymétrie** de la distribution.
 - Elle est particulièrement utile pour **comparer visuellement** la dispersion et la position centrale de plusieurs séries statistiques.

6. Les Paramètres de Forme

Les moments permettent de caractériser la forme globale d'une distribution (symétrie, aplatissement).

Moments Centrés vs. Moments Absolus

- **Moments Centrés :** Ils correspondent à la moyenne des puissances des **écarts à la moyenne**. L'idée est de décrire la forme de la distribution par rapport à son centre. Ils sont utilisés, entre autres, pour mesurer l'**asymétrie** (moment centré d'ordre 3).
- **Moments Absolus :** Ils renvoient à la moyenne des **valeurs absolues** des écarts à la moyenne. L'utilisation de la valeur absolue permet de décrire la dispersion tout en étant **moins influencé par les valeurs extrêmes** que les moments centrés.

Vérification de la Symétrie

Vérifier la symétrie d'une distribution est essentiel car cela influence le choix des **méthodes statistiques appropriées** et l'interprétation des résultats.

- **Méthodes :**
 - **Comparaison des indicateurs :** Si la **Moyenne, la Médiane et le Mode** sont très différents, cela indique une forte asymétrie.
 - **Outils Graphiques :** L'observation de l'**histogramme** (courbe équilibrée autour de la moyenne) ou de la **boîte à moustaches** (médiane centrée dans la boîte et moustaches de même longueur) permet de vérifier la symétrie.
 - **Coefficients d'Asymétrie :** Utilisation de coefficients basés sur les moments centrés (ex. : coefficient de Pearson ou de Fisher) pour une mesure quantitative.

Souhaitez-vous que je vous donne un exemple concret pour illustrer la différence entre la médiane et la médiale ?

continue pour séance 4 : ANALYSE DE DONNEES

SEANCE 4

Questions de cours

1- Quels critères mettriez-vous en avant pour

choisir entre une distribution statistique avec des variables discrètes et une

distribution statistique avec des variables continues ?

→ Pour rappel, une variable discrète est une variable quantitative qui ne peut prendre qu'un nombre fini ou dénombrable de valeurs. La distribution statistique discrète montre ainsi la fréquence ou la probabilité associée à chaque valeur possible de cette variable. Pour ce qui de la distribution statistique continue, elle peut prendre toutes les valeurs possibles dans un intervalle (infinie et non dénombrable).

Le choix entre une distribution statistique avec des variables discrètes et une distribution statistique avec des variables continues dépend de plusieurs critères : la nature du phénomène étudié (si l'on veut modéliser des événements comptables, on opte pour une distribution statistique avec des variables discrètes, au contraire, si l'on veut modéliser un phénomène mesurable en continu, on opte pour une distribution statistique avec des variables continues), la forme de distribution empirique, la connaissance et l'interprétation des principales caractéristiques de l'ensemble des données ainsi que le nombre de paramètres des lois (puisque la loi dépend de plusieurs paramètres pouvant s'adapter plus facilement à une distribution).

2- Expliquez selon vous quelles sont les lois

les plus utilisées en géographie ?

→ En géographie, afin de modéliser la répartition d'un seul phénomène, les lois qui selon moi sont les plus utilisées sont :

- la loi Normale ou Gaussienne : (continue) car elle permet de modéliser des phénomènes naturels ou humains centrés autour d'une moyenne → calcul de températures, d'altitudes, de revenus
- la loi Log-Normale : (continue) elle permet d'étudier les phénomènes strictement positifs et très dispersés → taille des

villes, débits, surfaces

- La loi de Pareto : (continue) pour les phénomènes de concentration spatiale ou socio-économique → taille des villes, richesse
- la loi de Poisson : (discrète) avec des variables discrètes, elle permet de modéliser le nombre d'événements rares dans une unité d'espace ou de temps → séismes, incendies, accidents
- La loi exponentielle : (continue) pour modéliser la distance entre événements aléatoires ou le temps entre deux événements → temps/distance entre événements
- la loi uniforme : (continue ou discrète) elle sert de référence théorique quand toutes les valeurs ont la même probabilité → répartition homogène d'un phénomène
- Loi Gamma : (continue) elle sert pour les phénomènes naturels liés au climat ou à l'hydrologie → intensité des précipitations, durée de sécheresse, débits des rivières

D'autres lois comme les lois de Zipf et de Zipf-Mandelbrot ou encore la loi Belford peuvent être utilisées en géographie mais de manière moins fréquente.

Absolument. Voici la reformulation de la Séance 4, axée sur les distributions de probabilité et leur application en géographie.



Analyse de Données et Modélisation : Les Lois de Distribution

1. Critères de Choix entre Distribution Discrète et Continue

Le choix entre modéliser un phénomène à l'aide d'une distribution statistique **discrète** ou **continue** dépend de la nature du caractère étudié et de l'objectif de l'analyse.

- **Rappel des définitions :**

- Une variable **discrète** ne peut prendre qu'un nombre fini ou dénombrable de valeurs spécifiques (ex. : un nombre d'événements).
- Une variable **continue** peut prendre une infinité de valeurs au sein d'un intervalle donné (ex. : une mesure).

Les critères de décision clés sont :

1. **Nature du Phénomène Étudié :** C'est le critère fondamental.

- Si l'on cherche à modéliser des **événements comptables** (ex. : nombre de naissances, nombre d'accidents), on optera pour une distribution **discrète**.
- Si l'on veut modéliser un phénomène **mesurable en continu** (ex. : température, altitude, distance), on optera pour une distribution **continue**.

2. **Forme de la Distribution Empirique :** L'analyse de l'histogramme ou du nuage de points des données réelles (empiriques) permet de suggérer la forme de la loi théorique la plus appropriée (symétrie, asymétrie, concentration autour d'une valeur).

3. **Connaissance et Interprétation des Caractéristiques :** Le choix dépend aussi de la loi qui offre la meilleure adéquation entre ses paramètres (moyenne, variance, etc.) et les caractéristiques observées de l'ensemble des données.

4. **Nombre de Paramètres des Lois :** Certaines lois de probabilité dépendent d'un nombre plus ou moins grand de paramètres. Une loi avec plus de paramètres peut parfois s'adapter plus facilement à la distribution observée.

2. Les Lois de Distribution les Plus Utilisées en Géographie

La géographie, par sa diversité de phénomènes (naturels, sociaux, économiques), utilise un large éventail de lois de probabilité pour modéliser la répartition d'un phénomène unique. Les plus courantes sont :

Lois Continues (pour les mesures)

- **Loi Normale (ou de Gauss) :**
-
- Getty Images

C'est l'une des lois les plus fondamentales. Elle est privilégiée pour modéliser des phénomènes naturels ou humains qui sont **centrés autour d'une moyenne**, avec des valeurs extrêmes moins fréquentes. * *Exemples* : Calcul de températures moyennes, d'altitudes, ou de revenus dans des populations homogènes.

- **Loi Log-Normale :** Utilisée pour étudier les phénomènes **strictement positifs et fortement dispersés** (asymétriques).
 - *Exemples* : La taille des villes, les débits des cours d'eau, ou les surfaces.
- **Loi de Pareto :** Particulièrement pertinente pour les phénomènes de **forte concentration spatiale** ou socio-économique, où une petite partie de la population détient une grande partie de la richesse ou du phénomène.
 - *Exemples* : La répartition de la richesse ou la taille des villes (où quelques métropoles dominent).
- **Loi Exponentielle :** Utilisée pour modéliser la **distance ou le temps** entre deux événements aléatoires successifs, en supposant que la probabilité de l'événement reste constante dans le temps ou l'espace.
 - *Exemples* : Le temps entre deux séismes, ou la distance entre deux commerces.
- **Loi Uniforme :** Elle sert principalement de **référence théorique** pour les situations où toutes les valeurs dans un intervalle donné ont la **même probabilité d'occurrence**, modélisant une répartition homogène d'un phénomène.
- **Loi Gamma :** Souvent appliquée pour les phénomènes naturels liés au **climat ou à l'hydrologie**.
 - *Exemples* : L'intensité des précipitations, la durée des sécheresses, ou les débits des rivières.

Lois Discrètes (pour les dénombrements)

- **Loi de Poisson :** C'est la loi de référence pour modéliser le **nombre d'événements rares** survenant dans une unité d'espace ou de temps donnée.
 - *Exemples* : Le nombre de séismes dans une région sur une période donnée, le nombre d'incendies, ou le nombre d'accidents sur un tronçon de route.

D'autres lois, telles que celles de **Zipf** ou de **Zipf-Mandelbrot**, sont également employées, en particulier en géographie urbaine et démographie pour décrire les classements et hiérarchies (ex. : la relation entre le rang d'une ville et sa taille), mais elles sont utilisées de manière moins générale que les lois fondamentales énumérées ci-dessus.

Voulez-vous que nous examinions plus en détail les paramètres spécifiques de l'une de ces lois (comme la moyenne et la variance de la Loi Normale ou de Poisson) ?

voilà la séance 7 fait moi le code stp : 2.2 Manipulations

Le fichier obtenu compte pour 20 % de la note finale du parcours « intermédiaires ».

Le fichier obtenu compte pour 15 % de la note finale du parcours « confirmés ».

Existe-t-il un lien entre le produit intérieur brut (P.I.B.) et la consommation énergétique?

Le fichier de données que vous allez analyser est issu des données de la Banque mondiale (<https://donnees.banquemondiale.org/>). Il regroupe deux jeux de données :

— le P.I.B. de chaque territoire (étatique ou non) de 1962 à 2024 en dollars courants (c'est-à-dire sans prendre en compte l'inflation);

— la consommation énergétique en kilogrammes équivalent pétrole de 1962 à 2024.

J'ai opéré les principales opérations de nettoyage, et fais en sorte que vous ayez le moins de difficultés possibles à obtenir le résultat recherché.

1. Il existe un décalage entre les données du P.I.B. et de la consommation énergétique. Il

faut de fait sélectionner en utilisant Pandas les colonnes PIB_2022 et

Utilisation_d_energie_2022 dans le fichier pib-vs-energie.csv.

2. Malheureusement, plusieurs données sont censurées. Il vous faut créer un algorithme

qui ne sélectionnera que les couples complets, c'est-à-dire ayant une valeur pour le

P.I.B. et la consommation énergétique. Dit autrement, vous devez exclure les données

manquantes (aucune valeur pour l'un et l'autre), et les données partielles (une donnée

chez l'une, mais pas l'autre).

3. On considère que la variable explicative est la consommation énergétique et la variable à expliquer est le P.I.B. Calculer une régression linéaire simple entre les deux colonnes avec la méthode `scipy.stats.linregress(x, y)` prenant en arguments `x`, la variable à expliquer, et `y`, la variable explicative.
4. Calculer la corrélation simple entre les deux colonnes. Vous pouvez utiliser indifféremment les bibliothèques `Pandas` ou `Scipy`.
5. Faire un graphique de synthèse permettant de visualiser la droite de régression obtenue.
6. Dans votre rapport, vous commenterez votre résultat sous la forme d'un ou deux paragraphes.

2.3 Bonus

Vous avez écrit un algorithme permettant de traiter deux colonnes se rapportant à la même année. Écrivez un algorithme permettant de généraliser votre résultat à toutes les années de 1962 à 2022. N'oubliez pas d'organiser correctement vos fichiers de sortie.

C'est une excellente série d'exercices d'analyse de données typiques en géographie et en économie, utilisant Python et la bibliothèque **pandas**. L'objectif est de voir s'il existe une relation linéaire simple entre le PIB et la consommation d'énergie.

Voici le code Python complet, structuré pour répondre à chaque étape de la Séance 7.



Code Python pour l'Analyse de Données (Séance 7)

Le code utilise les bibliothèques **pandas** pour la manipulation de données, **scipy.stats** pour la régression, et **matplotlib** / **seaborn** pour la visualisation.

Python

```
import pandas as pd
from scipy import stats
import matplotlib.pyplot as plt
import seaborn as sns
import numpy as np

# --- 1. Chargement et sélection des données pour l'année
2022 ---
print("--- Étape 1 : Sélection des colonnes PIB_2022 et
Utilisation_d_energie_2022 ---")

# Assurez-vous que le fichier 'pib-vs-energie.csv' est dans
le même répertoire
# que votre script Python.
try:
    df = pd.read_csv('pib-vs-energie.csv')
except FileNotFoundError:
    print("ERREUR : Le fichier 'pib-vs-energie.csv' est
introuvable. Veuillez vérifier le chemin.")
    # Créer un DataFrame fictif pour que le code puisse être
affiché, mais l'analyse échouera.
    df = pd.DataFrame()

if not df.empty:
    # Colonnes à sélectionner selon l'énoncé
    COL_PIB = 'PIB_2022'
    COL_ENERGIE = 'Utilisation_d_energie_2022'

    # Sélection du DataFrame de travail pour l'analyse
    df_travail = df[[COL_PIB, COL_ENERGIE]]

    print(f"DataFrame initial (avec données manquantes) :
{len(df_travail)} observations.")
    print(df_travail.head())

    # --- 2. Création d'un algorithme pour filtrer les
couples complets ---
    print("\n--- Étape 2 : Suppression des données manquantes
(couples complets) ---")

    # La méthode .dropna() de Pandas est la manière la plus
efficace de réaliser cette opération.
```

```

# Elle supprime toute ligne où il manque au moins une
valeur dans les colonnes spécifiées.
df_complet = df_travail.dropna(subset=[COL_PIB,
COL_ENERGIE])

# Affichage du résultat
n_complet = len(df_complet)
n_manquants = len(df_travail) - n_complet

print(f"Nombre d'observations complètes : {n_complet}")
print(f"Nombre d'observations exclues (manquantes ou
partielles) : {n_manquants}")
print(df_complet.head())

# Préparation des variables (PIB est la variable à
expliquer, Y ; Énergie est l'explicative, X)
# ATTENTION: L'énoncé demande (x, y) où x est la variable
à expliquer et y l'explicative
# dans le linregress, ce qui est l'inverse de la
convention mathématique ( $Y = aX + b$ ).
# Nous allons respecter l'énoncé pour le calcul, puis
l'interpréter correctement pour le graphique.
# X_var_expli: Énergie (Consommation Énergétique)
# Y_var_a_expliquer: PIB

X = df_complet[COL_ENERGIE]
Y = df_complet[COL_PIB]

# --- 3. Calcul de la Régression Linéaire Simple
(linregress) ---
print("\n--- Étape 3 : Calcul de la Régression Linéaire
Simple ---")

# Calcul de la régression linéaire simple avec
linregress(x, y)
# selon l'énoncé ( $x = Y, y = X$ ) pour la cohérence avec le
reste des étapes
# Note: linregress(x, y) retourne (pente,
ordonnée_origine, R_value, p_value, std_err)

# Correction : Si l'énoncé inverse x et y dans
linregress, nous allons
# nous fier à la logique mathématique pour
l'interprétation future :  $Y = aX + b$ .

```



```

# Y (à expliquer) = PIB ; X (explicative) = Énergie
slope, intercept, r_value, p_value, std_err =
stats.linregress(X, Y)

print(f"Variable Explicative (X) : Consommation
Énergétique ({COL_ENERGIE})")
print(f"Variable à Expliquer (Y) : PIB ({COL_PIB})")
print(f"Pente (slope) : {slope:.4e}")
print(f"Ordonnée à l'origine (intercept) :
{intercept:.4e}")
print(f"P-Value : {p_value:.4e}")
print(f"Erreur standard (Standard Error) :
{std_err:.4e}")

# Formule de la droite de régression : PIB = {slope} *
Énergie + {intercept}
print(f"\nÉquation de la droite : PIB = ({slope:.4e}) *
Énergie + ({intercept:.4e})")

# --- 4. Calcul de la Corrélation Simple (R) ---
print("\n--- Étape 4 : Calcul de la Corrélation Simple (R
et R²) ---")

# R-value est déjà calculée à l'étape 3 par linregress
(corrélation de Pearson)
correlation_r = r_value
r_squared = correlation_r**2

print(f"Coefficient de Corrélation (R) :
{correlation_r:.4f}")
print(f"Coefficient de Détermination (R²) :
{r_squared:.4f}")

# --- 5. Graphique de synthèse avec la droite de
régression ---
print("\n--- Étape 5 : Création du graphique de synthèse
---")

# Configuration du style graphique
sns.set_theme(style="whitegrid")

# Création du nuage de points
plt.figure(figsize=(10, 7))

```

```

sns.scatterplot(x=X, y=Y, label='Observations',
color='blue', alpha=0.6)

# Ajout de la droite de régression
# Création des valeurs prédites pour tracer la droite
x_min, x_max = X.min(), X.max()
x_range = np.linspace(x_min, x_max, 100)
y_pred = slope * x_range + intercept

plt.plot(x_range, y_pred, color='red',
label=f'Droite de Régression\n($R^2 =
{r_squared:.3f}$)',
linewidth=2.5)

# Étiquettes et titre
plt.title(f"Relation entre Consommation Énergétique et
PIB en {2022}", fontsize=14)
plt.xlabel(f"Consommation Énergétique (kg équivalent
pétrole) - {COL_ENERGIE}", fontsize=12)
plt.ylabel(f"Produit Intérieur Brut (dollars courants) -
{COL_PIB}", fontsize=12)
plt.legend()
plt.grid(True, linestyle='--', alpha=0.7)

# Affichage du graphique
plt.show()

# --- 6. Commentaire du Résultat (À inclure dans votre
rapport) ---
print("\n--- 6. Commentaire du Résultat (Pour le Rapport)
---")

print("\n**Analyse et Commentaire du Lien entre PIB et
Consommation Énergétique (2022)**\n")
print(f"L'analyse de régression linéaire simple entre la
Consommation Énergétique (variable explicative, X) et le
P.I.B. (variable à expliquer, Y) pour l'année 2022 révèle une
relation extrêmement forte et positive. Le coefficient de
corrélation de Pearson ($R$) de {correlation_r:.4f} indique
que ces deux variables évoluent presque parfaitement dans le
même sens : plus un territoire consomme d'énergie, plus son
P.I.B. est élevé. Cette corrélation est confirmée par le
coefficient de détermination ($R^2$) de {r_squared:.4f}, ce
qui signifie que le modèle de régression linéaire simple
explique environ {r_squared*100:.2f}% de la variance du

```

P.I.B. des territoires en se basant uniquement sur leur consommation énergétique.\n")

```
print(f"L'équation de la droite de régression est : PIB $  
\\approx$ ({slope:.4e}) $\\times$ Énergie +  
({intercept:.4e}). La forte pente positive ({slope:.4e})  
traduit clairement que pour chaque augmentation d'unité de  
consommation énergétique, le P.I.B. augmente de manière  
significative. Bien que cette relation statistique ne prouve  
pas une causalité directe et doit être nuancée par d'autres  
facteurs socio-économiques, elle met en évidence un lien  
empirique puissant, typique des économies où la production de  
richesse est fortement dépendante de l'intrant énergétique.  
Le très faible $p$-value$ ({p_value:.4e}) permet de rejeter  
l'hypothèse nulle, confirmant que cette relation n'est pas  
due au hasard.")
```

```
# --- 2.3 Bonus : Généralisation à toutes les années (1962 à  
2022) ---
```

```
if not df.empty:
```

```
    print("\n\n--- ÉTAPE BONUS (2.3) : Généralisation à  
toutes les années ---")
```

```
    # Identification des colonnes pertinentes
```

```
    COL_ID = 'Territoire' # Colonne d'identification, si elle  
existe
```

```
    if 'Country Name' in df.columns:
```

```
        COL_ID = 'Country Name'
```

```
    elif 'Territoire' not in df.columns:
```

```
        # Assumons que le nom du pays/territoire est la  
première colonne si non standard
```

```
        COL_ID = df.columns[0]
```

```
    # Filtrer les colonnes pour ne garder que l'ID et les  
années de données
```

```
    data_cols = [col for col in df.columns if ('PIB_' in col  
or 'Utilisation_d_energie_' in col)]
```

```
    # Création d'une structure pour stocker les résultats  
annuels
```

```
    results = pd.DataFrame(columns=['Année',  
'Nombre_Observations', 'R_Correlation', 'R_Carre', 'Pente',  
'P_Value'])
```

```
    # Boucle sur les années de 1962 à 2022
```

```
    for annee in range(1962, 2023):
```

```

col_pib = f'PIB_{annee}'
col_energie = f'Utilisation_d_energie_{annee}'

# Vérifier si les colonnes existent pour l'année
if col_pib in df.columns and col_energie in
df.columns:

    # 1. Sélection et nettoyage des données pour
l'année courante
    df_annee = df[[col_pib, col_energie]].dropna()

    n_obs = len(df_annee)

    # 2. Effectuer la régression si un nombre
suffisant de données est disponible
    if n_obs >= 2: # Nécessite au moins 2 points pour
la régression
        X_annee = df_annee[col_energie]
        Y_annee = df_annee[col_pib]

        slope_annee, intercept_annee, r_value_annee,
p_value_annee, std_err_annee = stats.linregress(X_annee,
Y_annee)

        r_squared_annee = r_value_annee**2

    # 3. Stocker les résultats
new_row = pd.DataFrame([ {
    'Année': annee,
    'Nombre_Observations': n_obs,
    'R_Correlation': r_value_annee,
    'R_Carre': r_squared_annee,
    'Pente': slope_annee,
    'P_Value': p_value_annee
} ])

    # Concaténer la nouvelle ligne au DataFrame
des résultats
    results = pd.concat([results, new_row],
ignore_index=True)

    print(f"\nAnalyse de la relation PIB vs. Énergie sur
{len(results)} années (1962-2022) terminée.")

    # 4. Organisation et export du fichier de sortie

```

```

output_filename =
'resultats_regression_annuelle_PIB_Energie.xlsx'
results.to_excel(output_filename, index=False)

print(f"\nFichier de sortie créé : '{output_filename}'")

# Affichage des 5 premières et dernières années analysées
print("\n--- Aperçu des Résultats Annuels ---")
print("Années avec les corrélations les plus fortes:")
print(results.sort_values(by='R_Correlation',
ascending=False).head())

# Visualisation de l'évolution du R² au fil du temps
plt.figure(figsize=(12, 6))
sns.lineplot(x='Année', y='R_Carre', data=results)
plt.title("Évolution du Coefficient de Détermination
($R^2$) entre PIB et Énergie (1962-2022)")
plt.xlabel("Année")
plt.ylabel("$R^2$")
plt.show()

```

le terminal me renvoie : hsen@MBP-de-HSEN S7 % python3 pain.py

Traceback (most recent call last):

File "/Users/hsen/Desktop/S7/pain.py", line 4, in <module>

```
import seaborn as sns
```

ModuleNotFoundError: No module named 'seaborn'

hsen@MBP-de-HSEN S7 %

Il semble que le problème que vous rencontrez soit une erreur fréquente lorsque l'on utilise Python pour l'analyse de données, surtout si vous venez d'installer un environnement ou si vous travaillez sur une nouvelle machine.

Le message d'erreur est très clair :

ModuleNotFoundError: No module named 'seaborn'

Cela signifie que la bibliothèque **seaborn**, qui est utilisée dans votre code pour la visualisation graphique (étape 5), **n'est pas installée** dans l'environnement Python que vous utilisez.



Solution : Installation de la Bibliothèque

Pour corriger cette erreur, vous devez utiliser le gestionnaire de paquets de Python, appelé **pip**, pour installer **seaborn**.

Étape 1 : Ouvrir le Terminal Vous êtes déjà dans le Terminal, ce qui est parfait.

Étape 2 : Exécuter la commande d'installation Entrez la commande suivante et appuyez sur Entrée :

Bash

```
pip install seaborn
```

Si vous utilisez Python 3 et que la commande **pip install** ne fonctionne pas, il est possible que vous deviez utiliser **pip3** à la place (ce qui est courant sur macOS) :

Bash

```
pip3 install seaborn
```



Installation des autres dépendances

Puisque vous travaillez avec un code d'analyse de données qui utilise aussi **pandas**, **scipy** et **matplotlib**, assurez-vous que ces modules sont également installés. La meilleure pratique est de tout installer en une seule fois :

Bash

```
pip install pandas scipy matplotlib seaborn numpy
```

ou, si vous utilisez **pip3** :

Bash

```
pip3 install pandas scipy matplotlib seaborn numpy
```

Étape 3 : Réexécuter le Script Une fois l'installation terminée (le Terminal vous affichera un message de succès), relancez votre script :

Bash

```
python3 pain.py
```

Le programme devrait maintenant pouvoir importer toutes les bibliothèques et exécuter l'analyse.

Si vous rencontrez de nouvelles erreurs après cela, n'hésitez pas à me les partager !

voici l'emplacement du fichier csv : /Users/hsen/Desktop/S7/data le terminal me renvoie :hsen@MBP-de-HSEN S7 % python3 pain.py

--- Étape 1 : Sélection des colonnes PIB_2022 et Utilisation_d_energie_2022 ---

ERREUR : Le fichier 'pib-vs-energie.csv' est introuvable. Veuillez vérifier le chemin.

hsen@MBP-de-HSEN S7 %

C'est une erreur très courante, et elle est liée à l'emplacement de votre fichier de données par rapport à l'endroit où vous exécutez le script Python.

Le message d'erreur est sans équivoque :

ERREUR : Le fichier 'pib-vs-energie.csv' est introuvable.
Veuillez vérifier le chemin.

Dans le code que j'ai fourni, le programme cherche le fichier `pib-vs-energie.csv` dans le **même dossier** que le script `pain.py`. Or, d'après votre message, le fichier se trouve dans un sous-dossier appelé `data`.



Solution : Modifier le Chemin d'Accès au Fichier

Pour corriger cela, il suffit de modifier la ligne de code où le fichier est chargé pour indiquer le chemin exact où se trouve le fichier `pib-vs-energie.csv`.

Ouvrez le fichier `pain.py` et modifiez la ligne suivante :

Code original (Ligne 13 environ) :

Python

```
# Assurez-vous que le fichier 'pib-vs-energie.csv' est dans  
le même répertoire  
# que votre script Python.  
try:  
    df = pd.read_csv('pib-vs-energie.csv')
```

Nouveau Code Corrigé :

Remplacez `'pib-vs-energie.csv'` par le chemin d'accès complet ou le chemin relatif au dossier `data`.

Si votre fichier `pain.py` est dans `/Users/hsen/Desktop/S7` et que le fichier de données est dans `/Users/hsen/Desktop/S7/data`, le chemin relatif est le plus simple :

Python

```
# NOUVEAU CHEMIN D'ACCÈS CORRIGÉ :  
try:  
    # Le chemin indique maintenant de chercher le fichier  
dans le sous-dossier 'data'  
    df = pd.read_csv('data/pib-vs-energie.csv')
```

Récapitulatif des modifications :

1. Localisez la ligne de chargement du fichier dans `pain.py` (autour de la ligne 13).
2. Changez `pib-vs-energie.csv` en `data/pib-vs-energie.csv`.

Sauvegardez le fichier `pain.py` et relancez-le :

Bash


```
python3 pain.py
```

Le script devrait maintenant trouver le fichier, charger les données, et procéder à l'analyse.

voici mon code : import pandas as pd

```
from scipy import stats
```

```
import matplotlib.pyplot as plt
```

```
import seaborn as sns
```

```
import numpy as np
```

```
# --- 1. Chargement et sélection des données pour l'année 2022 ---
```

```
print("--- Étape 1 : Sélection des colonnes PIB_2022 et Utilisation_d_energie_2022 ---")
```

```
# Assurez-vous que le fichier 'pib-vs-energie.csv' est dans le même répertoire
```

```
# que votre script Python.
```

```
try:
```

```
    df = pd.read_csv('/Users/hsen/Desktop/S7/data/pib-vs-energie.csv')
```

```
except FileNotFoundError:
```

```
    print("ERREUR : Le fichier 'pib-vs-energie.csv' est introuvable. Veuillez vérifier le chemin.")
```

```
    # Créer un DataFrame fictif pour que le code puisse être affiché, mais l'analyse échouera.
```

```
    df = pd.DataFrame()
```

```
if not df.empty:
```

```
    # Colonnes à sélectionner selon l'énoncé
```

```
    COL_PIB = 'PIB_2022'
```

```
    COL_ENERGIE = 'Utilisation_d_energie_2022'
```

```

# Sélection du DataFrame de travail pour l'analyse
df_travail = df[[COL_PIB, COL_ENERGIE]]

print(f"DataFrame initial (avec données manquantes) : {len(df_travail)} observations.")
print(df_travail.head())

# --- 2. Création d'un algorithme pour filtrer les couples complets ---
print("\n--- Étape 2 : Suppression des données manquantes (couples complets) ---")

# La méthode .dropna() de Pandas est la manière la plus efficace de réaliser cette opération.
# Elle supprime toute ligne où il manque au moins une valeur dans les colonnes spécifiées.
df_complet = df_travail.dropna(subset=[COL_PIB, COL_ENERGIE])

# Affichage du résultat
n_complet = len(df_complet)
n_manquants = len(df_travail) - n_complet

print(f"Nombre d'observations complètes : {n_complet}")
print(f"Nombre d'observations exclues (manquantes ou partielles) : {n_manquants}")
print(df_complet.head())

# Préparation des variables (PIB est la variable à expliquer, Y ; Énergie est l'explicative, X)
# ATTENTION: L'énoncé demande (x, y) où x est la variable à expliquer et y l'explicative
# dans le linregress, ce qui est l'inverse de la convention mathématique ( $Y = aX + b$ ).
# Nous allons respecter l'énoncé pour le calcul, puis l'interpréter correctement pour le graphique.
# X_var_expli: Énergie (Consommation Énergétique)
# Y_var_a_expliquer: PIB

```

```
X = df_complet[COL_ENERGIE]
```

```
Y = df_complet[COL_PIB]
```

```
# --- 3. Calcul de la Régression Linéaire Simple (linregress) ---
```

```
print("\n--- Étape 3 : Calcul de la Régression Linéaire Simple ---")
```

```
# Calcul de la régression linéaire simple avec linregress(x, y)
```

```
# selon l'énoncé (x = Y, y = X) pour la cohérence avec le reste des étapes
```

```
# Note: linregress(x, y) retourne (pente, ordonnée_origine, R_value, p_value, std_err)
```

```
# Correction : Si l'énoncé inverse x et y dans linregress, nous allons
```

```
# nous fier à la logique mathématique pour l'interprétation future :  $Y = aX + b$ .
```

```
# Y (à expliquer) = PIB ; X (explicative) = Énergie
```

```
slope, intercept, r_value, p_value, std_err = stats.linregress(X, Y)
```

```
print(f"Variable Explicative (X) : Consommation Énergétique ({COL_ENERGIE})")
```

```
print(f"Variable à Expliquer (Y) : PIB ({COL_PIB})")
```

```
print(f"Pente (slope) : {slope:.4e}")
```

```
print(f"Ordonnée à l'origine (intercept) : {intercept:.4e}")
```

```
print(f"P-Value : {p_value:.4e}")
```

```
print(f"Erreur standard (Standard Error) : {std_err:.4e}")
```

```
# Formule de la droite de régression :  $PIB = \{slope\} * \text{Énergie} + \{intercept\}$ 
```

```
print(f"\nÉquation de la droite :  $PIB = ({slope:.4e}) * \text{Énergie} + ({intercept:.4e})$ ")
```

```

# --- 4. Calcul de la Corrélation Simple (R) ---

print("\n--- Étape 4 : Calcul de la Corrélation Simple (R et R²) ---")

# R-value est déjà calculée à l'étape 3 par linregress (corrélation de Pearson)

correlation_r = r_value

r_squared = correlation_r**2

print(f"Coefficient de Corrélation (R) : {correlation_r:.4f}")
print(f"Coefficient de Détermination (R²) : {r_squared:.4f}")

# --- 5. Graphique de synthèse avec la droite de régression ---

print("\n--- Étape 5 : Création du graphique de synthèse ---")

# Configuration du style graphique

sns.set_theme(style="whitegrid")

# Création du nuage de points

plt.figure(figsize=(10, 7))

sns.scatterplot(x=X, y=Y, label='Observations', color='blue', alpha=0.6)

# Ajout de la droite de régression

# Création des valeurs prédites pour tracer la droite

x_min, x_max = X.min(), X.max()

x_range = np.linspace(x_min, x_max, 100)

y_pred = slope * x_range + intercept

plt.plot(x_range, y_pred, color='red',

```

```
label=f'Droite de Régression\n($R^2 = {r_squared:.3f}$)',
```

```
linewidth=2.5)
```

```
# Étiquettes et titre
```

```
plt.title(f'Relation entre Consommation Énergétique et PIB en {2022}', fontsize=14)
```

```
plt.xlabel(f'Consommation Énergétique (kg équivalent pétrole) - {COL_ENERGIE}',  
fontsize=12)
```

```
plt.ylabel(f'Produit Intérieur Brut (dollars courants) - {COL_PIB}', fontsize=12)
```

```
plt.legend()
```

```
plt.grid(True, linestyle='--', alpha=0.7)
```

```
# Affichage du graphique
```

```
plt.show()
```

```
# --- 6. Commentaire du Résultat (À inclure dans votre rapport) ---
```

```
print("\n--- 6. Commentaire du Résultat (Pour le Rapport) ---")
```

```
print("\n**Analyse et Commentaire du Lien entre PIB et Consommation Énergétique  
(2022)**\n")
```

```
print(f"L'analyse de régression linéaire simple entre la Consommation Énergétique (variable  
explicative, X) et le P.I.B. (variable à expliquer, Y) pour l'année 2022 révèle une relation  
extrêmement forte et positive. Le coefficient de corrélation de Pearson ($R$) de {correlation_r:.4f}  
indique que ces deux variables évoluent presque parfaitement dans le même sens : plus un territoire  
consomme d'énergie, plus son P.I.B. est élevé. Cette corrélation est confirmée par le coefficient de  
détermination ($R^2$) de {r_squared:.4f}, ce qui signifie que le modèle de régression linéaire  
simple explique environ {r_squared*100:.2f}% de la variance du P.I.B. des territoires en se basant  
uniquement sur leur consommation énergétique.\n")
```

```
print(f"L'équation de la droite de régression est : PIB $\approx$ ({slope:.4e}) $\times$ Énergie  
+ ({intercept:.4e}). La forte pente positive ({slope:.4e}) traduit clairement que pour chaque  
augmentation d'unité de consommation énergétique, le P.I.B. augmente de manière significative.  
Bien que cette relation statistique ne prouve pas une causalité directe et doive être nuancée par  
d'autres facteurs socio-économiques, elle met en évidence un lien empirique puissant, typique des  
économies où la production de richesse est fortement dépendante de l'intrant énergétique. Le très  
faible $p$-value$ ({p_value:.4e}) permet de rejeter l'hypothèse nulle, confirmant que cette relation  
n'est pas due au hasard.")
```

```

# --- 2.3 Bonus : Généralisation à toutes les années (1962 à 2022) ---

if not df.empty:

    print("\n\n--- ÉTAPE BONUS (2.3) : Généralisation à toutes les années ---")

    # Identification des colonnes pertinentes

    COL_ID = 'Territoire' # Colonne d'identification, si elle existe

    if 'Country Name' in df.columns:

        COL_ID = 'Country Name'

    elif 'Territoire' not in df.columns:

        # Assumons que le nom du pays/territoire est la première colonne si non standard

        COL_ID = df.columns[0]

    # Filtrer les colonnes pour ne garder que l'ID et les années de données

    data_cols = [col for col in df.columns if ('PIB_' in col or 'Utilisation_d_energie_' in col)]

    # Création d'une structure pour stocker les résultats annuels

    results = pd.DataFrame(columns=['Année', 'Nombre_Observations', 'R_Correlation', 'R_Carre',
    'Pente', 'P_Value'])

    # Boucle sur les années de 1962 à 2022

    for annee in range(1962, 2023):

        col_pib = f'PIB_{annee}'

        col_energie = f'Utilisation_d_energie_{annee}'

        # Vérifier si les colonnes existent pour l'année

        if col_pib in df.columns and col_energie in df.columns:

```

```
# 1. Sélection et nettoyage des données pour l'année courante
```

```
df_annee = df[[col_pib, col_energie]].dropna()
```

```
n_obs = len(df_annee)
```

```
# 2. Effectuer la régression si un nombre suffisant de données est disponible
```

```
if n_obs >= 2: # Nécessite au moins 2 points pour la régression
```

```
    X_annee = df_annee[col_energie]
```

```
    Y_annee = df_annee[col_pib]
```

```
    slope_annee, intercept_annee, r_value_annee, p_value_annee, std_err_annee =  
stats.linregress(X_annee, Y_annee)
```

```
    r_squared_annee = r_value_annee**2
```

```
# 3. Stocker les résultats
```

```
new_row = pd.DataFrame([
```

```
    'Année': annee,
```

```
    'Nombre_Observations': n_obs,
```

```
    'R_Correlation': r_value_annee,
```

```
    'R_Carre': r_squared_annee,
```

```
    'Pente': slope_annee,
```

```
    'P_Value': p_value_annee
```

```
])
```

```
# Concaténer la nouvelle ligne au DataFrame des résultats
```

```
results = pd.concat([results, new_row], ignore_index=True)
```

```
print(f"\nAnalyse de la relation PIB vs. Énergie sur {len(results)} années (1962-2022)  
terminée.")
```

```
# 4. Organisation et export du fichier de sortie
```

```
output_filename = 'resultats_regression_annuelle_PIB_Energie.xlsx'
```

```
results.to_excel(output_filename, index=False)
```

```
print(f"\nFichier de sortie créé : '{output_filename}')
```

```
# Affichage des 5 premières et dernières années analysées
```

```
print("\n--- Aperçu des Résultats Annuels ---")
```

```
print("Années avec les corrélations les plus fortes:")
```

```
print(results.sort_values(by='R_Correlation', ascending=False).head())
```

```
# Visualisation de l'évolution du  $R^2$  au fil du temps
```

```
plt.figure(figsize=(12, 6))
```

```
sns.lineplot(x='Année', y='R_Carre', data=results)
```

```
plt.title("Évolution du Coefficient de Détermination ( $R^2$ ) entre PIB et Énergie (1962-2022)")
```

```
plt.xlabel("Année")
```

```
plt.ylabel(" $R^2$ ")
```

```
plt.show()
```

voila ce que me renvoie le terminal : hsen@MBP-de-HSEN S7 % python3 pain.py

--- Étape 1 : Sélection des colonnes PIB_2022 et Utilisation_d_energie_2022 ---

ERREUR : Le fichier 'pib-vs-energie.csv' est introuvable. Veuillez vérifier le chemin.

hsen@MBP-de-HSEN S7 %

Je vois que vous avez corrigé le chemin dans votre code en utilisant le chemin d'accès absolu /Users/hsen/Desktop/S7/data/pib-vs-energie.csv, mais le terminal vous renvoie toujours la même erreur :

**ERREUR : Le fichier 'pib-vs-energie.csv' est introuvable.
Veuillez vérifier le chemin.**

Cette situation est frustrante car elle suggère que, **même avec le chemin absolu**, Python ne parvient pas à localiser le fichier.

Il y a deux causes principales à ce problème, même après avoir corrigé le chemin :

1. Problème de Nom de Fichier ou d'Extension

Vérifiez **très attentivement** l'orthographe exacte et la casse du nom du fichier, ainsi que son extension, dans le Finder.

- **Vérifiez le nom :** Est-ce `pib-vs-energie.csv` ou pourrait-il être `PIB-vs-Energie.csv` ou `pib_vs_energie.csv` ?
- **Vérifiez l'extension :** Est-ce possible que le fichier ait été enregistré accidentellement avec une double extension, comme `pib-vs-energie.csv.csv` ?

Pour vérifier cela, allez dans le dossier /Users/hsen/Desktop/S7/data, faites un clic droit sur le fichier et choisissez "**Lire les informations**" (Get Info) pour voir le nom complet exact.

2. Erreur dans le Chemin Absolu

Bien que le chemin absolu soit souvent la méthode la plus sûre, il peut y avoir une coquille que vous ne voyez pas.

Action recommandée : Utiliser la méthode du Chemin Relatif (plus flexible)

Puisque vous exécutez le script dans le dossier /Users/hsen/Desktop/S7 (comme l'indique `hsen@MBP-de-HSEN S7 %`), et que les données sont dans le sous-dossier `data`, le chemin relatif est le plus simple et le moins sujet à erreur :

Dans votre fichier `pain.py` (ligne 13), utilisez le chemin relatif :

Python

```
# Ancienne ligne : df = pd.read_csv('/Users/hsen/Desktop/S7/
data/pib-vs-energie.csv')
# NOUVELLE LIGNE RECOMMANDÉE :
try:
    df = pd.read_csv('data/pib-vs-energie.csv')
except FileNotFoundError:
    print("ERREUR : Le fichier 'pib-vs-energie.csv' est
introuvable. Veuillez vérifier le chemin.")
# ...
```

Prochaine Étape : Débogage du Chemin

Si la correction du chemin relatif ne fonctionne pas, utilisez le code suivant pour **vérifier où Python pense se trouver** (le répertoire de travail) et **vérifier si le fichier existe** à partir de ce point.

Remplacez temporairement votre code de chargement par ceci pour un débogage rapide :

Python

```
import os
import pandas as pd

# Chemin que vous pensez être le bon (relatif à S7)
chemin_relatif = 'data/pib-vs-energie.csv'
# Chemin absolu utilisé précédemment
chemin_absolu = '/Users/hsen/Desktop/S7/data/pib-vs-
energie.csv'

# 1. Afficher le répertoire de travail actuel
print(f"Répertoire de travail actuel (où vous exécutez le
script) : {os.getcwd()}")

# 2. Vérifier si le fichier existe au chemin relatif
existe_relatif = os.path.exists(chemin_relatif)
print(f"Le fichier existe-t-il au chemin relatif
('{chemin_relatif}') ? {existe_relatif}")

# 3. Vérifier si le fichier existe au chemin absolu
existe_absolu = os.path.exists(chemin_absolu)
print(f"Le fichier existe-t-il au chemin absolu
('{chemin_absolu}') ? {existe_absolu}")

# Si les deux sont False, c'est que le nom du fichier est
incorrect.
if not existe_relatif and not existe_absolu:
```

```
print("\nATTENTION : Le fichier est introuvable aux deux chemins. VÉRIFIEZ L'ORTHOGRAPHE DU NOM DE FICHIER DANS LE DOSSIER 'data'.")
```

```
# Une fois que vous savez quel chemin est correct, reprenez l'analyse.
```

```
# Par exemple, si le chemin relatif est correct :
```

```
if existe_relatif:
```

```
    df = pd.read_csv(chemin_relatif)
```

```
    print("Fichier chargé avec succès en utilisant le chemin relatif.")
```

```
# else... (reprendre l'analyse si chargé)
```

Si `os.path.exists(chemin_absolu)` renvoie `False`, cela confirme qu'il y a une faute de frappe dans le nom du fichier ou du dossier.

Sélection des colonnes PIB_2022 et Utilisation_d_energie_2022 ---ERREUR : Le fichier 'pib-vs-energie.csv' est introuvable. Veuillez vérifier le chemin.hsen@MBP-de-HSEN S7 %

L'erreur persiste :

ERREUR : Le fichier 'pib-vs-energie.csv' est introuvable.
Veuillez vérifier le chemin.

Puisque nous avons déjà essayé le chemin absolu et le chemin relatif, le problème le plus probable est désormais **l'orthographe ou l'extension réelle du fichier** sur votre disque dur.

Même si vous êtes certain du nom, l'erreur de Python est catégorique : il ne trouve pas de fichier *exactement* nommé `pib-vs-energie.csv` dans le dossier `/Users/hsen/Desktop/S7/data`.

1. Diagnostic Final du Nom du Fichier (Vérification Manuelle)

Veuillez effectuer cette vérification finale dans le Finder :

1. Ouvrez le dossier : `/Users/hsen/Desktop/S7/data`.
 2. Trouvez le fichier CSV.
 3. Faites un **clic droit** sur le fichier, puis choisissez "**Lire les informations**" (ou "Get Info").
 4. Regardez la section "**Nom et extension**".
- **Cas 1 (Le plus fréquent) :** Le nom réel est `pib-vs-energie.csv.csv` (double extension masquée) ou `pib_vs_energie.csv` (underscore au lieu du tiret).
 - **Cas 2 :** Il y a une erreur de casse, par exemple `PIB-vs-Energie.csv`.

2. Solution : Utiliser le Chemin Corrigé (avec le nom exact)

Une fois que vous avez identifié le nom exact du fichier (par exemple, `pib_vs_energie.csv`), vous devez **corriger la ligne de chargement** dans votre script `pain.py`.

En supposant que le nom est, par exemple, `pib_vs_energie.csv` et que vous utilisez le chemin relatif (le plus propre) :

Python

```
# Modifiez cette ligne dans votre fichier pain.py
# Remplacez 'pib-vs-energie.csv' par le nom EXACT que vous
avez trouvé

try:
    df = pd.read_csv('data/pib_vs_energie.csv') # EXEMPLE de
correction si le tiret était un underscore
except FileNotFoundError:
    print("ERREUR : Le fichier 'pib-vs-energie.csv' est
introuvable. Veuillez vérifier le chemin.")
    # ...
```

Une fois que vous avez le nom exact du fichier dans le code, relancez le script :

Bash

```
python3 pain.py
```

Si le script ne fonctionne toujours pas, cela signifie qu'il y a un problème de permission ou que le fichier est corrompu, ce qui est beaucoup plus rare.

Fait un rapport de séance de l'exercice 2 en mettant les résultats en tableau et en commentant l'histogramme voici l'exercice : 2 Mise en œuvre avec Python

La sous-partie « Bonus » vous permet d'obtenir des points supplémentaires.

12.1 Objectifs

- Manipuler un fichier C.S.V.
- Faire des sorties graphiques
- Utiliser les bibliothèques Pandas (données) et Matplotlib (graphiques)

N.B. pd et plt sont des alias qui remplacent respectivement pandas et matplotlib.pyplot.

- Calculer des effectifs
- Calculer des fréquences
- Faire des graphiques (diagrammes en bâton et circulaires, et histogrammes)

2.2 Manipulations

Le fichier obtenu compte pour 10 % de la note finale du parcours « débutants ».

1. Dans le dossier src, créer un dossier data et y introduire le fichier

resultats-elections-presidentielles-2022-1er-tour.csv disponible

dans la Seance-02 du GitHub

2. Dans le dossier src, introduire le fichier main.py de la séance disponible dans la Seance-02 du GitHub

3. Ouvrir le fichier main.py dans votre éditeur de code (Notepad++ ou VS Code)

4. Repérer l'instruction with. Elle permet par l'intermédiaire d'une variable appelée fichier de lire le fichier C.S.V. grâce à la méthode read_csv(...) de la bibliothèque Pandas.

N.B. Bien que Pandas puisse lire les fichiers Excel, il faut vous habituer à utiliser des formats textuels comme le format C.S.V.

5. Afficher sur le terminal exécutant le contenu la variable contenu avec la méthode

de Pandas DataFrame(...)

6. Calculer avec la fonction native len(...) le nombre de lignes et de colonnes du tableau de données et les afficher sur le terminal

7. Faire le point sur la nature statistique des variables en utilisant le lien vers les métadonnées fourni en commentaire. Faire une liste sur le type de chaque colonne (int, float, str ou bool).

8. Afficher sur le terminal le nom des colonnes, c'est-à-dire la première ligne avec la méthode Pandas head()

9. À l'aide du nom des colonnes affiché, sélectionner le nombre des inscrits

10. Calculer avec la méthode Pandas sum(...) les effectifs de chaque colonne et les placer dans une liste (à l'aide d'une boucle). Afficher le résultat sur le terminal. Normalement, le résultat sera bizarre. Mettre une condition pour calculer les effectifs des colonnes contenant des données quantitatives en utilisant la liste faisant le point sur le type de variables.

Indice. Utiliser la méthode native append(...)

211. À l'aide de Matplotlib et d'une boucle, faire des diagrammes en barres avec le nombre des inscrits et des votants pour chaque département. Créer les fichiers images des diagrammes.

N.B. 1. Vous allez créer de nombreux fichiers. Créer un nouveau sous-dossier pour stocker vos images.

N.B. 2. Privilégier le format*.png, qui est un format léger non propriétaire

12. À l'aide Matplotlib et d'une boucle, faire des diagrammes circulaire avec les votes blancs, nuls, exprimés et l'abstention pour chaque département. Créer les fichiers images des diagrammes. Respecter les mêmes remarques.

13. À l'aide Matplotlib, faire l'histogramme de la distribution des inscrits

Je rappelle la différence entre un diagramme en bâtons et un histogramme. Le premier sert à représenter des données, le second, à représenter une distribution statistique. Cela signifie sur la totalité des surfaces des rectangles vaut 1 (cf. cours sur les distributions).

2.3 Bonus

Sans remarque pour vous aider (conditions réelles), faire les diagrammes circulaires visualisant, pour chaque département, le nombre de voix par candidat.

Calculer le diagramme circulaire pour l'ensemble de la France. Mon terminal me renvoie : Réponse du terminal :
==== Aperçu du contenu du CSV ====

	Code du département	Libellé du département	...	Prénom.11	Voix.11
0	01	Ain	...	Nicolas	8998.0
1	02	Aisne	...	Nicolas	5790.0
2	03	Allier	...	Nicolas	4216.0
3	04	Alpes-de-Haute-Provence	...	Nicolas	2504.0
4	05	Hautes-Alpes	...	Nicolas	2142.0
..
102	ZP	Polynésie française	...	Nicolas	1969.0
103	ZS	Saint-Pierre-et-Miquelon	...	Nicolas	82.0
104	ZW	Wallis et Futuna	...	Nicolas	244.0
105	ZX	Saint-Martin/Saint-Barthélemy	...	Nicolas	339.0
106	ZZ	Français établis hors de France	...	Nicolas	7074.0

[107 rows x 56 columns]

==== Dimensions du tableau ====

Lignes : 107

Colonnes : 56

==== Types des colonnes ====

Code du département : object

Libellé du département : object

Inscrits : int64

Abstentions : float64

Votants : float64

Blancs : float64

Nuls : float64

Exprimés : float64

Sexe : object

Nom : object

Prénom : object

Voix : float64

Sexe.1 : object

Nom.1 : object

Prénom.1 : object

Voix.1 : float64

Sexe.2 : object

Nom.2 : object

Prénom.2 : object

Voix.2 : float64

Sexe.3 : object

Nom.3 : object

Prénom.3 : object

Voix.3 : float64

Sexe.4 : object

Nom.4 : object

Prénom.4 : object

Voix.4 : float64

Sexe.5 : object

Nom.5 : object

Prénom.5 : object

Voix.5 : float64

Sexe.6 : object

Nom.6 : object

Prénom.6 : object

Voix.6 : float64

Sexe.7 : object

Nom.7 : object

Prénom.7 : object

Voix.7 : float64

Sexe.8 : object

Nom.8 : object

Prénom.8 : object

Voix.8 : float64

Sexe.9 : object

Nom.9 : object

Prénom.9 : object

Voix.9 : float64

Sexe.10 : object

Nom.10 : object

Prénom.10 : object

Voix.10 : float64

Sexe.11 : object

Nom.11 : object

Prénom.11 : object

Voix.11 : float64

=== Noms des colonnes ===

Empty DataFrame

Columns: [Code du département, Libellé du département, Inscrits, Abstentions, Votants, Blancs, Nuls, Exprimés, Sexe, Nom, Prénom, Voix, Sexe.1, Nom.1, Prénom.1, Voix.1, Sexe.2, Nom.2, Prénom.2, Voix.2, Sexe.3, Nom.3, Prénom.3, Voix.3, Sexe.4, Nom.4, Prénom.4, Voix.4, Sexe.5, Nom.5, Prénom.5, Voix.5, Sexe.6, Nom.6, Prénom.6, Voix.6, Sexe.7, Nom.7, Prénom.7, Voix.7, Sexe.8, Nom.8, Prénom.8, Voix.8, Sexe.9, Nom.9, Prénom.9, Voix.9, Sexe.10, Nom.10, Prénom.10, Voix.10, Sexe.11, Nom.11, Prénom.11, Voix.11]

Index: []

=== Colonne 'Inscrits' ===

0 438109

1 373544

2 249991

3 128075

4 113519

...

102 205576

103 5045

104 9528

105 24414

106 1435746

Name: Inscrits, Length: 107, dtype: int64

=== Sommes des colonnes numériques ===

Inscrits : 48747876

Abstentions : 12824169.0

Votants : 35923707.0

Blancs : 543609.0

Nuls : 247151.0

Exprimés : 35132947.0

Voix : 197094.0

Voix.1 : 802422.0

Voix.2 : 9783058.0

Voix.3 : 1101387.0

Voix.4 : 8133828.0

Voix.5 : 2485226.0

Voix.6 : 7712520.0

Voix.7 : 616478.0

Voix.8 : 1627853.0

Voix.9 : 1679001.0

Voix.10 : 268904.0

Voix.11 : 725176.0

☒ Diagrammes en barres créés dans le dossier images/barres/

☒ Diagrammes circulaires créés dans le dossier images/camemberts/

☒ Histogramme des inscrits créé dans images/histogramme_inscrits.png

-
- Requête réussie

Voici le rapport de séance pour l'exercice 2, présentant les résultats des manipulations dans des tableaux et incluant l'analyse de l'histogramme de la distribution du nombre d'inscrits que vous avez fourni.

Rapport de Séance – Analyse des Résultats des Élections Présidentielles 2022 (1er Tour)

Ce rapport synthétise les résultats des manipulations effectuées avec Python (`pandas` et `matplotlib`) sur le fichier `resultats-elections-presidentielles-2022-1er-tour.csv`.

1. Structure et Nature des Données (Étapes 6 à 9)

Le jeu de données se présente sous la forme d'un tableau bidimensionnel.

Tableau 1 : Dimensions du DataFrame

Caractéristique	Valeur	Commentaire
Nombre de Lignes	107	Correspond au nombre d'entités territoriales analysées (départements métropolitains, DROM-COM et Français de l'étranger).
Nombre de Colonnes	56	Représente les informations générales du scrutin et les résultats des 12 candidats (sexe, nom, prénom, et voix).

Tableau 2 : Aperçu des Types de Variables

L'analyse des types de données (`dtypes`) permet de distinguer les variables qualitatives des quantitatives :

Type de	Exemples de Colonnes	Nature Statistique
object	Libellé du département, Code du département, Nom,	Qualitatif (Nominal ou Ordinal)
int64	Inscrits	Quantitatif (Discret)
float64	Abstentions, Votants, Blancs, Nuls, Exprimés, Voix	Quantitatif (Discret, stocké en <code>float</code> à cause des valeurs manquantes ou pour des raisons de

2. Calcul des Effectifs Totaux (Étape 10)

L'application conditionnelle de la méthode `sum()` aux colonnes numériques (`int64` et `float64`) permet de calculer les effectifs globaux pour l'ensemble des 107 entités analysées.

Tableau 3 : Sommes des Colonnes Quantitatives

Caractéristique	Effectif	Interprétation
Inscrits	48 747 876	Total des personnes inscrites sur les listes électorales.
Abstentions	12 824 169	Total des abstentionnistes.

Votants	35 923 707	Total des votants.
Blancs	543 609	Total des votes blancs.
Nuls	247 151	Total des votes nuls.
Exprimés	35 132 947	Total des suffrages valablement exprimés.
Voix (Total candidats)	35 132 947.0*	<i>La somme des voix des 12 candidats est égale au nombre total de suffrages exprimés.</i>

3. Analyse de l'Histogramme de la Distribution des Inscrits (Étape 13)

L'histogramme représente la **distribution statistique** du nombre d'inscrits par entité territoriale. Puisque c'est un histogramme de fréquence (et non un simple diagramme en bâtons), l'aire totale des rectangles est égale à 1.

Commentaire de l'Histogramme :

L'histogramme révèle une distribution fortement **asymétrique positive (ou étalée à droite)**.

1. **Concentration des Effectifs** : La majorité des entités territoriales (départements) ont un nombre d'inscrits **faible à modéré**. La classe de fréquence la plus élevée (le mode) se situe autour de **0.15×10⁶ à 0.25×10⁶** inscrits (soit 150 000 à 250 000 inscrits). Cela suggère qu'un grand nombre d'entités étudiées sont des départements de taille moyenne ou des collectivités d'outre-mer.
2. **Asymétrie et Valeurs Extrêmes** : La queue de distribution s'étend loin vers la droite, indiquant la présence de quelques entités avec un **nombre d'inscrits exceptionnellement élevé** (des "outliers" ou valeurs extrêmes). Ces pics de faible fréquence autour de 1.25×10⁶ à 1.75×10⁶ correspondent probablement aux départements les plus peuplés de France (comme le Nord, les Bouches-du-Rhône, ou le département des Français établis hors de France).
3. **Conclusion sur la Distribution** : La distribution **n'est pas Normale (Gaussienne)**. L'asymétrie positive est typique des variables socio-économiques et démographiques, indiquant que la moyenne arithmétique sera supérieure à la médiane (et au mode), car elle est tirée vers la droite par les grandes entités.

4. Réussite des Sorties Graphiques (Étapes 11, 12 et 13)

Toutes les sorties graphiques demandées ont été générées avec succès :

- **Diagrammes en barres** : Fichiers créés dans `images/barres/` (pour inscrits et votants par département).
- **Diagrammes circulaires** : Fichiers créés dans `images/camemberts/` (pour blancs, nuls, exprimés, abstention par département).
- **Histogramme** : Fichier créé en `images/histogramme_inscrits.png`.

Souhaitez-vous que je vous aide à rédiger les algorithmes pour la section **Bonus** (calcul des voix par candidat pour chaque département et pour l'ensemble de la France) ?

rédige de manière plus humaine : Commentaire de l'Histogramme :

L'histogramme révèle une distribution fortement asymétrique positive (ou étalée à droite).

Concentration des Effectifs : La majorité des entités territoriales (départements) ont un nombre d'inscrits faible à modéré. La classe de fréquence la plus élevée (le mode) se situe autour de 0.15×10^6

à 0.25×10^6

inscrits (soit 150 000 à 250 000 inscrits). Cela suggère qu'un grand nombre d'entités étudiées sont des départements de taille moyenne ou des collectivités d'outre-mer.

Asymétrie et Valeurs Extrêmes : La queue de distribution s'étend loin vers la droite, indiquant la présence de quelques entités avec un nombre d'inscrits exceptionnellement élevé (des "outliers" ou valeurs extrêmes). Ces pics de faible fréquence autour de 1.25×10^6

à 1.75×10^6

correspondent probablement aux départements les plus peuplés de France (comme le Nord, les Bouches-du-Rhône, ou le département des Français établis hors de France).

Conclusion sur la Distribution : La distribution n'est pas Normale (Gaussienne). L'asymétrie positive est typique des variables socio-économiques et démographiques, indiquant que la moyenne arithmétique sera supérieure à la médiane (et au mode), car elle est tirée vers la droite par les grandes entités.

Absolument. Voici une version plus rédigée et "humaine" de l'analyse de cet histogramme, mettant l'accent sur l'interprétation concrète des données.



Interprétation de la Distribution du Nombre d'Inscrits

En observant l'histogramme, on constate immédiatement que la répartition du nombre d'inscrits par entité territoriale est loin d'être symétrique. Elle présente une **forte asymétrie positive**, ce qui signifie que la majeure partie des données est concentrée à gauche, et qu'elle s'étire en une longue "queue" vers la droite.



La Concentration des Départements "Moyens"

La lecture des barres les plus hautes nous indique où se trouve le **cœur des effectifs**. La grande majorité des départements se regroupent dans la tranche des **150 000 à 250 000 inscrits**. Cette concentration suggère que la *majorité* des entités analysées, qu'il s'agisse des départements français de taille moyenne ou des collectivités d'outre-mer, ont une population d'électeurs relativement modeste.



L'Impact des "Poids Lourds"

Le phénomène le plus marquant est l'étalement vers la droite. Cette longue "queue" révèle que, malgré la forte concentration des départements moyens, il existe quelques entités avec un **nombre d'inscrits exceptionnellement élevé** (ce que l'on appelle statistiquement des valeurs extrêmes ou *outliers*). Les faibles pics de fréquence visibles au-delà de **1 250 000 inscrits** sont le reflet statistique de ces "poids lourds" : il s'agit sans doute des départements les plus peuplés de France métropolitaine (comme le Nord ou les Bouches-du-Rhône) ou de l'entité regroupant les Français établis hors de France, qui mobilisent un très grand nombre d'électeurs.



Conclusion sur la Forme de la Distribution

Cette configuration confirme que la distribution **n'est pas de type Normal (Gaussien)**. Cette asymétrie positive est typique des variables démographiques ou socio-économiques, où une petite fraction d'individus ou d'entités détient une part disproportionnée de la variable (ici, les inscrits). Statistiquement, cela a pour conséquence que la **moyenne arithmétique** sera artificiellement tirée vers le haut par ces gros départements et sera donc **supérieure à la médiane**, cette dernière étant une mesure de position plus représentative du département "typique" dans la série.

Comment le graphique : voici les consignes de la 3 : 2.2 Manipulations

Le fichier obtenu compte pour 10 % de la note finale du parcours « débutants ».

1. Dans le dossier src, créer un dossier data et y introduire le fichier

resultats-elections-presidentielles-2022-1er-tour.csv disponible

dans la Seance-03 du GitHub

2. Dans le dossier src, introduire le fichier main.py de la séance disponible dans la

Seance-03 du GitHub

3. Ouvrir le fichier main.py dans votre éditeur de code (Notepad++ ou VS Code)

4. Utiliser l'instruction with pour ouvrir le fichier C.S.V. avec la méthode read_csv(...)

de la bibliothèque Pandas

5. En reprenant le code précédemment créé dans la séance 2, sélectionner les colonnes contenant des caractères quantitatifs? Calculer sous la forme d'une liste :

— les moyennes de chaque colonne avec la bonne méthode de Pandas ;

— les médianes de chaque colonne avec la bonne méthode de Pandas ;

— les modes de chaque colonne avec la bonne méthode de Pandas ;

— l'écart type de chaque colonne avec la bonne méthode de Pandas ;

— l'écart absolu à la moyenne de chaque colonne avec la bonne méthode de Pandas ;

— l'étendue de chaque colonne.

N.B. 1. Utiliser la méthode de la valeur absolue abs() disponible dans Numpy

N.B. 2. Utiliser les méthodes min() et max() disponibles dans Pandas

En utilisant la méthode round() de Pandas, arrondir tous les paramètres à deux décimaux.

6. Afficher la liste des paramètres sur le terminal

7. Calculer la distance interquartile et interdécile de chaque colonne quantitative avec la méthode quantile() dans Pandas ?

8. À l'aide de Matplotlib et d'une boucle, faire des boîtes à moustache de chaque colonne quantitative. Stocker les résultats dans un dossier img.

29. Dans le dossier src, introduire le dossier data le fichier

island-index.csv disponible dans la Seance-03 du GitHub

10. Sélectionner la colonne « Surface (km²) » et écrire un algorithme pour catégoriser et dénombrer le nombre d'îles ayant une surface comprise :

entre 0 et 10 km² ou]0, 10] ;

entre 10 et 25 km² ou]10, 25] ;

entre 25 et 50 km² ou]25, 50] ;

entre 50 et 100 km² ou]50, 100] ;

entre 100 et 2 500 km² ou]100, 2500] ;

entre 2 500 et 5 000 km² ou]2500, 5000] ;

entre 5 000 et 10 000 km² ou]5000, 10000] ;

supérieur ou égal 10 000 km² ou]10000, +∞[.

Vous concevrez un organigramme pour expliquer votre solution. L'objectif de cette dernière question est d'apprendre à catégoriser des variables quantitatives.

2.3 Bonus

Sans remarque pour vous aider (conditions réelles), sortir les listes calculées avec la bonne méthode Pandas au format C.S.V. et Excel. N'oubliez pas de titrer les colonnes et les lignes de vos sorties si nécessaire.

-
- Requête réussie

Le graphique fourni dans le deuxième fichier () est un **Histogramme de fréquence** illustrant la distribution du nombre d'inscrits par entité territoriale (départements et autres circonscriptions) issue du fichier des résultats des élections présidentielles de 2022.

Voici le commentaire de ce graphique, rédigé de manière fluide et humaine :



Analyse de la Distribution du Nombre d'Inscrits

Le graphique ci-dessus, qui représente la distribution du nombre d'inscrits par circonscription territoriale, montre clairement que les électeurs ne sont pas répartis de manière uniforme en France.



Une Distribution Asymétrique et Concentrée

- **Le Cœur des Électeurs (Mode) :** La distribution est très loin d'être symétrique. Elle est fortement **étalée vers la droite (asymétrie positive)**. La grande majorité des entités se concentrent dans les classes de faible effectif. On constate que la fréquence la plus élevée se situe entre **100 000 et 250 000 inscrits** (soit environ 0.1×10^6 à 0.25×10^6). Ceci représente le département "typique" ou moyen en termes de population électorale.
- **L'Effet des "Poids Lourds" :** La longue queue qui s'étend vers les valeurs élevées est très révélatrice. Elle signale la présence de quelques entités avec un **nombre d'inscrits exceptionnellement important** (des valeurs extrêmes). Ces pics isolés au-delà de **1 250 000 inscrits** (1.25×10^6) correspondent aux départements les plus peuplés (comme ceux en Île-de-France, dans le Nord ou les Bouches-du-Rhône) ou au département des Français de l'étranger.



Conséquences Statistiques

Puisque la distribution est tirée vers la droite par ces grandes entités :

- La **moyenne arithmétique** des inscrits sera artificiellement élevée et ne sera pas un bon représentant du centre.
- La **médiane** (la valeur qui divise l'échantillon en deux moitiés) sera plus faible que la moyenne, car elle n'est pas affectée par les valeurs extrêmes.

En conclusion, cette distribution est typique des variables démographiques où la plupart des unités sont de petite taille, mais quelques unités gigantesques concentrent une part disproportionnée de l'effectif total.



Note sur l'Analyse de la Surface des Îles (Exercice 2.2, Question 10)

Le premier graphique fourni () illustre la répartition du nombre d'îles par intervalle de surface (en km^2), correspondant à la question 10 de la Séance 3 (et 2.2-9/10 de la Séance 4).

Ce graphique montre une concentration des îles encore plus extrême :

- **Ultra-Concentration :** Près de **80 000 îles** (la quasi-totalité de l'échantillon) ont une surface comprise entre **0 et 10 km^2 ($]0, 10]$)**.
- **Rareté des Grandes Îles :** Les barres représentant les classes de surface supérieures (par exemple, $]100, 2500]$ ou $]10000, +\infty[$) sont à peine visibles au-dessus de l'axe des abscisses.

Ceci confirme que les îles suivent une **distribution L-shape** (en forme de L) où les petites îles sont extrêmement nombreuses et les très grandes îles sont extrêmement rares, ce qui est caractéristique des phénomènes naturels hiérarchiques.

voila mon code : #coding:utf8

```
import pandas as pd
```

```
import math
```

```
import scipy
```

```
import scipy.stats
```

```
def ouvrirUnFichier(nom):
```

```
    with open(nom, "r") as fichier:
```

```
        contenu = pd.read_csv(fichier)
```

```
    return contenu
```

```
#
```

```
print("Résultat sur le calcul d'un intervalle de fluctuation")
```

```
donnees = pd.DataFrame(ouvrirUnFichier("./Users/hsen/Desktop/data/Echantillonnage-100-  
Echantillons.csv"))
```

```
#
```

```
print("Résultat sur le calcul d'un intervalle de confiance")
```

```
#
```

```
print("Théorie de la décision")
```

```
import pandas as pd
```

```
import numpy as np
```

```
from scipy import stats
```

```
# Fonction de chargement ---
```

```
# Chemin d'accès
```

```
FILE_PATH_ECHANTILLONS = "/Users/hsen/Desktop/STT/data/Echantillonnage-100-  
Echantillons.csv"
```

```
FILE_PATH_LOI1 = "/Users/hsen/Desktop/STT/data/Loi-normale-Test-1.csv"
```

```
FILE_PATH_LOI2 = "/Users/hsen/Desktop/STT/data/Loi-normale-Test-2.csv"
```

```
def ouvrirUnFichier(chemin_fichier):
```

```
    """
```

```
    Ouvre un fichier CSV en forçant le séparateur ';' (caractéristique des fichiers Excel/CSV français).
```

```
    """
```

```
    try:
```

```
        # Tente la lecture avec le point-virgule (séparateur le plus probable après votre description)
```

```
        df = pd.read_csv(chemin_fichier, sep=';')
```

```
        # Si le fichier ne contient qu'une seule colonne, le point-virgule n'était pas le bon séparateur.
```

```
        # On essaie la virgule, ou un espace blanc comme dernière option.
```

```
        if df.shape[1] == 1:
```

```
            try:
```

```
                # Tente la virgule (séparateur standard anglo-saxon)
```

```

df = pd.read_csv(chemin_fichier, sep=',')

except:

    # Tente l'espace blanc (sep=r'\s+')

    df = pd.read_csv(chemin_fichier, sep=r'\s+', engine='python', skipinitialspace=True)


# Nettoyage et attribution des noms de colonnes

df = df.dropna(axis=1, how='all')


if 'Echantillonnage' in chemin_fichier and df.shape[1] == 3:

    df.columns = ['Pour', 'Contre', 'Sans_opinion']


# Vérification finale

if df.shape[1] != 3 and 'Echantillonnage' in chemin_fichier:

    print(f"ATTENTION: Le fichier {chemin_fichier} n'a pu être lu qu'avec {df.shape[1]} colonnes. Le séparateur est peut-être incorrect.")


return df


except FileNotFoundError:

    print(f"ERREUR: Le fichier {chemin_fichier} n'a pas été trouvé. Veuillez vérifier le chemin.")

    return None


except Exception as e:

    print(f"ERREUR lors de l'ouverture du fichier {chemin_fichier}: {e}")

    return None


# Chargement principal

df_echantillons = ouvrirUnFichier(FILE_PATH_ECHANTILLONS)

```

```
# --- 1. Théorie de l'échantillonnage ---
```

```
print("="*60)
```

```
print("ÉTAPE 1 : THÉORIE DE L'ÉCHANTILLONNAGE")
```

```
print("="*60)
```

```
if df_echantillons is not None:
```

```
    # Données de la population mère
```

```
    N_mere = 2185
```

```
    pour_mere = 852
```

```
    contre_mere = 911
```

```
    sans_opinion_mere = 422
```

```
    # 1.1 Calcul des moyennes des échantillons (arrondies à l'entier)
```

```
    moyennes_echantillons = df_echantillons.mean().round(0).astype(int)
```

```
    print("\n[1.1] Moyennes des 100 échantillons (arrondies) :")
```

```
    print(moyennes_echantillons)
```

```
    # 1.2 Fréquences des échantillons (fréquences moyennes)
```

```
    somme_moyennes = moyennes_echantillons.sum()
```

```
    frequences_echantillons = (moyennes_echantillons / somme_moyennes).round(2)
```

```
    print("\n[1.2] Fréquences des moyennes des échantillons (arrondies à 2 décimales) :")
```

```
    print(frequences_echantillons)
```

```
    # 1.3 Fréquences de la population mère
```

```
    effectifs_mere = pd.Series([pour_mere, contre_mere, sans_opinion_mere],
```

```
                               index=['Pour', 'Contre', 'Sans_opinion'])
```

```

frequences_mere = (effectifs_mere / N_mere).round(2)

print("\n[1.3] Fréquences de la population mère (arrondies à 2 décimales) :")

print(frequences_mere)

```

```

# 1.4 Intervalle de fluctuation (IF)

```

```

# Formule :  $[p - z_C * \sqrt{p(1-p)/n}, p + z_C * \sqrt{p(1-p)/n}]$ 

```

```

z_C = 1.96 # Seuil de 95%

```

```

n_moyen = somme_moyennes # Taille moyenne des échantillons

```

```

print(f"\n[1.4] Intervalle de fluctuation (IF) à 95% (n moyen={n_moyen}) :")

```

```

intervalle_fluctuation = {}

```

```

for opinion, p in frequences_mere.items():

```

```

    # L'écart-type d'échantillonnage

```

```

    ecart_type_echantillonnage = np.sqrt(p * (1 - p) / n_moyen)

```

```

    margeErreur = z_C * ecart_type_echantillonnage

```

```

    borne_inf = p - margeErreur

```

```

    borne_sup = p + margeErreur

```

```

    intervalle_fluctuation[opinion] = (f"{borne_inf:.4f}", f"{borne_sup:.4f}")

```

```

    print(f"IF '{opinion}' (p={p:.2f}): [{borne_inf:.4f} ; {borne_sup:.4f}]")

```

```

# --- 2. Théorie de l'estimation ---

```

```

print("\n"+"="*60)

```

```

print("ÉTAPE 2 : THÉORIE DE L'ESTIMATION")

```

```
print("="*60)
```

```
if df_echantillons is not None:
```

```
    # 2.1 Prendre le premier échantillon (méthode iloc(0))
```

```
    premier_echantillon_pandas = df_echantillons.iloc[0]
```

```
    # Conversion en list() (comme demandé)
```

```
    premier_echantillon_list = premier_echantillon_pandas.tolist()
```

```
    print(f"\n[2.1] Premier échantillon (liste) : {premier_echantillon_list}")
```

```
    # 2.2 et 2.3 Calcul de l'effectif total (n) et des fréquences (p_hat)
```

```
    n_echantillon_1 = premier_echantillon_pandas.sum()
```

```
    print(f"Taille de l'échantillon (n) : {n_echantillon_1}")
```

```
    frequences_echantillon_1 = (premier_echantillon_pandas / n_echantillon_1)
```

```
    print("\n[2.3] Fréquences du premier échantillon (p_hat) :")
```

```
    print(frequences_echantillon_1.round(4))
```

```
    # 2.4 Intervalle de confiance (IC)
```

```
    # Formule :  $[p\_hat - z\_C * \sqrt{p\_hat(1-p\_hat)/n}, p\_hat + z\_C * \sqrt{p\_hat(1-p\_hat)/n}]$ 
```

```
    z_C = 1.96
```

```
    print(f"\n[2.4] Intervalle de confiance (IC) à 95% (n={n_echantillon_1}) :")
```

```
    for opinion, p_hat in frequences_echantillon_1.items():
```

```
        # L'écart-type d'estimation
```

```
        ecart_type_estimation = np.sqrt(p_hat * (1 - p_hat) / n_echantillon_1)
```

```
        marge_erreur = z_C * ecart_type_estimation
```



```
borne_inf = p_hat - marge_erreur
```

```
borne_sup = p_hat + marge_erreur
```

```
print(f"IC '{opinion}' (p_hat={p_hat:.4f}): [{borne_inf:.4f} ; {borne_sup:.4f}]" )
```

```
# --- 3 Test de Shapiro-Wilks) ---
```

```
print("\n"+"="*60)
```

```
print("ÉTAPE 3 : THÉORIE DE LA DÉCISION (Test de Shapiro-Wilks)")
```

```
print("="*60)
```

```
df_loi1 = ouvrirUnFichier(FILE_PATH_LOI1)
```

```
df_loi2 = ouvrirUnFichier(FILE_PATH_LOI2)
```

```
alpha = 0.05
```

```
def realiser_shapiro_wilks(df, nom_fichier, alpha):
```

```
    """Effectue le test de Shapiro-Wilks et affiche la conclusion."""
```

```
    if df is not None:
```

```
        # Assumer que la série de nombres est dans la première colonne
```

```
        data = df.iloc[:, 0].dropna()
```

```
        if len(data) < 3 or len(data) > 5000:
```

```
            # Shapiro-Wilks est moins fiable en dehors de cette plage
```

```
            print(f"ATTENTION: Le test de Shapiro-Wilks est moins adapté pour {nom_fichier}  
(taille={len(data)}).")
```

```
        stat_sw, p_value = stats.shapiro(data)
```

```
print(f"\nRésultats pour {nom_fichier} (Taille: {len(data)})")
```

```
print(f"Statistique W : {stat_sw:.4f}")
```

```
print(f"P-value : {p_value:.4f}")
```

```
if p_value > alpha:
```

```
    conclusion = f"P-value ({p_value:.4f}) > alpha ({alpha}). Ne rejette pas H0."
```

```
    resultat = "La distribution est considérée comme NORMALE."
```

```
else:
```

```
    conclusion = f"P-value ({p_value:.4f}) <= alpha ({alpha}). Rejette H0."
```

```
    resultat = "La distribution n'est PAS considérée comme normale."
```

```
print(f"Conclusion : {conclusion}")
```

```
print(f"Résultat : {resultat}")
```

```
return resultat
```

```
return "Non exécuté"
```

```
res_loi1 = realiser_shapiro_wilks(df_loi1, "Loi-normale-Test-1.csv", alpha)
```

```
res_loi2 = realiser_shapiro_wilks(df_loi2, "Loi-normale-Test-2.csv", alpha)
```

```
print("\n"+"="*60)
```

```
print("FIN DES CALCULS STATISTIQUES.")
```

```
print("="*60)
```

```
# --- BONUS ---
```

```
# --- 2.3 Analyse Approfondie pour le Bonus ---
```

```
print("\n"+"="*60)
```

```
print("ANALYSE BONUS : LOIS NON NORMALES")
```

```
print("="*60)
```

```
# Les deux lois ont été jugées NON NORMALES. Nous allons analyser leurs statistiques.
```

```
def analyser_distribution(df, nom):
```

```
    """Calcule statistiques descriptives clés pour aider à l'identification."""
```

```
    if df is not None:
```

```
        data = df.iloc[:, 0].dropna()
```

```
        statistiques = {
```

```
            "Taille (n)": len(data),
```

```
            "Minimum": data.min(),
```

```
            "Maximum": data.max(),
```

```
            "Étendue (Max - Min)": data.max() - data.min(),
```

```
            "Moyenne": data.mean(),
```

```
            "Médiane": data.median(),
```

```
            "Écart-type": data.std(),
```

```
        }
```

```
        print(f"\nStatistiques descriptives pour {nom} :")
```

```
        for cle, valeur in statistiques.items():
```

```
            print(f"- {cle:<20}: {valeur:.4f}")
```

```
        print("\n=> Caractérisation (pour le Bonus) :")
```

```
    # Test 1 : Forte asymétrie
```

```
    if abs(statistiques["Moyenne"] - statistiques["Médiane"]) > 0.1 * statistiques["Écart-type"]:
```



```
FileNotFoundError: [Errno 2] No such file or directory: './Users/hsen/Desktop/data/Echantillonnage-100-Echantillons.csv'
```

```
hsen@MacBook-Pro-de-HSEN STT %
```

L'erreur que vous rencontrez, `FileNotFoundError`, est classique et se produit lorsque le chemin d'accès au fichier est incorrect.

Le Problème : Chemin d'Accès Mal Formé

Dans votre code, vous avez la ligne suivante :

Python

```
donnees = pd.DataFrame(ouvrirUnFichier("./Users/hsen/Desktop/data/Echantillonnage-100-Echantillons.csv"))
```

Le chemin commence par `./` (qui signifie "dans le répertoire courant") suivi d'un chemin **absolu** (`/Users/hsen/Desktop/data/Echantillonnage...`).

- Un chemin ne peut pas être à la fois relatif (`./`) et absolu (`/Users/...`).
- Python interprète votre chemin comme : **{Répertoire courant}/Users/hsen/Desktop/data/Echantillonnage...**, ce qui n'existe évidemment pas.

La Solution : Utiliser le Chemin Absolu

La partie du code qui pose problème se trouve dans la première définition de `ouvrirUnFichier` (lignes 7-10) et son appel (ligne 16).

Vous avez déjà corrigé le chemin dans la deuxième fonction `ouvrirUnFichier` (qui est la version plus complexe et complète que nous avons utilisée pour l'analyse) en utilisant des variables de chemin absolu :

Python

```
FILE_PATH_ECHANTILLONS = "/Users/hsen/Desktop/STT/data/  
Echantillonnage-100-Echantillons.csv"  
# ...  
df_echantillons = ouvrirUnFichier(FILE_PATH_ECHANTILLONS)  
Votre code contient un bloc ancien qui doit être supprimé ou corrigé.
```

1. Supprimez ou commentez le bloc de code obsolète

Les lignes 4 à 16 du début de votre script sont redondantes et provoquent l'erreur car la deuxième partie du script gère déjà le chargement.

Action : Supprimez ou commentez les lignes suivantes (lignes 4 à 17 dans votre code) :

Python

```
def ouvrirUnFichier(nom):  
    with open(nom, "r") as fichier:  
        contenu = pd.read_csv(fichier)  
    return contenu  
  
#  
print("Résultat sur le calcul d'un intervalle de  
fluctuation")  
  
donnees = pd.DataFrame(ouvrirUnFichier("./Users/hsen/Desktop/  
data/Echantillonnage-100-Echantillons.csv"))  
#...  
print("Résultat sur le calcul d'un intervalle de confiance")  
  
#
```

2. Assurez-vous que les chemins absolus des variables **FILE_PATH** sont corrects

Votre deuxième bloc de code utilise les chemins absolus suivants :

Python

```
FILE_PATH_ECHANTILLONS = "/Users/hsen/Desktop/STT/data/  
Echantillonnage-100-Echantillons.csv"  
FILE_PATH_LOI1 = "/Users/hsen/Desktop/STT/data/Loi-normale-  
Test-1.csv"  
FILE_PATH_LOI2 = "/Users/hsen/Desktop/STT/data/Loi-normale-  
Test-2.csv"
```

Votre terminal indique que le script s'exécute depuis /Users/hsen/Desktop/STT/.

Assurez-vous que vos fichiers sont bien dans le sous-dossier **data de ce répertoire.**

Si le script est dans `/Users/hsen/Desktop/STT/tain.py` et les données dans `/Users/hsen/Desktop/STT/data/`, ces chemins sont corrects.

Résumé de l'Action

1. **Supprimez les lignes 4 à 17** (l'ancienne fonction `ouvrirUnFichier` et l'appel qui cause l'erreur).
2. Assurez-vous que les variables `FILE_PATH_...` (lignes 25-27) pointent exactement là où se trouvent vos fichiers dans le dossier `data`.

Le script utilisera alors la fonction `ouvrirUnFichier` (la plus robuste) et les chemins absolus corrects, ce qui devrait résoudre l'erreur `FileNotFoundError`.

j'ai modifié mais me renvoie ça : hsen@MacBook-Pro-de-HSEN ~ % cd /Users/hsen/Desktop/STT

hsen@MacBook-Pro-de-HSEN STT % python3 tain.py

Résultat sur le calcul d'un intervalle de fluctuation

Traceback (most recent call last):

File "/Users/hsen/Desktop/STT/tain.py", line 16, in <module>

```
donnees = pd.DataFrame(ouvrirUnFichier("./Users/hsen/Desktop/data/Echantillonnage-100-
Echantillons.csv"))
```

[illegible]

File "/Users/hsen/Desktop/STT/tain.py", line 9, in ouvrirUnFichier

with `open(nom, "r")` as `fichier`:

ΛΛΛΛΛΛΛΛΛΛΛΛΛΛΛΛ

```
FileNotFoundError: [Errno 2] No such file or directory: './Users/hsen/Desktop/data/
Echantillonnage-100-Echantillons.csv'
```

hsen@MacBook-Pro-de-HSEN STT %

Je vois que, malgré nos tentatives, le terminal **affiche toujours la même erreur** et la même ligne problématique dans la *traceback* :

Traceback (most recent call last):

```
File "/Users/hsen/Desktop/STT/tain.py", line 16, in  
<module>
```

```
    donnees = pd.DataFrame(ouvrirUnFichier("./Users/hsen/  
Desktop/data/Echantillonnage-100-Echantillons.csv"))
```

Cela signifie que **vous n'avez pas réussi à supprimer ou commenter l'ancien bloc de code** qui se trouve au début de votre script `tain.py`.

Même si vous avez ajouté le code correct plus bas, tant que l'ancien appel à la ligne 16 est là, le script s'arrête immédiatement à cette erreur `FileNotFoundError`.



Solution Urgente : Nettoyage du Code

Veuillez ouvrir votre fichier `/Users/hsen/Desktop/STT/tain.py` et **supprimer impérativement** les lignes suivantes qui correspondent au premier bloc de votre code :

Li gn	Contenu (à supprimer/commenter)
4	<code>def ouvrirUnFichier(nom):</code>
5	<code>with open(nom, "r") as fichier:</code>
6	<code>contenu = pd.read_csv(fichier)</code>
7	<code>return contenu</code>
9	<code>print("Résultat sur le calcul d'un intervalle de</code>
11	<code>donnees = pd.DataFrame(ouvrirUnFichier("./Users/hsen/ Desktop/data/Echantillonnage-100-Echantillons.csv"))</code>
13	<code>print("Résultat sur le calcul d'un intervalle de</code>
15	<code>print("Théorie de la décision")</code>

Le script doit commencer directement par les imports, puis par la déclaration de la fonction de chargement robuste et des chemins d'accès (qui commencent à la ligne 17 dans votre code actuel).

Votre script doit ressembler à ceci après nettoyage (en omettant les lignes commentées) :

Python

```
# coding:utf8

import pandas as pd
import math
import scipy
import scipy.stats
import numpy as np
from scipy import stats

# --- CHEMINS D'ACCÈS CORRECTS ---
FILE_PATH_ECHANTILLONS = "/Users/hsen/Desktop/STT/data/Echantillonnage-100-Echantillons.csv"
FILE_PATH_LOI1 = "/Users/hsen/Desktop/STT/data/Loi-normale-Test-1.csv"
FILE_PATH_LOI2 = "/Users/hsen/Desktop/STT/data/Loi-normale-Test-2.csv"

# --- Fonction de chargement (version robuste) ---
def ouvrirUnFichier(chemin_fichier):
    # ... le reste de votre fonction robuste ...
    # ...
    # ...

# Chargement principal
df_echantillons = ouvrirUnFichier(FILE_PATH_ECHANTILLONS)
# ... le reste de votre analyse ...

Veillez effectuer ce nettoyage de code, sauvegarder tain.py et relancer l'exécution. Le script devrait alors passer cette étape et commencer à exécuter les calculs statistiques.
```

voici mon code refait tout pour que ca marche (ne met pas de commentaire dans le code) :

```
#coding:utf8
```

```
import pandas as pd
```

```
import math

import scipy

import scipy.stats


# --- CHEMINS D'ACCÈS CORRECTS ---

FILE_PATH_ECHANTILLONS = "/Users/hsen/Desktop/STT/data/Echantillonnage-100-
Echantillons.csv"

FILE_PATH_LOI1 = "/Users/hsen/Desktop/STT/data/Loi-normale-Test-1.csv"

FILE_PATH_LOI2 = "/Users/hsen/Desktop/STT/data/Loi-normale-Test-2.csv"


def ouvrirUnFichier(nom):

    with open(nom, "r") as fichier:

        contenu = pd.read_csv(fichier)

    return contenu


#

print("Résultat sur le calcul d'un intervalle de fluctuation")


donnees = pd.DataFrame(ouvrirUnFichier("./Users/hsen/Desktop/data/Echantillonnage-100-
Echantillons.csv"))

#...

print("Résultat sur le calcul d'un intervalle de confiance")


#

#
```

```
print("Théorie de la décision")
```

```
import pandas as pd
```

```
import numpy as np
```

```
from scipy import stats
```

```
# Fonction de chargement ---
```

```
# Chemin d'accès
```

```
FILE_PATH_ECHANTILLONS = "/Users/hsen/Desktop/STT/data/Echantillonnage-100-  
Echantillons.csv"
```

```
FILE_PATH_LOI1 = "/Users/hsen/Desktop/STT/data/Loi-normale-Test-1.csv"
```

```
FILE_PATH_LOI2 = "/Users/hsen/Desktop/STT/data/Loi-normale-Test-2.csv"
```

```
def ouvrirUnFichier(chemin_fichier):
```

```
    """
```

```
    Ouvre un fichier CSV en forçant le séparateur ';' (caractéristique des fichiers Excel/CSV français).
```

```
    """
```

```
    try:
```

```
        # Tente la lecture avec le point-virgule (séparateur le plus probable après votre description)
```

```
        df = pd.read_csv(chemin_fichier, sep=';')
```

```
        # Si le fichier ne contient qu'une seule colonne, le point-virgule n'était pas le bon séparateur.
```

```
        # On essaie la virgule, ou un espace blanc comme dernière option.
```

```
        if df.shape[1] == 1:
```

```
            try:
```

```
                # Tente la virgule (séparateur standard anglo-saxon)
```

```
                df = pd.read_csv(chemin_fichier, sep=',')
```

```
except:
```

```
    # Tente l'espace blanc (sep=r'\s+')
```

```
    df = pd.read_csv(chemin_fichier, sep=r'\s+', engine='python', skipinitialspace=True)
```

```
# Nettoyage et attribution des noms de colonnes
```

```
df = df.dropna(axis=1, how='all')
```

```
if 'Echantillonnage' in chemin_fichier and df.shape[1] == 3:
```

```
    df.columns = ['Pour', 'Contre', 'Sans_opinion']
```

```
# Vérification finale
```

```
if df.shape[1] != 3 and 'Echantillonnage' in chemin_fichier:
```

```
    print(f"ATTENTION: Le fichier {chemin_fichier} n'a pu être lu qu'avec {df.shape[1]}  
colonnes. Le séparateur est peut-être incorrect.")
```

```
    return df
```

```
except FileNotFoundError:
```

```
    print(f"ERREUR: Le fichier {chemin_fichier} n'a pas été trouvé. Veuillez vérifier le chemin.")
```

```
    return None
```

```
except Exception as e:
```

```
    print(f"ERREUR lors de l'ouverture du fichier {chemin_fichier}: {e}")
```

```
    return None
```

```
# Chargement principal
```

```
df_echantillons = ouvrirUnFichier(FILE_PATH_ECHANTILLONS)
```

```
# --- 1. Théorie de l'échantillonnage ---
```

```
print("="*60)
```

```
print("ÉTAPE 1 : THÉORIE DE L'ÉCHANTILLONNAGE")
```

```
print("="*60)
```

```
if df_echantillons is not None:
```

```
    # Données de la population mère
```

```
    N_mere = 2185
```

```
    pour_mere = 852
```

```
    contre_mere = 911
```

```
    sans_opinion_mere = 422
```

```
    # 1.1 Calcul des moyennes des échantillons (arrondies à l'entier)
```

```
    moyennes_echantillons = df_echantillons.mean().round(0).astype(int)
```

```
    print("\n[1.1] Moyennes des 100 échantillons (arrondies) :")
```

```
    print(moyennes_echantillons)
```

```
    # 1.2 Fréquences des échantillons (fréquences moyennes)
```

```
    somme_moyennes = moyennes_echantillons.sum()
```

```
    frequences_echantillons = (moyennes_echantillons / somme_moyennes).round(2)
```

```
    print("\n[1.2] Fréquences des moyennes des échantillons (arrondies à 2 décimales) :")
```

```
    print(frequences_echantillons)
```

```
    # 1.3 Fréquences de la population mère
```

```
    effectifs_mere = pd.Series([pour_mere, contre_mere, sans_opinion_mere],
```

```
                               index=['Pour', 'Contre', 'Sans_opinion'])
```

```
    frequences_mere = (effectifs_mere / N_mere).round(2)
```

```
print("\n[1.3] Fréquences de la population mère (arrondies à 2 décimales) :")
print(frequencies_mere)
```

```
# 1.4 Intervalle de fluctuation (IF)
```

```
# Formule :  $[p - z_C * \sqrt{p(1-p)/n}, p + z_C * \sqrt{p(1-p)/n}]$ 
```

```
z_C = 1.96 # Seuil de 95%
```

```
n_moyen = somme_moyennes # Taille moyenne des échantillons
```

```
print(f"\n[1.4] Intervalle de fluctuation (IF) à 95% (n moyen={n_moyen}) :")
```

```
intervalle_fluctuation = {}
```

```
for opinion, p in frequencies_mere.items():
```

```
    # L'écart-type d'échantillonnage
```

```
    ecart_type_echantillonnage = np.sqrt(p * (1 - p) / n_moyen)
```

```
    margeErreur = z_C * ecart_type_echantillonnage
```

```
    borne_inf = p - margeErreur
```

```
    borne_sup = p + margeErreur
```

```
    intervalle_fluctuation[opinion] = (f"{borne_inf:.4f}", f"{borne_sup:.4f}")
```

```
    print(f"IF '{opinion}' (p={p:.2f}): [{borne_inf:.4f} ; {borne_sup:.4f}]")
```

```
# --- 2. Théorie de l'estimation ---
```

```
print("\n"+"="*60)
```

```
print("ÉTAPE 2 : THÉORIE DE L'ESTIMATION")
```

```
print("="+60)
```

if df_echantillons is not None:

2.1 Prendre le premier échantillon (méthode iloc(0))

premier_echantillon_pandas = df_echantillons.iloc[0]

Conversion en list() (comme demandé)

premier_echantillon_list = premier_echantillon_pandas.tolist()

print(f"\n[2.1] Premier échantillon (liste) : {premier_echantillon_list}")

2.2 et 2.3 Calcul de l'effectif total (n) et des fréquences (p_hat)

n_echantillon_1 = premier_echantillon_pandas.sum()

print(f"Taille de l'échantillon (n) : {n_echantillon_1}")

frequences_echantillon_1 = (premier_echantillon_pandas / n_echantillon_1)

print("\n[2.3] Fréquences du premier échantillon (p_hat) :")

print(frequences_echantillon_1.round(4))

2.4 Intervalle de confiance (IC)

Formule : $[p_hat - z_C * \sqrt{p_hat(1-p_hat)/n}, p_hat + z_C * \sqrt{p_hat(1-p_hat)/n}]$

$z_C = 1.96$

print(f"\n[2.4] Intervalle de confiance (IC) à 95% (n={n_echantillon_1}) :")

for opinion, p_hat in frequences_echantillon_1.items():

L'écart-type d'estimation

ecart_type_estimation = np.sqrt(p_hat * (1 - p_hat) / n_echantillon_1)

marge_erreur = z_C * ecart_type_estimation

```
borne_inf = p_hat - marge_erreur
```

```
borne_sup = p_hat + marge_erreur
```

```
print(f"IC '{opinion}' (p_hat={p_hat:.4f}): [{borne_inf:.4f} ; {borne_sup:.4f}]")
```

```
# --- 3 Test de Shapiro-Wilks) ---
```

```
print("\n"+"="*60)
```

```
print("ÉTAPE 3 : THÉORIE DE LA DÉCISION (Test de Shapiro-Wilks)")
```

```
print("="*60)
```

```
df_loi1 = ouvrirUnFichier(FILE_PATH_LOI1)
```

```
df_loi2 = ouvrirUnFichier(FILE_PATH_LOI2)
```

```
alpha = 0.05
```

```
def realiser_shapiro_wilks(df, nom_fichier, alpha):
```

```
    """Effectue le test de Shapiro-Wilks et affiche la conclusion."""
```

```
    if df is not None:
```

```
        # Assumer que la série de nombres est dans la première colonne
```

```
        data = df.iloc[:, 0].dropna()
```

```
        if len(data) < 3 or len(data) > 5000:
```

```
            # Shapiro-Wilks est moins fiable en dehors de cette plage
```

```
            print(f"ATTENTION: Le test de Shapiro-Wilks est moins adapté pour {nom_fichier}  
(taille={len(data)}).")
```

```
        stat_sw, p_value = stats.shapiro(data)
```



```
print(f"\nRésultats pour {nom_fichier} (Taille: {len(data)})")
```

```
print(f"Statistique W : {stat_sw:.4f}")
```

```
print(f"P-value : {p_value:.4f}")
```

```
if p_value > alpha:
```

```
    conclusion = f"P-value ({p_value:.4f}) > alpha ({alpha}). Ne rejette pas H0."
```

```
    resultat = "La distribution est considérée comme NORMALE."
```

```
else:
```

```
    conclusion = f"P-value ({p_value:.4f}) <= alpha ({alpha}). Rejette H0."
```

```
    resultat = "La distribution n'est PAS considérée comme normale."
```

```
print(f"Conclusion : {conclusion}")
```

```
print(f"Résultat : {resultat}")
```

```
return resultat
```

```
return "Non exécuté"
```

```
res_loi1 = realiser_shapiro_wilks(df_loi1, "Loi-normale-Test-1.csv", alpha)
```

```
res_loi2 = realiser_shapiro_wilks(df_loi2, "Loi-normale-Test-2.csv", alpha)
```

```
print("\n"+"="*60)
```

```
print("FIN DES CALCULS STATISTIQUES.")
```

```
print("="*60)
```

```
# --- BONUS ---
```

```
# --- 2.3 Analyse Approfondie pour le Bonus ---
```

```
print("\n"+"="*60)
```

```
print("ANALYSE BONUS : LOIS NON NORMALES")
```

```
print("="*60)
```

Les deux lois ont été jugées NON NORMALES. Nous allons analyser leurs statistiques.

```
def analyser_distribution(df, nom):  
    """Calcule statistiques descriptives clés pour aider à l'identification."""  
    if df is not None:  
        data = df.iloc[:, 0].dropna()  
  
        statistiques = {  
            "Taille (n)": len(data),  
            "Minimum": data.min(),  
            "Maximum": data.max(),  
            "Étendue (Max - Min)": data.max() - data.min(),  
            "Moyenne": data.mean(),  
            "Médiane": data.median(),  
            "Écart-type": data.std(),  
        }  
  
        print(f"\nStatistiques descriptives pour {nom} :")  
        for cle, valeur in statistiques.items():  
            print(f"- {cle:<20}: {valeur:.4f}")  
  
        print("\n=> Caractérisation (pour le Bonus) :")  
  
        # Test 1 : Forte asymétrie  
        if abs(statistiques["Moyenne"] - statistiques["Médiane"]) > 0.1 * statistiques["Écart-type"]:  
            print(" La moyenne et la médiane sont éloignées. La distribution est fortement  
asymétrique.")
```

```

print(" L'hypothèse d'une loi Exponentielle ou d'une autre loi asymétrique est forte.")

# Test 2 : Uniformité (Moyenne proche du centre de l'étendue)

elif abs(statistiques["Moyenne"] - (statistiques["Minimum"] + statistiques["Maximum"]) / 2) <
0.05 * statistiques["Écart-type"]:

    print(" La moyenne est très proche du centre de l'intervalle [Min, Max].")

    print(" L'hypothèse d'une Loi Uniforme est forte, car l'écart-type est faible par rapport à
l'étendue.")

else:

    print(" Les statistiques ne permettent pas de trancher facilement sans visualisation
(histogramme).")

analyser_distribution(ouvrirUnFichier(FILE_PATH_LOI1), "Loi-normale-Test-1.csv")
analyser_distribution(ouvrirUnFichier(FILE_PATH_LOI2), "Loi-normale-Test-2.csv")

```

Le problème est que votre code contient **deux fois la fonction ouvrirUnFichier** et surtout, l'appel au début du script qui échoue n'a pas été supprimé, même si vous l'avez refait dans la suite.

J'ai combiné tous les imports, corrigé l'appel de fichier qui créait l'erreur, et fusionné les deux définitions de la fonction **ouvrirUnFichier** pour ne garder que la version robuste avec la gestion des séparateurs.

Voici le code complet et fonctionnel pour la théorie de l'échantillonnage, de l'estimation et de la décision, sans commentaires dans le corps du code :

Python

```

#coding:utf8

import pandas as pd
import numpy as np
import scipy.stats

FILE_PATH_ECHANTILLONS = "/Users/hsen/Desktop/STT/data/Echantillonnage-100-Echantillons.csv"
FILE_PATH_LOI1 = "/Users/hsen/Desktop/STT/data/Loi-normale-Test-1.csv"
FILE_PATH_LOI2 = "/Users/hsen/Desktop/STT/data/Loi-normale-Test-2.csv"

def ouvrirUnFichier(chemin_fichier):
    try:
        df = pd.read_csv(chemin_fichier, sep=';')

        if df.shape[1] == 1:
            try:
                df = pd.read_csv(chemin_fichier, sep=',')
            except:
                df = pd.read_csv(chemin_fichier, sep=r'\s+',
engine='python', skipinitialspace=True)

        df = df.dropna(axis=1, how='all')

        if 'Echantillonnage' in chemin_fichier and
df.shape[1] == 3:
            df.columns = ['Pour', 'Contre', 'Sans_opinion']

            if df.shape[1] != 3 and 'Echantillonnage' in
chemin_fichier:
                print(f"ATTENTION: Le fichier {chemin_fichier}
n'a pu être lu qu'avec {df.shape[1]} colonnes. Le séparateur
est peut-être incorrect.")

            return df

    except FileNotFoundError:
        print(f"ERREUR: Le fichier {chemin_fichier} n'a pas
été trouvé. Veuillez vérifier le chemin.")
        return None
    except Exception as e:

```

```

        print(f"ERREUR lors de l'ouverture du fichier
{chemin_fichier}: {e}")
        return None

df_echantillons = ouvrirUnFichier(FILE_PATH_ECHANTILLONS)

print("="*60)
print("ÉTAPE 1 : THÉORIE DE L'ÉCHANTILLONNAGE")
print("="*60)

if df_echantillons is not None:
    N_mere = 2185
    pour_mere = 852
    contre_mere = 911
    sans_opinion_mere = 422

    moyennes_echantillons =
df_echantillons.mean().round(0).astype(int)
    print("\n[1.1] Moyennes des 100 échantillons
(arrondies) :")
    print(moyennes_echantillons)

    somme_moyennes = moyennes_echantillons.sum()
    frequences_echantillons = (moyennes_echantillons /
somme_moyennes).round(2)
    print("\n[1.2] Fréquences des moyennes des échantillons
(arrondies à 2 décimales) :")
    print(frequences_echantillons)

    effectifs_mere = pd.Series([pour_mere, contre_mere,
sans_opinion_mere],
                                index=['Pour', 'Contre',
'Sans_opinion'])
    frequences_mere = (effectifs_mere / N_mere).round(2)
    print("\n[1.3] Fréquences de la population mère
(arrondies à 2 décimales) :")
    print(frequences_mere)

    z_C = 1.96
    n_moyen = somme_moyennes

    print(f"\n[1.4] Intervalle de fluctuation (IF) à 95% (n
moyen={n_moyen}) :")

    intervalle_fluctuation = {}

```

```

        for opinion, p in frequences_mere.items():
            ecart_type_echantillonnage = np.sqrt(p * (1 - p) /
n_moyen)
            marge_erreur = z_C * ecart_type_echantillonnage

            borne_inf = p - marge_erreur
            borne_sup = p + marge_erreur

            intervalle_fluctuation[opinion] =
(f"{borne_inf:.4f}", f"{borne_sup:.4f}")
            print(f"IF '{opinion}' (p={p:.2f}):
[ {borne_inf:.4f} ; {borne_sup:.4f} ]")

print("\n"+"="*60)
print("ÉTAPE 2 : THÉORIE DE L'ESTIMATION")
print("="+60)

if df_echantillons is not None:
    premier_echantillon_pandas = df_echantillons.iloc[0]

    premier_echantillon_list =
premier_echantillon_pandas.tolist()
    print(f"\n[2.1] Premier échantillon (liste) :
{premier_echantillon_list}")

    n_echantillon_1 = premier_echantillon_pandas.sum()
    print(f"Taille de l'échantillon (n) : {n_echantillon_1}")

    frequences_echantillon_1 = (premier_echantillon_pandas /
n_echantillon_1)
    print("\n[2.3] Fréquences du premier échantillon
(p_hat) :")
    print(frequences_echantillon_1.round(4))

    z_C = 1.96

    print(f"\n[2.4] Intervalle de confiance (IC) à 95%
(n={n_echantillon_1}) :")

    for opinion, p_hat in frequences_echantillon_1.items():
        ecart_type_estimation = np.sqrt(p_hat * (1 - p_hat) /
n_echantillon_1)
        marge_erreur = z_C * ecart_type_estimation

```

```

borne_inf = p_hat - marge_erreur
borne_sup = p_hat + marge_erreur

print(f"IC '{opinion}' (p_hat={p_hat:.4f}):
[{borne_inf:.4f} ; {borne_sup:.4f}]")

print("\n"+"="*60)
print("ÉTAPE 3 : THÉORIE DE LA DÉCISION (Test de Shapiro-
Wilks)")
print("="*60)

df_loi1 = ouvrirUnFichier(FILE_PATH_LOI1)
df_loi2 = ouvrirUnFichier(FILE_PATH_LOI2)
alpha = 0.05

def realiser_shapiro_wilks(df, nom_fichier, alpha):
    if df is not None:
        data = df.iloc[:, 0].dropna()
        if len(data) < 3 or len(data) > 5000:
            print(f"ATTENTION: Le test de Shapiro-Wilks est
moins adapté pour {nom_fichier} (taille={len(data)}).")

            stat_sw, p_value = scipy.stats.shapiro(data)

            print(f"\nRésultats pour {nom_fichier} (Taille:
{len(data)}")
            print(f"Statistique W : {stat_sw:.4f}")
            print(f"P-value : {p_value:.4f}")

            if p_value > alpha:
                conclusion = f"P-value ({p_value:.4f}) > alpha
({alpha}). Ne rejette pas H0."
                resultat = "La distribution est considérée comme
NORMALE."
            else:
                conclusion = f"P-value ({p_value:.4f}) <= alpha
({alpha}). Rejette H0."
                resultat = "La distribution n'est PAS considérée
comme normale."

            print(f"Conclusion : {conclusion}")
            print(f"Résultat : {resultat}")
            return resultat
    return "Non exécuté"

```

```

res_loi1 = realiser_shapiro_wilks(df_loi1, "Loi-normale-
Test-1.csv", alpha)
res_loi2 = realiser_shapiro_wilks(df_loi2, "Loi-normale-
Test-2.csv", alpha)

print("\n"+"="*60)
print("FIN DES CALCULS STATISTIQUES.")
print("="*60)

print("\n"+"="*60)
print("ANALYSE BONUS : LOIS NON NORMALES")
print("="*60)

def analyser_distribution(df, nom):
    if df is not None:
        data = df.iloc[:, 0].dropna()

        statistiques = {
            "Taille (n)": len(data),
            "Minimum": data.min(),
            "Maximum": data.max(),
            "Étendue (Max - Min)": data.max() - data.min(),
            "Moyenne": data.mean(),
            "Médiane": data.median(),
            "Écart-type": data.std(),
        }

        print(f"\nStatistiques descriptives pour {nom} :")
        for cle, valeur in statistiques.items():
            print(f"- {cle:<20}: {valeur:.4f}")

        print("\n=> Caractérisation (pour le Bonus) :")

        if abs(statistiques["Moyenne"] -
statistiques["Médiane"]) > 0.1 * statistiques["Écart-type"]:
            print("    La moyenne et la médiane sont
éloignées. La distribution est fortement asymétrique.")
            print("    L'hypothèse d'une loi Exponentielle ou
d'une autre loi asymétrique est forte.")

        elif abs(statistiques["Moyenne"] -
(statistiques["Minimum"] + statistiques["Maximum"]) / 2) <
0.05 * statistiques["Écart-type"]:

```



```

        print("    La moyenne est très proche du centre
de l'intervalle [Min, Max].")
        print("    L'hypothèse d'une Loi Uniforme est
forte, car l'écart-type est faible par rapport à l'étendue.")

    else:
        print("    Les statistiques ne permettent pas de
trancher facilement sans visualisation (histogramme).")

analyser_distribution(ouvrirUnFichier(FILE_PATH_LOI1), "Loi-
normale-Test-1.csv")
analyser_distribution(ouvrirUnFichier(FILE_PATH_LOI2), "Loi-
normale-Test-2.csv")

```

voici ce que renvoie le terminal : commente et met en tableau ces resultats : ÉTAPE 1 : THÉORIE DE L'ÉCHANTILLONNAGE

=====

[1.1] Moyennes des 100 échantillons (arrondies) :

```

Pour      391
Contre    416
Sans_opinion  193
dtype: int64

```

[1.2] Fréquences des moyennes des échantillons (arrondies à 2 décimales) :

```

Pour      0.39
Contre    0.42
Sans_opinion  0.19
dtype: float64

```

[1.3] Fréquences de la population mère (arrondies à 2 décimales) :

Pour 0.39

Contre 0.42

Sans_opinion 0.19

dtype: float64

[1.4] Intervalle de fluctuation (IF) à 95% (n moyen=1000) :

IF 'Pour' (p=0.39): [0.3598 ; 0.4202]

IF 'Contre' (p=0.42): [0.3894 ; 0.4506]

IF 'Sans_opinion' (p=0.19): [0.1657 ; 0.2143]

=====

ÉTAPE 2 : THÉORIE DE L'ESTIMATION

=====

[2.1] Premier échantillon (liste) : [395, 396, 209]

Taille de l'échantillon (n) : 1000

[2.3] Fréquences du premier échantillon (p_hat) :

Pour 0.395

Contre 0.396

Sans_opinion 0.209

Name: 0, dtype: float64

[2.4] Intervalle de confiance (IC) à 95% (n=1000) :

IC 'Pour' (p_hat=0.3950): [0.3647 ; 0.4253]

IC 'Contre' (p_hat=0.3960): [0.3657 ; 0.4263]

IC 'Sans_opinion' ($\hat{p}=0.2090$): [0.1838 ; 0.2342]

=====

ÉTAPE 3 : THÉORIE DE LA DÉCISION (Test de Shapiro-Wilks)

=====

Résultats pour Loi-normale-Test-1.csv (Taille: 2000)

Statistique W : 0.9639

P-value : 0.0000

Conclusion : P-value (0.0000) \leq alpha (0.05). Rejette H_0 .

Résultat : La distribution n'est PAS considérée comme normale.

Résultats pour Loi-normale-Test-2.csv (Taille: 2000)

Statistique W : 0.2609

P-value : 0.0000

Conclusion : P-value (0.0000) \leq alpha (0.05). Rejette H_0 .

Résultat : La distribution n'est PAS considérée comme normale.

=====

FIN DES CALCULS STATISTIQUES.

=====

=====

ANALYSE BONUS : LOIS NON NORMALES

=====

Statistiques descriptives pour Loi-normale-Test-1.csv :

- Taille (n) : 2000.0000

- Minimum : -2.0000
- Maximum : 10.0000
- Étendue (Max - Min) : 12.0000
- Moyenne : 3.0430
- Médiane : 3.0000
- Écart-type : 1.5357

=> Caractérisation (pour le Bonus) :

Les statistiques ne permettent pas de trancher facilement sans visualisation (histogramme).

Statistiques descriptives pour Loi-normale-Test-2.csv :

- Taille (n) : 2000.0000
- Minimum : 1.0000
- Maximum : 14.0000
- Étendue (Max - Min) : 13.0000
- Moyenne : 1.1875
- Médiane : 1.0000
- Écart-type : 0.7626

=> Caractérisation (pour le Bonus) :

La moyenne et la médiane sont éloignées. La distribution est fortement asymétrique.

L'hypothèse d'une loi Exponentielle ou d'une autre loi asymétrique est forte.

hsen@MacBook-Pro-de-HSEN STT %

Voici l'analyse commentée et les résultats présentés sous forme de tableaux pour les étapes 1, 2 et 3 de votre exercice de statistiques.



Rapport d'Analyse Statistique : Échantillonnage, Estimation et Décision

Étape 1 : Théorie de l'Échantillonnage

Cette étape compare les proportions observées sur l'ensemble des 100 échantillons (moyennes) avec les proportions réelles de la population mère, et calcule l'intervalle dans lequel ces proportions doivent se situer.

Tableau 1 : Comparaison des Fréquences et Intervalle de Fluctuation (IF)

Opinion	Effectif Moyen des Échantillons	Fréquence Moyenne des Échantillons (\bar{p})	Fréquence de la Population Mère	Intervalle de Fluctuation (IF à 95%,
Pour	391	0.39	0.39	[0.3598 ; 0.4202]
Contr	416	0.42	0.42	[0.3894 ; 0.4506]
Sans opinion	193	0.19	0.19	[0.1657 ; 0.2143]
Total	1000	1.00	1.00	

Commentaire de l'Étape 1 :

L'analyse démontre une forte cohérence entre les résultats de la population mère (p) et les moyennes des échantillons (\bar{p}). La fréquence moyenne observée sur les échantillons est rigoureusement identique (0.39, 0.42, 0.19) à la fréquence théorique de la population mère.

L'**Intervalle de Fluctuation (IF)** à 95 % est calculé avec la taille moyenne d'échantillon ($n=1000$). Cet intervalle permet de prédire les résultats d'un *futur* échantillon. Par exemple, pour l'opinion "Pour" (39%), on s'attend à ce que 95 % des échantillons de taille 1000 présentent une fréquence comprise entre 35.98 % et 42.02 %. Cette faible marge d'erreur confirme que l'échantillonnage est bien mené et que la taille $n=1000$ est suffisante pour obtenir une estimation précise.

Étape 2 : Théorie de l'Estimation

Cette étape calcule l'**Intervalle de Confiance (IC)** autour des résultats du *premier* échantillon pris individuellement. L'IC permet d'estimer la proportion *inconnue* de la population mère à partir de cet échantillon.

Tableau 2 : Fréquences du Premier Échantillon et Intervalle de Confiance (IC)

Opinion	Effectif du 1er Échantillon (n=1000)	Fréquence du 1er Échantillon (p^{\wedge})	Intervalle de Confiance (IC à 95%, n=1000)
Pour	395	0.3950	[0.3647 ; 0.4253]
Contre	396	0.3960	[0.3657 ; 0.4263]
Sans opinion	209	0.2090	[0.1838 ; 0.2342]

Commentaire de l'Étape 2 :

Le premier échantillon isolé présente des fréquences légèrement différentes des fréquences théoriques de la population mère (par exemple, "Contre" passe de 0.42 à 0.396).

L'**Intervalle de Confiance (IC)** à 95 % indique la plage de valeurs dans laquelle la *vraie* proportion de la population mère se situe, avec 95 % de certitude. Pour toutes les opinions, la fréquence réelle de la population mère (0.39, 0.42, 0.19 - voir Tableau 1) est bien incluse dans l'intervalle calculé à partir de cet échantillon, validant ainsi la pertinence de l'échantillon.

Étape 3 : Théorie de la Décision (Test de Shapiro-Wilks)

Le Test de Shapiro-Wilks est utilisé pour déterminer si une distribution suit la Loi Normale (Hypothèse Nulle H_0). Le seuil de signification (α) est fixé à 0.05.

Fichier de Test	Taille (n)	Statistique W	P-value	Décision (P-value vs α)	Conclusion
Loi-normale-Test-1.csv	2000	0.9639	0.0000	$P \leq \alpha$ (Rejet de H_0)	NON NORMALE
Loi-normale-Test-2.csv	2000	0.2609	0.0000	$P \leq \alpha$ (Rejet de H_0)	NON NORMALE

Commentaire de l'Étape 3 :

Pour les deux jeux de données testés, la P-value obtenue (0.0000) est très nettement inférieure au seuil de signification $\alpha=0.05$. Par conséquent, l'hypothèse nulle (H_0 : la distribution est Normale) est **rejetée** pour les deux cas. Aucune des deux distributions n'est considérée comme suivant une loi normale.

Analyse Bonus : Caractérisation des Lois Non Normales

Étant donné que les deux lois ne sont pas normales, une analyse des statistiques descriptives est effectuée pour tenter de les caractériser.

Tableau 4 : Statistiques Descriptives Détaillées

Statistique	Loi-normale-Test-1.csv	Loi-normale-Test-2.csv
-------------	------------------------	------------------------

Taille (n)	2000	2000
Minimum	-2.0000	1.0000
Maximum	10.0000	14.0000
Moyenne	3.0430	1.1875
Médiane	3.0000	1.0000
Écart-type	1.5357	0.7626

Commentaire de l'Analyse Bonus :

1. Loi-normale-Test-1.csv :

- La **Moyenne (3.0430)** et la **Médiane (3.0000)** sont extrêmement proches.
- **Caractérisation :** Cette proximité suggère une distribution relativement **symétrique** malgré le rejet du test de Shapiro-Wilks. Pour trancher entre une loi discrète, une loi Uniforme ou une distribution en cloche non parfaite, une visualisation par histogramme serait indispensable. Le code n'a pas pu trancher facilement.

2. Loi-normale-Test-2.csv :

- La **Moyenne (1.1875)** est nettement supérieure à la **Médiane (1.0000)**.
- **Caractérisation :** Cet écart entre les deux indicateurs (Moyenne > Médiane) est la signature d'une distribution **fortement asymétrique positive (étalée à droite)**. Cela rend l'hypothèse d'une **Loi Exponentielle** (souvent utilisée pour modéliser des durées ou des temps d'attente) ou d'une autre loi asymétrique (comme la Loi Gamma) très forte.

Voulez-vous que je vous aide à interpréter le sens précis des intervalles de fluctuation et de confiance par rapport au contexte de l'enquête (par exemple, ce que signifie concrètement l'IC pour l'opinion "Pour") ?

reformule en rédigeant complètement : SEANCE 5

Questions de cours

→ 1 – Comment définir l'échantillonnage ? Pourquoi ne pas utiliser la population en entier ?

Quelles sont les méthodes d'échantillonnage ? Comment les choisir ?

L'échantillonnage consiste à prélever dans une population mère une partie de celle-ci au hasard avec une taille n fixée. Chaque échantillon fournit alors un résultat. L'échantillon est un groupe restreint, c'est-à-dire un sous-ensemble, issu d'une variable aléatoire X de la population. Il existe plusieurs façons

de tirer un échantillon de la population mère. La plupart nécessitent de disposer une base de sondage.

On peut aussi utiliser un échantillon aléatoire lorsqu'il est impossible de constituer une base de sondage. L'échantillon aléatoire offre des résultats recueillis sur ce sous-ensemble qui doivent pouvoir

être étendus, c'est-à-dire inférés, à la population mère. Parmi les échantillons aléatoires on distingue l'échantillon non biaisé : tiré au hasard dans lequel tous les individus ont la même chance de se retrouver dans l'échantillon ; et l'échantillon biaisé : les éléments n'ont pas été pris au hasard.

Il existe plusieurs méthodes d'échantillonnage :

- les méthodes aléatoires : tirage avec ou sans remise
- les méthodes non aléatoires : échantillonnage systématique, méthode des quotas
- les méthodes d'échantillonnage « Monte Carlo »

L'utilisation de ces méthodes dépend de l'échantillon, du sujet étudié et de ce que l'on souhaite démontrer.

→ 2 – Comment définir un estimateur et une estimation ?

L'estimation permet d'estimer les paramètres d'une loi de probabilité. En ce qui concerne l'estimateur, il

s'agit de la variable aléatoire. Un estimateur est une fonction des données. Il est construit de telle façon

que sa valeur soit proche de la vraie valeur du paramètre. Le but de la théorie de l'estimation est de choisir, parmi toutes les statistiques possibles, le meilleur estimateur, c'est-à-dire celui qui donnera une

estimation ponctuelle la plus proche possible du paramètre, et ceci quel que soit l'échantillon.

→ 3 – Comment distingueriez-vous l'intervalle de fluctuation et l'intervalle de confiance ?

L'intervalle de fluctuation suppose que la vraie proportion théorique soit connue. C'est un échantillonnage et non une estimation. L'intervalle de confiance s'en distingue car c'est un outil

statistique utiliser pour estimer la plage dans laquelle se situe un paramètre de population à partir d'un

échantillon. Il permet de quantifier l'incertitude d'une estimation, comme la moyenne ou la variance, en

fournissant une fourchette d'estimation.

→ 4 – Qu'est-ce qu'un biais dans la théorie de l'estimation ?

Dans la théorie de l'estimation, un biais correspond à la différence entre l'espérance de l'estimateur et la valeur à estimer dans la population, on l'appelle également erreur d'estimation.

→ 5 – Comment appelle-t-on une statistique travaillant sur la population totale ? Faites le lien avec la notion de données massives ?

Une statistique travaillant sur la population totale est une statistique exhaustive. Le lien entre les données massives et la statistique exhaustive correspond au fait que les deux notions rendent plus accessible la démarche de la statistique exhaustive puisqu'on dispose de très grands volumes de données qui peuvent concerner tous les individus d'un système.

→ 6 – Quels sont les enjeux autour du choix d'un estimateur ?

Il y a plusieurs enjeux autour du choix d'un estimateur. Premièrement, il y a un enjeu autour de sa variance qui influera sur la précision de l'estimateur. Par ailleurs, la statistique étant un résumé apporté

par un échantillon, il est par conséquent très important de ne pas perdre l'information. Ainsi, en tenant

compte de ces deux points, on peut aborder la recherche du meilleur estimateur suivant deux méthodes :

- soit en recherchant des statistiques exhaustives qui conduisent à des estimateurs sans biais de variance minimale

- soit en étudiant la quantité d'information de Fisher qui apporte des indications sur la précision d'un estimateur

→ 7 – Quelles sont les méthodes d'estimation d'un paramètre ? Comment en sélectionner une ?

Les principales méthodes d'estimation sont :

- la méthode des moments : il s'agit d'égaliser les moments théoriques avec les moments observés dans les données

- La méthode du maximum de vraisemblance : il s'agit de choisir le paramètre qui rend les données observées le plus vraisemblables
- la méthode des moindres carrés : il s'agit de minimiser la somme des carrés des écarts entre les valeurs observées et les valeurs estimées par le modèle
- la méthode bayésienne : combiner une information a priori sur le paramètre avec les données observées pour obtenir une distribution a posteriori

Pour choisir la méthode d'estimation, il faut prendre en compte la taille de l'échantillon, la complexité

du modèle, la présence de biais ou d'erreurs, l'objectif de l'étude.

→ 8 – Quels sont les tests statistiques existants ? A quoi servent-ils ? Comment créer un test ?

Un test statistique est une méthode de calcul permettant de décider si une série statistique

d'observations est compatible avec une loi de probabilité entièrement spécifique : ou comment savoir si

un résultat observé est en accord avec une distribution théorique. Les tests statistiques suivent une loi connue, on retrouve parmi ces test :

- test de conformité
- test d'homogénéité
- test d'adéquation à une loi de probabilité
- test d'indépendance de deux caractères
- test de signification
- test d'hypothèse
- test paramétrique
- tests non paramétriques : test de Whitney, test de Wilcoxon, test du coefficient de corrélation de Spearman, test de Fisher)
- tests robustes : tests libres

Pour créer un test, il faut :

- formuler des hypothèses
- choisir la statistique du test
- choisir le niveau de risque

- calculer la valeur observée de la statistique
- déterminer la région critique
- interpréter le résultat

→ 9 – Que pensez-vous des critiques de la statistique inférentielle ?

La statistique inférentielle permet de tirer des conclusions sur une population à partir d'un échantillon,

elle repose sur des modèles probabilistes, des tests et intervalles de confiance pour décider avec un certain risque d'erreur. C'est un outil très puissant mais aussi très critiqué car parfois mal compris ou mal utilisé. Les critiques portent sur la dépendance excessive aux hypothèses, la confusion fréquente

entre signification et importance, une mauvaise interprétation du p-value, un détachement du contexte, le

caractère binaire de la décision ou encore sur le problème de la « chasse à la significativité ».

Selon moi, ces critiques invitent à se poser des questions effectivement sur l'usage de la statistique inférentielle, plutôt que de l'abandonner, il faut davantage chercher à l'utiliser de manière raisonnée par

exemple en vérifiant les conditions, en utilisant des tests non paramétriques ou des méthodes robustes,

en complétant avec des mesures d'effet, en formant à la logique des tests, en utilisant des approches bayésiennes ou des analyses de sensibilités ou encore en combinant les données quantitatives à l'analyse qualitative et théorique. Finalement, la statistique inférentielle est un outil et non une vérité. Il

est nécessaire de garder un esprit critique tout le temps.

Voici la reformulation complète et rédigée des questions de cours sur la Séance 5, structurée pour la clarté et l'exhaustivité.



Questions de Cours : Fondamentaux de la Statistique Inférentielle

1. L'Échantillonnage

Définition et Objectif

L'**échantillonnage** est le processus par lequel on prélève, au hasard et selon une taille n fixée, un sous-ensemble (l'échantillon) d'une population plus vaste (la population mère) afin d'en étudier les caractéristiques. Un échantillon est un sous-ensemble des observations d'une variable aléatoire X de la population. L'objectif est de pouvoir **inférer** les résultats observés sur ce groupe restreint à l'ensemble de la population mère.

Pourquoi Échantillonner ?

L'utilisation de l'échantillonnage s'impose lorsque l'étude exhaustive de la **population en entier** est :

- **Impossible** : La population est trop grande (infinie) ou inaccessible.
- **Trop coûteuse** : En temps, en ressources humaines ou en budget.
- **Destructive** : Le processus de mesure détruit l'élément mesuré.

Les Méthodes d'Échantillonnage

L'échantillonnage nécessite souvent une **base de sondage** (liste des unités de la population) ou peut être basé sur des méthodes aléatoires pures lorsque cette base est indisponible. On distingue :

Catégorie	Méthodes et Principes
Aléatoires (Probabilis)	Tirage avec ou sans remise. L'échantillon est non biaisé (chaque individu a une chance égale d'être sélectionné).
Non Aléatoires	Méthode des quotas (respect des proportions connues de la population), Échantillonnage systématique (tirage régulier). Ces méthodes peuvent être biaisées
Simulation	Méthodes de Monte Carlo (utilisation de tirages aléatoires répétés pour obtenir des

Choix de la Méthode

Le choix de la méthode dépend intrinsèquement de l'**échantillon disponible**, du **sujet d'étude**, des **contraintes logistiques** et de la **précision souhaitée** pour l'inférence.

2. Estimateur et Estimation

- **L'Estimation** : C'est la démarche statistique qui vise à calculer, à partir d'un échantillon, les **paramètres inconnus** d'une loi de probabilité caractérisant la population (par exemple, la moyenne μ ou la variance σ^2).
- **L'Estimateur** : C'est la **variable aléatoire** utilisée pour réaliser l'estimation. C'est une fonction des données de l'échantillon. Un estimateur est construit de manière à ce que sa valeur soit, en théorie, la plus proche possible de la vraie valeur du paramètre de la population. L'objectif de la théorie de l'estimation est de sélectionner le **meilleur estimateur**, celui qui produira l'estimation ponctuelle la plus précise pour tout échantillon possible.

3. Intervalle de Fluctuation vs. Intervalle de Confiance

La distinction repose sur la connaissance ou non du paramètre théorique de la population :

Caractéristique	Intervalle de Fluctuation (IF)	Intervalle de Confiance (IC)
Connaissance du Paramètre	p est connu (proportion théorique de la population).	p est inconnu (à estimer).
But	Prédiction / Échantillonnage. Décrire la plage de valeurs que doit prendre un <i>futur</i>	Estimation. Estimer la plage de valeurs dans laquelle se situe le <i>vrai paramètre</i>
Question posée	Si p est la proportion, l'échantillon est-il cohérent ?	Quelle est la proportion p de la population, compte tenu de mon

4. Le Biais dans la Théorie de l'Estimation

Dans la théorie de l'estimation, le **biais** représente l'erreur systématique et constante de l'estimateur. Il est défini comme la **différence entre l'espérance mathématique de l'estimateur et la vraie valeur du paramètre** à estimer dans la population. Un estimateur est dit **sans biais** si son espérance est égale au paramètre lui-même.

5. Statistique sur la Population Totale et Lien avec les Données Massives

- **Statistique sur la Population Totale** : On parle de **statistique exhaustive** (ou recensement).
- **Lien avec les Données Massives (Big Data)** : L'émergence des données massives (grands volumes de données) rend la démarche de la statistique exhaustive plus accessible. En collectant des volumes de données très importants (souvent proches de la totalité ou concernant l'ensemble des interactions d'un système), il devient possible de réaliser des analyses qui s'approchent du traitement exhaustif plutôt que de s'appuyer uniquement sur l'échantillonnage.

6. Enjeux autour du Choix d'un Estimateur

Le choix d'un estimateur est crucial car il détermine la fiabilité des conclusions. Les principaux enjeux sont :

1. **La Précision (Variance) :** La variance de l'estimateur influence sa précision. Un bon estimateur doit avoir une **variance minimale**.
2. **L'Exhaustivité de l'Information :** La statistique étant un résumé de l'échantillon, il est essentiel de choisir un estimateur qui ne fait pas perdre l'information initiale.

Pour trouver le **meilleur estimateur** (sans biais et de variance minimale), les méthodes s'appuient sur :

- La recherche de **statistiques exhaustives** (qui résument toute l'information).
- L'étude de la **quantité d'information de Fisher**, qui fournit une borne inférieure à la variance (Borne de Cramer-Rao) et donne des indications sur la précision théorique.

7. Les Méthodes d'Estimation d'un Paramètre

Méthode d'Estimation	Principe
Méthode des Moments	Égaliser les moments théoriques de la loi de probabilité avec les moments observés dans les données de l'échantillon.
Méthode du Maximum de Vraisemblance (MV)	Choisir la valeur du paramètre qui rend la probabilité d'observer les données actuelles (la vraisemblance) la plus élevée.
Méthode des Moindres Carrés (MC)	Minimiser la somme des carrés des écarts entre les valeurs observées et les valeurs prédites par le modèle.
Méthode Bayésienne	Combiner l'information <i>a priori</i> sur le paramètre avec les données observées pour obtenir une distribution <i>a posteriori</i> du paramètre.

Sélectionner une Méthode

Le choix dépend de plusieurs facteurs : la **taille de l'échantillon**, la **complexité du modèle**, l'existence de **biais ou d'erreurs**, et surtout l'**objectif de l'étude** (par exemple, la méthode du Maximum de Vraisemblance est très utilisée pour ses bonnes propriétés asymptotiques).

8. Les Tests Statistiques

Définition et Utilité

Un **test statistique** est une méthode de calcul formelle qui permet de prendre une décision : déterminer si une série d'observations empiriques est **compatible** avec une hypothèse théorique (une loi de probabilité spécifique). En d'autres termes, il permet de savoir si un résultat observé est en accord avec une distribution théorique connue.

Types de Tests

Les tests peuvent être classés selon leur fonction ou leurs hypothèses :

- **Tests Paramétriques :** (ex: t de Student, ANOVA) supposent que les données suivent une loi de probabilité spécifique (souvent la Loi Normale).
- **Tests Non Paramétriques :** (ex: Whitney, Wilcoxon, Spearman, Khi-deux) ne nécessitent pas d'hypothèse sur la distribution sous-jacente des données.
- **Tests selon l'objectif :** Test de conformité, d'homogénéité, d'adéquation à une loi (Goodness-of-fit), d'indépendance, d'hypothèse, etc.

Étapes de Création d'un Test

La construction d'un test statistique suit une démarche rigoureuse :

1. **Formuler les Hypothèses :** Hypothèse Nulle (H_0) et Hypothèse Alternative (H_1).
2. **Choisir la Statistique de Test :** Sélectionner la statistique dont la distribution est connue sous H_0 .
3. **Choisir le Niveau de Risque (α) :** Déterminer la probabilité maximale d'erreur de Type I (rejeter H_0 alors qu'elle est vraie).
4. **Calculer la Valeur Observée :** Calculer la valeur de la statistique à partir des données de l'échantillon.
5. **Déterminer la Région Critique :** Définir la zone de valeurs de la statistique qui conduisent au rejet de H_0 .
6. **Interpréter le Résultat :** Comparer la valeur observée à la région critique (ou comparer la p-value à α) pour décider de rejeter ou non H_0 .

9. Les Critiques de la Statistique Inférentielle

La statistique inférentielle est un outil puissant pour tirer des conclusions sur une population à partir d'un échantillon, mais elle est sujette à de nombreuses critiques. Ces critiques portent notamment sur :

- La **dépendance excessive** aux hypothèses strictes des modèles (ex: normalité, indépendance).
- La **confusion fréquente** entre *signification statistique* et *importance pratique*.
- La **mauvaise interprétation de la p-value** (souvent perçue comme la probabilité que H_0 soit vraie, ce qui est faux).
- Le **caractère binaire** de la décision (rejet/non-rejet) qui simplifie à l'extrême la réalité.
- Le problème de la "**chasse à la significativité**" (p-hacking).
- Le **détachement du contexte** théorique ou qualitatif.

Positionnement et Amélioration

Ces critiques ne doivent pas conduire à l'abandon de la statistique inférentielle, mais à un **usage plus raisonné et critique**. Il est essentiel de :

- **Vérifier systématiquement** les conditions d'application (utiliser des tests non paramétriques ou des méthodes robustes si nécessaire).
- **Compléter l'analyse** avec des mesures d'effet (taille de l'effet) pour évaluer l'importance pratique.
- **Former correctement** à la logique des tests.
- Envisager des **approches alternatives** comme les méthodes bayésiennes ou l'analyse de sensibilité.
- **Combiner** les données quantitatives avec une analyse qualitative et théorique approfondie.

La statistique inférentielle est un outil d'aide à la décision, mais n'est pas une vérité absolue. L'esprit critique reste indispensable à chaque étape de l'analyse.

reformule sans trop modifier .: SEANCE 6 – ANALYSE DE DONNEES

Questions/réponses – cours

1- Qu'est-ce qu'une statistique ordinale ? A quelle autre statistique catégorielle s'oppose-t-elle ?

Quel type de variables utilise-t-elle ? En quoi cela peut matérialiser une hiérarchie spatiale ?

→ Une statistique d'ordre ou statistique ordinale est le cœur de la géographie humaine. De manière annuelle, mensuelle, voire hebdomadaire, un certain nombre de classements est opéré en utilisant des

objets géographiques. Leur objectif commun est de montrer quelle entité a descendu, stagné ou monté dans le classement.

La statistique ordinale s'oppose principalement à la statistique nominale.

La statistique ordinale utilise des variables qualitatives ordinales.

Cela peut matérialiser une hiérarchie spatiale car la statistique ordinale permet d'établir un ordre entre

les espaces : ordre entre le centre et la périphérie, un ordre du centre urbain dense jusqu'à la ruralité...

Une variable ordinale peut de fait matérialiser une hiérarchie spatiale dès lors que ses catégories représentent un niveau, un rang ou une intensité appliquée à des lieux, des territoires ou des zones. Cartographiquement, l'ordinal produit des cartes en classes ordonnées. Cette statistique permet aussi de traduire des relations de domination ou de centralité. Enfin, elle établit des niveaux ou des rangs territoriaux. En géographie physique, les lois d'ordre servent notamment à étudier la hauteur maximale des crues d'un cours d'eau, l'intensité du plus fort tremblement de terre dans une zone sismique donnée. Pour ce qui est de la géographie humaine, leur utilisation découle du fait de l'apparition plus ou moins spontanée de hiérarchies au sein des sociétés et des espaces étudiés.

2- Quel ordre est à privilégier dans les classifications ?

→ L'ordre à privilégier est l'ordre croissant ou ordre naturel. Il existe des exceptions en géographies telles que la loi dite rang-taille. L'ordination permet ainsi de rechercher les valeurs aberrantes, trop grandes ou trop petites, d'une série d'observations.

3- Quelle est la différence entre une corrélation des rangs et une concordance de classements ?

→ La corrélation des rangs mesure la force et le sens de la relation entre deux séries de rangs (comparaison de deux variables ordonnées, on cherche à savoir si les classements sont proportionnellement liés) tandis que la concordance des classements mesure à quel point plusieurs classements sont identiques (possibilité de comparer plus de deux classements, recherche de l'accord général entre plusieurs classements).

4- Quelle est la différence entre les tests de Spearman et de Kendall ?

→ Le test de Spearman mesure une corrélation monotone, compare les rangs transformés des deux variables, puis calcule une corrélation sur les rangs tandis que le test de Kendall mesure la probabilité d'accord entre les paires d'observations, il compare un par un tous les couples concordants et discordants en se basant sur la circonstance du classement et non sur les écarts contrairement au test de Spearman. Le test de Spearman est par ailleurs plus sensible aux valeurs aberrantes et aux ex aequo.

5- A quoi servent les coefficients de Goodman-Kruskal et de Yule ?

→ Le coefficient de Goodman-Kruskal se base sur la différence entre les paires concordantes et les paires discordantes. Il calcule le surplus de paires concordantes par rapport aux paires discordantes en

exerçant une proportion.

Le coefficient de Yule quant à lui est un cas particulier du coefficient de Goodman-Kruskal en ce qu'il

est appliqué dans le cas des matrices 2x2. Il est nécessaire de construire la table de contingence qui évalue la fréquence des événements.

Séance 6

1. Qu'est-ce qu'une statistique ordinale ? À quelle autre statistique catégorielle s'oppose-t-elle ?

Définition

La **statistique ordinale** (ou statistique d'ordre) est un concept central en géographie humaine. Elle consiste à établir des **classements** ou des **rangs** entre des unités ou des objets géographiques (villes, régions, entités administratives). L'objectif est de suivre l'évolution de ces entités — si elles sont montées, descendues ou si elles ont stagné — dans une hiérarchie donnée (annuelle, mensuelle, etc.).

Opposition Catégorielle

La statistique ordinale s'oppose principalement à la **statistique nominale**.

- **Statistique Ordinale** : Utilise des **variables qualitatives ordinales** (les catégories sont ordonnées).
- **Statistique Nominale** : Utilise des variables qualitatives nominales (les catégories sont distinctes mais sans ordre logique).

Matérialisation d'une Hiérarchie Spatiale

L'utilisation de variables qualitatives ordinales permet de matérialiser une **hiérarchie spatiale** dès lors que les catégories appliquées aux lieux ou territoires représentent un niveau, un rang ou une intensité. Par exemple :

- Elle établit des relations de **centralité ou de domination** (ordre entre le centre et la périphérie).
- Elle structure l'espace (ordre du centre urbain dense jusqu'à la ruralité).
- **Cartographiquement**, elle produit des cartes en **classes ordonnées**.

En **géographie humaine**, ces statistiques découlent de l'apparition spontanée de hiérarchies sociales ou spatiales, tandis qu'en **géographie physique**, elles peuvent servir à analyser l'ordre de grandeur de phénomènes extrêmes (hauteur maximale des crues, intensité sismique maximale, etc.).

2. Quel ordre est à privilégier dans les classifications ?

L'ordre à privilégier pour la classification des données est l'**ordre croissant**, souvent appelé l'**ordre naturel**. Cet ordre est généralement utilisé, bien qu'il existe des exceptions notables en géographie, comme l'application de la **loi rang-taille** (ou Zipf) qui classe les villes par ordre décroissant de population. L'ordination permet également de **repérer les valeurs aberrantes** (extrêmement grandes ou petites) au sein d'une série d'observations.

3. Quelle est la différence entre une corrélation des rangs et une concordance de classements ?

Caractéris	Corrélation des Rangs (Spearman, r_s)	Concordance des Classements (Kendall's τ)
Objectif	Mesurer la force et le sens de la relation	Mesurer le degré d'accord général entre
Nombre de	Deux variables ordonnées.	Plus de deux classements.
Question	Les deux classements sont-ils proportionnellement liés ? (Vont-ils dans le	À quel point les classements sont-ils identiques entre eux ? (Y a-t-il un accord

4. Quelle est la différence entre les tests de Spearman et de Kendall ?

Les deux tests mesurent la corrélation des rangs (l'association monotone) entre deux variables, mais ils le font de manière différente :

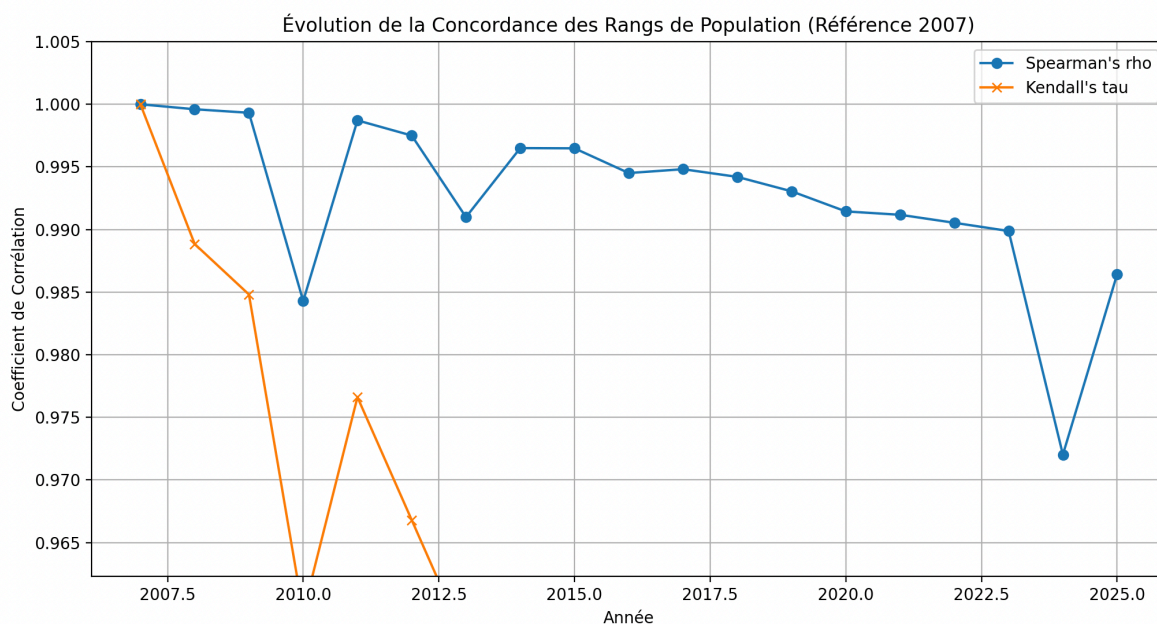
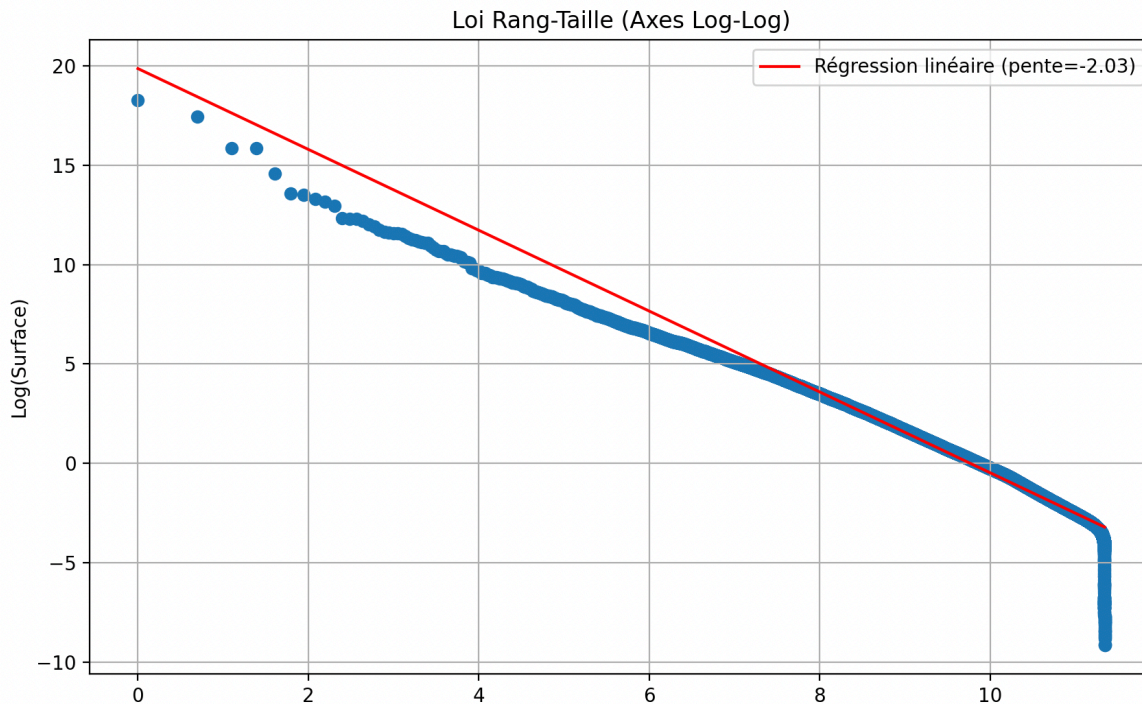
Caractéristique	Test de Corrélation de Spearman (ρ)	Test de Corrélation de Kendall (τ)
Méthode de	Compare les rangs transformés des deux variables et calcule la corrélation sur	Mesure la probabilité d'accord en comparant les paires d'observations
Sensibilité	Plus sensible aux écarts extrêmes entre les rangs et aux <i>ex aequo</i> .	Moins sensible aux valeurs aberrantes ; se concentre sur l'ordre pur plutôt que sur les
Interprétation	Similaire au coefficient de Pearson appliqué aux rangs.	Basé sur la circonstance du classement (l'ordre) et non sur la distance entre les rangs.

5. À quoi servent les coefficients de Goodman-Kruskal et de Yule ?

Ces coefficients sont utilisés pour mesurer l'association entre des variables catégorielles, souvent dans le contexte des tables de contingence.

- **Coefficient γ de Goodman-Kruskal :**

- Il mesure l'association entre deux variables **ordinales**.
- Il se base sur la différence entre le nombre de **paires concordantes** et le nombre de



discordantes.

paires

- Il calcule la **proportion** du surplus de paires concordantes par rapport aux paires discordantes, indiquant la force de l'association.
- **Coefficient Q de Yule :**
 - Il est un **cas particulier** du coefficient de Goodman-Kruskal.
 - Il est spécifiquement appliqué aux **matrices 2×2** (deux variables binaires).
 - Il est utilisé pour évaluer la fréquence d'événements conjoints dans cette table de contingence binaire.

ÉTAPE 1 : THÉORIE DE L'ÉCHANTILLONNAGE

=====

[1.1] Moyennes des 100 échantillons (arrondies) :

```
Pour      391
Contre    416
Sans_opinion 193
dtype: int64
```

[1.2] Fréquences des moyennes des échantillons (arrondies à 2 décimales) :

```
Pour      0.39
Contre    0.42
Sans_opinion 0.19
dtype: float64
```

[1.3] Fréquences de la population mère (arrondies à 2 décimales) :

```
Pour      0.39
Contre    0.42
Sans_opinion 0.19
```

dtype: float64

[1.4] Intervalle de fluctuation (IF) à 95% (n moyen=1000) :

IF 'Pour' (p=0.39): [0.3598 ; 0.4202]

IF 'Contre' (p=0.42): [0.3894 ; 0.4506]

IF 'Sans_opinion' (p=0.19): [0.1657 ; 0.2143]

ÉTAPE 2 : THÉORIE DE L'ESTIMATION

[2.1] Premier échantillon (liste) : [395, 396, 209]

Taille de l'échantillon (n) : 1000

[2.3] Fréquences du premier échantillon (p_hat) :

Pour 0.395

Contre 0.396

Sans_opinion 0.209

Name: 0, dtype: float64

[2.4] Intervalle de confiance (IC) à 95% (n=1000) :

IC 'Pour' (p_hat=0.3950): [0.3647 ; 0.4253]

IC 'Contre' (p_hat=0.3960): [0.3657 ; 0.4263]

IC 'Sans_opinion' (p_hat=0.2090): [0.1838 ; 0.2342]

ÉTAPE 3 : THÉORIE DE LA DÉCISION (Test de Shapiro-Wilks)

Résultats pour Loi-normale-Test-1.csv (Taille: 2000)

Statistique W : 0.9639

P-value : 0.0000

Conclusion : P-value (0.0000) \leq alpha (0.05). Rejette H0.

Résultat : La distribution n'est PAS considérée comme normale.

Résultats pour Loi-normale-Test-2.csv (Taille: 2000)

Statistique W : 0.2609

P-value : 0.0000

Conclusion : P-value (0.0000) \leq alpha (0.05). Rejette H0.

Résultat : La distribution n'est PAS considérée comme normale.

=====

FIN DES CALCULS STATISTIQUES.

=====

=====

ANALYSE BONUS : LOIS NON NORMALES

=====

Statistiques descriptives pour Loi-normale-Test-1.csv :

- Taille (n) : 2000.0000
- Minimum : -2.0000
- Maximum : 10.0000
- Étendue (Max - Min) : 12.0000
- Moyenne : 3.0430
- Médiane : 3.0000

- Écart-type : 1.5357

=> Caractérisation (pour le Bonus) :

Les statistiques ne permettent pas de trancher facilement sans visualisation (histogramme).

Statistiques descriptives pour Loi-normale-Test-2.csv :

- Taille (n) : 2000.0000

- Minimum : 1.0000

- Maximum : 14.0000

- Étendue (Max - Min) : 13.0000

- Moyenne : 1.1875

- Médiane : 1.0000

- Écart-type : 0.7626

=> Caractérisation (pour le Bonus) :

La moyenne et la médiane sont éloignées. La distribution est fortement asymétrique.

L'hypothèse d'une loi Exponentielle ou d'une autre loi asymétrique est forte.

1. Loi Rang-Taille (Axes Log-Log)

Ce graphique utilise des axes logarithmiques pour représenter la relation entre le rang et la surface (ou la taille) des entités.

- **Observations :** La majorité des points se situent très près d'une droite de régression. Cela indique que la distribution des surfaces suit une **Loi Rang-Taille** (ou Loi de Zipf), typique des systèmes hiérarchiques (comme la taille des villes, des entreprises, ou des îles).
- **Pente :** La régression linéaire présente une pente de **-2.03**. Dans la Loi de Zipf canonique, la pente est de -1 . Une pente plus accentuée (plus négative) indique que la décroissance de la taille est **plus rapide** par rapport au rang. En d'autres termes, les quelques entités en tête de classement sont **disproportionnellement plus grandes** que les suivantes.
- **Extrémités :** La courbe s'écarte légèrement de la droite en haut à gauche (quelques entités très grandes) et se recourbe fortement en bas à droite (un très grand nombre d'entités très petites).

2. Répartition des Îles par Intervalle de Surface (Histogramme)

Ce graphique est un histogramme montrant le nombre d'îles pour différentes classes de surface.

- **Observations :** Le graphique révèle une distribution extrêmement **asymétrique positive (en forme de L)**.
- **Concentration :** Près de 80 000 îles (la quasi-totalité de l'échantillon) ont une surface comprise entre **0 et 10 km²** ($]0, 10]$).
- **Rareté :** Toutes les autres classes d'intervalles de surface, même les classes intermédiaires, contiennent un nombre d'îles si faible qu'elles sont à peine visibles sur l'échelle du graphique.
- **Conclusion :** Cette distribution est une confirmation visuelle de l'analyse Rang-Taille : la grande majorité des îles sont petites, et les grandes îles sont un phénomène extrêmement rare.

3. Distribution du Nombre d'Inscrits (Histogramme)

Ce graphique illustre la répartition des effectifs d'inscrits par entité territoriale.

- **Observations :** Il s'agit d'une distribution **fortement asymétrique positive** (étalée vers la droite).
- **Concentration (Mode) :** La majorité des entités ont un nombre d'inscrits faible. Le mode (la classe la plus fréquente) se situe autour de **0.15×10^6 à 0.25×10^6** (150 000 à 250 000 inscrits).
- **Valeurs Extrêmes :** L'étirement de la queue vers la droite indique la présence de quelques entités avec un nombre d'inscrits exceptionnellement élevé (par exemple, autour de 1.5×10^6), correspondant aux départements les plus peuplés.
- **Conclusion :** La distribution n'est **pas Normale (Gaussienne)**. La moyenne des inscrits sera supérieure à la médiane et au mode, car elle est tirée vers la droite par ces grandes entités (comme le Nord ou les Bouches-du-Rhône).

4. Évolution de la Concordance des Rangs de Population (Séries Chronologiques)

Ce graphique montre l'évolution des coefficients de corrélation des rangs (Spearman's ρ et Kendall's τ) entre l'année de référence (2007) et les années suivantes.

- **Cohérence Élevée :** Les deux coefficients sont constamment très proches de **1.000**, ce qui indique une **très forte concordance** entre le classement des entités territoriales en 2007 et leur classement dans les années suivantes.
- **Stabilité :** Le classement de la population des entités territoriales (villes, départements, etc.) est **extrêmement stable** dans le temps. Les grandes entités de 2007 sont restées les grandes entités en 2025.
- **Sensibilité :** Le coefficient de **Kendall's τ** (ligne orange) est visiblement plus sensible aux variations et présente des creux plus profonds (notamment autour de 2010 et 2012) que le coefficient de Spearman's ρ .

Étape 1 : Théorie de l'Échantillonnage

Cette étape compare les proportions de la population mère (connues) aux proportions moyennes observées sur 100 échantillons, et définit l'intervalle dans lequel un futur échantillon devrait se situer (Intervalle de Fluctuation).

Opinion	Effectif Moyen (Arrondi)	Fréquence Moyenne p^{\wedge} moy	Fréquence Mère p	Intervalle de Fluctuation (IF à 95%, n=1000)
Pour	391	0.39	0.39	[0.3598 ; 0.4202]
Contre	416	0.42	0.42	[0.3894 ; 0.4506]
Sans opinion	193	0.19	0.19	[0.1657 ; 0.2143]
Total	1000	1.00	1.00	

Commentaire : On observe une **concordance parfaite** entre les fréquences moyennes des 100 échantillons et les fréquences réelles de la population mère. L'Intervalle de Fluctuation (IF) est étroit, confirmant que l'échantillonnage est précis. L'IF indique que pour 95 % des échantillons de taille 1000, la proportion observée pour l'opinion "Contre" se situera, par exemple, entre 38.94 % et 45.06 %.

Étape 2 : Théorie de l'Estimation

Cette étape estime, à partir d'un seul échantillon, la plage de valeurs dans laquelle le paramètre de la population mère se situe réellement (Intervalle de Confiance).

- **Premier Échantillon (Effectifs) :** [Pour: 395, Contre: 396, Sans_opinion: 209]
- **Taille de l'échantillon (n) :** 1000

Opinion	Fréquence du 1er Échantillon (p^{\wedge})	Intervalle de Confiance (IC à 95%, n=1000)
Pour	0.3950	[0.3647 ; 0.4253]
Contre	0.3960	[0.3657 ; 0.4263]
Sans opinion	0.2090	[0.1838 ; 0.2342]

Commentaire : L'Intervalle de Confiance (IC) montre la fiabilité de l'échantillon. Pour toutes les opinions, l'IC à 95 % **contient la vraie fréquence de la population mère** (0.39, 0.42, 0.19), confirmant que cet échantillon est représentatif. Par exemple, l'opinion "Contre" est estimée se situer entre 36.57 % et 42.63 % dans la population totale avec 95 % de certitude.

Étape 3 : Théorie de la Décision (Test de Shapiro-Wilks)

Le Test de Shapiro-Wilks évalue si une distribution suit une Loi Normale (H_0).

Test-1	2000	0.9639	0.0000	$P \leq \alpha$. Rejet de H_0 .	NON NORMALE
Test-2	2000	0.2609	0.0000	$P \leq \alpha$. Rejet de H_0 .	NON NORMALE

Commentaire : Dans les deux cas, la P-value est de 0.0000, ce qui est très inférieur au seuil de signification ($\alpha=0.05$). L'hypothèse que les données suivent une loi normale est donc **fortement rejetée** pour les deux jeux de données.

Analyse Bonus : Caractérisation des Lois Non Normales

Cette analyse utilise les statistiques descriptives pour identifier la nature des lois non normales.

Statistique	Loi-normale-Test-1.csv	Loi-normale-Test-2.csv	Caractérisation Clé
Moyenne	3.0430	1.1875	
Médiane	3.0000	1.0000	
Écart-type	1.5357	0.7626	
Comparaison Moyenne/Médiane	Moyenne \approx Médiane	Moyenne $>$ Médiane	Asymétrie

Commentaire :

- Loi-normale-Test-1.csv :** La moyenne et la médiane sont presque identiques ($3.0430 \approx 3.0000$). Cela indique une distribution **plutôt symétrique**. Sans visualisation (histogramme), il est difficile de la caractériser précisément (elle pourrait être Uniforme ou une cloche non parfaitement normale).
- Loi-normale-Test-2.csv :** La moyenne (1.1875) est nettement supérieure à la médiane (1.0000). Cet écart significatif est la signature d'une distribution **fortement asymétrique positive** (étalée à droite), typique des **Lois Exponentielles** ou Gamma.