

Méthodologie

Source des données

L'étude repose sur des données issues du Système National des Données de Santé (SNDS), couvrant la période 2018-2022 et portant sur environ 5000 communes. Ces données permettent d'analyser le nombre de consultations en médecine de ville en tenant compte des disparités territoriales et des caractéristiques locales. Nous avons trois bases essentielles. La première source est une base démographique contenant des données détaillées sur la répartition de la population par sexe et par tranche d'âge, ainsi que des indicateurs généraux tels que la population municipale et la structure des ménages. La deuxième source est une base généralisée intégrant des informations socio-économiques, notamment sur les statuts matrimoniaux, les catégories socio-professionnelles et le marché de l'emploi. Cette base permet d'étudier la composition sociale des communes et d'évaluer certains phénomènes tels que le taux d'activité ou la prévalence des unions libres. Enfin, une troisième base a été utilisée pour compléter les données géographiques, en particulier pour renseigner les latitudes et longitudes manquantes de certaines communes.

Les informations utilisées concernent principalement le volume des consultations médicales et leur répartition géographique. Des données complémentaires sur le contexte communal, telles que la densité de population et l'accessibilité aux soins, permettent d'affiner l'analyse. L'intégration de ces éléments facilite une approche spatiale de la modélisation, essentielle pour détecter d'éventuelles inégalités d'accès aux soins.

L'ensemble des données a été anonymisé et traité conformément aux normes en vigueur, garantissant ainsi la confidentialité des informations exploitées. Cette base constitue une ressource précieuse pour mieux comprendre les dynamiques d'accès aux soins en médecine de ville et proposer des modèles adaptés aux spécificités territoriales.

Traitement réalisés sur la base de données

Après la fusion des bases de données, un ensemble de traitements a été réalisé afin de structurer les informations et garantir la cohérence des analyses. La première étape a consisté à nettoyer et harmoniser les noms des variables pour faciliter leur manipulation. Ce travail a inclus la suppression des accents, le remplacement des espaces et caractères spéciaux par des underscores et la conversion en minuscules. De plus, certaines incohérences ont été corrigées, notamment des erreurs typographiques dans les intitulés des tranches d'âge.

Une fois cette normalisation effectuée, les différentes bases ont été fusionnées. Une jointure interne a été réalisée entre les bases démographiques et socio-économiques en utilisant le code unique des communes, ce qui a permis de conserver uniquement les communes présentes dans les deux sources. Ensuite, une jointure gauche avec la base des coordonnées géographiques a été effectuée pour compléter les informations manquantes sur la latitude et la longitude des communes concernées. Lors de cette étape, une attention particulière a été portée pour éviter la duplication des données et sélectionner les valeurs les plus pertinentes en cas de conflit entre plusieurs sources.

L'étape suivante a consisté à structurer et créer de nouvelles variables analytiques. En ce qui concerne la répartition de la population, les différentes tranches d'âge disponibles ont été regroupées en trois grandes catégories : 0-24 ans, 25-64 ans et 65 ans et plus. Ce regroupement permet d'obtenir une vision synthétique de la structure démographique tout en conservant des distinctions essentielles entre les jeunes, les actifs et les seniors. Pour chacune de ces nouvelles catégories, un indicateur de proportion a été calculé en rapportant la population de chaque groupe à la population totale de la commune.

Un travail spécifique a été réalisé pour enrichir l'analyse socio-économique. Le pourcentage des personnes en union libre a été calculé à partir des données disponibles sur les statuts matrimoniaux, et la part des ouvriers a été estimée en rapportant cette catégorie à la population totale. Le pourcentage des personnes sans emploi a été évalué en prenant comme référence la population âgée de 0 à 64 ans, excluant ainsi les retraités pour obtenir un indicateur plus pertinent sur la population active. De plus, un indicateur mesurant la proportion des familles avec trois enfants ou plus a été créé afin d'étudier la dynamique familiale dans les communes étudiées.

Un des traitements majeurs effectués a été la création d'une variable inédite : le taux de consultation. Cette variable, qui n'existait pas initialement dans la base, a été construite en rapportant le nombre total de consultations enregistrées dans chaque commune à la population municipale. Cette mesure permet d'évaluer la fréquence des consultations médicales ou administratives en fonction de la taille de la population et constitue un indicateur clé pour analyser l'accessibilité et l'utilisation des services sur le territoire étudié. Une version spécifique de cet indicateur a également été calculée pour la population de 19 ans et plus, en excluant les plus jeunes afin de mieux capter les tendances de consultation chez les adultes.

Enfin, un effort particulier a été porté sur la visualisation des données, notamment avec la construction d'une pyramide des âges. Pour cela, les effectifs masculins et féminins ont été extraits et réorganisés par tranche d'âge afin de respecter la convention des pyramides démographiques, affichant les populations masculines en valeurs négatives et les populations féminines en valeurs positives. Cette représentation permet d'identifier les déséquilibres entre les classes d'âge et de mieux comprendre la structure démographique des communes analysées.

Une fois ces traitements finalisés, la base de données enrichie a été exportée sous un format exploitable, intégrant l'ensemble des nouvelles variables créées ainsi que les indicateurs de taux de natalité et de mortalité. Ce travail de préparation assure une qualité optimale des données et permet de mener des analyses détaillées sur la dynamique démographique et socio-économique des communes étudiées.

Concepts fondamentaux en statistique spatiale

Autocorrélation spatiale

L'autocorrélation spatiale désigne la dépendance statistique entre des observations géographiquement proches. En d'autres termes, les valeurs prises par une variable en un lieu donné sont influencées par les valeurs observées dans les localisations voisines. Cette dépendance peut être positive, lorsque des valeurs similaires se regroupent, ou négative, lorsqu'une valeur élevée en un point est associée à une valeur faible dans les zones environnantes.

Diagramme de Moran

Le diagramme de Moran est un outil permettant d'analyser la structure spatiale d'une variable. Il représente un nuage de points où :

- L'axe des abscisses affiche les valeurs centrées de la variable d'intérêt y .
- L'axe des ordonnées affiche les valeurs moyennes de cette variable pour les observations voisines Wy , où W est la matrice de poids normalisée.

Interprétation En raison du centrage de y et de la normalisation de W , la moyenne empirique de Wy est égale à zéro.

- **Si les observations sont distribuées aléatoirement dans l'espace**, il n'existe pas de relation particulière entre y et Wy . La pente de la droite de régression linéaire est alors proche de zéro, et les observations sont réparties de manière uniforme dans les quadrants.
- **Si une structure spatiale existe**, la pente de la régression linéaire est différente de zéro, indiquant une corrélation entre y et Wy .

Le diagramme est divisé en **quatre quadrants** définis par les lignes $y = 0$ et $Wy = 0$, correspondant à différents types d'association spatiale :

1. Quadrant 1 (haut à droite, High-High)

- Valeurs élevées de la variable y entourées par des valeurs également élevées.
- Indique une **autocorrélation spatiale positive**.
- Peut refléter des zones homogènes où les valeurs sont naturellement élevées.

2. Quadrant 3 (bas à gauche, Low-Low)

- Valeurs faibles de y entourées par des valeurs également faibles.
- Indique une **autocorrélation spatiale positive**.
- Suggère la présence de zones où les valeurs sont homogènement basses.

3. Quadrant 2 (bas à droite, High-Low)

- Valeurs élevées de y entourées par des valeurs faibles.
- Indique une **autocorrélation spatiale négative**.
- Peut signaler une hétérogénéité spatiale locale, où certaines observations diffèrent fortement de leur voisinage.

4. Quadrant 4 (haut à gauche, Low-High)

- Valeurs faibles de y entourées par des valeurs élevées.
- Indique une **autocorrélation spatiale négative**.
- Peut traduire des zones de contraste spatial marquées.

Utilité du Diagramme de Moran Le diagramme de Moran permet d'identifier les structures spatiales dominantes en observant la répartition des points dans les quadrants. Il aide également à détecter les observations atypiques qui s'écartent du modèle spatial général et à confirmer l'existence d'une autocorrélation spatiale en complément de l'indice de Moran, qui quantifie cette relation. Cet outil visuel est ainsi essentiel pour explorer la dépendance spatiale et comprendre les structures spatiales sous-jacentes d'une variable.

Matrice des poids spatiaux

Pour quantifier la proximité spatiale entre unités géographiques, on utilise une matrice de poids spatiaux notée W . Cette matrice représente les relations de voisinage et permet d'introduire la structure spatiale dans les modèles économétriques. La matrice de poids peut être une matrice de contiguité binaire ou peut tenir compte de la distance entre les zones géographiques. Cette étude utilise une matrice de poids basée sur la distance et contient des pondérations inversement proportionnelles à la distance entre les régions.

Mesure de la corrélation spatiale

Indices de corrélation spatiale L'un des indicateurs les plus couramment utilisés est l'indice de Moran. Il évalue la similitude des valeurs d'une variable entre différentes entités géographiques (par exemple, des communes) en fonction de leur proximité spatiale. Il se base sur la matrice de poids spatiale (W), qui définit les relations entre ces entités. Il se calcule comme suit :

$$I = \frac{N}{\sum_i \sum_j W_{ij}} \times \frac{\sum_i \sum_j W_{ij} (y_i - \bar{y})(y_j - \bar{y})}{\sum_i (y_i - \bar{y})^2}$$

où y_i est la valeur de la variable d'intérêt en un point i , \bar{y} est la moyenne de cette variable et W_{ij} représente l'élément (i, j) de la matrice de poids.

D'autres indices existent, comme la **statistique de Geary**, qui est moins sensible aux valeurs extrêmes, et les **indicateurs locaux d'autocorrélation spatiale (LISA)**, qui permettent d'identifier des clusters spatiaux spécifiques.

Construction de la matrice de poids

Pour construire la matrice de poids, nous avons alors suivi ces étapes.

1. Calculer les distances de Haversine entre chaque paire d'entités.
2. Définir un seuil de distance maximale (d_{max}) :
 - Si $d_{ij} < d_{max}$, $w_{ij} = \frac{1}{d_{ij}}$.
 - Sinon, $w_{ij} = 0$.
3. Normaliser les poids pour que chaque ligne de la matrice ait une somme égale à 1 :

$$w_{ij}^{norm} = \frac{w_{ij}}{\sum_j w_{ij}}.$$

Test significativité de l'indice de Moran : Le test de significativité de l'indice de Moran permet d'évaluer si une variable présente une autocorrélation spatiale significative, c'est-à-dire si les valeurs observées dans des zones proches ont tendance à être similaires ou non.

Hypothèses du test

- **Hypothèse nulle H_0 :** Il n'y a **pas d'autocorrélation spatiale** significative. Les valeurs observées sont distribuées de manière aléatoire dans l'espace.
- **Hypothèse alternative H_1 :** Il existe une **autocorrélation spatiale** significative (positive ou négative).

Statistique de test

L'Indice de Moran standardisé suit approximativement une distribution normale sous H_0 . La statistique de test est donnée par :

$$Z = \frac{I - E(I)}{\text{Var}(I)}$$

où :

- I est l'Indice de Moran calculé sur les données,
- $E(I)$ est l'espérance théorique de I sous H_0 ,
- $\text{Var}(I)$ est la variance théorique de I .

L'espérance sous H_0 pour un échantillon de taille n est donnée par :

$$E(I) = -\frac{1}{n-1}$$

Règle de décision

On compare la statistique Z à une loi normale centrée réduite $\mathcal{N}(0, 1)$. Pour un seuil de significativité α (ex. 5 %), on utilise les quantiles de la loi normale :

- Si $Z > z_{1-\alpha/2}$ ou $Z < -z_{1-\alpha/2}$, on **rejette** H_0 et on conclut qu'il existe une autocorrélation spatiale significative.
- Si $Z \in [-z_{1-\alpha/2}, z_{1-\alpha/2}]$, on **ne rejette pas** H_0 , et on considère que la distribution spatiale est aléatoire.

Interprétation de l'Indice de Moran

- $I > 0$ et significatif : Autocorrélation spatiale **positive** (les zones proches ont des valeurs similaires).
- $I < 0$ et significatif : Autocorrélation spatiale **négative** (les zones proches ont des valeurs opposées).
- $I \approx 0$ et non significatif : Absence d'autocorrélation spatiale, les valeurs sont distribuées de manière aléatoire.

Ce test est couramment utilisé en analyse spatiale pour identifier des regroupements de valeurs similaires, par exemple dans les études de santé publique, d'aménagement du territoire ou d'économie régionale.

Modélisation en économétrie spatiale

Voici un rappel des différents éléments utilisés dans l'ensemble des modèles d'économétrie spatiale :

- Y : Il s'agit du vecteur des observations de la variable dépendante, c'est-à-dire la variable que l'on cherche à expliquer (par exemple, le taux de visites, le taux de chômage, etc).
- X : C'est la matrice des variables explicatives ou indépendantes. Elle regroupe toutes les caractéristiques observées qui sont supposées influencer \mathbf{Y} (comme des variables socio-économiques, démographiques ou structurelles).
- β : Ce vecteur de coefficients mesure l'effet direct des variables \mathbf{X} sur la variable dépendante \mathbf{Y} . Chaque coefficient indique l'impact d'une unité de variation dans la variable correspondante sur \mathbf{Y} , en l'absence d'effets spatiaux.
- W : La matrice des poids spatiaux définit la structure de voisinage entre les unités géographiques. Chaque élément W_{ij} quantifie l'influence ou la proximité de l'unité j par rapport à l'unité i . Le choix de cette matrice (par contiguïté, distance, ou K plus proches voisins) est crucial car il détermine la manière dont l'information spatiale est intégrée dans le modèle.
- WY : Le terme de décalage spatial de \mathbf{Y} , obtenu par le produit de la matrice \mathbf{W} par le vecteur \mathbf{Y} . Il représente l'influence moyenne pondérée des valeurs de \mathbf{Y} dans les zones voisines et permet de capturer la dépendance spatiale directe de la variable dépendante.
- WX : Il s'agit du terme de décalage spatial des variables explicatives. Concrètement, il représente une version pondérée des variables \mathbf{X} dans les zones voisines, où les pondérations sont définies par la matrice \mathbf{W} . Ce terme permet de mesurer l'effet indirect (ou spillover) des caractéristiques des voisins sur \mathbf{Y} .
- ε : C'est le terme d'erreur classique, qui capture les influences non observées ou aléatoires sur \mathbf{Y} . Il est généralement supposé être indépendant et identiquement distribué (iid).
- ρ : Utilisé dans les modèles qui intègrent directement l'effet des valeurs voisines de \mathbf{Y} (comme dans les modèles SAR et SDM). Ce paramètre mesure la force de l'interaction entre la valeur de \mathbf{Y} d'une unité et les valeurs de \mathbf{Y} des unités voisines. Un ρ positif indique une autocorrélation positive (les zones avec des valeurs élevées de \mathbf{Y} tendent à être entourées de zones à valeurs élevées, et inversement).

λ : Spécifique au modèle SEM (Spatial Error Model), ce paramètre quantifie la corrélation spatiale présente dans le terme d'erreur. Il mesure l'influence des erreurs des unités voisines sur l'erreur de l'unité considérée, suggérant que des facteurs non observés présentent une structure spatiale.

θ : Ce vecteur de coefficients est associé au terme \mathbf{WX} et apparaît dans les modèles SDM et SLX. Il mesure l'effet des variables explicatives des zones voisines sur la variable dépendante \mathbf{Y} , c'est-à-dire l'impact indirect des caractéristiques locales via leur diffusion spatiale.

Modèles principaux

Le modèle général est défini comme suit :

$$Y = \rho WY + X\beta + \theta WX + u, \quad u = \lambda Wu + \varepsilon$$

SAR (Spatial AutoRegressive Model) :

Le modèle SAR introduit une dépendance spatiale directement sur la variable dépendante Y . L'idée est que la valeur de Y en un lieu donné dépend des valeurs observées dans les zones voisines. Mathématiquement, il s'écrit :

$$Y = \rho WY + X\beta + \varepsilon$$

Interprétation :

- Si $\rho > 0$, les valeurs de Y ont tendance à être similaires entre voisins (autocorrélation positive).
- Si $\rho < 0$, on observe un effet de dispersion, où les valeurs de Y sont opposées dans les zones voisines (autocorrélation négative).
- Si $\rho = 0$, il n'y a pas de dépendance spatiale, et le modèle classique de régression linéaire est suffisant.

SEM (Spatial Error Model) :

Le modèle SEM est utilisé lorsque la dépendance spatiale affecte les erreurs du modèle plutôt que la variable dépendante elle-même. Il est défini par :

$$Y = X\beta + u, \quad u = \lambda Wu + \varepsilon$$

Interprétation :

- Contrairement au modèle SAR, le modèle SEM suppose que la dépendance spatiale est un effet de perturbation, provenant d'omissions de variables pertinentes qui suivent une structure spatiale.
- Il est utilisé lorsque la corrélation spatiale détectée dans un modèle classique provient d'erreurs spatialement autocorrélées, plutôt que d'une interaction directe entre observations.

SLX (Spatial Lag of X Model) :

Le modèle SLX est plus facile à estimer, car il suppose que la variable dépendante Y n'est pas directement influencée par les valeurs voisines, mais uniquement par les variables explicatives des zones voisines. Il est écrit comme suit :

$$Y = X\beta + \theta WX + \varepsilon$$

où WX capture l'effet des variables explicatives des unités voisines.

Interprétation :

- Il n'y a pas d'effet direct des valeurs voisines de Y .
- Il mesure uniquement l'effet de "spillover" (d'effet de débordement) des facteurs explicatifs.

SDM (Spatial Durbin Model) :

Le modèle SDM est une extension du modèle SAR. Il prend en compte non seulement la dépendance de Y aux observations voisines, mais aussi l'effet des variables explicatives des régions voisines. Il est défini par :

$$Y = \rho WY + X\beta + \theta WX + \varepsilon$$

Interprétation :

- Si $\theta = 0$, le modèle SDM devient un SAR classique.
- Si $\rho = 0$, il devient un modèle SLX (voir ci-dessous).
- Il permet de tester si des variables exogènes influencent Y au-delà des frontières administratives.

Comparaison des modèles

Modèle	Dépendance spatiale sur Y	Effet des X des voisins	Effet des erreurs
SAR	Oui	Non	Non
SEM	Non	Non	Oui
SLX	Non	Oui	Non
SDM	Oui	Oui	Non

Table 1: Comparaison des modèles spatiaux