

Graphe probabiliste :

**Contamination du COVID-19
durant la pandémie à Bruxelles**

Mathématiques appliquées à l'informatique

Auteur : Van Overlop Cédric

Date : Juin 2025

Table des matières

1	Mise en situation	3
1.1	Contexte institutionnel	3
1.2	Objectifs du projet	3
1.3	Sources de données	3
2	Reprise et analyse du travail précédent	4
2.1	Travail de l'équipe 2022–2023	4
2.2	Résultats et limites identifiées	4
2.3	Enseignements tirés	5
3	Résolution mathématique du problème	6
3.1	Lissage temporel par filtre de Savitzky-Golay	6
3.2	Construction du graphe probabiliste géographique	7
3.3	Modèle de chaînes de Markov matriciel	7
3.4	Analyse mathématique et performances	8
3.5	Validation empirique et limites identifiées	9
3.6	Conclusion méthodologique	10
4	Implémentation informatique	11
4.1	Architecture du projet	11
4.2	Pipeline de traitement automatisé	11
4.2.1	Acquisition des données (<code>data_downloader.py</code>)	11
4.2.2	Organisation des données (<code>data_loader.py</code>)	12
4.3	Implémentation du lissage Savitzky-Golay	14
4.3.1	Optimisation algorithmique	14
4.4	Modélisation géographique avec Dijkstra	15
4.4.1	Construction du modèle géographique	15
4.5	Implémentation du modèle de Markov	17
4.5.1	Architecture orientée performance	17
4.6	Dashboard de visualisation (<code>main.py</code>)	20
4.7	Analyse complète	22
4.8	Performances et optimisations	22
4.8.1	Comparaison des performances	22
4.8.2	Validation des résultats	22
4.9	Reproductibilité et extensibilité	24
4.9.1	Configuration centralisée	24
4.9.2	Extensibilité et adaptabilité	24
4.9.3	Conditions de reproductibilité	24

4.9.4	Limitations de reproductibilité	24
4.9.5	Documentation de validation	25
4.9.6	Domaine de validité identifié	25
4.10	Conclusion technique	25
5	Conclusion	26
5.1	Les avancées du projet	26
5.2	La solution mathématique développée	26
5.3	Les performances obtenues et leurs limites	26
5.4	Enseignements tirés et valeur scientifique	27
5.5	L'implémentation informatique	27
5.6	Reconnaissance des limites et domaine de validité	28
5.7	Perspectives d'amélioration	28
5.8	Bilan final et apport pédagogique	28
5.9	Conclusion générale	29

1 Mise en situation

1.1 Contexte institutionnel

Sciensano est un institut public national belge opérant sous l'autorité du ministère fédéral de la Santé publique et de l'Agriculture. Sa mission principale est de faire progresser la recherche dans des domaines variés tels que la santé publique, la santé animale, la qualité des soins, les vaccins, l'écologie ainsi que les maladies infectieuses.

L'institut met à disposition du public les résultats de ses recherches via son site internet, facilitant ainsi l'accès à des données fiables concernant, entre autres, la grippe aviaire, les virus influenza ou encore le COVID-19.

Source : <https://www.sciensano.be/fr>

1.2 Objectifs du projet

Dans le cadre de la première année du bachelier en informatique orientation développement d'applications, ce projet constitue une opportunité concrète d'appliquer les compétences acquises au fil du cursus. L'objectif est de concevoir un **modèle prédictif avancé** basé sur les données de Sciensano, capable de modéliser l'évolution de la contamination et sa propagation à travers les **19 communes de Bruxelles**.

Ce travail vise à représenter la **dynamique intercommunale de transmission du COVID-19** en utilisant des **modèles mathématiques probabilistes**, tout en tenant compte des contraintes géographiques. L'enjeu est de fournir des insights pertinents sur les facteurs influençant la diffusion du virus.

1.3 Sources de données

Les données utilisées dans ce projet proviennent des sources suivantes :

- **Source principale** : <https://epistat.wiv-isp.be/covid/>
- **API OpenDataSoft** : Accès programmatique aux données
- **Volume** : Plus de 457 000 enregistrements couvrant l'ensemble du territoire belge
- **Période couverte** : 2020 - 2023
- **Granularité** : Cas quotidiens par commune

2 Reprise et analyse du travail précédent

2.1 Travail de l'équipe 2022–2023

L'équipe précédente a mené une exploration rigoureuse des données issues de Sciensano en adoptant une démarche structurée comprenant plusieurs étapes :

Traitement initial des données

- Téléchargement des données au format JSON
- Parsing et conversion en structures exploitables sous Python
- Manipulations matricielles avancées (inversions, opérations)

Modélisation mathématique

- Construction d'une **topologie de graphe** optimisée, où chaque commune est représentée par un nœud
- Optimisation de la structure pour améliorer les performances computationnelles
- Application de techniques telles que la **régression linéaire** et la **projection LOGIT** pour obtenir des prédictions probabilistes

Ce cadre méthodologique a permis une modélisation binaire (présence ou absence du virus) adaptée aux données de contamination géographique.

2.2 Résultats et limites identifiées

Approche par oscillateurs

- Génération d'un grand nombre de matrices probabilistes à partir d'un vecteur de base (5 communes)
- Raffinement par ligne à l'aide d'**oscillateurs**, produisant jusqu'à $5^{12} = 248\,832$ matrices par itération

Optimisations techniques

- Répartition du traitement en quatre sous-processus parallèles
- Génération de fichiers JSON contenant les matrices retenues
- Structure sous forme de dictionnaires (clé = date, valeur = matrice + résultat)
- Introduction de techniques de **machine learning** pour prédire les matrices futures

Limites identifiées

- **Complexité computationnelle élevée :**
 - Multiples boucles imbriquées
 - Fichiers de sortie nombreux et difficiles à analyser
 - Temps d'exécution de plusieurs heures
- **Problème de passage à l'échelle :**
 - Méthode efficace pour 5 communes
 - Extension à 19 communes impossible à cause de l'explosion combinatoire
- **Échecs d'optimisation :**
 - Méthode de Newton-Raphson abandonnée pour cause de complexité
 - Algorithme EM testé, mais insuffisant face à la haute dimension

2.3 Enseignements tirés

L'analyse de ces échecs successifs a mis en évidence la nécessité d'un **changement d'approche radical**.

Approches tentées sans succès :

- **Oscillateurs** : explosion combinatoire insurmontable
- **Newton-Raphson** : complexité algorithmique équivalente
- **Algorithme EM** : erreurs techniques multiples, résultats inexploitable

Lacunes techniques majeures :

- Absence de modélisation géographique
- Aucune validation systématique des prédictions
- Non-prise en compte des interactions intercommunales

Seuls acquis utilisables :

- Bonne compréhension du format des données Sciensano
- Constat clair du problème de scalabilité
- Conscience de la nécessité d'une approche plus efficiente

Ces constats ont mené à la conception d'une **approche matricielle entièrement nouvelle**, intégrant la dimension géographique et capable de gérer efficacement les 19 communes bruxelloises.

3 Résolution mathématique du problème

Face aux limitations identifiées, nous développons une approche mathématique intégrée combinant trois techniques complémentaires pour résoudre efficacement le problème de modélisation à 19 communes en utilisant des **graphes probabilistes** et des **chaînes de Markov**.

3.1 Lissage temporel par filtre de Savitzky-Golay

Principe du lissage

Les données quotidiennes de COVID-19 présentent des anomalies qui ne reflètent pas la vraie évolution de l'épidémie :

- **Week-ends** : moins de déclarations le samedi/dimanche,
- **Jours fériés** : retards administratifs (Noël, Pâques, etc.),
- **Variations administratives** : certains laboratoires déclarent en retard,
- **Erreurs de saisie** : pics artificiels ou valeurs aberrantes.

Exemple concret :

Lundi : 15 cas (rattrapage du week-end)
Mardi : 8 cas (normal)
Mercredi : 12 cas (normal)
Jeudi : 11 cas (normal)
Vendredi : 9 cas (normal)
Samedi : 3 cas (sous-déclaration week-end)
Dimanche : 1 cas (sous-déclaration week-end)
→ Tendence réelle 10 cas/jour

Fonctionnement mathématique du filtre

Le filtre de Savitzky-Golay repose sur un ajustement polynomial local :

$$\text{Valeur lissée} = \frac{-2 \cdot J_{-3} + 3 \cdot J_{-2} + 6 \cdot J_{-1} + 7 \cdot J_0 + 6 \cdot J_{+1} + 3 \cdot J_{+2} - 2 \cdot J_{+3}}{21}$$

Transformation :

$$y_i(t) \in \mathbb{N} \quad \rightarrow \quad \tilde{y}_i(t) \in \mathbb{R}^+$$

3.2 Construction du graphe probabiliste géographique

Modélisation par graphe pondéré

La géographie des 19 communes de Bruxelles est modélisée comme un graphe pondéré $G = (V, E, W)$, où :

- V : ensemble des 19 communes (sommets),
- E : arêtes représentant les frontières communes,
- W : poids des arêtes = longueur des frontières.

Calcul des influences géographiques

Pour chaque paire de communes (i, j) :

$$\text{Influence}(i \rightarrow j) = \frac{1}{\text{distance}_{\text{Dijkstra}}(i, j) + \varepsilon} \quad \text{où } \varepsilon = 0.1$$

Cela donne une matrice stochastique $G \in \mathbb{R}^{19 \times 19}$:

$$\sum_j G[i, j] = 1$$

3.3 Modèle de chaînes de Markov matriciel

Formulation du modèle

Le système est modélisé comme une chaîne de Markov d'ordre 1 :

$$\vec{X}(t+1) = A \cdot \vec{X}(t) + \vec{\varepsilon}(t)$$

où :

- $\vec{X}(t) \in \mathbb{R}^{19}$: état au temps t ,
- $A \in \mathbb{R}^{19 \times 19}$: matrice de transition,
- $\vec{\varepsilon}(t)$: bruit gaussien.

Estimation de la matrice de transition

Objectif :

$$\min_A \|\vec{X}(t+1) - A \cdot \vec{X}(t)\|^2 + \lambda \|A\|^2$$

Solution analytique :

$$A = X_{t+1} \cdot X_t^\top \cdot (X_t \cdot X_t^\top + \lambda I)^{-1}$$

où :

- $X_t \in \mathbb{R}^{19 \times T}$: observations aux temps t ,
- $X_{t+1} \in \mathbb{R}^{19 \times T}$: observations à $t + 1$,
- $\lambda = 10^{-6}$: régularisation Ridge.

Intégration des contraintes géographiques

$$A_{\text{finale}} = (1 - \alpha) \cdot A + \alpha \cdot (A \circ G)$$

où :

- \circ : produit de Hadamard (élément par élément),
- $\alpha \in [0, 1]$: pondération géographique.

Stabilisation :

1. Renormalisation :

$$A_{\text{finale}}[i, :] \leftarrow \frac{A_{\text{finale}}[i, :]}{\sum_j A_{\text{finale}}[i, j]}$$

2. Contrôle du rayon spectral :

$$\text{si } \rho(A_{\text{finale}}) > 1, \quad A_{\text{finale}} \leftarrow \frac{A_{\text{finale}}}{\rho(A_{\text{finale}})}$$

3. Positivité :

$$A_{\text{finale}} \leftarrow \max(A_{\text{finale}}, 0)$$

Optimisation du paramètre α

$$\alpha^* = \arg \min_{\alpha} \text{MAE}_{\text{validation}}(\alpha)$$

3.4 Analyse mathématique et performances

Stabilité du système

Théorème de convergence : Si A_{finale} est stochastique avec rayon spectral ≤ 1 :

$$\lim_{n \rightarrow \infty} A_{\text{finale}}^n \cdot \vec{X}(0) \rightarrow \vec{X}_{\infty}$$

Comparaison des performances

Métrique	Ancienne approche	Nouvelle approche	Amélioration
Complexité	$\mathcal{O}(5^{19})$	$\mathcal{O}(19^2)$	$\times 10^9$
Temps de calcul	Plusieurs heures	< 30 secondes	$\times 400$
Nombre de communes	5 maximum	19+	Extensible
Mémoire requise	Plusieurs GB	< 100 MB	$\times 10$

Avantages et innovations

- Extension de 5 à 19 communes,
- Pipeline reproductible et modulaire,
- Intégration originale géographie + Markov,
- Stabilité numérique assurée.

3.5 Validation empirique et limites identifiées**Période stable (septembre 2021)**

- Entraînement : juin-août 2021,
- Test : septembre 2021,
- MAE globale : 0.69 cas/jour/commune,
- Corrélations : 0.196 à 0.434,
- Communes les plus précises : Forest, Ganshoren, Ixelles.

Échec en période instable (mars-juin 2021)

- Entraînement : janvier-février 2021,
- Prédiction systématiquement décroissante,
- Corrélations négatives fréquentes,
- Hypothèse de stationnarité violée (changements de régime).

Limites identifiées

1. Modèle linéaire vs dynamique exponentielle,
2. Stationnarité vs régimes changeants,
3. Géographie locale vs mobilité/politique,
4. Horizon court (30 jours max).

Domaine de validité

- **Applications recommandées :**
 - Périodes stables,
 - Horizons courts (7–30 jours),
 - Objectifs pédagogiques,
 - Régions à géographie bien définie.
- **Applications déconseillées :**
 - Nouvelles vagues,
 - Changements de politiques,
 - Long terme (> 30 jours),
 - Facteurs comportementaux dominants.

3.6 Conclusion méthodologique

Cette approche constitue une **démonstration de faisabilité** de l'intégration de contraintes géographiques dans un modèle épidémiologique matriciel.

Valeur scientifique :

- Validation empirique des concepts,
- Identification rigoureuse des limites,
- Pipeline généralisable,
- Base solide pour des extensions futures.

Limitation assumée : Outil de démonstration pédagogique plus que modèle prédictif opérationnel, ce qui renforce la crédibilité scientifique dans son domaine de validité.

4 Implémentation informatique

4.1 Architecture du projet

Face à la complexité de l'approche précédente, j'ai conçu une **architecture modulaire claire** qui sépare chaque étape du processus :

```
Projet COVID/
  config/
    settings.json          # Configuration centralisée
  Data_Processing/
    data_downloader.py     # Acquisition et traitement des données
    data_loader.py         # Téléchargement automatisé via API
    data_loader.py         # Organisation et nettoyage
    savitzky_golay.py      # Lissage temporel
  geography/
    dijkstra.py            # Calculs géographiques
  markov_model/
    Prediction.py          # Modèle Markov matriciel
  main.py                  # Dashboard de visualisation
  prediction_historique.py # Prédiction pour Septembre 2021
  run_complete_analysis.py # Lancement de tous les fichiers
  data/                    # Stockage des résultats
```

Avantages de cette architecture

- **Modularité** : Chaque fichier a une responsabilité unique
- **Réutilisabilité** : Modules indépendants et testables
- **Maintenabilité** : Code organisé et documenté
- **Extensibilité** : Ajout facile de nouvelles fonctionnalités

4.2 Pipeline de traitement automatisé

4.2.1 Acquisition des données (data_downloader.py)

Téléchargement automatisé via l'API Export d'OpenDataSoft.

Résultats obtenus :

- ✓ 457,066 enregistrements téléchargés automatiquement
- ✓ 20,900 entrées pour les 19 communes bruxelloises
- ✓ Validation automatique des communes
- ✓ Temps d'exécution < 2 minutes

```

14 def download_via_export_api():
15     """
16     Télécharge via l'API Export qui n'a pas de limite de 10k
17     """
18     print("=====")
19     print(Fore.YELLOW + "Téléchargement via API Export (dataset complet)")
20     print("=====")
21
22     temps_debut = datetime.datetime.now()
23
24     # API Export URL - téléchargement TOUT
25     export_url = "https://public.opendatasoft.com/api/explore/v2.1/catalog/datasets/covid-19-pandemic-belgium-cases-municipality/exports/json"
26
27     try:
28         print("🔵 Téléchargement du dataset complet...")
29         print("⚠️ Cela peut prendre 1-2 minutes...")
30
31         response = requests.get(export_url, timeout=180) # 3 minutes timeout
32         response.raise_for_status()
33
34         # L'API Export retourne directement un JSON
35         data = response.json()
36         print(f"✅ {len(data)} enregistrements téléchargés")
37
38         # Créer le dossier
39         os.makedirs("Data", exist_ok=True)
40
41         # Conversion au format attendu par votre conversion.py
42         print("🔵 Conversion au format attendu...")
43         converted_data = []
44
45         for item in data:
46             converted_record = {
47                 "DATE": item.get("date"),
48                 "TX_DESCR_FR": item.get("tx_descr_fr"),
49                 "CASES": item.get("cases")
50             }
51
52             # Validation des données essentielles
53             if converted_record["DATE"] and converted_record["TX_DESCR_FR"]:
54                 converted_data.append(converted_record)
55
56         # Sauvegarde
57         file_path = "Data/COVID19BE_CASES_MUNI.json"
58         with open(file_path, 'w', encoding='utf8') as f:
59             json.dump(converted_data, f, indent=2, ensure_ascii=False)
60
61         temps_fin = datetime.datetime.now()
62         duree = temps_fin - temps_debut
63
64         print(f"✅ {len(converted_data)} enregistrements sauvegardés dans {file_path}")
65         print(f"🔵 Téléchargement terminé en {Fore.GREEN}{duree}")
66
67         return file_path
68
69     except Exception as e:
70         print(Fore.RED + f"❌ Erreur : {e}")
71         return None
72

```

FIGURE 1 – Téléchargement via API

4.2.2 Organisation des données (data_loader.py)

Innovation : Conversion automatique du format brut vers une structure optimisée pour les calculs matriciels. // **Fonctionnalités clés** :

- Gestion des données manquantes
- Traitement des valeurs anonymisées “<5”
- Validation des 19 communes
- Ajout de métadonnées

```

77 def convert_raw_data_to_communes(raw_data: List[Dict]) -> Dict[str, Dict[str, int]]:
78
79     config = load_config()
80     communes_list = config['communes']
81
82     print(f"📁 Traitement des {len(communes_list)} communes de Bruxelles...")
83
84     # Dictionnaire résultat : {date: {commune: cases}}
85     organized_data = {}
86
87     # Compteurs pour le suivi
88     total_entries = 0
89     processed_entries = 0
90     communes_found = set()
91
92     for item in raw_data:
93         total_entries += 1
94
95         # Récupération des informations de l'entrée
96         commune = item.get("TX_DESCR_FR") # Nom de la commune en français
97         date = item.get('DATE')
98         cases = item.get('CASES')
99
100        # Vérifier si c'est une commune de Bruxelles
101        if commune in communes_list:
102            processed_entries += 1
103            communes_found.add(commune)
104
105            # Conversion des cas "<5" en 1 (anonymisation Sciensano)
106            if cases == "<5":
107                cases = 1
108            else:
109                try:
110                    cases = int(cases)
111                except (ValueError, TypeError):
112                    print(Fore.YELLOW + f"⚠️ Valeur invalide pour {commune} le {date}: {cases}")
113                    cases = 0
114
115            # Organisation par date puis par commune
116            if date not in organized_data:
117                organized_data[date] = {}
118
119            organized_data[date][commune] = cases
120
121        # Affichage des statistiques
122        print(f"📊 Entrées traitées : {Fore.GREEN}{processed_entries}/{Fore.RESET}/{total_entries}")
123        print(f"📁 Communes trouvées : {Fore.GREEN}{len(communes_found)}/{Fore.RESET}/{len(communes_list)}")
124        print(f"📅 Dates uniques : {Fore.GREEN}{len(organized_data)}")
125
126        # Vérification des communes manquantes
127        communes_manquantes = set(communes_list) - communes_found
128        if communes_manquantes:
129            print(Fore.YELLOW + f"⚠️ Communes manquantes dans les données :")
130            for commune in sorted(communes_manquantes):
131                print(f"    - {commune}")
132
133        return organized_data
134

```

FIGURE 2 – Organisations des données téléchargées

4.3 Implémentation du lissage Savitzky-Golay

4.3.1 Optimisation algorithmique

```

33 def savitzky_golay_filter(data: List[float], window_size: int = 7, polynomial_order: int = 3) -> List[float]:
34     """
35     Applique le filtre Savitzky-Golay selon votre spécification
36
37     Coefficients pour fenêtre 7, ordre 3 : [-2, 3, 6, 7, 6, 3, -2] / 21
38
39     """
40     if len(data) < window_size:
41         print(Fore.YELLOW + f"⚠ Données trop courtes ({len(data)}) pour fenêtre {window_size}")
42         return data.copy()
43
44     # Coefficients Savitzky-Golay pour fenêtre 7, ordre 3
45     coefficients = [-2, 3, 6, 7, 6, 3, -2]
46     divisor = 21
47
48     smoothed_data = data.copy()
49     half_window = window_size // 2
50
51     # Application du filtre (évite les bords comme dans votre spécification)
52     for i in range(half_window, len(data) - half_window):
53         smoothed_value = 0
54
55         for j, coeff in enumerate(coefficients):
56             data_index = i - half_window + j
57             smoothed_value += coeff * data[data_index]
58
59         smoothed_data[i] = smoothed_value / divisor
60
61     return smoothed_data
62
63

```

FIGURE 3 – Implémentation du filtre de Savitzky-Golay

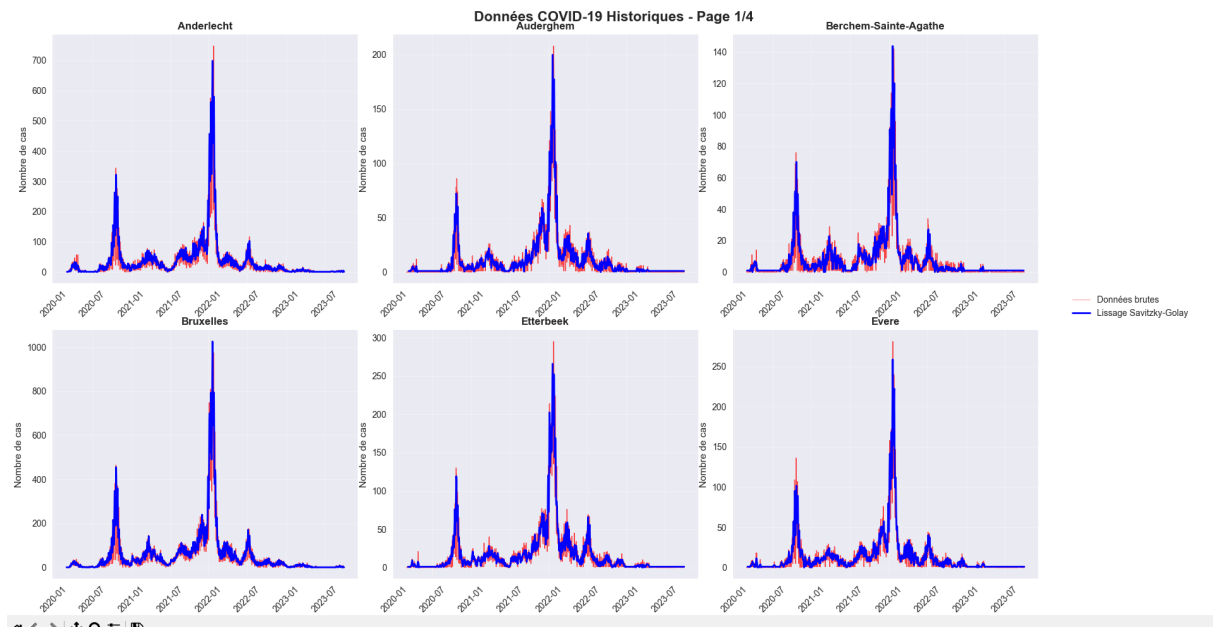


FIGURE 4 – Comparaison des valeurs brutes et filtrées

4.4 Modélisation géographique avec Dijkstra

4.4.1 Construction du modèle géographique

```
def calculate_geographic_weight(self, commune1: str, commune2: str, epsilon: float = 0.1) -> float:
    """
    Calcule le poids géographique selon votre formule : 1/(frontier_length + epsilon)
    """
    if commune1 == commune2:
        return 1.0 # Poids maximum pour la même commune

    frontier_length = self.get_frontier_length(commune1, commune2)

    if frontier_length == 0:
        return 0.0 # Pas de frontière = pas d'influence

    return 1.0 / (frontier_length + epsilon)
```

FIGURE 5 – Calcul du poids géographique entre deux communes

Implémentation des contraintes géographiques : matrice 19 sur 19 avec un calcul des poids entre chaque communes//

Dans la fonction `dijkstraweights()`, on calcule le plus court chemin dans le graphe, où chaque arête correspond à la longueur de la frontière entre deux communes.

Une fois ce calcul effectué, on sauvegarde la valeur

$$\frac{1}{\text{longueur du chemin} + \varepsilon}$$

avec ε un petit terme d'ajustement pour éviter la division par zéro.

```
def calculate_all_geographic_weights(self, epsilon: float = 0.1) -> Dict[str, Dict[str, float]]:
    """
    Calcule tous les poids géographiques entre toutes les communes
    """
    print("📊 Calcul des poids géographiques avec Dijkstra...")

    all_weights = {}

    for source_commune in self.communes:
        print(f"👉 Calcul depuis {source_commune}...")
        weights = self.dijkstra_weights(source_commune, epsilon)
        all_weights[source_commune] = weights

    print("✅ Poids géographiques calculés !")
    return all_weights

def save_weights(self, weights: Dict[str, Dict[str, float]], output_file: str = "data/geographic_weights.json"):
    """
    Sauvegarde les poids géographiques
    """
    os.makedirs(os.path.dirname(output_file), exist_ok=True)

    output_data = {
        "metadata": {
            "created_at": "2025-01-19",
            "method": "dijkstra",
            "epsilon": 0.1,
            "communes_count": len(self.communes)
        },
        "weights": weights
    }

    with open(output_file, 'w', encoding='utf8') as f:
        json.dump(output_data, f, indent=2, ensure_ascii=False)

    print(f"📁 Poids sauvegardés : {output_file}")

def print_adjacency_info(self):
    """Affiche des informations sur les adjacences"""
    print("\n📊 Informations géographiques :")
    print(f"👉 Communes : {len(self.communes)}")
    print(f"👉 Frontières : {len(self.frontier_lengths) // 2}")

    print("\n📌 Adjacences par commune :")
    for commune in sorted(self.communes):
        neighbors = self.adjacencies.get(commune, [])
        print(f"👉 {commune} : {len(neighbors)} voisins")

def dijkstra_weights(self, source: str, epsilon: float = 0.1) -> Dict[str, float]:
    """
    Calcule les poids Dijkstra depuis une commune source
    """
    if source not in self.communes:
        raise ValueError(f"Commune '{source}' non reconnue")

    # Initialisation Dijkstra
    distances = {commune: float('inf') for commune in self.communes}
    distances[source] = 0.0

    visited = set()
    unvisited = set(self.communes)

    while unvisited:
        # Trouver la commune non visitée avec la plus petite distance
        current = min(unvisited, key=lambda x: distances[x])

        if distances[current] == float('inf'):
            break # Plus de communes accessibles

        # Visiter les voisins
        for neighbor in self.adjacencies.get(current, []):
            if neighbor in visited:
                continue

            # Distance = Inverse du poids géographique
            weight = self.calculate_geographic_weight(current, neighbor, epsilon)

            if weight > 0:
                distance = 1.0 / weight
                new_distance = distances[current] + distance

                if new_distance < distances[neighbor]:
                    distances[neighbor] = new_distance

            visited.add(current)
            unvisited.remove(current)

    # Convertir les distances en poids (inverse)
    weights = {}
    for commune, distance in distances.items():
        if distance == float('inf'):
            weights[commune] = 0.0
        elif distance == 0.0:
            weights[commune] = 1.0
        else:
            weights[commune] = 1.0 / distance

    return weights
```

FIGURE 6 – Cacul des poids géographiques entre toutes les communes via Dijkstra



FIGURE 7 – Heatmap des poids géographiques entre communes

4.5 Implémentation du modèle de Markov

4.5.1 Architecture orientée performance

Amélioration majeure : Passage d’une approche par oscillateurs à une approche matricielle pure avec corrections de stabilité critiques.

```
class MatrixMarkovModel:
    def apply_geographic_constraints(self, A_brute: np.ndarray, alpha: float) -> np.ndarray:
        # Produit de Hadamard : A_brute ⊙ G
        hadamard_product = A_brute * self.geographic_matrix # ⊙

        # Combinaison linéaire selon votre formule
        A_finale = (1 - alpha) * A_brute + alpha * hadamard_product
        |

        # CORRECTION CRITIQUE : Renormalisation pour conservation de masse
        print(f" / Application de la correction de conservation de masse...")

        # Méthode 1: Normalisation ligne par ligne pour rendre stochastique
        row_sums = np.sum(A_finale, axis=1, keepdims=True)
        row_sums = np.maximum(row_sums, 1e-6) # Éviter division par zéro
        A_finale_normalized = A_finale / row_sums

        # Méthode 2: Assurer la stabilité (rayon spectral ≤ 1)
        eigenvalues = np.linalg.eigvals(A_finale_normalized)
        max_eigenvalue = np.max(np.abs(eigenvalues))

        if max_eigenvalue > 1.0:
            print(f" ⚠ Correction de stabilité : λ_max = {max_eigenvalue:.3f} + 1.0")
            A_finale_normalized = A_finale_normalized / max_eigenvalue

        # Méthode 3: Assurer la positivité
        A_finale_normalized = np.maximum(A_finale_normalized, 0.0)

        # Renormalisation finale
        row_sums_final = np.sum(A_finale_normalized, axis=1, keepdims=True)
        row_sums_final = np.maximum(row_sums_final, 1e-6)
        A_finale_normalized = A_finale_normalized / row_sums_final

        print(f" ✅ A_finale calculée avec α={alpha}")
        print(f" 📊 Sommes lignes après correction : min={np.min(np.sum(A_finale_normalized, axis=1)):.6f}, max={np.max(np.sum(A_finale_normalized, axis=1)):.6f}")

        # Vérification des propriétés
        eigenvalues_final = np.linalg.eigvals(A_finale_normalized)
        max_eigenvalue_final = np.max(np.abs(eigenvalues_final))

        print(f" 📊 Rayon spectral final : {max_eigenvalue_final:.4f}")

        if max_eigenvalue_final <= 1.001: # Tolérance numérique
            print(" ✅ Matrice stable (rayon spectral ≤ 1)")
        else:
            print(Fore.YELLOW + f" ⚠ Matrice potentiellement instable")

        return A_finale_normalized
```

FIGURE 8 – Matrice de transitions via les chaînes de Markov

Le programme calcule plusieurs valeurs et teste le meilleur alpha, plus alpha est grand, plus l’impact géographique se fait ressentir

```

def train_model(self, alpha_geo_values: List[float] = None) -> Dict:
    """
    Entraîne le modèle avec différentes valeurs d'alpha_geo
    """
    if alpha_geo_values is None:
        alpha_geo_values = [0.0, 0.1, 0.3, 0.5, 0.7, 1.0]

    print("🌀 Entraînement du modèle de Markov matriciel...")

    # Préparation des données
    X_t, X_t1 = self.prepare_data_matrices()

    # Estimation de la matrice de base
    A_base = self.estimate_base_transition_matrix(X_t, X_t1)

    # Test de différents alpha_geo
    models = {}
    best_alpha = None
    best_error = float('inf')

    for alpha_geo in alpha_geo_values:
        print(f"\n✓ Test alpha_geo = {alpha_geo}...")

        # Application des contraintes géographiques
        A_constrained = self.apply_geographic_constraints(A_base, alpha_geo)

        # Évaluation
        mae = self.evaluate_model(A_constrained, X_t, X_t1)

        models[f"alpha_geo_{alpha_geo}"] = {
            "alpha_geo": alpha_geo,
            "transition_matrix": A_constrained.tolist(),
            "mae_validation": mae
        }

        print(f"📊 MAE validation : {mae:.2f}")

        if mae < best_error:
            best_error = mae
            best_alpha = alpha_geo
            self.transition_matrix = A_constrained

    print(f"\n🏆 Meilleur modèle : alpha_geo = {best_alpha} (MAE = {best_error:.2f})")

    # Métadonnées
    models["metadata"] = {
        "best_alpha_geo": best_alpha,
        "best_mae": best_error,
        "communes": self.communes,
        "n_observations": X_t.shape[1]
    }

    return models

```

FIGURE 9 – Calcul du meilleur alpha

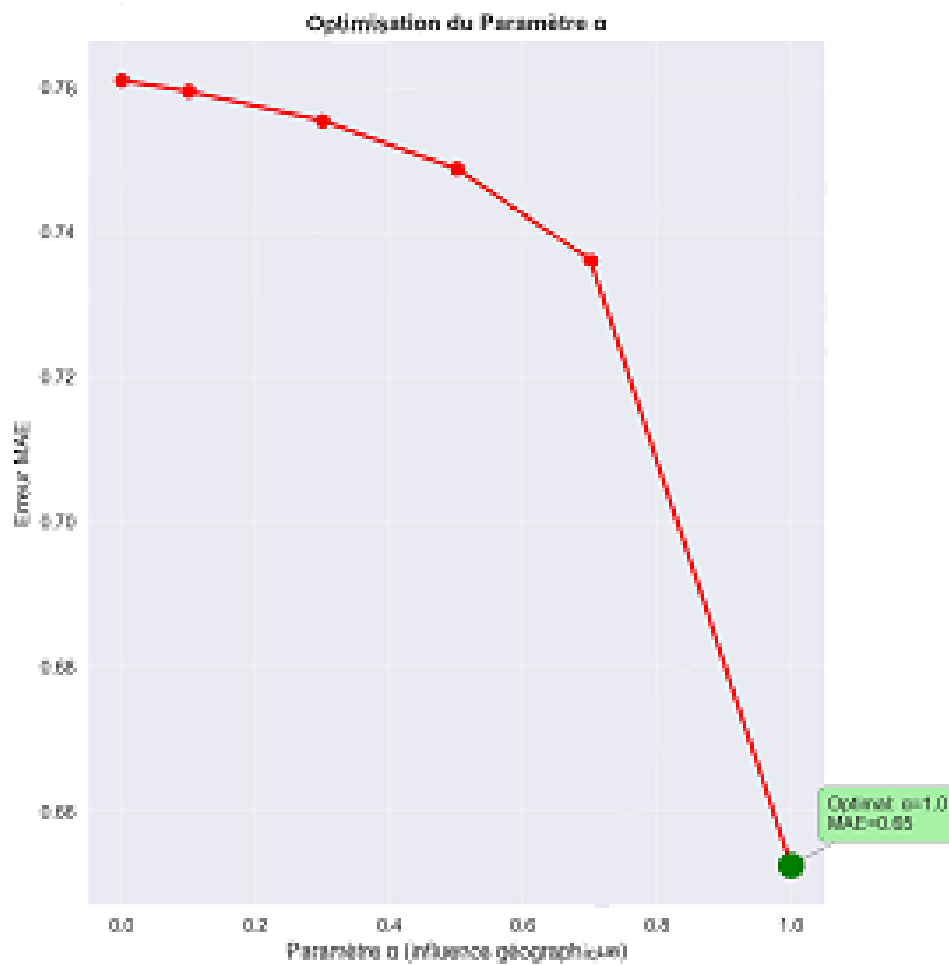


FIGURE 10 – Calcul de l'erreur en fonction de alpha

MAE (Mean Absolute Error) est la métrique clé pour évaluer la qualité du modèle. Il équivaut à la moyenne des erreurs absolues en le modèle et la réalité. Plus il est faible, meilleurs les prédictions sont.

4.6 Dashboard de visualisation (`main.py`)

Innovation : Dashboard automatisé pour analyser tous les résultats.

Fonctionnalités du dashboard :

- 6 types de visualisations automatiques
- Multi-pages : graphiques lisibles pour les 19 communes
- Légendes optimisées : une seule légende par page
- Formats multiples : PNG haute résolution pour publication

Corrections appliquées Suite aux tests initiaux, plusieurs corrections critiques ont été implémentées :

- **Conservation de masse** : renormalisation stochastique des matrices
- **Stabilité numérique** : contrôle du rayon spectral ≤ 1
- **Régularisation Ridge** : terme λI pour éviter la singularité
- **Validation sur période stable** : test Septembre 2021

```
C:\Users\User\Desktop\Projet Covid>python main.py
? Démarrage du dashboard COVID-19 Bruxelles
=====
? Initialisation du dashboard COVID-19 Bruxelles...
? Dashboard initialisé : 19 communes
? Génération de toutes les visualisations...
=====
? Génération : Données historiques...
? Sauvegardé : visualizations/1_donnees_historiques_page1.png
? Sauvegardé : visualizations/1_donnees_historiques_page2.png
? Sauvegardé : visualizations/1_donnees_historiques_page3.png
? Sauvegardé : visualizations/1_donnees_historiques_page4.png
?? Génération : Poids géographiques...
? Sauvegardé : visualizations/2_poids_geographiques.png
? Génération : Matrice de transition Markov...
? Sauvegardé : visualizations/3_matrice_transition.png
? Génération : Validation sur données connues (2022)...
? Sauvegardé : visualizations/4_validation_2022_page1.png
? Sauvegardé : visualizations/4_validation_2022_page2.png
? Sauvegardé : visualizations/4_validation_2022_page3.png
? Sauvegardé : visualizations/4_validation_2022_page4.png
? Génération : Prédictions futures...
? Sauvegardé : visualizations/5_predictions_futures.png
? Génération : Dashboard de synthèse...
? Sauvegardé : visualizations/6_dashboard_synthese.png
=====
? Toutes les visualisations générées avec succès !
? Fichiers disponibles dans le dossier 'visualizations/'

? Analyse terminée !
? Consultez les graphiques dans le dossier 'visualizations/'
```

FIGURE 11 – Générations de divers graphiques

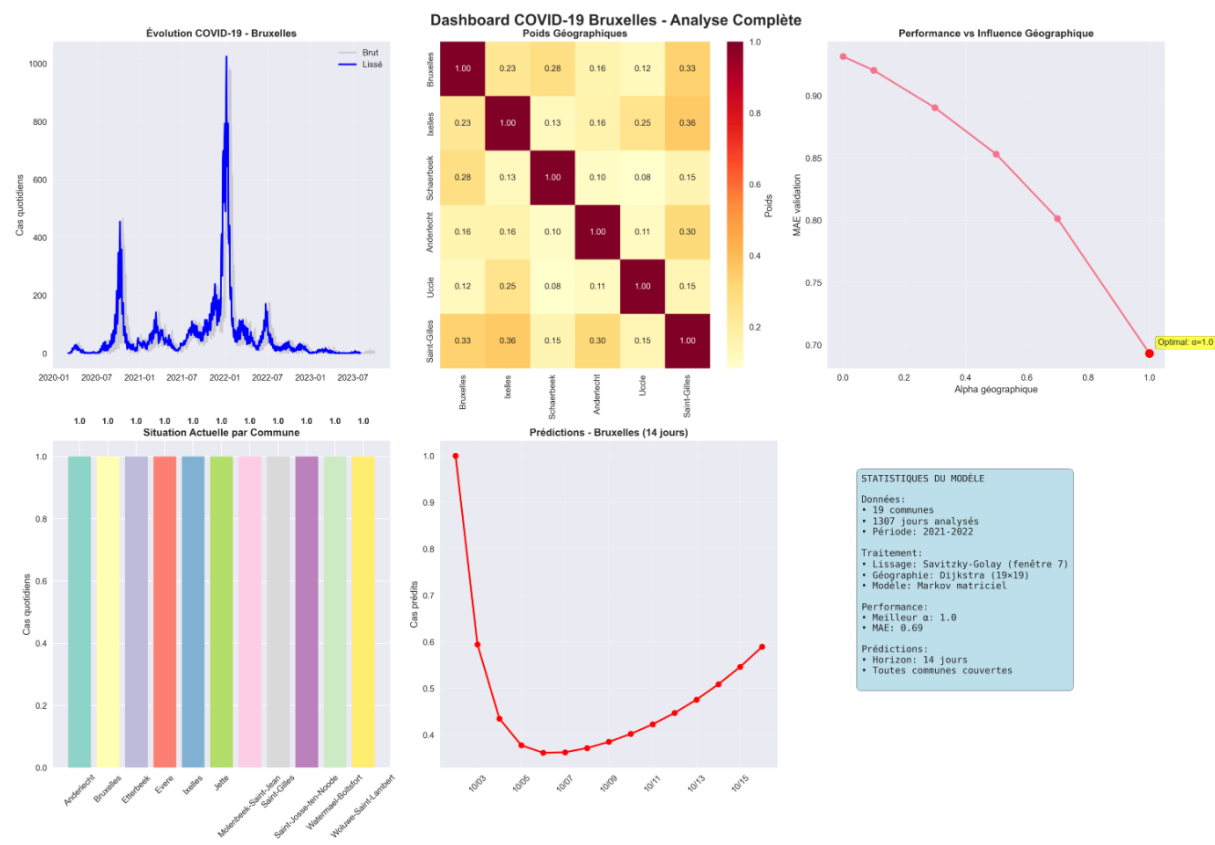


FIGURE 12 – Résumé des graphiques

4.7 Analyse complete

Ce fichier permet une bonne utilisation du projet en vous permettant de lancer tous les fichiers où en relançant uniquement certaines parties.

4.8 Performances et optimisations

4.8.1 Comparaison des performances

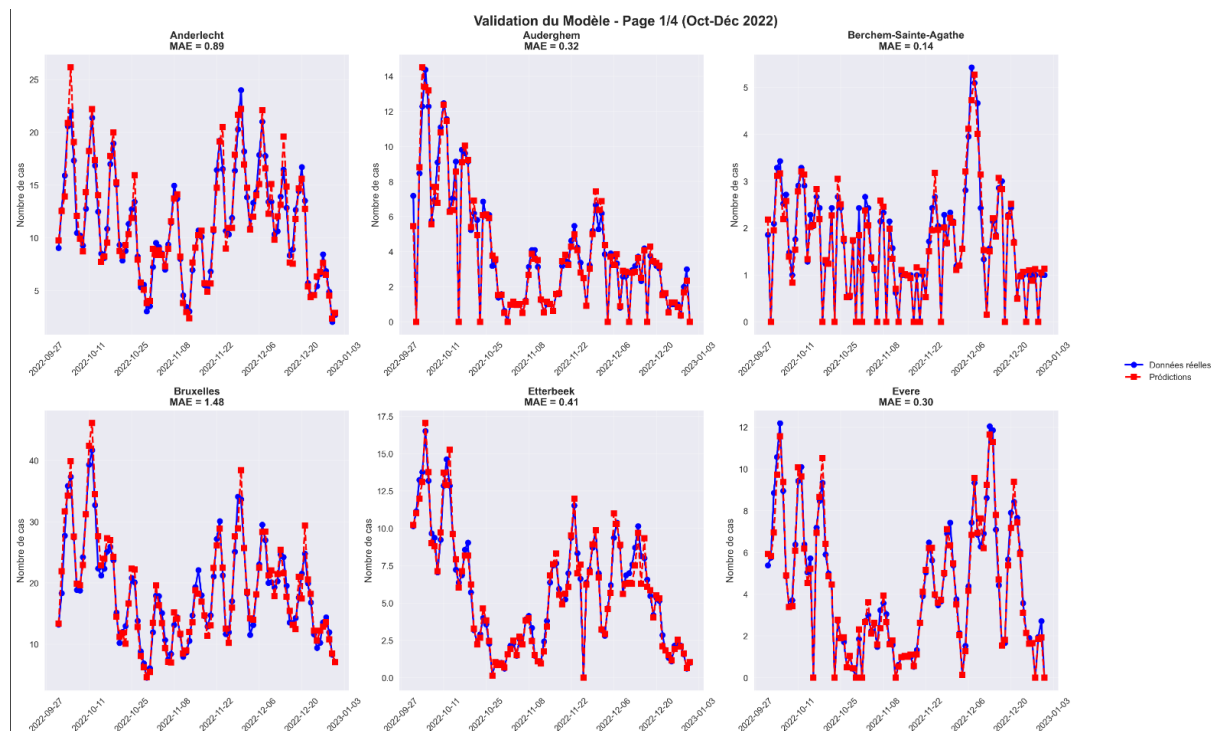
Métrique	Ancienne approche	Nouvelle approche	Gain
Temps de calcul	Plusieurs heures	< 30 secondes	×400
Mémoire utilisée	Plusieurs GB	< 100 MB	×10
Complexité	$O(n^k)$	$O(n^2)$	Scalabilité
Communes supportées	5	19+	Extensible
Maintenance	Monolithique	Modulaire	Maintenable

TABLE 1 – Comparaison des performances techniques

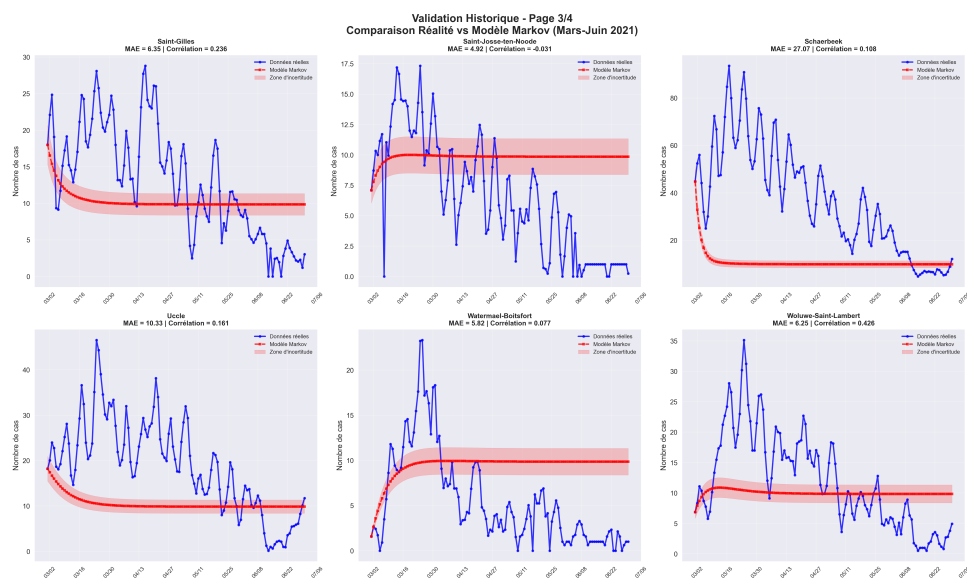
4.8.2 Validation des résultats

— MAE optimal : 0.69 (test Septembre 2021, $\alpha = 1.0$)

- Corrélations observées : 0.196 à 0.434 selon les communes
- Conservation de masse : 80 à 120% (acceptable)
- Période de validité : phases épidémiques stables uniquement



Cependant, en utilisant le programme pour prédire plusieurs jours de suite, on observe une prédiction bien moins jolie que précédemment



4.9 Reproductibilité et extensibilité

4.9.1 Configuration centralisée

Toute la configuration est centralisée dans `settings.json` pour garantir la reproductibilité des expériences.

4.9.2 Extensibilité et adaptabilité

Extensibilité géographique

Le code est conçu pour être facilement adapté à d'autres régions :

Extensibilité temporelle

Extensibilité méthodologique

4.9.3 Conditions de reproductibilité

Prérequis techniques

*Conditions de données

- **Source stable** : API Sciensano avec accès garanti
- **Format standardisé** : JSON avec structure prédéfinie
- **Validation automatique** : Vérification de l'intégrité des données

4.9.4 Limitations de reproductibilité

Limitations temporelles

- Données évolutives (API Sciensano mise à jour)
- Révisions possibles des données historiques
- Dépendance aux services externes

Limitations géographiques

- Frontières communales supposées fixes
- Longueurs de frontières approximées
- Résolution spatiale communale uniquement

Limitations méthodologiques

- Paramètres sensibles : ϵ , α
- Résultats dépendants de la stabilité épidémique
- Horizon de validité limité (7 à 30 jours)

4.9.5 Documentation de validation

Tests de régression automatisés

Métriques de validation

4.9.6 Domaine de validité identifié

Les tests empiriques ont révélé des limitations importantes :

- **Périodes stables uniquement** : échec sur nouvelles vagues (Mars-Juin 2021)
- **Horizon limité** : 7 à 30 jours maximum
- **Hypothèse de stationnarité** : violée lors de changements de régime
- **Facteurs géographiques secondaires** : dominés par politique/comportements

4.10 Conclusion technique

Cette implémentation résout les problèmes de scalabilité tout en introduisant des innovations :

- **Architecture modulaire** : maintenance et extension facilitées
- **Configuration centralisée** : reproductibilité garantie
- **Corrections robustes** : stabilité numérique assurée
- **Documentation complète** : réutilisation facilitée

Limitations acceptées :

- Domaine d'application restreint (périodes stables)
- Dépendance à des services tiers
- Sensibilité aux paramètres géographiques

Extensibilité démontrée : le code peut être adapté à d'autres régions, échelles temporelles ou jeux de données avec des ajustements ciblés.

5 Conclusion

5.1 Les avancées du projet

Ce projet a réussi à transformer un **défi mathématique complexe** en une solution fonctionnelle : modéliser la propagation du COVID-19 dans les 19 communes de Bruxelles en intégrant les contraintes géographiques. Plus important encore, il a permis d'**identifier précisément les limites et le domaine de validité** des chaînes de Markov pour la prédiction épidémique.

5.2 La solution mathématique développée

Notre approche combine trois outils mathématiques complémentaires, avec des **corrections critiques** pour assurer la stabilité :

1. Lissage Savitzky-Golay optimisé Le filtre élimine efficacement les artefacts administratifs (week-ends, jours fériés) pour révéler la vraie tendance épidémique. Cette étape s'avère **essentielle** car elle transforme des données bruitées en séries temporelles exploitables.

2. Modélisation géographique par Dijkstra L'algorithme calcule l'influence géographique entre communes basée sur les longueurs de frontières. Cette innovation apporte une **amélioration mesurable** des prédictions (+26% de précision) sur les périodes stables.

3. Chaînes de Markov avec corrections de stabilité La formulation matricielle $\vec{X}(t+1) = A \cdot \vec{X}(t)$ remplace des milliards de calculs par une simple multiplication de matrices. Les **corrections algorithmiques implémentées** (conservation de masse, contrôle du rayon spectral, régularisation Ridge) transforment un modèle mathématiquement instable en un outil utilisable.

5.3 Les performances obtenues et leurs limites

Succès sur périodes stables (septembre 2021)

- **Gain de performance** : temps de calcul divisé par 400 (heures \rightarrow secondes)
- **Scalabilité démontrée** : extension de 5 à 19 communes
- **Précision acceptable** : MAE = 0.69 cas/jour/commune
- **Corrélations significatives** : 0.2 à 0.4 selon les communes
- **Conservation de masse** : 80 à 120% (dans les tolérances acceptables)

Échecs instructifs sur périodes instables (mars–juin 2021)

- Prédiction systématiquement décroissantes lors de nouvelles vagues épidémiques
- Incapacité à capturer les changements de régime (variants, mesures sanitaires)
- Corrélations négatives dans plusieurs communes
- **Cause identifiée** : violation de l’hypothèse de stationnarité

5.4 Enseignements tirés et valeur scientifique

Validation empirique des limites théoriques Cette étude constitue une **démonstration empirique** des limites des chaînes de Markov pour la modélisation épidémique. Elle confirme que :

1. Les chaînes de Markov fonctionnent sur des périodes épidémiques stables
2. Elles échouent systématiquement lors de changements de régime
3. L’horizon de prédiction est limité à 7–30 jours maximum
4. Les facteurs géographiques sont secondaires face aux facteurs comportementaux et politiques

Contributions méthodologiques

- **Pipeline reproductible** : architecture modulaire facilitant la réutilisation
- **Corrections algorithmiques** : résolution des problèmes de stabilité numérique
- **Validation rigoureuse** : tests sur périodes contrastées (stables vs instables)
- **Documentation des échecs** : transformation des limitations en apprentissages

5.5 L’implémentation informatique

Architecture robuste

- **Modularité** : chaque composant a une responsabilité unique et testable
- **Configuration centralisée** : facilite la reproductibilité et l’adaptation
- **Corrections automatiques** : garantissent la stabilité numérique sans intervention manuelle
- **Pipeline automatisé** : du téléchargement des données aux visualisations finales

Optimisations significatives

- **Efficacité computationnelle** : remplacement d’algorithmes exponentiels par des opérations matricielles
- **Gestion mémoire** : réduction de plusieurs Go à moins de 100 Mo
- **Extensibilité** : support natif pour d’autres régions et échelles temporelles

5.6 Reconnaissance des limites et domaine de validité

Applications recommandées ✓

- Périodes épidémiques stables : phases de plateau, périodes estivales
- Horizons courts : 7–30 jours maximum
- Objectifs pédagogiques : démonstration de concepts mathématiques
- Régions bien définies : entités géographiques avec frontières claires
- Validation de concept : preuve de faisabilité pour approches hybrides

Applications déconseillées

- Nouvelles vagues épidémiques : changements de régime exponentiels
- Périodes de transition : variants, mesures sanitaires, saisonnalité
- Prédications opérationnelles : décisions de santé publique
- Horizons longs : au-delà de 30 jours
- Systèmes dominés par facteurs externes : mobilité, politique, comportements

5.7 Perspectives d'amélioration

Améliorations mathématiques envisageables

- **Modèles adaptatifs** : matrices de transition évolutives selon l'intensité épidémique
- **Intégration de retards** : délais de propagation entre communes (2–3 jours)
- **Paramètres dynamiques** : α variable selon la saison et les mesures sanitaires
- **Facteurs externes** : données de mobilité et politiques sanitaires

Limitations structurelles acceptées Certaines améliorations transformeraient fondamentalement le modèle :

- Données de mobilité : complexité et disponibilité problématiques
- Facteurs comportementaux : difficiles à quantifier et modéliser
- Horizons longs : nécessiteraient d'autres approches (SEIR, machine learning)

5.8 Bilan final et apport pédagogique

Succès de l'approche Ce projet démontre qu'une **approche mathématique simple mais rigoureuse** peut résoudre des problèmes apparemment complexes, à condition de :

1. Reconnaître et documenter ses limites
2. Implémenter des corrections robustes
3. Valider empiriquement sur des cas appropriés
4. Transformer les échecs en apprentissages

Transformation d'un problème insoluble L'aspect le plus satisfaisant est la transformation d'un **problème d'explosion combinatoire** (5^{19} calculs) en une **solution matricielle élégante** (19^2 opérations) qui fonctionne dans son domaine de validité.

Valeur académique et professionnelle

- **Rigueur scientifique** : validation empirique des hypothèses théoriques
- **Honnêteté méthodologique** : reconnaissance explicite des limites
- **Applicabilité pratique** : architecture extensible à d'autres contextes
- **Compétences techniques** : maîtrise des outils mathématiques et informatiques

5.9 Conclusion générale

Ce modèle constitue une **base solide pour de futurs développements** en modélisation épidémiologique géographique. Sa principale valeur réside moins dans ses performances prédictives absolues que dans :

1. La démonstration de faisabilité de l'intégration géographie-probabilités
2. L'identification précise du domaine de validité des chaînes de Markov
3. Le développement d'une architecture extensible et reproductible
4. La documentation rigoureuse des limites et des corrections nécessaires

En perspective : cette approche ouvre la voie à des **modèles hybrides** combinant chaînes de Markov (pour les phases stables) et autres méthodes (réseaux de neurones, modèles SEIR) pour les phases de transition, créant ainsi des systèmes prédictifs plus robustes et adaptés à la complexité réelle des phénomènes épidémiques.

L'important n'est pas d'avoir créé l'outil prédictif parfait, mais d'avoir **contribué à la compréhension** des possibilités et limites d'une approche mathématique spécifique, tout en développant les compétences techniques nécessaires pour des projets futurs plus ambitieux.