



**HAL**  
open science

# Compared performance of Covid19 reproduction number estimators based on realistic synthetic data

Juliana Du, Barbara Pascal, Patrice Abry

## ► To cite this version:

Juliana Du, Barbara Pascal, Patrice Abry. Compared performance of Covid19 reproduction number estimators based on realistic synthetic data. Colloque Francophone de Traitement du Signal et des Images (GRETSI), GRETSI, Aug 2023, Grenoble, France. hal-04032614v4

**HAL Id: hal-04032614**

<https://hal.science/hal-04032614v4>

Submitted on 17 Apr 2024

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution 4.0 International License

# Performances comparées d’estimateurs du coefficient de reproduction de la Covid19 à l’aide de données synthétiques réalistes<sup>†</sup>

Juliana DU<sup>1</sup> Barbara PASCAL<sup>2</sup> Patrice ABRY<sup>1</sup>

<sup>1</sup>ENSL, CNRS, Laboratoire de physique, F-69342 Lyon, France

<sup>2</sup>Nantes Université, École Centrale Nantes, CNRS, LS2N, UMR 6004, F-44000 Nantes, France

**Résumé** – Estimer précisément le coefficient de reproduction est un enjeu majeur pour surveiller et contrôler une épidémie. Afin d’évaluer et de comparer plusieurs estimateurs proposés durant la pandémie de Covid19, une procédure de génération de données épidémiologiques synthétiques réalistes est élaborée, reflétant à la fois la propagation épidémique, suivant le modèle épidémiologique proposé par Cori et col., et un processus de collecte de nombre de cas dégradé (erreur de report, données manquantes). Des simulations de Monte Carlo sur un large jeu de données synthétiques aux dynamiques variées fournissent une évaluation quantitative des performances d’estimation du coefficient de reproduction.

**Abstract** – Accurate estimation of the reproduction number is a major challenge for the surveillance and control of epidemics. To evaluate and compare several estimators proposed during the Covid19 pandemic, a procedure for generating realistic synthetic epidemiological data is designed, reflecting both the epidemic spread, following the epidemiological model proposed by Cori et al., and the degraded collection of infection counts (erroneous figures, missing data). Monte Carlo simulations over a large set of synthetic data with different dynamics yield quantitative assessment of reproduction number estimation performance.

## 1 Introduction

**Contexte.** Durant la pandémie de Covid19 des efforts conséquents ont été dévolus à la construction d’outils précis pour le suivi de la situation sanitaire au jour le jour. L’indicateur de l’intensité épidémique le plus populaire est le *coefficient de reproduction*  $R_t$ , défini comme le nombre moyen de personnes contaminées par un individu infectieux, quantifiant la transmissibilité du virus au jour  $t$  [5, 7]. Plusieurs estimateurs de  $R_t$  sont documentés [9, 5, 1, 8], dont certains ont été développés pour répondre à la crise engendrée par la pandémie de Covid19. Néanmoins, faute de données épidémiologiques annotées, par exemple *synthétiques*, aucun de ces estimateurs n’a reçu de validation systématique de ses performances et aucune comparaison quantitative n’est possible. Or, le contrôle d’une pandémie, aux enjeux sanitaires, sociaux et économiques majeurs, requiert des outils validés scientifiquement.

**État-de-l’art.** Une revue de la littérature récente montre que peu de travaux se sont intéressés à la génération de données synthétiques pour la Covid19. Les modèles compartimentaux [7, 3], décrivant directement la dynamique épidémique, sont bien adaptés à la génération de données, et ont été mis à profit, par exemple pour la prédiction [2]. Cependant, la qualité des données synthétiques produites dépend de la précision dans l’estimation des paramètres du modèle, non seulement coûteuse en ressources de calcul, mais également fortement dégradée par la faible qualité des données relatives à la Covid19 rapportées par les autorités sanitaires (cf. exemple Fig. 1).

Alternativement, dans le modèle statistique de Cori [5], le nombre de nouvelles infections  $Z_t$  au jour  $t$  dépend du coefficient de reproduction  $R_t$  et des nombres d’infections passés  $\{Z_1, Z_2, \dots, Z_{t-1}\}$ , pondérés par la fonction d’intervalle de série  $\phi$  modélisant le délai aléatoire entre le début des symp-

tômes dans l’infection primaire et l’infection secondaire. Le schéma bayésien privilégié par [5] pour estimer  $R_t$  à partir des comptes  $Z_t$ , approprié pour des données consolidées *a posteriori*, est sévèrement mis en difficulté par les données quotidiennes liées à la Covid19 qui sont de qualité limitée.

**Objectifs et contributions.** Le but est de comparer les estimateurs du coefficient de reproduction introduits par [1, 8] sur la base de performances d’estimation quantitatives. Une première contribution consiste en la construction d’une stratégie de génération de données synthétiques réalistes, fournissant des comptes d’infection accompagnés de la vérité terrain sur le coefficient de reproduction. Une seconde contribution est la comparaison quantitative systématique des performances d’estimation et du temps de calcul, menée *via* des simulations numériques de Monte Carlo intensives.

## 2 Estimation de $R_t$

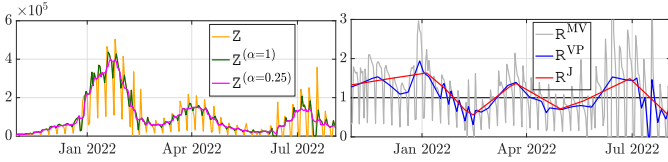
**Modèle épidémiologique.** Le modèle de Cori [5] fait l’hypothèse que  $Z_t$  suit une loi de Poisson dont l’intensité  $p_t$  varie dans le temps, de sorte que la log-vraisemblance s’écrit

$$\ln \mathbb{P}(Z_t | Z_1, \dots, Z_{t-1}; p_t) = Z_t \ln(p_t) - p_t - \ln(Z_t!). \quad (1)$$

**Données relatives à la Covid19.** Depuis le début de la pandémie, l’Université Johns Hopkins (JHU) a collecté et mis en ligne le nombre quotidien d’infections par la Covid19 dans plus de 200 pays<sup>1</sup>. À titre d’exemple, la Fig. 1a présente les comptes de nouveaux cas  $Z_t$  en France sur une durée de dix mois, qui présentent des fluctuations rapides, parfois très importantes (cf. courbe en orange, Fig. 1a), sans lien avec la dynamique épidémique et reflétant avant tout l’impact des jours chômés sur les stratégies de test et de report des cas. Une

<sup>†</sup>J. Du est financée par le projet 80PRIME-2021 CNRS « CoMoDécartes ».

<sup>1</sup><https://coronavirus.jhu.edu/>



(a) Comptes bruts  $\mathbf{Z}$  (en orange) et pré-traités par médian glissant  $\mathbf{Z}^{(\alpha)}$ , avec  $\alpha$  grand (resp. petit) en vert (resp. en rose) (b) Coefficient de reproduction estimé par Max. Vrais. (MV), Vrais. Pén. (VP) et estimation jointe (J).

FIGURE 1 : Infections par la Covid19 en France et coefficient de reproduction entre le 1<sup>er</sup> novembre 2021 et le 3 août 2022.

estimation précise du coefficient de reproduction sous-jacent nécessite de gérer la mauvaise qualité de ces données réelles.

**Maximum de Vraisemblance.** Originellement, l'intensité  $p_t$  d'une épidémie est le produit du coefficient de reproduction et de la *contagiosité effective* dans la population [5] :

$$p_t^{(1)} = R_t \Phi_t^Z, \quad \text{où } \Phi_t^Z = \sum_{s=1}^{\tau_\phi} \phi(s) Z_{t-s}. \quad (2)$$

La fonction d'intervalle de série  $\phi$  encode les caractéristiques de l'épidémie et est modélisée, pour le virus de la Covid19, par une distribution Gamma de moyenne (resp. écart-type) 6,6 (resp. 3,5) jours tronquée à  $\tau_\phi = 25$  jours. L'estimateur de Maximum de Vraisemblance (MV) est obtenu par maximisation directe de (1) :

$$\widehat{R}_t^{\text{MV}} = Z_t / \Phi_t^Z, \quad (3)$$

et tracé en gris sur la Fig. 1b, dans l'exemple susmentionné des données françaises. La grande variabilité des décomptes d'infections induit un comportement très irrégulier de  $\widehat{R}_t^{\text{MV}}$ , non réaliste du point de vue épidémiologique.

**Débruitage puis estimation régularisée.** Les fluctuations erratiques des comptes reportés proviennent majoritairement du processus de collecte fortement impacté par les jours non-travaillés entraînant des erreurs de report, des données aberrantes, des décomptes nuls le week-end et des pseudo-périodicités, et seulement dans une moindre mesure du caractère aléatoire de la propagation de la pandémie.

i) *Pré-traitement des comptes erronés.* Une stratégie naturelle est d'appliquer aux données corrompues une étape de débruitage en amont de l'estimation de  $R_t$ . Pour retirer les décomptes anormalement bas/élevés tout en préservant autant que possible la dynamique sous-jacente, [1] applique un médian glissant aux données brutes de nouvelles infections. Chaque jour  $t$ , le médian de  $Z_s$  sur une fenêtre  $W_t$  de deux semaines centrée en  $t$  est comparé au compte observé  $Z_t$  : pour un seuil fixé  $\alpha > 0$ ,  $Z_t$  est remplacé par le médian local si

$$|Z_t - \text{med}\{Z_s, s \in W_t\}| \geq \alpha \times \text{mad}\{Z_s, s \in W_t\} \quad (4)$$

où med (resp. mad) désigne le médian (resp. la déviation au médian en valeur absolue) d'une série temporelle ; sinon  $Z_t$  est laissé inchangé. Les données françaises pré-traitées via la procédure de médian glissant, notées  $\mathbf{Z}^{(\alpha)}$ , sont tracées en Fig. 1a : lorsque  $\alpha$  est grand seules les valeurs extrêmement anormales sont corrigées, tandis qu'un petit  $\alpha$  induit un débruitage plus fort, au prix d'une dégradation de la dynamique sous-jacente, (voir discussion en Sec. 4).

ii) *Estimation régularisée.* Une fois les données corrigées de leurs valeurs aberrantes par le médian glissant, [1] estime  $R_t$  via la minimisation d'une Vraisemblance Pénalisée (VP) favorisant un comportement linéaire par morceaux qui permet

de saisir la tendance globale de la dynamique épidémique :

$$\widehat{\mathbf{R}}^{\text{VP}} = \underset{\mathbf{R} \in \mathbb{R}_+^T}{\text{argmin}} \sum_{t=1}^T d_{\text{KL}}(Z_t | p_t^{(1)}) + \mu_R \|\mathbf{D}_2 \mathbf{R}\|_1 \quad (5)$$

où la divergence de Kullback-Leibler, définie comme

$$d_{\text{KL}}(Z|p) = \begin{cases} Z \ln(Z/p) + p - Z & \text{si } Z > 0 \text{ et } p > 0 \\ p & \text{si } Z = 0 \text{ et } p \geq 0 \\ \infty & \text{sinon,} \end{cases}$$

coïncide avec l'opposé de la log-vraisemblance de Poisson. La parcimonie du laplacien discret de l'estimée, c.-à-d., le caractère linéaire par morceaux de  $\widehat{\mathbf{R}}^{\text{VP}}$ , est favorisé par  $\|\mathbf{D}_2 \mathbf{R}\|_1 = \sum_{t=3}^T |R_t - 2R_{t-1} + R_{t-2}|$ . Le paramètre  $\mu_R > 0$  contrôle le niveau de régularisation (voir Sec. 4).

**Modèle de Cori étendu et procédure jointe.** Bien que fournissant des estimées régularisées satisfaisantes sur certaines données peu corrompues, cette méthode en deux étapes est plus difficile à mettre en œuvre lorsque la qualité des données diminue (cf. courbe en bleu en Fig. 1b). Pour palier cette limitation et gérer efficacement les comptes erronés et les données manquantes dans les données récentes, un modèle étendu de l'intensité de l'épidémie a été proposé dans [8] :

$$p_t^{(2)} = R_t \Phi_t^Z + O_t, \quad (6)$$

où la variable  $O_t$  modélise l'écart entre les comptes reportés et les nombres réels de nouvelles infection et est estimée conjointement avec  $R_t$  via un algorithme joint résolvant

$$\widehat{\mathbf{R}}, \widehat{\mathbf{O}}^J = \underset{\mathbf{R} \in \mathbb{R}_+^T, \mathbf{O} \in \mathbb{R}^T}{\text{argmin}} \sum_{t=1}^T d_{\text{KL}}(Z_t | p_t^{(2)}) + \lambda_R \|\mathbf{D}_2\|_1 + \lambda_O \|\mathbf{O}\|_1. \quad (7)$$

Le terme  $\|\mathbf{O}\|_1$  promeut la parcimonie de la variable d'erreur  $O_t$ , dont le niveau est contrôlé par un paramètre  $\lambda_O > 0$ . L'estimée de  $R_t$  via la procédure jointe apparaît significativement plus robuste à la qualité limitée des données (voir courbe en rouge, Fig. 1b) ce qui sera quantifié en Sec. 4.

### 3 Données synthétiques

L'évaluation et la comparaison quantitatives des stratégies d'estimation de  $R_t$  décrites en Sec. 2 requièrent des comptes  $\mathbf{Z}$  accompagnés des coefficients de reproduction  $\mathbf{R}$  constituant la vérité terrain. L'Alg. 1 est construit pour générer efficacement des nombres de cas reportés synthétiques sous le modèle (6) à partir de vérités terrain  $\mathbf{R}, \mathbf{O}$ . La production de données synthétiques *réalistes* repose sur la capacité à fabriquer des vérités terrains appropriées pour  $R_t$  et  $O_t$  ; pour cela, une procédure s'appuyant sur les données réelles disponibles est proposée.

**Échantillonnage itératif sous le modèle de Cori étendu.** Soit  $Z_1, \dots, Z_{t-1}$  les comptes passés, échantillonner  $Z_t$  sous le modèle de Cori étendu (1) revient à calculer l'intensité de Poisson instantanée  $p_t^{(2)}$  à partir des nombres de cas passés selon (6) (cf. l. 2 de l'Alg. 1), puis de tirer une variable aléatoire sous la loi de Poisson d'intensité  $p_t^{(2)}$  (cf. l. 3 de l'Alg. 1). La contagiosité effective  $\Phi_t^Z$  intervenant dans le calcul de  $p_t^{(2)}$  dépend des  $\tau_\phi$  précédents comptes. Pour  $t < \tau_\phi$  les effets de bords dans le calcul de  $\Phi_t^Z$  sont gérés en imposant que la somme des poids appliqués vaille un via une normalisation par  $n_t^\phi$  (cf. l. 1 de l'Alg. 1).

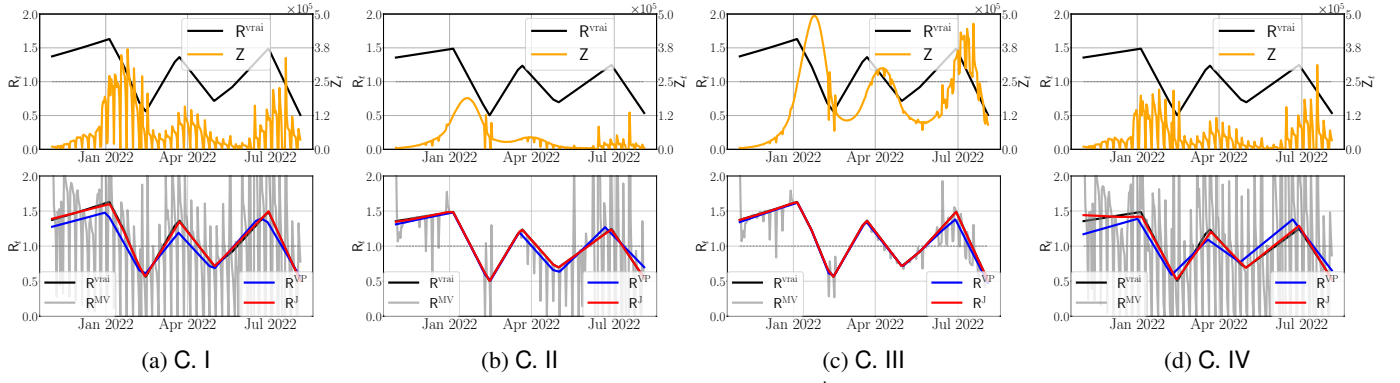


FIGURE 2 : Haut : données synthétiques  $\mathbf{Z}$  (en orange) et vérité terrain  $\mathbf{R}^{\text{vrai}}$  (en noir). Bas : estimations de  $\mathbf{R}$  par maximum de vraisemblance (MV, en gris), médian glissant puis vraisemblance pénalisée (VP, en bleu) et procédure jointe (J, en rouge) sur des données synthétiques de vérité terrain  $\mathbf{R}^{\text{vrai}}$  (en noir).

---

### Algorithme 1 : Génération de comptes synthétiques.

---

**Données :**  $\mathbf{R}^{\text{vrai}} \in \mathbb{R}_+^T$ ,  $\mathbf{O}^{\text{vrai}} \in \mathbb{R}^T$ ,  
 $Z_0 \in \mathbb{N}$  # nombre d'infections initial

**Résultat :**  $\mathbf{Z} = (Z_1, Z_2, \dots, Z_T) \in \mathbb{N}^T$

```

1 pour  $t = 1, 2, \dots, T$  faire
2    $n_t^\phi = \sum_{s=1}^{\min(t, \tau_\phi)} \phi(s)$  # effets de bords
3    $p_t = \mathbf{R}_t^{\text{vrai}} \times \sum_{s=1}^{\min(t, \tau_\phi)} \frac{\phi(s)}{n_t^\phi} Z_{t-s} + \mathbf{O}_t^{\text{vrai}}$  # intensité
4    $Z_t \sim \text{Pois}(\max(p_t, 0))$  # nombre d'infections
5 fin

```

---

**Vérités terrains réalistes pour  $\mathbf{R}$  et  $\mathbf{O}$ .** Obtenir des comptes de nouvelles infections synthétiques *réalistes*, nécessite que  $\mathbf{R}^{\text{vrai}}$  et  $\mathbf{O}^{\text{vrai}}$  imitent fidèlement à la fois la dynamique de propagation sous-jacente, et la dégradation des données par le processus de collecte. Or, ni le coefficient de reproduction  $R_t$  évoluant au cours du temps, ni la variable de corruption non stationnaire  $O_t$ , qui dépendent fortement de l'épidémie considérée, ne sont accessibles à une modélisation exhaustive. Afin de contourner cet écueil, des  $R_t, O_t$  réalistes sont construits à partir d'échantillons de données réelles extraits de la base JHU. Pour un pays et une période de temps donnés, les comptes d'infections reportés sont analysés au moyen de la stratégie jointe (7) pour des paramètres de régularisation fixés ; les estimées obtenues  $\hat{\mathbf{R}}$  et  $\hat{\mathbf{O}}$  fournissent alors des vérités terrains réalistes. En considérant différentes configurations, c.-à-d. différents réglages des paramètres de régularisation ( $\lambda_R, \lambda_O$ ), cette méthode permet de construire des vérités terrains réalistes, ensuite combinées pour générer, *via* l'Alg. 1, des données synthétiques diversifiées (cf. Fig. 2a à 2d, haut).

## 4 Performances d'estimation

**Optimisation convexe non lisse.** Les estimées par vraisemblance pénalisée (5) et par la procédure jointe (7) sont calculées par les schémas primaux-duaux de Chambolle-Pock [4] développés dans [1, 8], recourant à des opérateurs proximaux pour la minimisation de fonctionnelles non-lisses. Les pas de descente sont choisis de la même manière que dans [1, 8]. Cependant, pour renforcer la robustesse de l'estimation vis-à-vis de dynamiques épidémiques très différentes et de différents niveaux de régularisation, l'algorithme est stoppé lorsque le maximum, sur les  $T$  composantes, des incréments relatifs du

coefficient de reproduction entre deux itérations successives  $\mathbf{R}^{[k]}$  et  $\mathbf{R}^{[k+1]}$ , lissé sur 500 itérations, est inférieur à  $10^{-6}$ , c.-à-d. quand

$$\max_{k-500 < k' \leq k} \max_{t=1, \dots, T} \frac{|\mathbf{R}_t^{[k'+1]} - \mathbf{R}_t^{[k']}|}{\mathbf{R}_t^{[k']}} \leq 10^{-6} \quad (8)$$

ou bien après un maximum de  $7 \cdot 10^5$  itérations.

**Données synthétiques.** Les performances des estimateurs décrits en Sec. 2 sont évaluées sur des données synthétiques produites *via* la procédure exposée en Sec. 3 : les comptes de nouvelles infections reportés en France entre le 1<sup>er</sup> novembre 2021 et le 3 août 2022 (cf. Fig. 1a) sont analysés au moyen de (7), avec deux réglages de paramètres de régularisation très différents<sup>2</sup>, fournissant une collection de  $\mathbf{R}^{\text{vrai}}, \mathbf{O}^{\text{vrai}}$ . Ces vérités terrains sont ensuite combinées en quatre configurations :  
 C. I : beaucoup de changements de pente de  $R_t$ ,  
 beaucoup de valeurs non nulles de  $O_t$  (cf. Fig. 2a haut) ;  
 C. II : peu de changements de pentes de  $R_t$ ,  
 peu de valeurs non nulles de  $O_t$  (cf. Fig. 2b haut) ;  
 C. III : beaucoup de changements de pentes de  $R_t$ ,  
 peu de valeurs non nulles de  $O_t$  (cf. Fig. 2c haut) ;  
 C. IV : peu de changements de pentes de  $R_t$ ,  
 beaucoup de valeurs non nulles de  $O_t$  (cf. Fig. 2d haut).

Ces vérités terrains réalistes et diversifiées, avec des dynamiques de  $R_t$  et  $O_t$  similaires pour C. I et II et « croisées » pour C. III et IV, sont ensuite mises à profit pour générer des comptes de nouvelles infections synthétiques  $\mathbf{Z}$  au moyen de l'Alg. 1. Pour chaque couple  $(\mathbf{R}^{\text{vrai}}, \mathbf{O}^{\text{vrai}})$ , c.-à-d. pour chaque C I à IV, 20 réalisations de  $\mathbf{Z}$  sont générées.

**Évaluation des performances.** Pour mesurer la qualité des estimées du coefficient de reproduction, deux critères sont utilisés. Tout d'abord, le rapport signal à bruit (SNR), indicateur standard mesuré en décibels (dB), défini comme

$$\text{SNR} := 10 \times \log_{10} \left( \frac{\|\mathbf{R}^{\text{vrai}}\|_2^2}{\|\hat{\mathbf{R}} - \mathbf{R}^{\text{vrai}}\|_2^2} \right) \quad (9)$$

quantifie l'adéquation entre l'estimée  $\hat{\mathbf{R}}$  et la vérité terrain  $\mathbf{R}^{\text{vrai}}$ . La qualité des estimées fournie par les méthodes en deux vs. une étape(s) dépend fortement du réglage des hyperparamètres, c.-à-d. du choix du seuil  $\alpha$  dans (4), du paramètre  $\mu_R$  dans (5), et de  $\lambda_R, \lambda_O$  dans (7). Pour chaque estimateur, chaque configuration et chaque réalisation, le SNR est maximisé sur une grille de  $20 \times 20$  paramètres espacés logarithmiquement.

<sup>2</sup> $(\lambda_R, \lambda_O) = (3,7 \cdot 10^5; 0,03)$  et  $(\lambda_R, \lambda_O) = (5,3 \cdot 10^6; 0,75)$

	Max. Vrais. (MV)	Vrais. Pén. (VP)	Estim. Jointe (J)
SNR (dB)			
C. I	0,90 ± 0,03	20,65 ± 0,06	<b>35,68 ± 0,39</b>
C. II	-2,31 ± 0,06	26,55 ± 0,15	<b>48,56 ± 0,90</b>
C. III	19,14 ± 0,08	29,79 ± 0,09	<b>56,34 ± 0,77</b>
C. IV	-2,32 ± 0,01	18,79 ± 0,02	<b>31,06 ± 0,20</b>
Indice de Jaccard sur le laplacien (%)			
C. I	0,79 ± 0,00	57,63 ± 1,33	<b>84,28 ± 1,48</b>
C. II	0,57 ± 0,00	47,71 ± 1,10	<b>96,46 ± 0,77</b>
C. III	2,28 ± 0,01	73,76 ± 0,73	<b>98,03 ± 0,90</b>
C. IV	0,67 ± 0,00	13,58 ± 0,65	<b>76,79 ± 0,57</b>

TABLE 1 : SNR et indice de Jaccard pour trois estimateurs et quatre configurations, moyennées sur  $N = 20$  réalisations.

Pour chaque estimateur et chaque configuration, le SNR maximal moyen et son intervalle de confiance gaussien à 95% sont reportés dans la Tab. 1, l. 3 à 6. Puis, pour la meilleure estimation au sens du SNR, l'indice de Jaccard<sup>3</sup> [6] entre le laplacien discret binarisé de  $\hat{\mathbf{R}}$  et celui de  $\mathbf{R}^{\text{vrai}}$ , avec une fenêtre  $\mathbf{g}$  gaussienne d'écart-type 1 jour, reporté dans la Tab. 1, l. 8 à 11, permet d'évaluer si les points de rupture dans le comportement linéaire de  $R_t$ , correspondant aux changements brutaux de dynamique épidémique, sont détectés correctement.

**Comparaison entre les estimateurs.** Pour chaque configuration C. I à IV, les Fig. 2a à 2d, en bas, présentent les estimées de  $\mathbf{R}$  obtenues par maximum de vraisemblance (en gris), médian glissant puis vraisemblance pénalisée (en bleu), et procédure jointe (en rouge) sur une réalisation de données synthétiques. Quelle que soit la configuration C. I à IV, l'estimée  $\mathbf{R}^{\text{J}}$  est systématiquement la plus fidèle à  $\mathbf{R}^{\text{vrai}}$ . L'étude quantitative du SNR sur 20 réalisations de chaque configuration, l. 3 à 6 du Tab. 1, montre la qualité nettement supérieure des estimées régularisées, col. 3 et 4, par rapport au maximum de vraisemblance, col. 2. En outre, la procédure jointe, col. 4 du Tab. 1, atteint des SNR significativement plus élevés, de plus de 10 dB, que la vraisemblance pénalisée, col. 3 du Tab. 1.

**Robustesse de la procédure d'estimation.** La diversité des données synthétiques C. I à IV permet de sonder la robustesse des estimateurs de  $R_t$  vis-à-vis de différentes combinaisons de dynamiques épidémiques et de niveaux de corruption des données. Pour toutes les configurations, l'indice de Jaccard entre le laplacien de l'estimée jointe et celui de la vérité terrain, Tab. 1 col. 8 à 11, varie entre 76% et 97%, attestant à la fois d'une détection efficace et d'une localisation correcte des changements de dynamique épidémique, plaidant pour l'utilisation de la procédure jointe [8] pour surveiller une épidémie, préférentiellement aux estimateurs de [5, 1].

**Temps de calcul.** Les meilleures performances d'estimation de la procédure jointe se font au prix d'un coût de calcul plus élevé (cf. Fig. 3). Pour des paramètres de régularisation équivalents pour les deux méthodes, le critère de convergence (8) atteint la précision requise ( $10^{-6}$ ) en 20 fois plus d'itérations pour la procédure jointe (J, en rouge) que pour la méthode par vraisemblance pénalisée (VP, en bleu), et ce de manière robuste, comme l'attestent les barres d'erreur très étroites calculées sur  $N = 20$  réalisations (cf. encart de droite, Fig. 3).

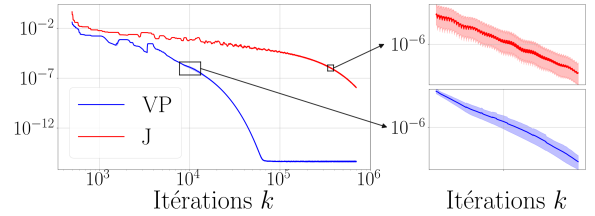


FIGURE 3 : Incréments relatifs (8) lors de l'estimation de  $R_t$  via VP (bleu) ou J (rouge) sur  $N = 20$  réalisations de C. I.

## 5 Conclusion

Une procédure de génération de comptes d'infections synthétiques *réalistes*, reproduisant à la fois la dynamique épidémique et les erreurs dans le report des cas, a été proposée et mise à profit pour construire des données synthétiques diversifiées sur lesquelles plusieurs estimateurs du coefficient de reproduction ont été évalués et comparés *via* des simulations de Monte Carlo intensives, démontrant la supériorité d'une approche *jointe*, corrigeant les comptes et estimant le coefficient de reproduction simultanément. La génération de données synthétiques constitue une étape cruciale pour évaluer systématiquement les estimateurs d'indicateurs épidémiologiques. En outre, cette étude ouvre la voie à l'élaboration de méthodes automatiques et pilotées par les données pour le réglage des paramètres de régularisation de la procédure jointe.

## Références

- [1] P. ABRY, N. PUSTELNIK, S. ROUX, P. JENSEN, P. FLANDRIN, R. GRIBONVAL, C.-G. LUCAS, É. GUICHARD, P. BORGAT et N. GARNIER : Spatial and temporal regularization to estimate COVID-19 reproduction number  $R(t)$  : Promoting piecewise smoothness via convex optimization. *PLOS One*, 15(8):e0237901, 2020.
- [2] N. BANNUR, V. SHAH, A. RAVAL et J. WHITE : Synthetic Data Generation for Improved covid-19 Epidemic Forecasting. *medRxiv*, pages 2020–12, 2020.
- [3] F. BRAUER, C. CASTILLO-CHAVEZ et Z. FENG : *Mathematical models in epidemiology*. Springer, New York, 2019.
- [4] A. CHAMBOLLE et T. POCK : A first-order primal-dual algorithm for convex problems with applications to imaging. *J. Math. Imaging Vis.*, 40(1):120–145, 2011.
- [5] A. CORI, N. M. FERGUSON, C. FRASER et S. CAUCHEMEZ : A new framework and software to estimate time-varying reproduction numbers during epidemics. *Am. J. Epidemiol.*, 178(9):1505–1512, 2013.
- [6] P. JACCARD : Distribution de la flore alpine dans le bassin des Dranses et dans quelques régions voisines. *Bull. Soc. Vaudoise Sci. Nat.*, 37:241–272, 1901.
- [7] Q.-H. LIU, M. AJELLI, A. ALETA, S. MERLER, Y. MORENO et A. VESPIGNANI : Measurability of the epidemic reproduction number in data-driven contact networks. *PNAS*, 115(50):12680–12685, 2018.
- [8] B. PASCAL, P. ABRY, N. PUSTELNIK, S. ROUX, R. GRIBONVAL et P. FLANDRIN : Nonsmooth convex optimization to estimate the Covid-19 reproduction number space-time evolution with robustness against low quality data. *IEEE Trans. Signal Process.*, 70:2859–2868, 2022.
- [9] J. WALLINGA et P. TEUNIS : Different Epidemic Curves for Severe Acute Respiratory Syndrome Reveal Similar Impacts of Control Measures. *Am. J. of Epidemiol.*, 160(6):509–516, 09 2004.

<sup>3</sup>Pour  $\mathbf{X}, \mathbf{Y} \in \{0,1\}^T$ ,  $\text{Jac}^g(\mathbf{X}, \mathbf{Y}) := \frac{\sum_{t=1}^T \sqrt{\mathbf{X}_t^g \times \mathbf{Y}_t^g}}{\sum_{t=1}^T \mathbf{X}_t^g + \mathbf{Y}_t^g - \sqrt{\mathbf{X}_t^g \times \mathbf{Y}_t^g}}$   
où  $\mathbf{X}^g$  est la convolution de  $\mathbf{X}$  avec la fenêtre de lissage  $\mathbf{g}$ .