



Clustering des trajectoires thérapeutiques pour la bronchopneumopathie chronique obstructive (BPCO) : comparaison des méthodes d'encodage de séquences

Romane Péan\*<sup>1</sup>, Nina Temam<sup>1</sup>, Marie Génin<sup>1</sup>, Diane Vincent<sup>1</sup>, Rachel Nadif<sup>2</sup>, Sofiane Kab<sup>3</sup>, Nicolas Roche<sup>4</sup>, and Pauline Guilmin<sup>1</sup>

<sup>1</sup>Quinten Health – Quinten – 8 rue Vernier, 75017, Paris, France

<sup>2</sup>Centre de recherche en épidémiologie et santé des populations (CESP) – Assistance publique - Hôpitaux de Paris (AP-HP), Institut National de la Santé et de la Recherche Médicale – 16 avenue Paul Vaillant Couturier 94807 Villejuif Cedex, France, France <sup>3</sup>Unité Cohortes épidémiologiques en population (Constances) – Institut National de la Santé et de la Recherche Médicale - INSERM – 16 Av. Paul Vaillant Couturier, 94800 Villejuif, France

<sup>4</sup>Institut Cochin (IC UM3 (UMR 8104 / Ú1016)) – Institut National de la Santé et de la Recherche Médicale, Université Paris Cité – 22 rue Méchain, 75014 Paris, France

Cette étude a été réalisée dans le cadre d'un partenariat multiprojets avec Constances par le biais d'une convention Inserm « Inserm transfert ».



# Analyse de séquences thérapeutiques : défis méthodologiques

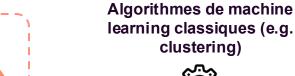
Comprendre la prise en charge des patients





Utilisation des RWD (séquences thérapeutiques / événements médicaux)









Les algorithmes de ML ne savent pas gérer la notion de séquence ou ordre temporel entre les données

# Analyse de séquences thérapeutiques : défis méthodologiques

Comprendre la prise en charge des patients



Utilisation des RWD (séquences thérapeutiques / événements médicaux)



Encodage des données longitudinales en vecteurs numériques



Algorithmes de machine learning classiques (e.g. clustering)





OBJECTIF: Comparer trois méthodes d'encodage de séquences combinées au clustering & évaluer leur capacité à regrouper les trajectoires thérapeutiques des patients

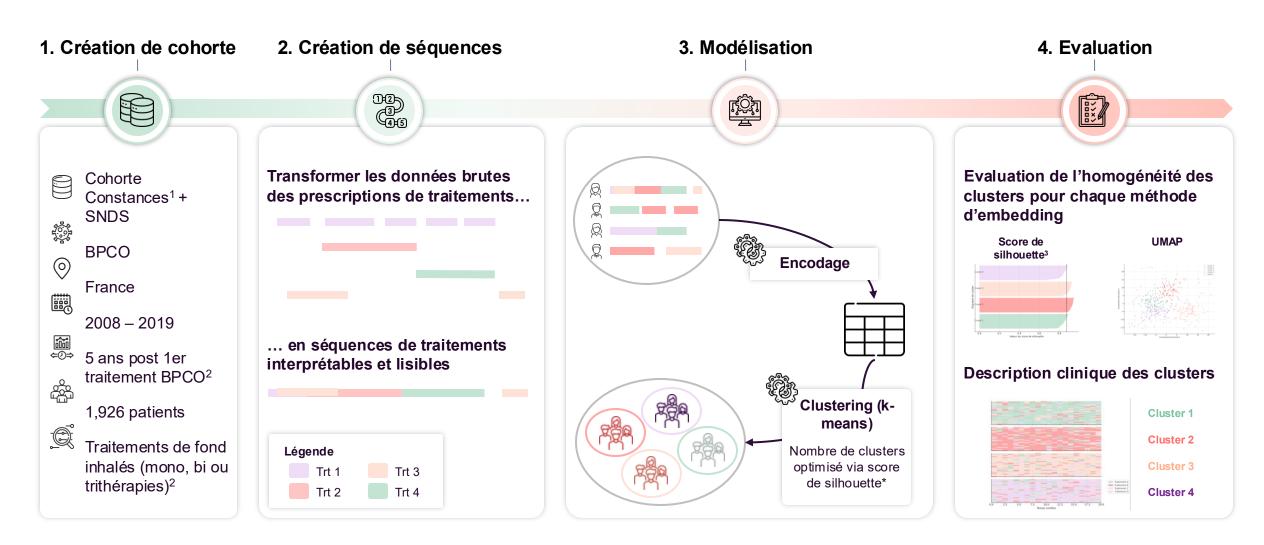








Résultats



<sup>1</sup> Zins, M.; Goldberg, M. The French CONSTANCES Population-Based Cohort: Design, Inclusion and Follow-Up. Eur J Epidemiol 2015, 30, 1317–1328. https://doi.org/10.1007/s10654-015-0096-4 (Cette étude a été réalisée dans le cadre d'un partenariat multiprojets avec Constances par le biais d'une convention Inserm Inserm transfert.)

<sup>2</sup> Mirza, S.; Clav, R. D.; Koslow, M. A.; Scanlon, P. D. COPD Guidelines: A Review of the 2018 GOLD Report. Mayo Clin Proc 2018, 93 (10), 1488-1502. https://doi.org/10.1016/j.mayocp.2018.05.026 (Cortico stéroïdes inhalés (ICS), Bêta-ago niste à longue durée d'action (LABA). Anticholiner gique à longue du rée d'action (LAMA))

<sup>3</sup> Rousseeuw, Peter. (1987). Rousseeuw, P.J.: Silhou ettes: A Graphical Aid to the Interpretation and Validation of Cluster Analysis. Comput. Appl. Math. 20, 53-65. Journal of Computational and Applied Mathematics. 20. 53-65. 10.1016/0377-0427(87)90125-7. \* en utilisant la méthode du coude

### Création de cohorte et séquences







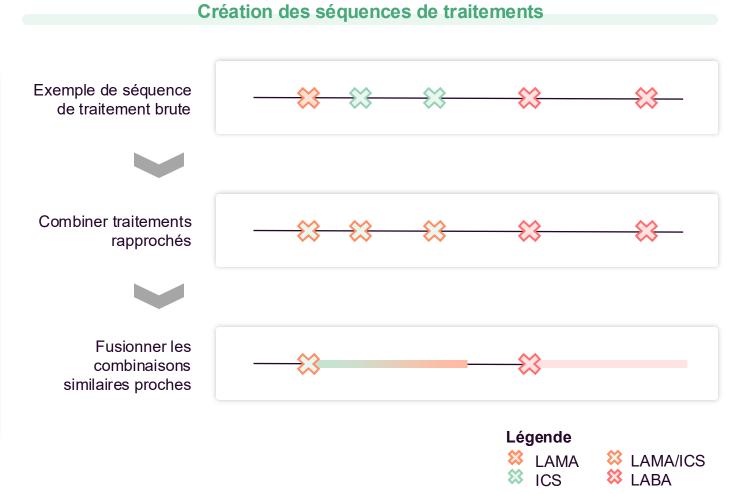


#### Patients BPCO inclus dans la cohort Constances<sup>1</sup> via OU Questionnaire de Test de diagnostic auto **Spirométrie** rempli Patients suivis pendant au moins 5 ans après leur premier traitement 1ère date Demière date active dans le active dans le Inclusion dans 1er trt BPCO2 **SNDS SNDS** Constances

≥5 ans

Aucun traitement pendant ≥1 an

Création de la cohorte



<sup>1</sup> Zins, M.; Goldberg, M. The French CONSTANCES Population-Based Cohort: Design, Inclusion and Follow-Up. Eur J Epidemiol 2015, 30, 1317–1328. https://doi.org/10.1007/s10654-015-0096-4 (Cette étude a été réalisée dans le cadre d'un partenariat multiprojets avec Constances par le biais d'une convention Inserm Inserm transfert.) 2 Traitements de fond inhalés (mono, bi ou trithérapies): Cortico stéroïdes inhalés (ICS), Bêta-agoniste à longue durée d'action (LABA), Anticholinergique d'action (LABA), anticholinergique d'action (LABA), anticholinergique d'action (LABA), anticholinergique d'action (LABA), anticholinerg

# 3 méthodes d'encodage des séquences



Données temporelles brutes







Vecteurs numériques utilisables par les algorithmes de MI





### (A) SegMining

Extraction de sous-séquences récurrentes via l'algorithme SPADE<sup>1</sup>.

Génération de vecteurs binaires indiquant la présence ou l'absence de motifs spécifiques.

#### Représentation numérique

PTID 1	1	1	1	1
PTID 2	0	1	0	0
PTID 3	1	0	0	0
PTID 4	0	1	1	0



#### (B) SeqToChar

Représentation des séquences sous forme de chaînes de caractères.

PTID 1	A – F – J – B – F	
PTID 2	B — J — J	
PTID 3	A – B	
PTID 4	J — F	

Calcul d'une matrice de similarité (Jaro<sup>2</sup>) par paires.

#### Représentation numérique

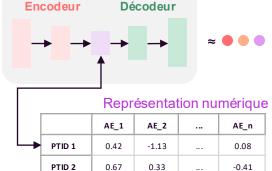
	PTID 1	PTID 2	PTID 3	PTID 4
PTID 1	1	0.63	0.43	0.97
PTID 2	0.63	1	0.24	0.33
PTID 3	0.43	0.24	1	0.12
PTID 4	0.97	0.33	0.12	1



### (C) Autoencodeur

Modèle d'apprentissage profond apprenant des représentations continues abstraites (vecteurs) des séquences en prenant en compte leur durée, dans un espace latent de dimension réduite<sup>3</sup>





<sup>1</sup> Zaki, M. J. SPADE: An Efficient Algorithm for Mining Frequent Sequences. 2021.

# Résultats – SeqToChar présente les clusters les plus homogènes



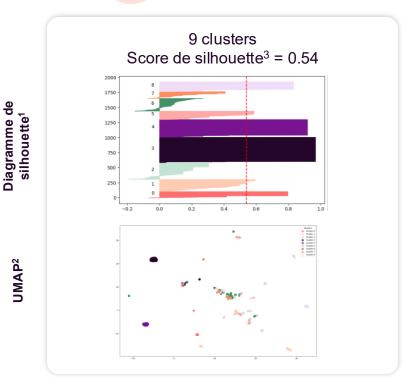






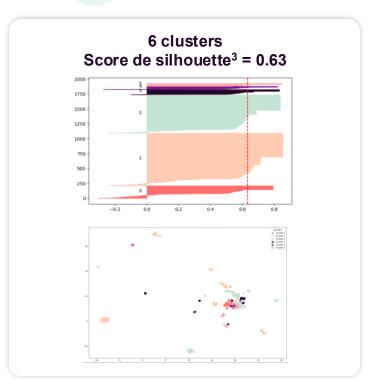
Interprétation du score de silhouette : (Excellent : 0,71-1 ; Bon : 0,51-0,7 ; Moyen : 0,26-0,5 ; Faible : 0-0,25 ; Mauvais : < 0)



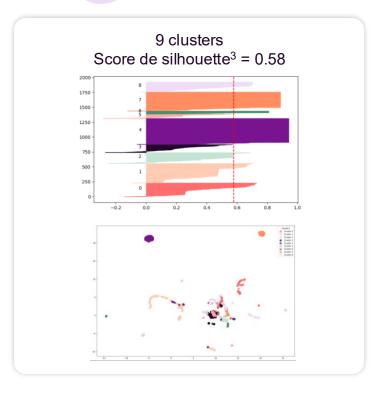


fix learning rate=0.01; starter learning rate=None; decay rate=0.1; decay steps=5)









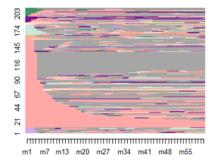
<sup>1</sup> Rousseeuw, Peter. (1987). Rousseeuw, P.J.: Silhou ettes: A Graphical Aid to the Interpretation and Validation of Cluster Analysis. Comput. Appl. Math. 20, 53-65. Journal of Computational and Applied Mathematics. 20. 53-65. 10.1016/0377-0427(87)90125-7. 2 McInnes, L., Healy, J., Melville, J. UMAP: Uniform Manifold Approximation and Projection for Dimension Reduction. arXiv September 18, 2020. https://doi.org/10.48550/arXiv.1802.03426 3 Kaufman & Rousseeuw (1990), Finding Groups in Data, pp. 79-81 - (Excellent: 0.71-1: Bon: 0.51-0.7: Moven: 0.26-0.5: Faible: 0-0.25: Mauvais: < 0) Hyperparametres: Seq Mining (min support=0.3); Seq ToChar (event same visit strategy="alphabetical order"; edit distance="jaro"); Autoencoder object (hidden dim2=32; hidden dim2=32); Autoencoder embeddings (epochs=10; val split=0.2;

### Résultats - description clinique des clusters

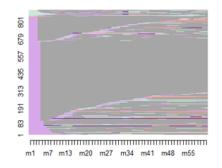




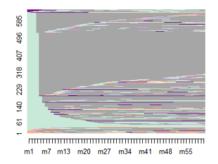
### Cluster (N=211)



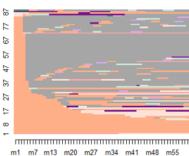
#### Cluster (N=882)



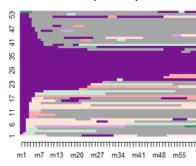
### Cluster (N=642)



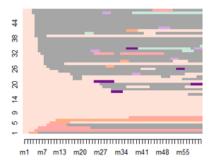
#### **Cluster4** (**N**=88)



Cluster **(N=54)** 



Cluster **6** (N=49)



- Légende

  ICS ICS/LABA ICS/LAMA

  LAMA LABA/LAMA Not treated

  LABA ICS/LABA/LAMA
- Trajectoires de traitement hétérogènes, nombreux arrêts précoces sans relais thérapeutique → problème d'adhérence connu
- Clusters principalement regroupés selon le traitement initial
- Cluster 2: majoritairement initié par une monothérapie ICS (non recommandée en usage isolé), semblant correspondre à un usage ponctuel
- Cluster **6** et Cluster **6** : débutant par une trithérapie ICS/LABA/LAMA ou une bithérapie LABA/LAMA, présentant les trajectoires les plus durables dans le temps
- Cluster 3: initié par ICS/LABA, moindre persistance au traitement (mais avec relais partiel vers d'autres catégories de traitement)
- Cluster 1 et Cluster 2 : débutant par des monothérapies LABA ou LAMA présentant une faible continuité dans le temps, mais transitions fréquentes vers d'autres catégories de traitement

Corticostéroïdes inhalés (ICS), Bêta-agoniste à longue durée d'action (LABA), Anticholinergique à longue durée d'action (LAMA)

### Discussion

#### Des résultats prometteurs

- √ Groupes homogènes sur le traitement initial
- ✓ Pas de différence majeure de performance observée entre les 3 méthodes d'encodage
- √ Cohérence (score de silhouette de 0.63) et séparation visible des clusters
- ✓ Reflètent problèmes d'adhérence connus dans la BPCO

#### Limites

- ~ Durée de traitement imputée → biais
- Diagnostic de BPCO posé avant ou après le début du traitement (n.b. exclusion des patients asthmatiques)
- ~ Clustering sur le premier événement

#### Pour aller plus loin

- > Explorer d'autres méthodes de clustering et d'encodage
- Enrichir le clustering avec des variables additionnelles
- > Trouver un **équilibre** entre performance statistique et pertinence médicale dans le choix de la méthode



Des questions?

### Bibliographie

- Zins, M.; Goldberg, M. The French CONSTANCES Population-Based Cohort: Design, Inclusion and Follow-Up. Eur J Epidemiol 2015, 30, 1317–1328. https://doi.org/10.1007/s10654-015-0096-4
- Mirza, S.; Clay, R. D.; Koslow, M. A.; Scanlon, P. D. COPD Guidelines: A Review of the 2018 GOLD Report. Mayo Clin Proc 2018, 93 (10), 1488–1502. <a href="https://doi.org/10.1016/j.mayocp.2018.05.026">https://doi.org/10.1016/j.mayocp.2018.05.026</a>
- 3. Rousseeuw, Peter. (1987). Rousseeuw, P.J.: Silhouettes: A Graphical Aid to the Interpretation and Validation of Cluster Analysis. Comput. Appl. Math. 20, 53-65. Journal of Computational and Applied Mathematics. 20. 53-65. 10.1016/0377-0427(87)90125-7
- 4. Zaki, M. J. SPADE: An Efficient Algorithm for Mining Frequent Sequences. 2021.
- 5. Jaro, M. A. Advances in Record-Linkage Methodology as Applied to Matching the 1985 Census of Tampa, Florida. Journal of the American Statistical Association 1989, 84 (406), 414–420. https://doi.org/10.1080/01621459.1989.10478785
- Baytas, I. M.; Xiao, C.; Zhang, X.; Wang, F.; Jain, A. K.; Zhou, J. Patient Subtyping via Time-Aware LSTM Networks. In Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining; ACM: Halifax NS Canada, 2017; pp 65–74. <a href="https://doi.org/10.1145/3097983.3097997">https://doi.org/10.1145/3097983.3097997</a>.
- 7. Rousseeuw, Peter. (1987). Rousseeuw, P.J.: Silhouettes: A Graphical Aid to the Interpretation and Validation of Cluster Analysis. Comput. Appl. Math. 20, 53-65. Journal of Computational and Applied Mathematics. 20. 53-65. 10.1016/0377-0427(87)90125-7.
- 8. McInnes, L.; Healy, J.; Melville, J. UMAP: Uniform Manifold Approximation and Projection for Dimension Reduction. arXiv September 18, 2020. <a href="https://doi.org/10.48550/arXiv.1802.03426">https://doi.org/10.48550/arXiv.1802.03426</a>
- 9. Kaufman & Rousseeuw (1990), Finding Groups in Data, pp. 79–81
- 10. Guo, Y.; Guo, S.; Jin, Z.; Kaul, S.; Gotz, D.; Cao, N. Survey on Visual Analysis of Event Sequence Data. arXiv June 25, 2020. http://arxiv.org/abs/2006.14291 (accessed 2023-11- 23)

ICS/LABA LABA/LAMA

ICS/LABA/LAMA

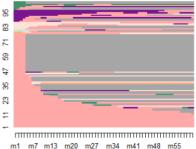
ICS/LAMA

Not treated

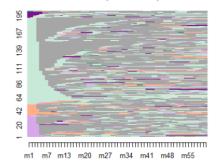
### Résultats - description clinique des clusters



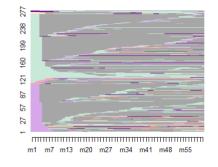




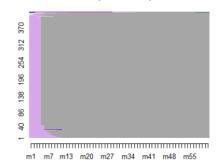
Cluster (N=202)

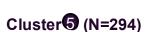


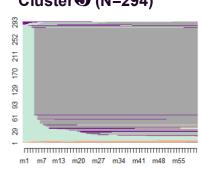
Cluster (N=278)



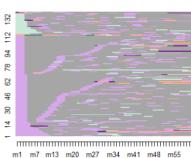
Cluster (N=422)



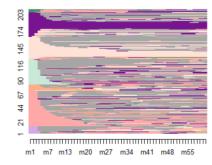




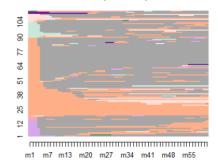
Cluster 6 (N=140)



Cluster (N=49)



Cluster (N=49)



Cluster (N=49)

Légende

