



Re-identifying cancer treatment lines in real-world data
an innovative sequence alignment algorithm



MYLORD

Objective

Case study in
hematology:
Multiple Myeloma (MM)

HEVA



MM disease history

succession of remission
and relapse, which form
treatment line

**How many patients
under each line?**



MM treatment

treatment line are treated
by protocols

**In which line is a
protocol given?**



MM care evolution

personalized medicine
thanks to fast evolving
recommendations

**How does the patient
management change
year after year?**

Automatically describe the treatment lines received
by MM patients in France from the SNDS

Why do we need AI?

To learn the matches between medical theory and Real-World Data, despite the variations



Medical theory

≠



Patient reality

≠



Real-World Data

VTD cycle

One part of a patient's treatment sequence

SNDS: French claim databases

 D1 - D4 - D8 - D11

 D1 - D4 - D8 - D13

 68 million people


 D1

 D1 - D13

Completeness of reimbursement

 D1

 D1 - D13

 Unnamed injections
For drugs, only date of sale is available

Complex protocols

Deviations from theoretical protocols

Why do we need AI?

Because of the number of patients and of protocols

HEVA



Identify theoretical protocols

Numerous protocols

41 cycles

VS.



**Cohort selection
(2014 - 2017)**

Thousands of patients

17 442 patients



ATLAS

Analysis of Treatment Lines
using Alignment of Sequences

By HEVA

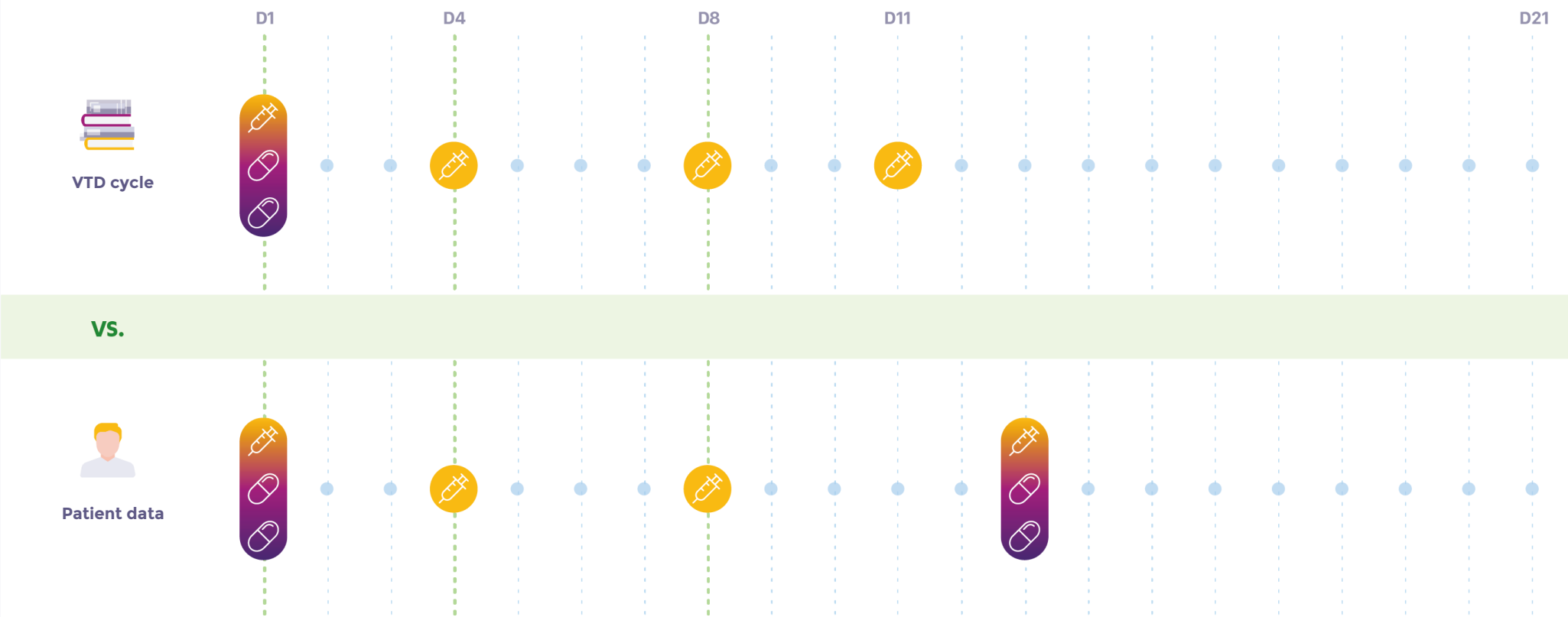


1 cycle  1 vs.  2 1 part of the patient's sequence

HEVA

3 4

Alignment and score computation



 Velcade + Thalidomide + Déxaméthasone  Velcade

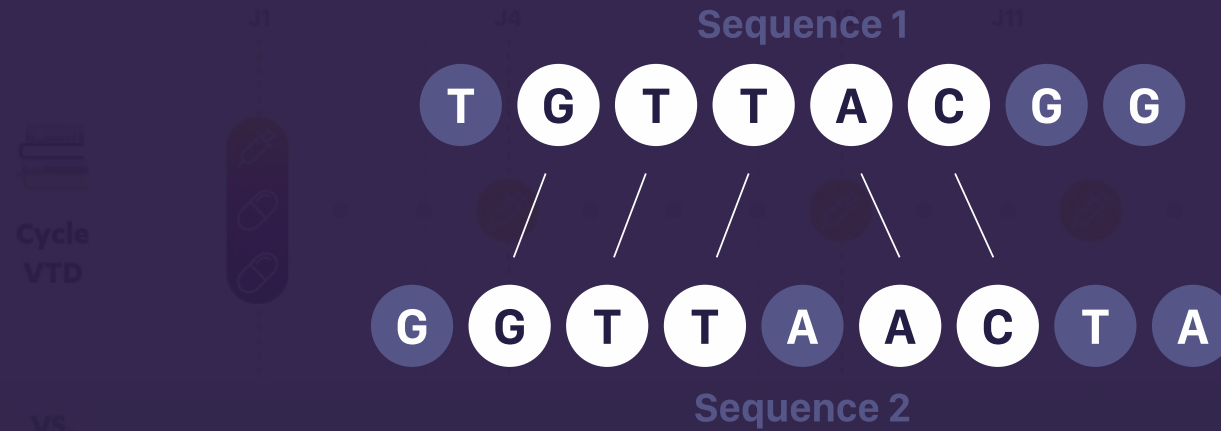


1 cycle vs. 1 part of the patient's sequence

Method

Which sections are the most similar?

Alignement et calcul du score



Adaptation of the Smith Waterman algorithm,
used for DNA sequences alignment



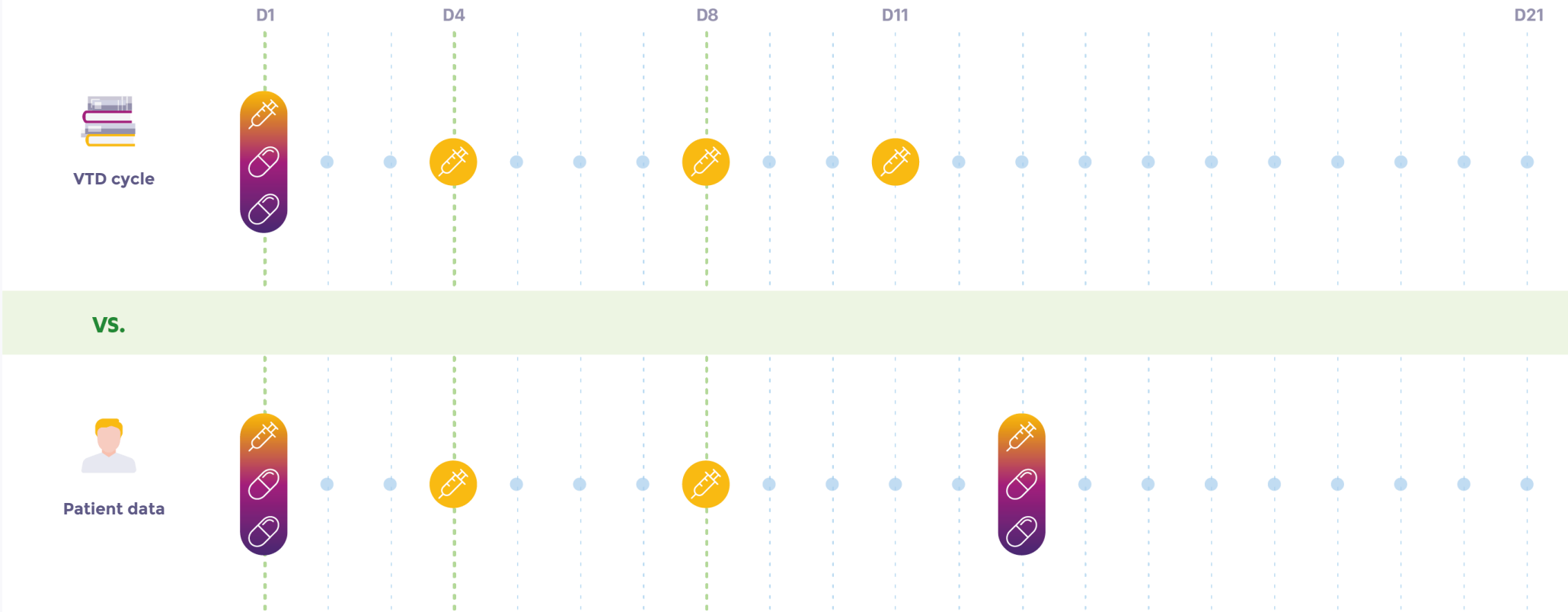
1 cycle  1 vs.  2 1 part of the patient's sequence

HEVA

3

4

Alignment and score computation



Velcade + Thalidomide + Déxaméthasone



Velcade

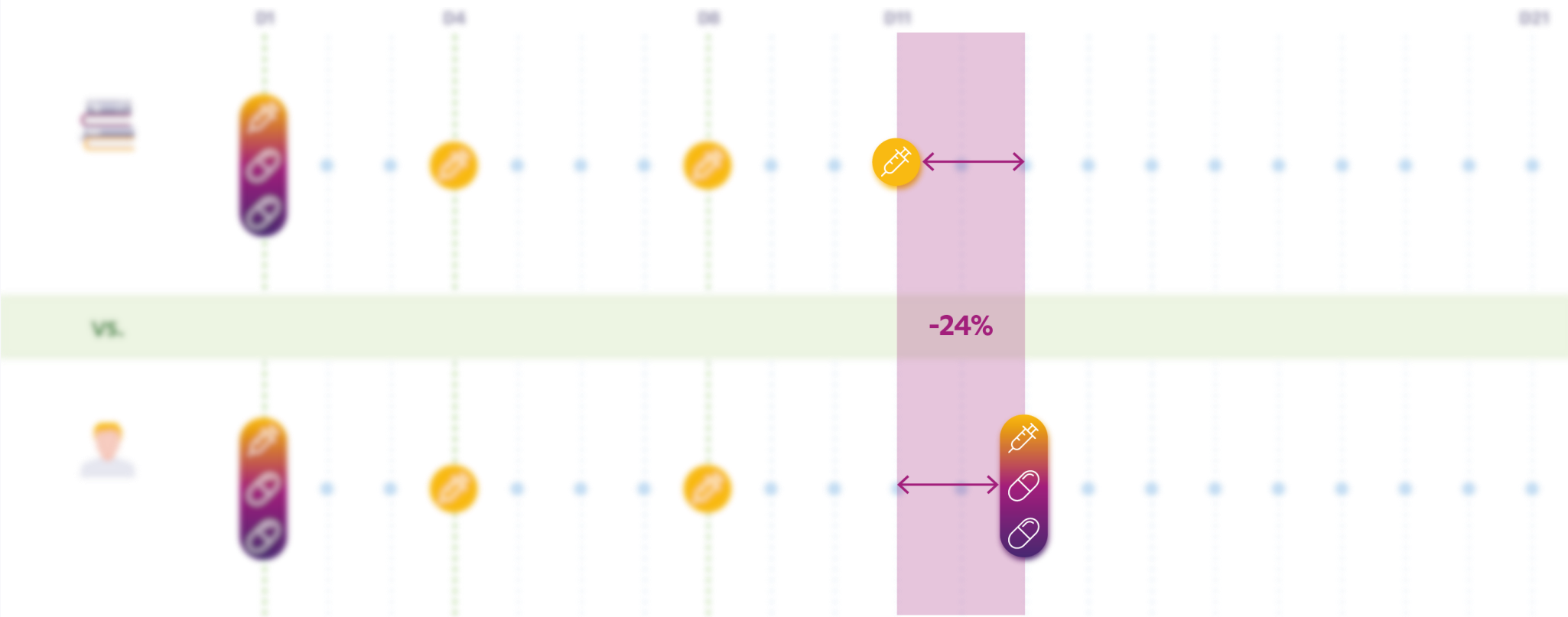


1 cycle  1 vs.  2 1 part of the patient's sequence

HEVA

3 4

Alignment and score computation



Velcade + Thalidomide + Dexaméthasone



Velcade



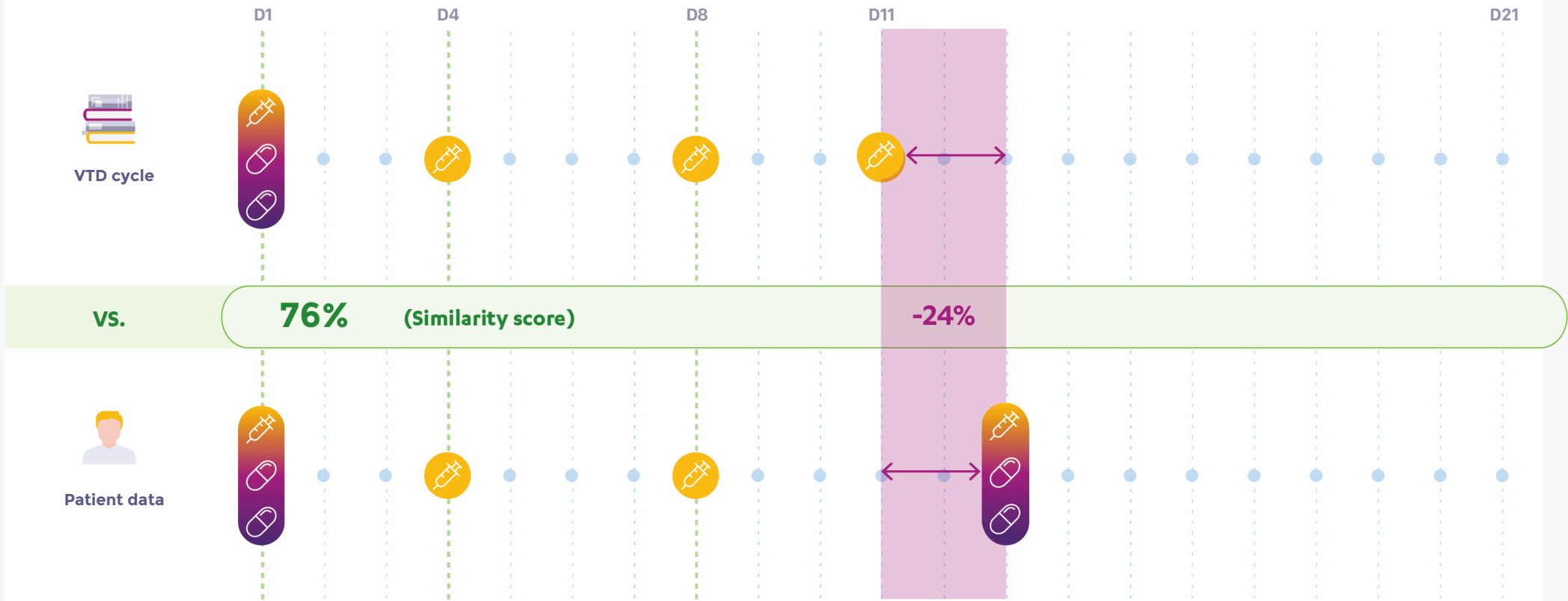
1 cycle  1 vs.  2 1 part of the patient's sequence

HEVA

3

4

Alignment and score computation





Données patient

Données patient



All cycles



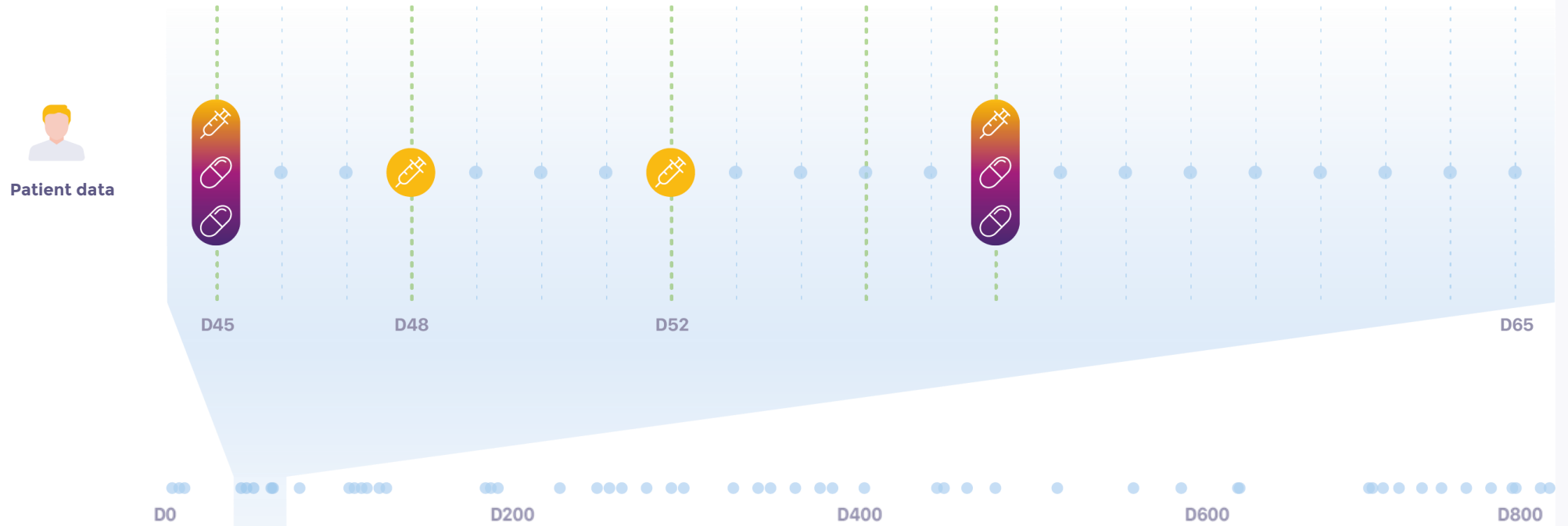
vs.



1 patient

HEVA

Scan the entire patient's history (800 days)



Analyze of the repetition of cycles over time
→ Reconstruction of complete treatment lines



All cycles



vs.



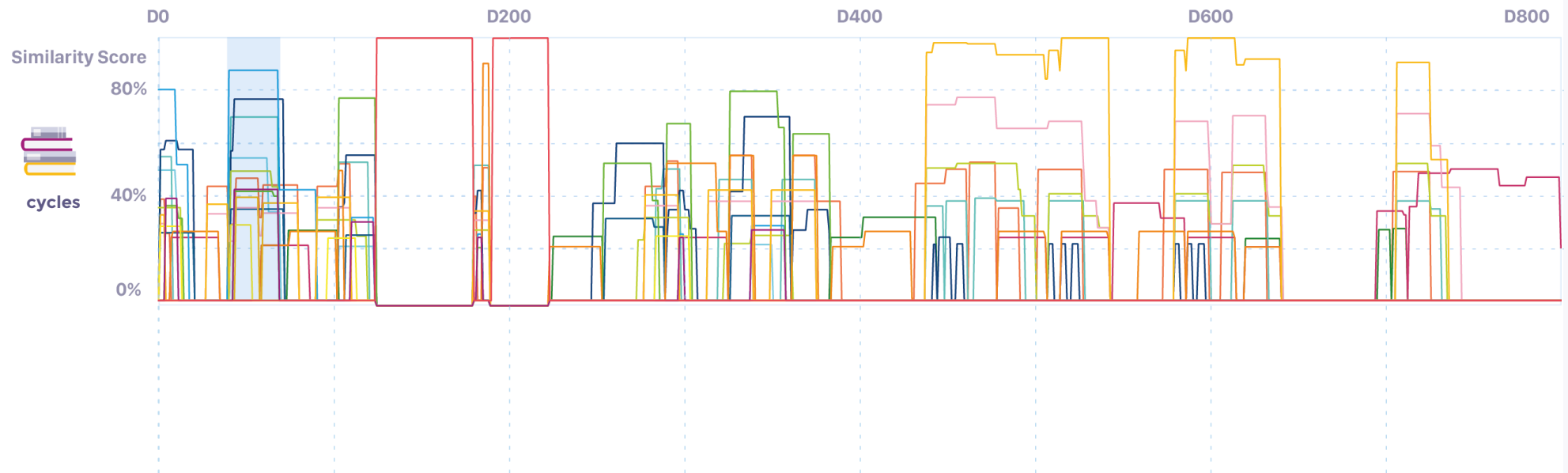
1 patient

HEVA

5

Display all the similarity scores between a patient and each of the theoretical cycles

- VTD
- Pom Dex
- Dara Dex
- Dara Vd
- VTD Dara
- Dara Rd
- Pom Dex Dara
- Pom Vel Dex
- MPV
- VRD
- BVD
- Rd
- MP
- V
- Pano Vd
- ASCT

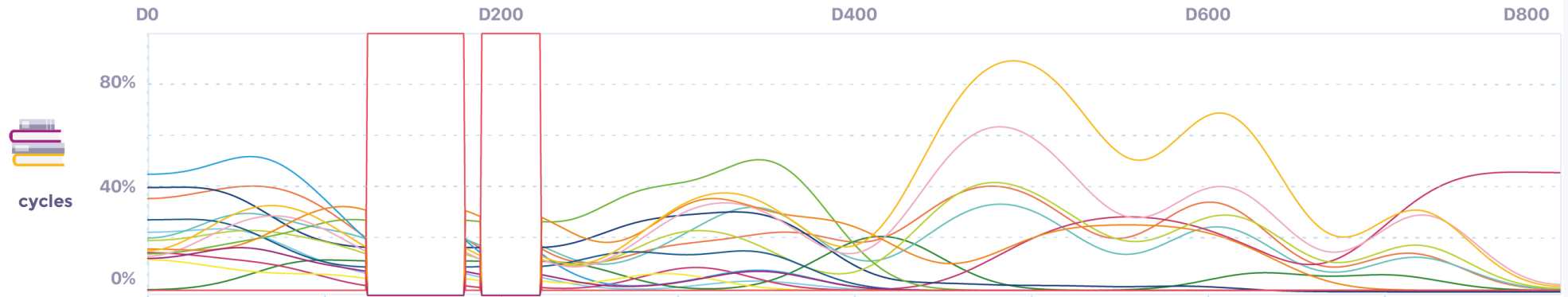




6

Smooth the scores over time

- VTD
- Pom Dex
- Dara Dex
- Dara Vd
- VTD Dara
- Dara Rd
- Pom Dex Dara
- Pom Vel Dex
- MPV
- VRD
- BVD
- Rd
- MP
- V
- Pano Vd
- ASCT

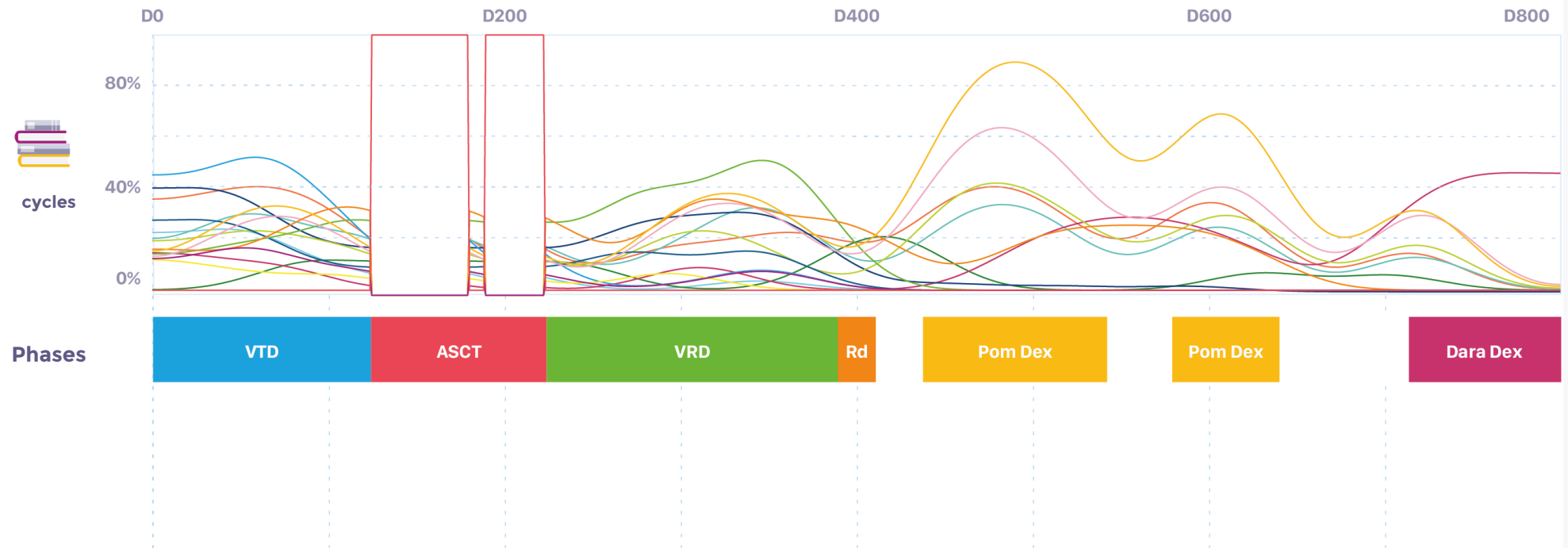




6

Smooth the scores over time and identify the best phase at each time of the follow-up

- VTD
- Pom Dex
- Dara Dex
- Dara Vd
- VTD Dara
- Dara Rd
- Pom Dex Dara
- Pom Vel Dex
- MPV
- VRD
- BVD
- Rd
- MP
- V
- Pano Vd
- ASCT

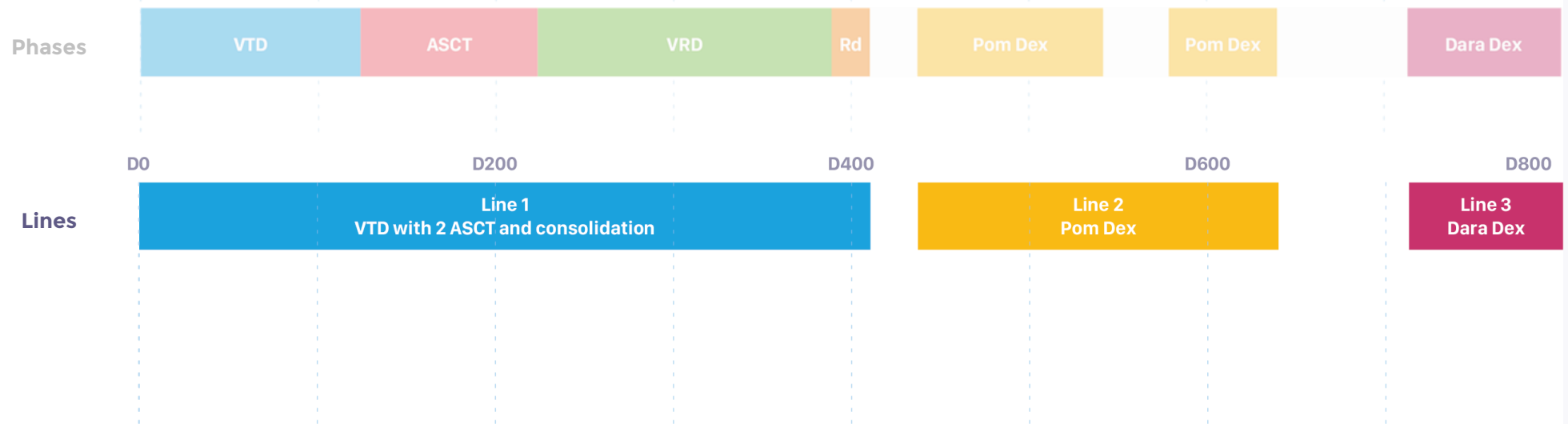




7

Gather phases into medical lines

- VTD
- Pom Dex
- Dara Dex
- Dara Vd
- VTD Dara
- Dara Rd
- Pom Dex Dara
- Pom Vel Dex
- MPV
- VRD
- BVD
- Rd
- MP
- V
- Pano Vd
- ASCT





Résultats 7 lignes pour le patient

Go from each individual patient result

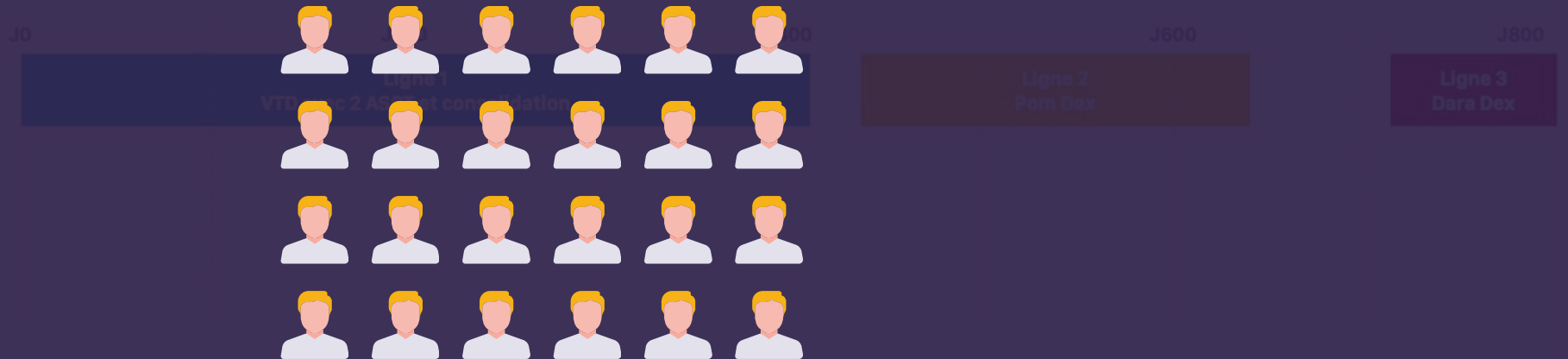
to

smart data-visualizations of the entire cohort

- VTD
- Pom Dex
- Dara Dex
- Dara Vd
- VTD Dara
- Dara Rd
- Pom Dex Dara
- Pom Vel Dex
- MPV
- VRD
- BVD
- Rd
- MP
- V
- Pano Vd
- ASCT

Patients

Lignes



Cohort-level results

Empower dataviz'

Restitution



Tables



Histograms



Sunburst



Kaplan-Meier

Confrontation with the clinical knowledge

to improve the algorithm



Counts of patients, by lines, by protocols, by year



Duration of lines, of inter-lines

Results

Re-identification
power of the algorithm



1st execution



77 %
of the patients

Have all their lines
identified



2 cycles identified
a posteriori were added

Improvement
of the parameterization

Improvement of medical
understanding



2nd execution



92 %
of the patients

Have all their lines
identified

Advantages of the AI approach

HEVA



**Flexibility and
performant on RWD**



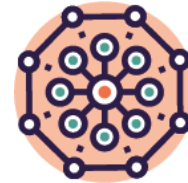
**Custom-made
for each disease**



**Confidence scores
allow to work with
several iterations, to
obtain high-quality
results**



Therapeutic **innovations** arise from new treatment sequences



Pharmaceutical companies and public institutions face the challenge to find **such sequences**



In our experience, **algorithms and data** have a role to play



Thank you for your attention



MYLORD

Scientific article on the ATLAS algorithm

ATLAS: a Robust Algorithm for Temporal Sequence Alignment of Treatment Lines using Claim Databases

Martin Prodel
dept. of Data Analytics
HEVA
Lyon, France
mprodel@hevaweb.com

Ludovic Lamarsalle
dept. of Data Analytics
HEVA
Lyon, France
llamarsalle@hevaweb.com

Vincent Augusto
Mines Saint-Etienne, Univ. Clermont Auvergne
CNRS, UMR 6158 LIMOS
Center for Healthcare Engineering
F-42023 Saint-Etienne, France
augusto@emse.fr

Abstract—Comparison of cancer treatment protocols against patient history is important to assess the efficiency of chemotherapy protocols. However, manual identification of protocols is not possible when considering a large population of patients. This paper proposes a new method called ATLAS (Analysis of Treatment Lines using Alignment of Sequences) to tackle the problem related to protocol identification using claim databases as data input. The proposed algorithm is an extension of the Smith-Waterman algorithm and allows to compare any patient's medical history against a list of theoretical protocols taking into account temporal information in sequences as well as missing data. Numerical experiments show that the proposed method could identify the right protocol in over 95% of the cases for realistically noised sequences (15,000 generated patients aligned with 15 protocols). The method is meant to be used as a decision aid tool for practitioners.

Index Terms—sequence alignment, temporal Smith-Waterman, claim data, healthcare, oncology, treatment lines

Claim databases, Electronic Health Records (EHR) or Electronic Health Data (EHD), contain the information found on medical bills or claims, and are collected by insurers (public or private). Although claim databases are collected for reimbursement purposes, they contain relevant information on patients' medical condition and on delivered care. The increasing availability of claim data seems promising to lead epidemiological studies on how well treatment protocols are followed [5]. However, using data to identify which protocol a patient has followed is not straightforward. First, the exact administration date of some drugs might be unknown, only the delivery date in a pharmacy is known. The correct matching of such events to their theoretical administrations is uncertain. The second hindrance is missing data. Databases sporadically lack the information about a drug given to a patient, even if a chemotherapy session is registered. Those omitted administration blur the identification of the protocol, for which all the drugs are supposedly given at precise dates.

I. INTRODUCTION

A. Context

Chemotherapy is a type of cancer treatment that uses one or more anti-cancer drugs as part of a standardized chemotherapy regimen. Chemotherapy may be given with a curative intent, or it may aim to prolong life or to reduce symptoms (palliative chemotherapy). A regimen is a systematic plan of treatment administration designed to improve the health of the patient. The regimens recommended by the national guidelines are called protocols. These protocols are evaluated through clinical trials, but also a posteriori using historical patient data.

To assess the efficiency of chemotherapy regimens, epidemiologists need to identify the actual treatments received by a cohort of patients in large databases [1]–[4]. It helps improve existing protocols and thus the quality of care. Protocol identification is possible thanks to the increasing availability of medical and claim data in recent years. Still, such identification is tedious because drug protocols are complex and numerous. Scientific challenges related to protocol identification are: co-occurrence of multi-drug treatments, data quality, temporal nature of the problem.

978-1-7281-1462-0/19/\$31.00 ©2019 IEEE

Existing works on protocol identification from claim databases mainly rely on sequence alignment algorithms. Those algorithms were made popular and widespread in the field of bioinformatics [6]. Sequence alignment is used to determine similar regions between two strings of DNA also named sequences. When aligning two sequences, one is the sequence in which we search similarities and the other is the "reference" sequence. Two classical and proven algorithms for sequence alignments are the Needleman-Wunsch [7] and Smith-Waterman [8] algorithms. Both are based on dynamic programming and use a scoring matrix. Each value of this matrix represents the similarity between two elements, one from the sequence to align and the other from the reference sequence. The higher the score, the more the elements are similar.

These alignment methods are meant to identify shared patterns in two sequences, but they do not consider time between elements of a sequence, nor the uncertain position of such elements (e.g. uncertain administration date of drugs).

Poster (in French) on the ATLAS algorithm

MYLORD ATLAS, nouvelle méthode d'analyse des lignes de traitements à partir du SNDS
Exemple de l'étude MYLORD, sur les patients français atteints du Myélome Multiple

Martin Prodel*, Ludovic Lamarsalle*, Matthieu Javelot*, Gabriel Gagané*, Marie Perre*, Vincent Augusto*, Ludovic Lamarsalle*, Benny Raguideau*, Martin Prodel*
*Recherche associée à l'Université Clermont Auvergne, CNRS, UR 2016 - Mines Saint-Etienne, Univ. Clermont Auvergne, CNRS, UMR 6158 LIMOS, Centre CIC, F-42023 Saint-Etienne, France, **Mines Saint-Etienne, Univ. Clermont Auvergne, CNRS, UMR 6158 LIMOS, Centre CIC, F-42023 Saint-Etienne, France

Objectif
Décrire de manière automatique les lignes de traitements reçues par les patients atteints d'un Myélome Multiple (MM) en France à partir du SNDS.

Contexte
Le MM est une hémopathie maligne qui touche environ 15000 personnes par an en France. Au cours de cette maladie, on recense des épisodes de thérapies et de rechutes prises en charge par des protocoles de soins consistant de combinaisons de médicaments et/ou greffe de cellules souches hématopoïétiques. Ces protocoles sont constitués d'une ou plusieurs phases. Chaque phase étant une répétition de cycles. A chaque nouvelle ligne de traitement, le protocole de soins est adapté par les cliniciens en regard du profil du patient et/ou des effets secondaires.
Les données du SNDS, couplées à l'algorithme ATLAS, permettent d'identifier ces lignes de traitements.

Challenges
Détection et sous-identification complexes

Théorie
Protocoles complexes (plusieurs médicaments combinés, avec des timings différents)
Protocoles nombreux

Pratique
Distinction par rapport aux protocoles théoriques (sédation des traitements, effets secondaires, etc.)
Des milliers de patients

Data
Les données de remboursements sont volumineuses
Traitements sous CHS nombreuses
Essais cliniques

Méthode : ATLAS, un algorithme d'intelligence artificielle en 8 étapes

1. Reconnaissance de tous les cycles constituant les protocoles théoriques, ainsi que les données de cycles reçus par 20 protocoles théoriques
2. Sélection de la cohorte de patients atteints d'un MM à partir du SNDS, recherche de l'historique médical de chaque patient (traitements, hospitalisations et sécs) comparant les cycles théoriques à 15 000 patients réels (CHS - 2017)
3. Pour chaque patient :
 - 3.1. Alignement des données médicales administratives du patient avec les cycles théoriques existants
 - 3.2. Calcul d'un score de similarité pour chacun de ses alignements
 - 3.3. Sélection d'un alignement de l'algorithme d'alignement de séquences de Smith-Waterman
 - 3.4. Visualisation des scores de similarité retrouvés au cours du temps
 - 3.5. Choix des scores optimaux pour définir des phases
4. Phase consistant de tester plusieurs dates de début
5. Utilisation d'une estimation de durée à chaque session sur chaque cycle, afin de fixer les scores dans le temps
6. Passage des phases aux lignes
7. Ligne consistant de passer les phases alignées à un état d'alignement compatible avec les données de sécs
8. Utilisation de la connaissance médicale des experts cliniciens

Exemple : des cycles V1D peuvent former l'induction pour préparer une greffe (SNDS), suite de cycles V1D et D.

Sur l'ensemble de la cohorte :

1. Agrégation des résultats
2. Les étapes 3 à 7 se font indépendamment pour chaque patient. Les résultats sont finalement agrégés sur l'ensemble de la cohorte (nombre d'hospitalisations, distributions, Kaplan-Meier)

L'ensemble a été réduit à 8 fois, afin d'améliorer chaque étape pour qu'elle soit adaptée à la prise en charge réelle, afin que les résultats soient les plus précis possible.

Résultats
Vie exécution de l'algorithme

- 95% des patients ont des lignes de traitements identifiées
- Ajout des 2 cycles identifiés à la prescription
- Amélioration de la prescription
- Amélioration de la prescription médicale

2ème exécution de l'algorithme

- 92% des patients ont des lignes de traitements identifiées
- Amélioration de la prescription
- Amélioration de la prescription médicale
- Amélioration de la prescription

Conclusion

Méthode innovante
Algorithme d'alignement de séquences d'ADN adapté à la recherche de séquences de soins

Verrous scientifiques levés
Flexibilité maximale aux étapes
Automatisation

Application médicale
Vieilles données de remboursements dans le domaine du myélome multiple (hospitalisations levées)

Validité de la prise en charge de MM peut être décrite, ouvrant la possibilité d'une meilleure description des séquences thérapeutiques sur un grand nombre de données patients en vie réelle. Cette méthode peut également s'appliquer à la description des lignes thérapeutiques d'autres pathologies.

ATLAS: Colloque Doctoral de Santé en vie réelle, 9 septembre 2019, Cité Internationale Universitaire de Paris 10^{ème} France



Our website



Our videos